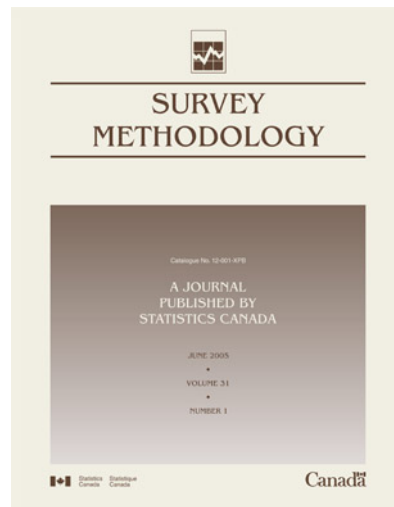




Catalogue no. 12-001-XIE

# Survey Methodology

June 2005



## How to obtain more information

Specific inquiries about this product and related statistics or services should be directed to: Business Survey Methods Division, Statistics Canada, Ottawa, Ontario, K1A 0T6 (telephone: 1 800 263-1136).

For information on the wide range of data available from Statistics Canada, you can contact us by calling one of our toll-free numbers. You can also contact us by e-mail or by visiting our website.

National inquiries line	1 800 263-1136
National telecommunications device for the hearing impaired	1 800 363-7629
Depository Services Program inquiries	1 800 700-1033
Fax line for Depository Services Program	1 800 889-9734
E-mail inquiries	<a href="mailto:infostats@statcan.ca">infostats@statcan.ca</a>
Website	<a href="http://www.statcan.ca">www.statcan.ca</a>

## Information to access the product

This product, catalogue no. 12-001-XIE, is available for free. To obtain a single issue, visit our website at [www.statcan.ca](http://www.statcan.ca) and select Our Products and Services.

## Standards of service to the public

Statistics Canada is committed to serving its clients in a prompt, reliable and courteous manner and in the official language of their choice. To this end, the Agency has developed standards of service that its employees observe in serving its clients. To obtain a copy of these service standards, please contact Statistics Canada toll free at 1 800 263-1136. The service standards are also published on [www.statcan.ca](http://www.statcan.ca) under About Statistics Canada > Providing services to Canadians.



Statistics Canada  
Business Survey Methods Division

# Survey Methodology

June 2005

Published by authority of the Minister responsible for Statistics Canada

© Minister of Industry, 2006

All rights reserved. The content of this electronic publication may be reproduced, in whole or in part, and by any means, without further permission from Statistics Canada, subject to the following conditions: that it be done solely for the purposes of private study, research, criticism, review or newspaper summary, and/or for non-commercial purposes; and that Statistics Canada be fully acknowledged as follows: Source (or "Adapted from", if appropriate): Statistics Canada, year of publication, name of product, catalogue number, volume and issue numbers, reference period and page(s). Otherwise, no part of this publication may be reproduced, stored in a retrieval system or transmitted in any form, by any means—electronic, mechanical or photocopy—or for any purposes without prior written permission of Licensing Services, Client Services Division, Statistics Canada, Ottawa, Ontario, Canada K1A 0T6.

May 2006

Catalogue no. 12-001-XIE  
ISSN 1492-0921

Frequency: semi-annual

Ottawa

Cette publication est disponible en français sur demande (n° 12-001-XIF au catalogue).

---

## Note of appreciation

Canada owes the success of its statistical system to a long-standing partnership between Statistics Canada, the citizens of Canada, its businesses, governments and other institutions. Accurate and timely statistical information could not be produced without their continued cooperation and goodwill.

# Approximations to $b^*$ in the Prediction of Design Effects Due to Clustering

Peter Lynn and Siegfried Gabler<sup>1</sup>

## Abstract

Kish's well-known expression for the design effect due to clustering is often used to inform sample design, using an approximation such as  $\bar{b}$  in place of  $b$ . If the design involves either weighting or variation in cluster sample sizes, this can be a poor approximation. In this article we discuss the sensitivity of the approximation to departures from the implicit assumptions and propose an alternative approximation.

Key Words: Complex sample design; Intracluster correlation coefficient; Selection probabilities; Weighting.

