

Small Area Estimation Strategy for the 2011 Census in England and Wales

Bernard Baffour^{a*}, Denise Silva^b, Alinne Veiga^b, Christine Sexton^c, James J. Brown^d

^a *School of Demography, Australian National University, Canberra, ACT 2601, Australia*

^b *National School of Statistical Sciences (ENCE/IBGE), Rio de Janeiro, Brazil*

^c *Office for National Statistics, Segensworth Road, Titchfield, PO15 5RR, UK*

^d *School of Mathematical & Physical Sciences, University of Technology Sydney, 15 Broadway, Ultimo NSW 2007, Australia*

*corresponding author: School of Demography, Australian National University, Canberra, ACT 2601, Australia. Tel: +61 26125 9030, email: bernard.baffour@anu.edu.au

Abstract. The use of model-based small area estimation methods for adjusting census results in the UK was first introduced in the 2001 Census. The aim was to obtain local level population estimates (local authority totals) by age-sex groups, adjusted for the level of undercount based on regression models combining results from the census and the Census Coverage Survey. A similar approach was adopted for the 2011 Census but with new features and this paper describes the work carried out to arrive at the chosen small area strategy. Simulation studies are used to investigate three proposed small area estimation methods: a local fixed effects model (the 2001 Census approach), a direct estimator and a synthetic estimator. The results indicate that both the synthetic and the local fixed models constitute good options to produce accurate and reliable local authority population estimates. A proposal is made to implement a small area estimation procedure that accommodates both the synthetic and local fixed models, as in some selected areas with differing local authority under-coverage rates a local fixed effects model may perform best. We examine this strategy under real census conditions based on the final results from the 2011 census.

Keywords: Census coverage, small area estimation, synthetic estimator, direct estimator

1. Introduction

The key purpose of a census is to produce accurate and reliable estimates of the population, not just at the national level but also, more importantly, for small areas. However, it is widely known that despite all the efforts of the census, some people will be missed [1] and it is standard practice to include an assessment of coverage within the census process. This is usually accomplished through a post-enumeration survey [2]. In the 2001 Census of England and Wales the Office for National Statistics (ONS) re-designed the post-enumeration survey, referred to as the Census Coverage Survey (CCS), to dramatically increase the sample size with a focus on coverage. The result was a large-scale survey designed to provide information that could be matched with the Census in order to estimate directly the age-sex structure of estimation areas (EAs), consisting of populations around 0.5 million individuals [3]. EAs were either a single large local authority (LA) or a contiguous group of smaller LAs. LAs are administrative units of local government and are primarily in charge of key services such as education, housing and social services. At the time of the 2011 Census, there were 348 local authorities in England and Wales and the census is often the main source of information about the population at such small geographies [4]. The same basic strategy was also implemented for Scotland and Northern Ireland within their EA and hard-to-count structures. The units of local administration in Scotland are known as council areas, of which there were 32 for the 2011 Census and in Northern Ireland they are known as districts, of which there were 26 for the 2011 Census. We refer to the ‘UK census’ as shorthand for the censuses in England & Wales, Scotland and Northern Ireland.

Population size and structure are key drivers in the allocation of funding to LAs from central government. Hence it is important that the census counts are adjusted for the estimated undercount to enable a fair and accurate allocation of resources. To facilitate this, the ideal would be a CCS designed to estimate the coverage of the age-sex population directly at LA level. However, like any other national statistical institute, the ONS faces the challenge of producing comprehensive, accurate and reliable information in a timely and cost-efficient manner. A CCS with sufficient sample size for direct estimation of all LAs would not only increase costs, but its size would potentially reduce the overall quality, as undertaking such a large data collection exercise very close to the census would be problematic. Therefore, it is necessary to turn to small area techniques [5] that allow the age-sex estimates for an individual LA to borrow strength from neighbouring LAs or neighbouring age-sex categories within the LA, while still attempting to reflect localised effects. In general, direct estimators (based only on the small CCS sample from within an LA) will be unbiased, but have large standard errors and so are imprecise. On the other hand, indirect methods, although more precise, can have large biases ([6]; [7]). For the 2001 Census, borrowing strength was achieved with the inclusion of LA specific fixed-effects within a collapsed version of the main estimation model used for EAs. Such an approach combined direct information from the specific LA with pooled information across the LAs within their EA.

Following reviews of the 2001 Census adjustment approach (see [8]; [9]), ONS adopted broadly the same strategy for the 2011 Census [2]. However, the 2001 Census provided substantially more data from which to develop the 2011 approach. This led to a change in the CCS design structure so that allocation to LAs was directly controlled in the design, stratification within LAs was based on more up-to-date information on the population structure and the allocation was driven by variation in coverage patterns actually observed in 2001 [10].

The result is that many of the city LAs, Coventry for example, that did not have a big enough population to count as an estimation area in 2001 are a single LA estimation area in the 2011 design. Conversely, the estimation areas that are aggregates of LAs tend to contain more LAs than in 2001 but with a stronger expectation that within estimation area homogeneity across the LAs can be achieved at estimation [11]. First, this is because the EAs are formed after the design stage so LAs can be aggregated, albeit still reflecting geographical contiguity, to take account of the observed patterns in coverage from 2001. Second, the move to post-out as well as post-back with flexible allocation of staff for non-response follow-up is expected to smooth out census coverage patterns across local geography more than was seen in 2001 [12]. Therefore, in this paper we outline the development of the strategy for applying small area techniques to produce LA population estimates for the 2011 census in the light of the updated design of the CCS [10] and the overall estimation strategy for the estimation area level. The discussion in this paper has been made necessarily concise to focus on the small area estimation strategy to provide local authority estimates. Interested readers can refer to the partner paper [11] which provides the background, context and details of the coverage assessment process of the 2011 census.

The outline of the paper is as follows. Section 2 gives some background on the UK census coverage assessment focusing on the small area estimation, linking it to the main estimation strategy for EAs. Section 3 describes the various small area models considered. A simulation study was used to determine how different small area methods performed under a number of scenarios, and Section 4 gives the results of the simulation study. Two estimators are found to be the most suitable and these are further investigated in the last part of this section. The paper concludes with an evaluation of the implementation of the small area strategy in 2011 showing that in general the single estimator that was actually applied performed well, but in a few

instances there is some evidence that would have supported consideration of the proposed alternative approach.

