

On sample allocation for effective EBLUP estimation of small area totals

Valmiera workshop, August 26, 2012

Basic assumptions

This research uses stratified sampling (stratum = area).

Sampling units inside strata have been selected with SRSWOR sampling method.

Overall sample size (n) is fixed and small (limited resources).

Effective estimation

Selected estimation method (direct – indirect with different variations) should produce such area estimates for means, totals, proportions etc. which have as low sampling errors as possible. The objectives can be determined on area or population level.

Measuring sampling errors:

- Variance, MSE, CV
- Quality measures (ARE, ARB, ASE, RRMSE, EFF)
- Coverage of confidence intervals

Optimal allocation?

Overall sample size (n) is allocated for d areas in order to reach pre-set optimization criteria concerning sampling errors. In many cases it is a question of minimizing a mathematical expression as a function of sample sizes of areas under certain constraints.

The expression can contain variances, MSE's etc.

Constraints: 1) $\sum_d n_d = n$
2) sampling error of each area < given limit
etc.

Earlier approaches to reach optimal allocation in analytical way

Main attention in optimal allocation has been focused on direct estimation so far.

Longford (2006): areas have different priorities (weights).

Falorsi and Righi (2008, 2011): basic domains divided into different partitions and balanced sampling technique.

Khan et al. (2010): several response variables and one auxiliary variable for each; minimization of increment of variance mean.

Keto and Pahkinen (Katowice 2009) have used experimental allocation in model-based EBLUP estimation. The idea was to find out topics for further research.

Example of allocation problem

A big reform concerning local administration (municipalities) has been started in Finland. Among other things the number of municipalities will be reduced from present value 336 as down as 70-100. What do the people in municipalities think about this reform?

Suppose that a nationwide survey research is carried out. If overall sample size is "normal" 2 000 (avg. 6/municipality) and sample allocation is proportional, what would it mean? Helsinki would take 10 % of 2 000, and many small municipalities will have zero sample size!

Basic problem in optimal allocation

Basic problem in area statistics is that area level in general has not been taken into account in sampling design, and there may appear "zero" areas ($n_d = 0$). This forces to apply model-based estimation. Well-known are hierarchical models which use EBLUP estimation. In this research sampling design should lead to optimal area estimation from the point of view of selected model.

This research searches for analytical solution of optimal allocation problem conditional to selected model.

One example of analytical solution in a simple case (regression model) is presented in CP (by M Keto).

Conventional, widely used allocations

In the formulas we assume usage of auxiliary variable (x).

General restriction:
$$\sum_d n_d = n$$

Equal allocation:
$$n_{d, equ} = n/D \text{ (} D = \text{number of areas)}$$

Proportional allocation:
$$n_{d, pro} = (N_d/N)n$$

Optimal (Neyman) allocation:
$$n_{d, opt} = (N_d S_{dx} / \sum_{d=1}^D N_d S_{dx}) n$$

Power allocation:
$$n_{d, pow} = (X_d^a CV(x)_d / \sum_{d=1}^D X_d^a CV(x)_d) n,$$

where $X_d^a = \text{sum of } x\text{-values in area } d$

X – total allocation:
$$n_{d, tot} = (\sum_{k=1}^{N_d} x_{dk} / \sum_{d=1}^D \sum_{k=1}^{N_d} x_{dk}) n$$

None of these allocations is based on a specific model.

Used model

This research uses nested-error regression, basic unit level model which is a special case of general linear model:

$$y_{dk} = \mathbf{x}'_{dk} \boldsymbol{\beta} + v_d + e_{dk}; \quad k = 1, \dots, N_d; d = 1, \dots, D$$

v_d : *random area effect, mean = 0, variance = σ_v^2*

e_{dk} : *random area effect, mean = 0, variance = σ_e^2*

General theory of this model (estimation of variance components, regression coefficients, area effects etc.) is well-known and many times applied.

EBLUP estimates and MSE

EBLUP estimate for area total of response variable y :

$$\hat{Y}_{d,EBLUP} = \sum_{k \in s_d} y_{dk} + \sum_{k \in \bar{s}_d} \hat{y}_{dk} = \sum_{k \in s_d} y_{dk} + \sum_{k \in \bar{s}_d} \mathbf{x}'_{dk} \hat{\boldsymbol{\beta}} + (N_d - n_d) \hat{v}_d$$

EBLUP estimator is biased \rightarrow MSE is used instead of variance:

$$MSE(\hat{Y}_{d,Eblup}) = E(\hat{Y}_{d,Eblup} - Y_d)^2 = Var(\hat{Y}_{d,Eblup}) + (\hat{Y}_{d,Eblup} - Y_d)^2$$

Prasad-Rao approximation of MSE for finite populations:

$$mse(\hat{Y}_{d,EBLUP}) = g_{1d}(\hat{\sigma}_v^2, \hat{\sigma}_e^2) + g_{2d}(\hat{\sigma}_v^2, \hat{\sigma}_e^2) + 2g_{3d}(\hat{\sigma}_v^2, \hat{\sigma}_e^2) + g_{4d}(\hat{\sigma}_v^2, \hat{\sigma}_e^2)$$

Four components of MSE approximation:

$$g_{1d}(\hat{\sigma}_v^2, \hat{\sigma}_e^2) = (N_d - n_d)^2 (1 - \gamma_d) \hat{\sigma}_v^2$$

$$g_{2d}(\hat{\sigma}_v^2, \hat{\sigma}_e^2) = (N_d - n_d)^2 (\bar{\mathbf{x}}_d^* - \gamma_d \bar{\mathbf{x}}_d)' (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} (\bar{\mathbf{x}}_d^* - \gamma_d \bar{\mathbf{x}}_d)$$

$$g_{3d}(\hat{\sigma}_v^2, \hat{\sigma}_e^2) = (N_d - n_d)^2 n_d^{-2} (\hat{\sigma}_v^2 + \hat{\sigma}_e^2 n_d^{-1})^{-3} [\hat{\sigma}_e^4 \text{Var}(\hat{\sigma}_v^2) + \hat{\sigma}_v^4 \text{Var}(\hat{\sigma}_e^2) - 2\hat{\sigma}_e^2 \hat{\sigma}_v^2 \text{Cov}(\hat{\sigma}_e^2, \hat{\sigma}_v^2)]$$

$$g_{4d}(\hat{\sigma}_v^2, \hat{\sigma}_e^2) = (N_d - n_d) \hat{\sigma}_e^2$$

Ratio γ_d : $\gamma_d = \hat{\sigma}_v^2 / (\hat{\sigma}_v^2 + \hat{\sigma}_e^2 n_d^{-1}) = n_d \hat{\sigma}_v^2 / (n_d \hat{\sigma}_v^2 + \hat{\sigma}_e^2)$

Common intra-area correlation ρ :

$$\hat{\rho} = \hat{\sigma}_v^2 / (\hat{\sigma}_v^2 + \hat{\sigma}_e^2) = 1 / (1 + \hat{\sigma}_e^2 / \hat{\sigma}_v^2).$$

Optimization criterion

Basic criterion for optimization: minimize the arithmetic mean of areal MSE approximations

$$1/D \sum_{d=1}^D mse(\hat{Y}_{d,EBLUP})$$

subject to the constraint of fixed overall sample size

$$\sum_{d=1}^D n_d = n .$$

The model is used as given information when searching for optimum.

Use of component g_1

Because of the complexity of whole MSE approximation the optimum is impossible to reach. We turn our attention to the first and most important component g_1 :

$$g_{1d}(\hat{\sigma}_v^2, \hat{\sigma}_e^2) = (N_d - n_d)^2 (1 - \gamma_d) \hat{\sigma}_v^2$$

If variation between areas is strong enough and the model is suitable for estimation, then the proportion

$$100 \times (g_{1d}(\hat{\sigma}_v^2, \hat{\sigma}_e^2) / \text{mse}(\hat{Y}_{d,EBLUP})) \%$$

reaches easily 85-90 %, often as much as 95 % according to for ex. Nissinen (2009). Now it is reasonable to find minimum for the mean of area g_1 values.

Minimization problem

Minimize expression

$$1/D \sum_{d=1}^D g_{1d}(\hat{\sigma}_v^2, \hat{\sigma}_e^2) = 1/D \sum_{d=1}^D (N_d - n_d)^2 (1/\hat{\sigma}_e^2 \times n_d + 1/\hat{\sigma}_v^2)^{-1}$$

with respect of n_d subject to constraint

$$\sum_d n_d = n .$$

Minimum is obtained by using Lagrange's multiplier method.