## 1. Alternative Functions of Cluster Size

Kish (1965) used an expression for the design effect (variance inflation factor) due to sample clustering,  $\text{deff} = 1 + (b - 1)\rho$ , where  $b$  is the number of observations in each cluster (primary sampling unit) and  $\rho$  is the intracluster correlation coefficient. This expression is well-known, is taught on courses on sampling theory, and is used by survey practitioners in designing and evaluating samples.

The expression holds when there is no variation in cluster sample size and the design is equal-probability (self-weighting). We can express these two criteria formally:

$$b_c = b \quad \forall c \quad (1)$$

where  $c = 1, \dots, C$  denote the clusters, and

$$w_i = w \quad \forall i \quad (2)$$

where  $i = 1, \dots, I$  denote the weighting classes, with  $w_i$  the associated design weights.

However, most surveys involve departures from (1) and (2). In the general case, *i.e.*, removing restrictions (1) and (2), Gabler, Häder and Lahiri (1999) showed that under an appropriate model,  $\text{deff}_c = 1 + (b^* - 1)\rho$ , where

$$b^* = \frac{\sum_{c=1}^C \left( \sum_{i=1}^I w_i b_{ci} \right)^2}{\sum_{i=1}^I w_i^2 b_i} = \frac{\sum_{c=1}^C \left( \sum_{j=1}^{b_c} w_{cj} \right)^2}{\sum_{c=1}^C \sum_{j=1}^{b_c} w_{cj}^2} \quad (3)$$

and  $b_{ci}$  is the number of observations in weighting class  $i$  in cluster  $c$ ,  $b_i = \sum_{c=1}^C b_{ci}$  (we have changed the notation from that of Gabler *et al.* (1999), to provide consistency) and  $w_{cj}$  is the weight associated with the  $j^{\text{th}}$  observation in cluster  $c$ ,  $j = 1, \dots, b_c$ .

The quantity  $b^*$  can be calculated from survey micro-data, provided the design weight and cluster membership is known for each observation. However, at the sample design stage it is not clear how  $b^*$  can be predicted. Gabler *et al.*

(1999) interpreted Kish's  $b$  as a form of weighted average cluster size:

$$\begin{aligned} \bar{b}_w &= \frac{\sum_{c=1}^C b_c \left( \sum_{i=1}^I w_i^2 b_{ci} \right)}{\sum_{c=1}^C \sum_{i=1}^I w_i^2 b_{ci}} \\ &= \frac{\sum_{c=1}^C \left( b_c \sum_{j=1}^{b_c} w_{cj}^2 \right)}{\sum_{c=1}^C \sum_{j=1}^{b_c} w_{cj}^2} \quad (4) \end{aligned}$$

where  $b_c$  is the number of observations in cluster  $c$ ,  $b_c = \sum_{i=1}^I b_{ci}$ . However, (4) is no easier than (3) to predict at the sample design stage. A simpler interpretation, perhaps commonly used in sample design, is the unweighted mean cluster size:

$$\bar{b} = \frac{\sum_{c=1}^C b_c}{C} = m/C. \quad (5)$$

It is much easier to predict  $\bar{b}$  at the sample design stage than either  $\bar{b}_w$  or  $b^*$ , as it requires knowledge only of the total number of observations,  $m$ , and total number of clusters,  $C$ .

## 2. Relationship Between $b^*$ , $\bar{b}_w$ and $\bar{b}$ Under Alternative Assumptions

Let

$$\bar{w}_c = \frac{1}{b_c} \sum_{j=1}^{b_c} w_{cj} = \sum_{i=1}^I w_i \frac{b_{ci}}{b_c},$$

$$\text{Cov}(b_c, b_c \bar{w}_c^2) = \frac{1}{C} \sum_{c=1}^C b_c^2 \bar{w}_c^2 - \frac{m}{C^2} \sum_{c=1}^C b_c \bar{w}_c^2$$

and

1. Peter Lynn, Institute for Social and Economic Research, University of Essex, Wivenhoe Park, Colchester, Essex CO4 3SQ, United Kingdom. E-mail: p.lynn@essex.ac.uk; Siegfried Gabler, Zentrum für Umfragen, Methoden und Analysen (ZUMA), Postfach 12 21 55, 68072 Mannheim, Germany. E-mail: gabler@zuma-mannheim.de.

