# Small Area Estimation

**1** - Small area estimation problem

**2** - Estimation for domains - **Direct estimators** – estimation for planned domains

**3** – Coefficient of Variation and Minimum level of precision

**4** - Estimation for unplanned domains and/or where the sample size is not enough for the minimum level of precision – **Indirect estimators**

# Recap

- Target parameters in the domain:

- Total of the study variable $\quad t = \sum_{k \in U} y_k$

- Mean of the study variable $\quad \bar{y} = \sum_{k \in U} y_k / N$

- At risk of poverty rate $\quad P_0 = \dfrac{1}{N} \sum_{i=1}^{N} I(y_i < z).$

- Poverty gap $\quad P_1 = \dfrac{1}{N} \sum_{i=1}^{N} \dfrac{G_i}{z}.$

# Inference framework

Further we have, depending on the "reference framework" for inference:

— **Design Based Approach:** Estimator properties are assessed with respect to the sampling design (see previous example). This framework is used for small area estimation, mainly because of its simplicity.

— **Model Assisted Approach:** In practice, the values of Y are typically defined by assuming a model for the distribution of Y given X. That is, practitioners have been willing to use models in order to identify optimal strategies for estimating $T_Y$. However, their assessment of these strategies remain design-based (Särndal, Swensson and Wretman, 1992).

— **Model Based Approach:** design-unbiasedness is no longer a requirement, the alternative property we require of the estimator under this approach is that it be model-unbiased $E\left(\hat{T}_Y - T_Y \mid \mathbf{S}, \mathbf{X}\right) = 0$. given the sample S and aux info X.

# What are we modelling?

We are modelling the relationship between an outcome and the auxiliary variables

note that:

- unplanned domains=geographical domains= areas

- notation:

$Y_{ij}$  outcome= the value of the study variable (*income survey data unit j, individual or household, in   area i*)

$\hat{\theta}_i^{dir}$  outcome= survey direct estimator (*per capita income in area i, total income in area i*)