2. Background

The output from the census coverage adjustment process is a complete database with individual and household level records for the entire population, taking full account of any estimated under-coverage. The process begins with the census, which attempts to enumerate the whole population. This is followed by the CCS which undertakes an intensive re-enumeration of a sample of the population. The CCS is a nationally representative sample of over 300,000 households (grouped into postcodes, which are small geographical units made up of 15 to 20 households) and the design is described in [10]. The CCS responding households are matched to the census responses and, for the sampled postcodes, estimates of the missed households and persons are calculated through the application of dual-system estimation [13]. The dual-system estimates are used as inputs to a ratio estimation using census counts as an auxiliary variable to produce estimates of the population for estimation areas. Where an estimation area consists of more than one local authority the estimation area totals then need to be allocated to the constituent local authorities through small area techniques. There are additional stages in the census coverage process, such as quality assurance using administrative datasets and demographic analysis, that often involve inspecting the implied sex ratios of the population as well as birth and death rates. The resulting LA level estimates are used as control totals for the imputation system that produces the fully adjusted database, as outlined in [14]. This paper focuses on the small area estimation part of the coverage process and complements [11] which describes the framework for estimation at the EA level.

The small area approach outlined here builds on the approach used in 2001 but accommodating the adjustments to the CCS design for 2011 outlined in [10]. The CCS design in 2001 created estimation areas by grouping contiguous LAs together with the aim of having a population of around 0.5 million. This was done at the design stage and then there was a further stratification by a Hard-to-Count (HtC) index before allocating the sample [3]. LAs were not explicitly accounted for in the design, and there was no historical data to provide evidence of variation in census coverage to drive the formation of the estimation areas. Therefore, it was important that the small area technique used could directly reflect LA specific variation in coverage remaining after controlling for age-sex and HtC at the estimation stage.

The small area level estimates are contingent on the results of the dual-system estimation, which in turn are reliant on the accuracy of the matching of the census and the CCS. This matching process produces a contingency table with the number of individuals that were in both the census and CCS (n_{11}), in the census but not in the CCS (n_{10}) and those not in the census but in the CCS (n_{01}). By definition, the individuals that are counted neither by the census nor CCS (n_{00}) are unknown, and are referred to as the undercount. In order to estimate the total population it is required to adjust for this undercount by finding an estimate of those missed by both the census and CCS. This is achieved through the assumption that there is independence between the census and CCS. Thus the estimate of those missed by both the census and CCS can be found by the expression

$$\hat{n}_{00} = \frac{n_{01}n_{10}}{n_{11}}.$$

Dual-system estimation also relies on the assumption that individuals have the same chance of being counted by either the census or CCS. The homogeneity assumption does not in fact hold

across the entire population, unless the population is subdivided into groups of similar individuals through post-stratification [13]. In the UK, this is achieved firstly by dividing the country broadly along regional lines into estimation areas. If the local authority is particularly large – for example Manchester – the local authority comprises an estimation area of its own. On the other hand, London has several estimation areas based on grouping contiguous local authorities within the metropolitan area.

The population is further stratified by age and sex, and a ‘hard-to-count’ index. The 2001 HtC index (see [3]) was constructed from household characteristics known to be associated with under-coverage, such as high levels of multi-occupancy and private rented accommodation, based on information from previous censuses and social surveys. It had three strata – easy, medium and hard – and it was assumed that post-stratification using age, sex and hard-to-count index gave reasonable assurance that within each post-stratum there was homogeneity of being counted in the census or CCS. (For the 2011 census the hard-to-count index described by [12] was extended to five strata.) Then for each of the post-strata, those missed in both the census and CCS (n_{00}) can be reasonably estimated with the dual-system estimator (DSE). The dual-system estimator is applied at low levels of geography consisting of three to five postcodes, which provide sufficient data to yield stable estimates as well as forming the primary sampling unit for the design of the CCS [10].

It is possible to produce direct estimates of the local authority totals based on information from the CCS. However, these have unacceptably large standard errors due to small sample sizes, particularly after stratifying by the CCS design variables (such as age and sex). Sample sizes for the local authorities are small partly to keep the survey manageable, and also because the overall sample size was determined to provide specific accuracy at the estimation area level.

Research was carried out to ascertain if it were possible to increase the sample size in order to facilitate direct estimation of the local authority totals from the CCS. However, this was deemed not feasible [10]. The CCS, in addition to being nationally representative, is already a large survey. It is eight times the size of the quarterly Labour Force Survey, which has a responding sample of approximately 40,000 households per quarter [15].

Indirect estimates of the small area population can be produced which increase the effective sample sizes of the local authorities using information from related areas and thereby reducing standard errors. The drawback of these indirect techniques, however, is that they rely on strong assumptions about the relationship between the small areas themselves, in addition to the relationship between the small area and the larger area. Thus, while the estimators may have low variances, they tend to be biased. Therefore, the small area strategy has to strike a balance between the potential bias of an indirect estimator and the imprecision of the direct estimator.

In 2001 a number of different approaches were considered on the basis of available literature and the suitability of the underpinning model assumptions. The small area models were then assessed to find the model that was capable of delivering accurate estimates of the population under various coverage scenarios. In the final model selected, information from all the local authorities within an estimation area was used to model the undercount, but the model coefficients (i.e. the slopes of the regression lines) were allowed to vary by local authority. As a consequence, the heterogeneity of the slopes accounted for the differences in coverage between local authorities and within the specific estimation area [16]. In the next section we discuss approaches investigated when developing the small area strategy for the 2011 UK census, which builds upon the research and lessons learnt from 2001.

3. Small area estimation for local authorities in the 2011 Census

The main objective of the small area estimation strategy is to produce reliable population estimates, with corresponding precision measures, by HtC strata and age-sex groups within each local authority. The age-sex categories used were similar to those used in 2001. There were 35 age-sex groups given by males and females under 1 year old, males from 1 to 4 years old, females 1 to 4 years old, then 5 year age groups for males and for females up to 79 years old, males over 80 years old and females over 80 years old. The small area estimation procedure implemented for the 2011 census apportions the estimation area estimates to the local authorities by assuming a relationship between the undercount pattern at the local authority (small area) level and the broader area (i.e. the estimation area). The starting point is an LA by HtC age-sex specific model and we then explore how to estimate that model by borrowing strength in various dimensions.

To specify a model we start by defining some notation using the same structure as [11]. We assume that modelling takes place within an estimation area, and drop any subscript to distinguish EAs (although we use a subscript e to show statistics calculated over the whole EA). Let Y_{oa} be the true count for age-sex group a from sampled postcodes in output area o , assumed to be within HtC-within-LA stratum h . In reality, this is the dual-system estimate (see [11]) at the cluster level combining across sampled postcodes within output area. Also, let X_{oa} be the corresponding unadjusted census count. A simple model that links the true counts to the census counts as an auxiliary is the ratio model

$$Y_{oa} = R_{ha} X_{oa} + \varepsilon_{ha} \sqrt{X_{oa}}$$

$$\begin{aligned} \text{Var}(Y_{oa} | X_{oa}) &= \sigma_{ha}^2 X_{oa} \quad \text{with } \varepsilon_{ha} \sim N(0, \sigma_{ha}^2) \\ \text{Cov}(Y_{oa}, Y_{o^*a} | X_{oa}, X_{o^*a}) &= 0 \quad \text{for all } o \neq o^* \end{aligned} \quad (1)$$

It is essentially a set of independent ratio models for each age-sex group by HtC-within-LA strata, that is with ratios R_{ha} at the level of the individual LA.