$$\begin{aligned} \text{Var}(w_{cj}) &= \frac{1}{b_c} \sum_{j=1}^{b_c} (w_{cj} - \bar{w}_c)^2 \\ &= \sum_{i=1}^I \frac{b_{ci}}{b_c} (w_i - \bar{w}_c)^2 \quad \forall c. \end{aligned}$$

Then

$$b^* = \frac{C \cdot \text{Cov}(b_c, b_c \bar{w}_c^2) + \bar{b} \sum_{c=1}^C b_c \bar{w}_c^2}{\sum_{c=1}^C b_c \cdot \text{Var}(w_{cj}) + \sum_{c=1}^C b_c \bar{w}_c^2}. \quad (6)$$

If (1) holds, then (6) becomes:

$$b^* = \bar{b} \left( \frac{\sum_{c=1}^C \bar{w}_c^2}{\sum_{c=1}^C \text{Var}(w_{cj}) + \sum_{c=1}^C \bar{w}_c^2} \right). \quad (7)$$

So, in that circumstance,  $b^* \leq \bar{b}$ . If, additionally, weights are equal within clusters, *viz*:

$$w_{cj} = w_c \quad \forall j \in c \quad (8)$$

then  $b^* = \bar{b}$ .

If (8) holds, but not (1), then

$$\begin{aligned} b^* &\geq \bar{b} \text{ if and only if } \text{Cov}(b_c, b_c \bar{w}_c^2) \geq 0 \\ \text{since } b^* - \bar{b} &= \frac{C \cdot \text{Cov}(b_c, b_c \bar{w}_c^2)}{\sum_{c=1}^C b_c \bar{w}_c^2}. \end{aligned}$$

The covariance would be negative only if small cluster sizes coincide with large average weights within the clusters and *vice versa*. In section 4 below, we observe that this did not occur in any country on round 1 of the European Social Survey. Furthermore, from (3) and (4), we have:

$$b^* = \bar{b}_w = \sum_{c=1}^C (w_c b_c)^2 / \sum_{c=1}^C w_c^2 b_c. \quad (9)$$

If we additionally impose the restriction (1), then we have the obvious result  $b^* = \bar{b}_w = \bar{b} = b_c \quad \forall c$ .

The result in (9) would apply to surveys where the only variation in selection probabilities was due to disproportionate sampling between domains that did not cross-cut clusters. A common example would involve disproportionate stratification by region, with PSUs consisting of geographical areas hierarchical to regions.

A practical relaxation of the restriction on the variation in weights is:

$$b_{ci} = b_c \left( \frac{b_i}{m} \right) \quad \forall i, c. \quad (10)$$

In other words, we allow variation in weights within clusters, but we constrain the weights to have the same relative frequency distribution in each cluster, *i.e.*, the means and the variances of the weights within clusters do not depend on the clusters.

Now, (3) simplifies as follows:

$$\begin{aligned} b^* &= \sum_{c=1}^C \left( \sum_{i=1}^I w_i b_c \frac{b_i}{m} \right)^2 / \sum_{i=1}^I w_i^2 b_i \\ &= \sum_{c=1}^C \left( b_c^2 \left( \sum_{i=1}^I w_i b_i \right)^2 \right) / m^2 \sum_{i=1}^I w_i^2 b_i \\ &= \frac{\left( \sum_{i=1}^I w_i b_i \right)^2 \sum_{c=1}^C b_c^2}{\sum_{i=1}^I w_i^2 b_i m^2}. \quad (11) \end{aligned}$$

