# Evaluation of the precision of systematic sampling: nugget effect and covariogram modelling

J. THIOULOUSE,* J. P. ROYET,† H. PLOYE* & F. HOULLIER‡

*\* Laboratoire de Biométrie, Génétique et Biologie des Populations, (U.R.A.CNRS 243), Université Claude Bernard-Lyon 1, F-69622 Villeurbanne Cedex, France*

*† Laboratoire de Physiologie Neurosensorielle, (U.R.A. CNRS 180), Université Claude Bernard-Lyon 1, F-69622 Villeurbanne Cedex, France*

*‡ Ecole Nationale du Génie Rural, des Eaux et Forêts, Laboratoire de Recherches en Sciences forestières, 14, rue Girardet, F-54042 Nancy Cedex, France*

## Summary

Systematic sampling designs are widely used in stereology. When an estimator of the total amount, $Q$, of the sampled variable is evaluated by such a procedure, the coefficient of error can be predicted by applying Matheron's theory of regionalized variables. To evaluate the accuracy of the estimate of $Q$, it is necessary to study the behaviour of the regionalized variable and to model its covariogram. Histological data with a low short-range variability and agronomic data with a pronounced nugget effect provided the biological material for extreme case studies. Results show that the short-range variability, if present, cannot be detected when only small samples are available. An underestimation of the coefficient of error is then to be expected. We propose several models of the covariogram, which can be used to test for the presence of a nugget effect. If a nugget effect is present, these models will provide better estimates of the coefficient of error. If there is no nugget effect a simplified method can be used and will provide reliable estimates of the coefficient of error.

## 1. Introduction

Systematic sampling is frequently used in biology and materials science (Matérn, 1960; Weibel, 1979, 1980). It has been found to be much more efficient than simple random sampling (Cochran, 1953; Ebbesson & Tang, 1967; Gundersen & Jensen, 1987; Mattfeldt, 1989). Several recent studies confirm the interest in the problem of estimating the precision of systematic sampling (Cruz-Orive, 1989; Kellerer, 1989; Matérn, 1989; Mattfeldt, 1989). The transitive theory of regionalized variables provides a consistent framework for solving this problem (Matheron,

1965, 1970; Journel & Huijbregts, 1978). This theory was originally developed in the context of geostatistics (Matheron, 1965, 1970), but has also been applied to forestry (Bachacou & Decourt, 1976; Bachacou et al., 1978), agronomy (Thioulouse et al., 1985a) and stereology (Digabel, 1971; Thioulouse et al., 1985b; Cruz-Orive, 1987, 1989; Michel & Cruz-Orive, 1988; Gundersen & Jensen, 1987; Mattfeldt, 1989).

A crucial aspect of this theory consists in analysing and modelling the structure of the regionalized variable, and representing it globally by its covariogram. To predict accurately the precision of the systematic sampling estimators, the model of the covariogram must account for the long-range variability (i.e. the structure of the regionalized variable for distances greater than the interval between adjacent sampling points) and the short-range variability (i.e. for distances smaller than the interval between adjacent sampling points) of the phenomenon under study. In some cases, the short-range variability is pronounced (due either to the measurement process or to the high irregularity of the biological object), and the covariogram is consequently not continuous at the origin. This discontinuity was called the 'nugget effect' by Matheron (1965).

The scope of this paper is to show (i) the key role of a good description of short-range variability in the covariogram model and (ii) the influence of the nugget effect on the estimation of precision of the systematic sampling. As a consequence, it aims at pointing out the usefulness of large reference samples (i.e. systematic samples with a small distance between adjacent points) prior to operational small samples in order to model correctly the covariogram, detect potential nugget effects and compute precisely the coefficient of error (CE) for different sampling intensities.

Thioulouse, J., J. P. Royet, H. Ploye, and F. Houllier. 1993. Evaluation of the precision of systematic sampling: nugget effect and covariogram modelling. Journal of Microscopy 172:249-256.

Two different biological materials were used: (i) histological data with a low short-range variability and (ii) agronomic data exhibiting a pronounced short-range variability (Thioulouse *et al.*, 1985a). The quantities to be estimated were, in the first case, the volume of the olfactory bulb (obtained from areal measurements on serial sections) and, in the second case, the total number of insects along a row of tomato plants. Agronomic data were included in this article because they provide a biological example in which the covariogram model must include a nugget effect. In both cases, we compared the predictions of the coefficients of error computed from the covariograms modelled using either a large reference sample or a small systematic sample.