# What are we modelling?
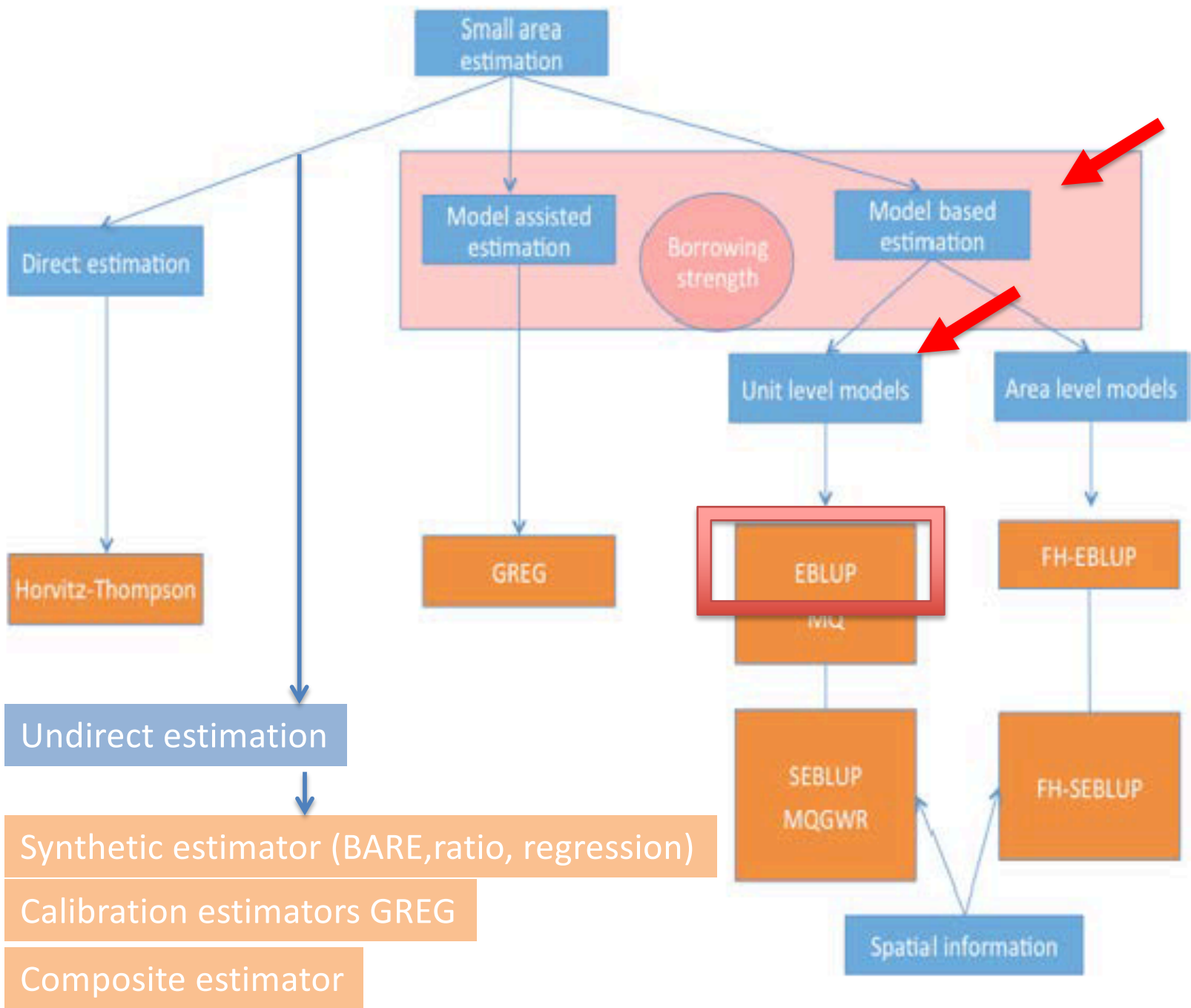
The models are classified into two broad types:

1   Aggregate level (or area-level ) models that relate the small area outcome (means, totals) to area-specific auxiliary variables. Such models are essential if unit level data are not available