An optimal estimator for (1) follows from [17] and uses the weighted least squares estimator

for R_{ha} given by $\frac{\sum_{o \in s_h} Y_{oa}}{\sum_{o \in s_h} X_{oa}}$, where Y_{oa} , the sum across the sampled postcodes in OA o , is then

replaced by the cluster level dual-system estimator and s_h represents the OAs sampled from HtC-within-LA stratum h . An estimator of the total is then given by $\hat{T}_{ha} = \hat{R}_{ha} X_{ha}$. This is just applying the ratio adjustment to the total unadjusted census count; or more correctly it sums the estimated true counts, observed for the sample data, and then predicts using the estimated ratio applied to the unadjusted census counts for the non-sampled postcodes. This is the model and estimator that is used for an EA containing a single LA with the Y 's replaced with cluster level dual-system estimates to estimate the individual ratios. We now explore ways to 'borrow strength' to estimate for LAs when the sample size is too small to support directly estimating model (1).

Various regression type models that collapsed (1) across different dimensions were considered in a simulation study with the objective of finding an estimator that balanced the trade-off between variance and bias, yielding estimates with good precision and as little bias as possible. As the CCS was stratified by the HtC index, and this was expected to be a good proxy for variation in census coverage, the small area models produce HtC-specific estimates of the local

authority population totals. The general objective is, therefore, to produce model-based estimators for the population total by HtC-within-LA stratum and age-sex group, \hat{T}_{ha} . Here we focus on three alternatives; one direct estimator and two indirect estimators. In the 2001 Census, and again in 2011, the final model-based estimates \hat{T}_{ha} were scaled to the estimation area age-sex population total. This calibration ensured that estimates produced by the small area modelling would be consistent with the sub-national and national population estimates. Variance estimation for the LA estimates within an EA was undertaken using a bootstrap approach developed by [18] in application to population total estimation with a finite sampling population correction (see Chapter 5 of [19]) to ensure that the lower level LA estimates aligned to the (higher-level) EA estimates.

3.1 *The direct estimator*

In this context, the small area *direct estimator* of the local authority total population is one that still relies only on data from the LA, but looks to borrow strength by collapsing (1) within the LA. To do this we fit the model in broader age-sex groups, exploiting the similarity in the age and sex categories. Thus the 35 groups are collapsed into 16 groups indexed by c (therefore with $a \in c$) for estimating model parameters, although the input data still reflect the full 35 groups. These collapsed categories were 0-4 year olds, 5-14 year olds, 15-19 year old males, 15-19 year old females, 20-24 year old males, 20-24 year old females, 25-29 year old males, 25-29 year old females, 30-39 year old males, 30-39 year old females, 40-49 year olds, 50-59 year olds, 60-69 year olds, 70-79 year olds, over 80 year old males and over 80 year old females. Therefore, the adjustment ratios are smoothed across the collapsed age-sex groups requiring fewer ratios to be estimated. This leads to a model for Y_{oa} given by

$$Y_{oa} = R_{hc} X_{oa} + \varepsilon_{hc} \sqrt{X_{oa}} \quad (2)$$

with a variance structure that is specific to the collapsed groupings with $\varepsilon_{hc} \sim N(0, \sigma_{hc}^2)$. The population estimate for age-sex group a , HtC stratum h , and local authority l in a given estimation area is calculated as

$$\hat{T}_{ha}^{dir} = \frac{\sum_{o \in s_h} \sum_{a \in c} Y_{oa}}{\sum_{o \in s_h} \sum_{a \in c} X_{oa}} X_{ha} = \hat{R}_{hc} X_{ha} \quad (3)$$

where s_h are the sample areas from HtC-within-LA stratum h and Y_{oa} is replaced by the cluster level dual-system estimator. The ratio \hat{R}_{hc} is an adjustment factor applied to each age-sex group and HtC stratum within a local authority, with the collapsed category levels satisfying $a \in c$. Therefore, distinct local authorities within the estimation area will have different adjustment factors but there will be less variation amongst the direct estimates by age-sex than at the EA level. However, although the estimates in (3) of the coverage ratio will not vary by age-sex group a within collapsed grouping c , the individual LA estimates are calibrated to the overall EA estimate which will then impose the EA variation in coverage ratios by a within c .

3.2 The synthetic estimator

In this context, the *synthetic estimator* uses data from all the local authorities within a specified estimation area when estimating the coverage of a specific LA. The underlying assumption is that there is a common undercount pattern (observed in the whole estimation area) for all local authorities after controlling for HtC and age-sex. In this way it simplifies (1) by borrowing strength across the LAs within an EA using the level of undercount in each age-sex category by HtC stratum in the estimation area to adjust the local authority census populations. This leads to a model for Y_{oa} given by

$$Y_{oa} = R_{eha} X_{oa} + \varepsilon_{eha} \sqrt{X_{oa}} \quad (4)$$

with a variance structure that is specific to the collapsed groupings with $\varepsilon_{eha} \sim N(0, \sigma_{eha}^2)$. The population estimate for age-sex group a in stratum h in a given estimation area is calculated as

$$\hat{T}_{ha}^{synth} = \frac{\sum_{HtC(h')=HtC(h)} \sum_{o \in s_{h'}} Y_{oa}}{\sum_{HtC(h')=HtC(h)} \sum_{o \in s_{h'}} X_{oa}} X_{ha} = \hat{R}_{eha} X_{ha} \quad (5)$$

where the first sum is over strata with the same HtC level as the target estimator (but varying LAs) and Y_{oa} is replaced by the cluster level dual-system estimator. Comparing the model (4) and estimator (5) with the *direct estimator* given by (2) and (3), we see that the *direct estimator* keeps the full geography by collapsing \hat{R}_{ha} to \hat{R}_{hc} while the synthetic estimator keeps the full age-sex profile by collapsing \hat{R}_{ha} to \hat{R}_{eha} .

3.3 The local fixed effects model

The *local fixed effects model* is another indirect estimator and was the approach implemented in 2001. It is similar to the *synthetic estimator* in that a simple ratio model is fitted that relates the dual-system estimates to the unadjusted census counts using data from the whole EA. The differences are that the regression coefficients vary according to the local authorities, and the age-sex coefficients are for the collapsed groups as in the *direct estimator*. Again the model is fitted to each HtC stratum within each EA using age-sex group by postcode level data and is given by

$$Y_{oa} = (R_{ehc} + \gamma_h) X_{oa} + \varepsilon_{eh} \sqrt{X_{oa}}$$

$$Var(Y_{oa} | X_{oa}) = \sigma_{eh}^2 X_{oa} \quad \text{with} \quad \varepsilon_{eh} \sim N(0, \sigma_{eh}^2)$$

$$Cov(Y_{oa}, Y_{o'a} | X_{oa}, X_{o'a}) = 0 \quad \text{for all } o \neq o' \quad (6)$$

with the collapsed category levels satisfying $a \in c$ and the HtC-within-LA specific effects γ_h in each estimation area assumed to sum to zero within each HtC stratum $\sum_{HtC(h')=HtC(h)} \gamma_{h'} = 0$.