Solution (is not shown, but can be proved):

$$\begin{aligned}
 n_{d,opt} &= -\hat{\sigma}_e^2 / \hat{\sigma}_v^2 + \frac{(N_d + \hat{\sigma}_e^2 / \hat{\sigma}_v^2)(n + D(\hat{\sigma}_e^2 / \hat{\sigma}_v^2))}{N + D(\hat{\sigma}_e^2 / \hat{\sigma}_v^2)} = -\delta + \frac{(N_d + \delta)(n + \delta D)}{N + \delta D} \\
 &= \frac{N_d n - (N - N_d D - n)\delta}{N + D\delta} = \frac{N_d n - (N - N_d D - n)(1/\hat{\rho} - 1)}{N + D(1/\hat{\rho} - 1)},
 \end{aligned}$$

where ratio of variance components is $\delta = \hat{\sigma}_e^2 / \hat{\sigma}_v^2$

and intra-area correlation is $\hat{\rho} = \hat{\sigma}_v^2 / (\hat{\sigma}_v^2 + \hat{\sigma}_e^2) = 1 / (1 + \hat{\sigma}_e^2 / \hat{\sigma}_v^2)$.

Because intra-area correlation depends on values of response variable y , we have to replace this correlation with a value produced from x -values and which measures the proportion of variation between areas and total variation. The reasoning is that same variation in x -values transfers to the sample.

Homogeneity measure

We know from cluster sampling with unequal clusters: First: simple ANOVA for auxiliary variable x and then the measure

$$R_a^2 = 1 - R^2 = 1 - \frac{MSW}{S^2},$$

where R^2 is coefficient of determination (regr. analysis), MSW is mean SS of clusters (strata) and S^2 is variance of x .

Remark: ratio SS_B / SS_{tot} is very close to homogeneity measure
Also that can be used.

Computational values and extreme case

Computational sample sizes are not integers (very likely). They are normally rounded to nearest integer (sometimes compromises have to be made).

If overall sample size (n) is small or/and size of area (N_d) is small computational sample size can become negative. This is of course a restriction.

If all variation is between areas, the result is proportional allocation, because ratio of variances $\delta = 0$ (and intra-area correlation $\rho = 1$).

Research data

Population:	9815 apartments
Areas:	34 Finnish towns (small – large)
Response variable (y):	Price of apartment (1 000 €)
Auxiliary variable (x):	Size of apartment (m^2)
xy-correlation in population:	0,674
Sizes of areas:	111 – 833
Homogeneity measure (of x):	0,33 (quite high → strong variation between areas)

Testing the performance of g_1 allocation: results vs results of "conventional" allocations

"Competing" allocations:

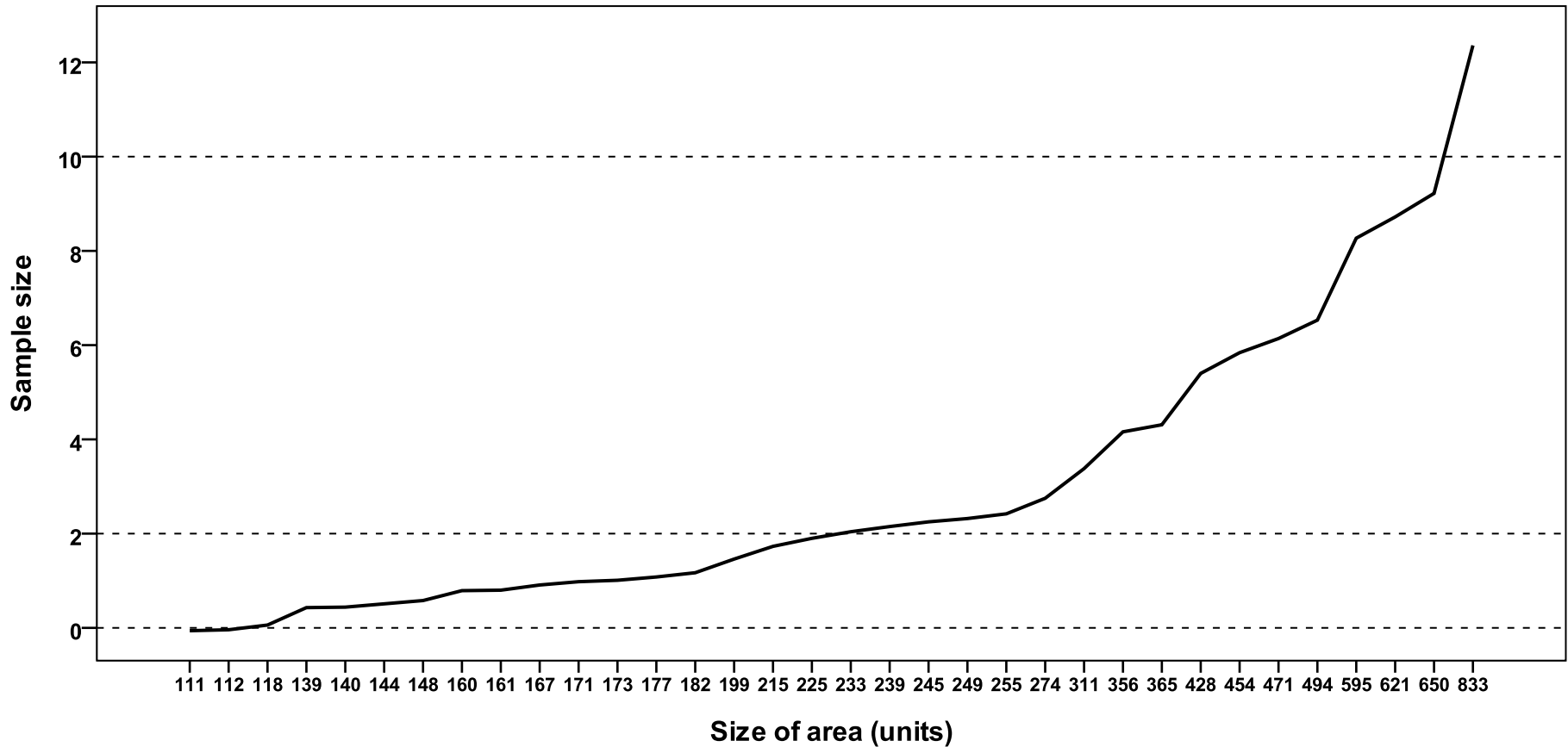
- equal, proportional and power
- g_1 -allocation

Results of Neyman allocation are not presented because its performance is clearly the poorest.

1500 random samples were simulated (with SAS program) for each allocation alternative, sampling method was SRSWOR inside strata (=area) and necessary statistics and quality measures were computed. Overall sample size was 102, 170 (original 34 areas) and 180 (15 combined areas).

Phase 1: All 34 original areas, $N = 9815$, $n = 102$ ($E(n_d) = 3$)
- 3 smallest areas: $n_d = 0$