2   Unit level models that relate the outcome (unit values of the study variable) to unit-specific auxiliary variables

The use of *explicit models* offers several advantages

# What are we modelling?
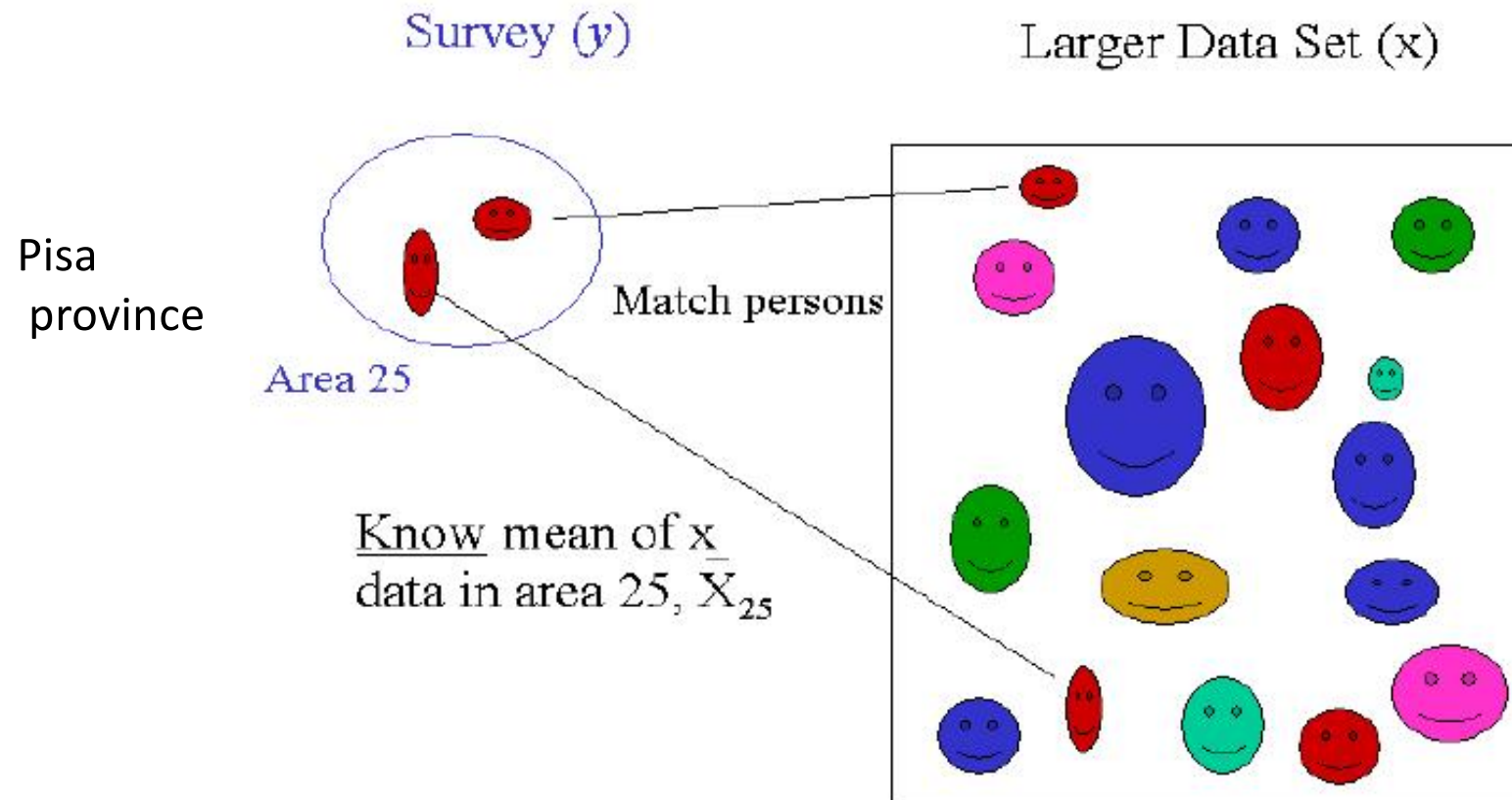
Advantages both for are-level and unit-level models:


1  Model diagnostics can be used to find suitable models that fit the data well

2  Area-specic measures of precision can be associated with each small area estimate, solving the problem of instability seen for synthetic and composite estimators

3  Linear mixed models as well as nonlinear models can be used.
4  Complex data structures, such as spatial dependence and time series structures, can also be handled

5  Methodological developments for random effects models can be utilized to achieve accurate small area inferences

# Unit Level Approach

1 - outcome : $y_{ij}$ , auxiliary variable *X* available at unit level (j) from a larger data set
Ex:   unit=individual, area=province: $y_{ij}$ income of unit, X household size



Survey (y)

Larger Data Set (x)

Pisa province

Area 25

Match persons

Know mean of x data in area 25, $\bar{X}_{25}$

# EBLUP: Unit level approach

- $\boldsymbol{y}$ the vector for the $y$ variable for the population $\Omega$
- $\boldsymbol{y} = [\boldsymbol{y}_s', \boldsymbol{y}_r']'$, where $\boldsymbol{y}_s$ is the vector of the observed units (the sampled ones) and $\boldsymbol{y}_r$ is the vector of the non observed units ($N-n$, $r = 1, \ldots, N-n$)
- $\boldsymbol{X}$ is the covariates matrix and is considered know for all the population units
- Subscript $i$ refers to small areas (e.g. $\boldsymbol{y}_{s_i}$ is the vector of observed variables in area $i$)

$$\boldsymbol{y} = [\boldsymbol{y}_s', \boldsymbol{y}_r']'$$

EUSILC: ...you observe income of the units in the sample

There are other units in the area i ...but they are not included in the sample: you do not observe their income

# EBLUP: Unit level approach

- Model for the *y* variable (known as superpopulation model)

$$y = X\beta + Zu + e$$

- that can be alternatively write as

$$y_{ij} = x_{ij}\beta + u_i + e_{ij}$$

In addition to the assumptions already made, we require that the model holds for both the population and the sample.

