

Some Model-based Estimator

Danute Krabavickaitė and Vilma Nekrasaitė

Vilnius, Lithuania

Survey sampling:

Finite population is a fixed, mammoth parameter

Sampler does a dance around it: picks this, omits that, introducing randomness

Survey design and inference – the study of that dance

Main branches of statistics:

Data is a realization of random variables

Statistician makes inference about probability laws of these random variables

Something is wrong with this principle in survey sampling!

not of the outcome of that process.

Sampling variance is a property of process of selecting the sample and

and not dependent on parameters about which we must make inferences.

on observed random variables whose probability distribution is known

statistical inference is made conditionally

Conditionality principle used often in statistical inference:

They differ much if $\frac{n}{N} = \frac{1}{2}$

$$\text{Var}_{SRS}(t_y) = N^2 \left(1 - \frac{n}{N}\right) \frac{s^2}{n}, \quad \text{Var}_{SRW}(t_y) = N^2 \frac{s^2}{N}$$

$$t_y = \frac{1}{N} \sum_{k=1}^N y_k$$

Sample $i \in U$ size n drawn according to SRS and SRW designs:

Parameter of a population, for example, total $t_y = \sum_{k=1}^N y_k$, fixed, unknown

Study variable y : $\{y_1, y_2, \dots, y_N\}$,

Finite population $u = \{1, 2, \dots, N\}$

Nothing is fixed in the world.

What statistical model can describe the process
which has generated the finite population?

Other statisticians are interested in estimating the parameters of the model

For the survey statistician when the model is chosen
estimation of model parameters plays some assistant role for the main task:

to make inference about the realized finite population itself.

- observational studies
- clinical trials
- designed experiments

Question: whether inferences should be based on distributions created by randomization or on probability models arises when interpreting data from

Valid inference can be done with and without it, and when randomisation is present, it does not necessarily creates the best probabilistic approach for inference.

Randomisation is desirable but not necessary nor sufficient for statistical inference.

- Bayesian prediction techniques (the posterior distribution of t_y is calculated when y_k , $k \in i$ are given) Bolafriene and Zacks, 1992; Ericson, Ghosh, Pfeferman, Royal
- Fiducial inference (palyginomial analyse), Kalbfleisch and Sprott, 1992
- Likelihood prediction, Björnstad, Royal
- Prediction based on general linear models. Valliant, Dorfman, Royall, 2000; Cham-
- bers

Model-based approaches:

Prediction Theory and the General Linear Model

Finite population $u = \{1, 2, \dots, N\}$

The population vector of values of the variable y : $y = (y_1, \dots, y_N)'$

It is treated as realization of the **random vector** $\mathbf{Y} = (Y_1, \dots, Y_N)'$

The goal – to estimate a linear combination of y_i 's:

$\gamma' \mathbf{Y}$, $\gamma = (\gamma_1, \dots, \gamma_N)'$ – vector of constants

Separate cases:

$\gamma = (1, \dots, 1)'$ $\iff \gamma' \mathbf{y} = t_y$ population total

$\gamma = (\frac{1}{N}, \dots, \frac{1}{N})'$ $\iff \gamma' \mathbf{y} = \mu_y$ population mean

Let i consists of the first n population elements. Then:

Any sample $i \in \mathcal{U}$

$$\boldsymbol{y} = (\boldsymbol{y}_i, \boldsymbol{y}_{u|i})'$$

$$\boldsymbol{Y} = (\boldsymbol{Y}_i, \boldsymbol{Y}_{u|i})'$$

$$\boldsymbol{\gamma} = \boldsymbol{\gamma}_i + \boldsymbol{\gamma}_{u|i} y_{u|i}$$

$$\boldsymbol{\gamma}_i = \boldsymbol{\gamma}_i^i y_i + \boldsymbol{\gamma}_{u|i}^i y_{u|i}$$

is a realization of

$\boldsymbol{\gamma}_i = (\boldsymbol{\gamma}_i^i, \boldsymbol{\gamma}_{u|i}^i)'$

Estimation target

$\boldsymbol{\gamma} = (\boldsymbol{\gamma}_i, \boldsymbol{\gamma}_{u|i})'$

Estimation of $\boldsymbol{\gamma}_i$ \iff prediction of $\boldsymbol{\gamma}_{u|i}$.

The term „**prediction**“ is used in the sense of making statistical guesses about the values of \boldsymbol{Y}_k , which we have not seen – not in the sense of forecasting the future values.

its expression and error variance is known.

$\text{Var}_M(\hat{\theta})$ is called best linear unbiased predictor (BLU) of θ .

Def. 6. A linear model unbiased estimator $\hat{\theta}$ of θ which minimizes the error variance

Def. 5. The error variance of $\hat{\theta}$ under a model M is $\text{Var}_M(\hat{\theta}) = \mathbb{E}^M(\hat{\theta} - \theta)^2$.

Def. 4. The estimator $\hat{\theta}$ is unbiased for θ under a model M if $\mathbb{E}^M(\hat{\theta} - \theta) = 0$.

Λ – positive definite variance-covariance matrix of \mathbf{Y} .

$\boldsymbol{\theta}$ – vector of unknown parameters,

$\mathbf{X} = (\mathbf{X}_{(1)}, \dots, \mathbf{X}_{(J)})$ – matrix of auxiliary variables,

$$\begin{aligned} \mathbb{E}^M(\mathbf{Y}) &= \mathbf{X}\boldsymbol{\theta}, \\ (1) \end{aligned} \quad \begin{aligned} \text{Var}_M(\mathbf{Y}) &= \Lambda, \end{aligned} \quad (2)$$

Def. 3. General linear model M for \mathbf{Y} :

$$\hat{\theta} - \theta = \mathbf{u}_i^T \mathbf{Y}_i - \gamma_i \mathbf{Y}.$$

Def. 2. The estimation error of an estimator $\hat{\theta} = \mathbf{u}_i^T \mathbf{Y}_i$ is

$\mathbf{u}_i = (u_{i1}, \dots, u_{in})'$ – vector of coefficients.