Computational sample sizes for 34 areas in g1-allocation

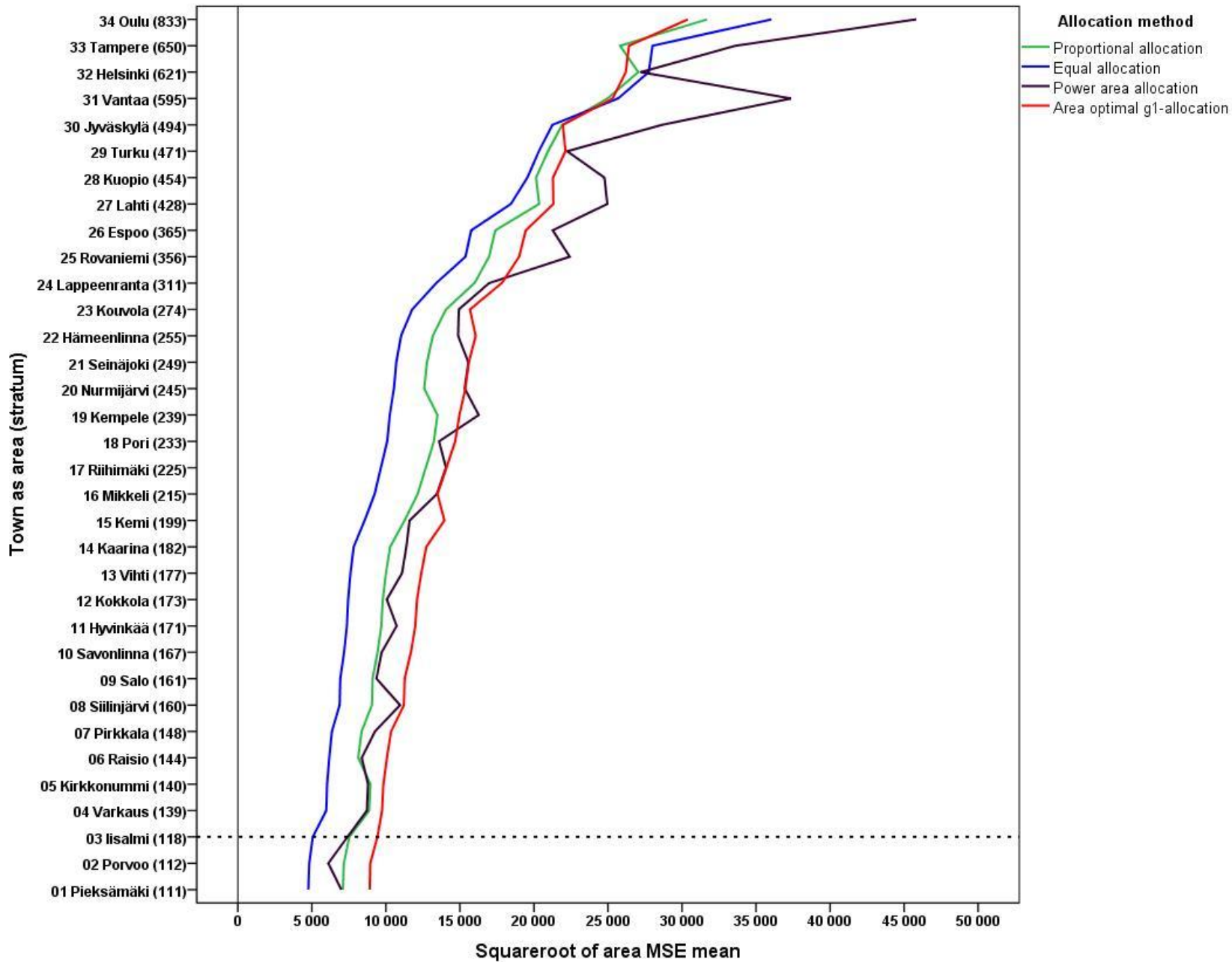


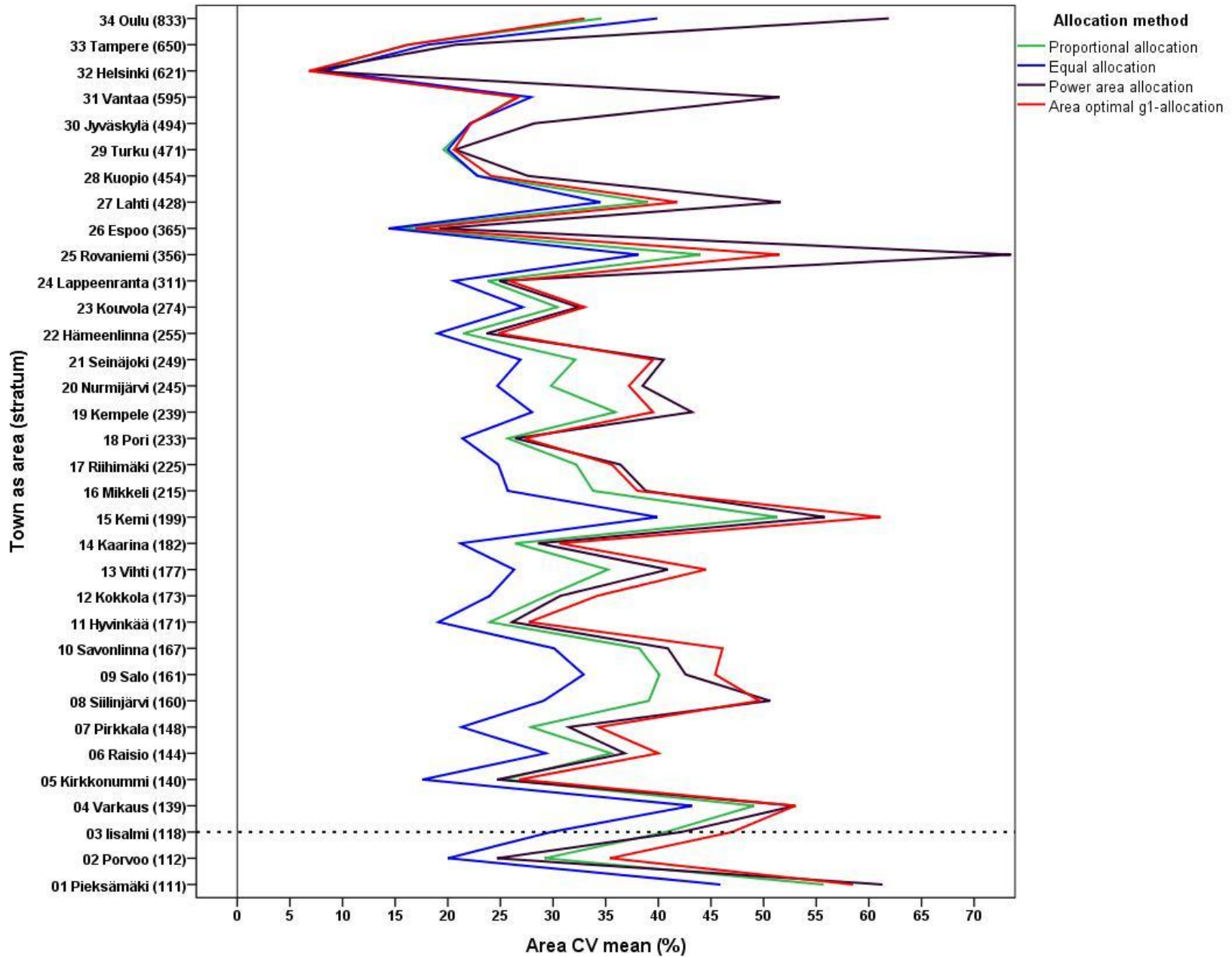
Some results computed from sample simulations (1500 samples / allocation):

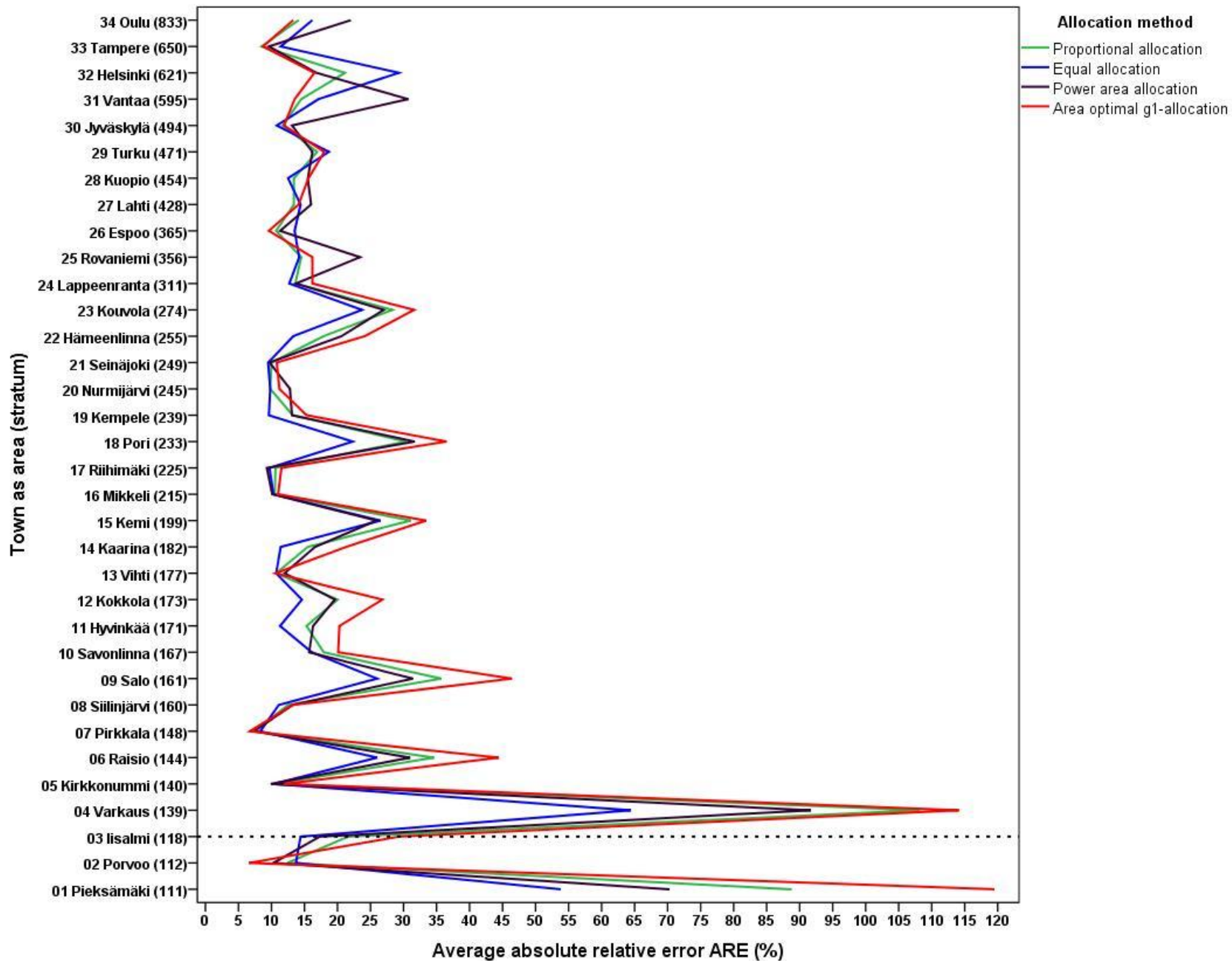
- dash line: limit of "zero" areas
- MSE means of areas
- CV means of areas
- ARE means of areas (average absolute relative error)