$X\beta$

$Zu$

| Base-line part of the model | Differences among areas |

# EBLUP: Unit level approach

Starting with the linear regression model for grouped individuals:

$$y_{dj} = \beta_{0d} + \beta_{1d} \cdot x_{dj} + \varepsilon_{dj}$$

for the groups $d = 1, \ldots, D$ and the individuals $j$, we assume that $\beta_{0d} = \beta_0 + u_{0d}$ and $\beta_{1d} = \beta_1 + u_{1d}$. For the random effects $u_{0d}$ and $u_{1d}$ we further assume

$$E(u_{0d}) = E(u_{1d}) = 0 \quad \text{and}$$
$$V(u_{0d}) = \sigma_{u0}^2, \quad V(u_{1d}) = \sigma_{u1}^2, \quad \text{Cov}(u_{0d}, u_{1d}) = \sigma_{u01}.$$
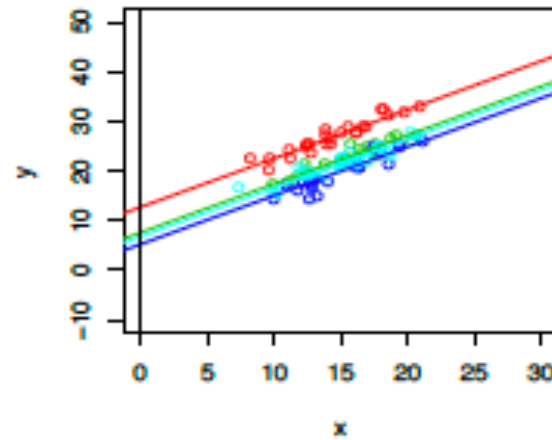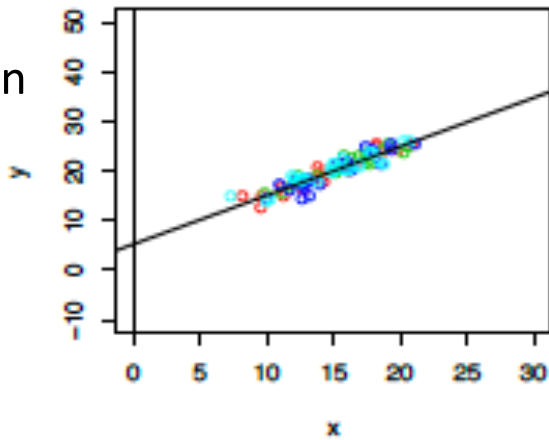
This yields

$$\boxed{X\beta} \qquad \boxed{Zu}$$

$$y_{dj} = \beta_0 + \beta_1 \cdot x_{dj} + (u_{0d} + u_{1d} \cdot x_{dj} + \varepsilon_{0dj})$$
$$V(\varepsilon_{0dj}) = \sigma_{\varepsilon 0}^2 \quad .$$

Moreover, we require independence between $\varepsilon_{0dj}$ and both random effects.
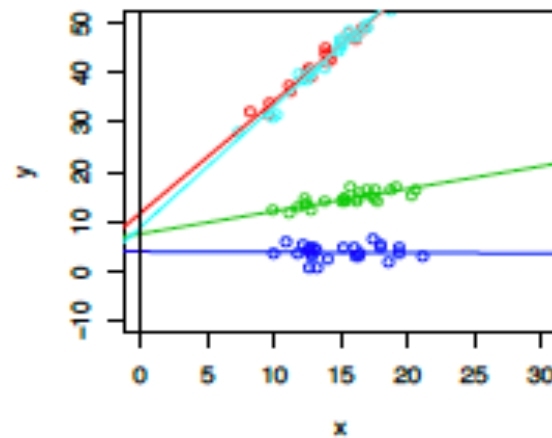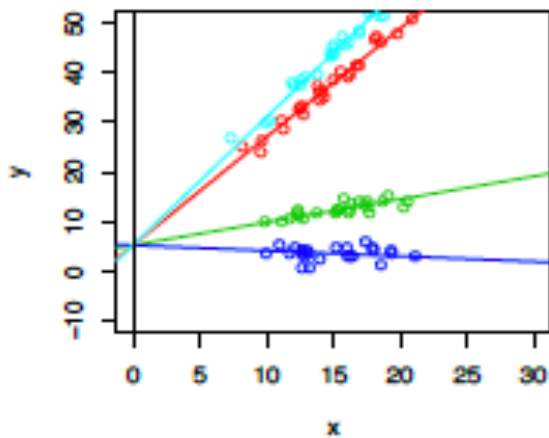
# EBLUP: Unit level approach