## 2. Covariogram models and example data sets

Let $f(x)$ be the value of a one-dimensional regionalized variable defined on a domain $\mathscr{D}$ and measured at a point $x$, and $Q$ be the 'total amount' of this variable:

$$Q = \int_{\mathscr{D}} f(x)dx \qquad (1)$$

with

$$f(x) = 0 \text{ if } x \notin \mathscr{D}$$

Let a systematic sample be defined by $t$, the sampling interval (i.e. the distance between adjacent sampling points), $m$ the total number of sampling points and $x_1$ the abscissa of the first sampling point, uniformly randomly chosen between 0 and $t$. The abscissae $x_1, x_2, \ldots, x_m$ of the sampling points are defined by $x_{j+1} = x_j + t$ with $(j = 1, 2, \ldots, m-1)$. For a given systematic sample, an estimator of $Q$ is $[Q_t]$:

$$[Q_t] = t \sum_{j=1}^{m} f(x_j). \qquad (2)$$

As demonstrated by Matheron (1965), the estimation of the variance of $[Q_t]$ may use the transitive covariogram, $g(h)$:

$$g(h) = \int_{\mathscr{D}} f(x)f(x+h)dx. \qquad (3)$$

$g(h)$ can be calculated only for distances $h$ which are multiples of the sampling interval, and it is therefore necessary to model the covariogram. This model must take into account two different types of variability: (i) the short-range variability, i.e. the dependence between $f(x)$ and $f(x+h)$ when $h$ varies between 0 and $t$; and (ii) the long-range variability, which reflects the behaviour of the regionalized variable throughout the sample (i.e. when $h$ is greater than $t$), also called 'functional regionalization'. Two situations are distinguished below: (i) a large reference sample is available and may be used to infer the *CE* for subsequent smaller systematic samples, and (ii) only a small systematic sample is available. The *CE* is computed from the

variance, according to Eq. (4):

$$CE([Q_t]) = \frac{\sqrt{Var([Q_t])}}{[Q_t]}. \qquad (4)$$

### 2.1. Covariogram models for large reference samples

1 The first covariogram model ($L_1R$, where L stands for linear, and R for reference) is linear and takes into account only the first two points, $h = 0$ and $h = t$, which gives the slope of the tangent at the origin:

$$g_1(h) = C_{01} + C_{11}h \qquad (5)$$

with

$$C_{01} = g_1(0)$$

and

$$C_{11} = \frac{g_1(t) - g_1(0)}{t}.$$

The evaluation of the sampling variance of $[Q_t]$ is

$$Var([Q_t]) = -C_{11}\frac{t^2}{6}. \qquad (6)$$

2 The second covariogram model ($L_2R$) is also linear but with a discontinuity at the origin (i.e. a nugget effect), the linear trend being adjusted over the following points, excluding the first one ($h = 0$). The number of points used to adjust the linear trend must be chosen according to the covariogram shape.

$$g_2(h) = C_{02} + C_{12}h \qquad \text{if } h > 0 \qquad (7)$$

and

$$g_2(0) = C_{02} + C_{\Delta 2}.$$

The term $C_{\Delta 2}$ is used to model the nugget effect. The variance of $[Q_t]$ is estimated by

$$Var([Q_t]) = tC_{\Delta 2} - \frac{C_{12}}{6}t^2. \qquad (8)$$

3 The third covariogram model (QR) is quadratic with a nugget effect:

$$g_3(h) = C_{03} + C_{13}h + C_{23}h^2 \qquad \text{if } h > 0 \qquad (9)$$

and

$$g_3(0) = C_{03} + C_{\Delta 3}.$$

The quadratic trend is adjusted by ordinary least squares, using the first points but excluding $h = 0$. The number of points to be used can be chosen on an empirical basis, in order to improve the accuracy of the adjustment. The first point provides an estimate of $C_{\Delta 3}$. With this model, the variance of $[Q_t]$ is estimated by