$$ARE \% = 100 \times (1/r) \sum_{k=1}^r \left| \hat{Y}_{dk,EBLUP} - Y_d \right| / Y_d ,$$

where r = number of sample simulations.

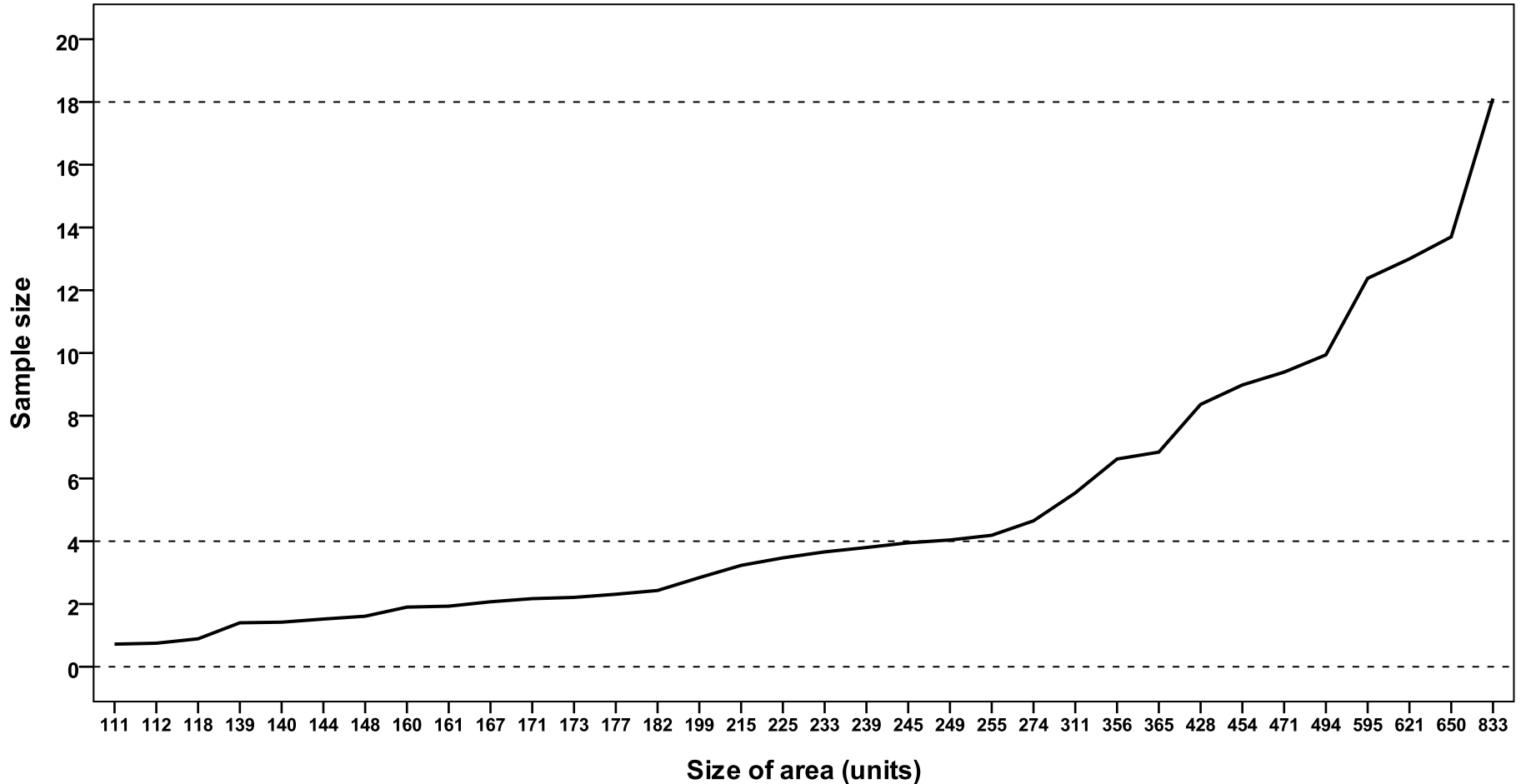






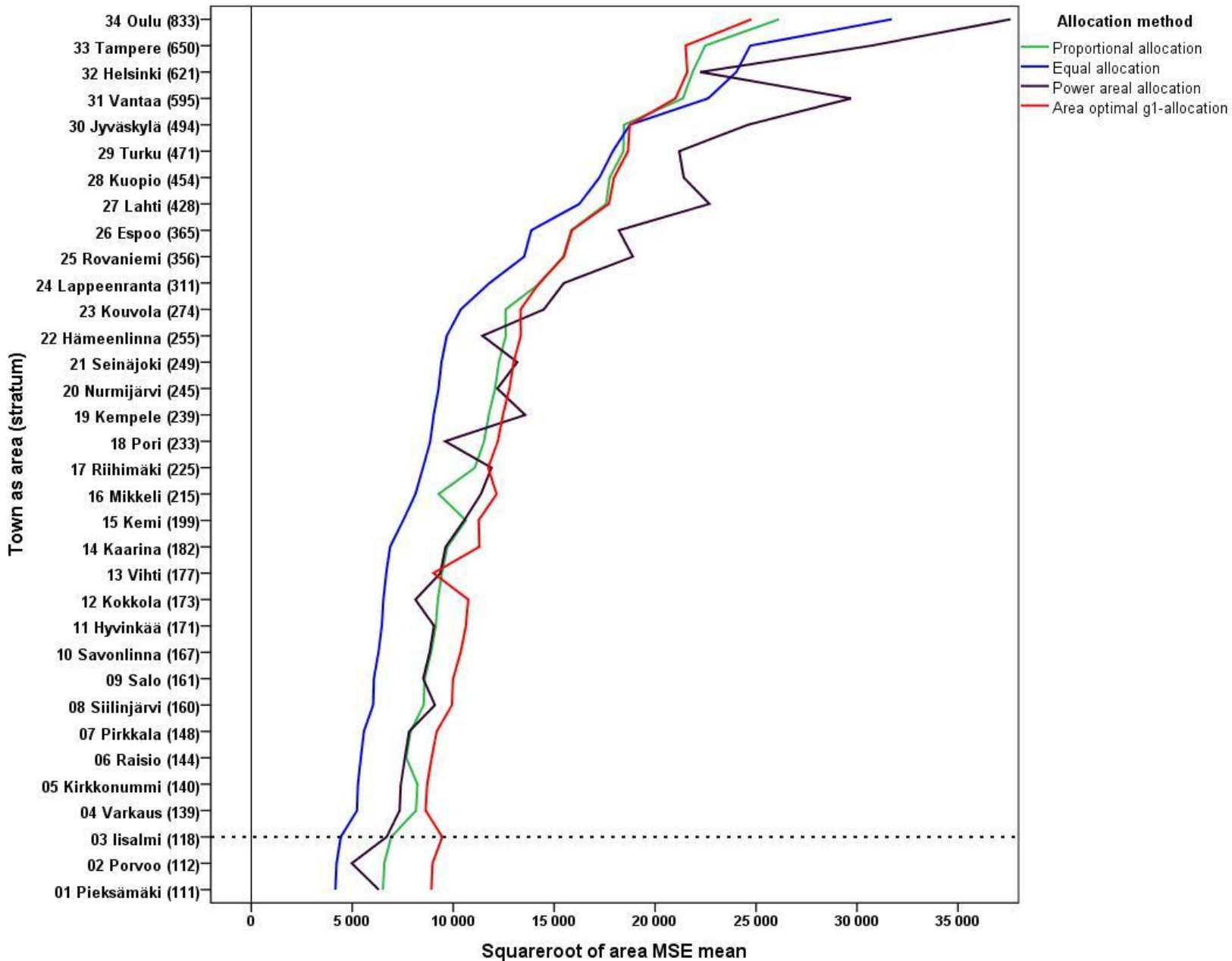
Phase 2: All 34 original areas, $N = 9815$, $n = 170$ ($E(n_d) = 5$)
- 3 smallest areas: $n_d = 0$