Def. 1. A linear estimator of $\theta = \gamma_i \mathbf{Y}_i$ is defined as $\hat{\theta} = \mathbf{u}_i^T \mathbf{Y}_i$,

Valiant, R., Dorfman A. H., and Royal R. M. (2000) *Finite Population sampling and inference: A prediction Approach*. John Wiley & Sons

.....

stratified expansion estimator,

linear regression estimator,

expansion estimator,

Taking various matrices of auxiliary variables X , matrices V and vectors β we obtain BLU predictors coinciding with the estimators of totals known in the design-based theory:

$$\text{Estimating } \gamma_y = t_y \iff \text{predicting value } \sum_{k \in U \setminus i} y_k \text{ of unobserved random variable } \sum_{k \in U \setminus i} y_k + \sum_{k \in i} y_k.$$

is equivalent to

$$\text{Let } \gamma = (1, \dots, 1) : \quad \gamma_y = \gamma_i y_i + \gamma_{U \setminus i} y_{U \setminus i}.$$

When X are quantitative variables the estimator $\hat{\beta}$ for β which gives BLU predictor of θ , is the least squares estimator.

Prediction theory under the nonlinear model

R. Valliant, JASA, 1985

Nonlinear model:

$$(3) \quad \mathbb{E}^M(Y_k) = f(\mathbf{X}_k; \boldsymbol{\varphi}), \quad k = 1, \dots, N$$

$$(4) \quad \text{Var}^M(Y) = V,$$

$\boldsymbol{\varphi}$ – vector of J unknown coefficients, model parameter,

$f(\cdot, \cdot)$ – nonlinear function, has at least 3 partial derivatives with respect to $\boldsymbol{\varphi}$.

Remark: if model is such that

Y_1, \dots, Y_N are independent random variables \iff

V is diagonal \iff

$\text{Var}^M(Y_k)$ is only needed for it in (4).

Population total

When ϕ is known, the BLU estimator of t_y is obtained by adding $\sum_{k \in U \setminus i} y_k$ with BLU predictor of $\sum_{k \in U \setminus i} y_k$.

$$t_y = \sum_{k=1}^N y_k + \sum_{k \in U \setminus i} y_k.$$

When ϕ is not known – its estimator is inserted into BLU predictor of $\sum_{k \in U \setminus i} y_k$.
Valliant: If asymptotically normal and consistent estimator of the model parameters ϕ is available
 t_y becomes biased.

1) Estimator \hat{t}_y is consistent

2) Expressions for

$$\begin{aligned} ABias(\hat{t}_y) &\approx E^M(\hat{t}_y - t_y), \\ Var(\hat{t}_y) &\approx Var(t_y - t_y). \end{aligned}$$

3) Consistent estimator of $Var M(\hat{t}_y)$ constructed through the consistent estimator of ϕ .

With the study variable because zero values can not be known in advance.

Model-assisted estimators are not useful – there are no auxiliary variables correlated

Design-based estimators are not useful

Investment of enterprise for an environmental protection

Area under some kind of crop in the farm, which is not grown up often (rape),

Examples:

$$\left. \begin{array}{l} y_k < 0, \quad k = 1, \dots, N \\ y_k = 0 \quad \text{or} \end{array} \right\} \Leftrightarrow \text{Var}(y) \text{ large} \Leftrightarrow \text{Var}(t_y) \text{ large.}$$

Let us consider a study variable y , values of which have the properties:

Models for study variable with positive values and many zeroes

$$\hat{y}_k = \begin{cases} 0, & 1 - p_k, \\ p_k, & \text{otherwise} \end{cases} \Leftrightarrow \hat{y}_k = \mathbb{E}_{M^k} Y_k$$

$$Y_k = \hat{y}_k Y_k, \quad Y_k = x_k' g + e_k, \quad e_k \sim N(0, \sigma^2)$$

I. F. Kalberer, *Journal of Official Statistics*, 2000

Nonlinear prediction $\hat{y}_k = ?$

$$\hat{y}_k = \sum_{k \in U \setminus i} y_k + \sum_{k \in i} y_k, \quad t_y = \sum_{k \in U \setminus i} y_k + \sum_{k \in i} y_k,$$

$y_k, k = 1, 2, \dots, N, t_y$ – random variables

U_* – superpopulation or a model for y , U – its *random realization*,

Model-based estimator

Econometric models

2. Censored regression or Tobit model for a study variable y .

Let there are unobserved random variables

$$Y_*^k = \mathbf{x}_k' \boldsymbol{\beta} + \epsilon_k^k, \quad \epsilon_k^k \sim N(0, \sigma^2), \quad k = 1, 2, \dots, N$$

ϵ_k^k , ϵ_l^k independent for $k \neq l$.

$$Y_k = \begin{cases} Y_*^k, & \text{if } Y_*^k < 0, \\ 0, & \text{if } Y_*^k \geq 0, \end{cases} \quad k = 1, 2, \dots, N$$

Define random variables

W.H. Greene. Econometric Analysis. Prentice Hall, Upper Saddle River, 2003.