$$Var([Q_t]) = tC_{\Delta 3} - \frac{C_{13}}{6}t^2. \qquad (10)$$

### 2.2. Covariogram models for small operational samples

Three models were used: model $Q_1S$, where Q stands for quadratic, and S for small, (Gundersen & Jensen, 1987),
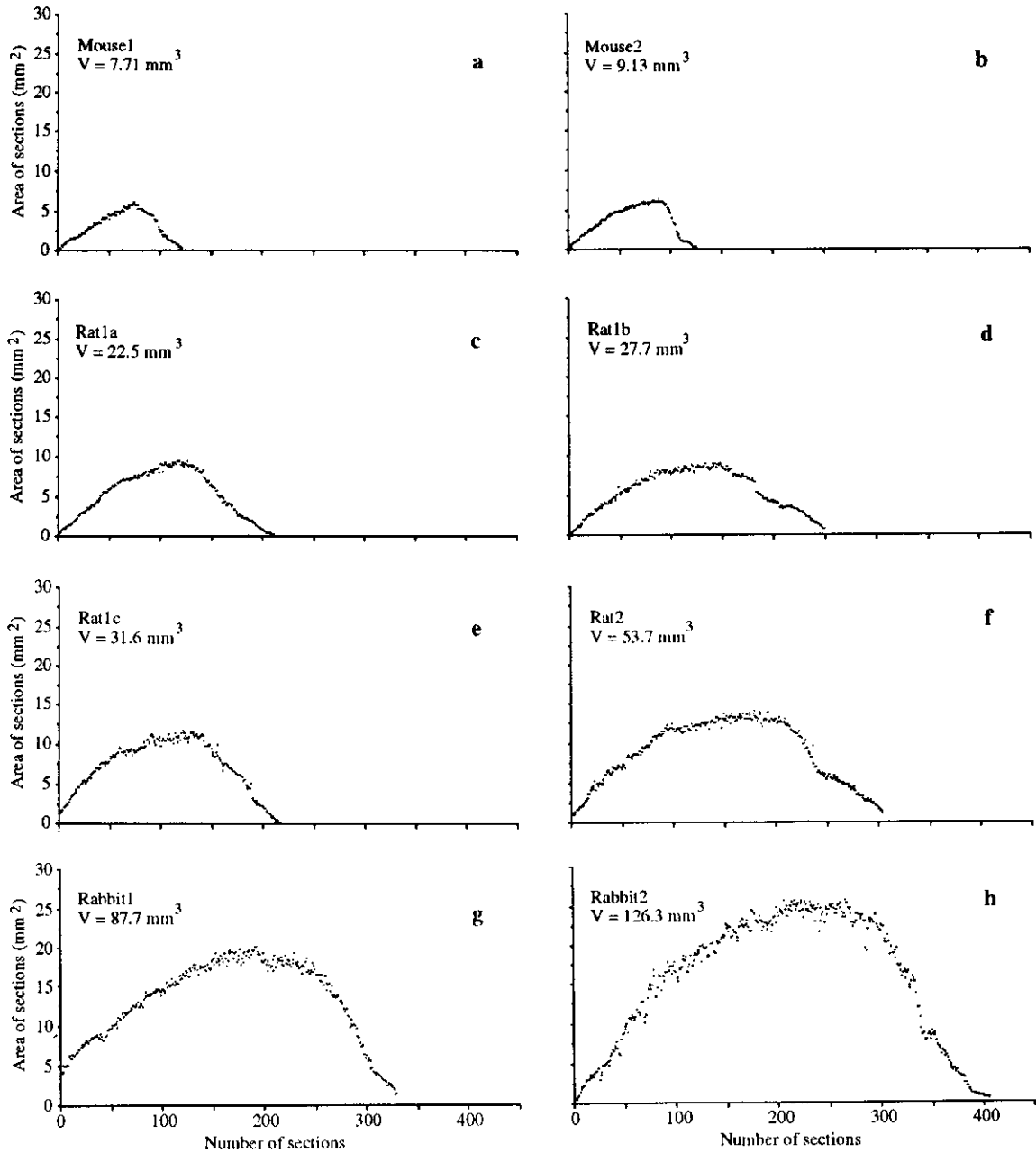
Fig. 1. Section areas from whole stacks of consecutive sections of 20 μm thickness through the olfactory bulb (OB) for different species of rodents (a: Mouse 1; b: Mouse 2; c: Rat 1a; d: Rat 1b; e: Rat 1c; f: Rat 2; g: Rabbit 1; h: Rabbit 2). V denotes the total volume of the OB.

model $Q_2S$ (Roberts *et al.*, 1993) and a linear model (LS) deduced from the first model ($L_1R$) described above. In this case, the CE can be computed from simplified formulas, in which the following variables are used:

$$A = [g(0)]/t = \sum_{j=1}^{m} f(x_j)^2 \qquad (11)$$

$$B = [g(t)]/t = \sum_{j=1}^{m-1} f(x_j)f(x_j + t) \qquad (12)$$

$$C = [g(2t)]/t = \sum_{j=1}^{m-2} f(x_j)f(x_j + 2t) \qquad (13)$$

$$D = [g(3t)]/t = \sum_{j=1}^{m-3} f(x_j)f(x_j + 3t) \qquad (14)$$

$$S = \sum_{j=1}^{m} f(x_j). \qquad (15)$$
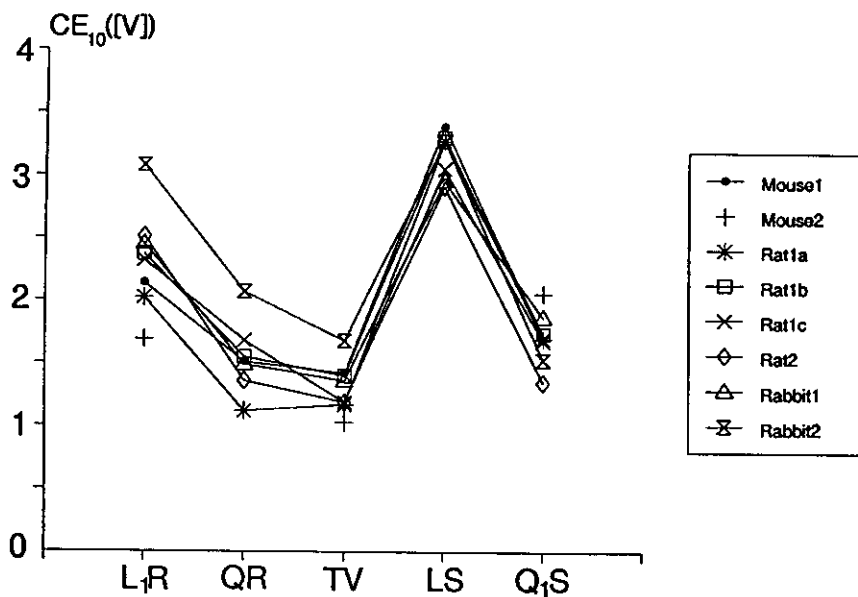
1 The LS covariogram model is linear, without nugget

Thioulouse, J., J. P. Royet, H. Ploye, and F. Houllier. 1993. Evaluation of the precision of systematic sampling: nugget effect and covariogram modelling. Journal of Microscopy 172:249-256.