Different models (colours are identifying groups of units:

Regression



Different intercepts

Different slopes

Different Intercepts and slopes

# EBLUP: Unit level approach

Random intercept model:

- ▶ Assume that only the intercept is a random component on the second level.
- ▶ This is the by far most common choice in small area estimation.
- ▶ It implies $u_{1d} \equiv 0$ and hence

$$y_{dj} = \mathbf{x}'_{dj}\beta + u_d + \varepsilon_{dj}$$

- ▶ The covariance of any two units follows as:

$$\text{Cov}\,(y_{dj}, y_{d'j'}) = \begin{cases} \sigma_{u0}^2 + \sigma_{\varepsilon 0}^2, & \text{if } d = d' \text{ and } j = j', \\ \sigma_{u0}^2, & \text{if } d = d' \text{ and } j \neq j', \\ 0,, & \text{if } d \neq d' \text{ and } j \neq j' \end{cases}$$

- ▶ Hence, two units from the same group $d$ will be correlated, whereas units from different groups are independent.

# Battese Harter Fuller Model

Battese, G.E., Harter, R.M., and Fuller, W.A. (1988) *An error-components model for prediction of county crops using survey and satellite data*. Journal American Statistical Association 83 28-36.

► This yields

$$y_{dj} = x'_{dj}\beta + u_d + \varepsilon_{dj}, \quad d = 1, \ldots, D, j = 1, \ldots, N_d, \quad (1)$$

$$u_d \overset{iid}{\sim} N(0, \sigma_u^2)$$

$$\varepsilon_{dj} \overset{iid}{\sim} N(0, \sigma_\varepsilon^2)$$

with $x_{dj} = (x_{dj1}, \ldots, x_{djp})'$ as a $p \times 1$ column vector of the covariates for the $j$-th unit within the $d$-th area.

► We estimate $\beta$ by GLS, i.e

$$\widehat{\beta} = \left[ \sum_{d=1}^{D} \sum_{j=1}^{n_d} x_{dj}(x_{dj} - \gamma_d \bar{x}_d)' \right]^{-1} \left[ \sum_{d=1}^{D} \sum_{j=1}^{n_d} (x_{dj} - \gamma_d \bar{x}_d)' y_{dj} \right]$$

# Battese Harter Fuller Model

▶ Under model (1) the small area means are $\boxed{\mu_d = \overline{\mathbf{X}}'_d \beta + u_d + \bar{\varepsilon}_d}$ (cf. Pfeffermann (2013))

▶ This may be approximated as $\mu_d \approx \overline{\mathbf{X}}'_d \beta + u_d$

▶ For small sampling fractions the empirical best linear unbiased predictor follows as

$$\hat{\mu}_d^{BHF} = \overline{\mathbf{X}}'_d \hat{\beta} + \hat{u}_d$$
$$\hat{u}_d = \hat{\gamma}_d \left( \bar{y}_d - \bar{\mathbf{x}}'_d \hat{\beta} \right) \qquad (2)$$
$$\hat{\gamma}_d = \frac{\hat{\sigma}_u^2}{\hat{\sigma}_u^2 + \frac{\hat{\sigma}_\varepsilon^2}{n_d}}.$$

▶ Upper-case notation refers to population values and lower-case notation to sample values; a hat indicates that the variable is estimated and a bar denotes the mean

# Battese Harter Fuller Model

$$\widehat{\mu}_d^{\text{BHF}} = \widehat{\gamma}_d \left( \overline{y}_d + (\overline{\mathbf{X}}_d - \overline{\mathbf{x}}_d)' \widehat{\beta} \right) + (1 - \widehat{\gamma}_d) \overline{\mathbf{X}}_d' \widehat{\beta} \qquad (3)$$

▶ Expressions (2) and (3) are equivalent
▶ As (3) indicates, the BHF-estimator may be viewed as a composite estimator of

> ▶ the survey regression estimator (multilevel-GREG under SRS)
> $$\overline{y}_d + (\overline{\mathbf{X}}_d - \overline{\mathbf{x}}_d)' \widehat{\beta}$$
> ▶ and the regression-synthetic component $\overline{\mathbf{X}}_d' \widehat{\beta}$
>
> $$\widehat{\gamma}_d = \frac{\widehat{\sigma}_u^2}{\widehat{\sigma}_u^2 + \frac{\widehat{\sigma}_\varepsilon^2}{n_d}}.$$
>
> with weights $\widehat{\gamma}_d$ and $(1 - \widehat{\gamma}_d)$.