$$\mathbb{E}^M(Y_k) = f(\mathbf{x}_k; \boldsymbol{\phi}), \quad \Phi(\boldsymbol{\phi}) = \left(\frac{\partial}{\partial \mathbf{x}_k' \boldsymbol{\beta}} \right) \Phi(\boldsymbol{\phi}) + \left(\frac{\partial}{\partial \mathbf{x}_k' \boldsymbol{\beta}} \right) \Phi(\mathbf{x}_k) f(\mathbf{x}_k; \boldsymbol{\phi})$$

$$Var^M(Y_k) = \sigma^2 \Phi(\boldsymbol{\phi}) (\alpha_0 + \alpha_1 x_k) g_k(\alpha_0, \alpha_1), \quad k = 1, \dots, N$$

$$\boldsymbol{\phi} = (\sigma, \alpha_0, \alpha_1), \quad \alpha_0 = \beta_0/\sigma, \quad \alpha_1 = \beta_1/\sigma.$$

If $\mathbf{x}_k = (1, x_k)$:

T. Amemiya, *Econometrica*, 1973.

Maximum likelihood estimators $\hat{\alpha}_0, \hat{\alpha}_1$ of tobit model parameters α_0, α_1 are consistent and asymptotically normal.

We estimate

$$\hat{y}_{(tobit)} = \hat{E}_{(tobit)} Y_i = \hat{\beta}'_i X_i + \hat{\epsilon}_i$$

Using results of Valiant we obtain expressions for

SAS procedure LIFEREG: $\hat{\beta}, \hat{\sigma}$.

$$\hat{Var}(\hat{t}_{(tobit)}^y)$$

$$AVar(\hat{t}_{(tobit)}^y)$$

$$ABias(\hat{t}_{(tobit)}^y)$$

If both equations in (5) coincide \iff Heckman \equiv Tobit.

$$I^k = \begin{cases} 0 & \text{otherwise.} \\ 1, & \text{if } Z_*^k < 0, \\ Y_*^k, & \text{if } I^k = 1, \end{cases}$$

$$\left(\begin{pmatrix} 1 & \varrho_{\varepsilon n} \\ \varrho_{\varepsilon n} & \varrho^2 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \end{pmatrix} \right) \mathcal{N} \sim \begin{pmatrix} n \\ \varepsilon \end{pmatrix}$$

$$(6) \quad Y_*^k = \mathbf{x}_{(2)}' \boldsymbol{\beta}_{(2)} + u^k, \quad \mathbf{x}_{(2)}' \boldsymbol{\beta}_{(2)} = Z^k$$

$$(5) \quad Y^k = \mathbf{x}_{(1)}' \boldsymbol{\beta}_{(1)} + \varepsilon^k, \quad \mathbf{x}_{(1)}' \boldsymbol{\beta}_{(1)} = 0$$

Two unobserved random variables:

Maddala G. S. Limited-dependent and qualitative variables in econometrics. Cambridge University Press, 1983.

3. Censored regression model with the unobserved stochastic threshold (Heckman model)

how big will be Y^k if "yes"?

to have positive values of Y^k or not to have?

The decisions can be separated:

Nonlinear regression

Other generalizations of nonlinear model (3), (4) when residuals are not normally distributed – skew distributions in economics.

$$\sum_{k \in \mathcal{E} \setminus i} y_k + \sum_{k \in \mathcal{E} \setminus i} = \underline{\mathbb{E}}^M_{(Heckman)}(Y_K), \quad t_{(Heckman)} y_k$$

Estimation of the model parameters

$$\left(\frac{(\mathbf{x}_{(2)}' \boldsymbol{\theta}_{(2)} \mathbf{x}) \boldsymbol{\phi}}{(\mathbf{x}_{(2)}' \boldsymbol{\theta}_{(2)})} + \mathbf{x}_{(1)}' \boldsymbol{\theta}_{(1)} \right) (\mathbf{x}_{(2)}' \boldsymbol{\theta}_{(2)} \mathbf{x}) \boldsymbol{\phi} = \underline{\mathbb{E}}^M(Y) = f(\mathbf{x}_{(1)}, \mathbf{x}_{(2)}, \boldsymbol{\theta}_{(1)}, \boldsymbol{\theta}_{(2)})$$

Simulation

$N = 1000$ size finite population generated

$$\mathbf{x}^k = (1, x_k)', \quad x_k \sim N(\mu_k, \sigma_k^2), \quad \varepsilon_k \sim N(0, \sigma_k^2)$$

$$Y_* = \beta_0 + \beta_1 x_k + \varepsilon_k$$

$M = 1000$ replicates:

Simple random sample of size n

Estimates: design-based and tobit model-based

$$y_{(tobit)} = f(x_k; \beta_0, \beta_1), \quad k \in \mathcal{U} \setminus i$$

$$\hat{\beta}_0^{(tobit)} = \sum_{k \in I} y_k + \sum_{k \in U \setminus I} \frac{u}{N}$$

Dependence of $RMS\hat{e}$ on n , α , percentage of zeroes in the population.