Fig. 2. Coefficients of error computed for $L_1R$, QR, LS and $Q_1S$ models and for the true variance (TV) when the number ($n$) of analysed sections is equal to 10. In all cases, the estimate of CE decreases from $L_1R$ to QR. For QR, the CE is approximately equal to that of the TV. Estimates of the CE corresponding to the true variance are always smaller than estimates of CE obtained from the other models.

effect, and is adjusted with only the first two points of the covariogram: $[g(0)]$ and $[g(t)]$. The CE is

$$CE([Q_t]) = 1/S\sqrt{\frac{A-B}{6}}. \tag{16}$$

2 The $Q_1S$ model is based on a quadratic interpolation of the first three points of the covariogram. This model is similar to the QR model above, but without nugget effect. The CE is

$$CE([Q_t]) = 1/S\sqrt{\frac{3A-4B+C}{12}}. \tag{17}$$

3 The $Q_2S$ model is also based on a quadratic interpolation of the estimated covariogram but the three points used are the second to fourth, while the first one provides an estimate of the nugget effect. The CE is

$$CE([Q_t]) = 1/S\sqrt{A - \frac{31}{12}B + \frac{7}{3}C - \frac{3}{4}D}. \tag{18}$$

All these estimators of the CE are finally compared to the CE derived from the true variance (TV), obtained by taking all possible sets of $m$ equidistant points with $x_1$ varying between 0 and $t$.

### 2.3. Histological data set

Estimation of the olfactory bulb (OB) volume was carried out for males of two rodent and one lagomorph species: two C57B1/6J mice, four Wistar SPF rats and two rabbits. For the same species, subjects differed according to age or experimental conditions (Royet et al., 1989a,b). For each section, the area of the main OB was computed with the aid of a Quantimet 900 (Cambridge Instruments). The OB volume estimation from serial histological sections was

performed according to Cavalieri's principle (Gundersen & Jensen, 1987).

In order to study the influence of the nugget effect on the precision of OB volume estimates, a simulated noise (i.e. a random error) was added to observed values for mouse 1:

$$\tilde{f}(x) = f(x) + e(x)$$

$$e(x) = a(x)\sqrt{f(x)} \tag{19}$$

where $e(x)$ is the simulated noise, $a(x)$ a uniform random variable and $f(x)$ the measured area. Four simulations were carried out with different values for the variance of $a(x)$: 0·1, 0·5, 1, and 2.

### 2.4. Agronomic data set

Data come from a study on sampling procedures in population biology (Thioulouse et al., 1985a). The aim was to estimate the total number of adult whiteflies (Trialeurodes vaporariorum) in nine rows of tomato plants. One row with 380 plants and 778 insects was chosen. This row had been exhaustively measured.

### 3. Results

### 3.1. Histological data set

3.1.1. Original data. OB section areas are presented in Fig. 1 (each point represents a section). The variability around the general trend appears to be low. The corresponding covariograms (not shown here) are indeed very smooth, and the nugget effect is negligible.

Figure 2 shows the CEs computed for a given sampling intensity (10 sections) with the $L_1R$ and QR models when a