The model is actually fitted using weighted least squares applied to data based on the cluster of sampled postcodes within an OA to get estimates \hat{R}_{ehc} and $\hat{\gamma}_h$ of the model parameters. Given these estimated parameters, it follows that a model based estimator for the population total by local authority, HtC stratum and age-sex group can be defined as $\hat{T}_{ha}^{LFE} = (\hat{R}_{ehc} + \hat{\gamma}_h) X_{ha}$. We can see that this estimator has age-sex effects that are common to all LAs within the EA but also allows for LA specific coverage adjustments that apply to all age-sex groups by collapsing \hat{R}_{ha} to $(\hat{R}_{ehc} + \hat{\gamma}_h)$. This then allows for local factors that might be expected to have a universal impact on census coverage for the whole LA, while recognising that the main coverage patterns will be driven by general age-sex and HtC effects for the whole EA. Such an approach was important in 2001 where there was little historical information on coverage to use when combining LAs, and the census fieldwork was still locally organised and managed, with individual enumerators directly responsible for small areas making localised census failures more possible.

4. Evaluation of the small area methods

Section 3 outlined three estimators that can be applied at local authority level to produce population estimates. Their relative performance depends on the strength of localised census enumeration effects that cannot be controlled for using a combination of age-sex and hard-to-count classifiers within an estimation area. To get an idea of the trade-offs, a simulation study

was used to evaluate the three competing estimators. A series of censuses and CCSs were simulated using predicted coverage probabilities obtained through modelling of the under coverage in the 2001 census and CCS data. Simulations were produced for a number of estimation areas with a variety of coverage patterns. For each estimation area in the simulation, 400 censuses and 400 CCSs were used. The first step in the estimation procedure was to produce estimates of the population totals for the larger domains, here the estimation areas. For each simulated census and CCS combination, dual-system estimation and ratio estimation were used to produce estimates of the estimation area totals for the detailed age-sex groups by hard-to-count stratum. After this was completed, the local authority estimates by age-sex group and HtC stratum were obtained for each of the 400 simulations within an estimation area using the three competing estimators.

As outlined in section 2, it is known that, although more precise (that is, with lower variance), the indirect estimators have a tendency to be biased in comparison with the direct estimators. As such the aim of the evaluation process was to weigh the reduction in variance against potentially larger biases. Therefore, based on the 400 simulation results the relative bias and the relative root mean squared error were calculated as suitable measures of performance that could be used to investigate the bias and variance. The mean squared error is a function of both the variance and bias, and is consequently a good measure of the overall accuracy of the different estimators (see page 253 of [20]). Therefore, for our application, the relative root mean squared error (RRMSE) and the relative bias (RB) for each domain (HtC by age-sex) in a given local authority are respectively calculated as

$$\text{RRMSE}(\hat{T}_{ha}) = \frac{1}{T_{ha}} \sqrt{\frac{\sum_{j=1}^{400} (\hat{T}_{ha}^{(j)} - T_{ha})^2}{400}} \quad \text{and} \quad \text{RB}(\hat{T}_{ha}) = \frac{1}{T_{ha}} \frac{\sum_{j=1}^{400} (\hat{T}_{ha}^{(j)} - T_{ha})}{400} \quad (7)$$

where:

T_{ha} is the true population count for the age-sex group a in HtC-within-LA stratum h ; and

$\hat{T}_{ha}^{(j)}$ is the corresponding model based population estimate obtained from the j^{th} simulation,

with $j = 1, \dots, 400$.

4.1 Results of the simulations

Simulated census and CCS data were obtained for some estimation areas which were selected because they had different levels of coverage in the 2001 census. As the investigation sought to determine how each of the different small area models fared under a range of coverage scenarios, estimation areas were chosen to exhibit diverse census coverage characteristics. This paper presents results from four estimation areas, to show the methodological development of the small area strategy for the 2011 UK census. The chosen areas are KK and KO from the Midlands, LB from Inner London, and LJ from Outer London, which cover a range of observed census coverage patterns for the 2001 Census. These estimation areas consist of two or three constituent local authorities and showcase the issues that had to be considered when choosing a suitable small area methodology to produce reliable estimates of the local authority totals.

Table 1 gives the 2001 Census coverage rates by local authority and estimation area. It shows that higher coverage is achieved in KK and KO but lower coverage in LB and LJ. In addition, there are some differences in coverage by local authority within estimation areas reflecting the fact that 2001 estimation areas were based on geography and population size with little available evidence relating to localised variation in census coverage. However, this variation may also be related to differing age-sex and HtC structures within the local authorities of each estimation area.

----- TABLE 1 ABOUT HERE -----

For each of the estimation areas, the RRMSEs and RBs were calculated for the 3 competing small area estimation techniques (namely *direct estimator* \hat{T}_{ha}^{dir} , *synthetic estimator* \hat{T}_{ha}^{synth} and *local fixed effects* \hat{T}_{ha}^{LFE}). We are interested in exploring the behaviour of the different small area estimators and looking to determine which estimator produces the most robust estimates of the local authority population totals. Table 1 shows the RRMSE and RB for the local authority population totals in each estimation area. The results in the table for the three small area model-based estimates are found by summing across the age-sex groups and the hard-to-count strata. This gives an indication of the variability of the different local authority population totals produced by the different small area strategies. From Table 1, when the target parameter is the local authority population total, the *synthetic estimator* produces the lowest RRMSE in five of the 11 local authorities; and is very similar to the lowest in a further three. The estimates where it is lowest all occur in the two London EAs where the observed coverage patterns for the LAs in the 2001 Census are relatively similar within each EA. *Local fixed effects* is also the lowest in five local authorities and these occur in the other two EAs which tend to have higher coverage but greater variation across the LAs within each EA.

In terms of RRMSE, Table 1 suggests the choice is between a *synthetic estimator* that is likely to have smaller variance but more potential for bias and *local fixed effects* with potentially higher variance but less bias. This is confirmed by the bias results in Table 1, where the *synthetic estimator* typically has larger absolute bias with either the *local fixed effects* or *direct estimator* having the smaller absolute biases. However, it is worth noting that in the design for

the 2011 CCS [10], the direct use of local authority in the design results in KK1, KO1 and all of LB being treated as estimation areas with a single local authority at estimation [11] due to their more extreme coverage patterns relative to neighbouring local authorities. Therefore, taking the results in Table 1 with the changing structure of the CCS, the *synthetic estimator* would be expected to perform better in terms of RRMSE but there may be a small bias if the estimation areas combine local authorities that then experience localised coverage effects in 2011.