Investigated:

$$\begin{aligned}
 & MSE(\hat{\eta}) = \sqrt{\frac{t}{MSE(\hat{\eta})}} \\
 & MSE(\hat{\eta}) = Bias(\hat{\eta})^2 + Var(\hat{\eta}) \\
 & MSE_{emp}(\hat{\eta}) = Bias_{emp}(\hat{\eta})^2 + Var_{emp}(\hat{\eta}) \\
 & Var_{emp}(\hat{\eta}) = \frac{1}{M} \sum_{m=1}^{M-1} (\hat{\eta}_m - \bar{\hat{\eta}})^2, \\
 & Bias_{emp}(\hat{\eta}) = \bar{\hat{\eta}} - t, \quad \text{for } \hat{\eta} = t(SRS), \\
 & Bias_{emp}(\hat{\eta}) = \bar{\hat{\eta}} - t, \quad \text{for } \hat{\eta} = t(tobit), \\
 & \hat{\eta} = \frac{1}{M} \sum_{m=1}^M \hat{\eta}_m, \\
 & \hat{\eta} = t_{(SRS)} \quad \text{or} \quad \hat{\eta} = t_{(tobit)}
 \end{aligned}$$

Notations: $\hat{\eta}_m$ and $Var(\hat{\eta}_m)$ – replicates of the estimates $\hat{\eta}$:

$$\mathbf{x}_k = (1, x_k)', \quad x_k \sim N(3.49, 1.56), \quad \epsilon_k \sim N(0, 0.8), \quad \beta = (-2.42, 1.1), \quad 20\% \text{ zeroes}$$
$$Y_k^* = -2.42 + 1.1x_k + \epsilon_k \quad \Leftarrow \quad Y_k$$

Fig. 1. Population

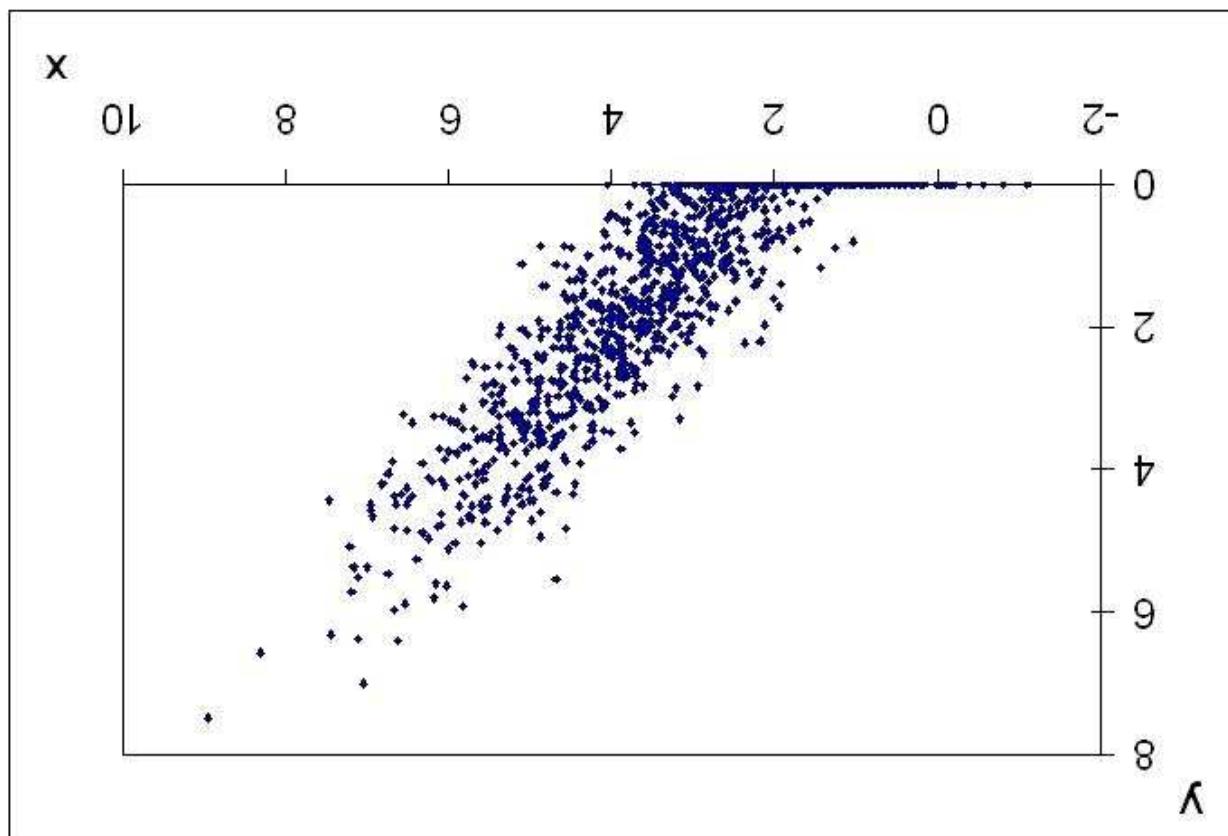


Fig. 2. Dependence of $\hat{t}_y(SRS)$ (left) and $\hat{t}_y(tobit)$ (right) on the sample size n