Thioulouse, J., J. P. Royet, H. Ploye, and F. Houllier. 1993. Evaluation of the precision of systematic sampling: nugget effect and covariogram modelling. Journal of Microscopy 172:249-256.
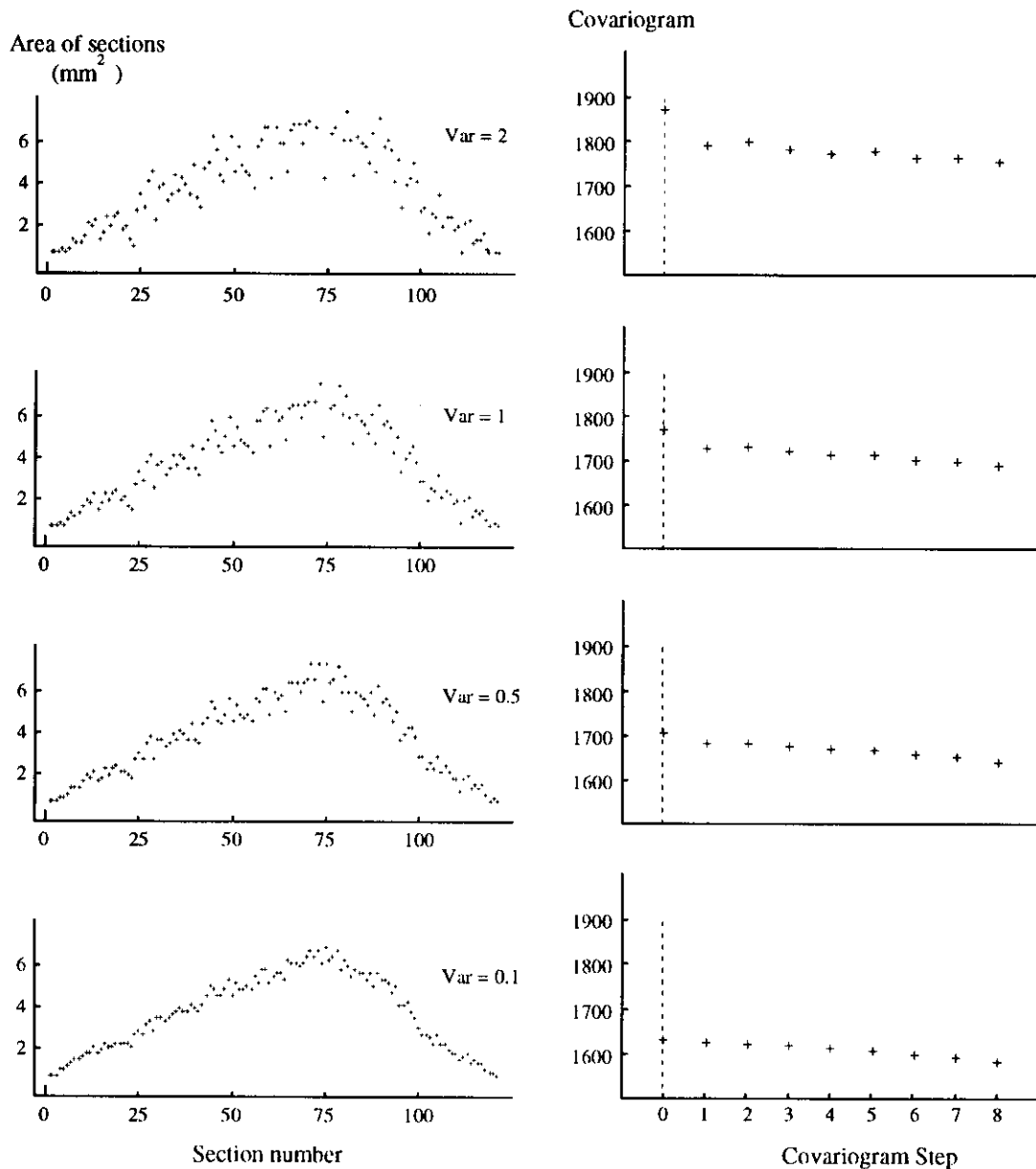
Fig. 3. Left: simulated data corresponding to section areas through the whole set of systematic sections of the olfactory bulb for Mouse 1. The variance of the simulated noise is equal to 0·1, 0·5, 1, and 2 (from bottom to top). Abscissa: number of sections analysed. Ordinate: area of sections. Right: the first 10 points of the corresponding covariograms. The importance of the nugget effect depends on the magnitude of the variance of the simulated noise.

reference sample is available, with the LS and $Q_1 S$ models when only measurements from a small sample are available and in the true variance case. The QR model is very well adjusted to the covariogram since it gives a CE similar to the one of the true variance method. Similarly, the $Q_1 S$ model predicts a CE very similar to that of the true variance. These results can be explained by the absence of short-range variability. In such conditions, the number of successive sections necessary to estimate the OB volume with a CE equal to 5% is small and approximately constant (from 2 to

7 sections only are necessary to estimate the OB volume). In contrast, it is clear that the sampling intensity is much higher when a sampling accuracy of 1% is required (from 12 to 31 sections).

*3.1.2. Simulated data.* Figure 3 (left part) illustrates the set of simulated data for Mouse 1 with different variances. The corresponding covariograms are illustrated on the right side of the figure. The augmentation of the nugget effect when the variance of the noise increases is very clear. Finally, the
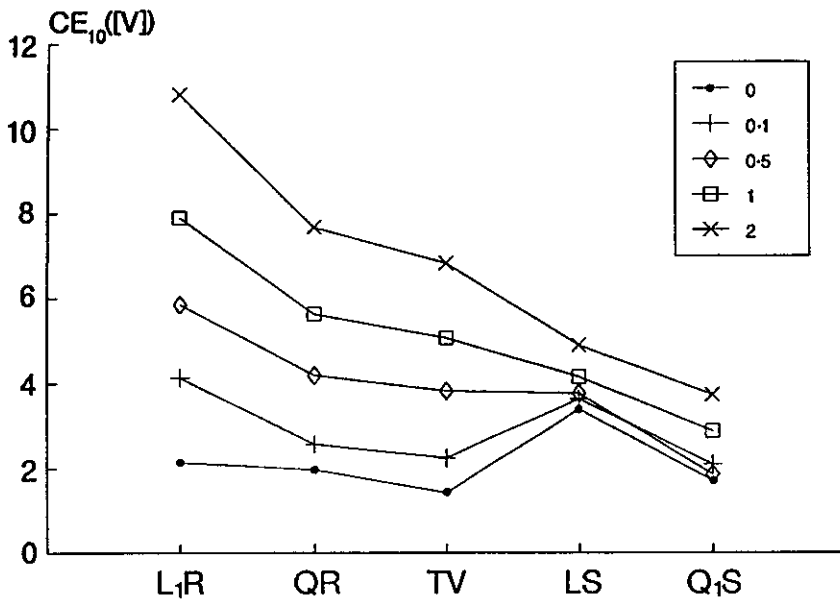
Thioulouse, J., J. P. Royet, H. Ploye, and F. Houllier. 1993. Evaluation of the precision of systematic sampling: nugget effect and covariogram modelling. Journal of Microscopy 172:249-256.