Note that  $(\sum_{i=1}^I w_i b_i)^2 / \sum_{i=1}^I w_i^2 b_i = m / (1 + c_w^2)$ , where  $c_w^2$  is the squared coefficient of variation, across all observations, of the weights. Also,  $(\sum_{c=1}^C b_c^2) / m^2 = (1 + c_b^2) / C$ , where  $c_b^2$  is the squared coefficient of variation, across all clusters, of the cluster sample sizes. Thus, (11) becomes:

$$b^* = \frac{m}{(1 + c_w^2)} \frac{(1 + c_b^2)}{C} = \bar{b} \frac{(1 + c_b^2)}{(1 + c_w^2)} = \bar{b}, \text{ say.} \quad (12)$$

So,  $\bar{b}$  will underestimate  $b^*$  if  $c_b^2 > c_w^2$  and *vice versa*. In particular, if  $w_{cj} = w \quad \forall j, c$  and  $c_b^2 > 0$ , then  $b^* > \bar{b}$ . The greater the variation in  $b_c$ , the greater the extent to which  $\bar{b}$  will under-estimate  $b^*$ .

Assumption (10) will rarely hold exactly, but this result might be useful in situations where the distribution of weights is expected to be similar across clusters. An example might be address-based samples where one person is selected per address. If the distribution of the number of persons per address is approximately constant across PSUs (in the population), then the distribution of weights will vary across clusters in the sample only due to sampling variation and disproportionate nonresponse (the effect of this could, of course, be substantial if cluster sample sizes are small).

If no restriction is imposed on the variation in weights, but  $\text{Var}(w_{cj}) > 0$  for at least one  $c$ , then, from (6),

$$b^* \geq \bar{b} \text{ if and only if } \zeta = \frac{C^2 \text{Cov}(b_c, b_c \bar{w}_c^2)}{m \sum_{c=1}^C b_c \text{Var}(w_{cj})} \geq 1. \quad (13)$$

If (10) holds, then  $\zeta = c_b^2 / c_w^2$ .

### 3. Implications for Sample Design

Expression (12) suggests that  $b^*$  may be predicted by predicting the relative magnitudes of  $c_b^2$  and  $c_w^2$ . However, this result applies to a special situation, where

$$\begin{aligned} \text{Cov}(w_{cj}, b_c) &= \frac{1}{m} \sum_{c=1}^C \sum_{j=1}^{b_c} (w_{cj} - \bar{w})(b_c - \bar{b}) \\ &= \frac{1}{m} \sum_{c=1}^C (b_c - \bar{b}) \left( \sum_{i=1}^I w_i b_{ci} - b_c \bar{w} \right) \\ &\stackrel{\text{from (10)}}{=} \frac{1}{m^2} \sum_{c=1}^C (b_c - \bar{b}) b_c \left( \sum_{i=1}^I w_i b_i - m \bar{w} \right) \\ &= 0 \end{aligned}$$

where

$$\begin{aligned} \bar{w} &= \frac{1}{m} \sum_{c=1}^C \sum_{j=1}^{b_c} w_{cj} = \frac{1}{m} \sum_{c=1}^C b_c \bar{w}_c \\ \bar{b} &= \frac{1}{m} \sum_{c=1}^C \sum_{j=1}^{b_c} b_c = \frac{1}{m} \sum_{c=1}^C b_c^2 = \frac{m}{C} (1 + c_b^2). \end{aligned}$$

When this covariance is expected to be small, it may be appropriate to predict  $b^*$  thus:

$$\hat{b}^* = \hat{b} = \hat{b} \frac{(1 + \hat{c}_b^2)}{(1 + \hat{c}_w^2)}. \tag{14}$$

Both coefficients of variation can be estimated from knowledge of the proposed sample design. In the following section, we investigate sensitivity of predictions obtained in this way to assumption (10) using real data from different sample designs with  $\text{Cov}(w_{cj}, b_c) > 0$ .