While Table 1 presents results for the total population, it is important to consider the age-sex by HtC estimates as this is the level at which the estimators operate. Boxplots of the distributions of RRMSEs and RBs for the 105 (i.e. 35 x 3) age-sex by HtC model-based population estimates for each local authority are shown in Figure 1. Small area techniques that perform well should produce an RRMSE distribution with lower median and a smaller spread. In the case of bias, a good technique should produce an RB distribution that is centred around zero with small spread. For both RB and RRMSE distributions outliers are indicative of possible model failure, therefore any outlying observations are highlighted in the boxplots.

----- FIGURE 1 ABOUT HERE -----

The boxplots for KO, KK, and LJ are less skewed and exhibit smaller variability in comparison to LB. These boxplots provide evidence that in general the *synthetic estimator* has lower RRMSEs and performs best in comparison to the *local fixed model* and the *direct estimator*.. Furthermore, the distributions have smaller spread within local authorities for each of the estimation areas. However, when examining the RBs, the *local fixed effects model* produces better behaved distributions, which are mostly centred around zero and are therefore

approximately unbiased. The reasoning behind the *local fixed effects estimator* is to capture any difference in coverage due to local authority effects. Although no improvement in the RRMSE was found, the model containing local authority effects may protect the estimation procedure against failure when local authority differentials are observed. This motivated the use of the *local fixed effects* model in estimation areas where there was evidence of coverage variation between LAs within the EA.

The analysis shows that the *synthetic estimator* seems to be doing the best overall. An explanation of why the *synthetic estimator* does better than the *local fixed effects estimator* might simply be that the simpler model behind the estimator is sufficient to capture the likely coverage patterns. The *local fixed effects model* includes a fixed effect for each local authority, however if there are no (or only small) local authority differentials in undercoverage, then additional modelling error is being introduced, with little benefit. Furthermore, the results do make some sense in the context of the coverage rates in Table 1. Most of the local authorities have similar coverage rates to the overall estimation area coverage. Even in estimation areas with relatively poor coverage, such as the inner London boroughs of LB, all the local authorities exhibit similar coverage patterns. The *local fixed effects model* becomes useful if the different local authorities in the estimation area have varying coverage rates. Nonetheless, the *local fixed effects model* does have some definite benefits with regards to its intuitive appeal; it can also offer more protection against model failure than the *synthetic estimator*. Notice that the *direct estimator*, which is typically less efficient than the *synthetic* and *local fixed* models since it does not borrow strength outside the estimation domain, still performs well; and can perform as well as the other two, as is evidenced in KO.

The results indicate that both the *synthetic estimator* and *local fixed effects estimator* are reasonable options to produce local authority population estimates. The first performs better in terms of RRMSE whereas the latter produces estimates with smaller biases. The *synthetic estimator*, however, seems more stable as it shows less variability in performance across local authorities (as shown earlier in Figure 1). The use of a *local fixed effects model* could represent a safeguard for local authority undercoverage differentials. However, as demonstrated in some of the results, the *local fixed effects model* may add unnecessary noise into the estimates if there are no local authority effects to be observed. The compromise solution for the 2011 census was to implement a small area estimation procedure that accommodated both options. That is, the *synthetic estimator* was the default option for each estimation area, thereby assuming the local authority effects were not important. Then, if the quality assurance procedure found evidence of a localised failure in coverage, fit a *local fixed effects model* and test the significance of the areal effects.

4.2 *Assessing the Performance in 2011*

Based on the simulation results and the change in structure to the CCS, the standard approach implemented in the 2011 Census utilised synthetic estimation for local authorities within an estimation area. The use of *local fixed effects* would be explored only if quality assurance identified evidence of localised coverage effects that needed to be accounted for. No such situations occurred, so all local authority outputs were either for a single LA making up an EA by itself, or *synthetic estimates* within the EA. However, we can now explore the models in a little more detail to assess the robustness of this approach using the actual 2011 data.

For the 70 EAs that contain more than a single LA, we compare the *synthetic* model with the full set of age-sex categories to a *synthetic* model with the collapsed age-sex categories and then the *local fixed effects* model (with the same collapsed age-sex categories). Having the synthetic approach for both the full and collapsed age-sex groups allows us to assess the cost of reducing the number of groups prior to assessing the potential benefit of adding the local fixed effects. The approach used to assess the strength of the local authority effects in a given estimation area was to compare the different models using two goodness-of-fit measures: the Schwarz Bayesian Information Criterion (BIC) and the adjusted R² value. In both cases the measures are based on the variation explained by the model but with penalties for the number of parameters, making them suitable to compare non-nested models. In the case of the BIC smaller values represent better fit, while for the adjusted R² larger values imply better fit.

The BIC for the *local fixed effects* model was found to be smaller than that for either of the *synthetic* models in just six of the 70 estimation areas considered. This indicates that for the vast majority of estimation areas there was no evidence of strong local authority effects. The six estimation areas where there was some indication of stronger local authority effects were examined in greater detail. The model goodness of fit statistics for these estimation areas are given in Table 2.

----- TABLE 2 ABOUT HERE -----

In all but one of these six estimation areas in Table 2, just one of the hard to count strata had the smallest BIC for the *local fixed effects* model. The exception is the EA coded SW04 from the South-West, where both hard to count strata 2 and 3 have smaller BIC values for the *local fixed effects* models. In Table 2 it can also be seen that the difference in BIC values between

the *local fixed effects* model and the collapsed age-sex group *synthetic* model is small for these six areas, regardless of which model has the actual lowest value. This implies that the addition of fixed effects over broader age-sex groups has little advantage. The BIC values for both the collapsed age-sex group *local fixed effects* model and the collapsed age-sex group *synthetic* model are smaller than the corresponding values for the full age-sex group *synthetic* model. This implies there is some potential efficiency gain from collapsing age-sex groups, but the requirement to produce estimates for the five-year age-sex groups means we would not want to collapse unless it was needed to allow the inclusion of the *local fixed effects*. The adjusted R^2 values are generally largest for the *local fixed effects* model, but there is really very little improvement in adjusted R^2 from including the local authority effects or collapsing the age-sex groups.