Fig. 4. Coefficients of error computed for Mouse 1 with $L_1R$, QR, LS and $Q_1S$ models and the true variance (TV), from four sets of simulated data with different variances (0·1, 0·5, 1 and 2). The variance 0 data correspond to those observed for Mouse 1, as illustrated in Fig. 1. These results show that model QR gives an estimation for the CE which is fairly close to that obtained from the true variance. In contrast, the $Q_1S$ model leads, in the four cases of simulation, to an underestimation of the CE. In other words, this model is not adequate to estimate the CE when areal measurements present a pronounced short-range variability. The LS model performs more poorly, except in the case where the variance of the simulated noise is equal to 0·1.

CEs computed for a sampling intensity of 10 sections with $L_1R$, QR, LS and $Q_1S$ models and the CE computed from the true variance are illustrated in Fig. 4. Analysis of this figure reveals that whichever model is considered, the simulated noise clearly influences the CE. However, the QR model gives results which are very similar to that of the true variance. The performances of $Q_1S$ and LS models are much poorer: for $a(x)$ equal to 1 and 2, the CE is underestimated. This results from the fact that these models do not take into account the nugget effect and only consider the first two or three points of the particular covariogram established at a low sampling intensity. Model $Q_2S$, which does include a nugget effect, was applied only to the simulated data with variance $a(x)$ equal to 2. Over the 12 possible subsamples consisting of 10 sections, only three (numbers 4, 7 and 10) provided positive estimations of the variance, while in the nine other cases, the estimation of $g(0)$ was actually higher than the computed $g(0)$, hence leading to a negative nugget effect. The three corresponding values for the CE are 11·3, 9·4, and 12·3%. These estimates are correct and in good agreement with the results of $L_1R$ and QR models and TV. However, the small number of cases in which model $Q_2S$ can be used without problem seems to preclude its use as a 'ready made' model.

### 3.2. Agronomic data set

The number of whiteflies on the 380 tomato plants is illustrated in Fig. 5(a). The corresponding covariogram is shown in Fig. 5(b). We observe a very strong nugget effect. Coefficients of error determined for $L_2R$, LS and $Q_1S$ models and for the true variance are depicted in Fig. 6. The $L_2R$ model gives results very close to the true variance. By contrast, the LS and $Q_1S$ models present a mean CE which is less than half the CE of the TV method. Therefore, it is

evident that LS and $Q_1S$ methods are over-optimistic when the short-range variability is high and that they lead to a sampling intensity which is too low.

Model $Q_2S$ was also applied to this data set. Over the 38 possible subsamples of size 10, only 18 provided positive estimates of the variance. But even in these 18 cases, the variability of the subsamples was so high that the model was not usable, for example when $g(1)$ was lower than $g(2)$ [and sometimes lower than $g(3)$], hence leading to aberrant (possibly negative) estimates of $g(0)$. One solution in this case is to use a model adjusted over a higher number of points (50 points were used here for model $L_2R$) in order to improve the reliability of the estimation of $g(0)$. But this solution is valid only if the number of points in the subsample is sufficient: in our example, with only 10 points, the nugget effect can be randomly estimated as either negligible or very high.

### 4. Discussion

The theory of regionalized variables requires a model of the covariogram. If only operational small samples are available, there is no way to describe accurately the behaviour of the covariogram near the origin. In this case, estimates of the CE can be unreliable (e.g. LS and $Q_1S$ models for simulated histological data and agronomic material).

One way to avoid this difficulty is to have additional information on the structure of the phenomenon. Such information may come from empirical knowledge on similar situations (e.g. OB data from other species with the same histological technique). But the best way to get it is to make some preliminary large reference samples for a few cases. These samples enable us to analyse the structure and the