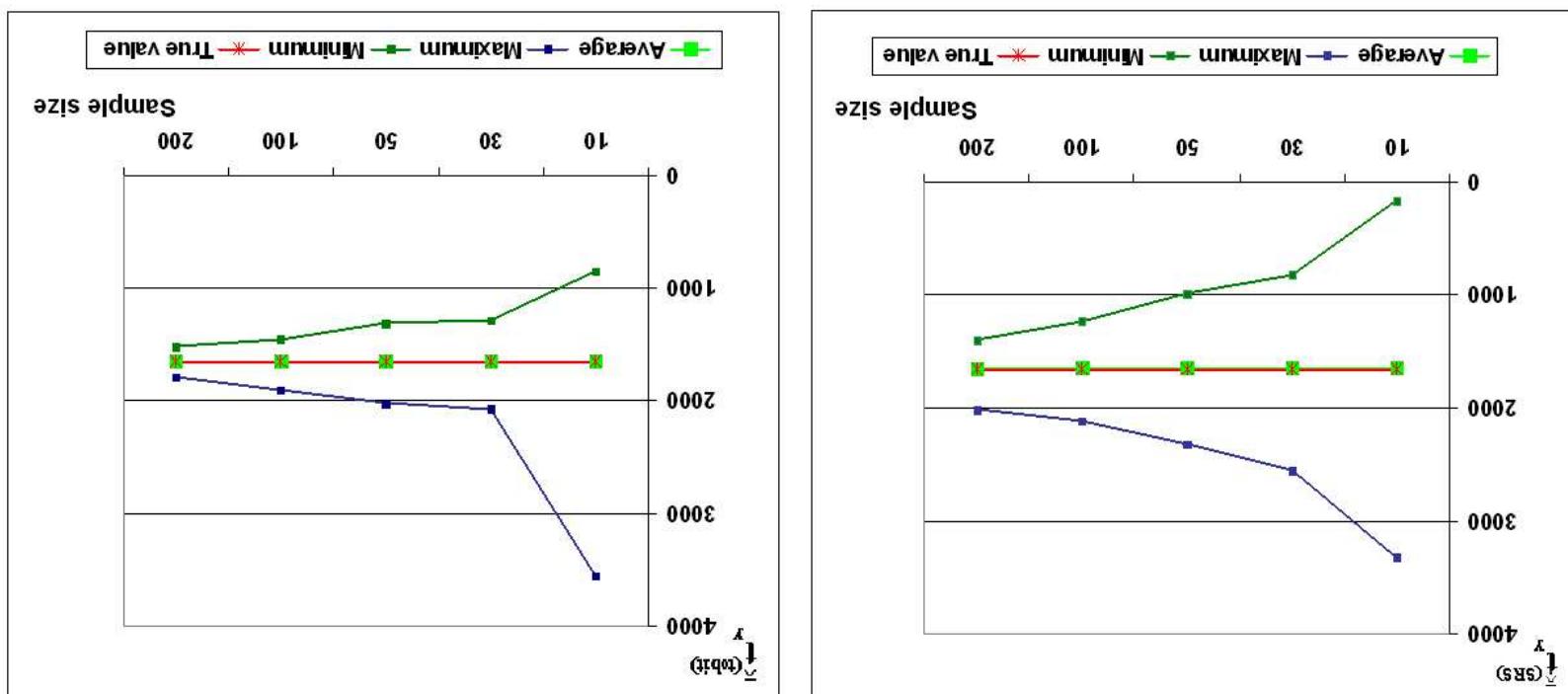


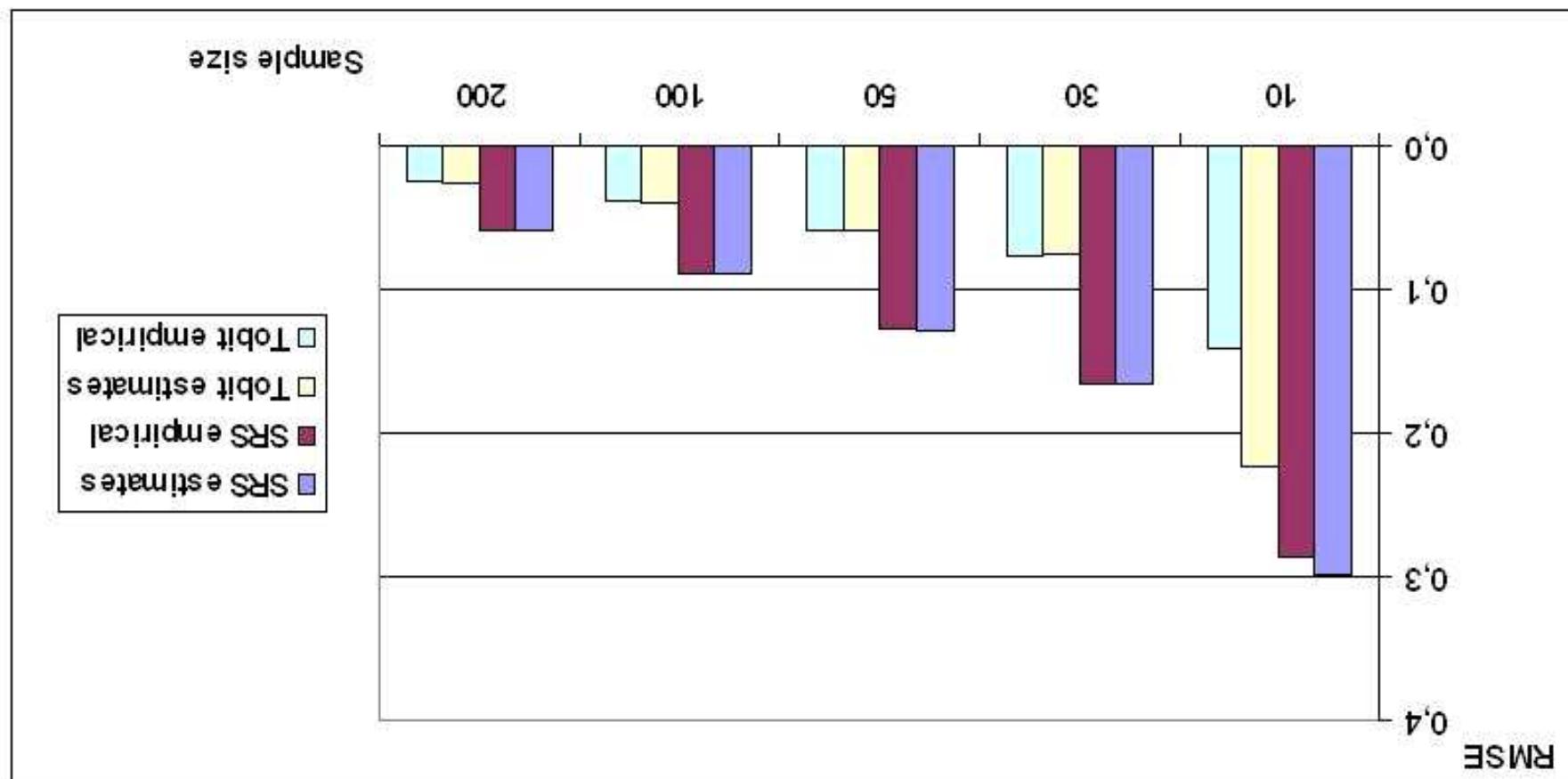
Fig. 3. Dependence of $RMS{E}(\hat{t}_y)$ on the sample size n 