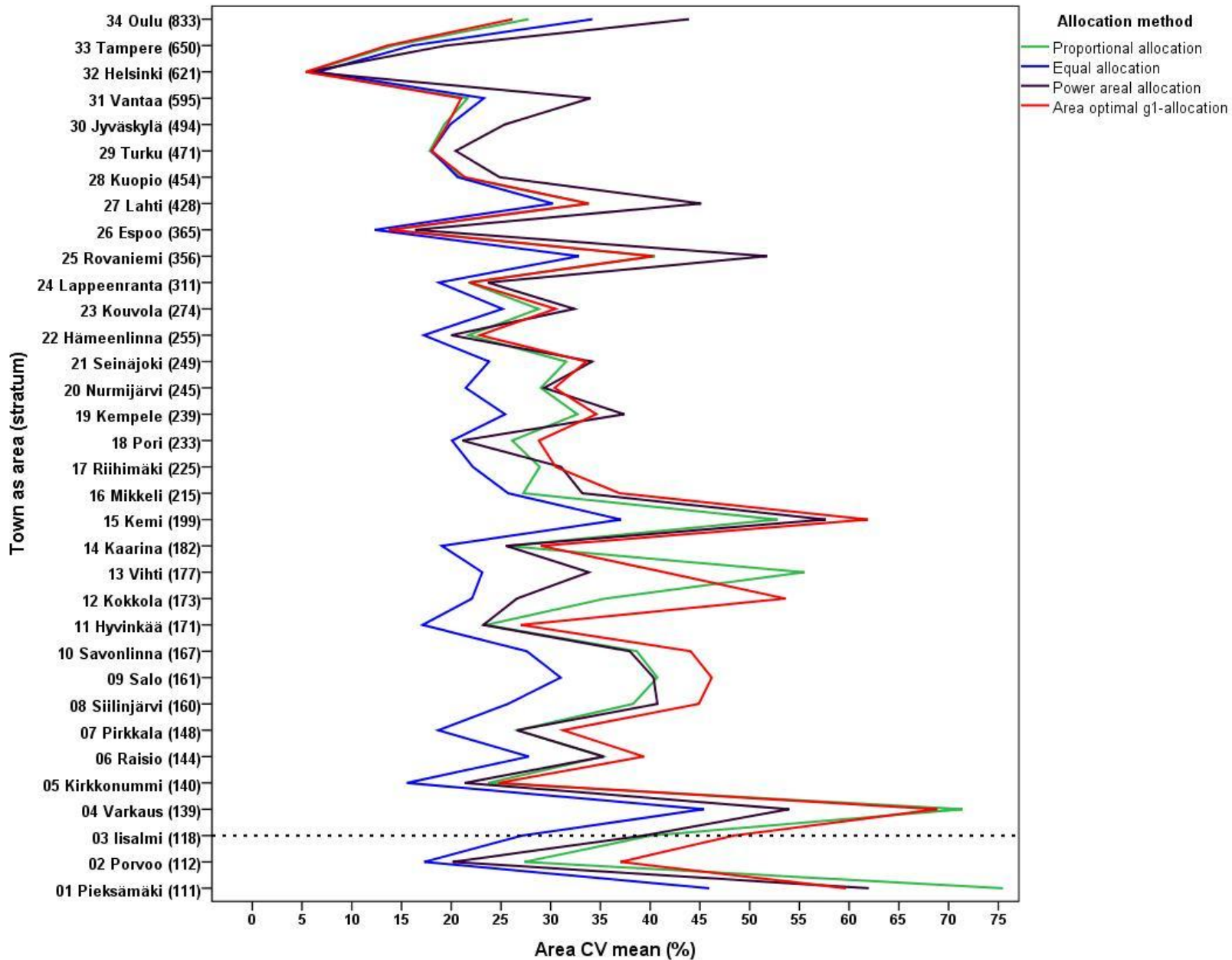
Computational sample sizes for 34 areas in g1-allocation

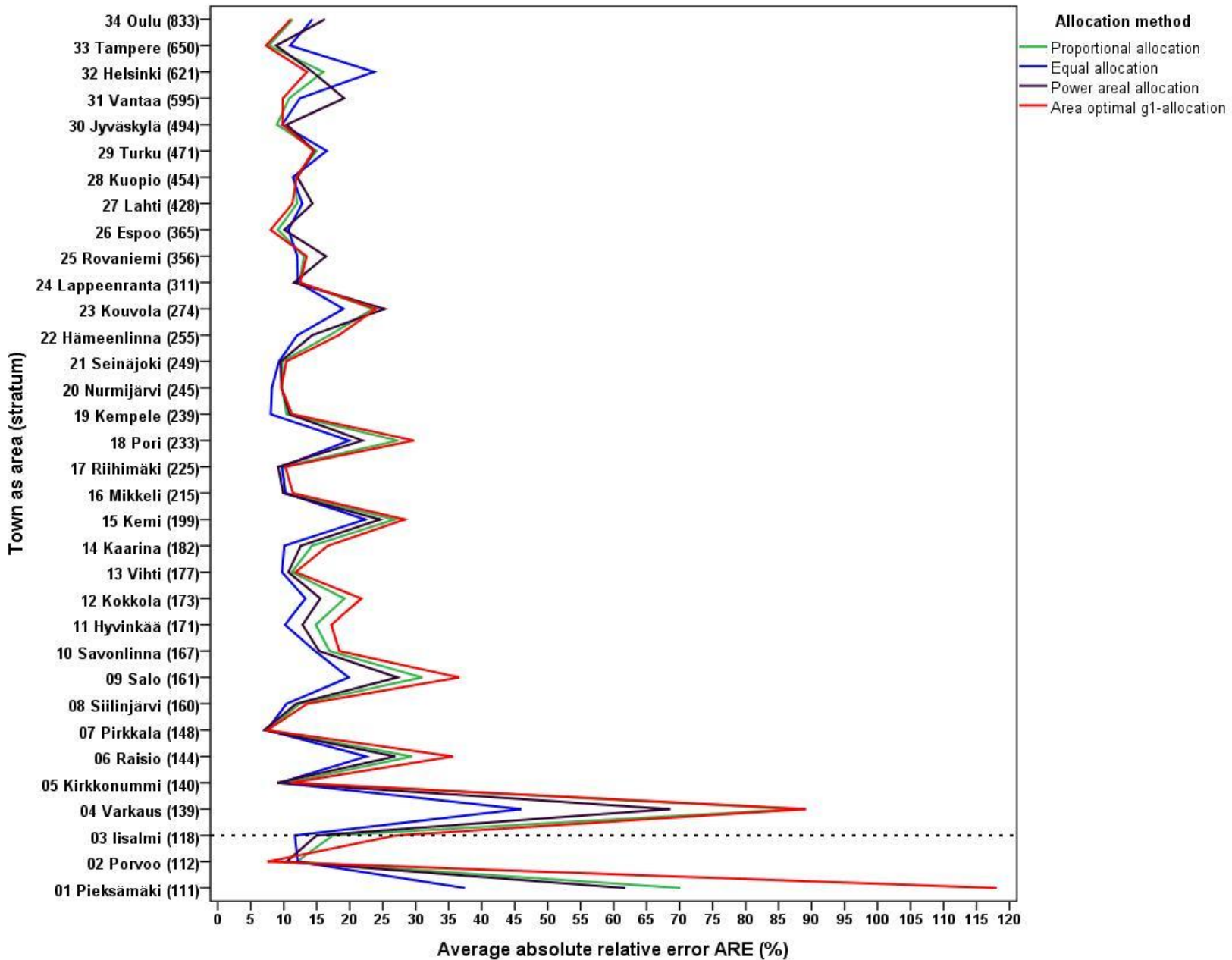


Results presented:

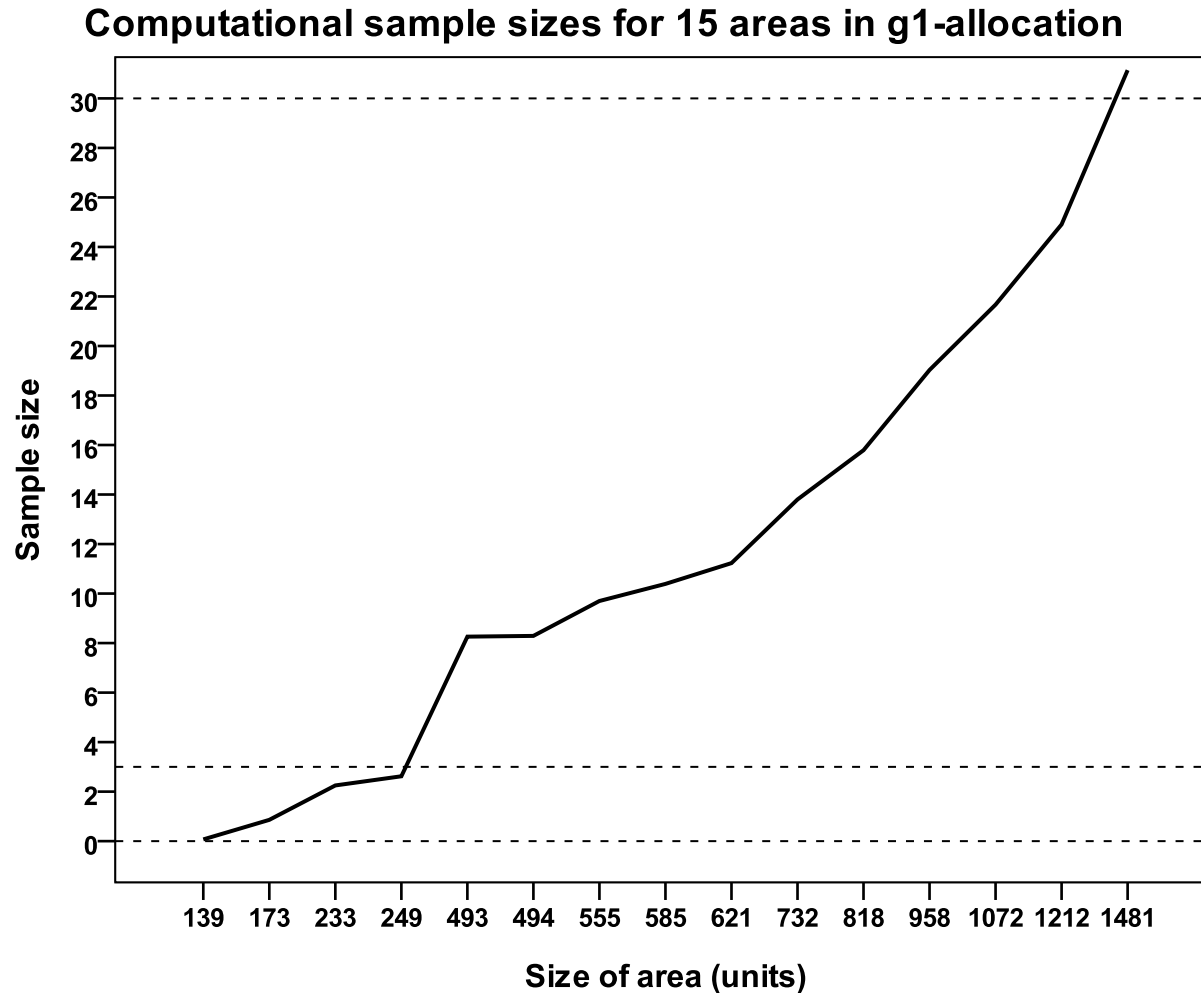
- dash line: limit of "zero" areas**
- MSE means of areas**
- CV means of areas**
- ARE means of areas (average absolute relative error)**