Thioulouse, J., J. P. Royet, H. Ploye, and F. Houllier. 1993. Evaluation of the precision of systematic sampling: nugget effect and covariogram modelling. Journal of Microscopy 172:249-256.
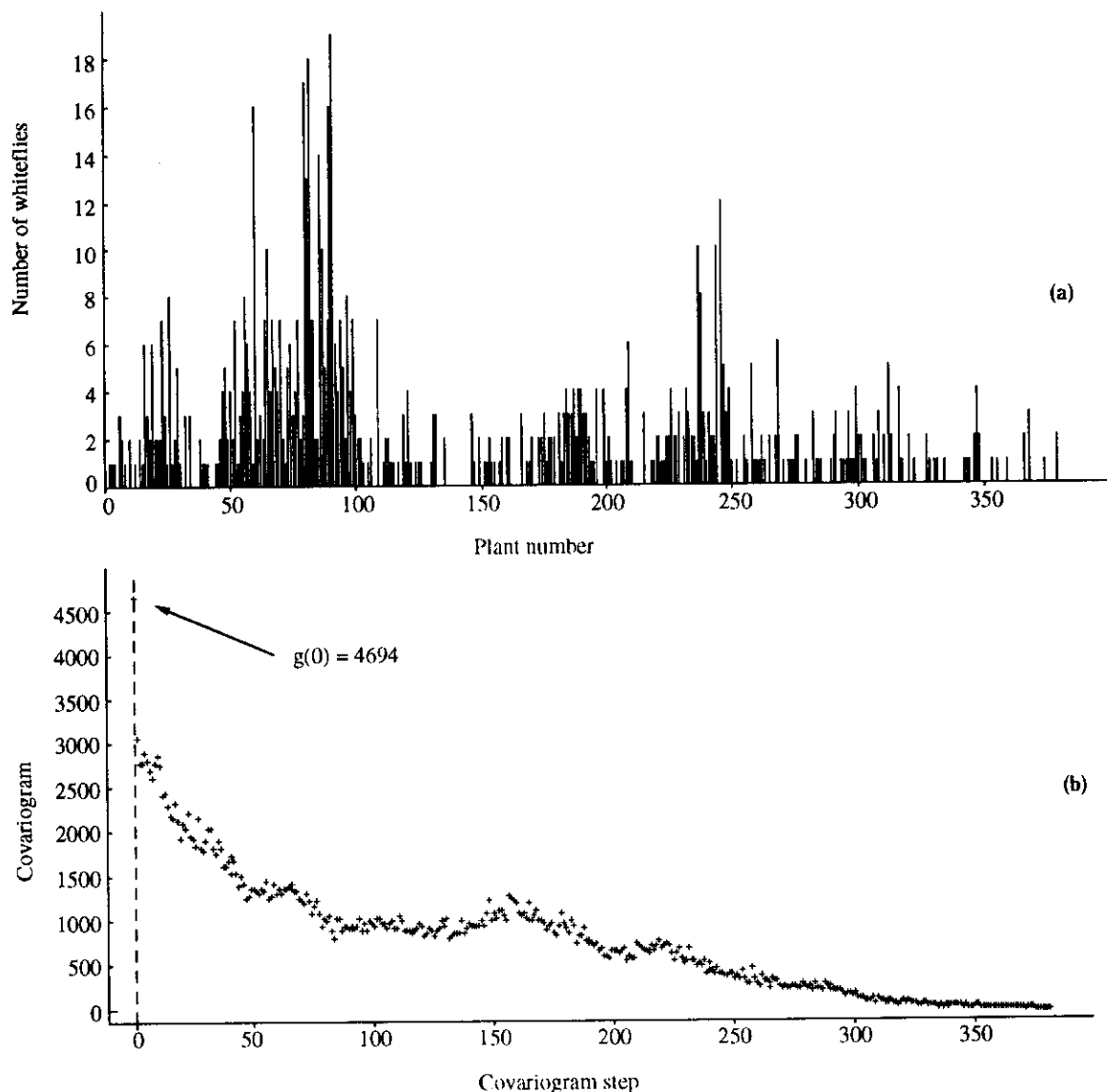
Fig. 5. (a) Number of adult whiteflies along the row of tomato plants. Abcissa: position of plants along the row (from 1 to 380). Ordinate: frequency of insects for each plant. (b) Corresponding covariogram. Abcissa: covariogram step. Ordinate: covariogram values. The linear model ($L_2R$) was fitted to points 2 to 51. The nugget effect is very high (Thioulouse et al., 1985a).

variability of the phenomenon and to propose reliable covariogram models. Once the absence of a strong nugget effect has been proven, it is possible to use the $Q_1S$ covariogram model and to use this model for subsequent small operational samples. But if it appears that the short-range variability is pronounced, a covariogram model with a nugget effect must be chosen (e.g. see the agronomic data presented here). By using this model, it is then possible to optimize the sampling intensity of subsequent small operational samples on a reliable basis.

Our study confirms that the simplified method (i.e. the quadratic model $Q_1S$ and the procedure to estimate it) proposed by Gundersen & Jensen (1987) is reliable for

histological data. If a high level of short-range variability is expected in the data set, one can, as a first step, try to apply model $Q_2S$. If the estimate of the variance is positive, then this value can be used, and will provide a better estimate of $CE([V])$ than the one computed with model $Q_1S$. If the estimate of the variance is negative, this means that the nugget effect is actually negligible, and one can therefore switch to model $Q_1S$.