▶ Unlike the composite estimators in the previous lecture, the weights of the EBLUP emerge as a simple ratio of the estimated variance components.

# Battese Harter Fuller Model

For finite populations with non-negligible sampling fractions equation (2) or (3) have to be replaced by:

$$\hat{\mu}_d^{BHF} = \frac{1}{N_d} \left[ \sum_{j \in S_d} y_{dj} + \sum_{j \notin S_d} \hat{y}_{dj} \right] \quad \text{with} \tag{4}$$

$$\hat{y}_{dj} = \mathbf{x}'_{dj}\hat{\boldsymbol{\beta}} + \hat{u}_d \tag{5}$$

Equation (4) is a general representation of an empirical best predictor (EBP) as well. The EBP generally comprises two parts

1. The sum of the observations for the sampled units and
2. The sum of the predictions for the non-sampled units.

# BHF model: MSE

$\hat{\theta}_i$   statistics of interest at area *i* level ( or group *d* level):
say the area mean, area poverty rates , estimated by BHF model

Next step is to derive an MSE estimator

- $MSE(\hat{\theta}_i) \approx g_{1i}(\sigma) + g_{2i}(\sigma) + g_{3i}(\sigma))$

- $g_{1i}(\sigma) = \alpha'_r Z_r T_s Z'_r \alpha_r$

- $g_{2i}(\sigma) =$
  $[\alpha'_r bX_r - \alpha'_r Z_r T_s Z'_s R'_s X_s](X'_s V^{-1} X_s)^{-1}[X'_r \alpha_r - X'_s R_s^{-1} Z_s T_s Z'_r \alpha_r]$

- $g_{3i}(\sigma) = tr\{(\nabla(\alpha'_r Z_r \Sigma_u Z'_s V_s^{-1})')')V_s(\nabla(\alpha'_r Z_r \Sigma_u Z'_s V_s^{-1})')')'E[(\hat{\sigma} - \sigma)(\hat{\sigma} - \sigma)']\}$

- $T = \Sigma_u - \Sigma_u Z'_s(\Sigma_{es} + Z_s \Sigma_u Z'_s)^{-1} Z_s \Sigma_u$

- $\sigma = (\sigma_e^2, \sigma_u^2/\sigma_e^2)'$

# BHF model: MSE

Finally, the estimator for the MSE of $\hat{\theta}_i$ is

$$\widehat{MSE}(\hat{\theta}_i) = g_{1i}(\hat{\sigma}) + g_{2i}(\hat{\sigma}) + 2g_{3i}(\hat{\sigma})$$

- $\hat{\sigma}$ is an unbiased estimator for $\sigma$

Remark: it is possible to obtain an estimate of the MSE using alternative techniques, such as bootstrap and jackknife

# BHF model: recap

- If $\hat{\sigma}_u^2$ is small, a small value $\hat{\gamma}_d$ results as well and more weight will be given to the regression synthetic component.
  The same holds if the area specific sample sizes $n_d$ is very small.

- The weighting factor will tend to 1 and more weight will be given to the survey regression estimators, if the random effects variance is large
  The same holds for large sample sizes $n_d$.

- The BHF-estimator is not design-consistent for general survey designs
  $\longrightarrow$ Survey regression estimator does not consider design-weights

- The EBLUP is not model unbiased when conditioning on $u_d$, as this implies assuming fixed intercepts in the different areas (see Pfeffermann 2013).

- Examples  of application of the BHF EBLUP estimator during the R lab

Pros and cons of the model will be discussed after examples and applications to real data