#### 4. Example: European Social Survey

The European Social Survey (ESS) is a cross-national survey for which great efforts have been made to achieve approximate functional equivalence in sample design between participating nations (Lynn, Häder, Gabler and Laaksonen 2004). Nevertheless, there is considerable variety in the types of design used, primarily due to variation in the nature of available frames and in local objectives, such as a desire for sub-national analysis which may lead to disproportionate stratification by domain. We use here data from the first round of the ESS, for which fieldwork was carried out in 2002 – 2003. Of the 22 participating nations, 17 had a clustered sample design. Of these, two had not yet provided useable sample data at the time of writing. In Table 1 we present the sample values of  $b^*$ ,  $\bar{b}$ ,  $c_b^2$ ,  $c_w^2$ ,  $\bar{b}$ ,  $|\bar{b} - b^*|$ ,  $|\bar{b} - b^*|$ ,  $\text{Corr}(w_{cj}, b_c)$ , and  $\zeta$  for the remaining 15. Note that the United Kingdom and Poland both had a 2 – domain design with the sample clustered only in one domain, namely Great Britain (*i.e.*, excluding Northern Ireland) and less densely-populated areas (*i.e.*, all except the largest 42 towns) respectively. Figures presented in table 1 relate only to the clustered domain.

**Table 1**  
Sample Values of  $b^*$ ,  $\bar{b}$ ,  $c_b^2$ ,  $c_w^2$ ,  $\bar{b}$ ,  $|\bar{b} - b^*|$ ,  $|\bar{b} - b^*|$ ,  $\text{Corr}(w_{cj}, b_c)$ , and  $\zeta$ , for 15 Surveys

Country		$b^*$	$\bar{b}$	$c_b^2$	$c_w^2$	$\bar{b}$	$ \bar{b} - b^* $	$ \bar{b} - b^* $	$\text{Corr}(w_{cj}, b_c)$	$\zeta$
Austria	AT	6.49	7.08	0.08	0.25	6.15	0.34	0.58	0.0036	0.4549
Belgium	BE	6.56	5.79	0.13	0.00	6.56	0.00	0.77	.	.
Switzerland	CH	8.83	9.23	0.12	0.21	8.50	0.34	0.40	0.0223	0.7060
Czech Republic	CZ	2.94	2.70	0.24	0.25	2.68	0.26	0.24	0.0225	1.7350
Germany	DE	18.85	18.13	0.07	0.11	17.42	1.43	0.72	-0.2287	.
Spain	ES	4.96	5.04	0.17	0.22	4.80	0.15	0.08	-0.0767	0.8757
Great Britain	GB	11.11	12.27	0.08	0.22	10.90	0.21	1.16	0.0114	0.4198
Greece	GR	5.47	5.86	0.09	0.22	5.25	0.22	0.39	-0.0280	0.5207
Hungary	HU	8.68	8.18	0.06	0.00	8.68	0.00	0.50	.	.
Ireland	IE	12.09	11.18	0.13	0.04	12.05	0.05	0.91	0.0006	3.1054
Israël	IL	11.79	12.82	0.12	0.56	9.27	2.53	1.02	-0.1271	0.4401
Italy	IT	10.98	10.87	0.26	0.16	11.80	0.83	0.10	-0.5589	1.3018
Norway	NO	44.09	18.68	1.33	0.01	43.32	0.77	25.41	0.0807	.
Poland (rural)	PL	10.07	9.45	0.06	0.01	9.88	0.19	0.62	0.2923	.
Slovenia	SI	10.76	10.13	0.06	0.00	10.76	0.00	0.63	.	.

From (12), we would expect to observe  $\bar{b} > b^*$  when  $\hat{c}_w^2 > \hat{c}_b^2$ . A common sample design for which this inequality can be anticipated is one where, a) the selected cluster sample size is constant, so variation in  $b_c$  will be limited to that caused by differential non-response; and b) the samples are equal-probability samples of addresses, with subsequent random selection of one person per address, leading to variation in design weights reflecting the variation in household size. There are six nations with sample designs of this type (AT, CH, ES, GB, GR, IL). It is indeed the case that for all of these nations,  $\zeta < 1$  and  $\bar{b} > b^*$ . Furthermore, for 5 of these 6 nations (AT, CH, ES, GB, GR,  $h = 1, \dots, 5$ ) we might expect (10) to be a reasonable approximation as the only variation in weights is that due to selection within a household/address. For these, we might expect  $\hat{b}$  to perform better than  $\bar{b}$ . Indeed,  $|\bar{b} - b^*| < |\hat{b} - b^*|$  for 4 of the 5, and  $(\sum_{h=1}^5 |\bar{b} - b^*|) / \sum_{h=1}^5 |\hat{b} - b^*| = 0.48$ . The one nation where  $\hat{b}$  would not provide an improvement is Spain and this is to be expected as  $\bar{b}$  is small. Small cluster sample sizes leave them relatively more susceptible to the effects of nonresponse and also sampling variance, which will lead to violation of (10). In Israel, there was a further source of variation in design weights as there was disproportionate stratification by geographical areas. This too causes violation of (10), so we would not expect  $\hat{b}$  necessarily to provide an improvement on  $\bar{b}$  as a predictor of  $b^*$ .