----- FIGURE 2 ABOUT HERE -----

In Figure 2 the BIC values for all areas obtained from fitting both *synthetic* models are plotted against the BIC value from the corresponding *local fixed effects* model, together with the fitted lines. Also plotted is the $y=x$ line to demonstrate how close the values from the *synthetic* models are to the *local fixed effects* model. In this figure the signs of BIC values have been changed so that the larger the BIC value the better. In Figure 2, the fitted line of the *local fixed effects* against *synthetic* with collapsed age-sex groups is very close to the $y=x$ line showing, in general, that adding the local authority effects does not improve the fit of the model compared to a *synthetic* model with the same age-sex groups. However, the fitted line for the BIC values from the comparison of local fixed effects to the *synthetic* model with the full age-sex categories is slightly below the $y=x$ line, which indicates that having a greater number of age-sex groups in the model generally results in an improved fit over the inclusion of the local

authority effects and a reduced age-sex categorisation. From these overall results in Figure 2, combined with the small number of EAs highlighted in Table 2, we can see that the small area strategy for 2011 performed well in that the *synthetic* approach did well in the vast majority of cases. Even when the *local fixed effects* model gave an improved fit, the gain was marginal; and this shows why these impacts were not detected in the quality assurance process.

5. Conclusions

Small area estimation techniques are useful in overcoming the problem of small sample sizes since direct estimates using data from the CCS would have correspondingly large standard errors and be imprecise. However, although they are precise, these (indirect) model based estimators may be more biased than the direct estimators. Therefore, the aim of the evaluation of different estimators was to balance the trade-off between variance and bias in order to find the estimator that produced estimates with good precision and as little bias as possible. The small area models work by incorporating auxiliary information by assuming relationships between the undercount pattern in the local authority and broader areas such as the estimation area. The underlying idea was to exploit the similarities in the undercount patterns so as to borrow strength over the areas through the use of regression models relating the dual-system estimates to the census counts.

Indirect estimators such as the *synthetic estimator* and the *local fixed effects estimator* can realistically be used to improve the precision of the local authority estimates. However, this improvement in precision is wholly dependent on being able to exploit the similarities between local authorities either within the estimation area or within the region. In the simulation

exercise, it was demonstrated that the choice of the indirect model can be complicated because it is often not easy to know how to exploit these similarities. Generally speaking, if the local authorities within an estimation area are broadly similar, the *synthetic estimator* will often be the most appropriate indirect model; and this is supported by the results of the simulation. However, in 2001 the approach was based on *local fixed effects* for local authorities as the CCS did not directly control local authorities and the formation of estimation areas was based on population size rather than historical coverage patterns. The re-design of the CCS for 2011 outlined in [10] brings in the historical coverage of local authorities. Consequently estimation areas were formed to be more internally homogenous with respect to coverage; and several smaller local authorities became estimation areas by themselves rather than being grouped with neighbours. However, when there is localised failure of the census (and/or the CCS) - for example, a specific local authority behaves differently to the estimation area within which it is found - then the *synthetic estimator* can be less precise than the *local fixed model*.

The main reason for using indirect estimation for the local authority population totals is to improve precision by combining information from the broader estimation area to increase the effective sample size. In this paper we have explored two indirect approaches, the *synthetic estimator* and *local fixed effects estimator*, both applied to an estimation area. In preparation for the 2011 Census, additional research was carried-out to assess more complex indirect estimators based on models using random effects but fitted to larger areas, in our case government office region (GOR). The underlying assumption here was that the undercount pattern in the GOR was similar to the undercount pattern in the local authority. Obviously, this is not necessarily true but the inclusion of random effects helps account for local authority differentials. However, results from this additional research, not presented here, found that these more complex (random effects) models did not do any better than the *synthetic model* or

the *local fixed model*. In addition, we considered composite models which took a weighted combination of the synthetic estimator and the local fixed model. These composite estimators tended to increase the variability, and were therefore found to be relatively inefficient.

Thus the recommendation made from these simulation results was that the most appropriate small area strategy involved accommodating both *synthetic* estimation and *local fixed effects* regression. The *synthetic estimator* was the default technique, and could cope with some local authority differentials provided they could be explained by hard-to-count and age-sex. However, in the case that there were unanticipated problems in the census and the CCS leading to greater differences in the observed local authority coverage levels, this would be detected by the quality assurance process and the *local fixed model* would be better placed to produce more robust population estimates.

During the estimation for the 2011 Census, the quality assurance did not trigger the use of *local fixed effects*, as the default *synthetic* estimates were accepted. However, here we present the results from a modelling exercise that compared the two approaches for all 70 EAs. The results of this confirm that the synthetic model was generally a better fit than the local fixed effects model. However, it also highlighted how little difference there was between the approaches which all had very high values for the adjusted R^2 showing how well the models explained the variation in coverage using the census counts. This demonstrates that an initial population count that manages to count everyone well, with very little undercount, will ensure a more robust small area adjustment with accurate local authority population estimates. Conversely, any small area technique will struggle to adjust a poorly performing census. Looking ahead for the next censuses in 2021 and beyond, the small area estimation strategy can be enhanced with the use of administrative register data, specifically during the final quality assurance of

the estimates, to ensure more robust and reliable adjusted population counts at a local authority level.

Acknowledgements

The authors thank the members of the various census committees that have commented on this work as it has developed. They would specifically like to acknowledge the contribution and support of Dr Frank Nolan from the Office for National Statistics, who passed away unexpectedly in 2012, in the development of the coverage assessment plans for the 2011 Census.

References

- [1] Diamond I. The Census. In: Dorling D, Simpson L, editors. *Statistics in Society: the arithmetic of politics*. London: Arnold; 1999; 9-18.
- [2] Abbott O. 2011 UK Census Coverage Assessment and Adjustment Methodology. *Population Trends* 2009; 137: 25-32.
- [3] Brown JJ, Diamond ID, Chambers RL, Buckner LJ, Teague AD. A methodological strategy for a one-number census in UK. *Journal of the Royal Statistical Society: Series A* 1999; 162: 247-267.
- [4] Rao JNK, Molina I. *Small area estimation*, 2nd Edition. New York: Wiley; 2015.
- [5] Martin D. Editorial: census present and future. *Journal of the Royal Statistical Society: Series A* 2007; 170: 263-266.
- [6] Pfeffermann D. Small area estimation – new developments and directions. *International Statistical Review* 2002; 70: 125-143.
- [7] Ghosh M, Rao JNK. Small area estimation: an appraisal. *Statistical Science* 1994; 9: 55-93.
- [8] Office for National Statistics. 2001 census: Manchester and Westminster matching studies full report. London: Office for National Statistics 2004 [cited 2017 Oct 19]. Available from <http://www.ons.gov.uk/ons/guide-method/method-quality/specific/population-and-migration/pop-ests/local-authority-population-studies/2001-census---manchester-and-westminster-matching-studies-full-report.pdf>
- [9] Local Government Association. *The 2001 One Number Census and its quality assurance: a review*. Research Briefing 6.03. London: Local Government Association; 2003.
- [10] Brown J, Abbott O, Smith PA. Design of the 2001 and 2011 census coverage surveys for England and Wales. *Journal of the Royal Statistical Society Series A* 2011; 174: 881-906.