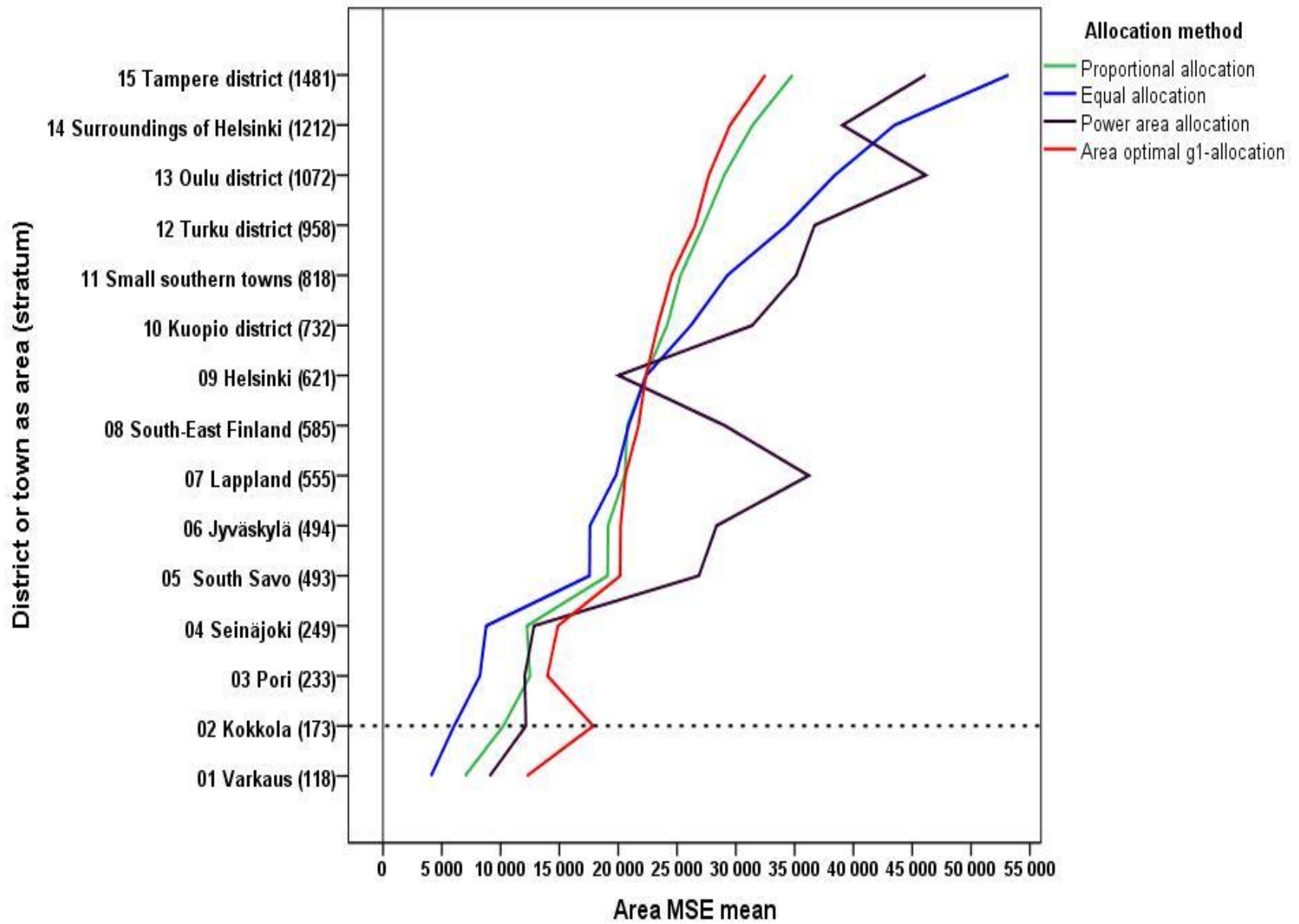


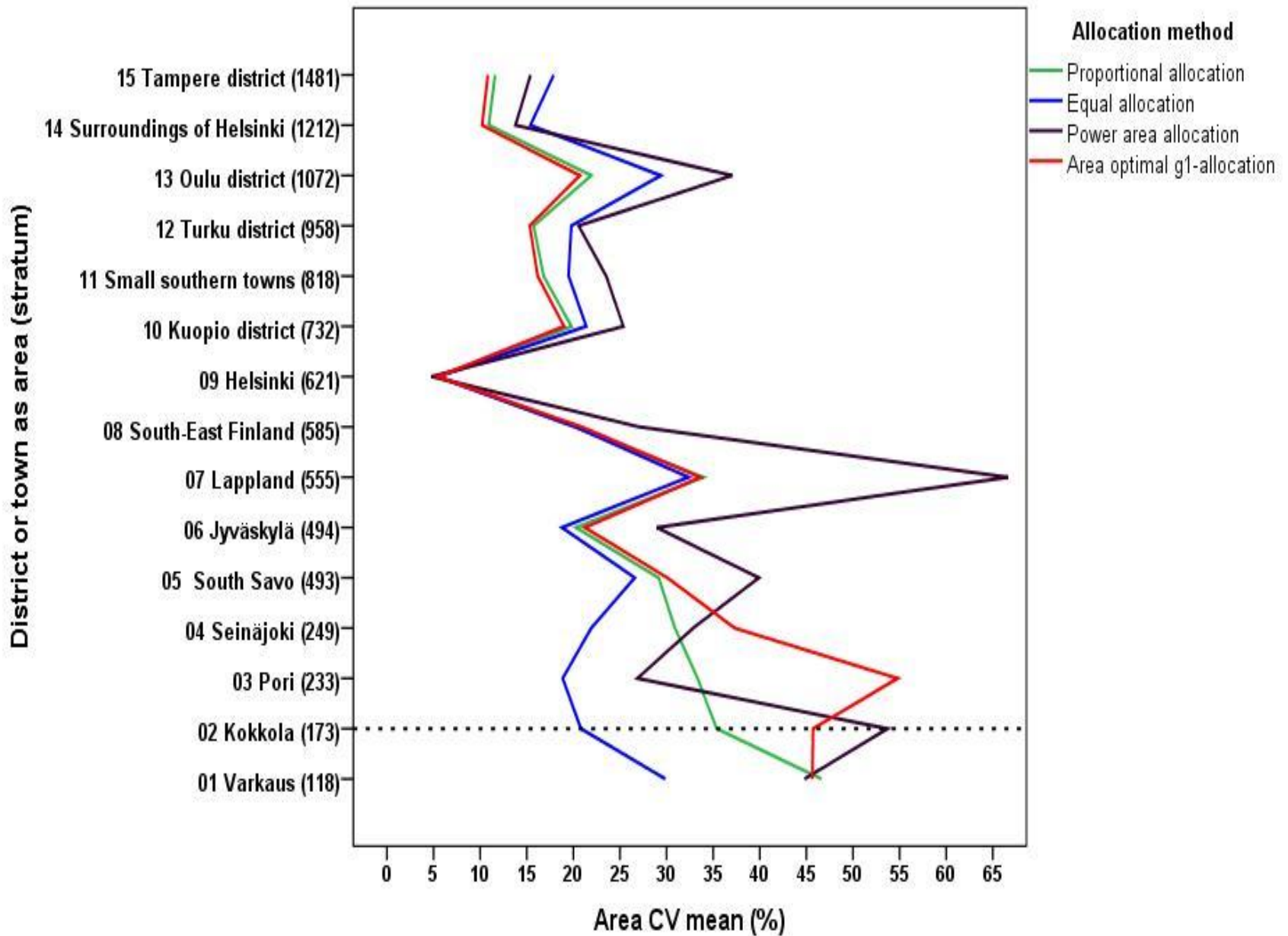
Phase 3: Combined 15 areas, $N = 9815$, $n = 180$ ($E(n_d) = 12$), homogeneity measure of $x = 0,237$, 2 smallest areas: $n_d = 0$

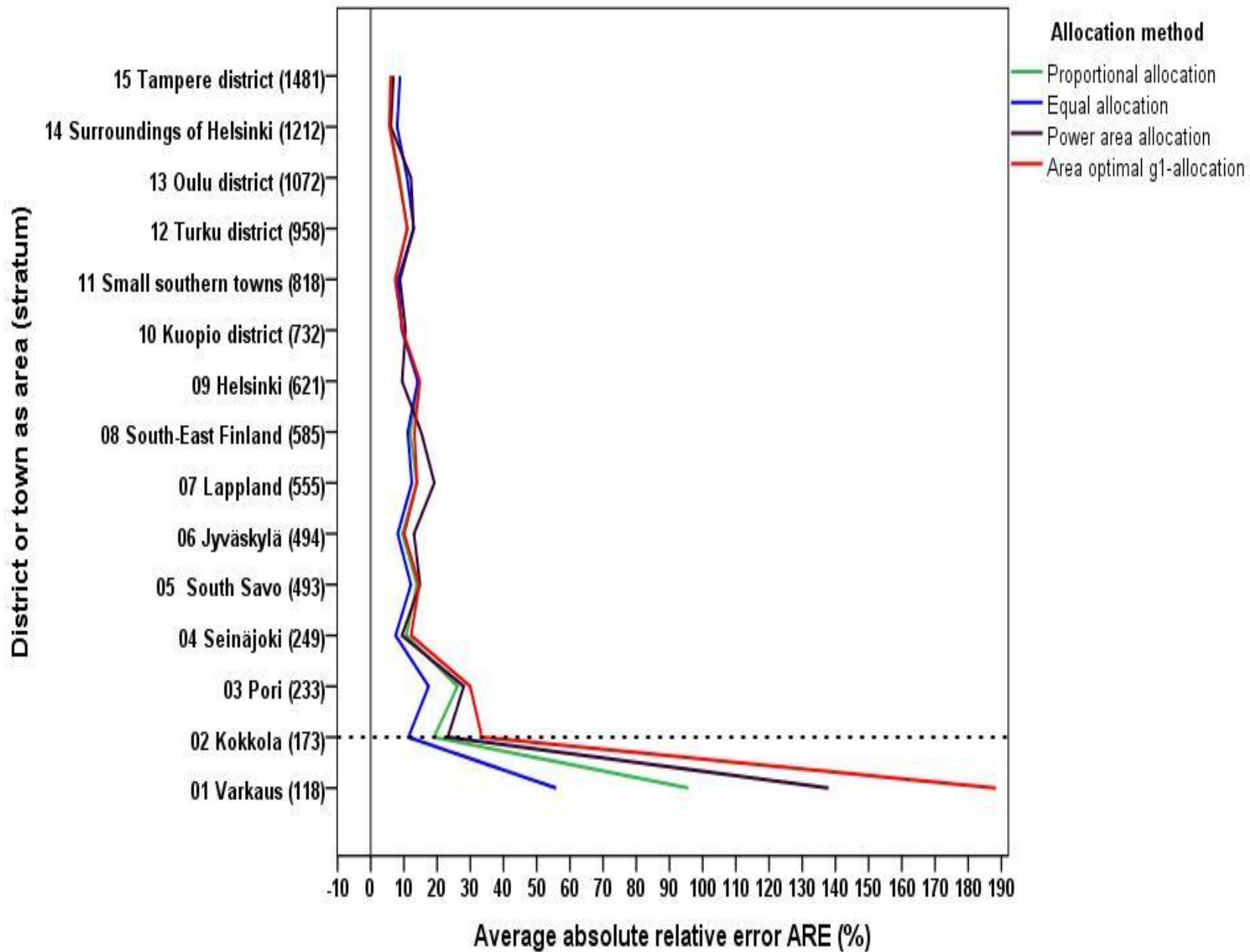


Results presented:

- dash line: limit of "zero" areas**
- MSE means of areas**
- CV means of areas**
- ARE means of areas (average absolute relative error)**







Results measuring the performance of each allocation by using MSE and CV means:

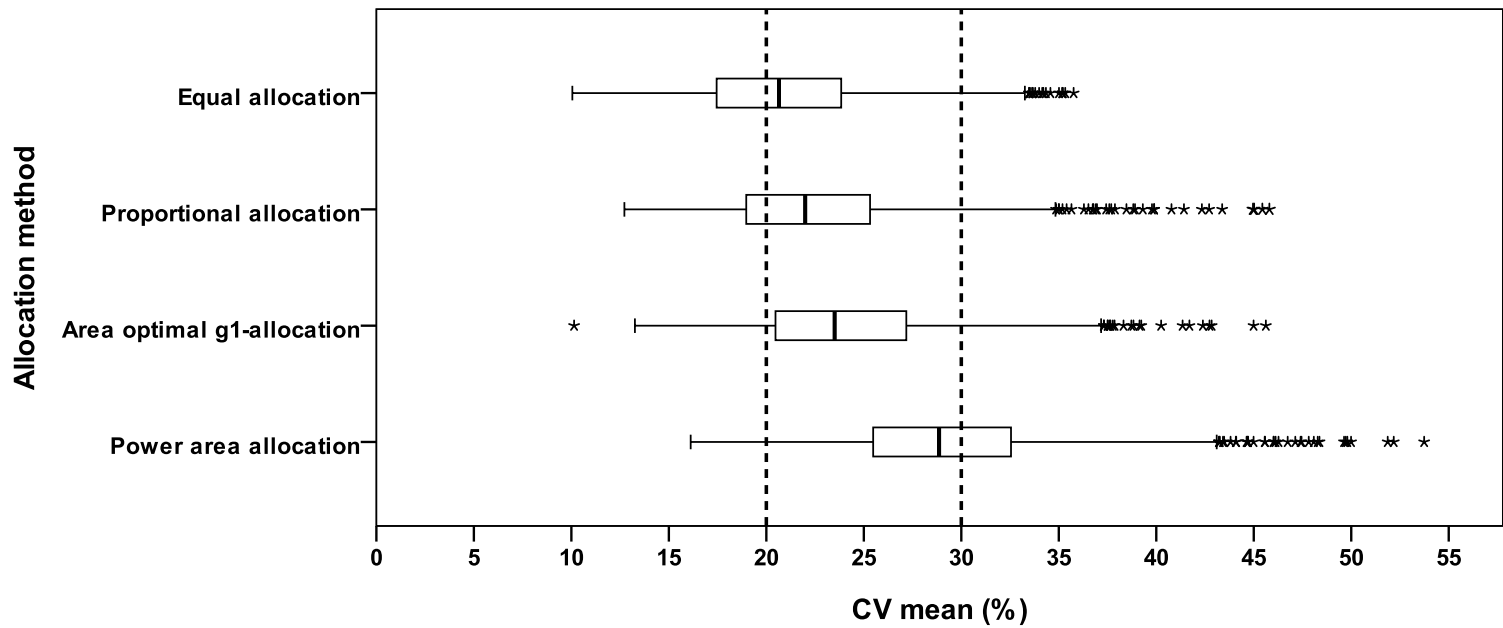
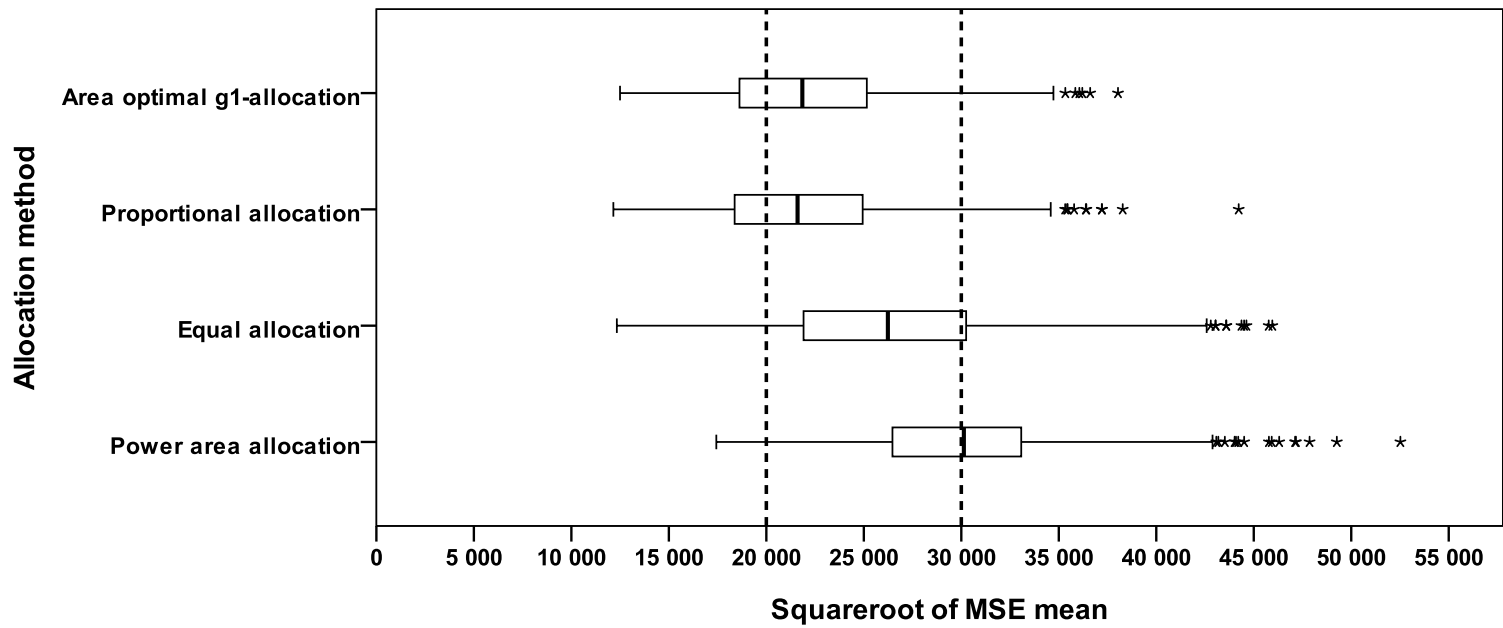
Distributions (as boxplot graphs) of

- MSE means of samples (100 %)

- CV means of samples (99 % of samples are presented because of a few very large values)

MSE mean in one sample = mean of area MSE's

CV mean in one sample = mean of area CV's.



Conclusions

Necessary condition for the use of g_1 -allocation in general is that variation between areas is strong enough. This can be confirmed through examination of auxiliary variable.

If the sizes of areas vary little (only a few large areas), equal allocation seems to be best among tested allocations in this situation, but g_1 -allocation has better performance compared with power and Neyman allocation and clearly better compared with latter. Proportional allocation has better performance than g_1 -allocation.

When the sizes of areas vary strongly and number of very small areas is low, use of g_1 -allocation is justified.

Factors affecting area MSE-mean: size of area, area mean, range and CV of x .

Factors affecting areal CV-mean: area mean, range and CV of x .

Factors affecting accuracy (ARE, ARB, RRMSE) of areal estimate: area CV of x , zero sample size or low sample size.

Also a "zero"-area can have good estimation results if its x -characteristics are close to corresponding x -characteristics in the whole population.

Compared with other allocations, area estimation results of g_1 -allocation improve when size of area grows.

Summary

g_1 -allocation seems to be an allocation alternative worth considering. It seems to work well in certain situations. It is better than Neyman and power allocation, and it is slightly better than proportional allocation when between-area variation is strong enough and area size is large enough.

Low overall sample size can lead to zero sample sizes for smallest areas in g_1 -allocation. This research found out that in spite of this estimation results can be moderately good if the area has properties which are near population properties.

Very small areas should be united into larger areas with similar properties before sampling.

References

Falorsi, P.D. and Righi, P. (2008). A balanced sampling approach for multi-way stratification for small area estimation. *Survey Methodology* **34**, 223-234.

Keto, M. and Pahkinen, E. (2009). On sample allocation for effective EBLUP estimation of small area totals – “Experimental Allocation”. In: J. Wywiał and W. Gamrot (eds.). (2010). *Survey Sampling Methods in Economic and Social Research*. Katowice: Katowice University of Economics.

Khan, M.G.M., Maiti, T. and Ahsan, M.J. (2010). An Optimal Multivariate Stratified Sampling Design Using Auxiliary Information: An Integer Solution Using Goal Programming Approach. *Journal of Official Statistics* **26**, 695-708.

Longford, N. T. (2006). Sample Size Calculation for Small-Area Estimation. *Survey Methodology* **32**, 87 - 96.

Nissinen, K. (2009). *Small Area Estimation With Linear Mixed Models From Unit-Level Panel and Rotating Panel Data*. University of Jyväskylä, Department of Mathematics and Statistics, Report **117**. (Dissertation).

Liels paldies!
Kiitos paljon!
Tack så mycket!
Thank you very much!
Danke schön!