Fig. 4. Dependence of $\hat{\sigma}_2^2$, \hat{B}_0 , \hat{B}_1 on the sample size n

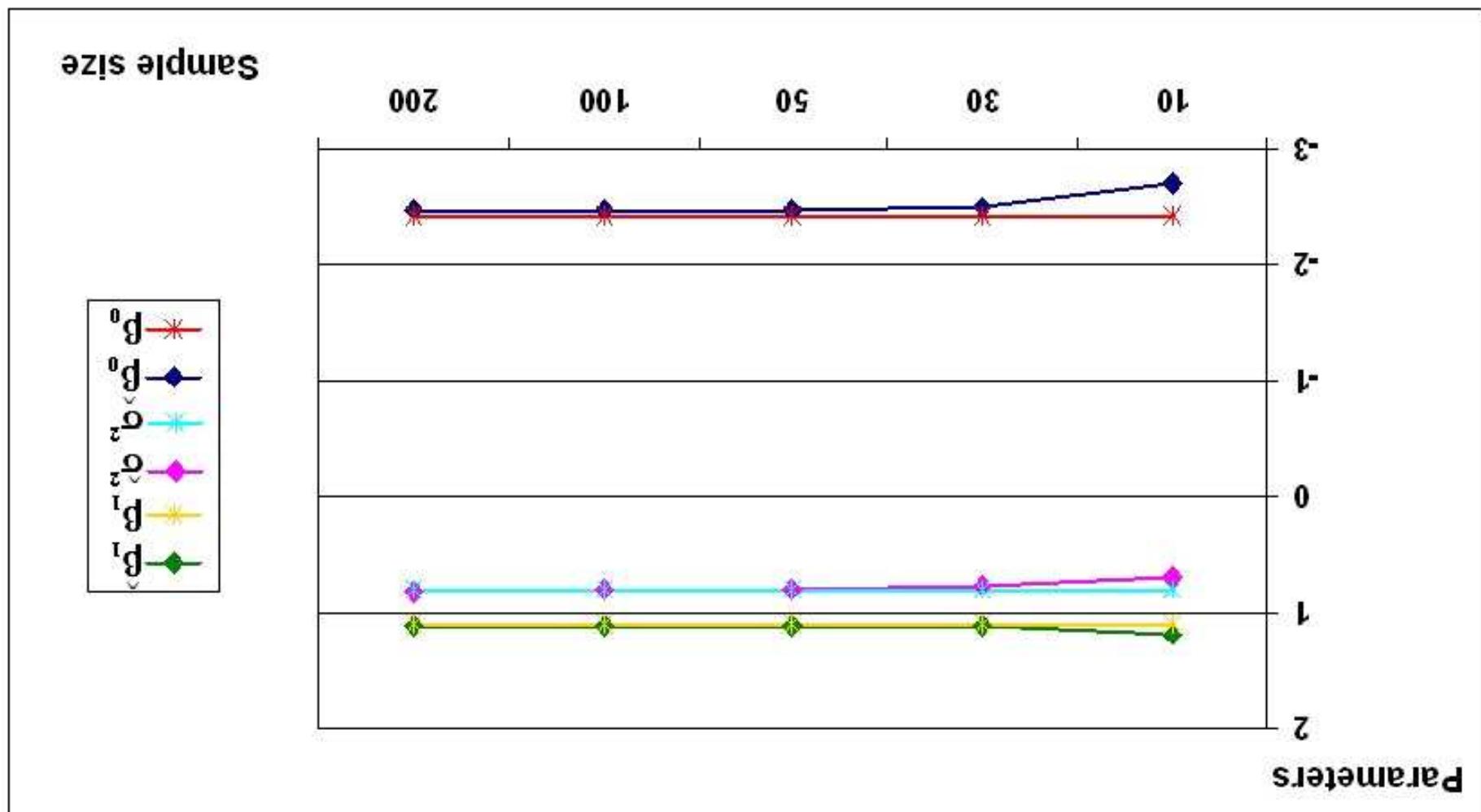
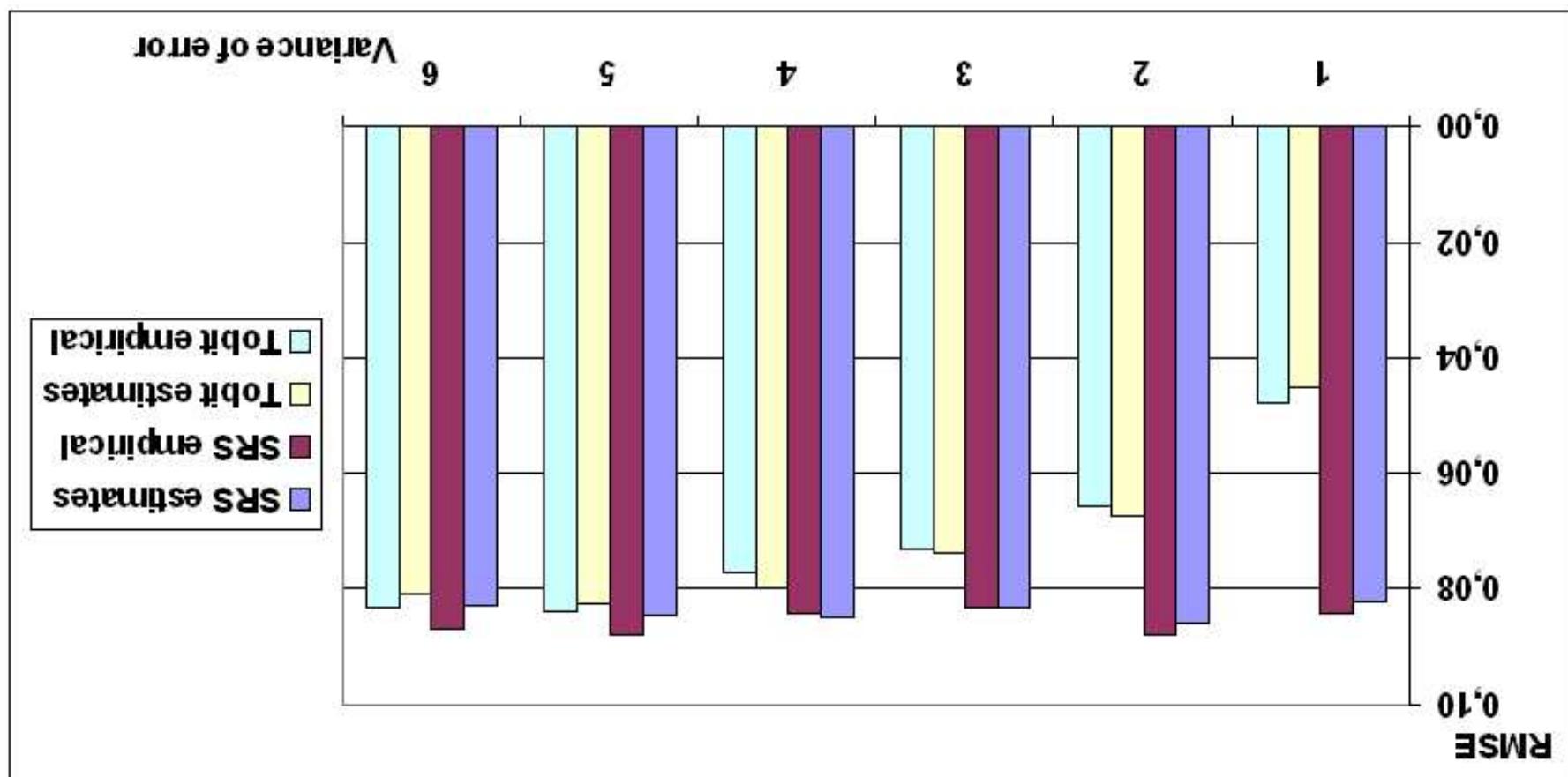
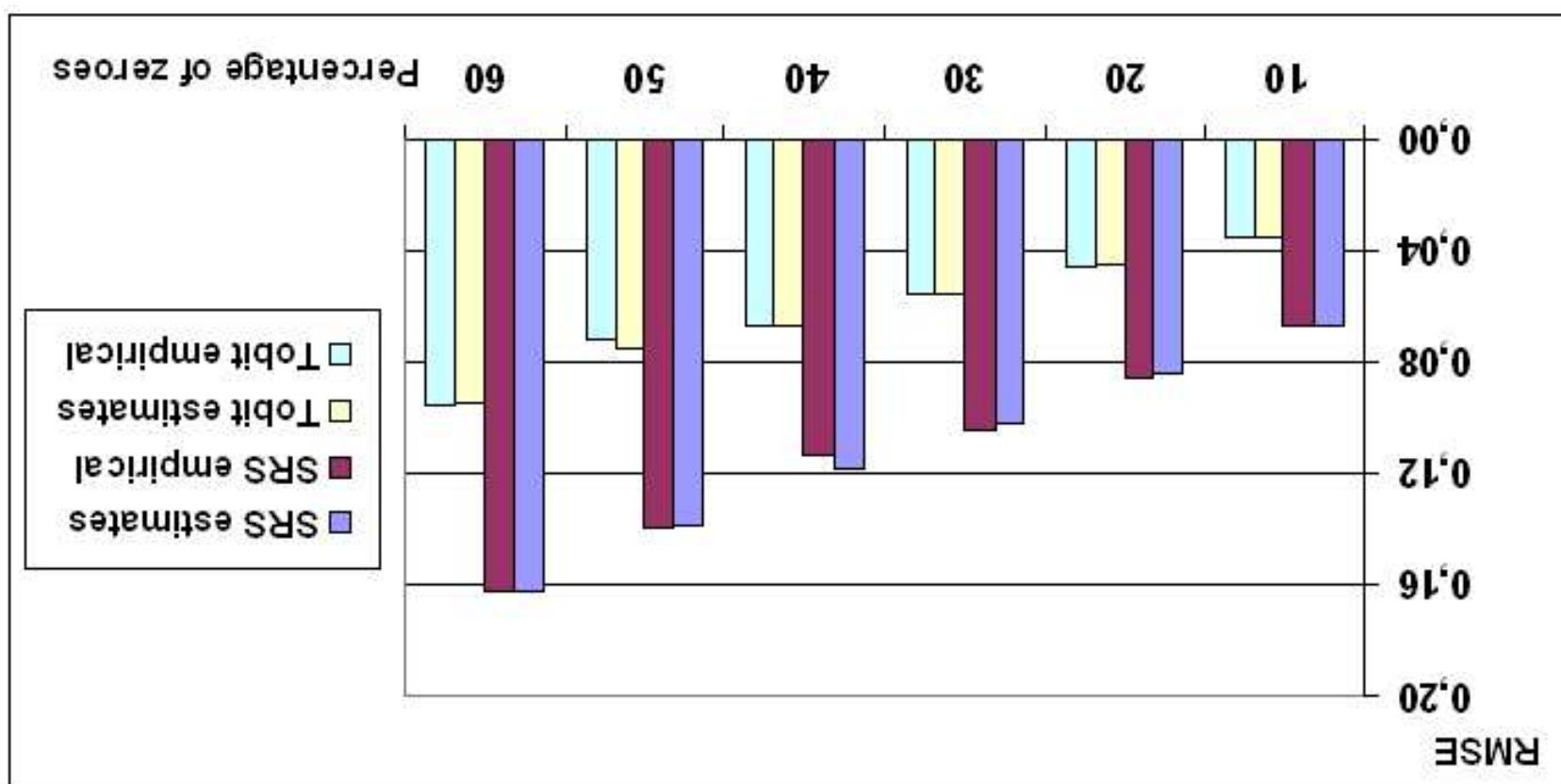