Of the nations where  $c_b^2 < c_w^2$ , there is only one (CZ) for which  $\bar{b} < b^*$  and  $\zeta > 1$ . This is also the nation with the smallest value of  $\bar{b}$ . When cluster sample sizes are particularly small, deff will be small and the choice between estimators of  $b^*$  may be less important.

There are five nations where sample units were individuals selected with equal probabilities (within clusters) from population registers (BE, DE, HU, PL, SI). In this case (8) (and, therefore, (10)) holds strictly, so we have  $\bar{b} < b^*$ . For three of these nations (BE, HU, SI) the sample is equal-probability, so we observe  $\bar{b} = b^*$ . It is clear that  $\hat{b}$  is superior to  $\bar{b}$  for equal-probability samples. For Germany and Poland, there is some variation in design weights between clusters (but not within). This variation is modest in Poland, and  $|\bar{b} - b^*| < |\hat{b} - b^*|$ , but the same is not true in Germany, where the ex-East Germany was sampled at a considerably higher rate than the ex-West Germany.

The Norwegian sample design was the only one that resulted in considerable variation in cluster sample sizes at the selection stage. The dramatic impact of this on  $\bar{b} \approx b^*$  can clearly be seen. Again, this is a situation in which  $\hat{b}$  is likely to be preferable to  $\bar{b}$  as a predictor of  $b^*$ .

The designs in Ireland and Italy both involved selecting addresses from the electoral registers with probability

proportional to number of electors and then selecting one resident at random from each selected address. Such designs are not equal-probability, but are likely to result in considerably less variation in design weights than the address-based sample designs discussed earlier (Lynn and Pisati 2005). In both these cases,  $\hat{c}_w^2 < \hat{c}_b^2$ , the difference being greater in the case of Italy where some cluster sample sizes (in the largest municipalities) were considerably larger than the others (in Ireland, all were equal at the selection stage). Aside from the Czech Republic, these are the only two nations with  $\zeta > 1$ .

## 5. Conclusion

To aid prediction of the design effect due to clustering, we believe that  $\hat{b}$  is likely to be a better choice than  $\bar{b}$  as a predictor of  $b^*$  in situations where it can reasonably be expected that (10) will approximately hold. This includes, but is not restricted to, the following common types of sample design:

- Equal-probability designs where cluster sample sizes vary by design;
- Equal-probability designs where clusters do not vary by design but are likely to vary due to nonresponse;
- Address-based samples where one person is selected at each address, there is no other significant source of variation in selection probabilities, and cluster sizes do not vary by design.

## Acknowledgement

This research was carried out while the first author was Guest Professor at the Center for Survey Research and Methodology (ZUMA), Mannheim, Germany.

## References

- Gabler, S., Häder, S. and Lahiri, P. (1999). A model based justification of Kish's formula for design effects for weighting and clustering. *Survey Methodology*, 25, 105-106.
- Kish, L. (1965). *Survey Sampling*. New York: John Wiley & Sons, Inc.
- Lynn, P., Häder, S., Gabler, S. and Laaksonen, S. (2004). Methods for Achieving Equivalence of Samples in Cross-National Surveys. ISER Working Paper 2004-09. Available at <http://www.iser.essex.ac.uk/pubs/workpaps/pdf/2004-09.pdf>.
- Lynn, P., and Pisati, M. (2005). Improving the quality of sample design for social surveys in Italy: Lessons from the European Social Survey. Forthcoming.