- [11] Brown J, Sexton C, Abbott O, Smith PA. The framework for estimating coverage in the 2011 Census of England and Wales: combining dual-system estimation with ratio estimation. Submitted to *Statistical Journal of the International Association of Official Statistics*; 2017.
- [12] Abbott O, Compton G. Counting and estimating hard-to-survey populations in the 2011 Census. In: Tourangeau R, Edwards B, Johnson TP, Wolter KM, Bates NA, editors. *Hard-to-Survey Populations*. Cambridge: Cambridge University Press; 2014.
- [13] Sekar CC, Deming WE. On a method of estimating birth and death rates and the extent of registration. *Journal of the American Statistical Association* 1949; 44: 101-115.
- [14] Steele F, Brown J, Chambers R. A controlled donor imputation system for a one-number census. *Journal of the Royal Statistical Society, Series A* 2002; 165: 495-522.
- [15] Office for National Statistics. Quality and methodology information (LFS). Information Paper. Newport: Office for National Statistics; 2015.
- [16] Office for National Statistics. One number census local authority estimation. London: Office for National Statistics; 2000 [cited 2017 Oct 19]. Available from <http://www.ons.gov.uk/ons/guide-method/census/census-2001/design-and-conduct/the-one-number-census/methodology/steering-committee/key-papers/local-authority-estimation.pdf>.
- [17] Royall RM. On finite population sampling under certain linear regression models. *Biometrika* 1970; 57: 377-387.
- [18] Efron B, Tibshirani RJ. *An introduction to the bootstrap*. Boca Raton: Chapman & Hall/CRC; 1993.
- [19] Wolter K. *Introduction to variance estimation*. 2nd edition. New York: Springer; 2007.
- [20] Cox DR, Hinkley DV. *Theoretical statistics*. London: Chapman and Hall; 1974.

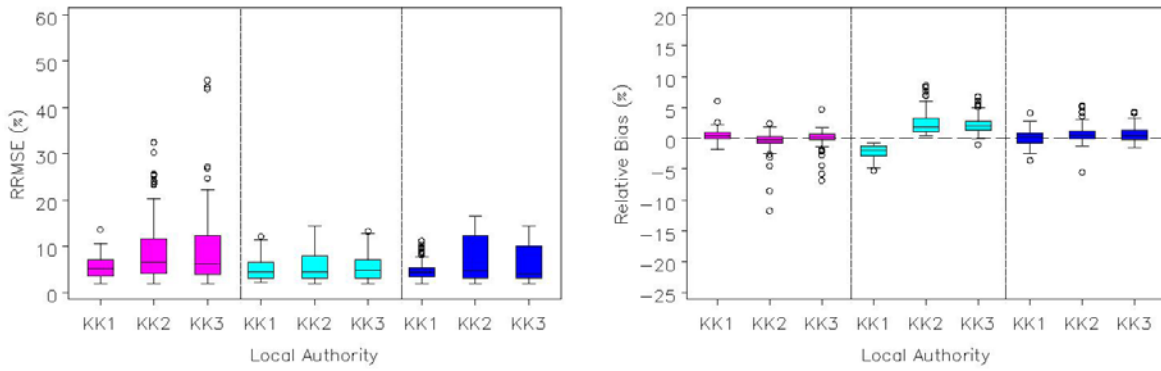
Table 1. Performance (RRMSE and Relative Bias) for local authority total population estimates by small area model

		Small Area Estimation Models/Estimators					
Estimation Area ¹	Local Authority ¹	RRMSE (%)			Relative Bias (%)		
		Direct	Synthetic	Local Fixed	Direct	Synthetic	Local Fixed
	KK1 (91.42)	1.97	1.96	1.78	0.47	-1.37	0.12
KK (95.5)	KK2 (98.00)	2.03	2.48	2.05	-0.12	2.24	0.38
	KK3 (97.17)	1.79	2.48	1.67	0.10	2.23	0.45
	KO1 (92.39)	1.32	1.36	1.30	-0.09	-0.75	-0.15
KO (95.2)	KO2 (98.02)	1.01	1.34	1.00	0.10	1.08	0.19
	LB1 (73.28)	3.81	3.15	3.66	-0.97	-2.14	-0.60
LB (76.5)	LB2 (79.32)	3.62	4.32	3.50	-0.94	3.65	-1.01
	LB3 (76.93)	4.79	3.60	4.69	0.14	-2.72	-0.17
	LJ1 (87.80)	2.40	1.53	2.21	-0.14	-0.74	0.12
LJ (88.4)	LJ2 (88.38)	2.46	1.63	2.35	0.06	0.66	-0.03
	LJ3 (88.93)	2.75	1.94	2.67	-0.18	-0.38	-0.27

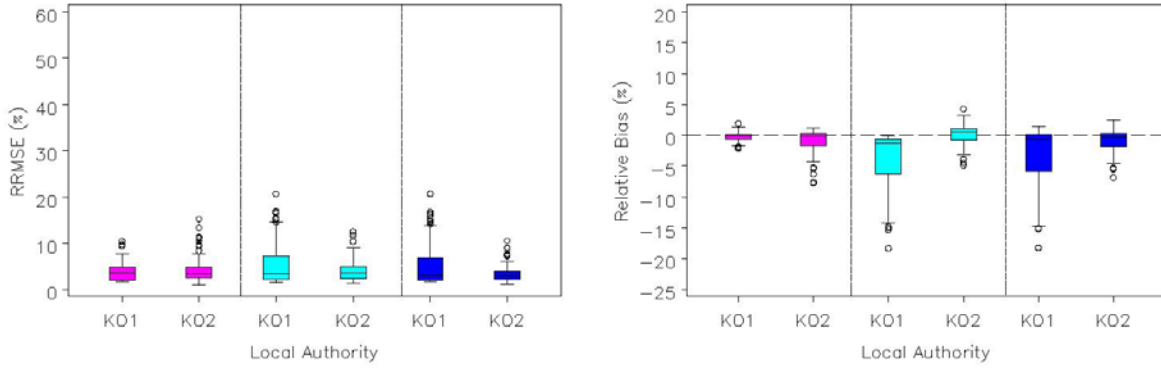
1. Estimated coverage percentage for 2001 Census in brackets.

Table 2. A comparison of model ‘goodness of fit’ for estimation areas and hard to count strata where the BIC (Schwarz Bayesian Information Criterion) goodness of fit measure for the fixed effects model is smaller than that for the synthetic models.