Fig. 5. Dependence of $\text{RMSE}(t_y)$ on the variance of error σ ($\sigma \downarrow \Leftrightarrow G_0 \downarrow$)



(zeroes $\downarrow \Leftarrow u \uparrow$)

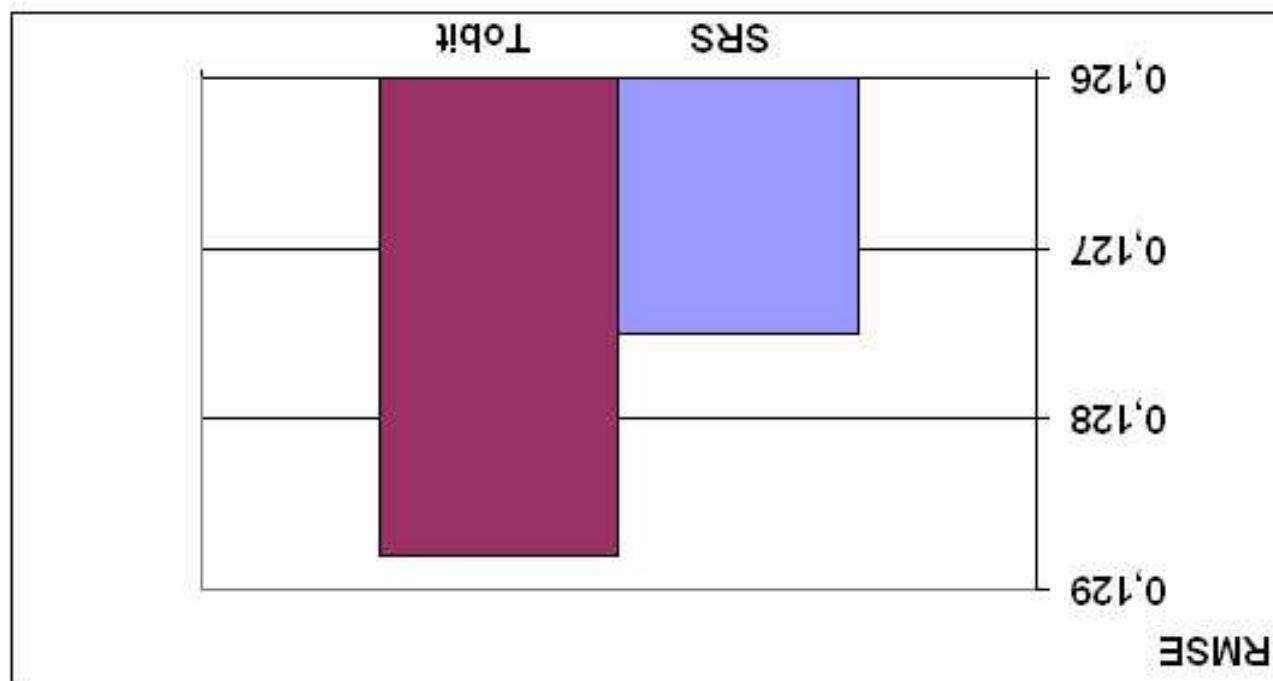
Fig. 6. Dependence of $RMS{E}(t_y)$ on the percentage of zeros in the population



$$y_k^* = 0.5 + 0.3x_k + \varepsilon_k \quad \Leftarrow \quad y_k$$

$$\mathbf{x}_k = (1, x_k)', \quad x_k \sim N(1, 2), \quad \varepsilon_k \sim N(0, 6), \quad \beta = (0.5, 0.3), \quad 45\% \text{ zeroes}$$

Fig. 7. The case when $RMS{E}(t_{Tobit}^y) > RMS{E}(t_{SRS}^y)$



Conclusions

For small sample size ($n = 10$) average $RMS{E}$ and empirical $RMS{E}$ of model-based estimator can differ much.
The efficiency of the model-based estimator does not show any dependence on the sample size or percentage of zeroes in the population.

Acknowledgements

Support: Nordic Council of Ministers

Library facilities: Umeå University

Thank you for your attention!