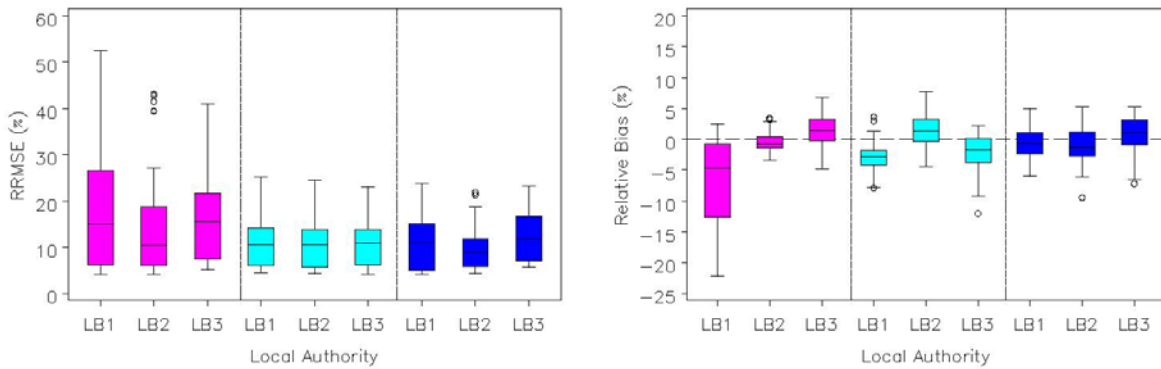
EA code	HtC stratum	Number of LAs	Fixed effects – collapsed age-sex groups		Synthetic model – collapsed age-sex groups		Synthetic model – full age-sex groups	
			BIC	AdjR ²	BIC	AdjR ²	BIC	AdjR ²
EE05	1	6	-996.9	0.9855	-957.1	0.9834	-897.2	0.9833
	2	7	-1712.0	0.9892	-1748.5	0.9893	-1681.8	0.9892
SE03	2	3	-1268.2	0.9858	-1270.4	0.9856	-1220.91	0.9855
	3	2	-402.4	0.9776	-376.2	0.9752	-326.4	0.9745
SW04	1	3	-441.4	0.9850	-450.2	0.9850	-401.5	0.9846
	2	3	-681.6	0.9845	-664.8	0.9833	-608.8	0.9829
	3	2	-38.9	0.9368	-37.8	0.9345	8.9	0.9305
WA02	1	3	-1098.7	0.9775	-1093.0	0.9769	-1024.6	0.9766
	2	3	-466.8	0.9750	-472.0	0.9747	-420.1	0.9746
WM03	2	2	-1436.4	0.9809	-1441.6	0.9808	-1366.8	0.9806
	3	2	-304.5	0.9449	-303.1	0.9441	-241.1	0.9437
YH07	1	2	-719.3	0.9953	-713.6	0.9951	-665.7	0.9950
	2	2	-1417.3	0.9839	-1423.9	0.9839	-1361.4	0.9837



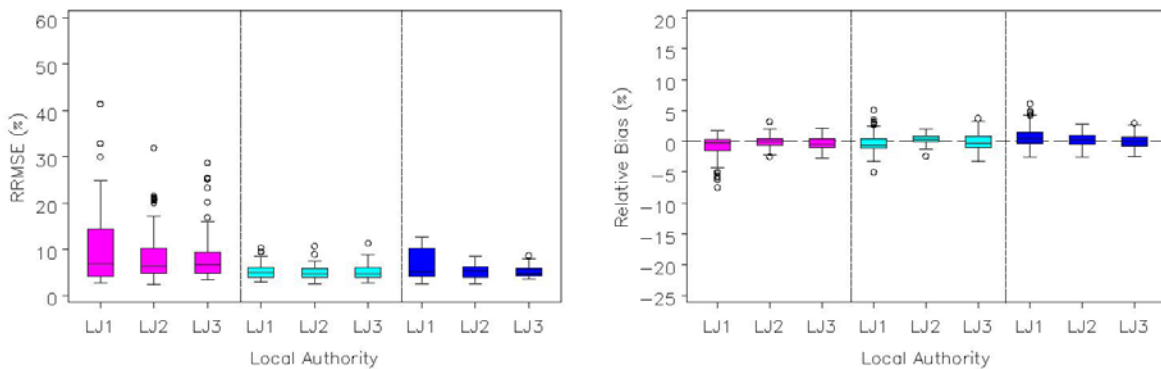
(a) Estimation Area KK



(b) Estimation Area KO

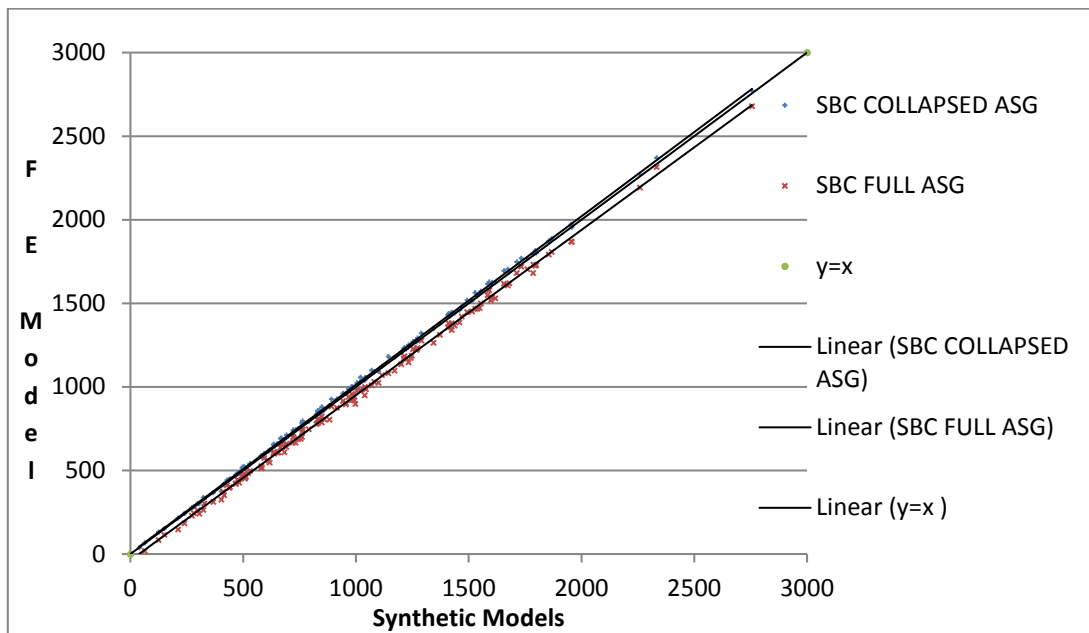


(c) Estimation Area LB



(d) Estimation Area LJ

Figure 1. Boxplots showing the RB and RMSE distribution of the different small area estimators for the selected four estimation areas. For each plot the left panel represents the direct estimator, middle panel represents the synthetic estimator, and the right panel represents the local fixed effects model.



*The BIC values have been multiplied by -1 so that in this figure the larger the BIC value the better

Figure 2. Adjusted Schwarz Bayesian Information Criterion (BIC*) values for fixed effects models against synthetic models