## Acknowledgments

Thioulouse, J., J. P. Royet, H. Ploye, and F. Houllier. 1993. Evaluation of the precision of systematic sampling: nugget effect and covariogram modelling. Journal of Microscopy 172:249-256.
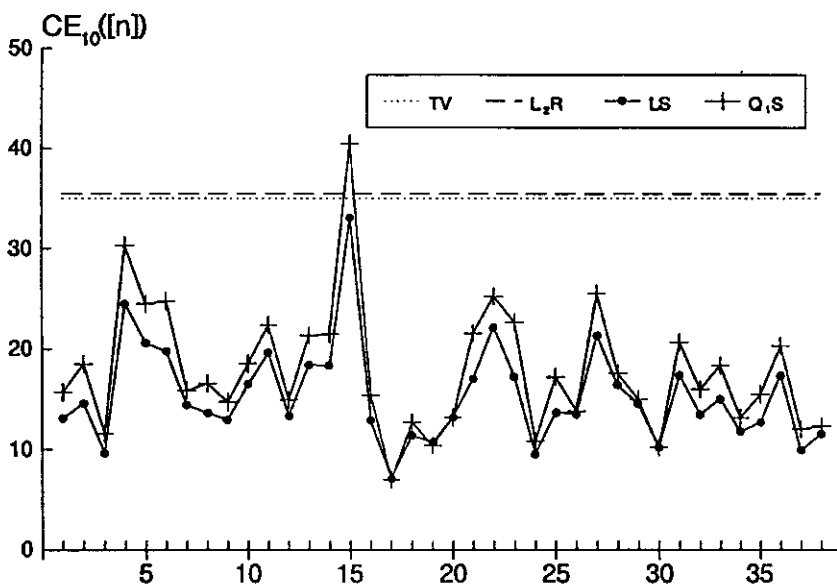
Fig. 6. Coefficients of error computed for whiteflies with LS and $Q_1S$ models for the 38 possible subsamples made of 10 plants. The great variance of the different estimates comes from the high spatial irregularity of the regionalized variable and underlines the instability of the estimated coefficients of error. Dotted and dashed horizontal lines represent, respectively, the true variance and the $L_2R$ model. These values are more than twice those of LS and $Q_1S$ models.

# References

Bachacou, J., Bouchon, J. & Marbeau, P. (1978) Etude structurale d'une régénération naturelle en forêt. Biométrie et Ecologie, Société Francaise de Biométrie, no. 1 (ed. by J.M. Legay and R. Tomassone), pp. 1–16. Société Française de Biométrie.

Bachacou, J. & Decourt, N. (1976) Etude de la compétition dans des plantations régulières à l'aide de variogrammes. Ann. Sci. Forest. 33, 177–198.

Cochran, W.G. (1953) Sampling Techniques. Wiley and Sons, New York.

Cruz-Orive, L.M. (1987) Precision of stereological estimators from systematic probes. Proc. Seventh Int. Congr. Stereol. (ed. by J.L. Chermant), Acta Stereol. 6/III, 153–158.

Cruz-Orive, L.M. (1989) On the precision of systematic sampling: a review of Matheron's transitive methods. J. Microsc. 153, 315–333.

Digabel, H. (1971) Estimation du volume ovogonial d'une gonade. Internal Report no. 259, Ecole Nationale Supérieure des Mines de Paris, F-Fontainebleau.

Ebbesson, S.O.E. & Tang, D.B. (1967) A comparison of sampling procedures in a structured cell population. Proc. Second Int. Congr. Stereol. (ed. by H. Elias), pp. 131–132. Springer, Berlin.

Gundersen, H.J.G. & Jensen, E.B. (1987) The efficiency of systematic sampling in stereology and its prediction. J. Microsc. 147, 229–263.

Journel, A.G. & Huijbregts, CH.J. (1978) Mining Geostatistics. Academic Press, New York.

Kellerer, A.M. (1989) Exact formulae for the precision of systematic sampling. J. Microsc. 153, 285–300.

Matérn, B. (1960) Spatial variation. Meddelanden From Statins, Skoqsforskninstitnt, 49(5).

Matérn, B. (1989) Precision of area estimation: a numerical study. J. Microsc. 153, 269–284.

Matheron, G. (1965) Les variables régionalisées et leur estimation. Masson et Cie, Paris.

Matheron, G. (1970) La théorie des variables régionalisées et ses applications. Les Cahiers du Centre de Morphologie Mathématique de Fontainebleau. No. 5, Ecole Nationale Supérieure des Mines de Paris, F-Fontainebleau.

Mattfeldt, T. (1989) The accuracy of one dimensional systematic sampling. J. Microsc. 153, 301–313.

Michel, R.P. & Cruz-Orive, L.M. (1988) Application of the Cavalieri principle and vertical sections method to lung: estimation of volume and pleural surface area. J. Microsc. 150, 117–136.

Roberts, N., Cruz-Orive, L.M., Reid, N., Brodie, D., Bourne, M. & Edwards, R.H.T. (1993) Unbiased estimation of whole body composition in man by the Cavalieri method using magnetic resonance imaging. J. Microsc. 171, 239–253.

Royet, J.P., Jourdan, F. & Ploye, H. (1989a) Morphometric modifications associated with early sensory experience in the rat olfactory bulb: I. volumetric study of the bulbar layers. J. Comp. Neurol. 289, 586–593.

Royet, J.P., Jourdan, F., Ploye, H. & Souchier, C. (1989b) Morphometric modifications associated with early sensory experience in the rat olfactory bulb: II. stereological study of the glomerular population. J. Comp. Neurol. 289, 594–609.

Thioulouse, J., Houllier, F. & Onillon, J.C. (1985a) Variables régionalisées et dénombrements d'insectes: cas unidimensionnel. C. R. Acad. Sci. Paris, Série III, 301, 423–428.

Thioulouse, J., Mathy, B. & Ploye, H. (1985b) Estimation d'un volume à partir de coupes sériées: sous-échantillonnage, covariogramme transitif et calcul de précision. Mikroskopie, 42, 215–224.

Weibel, E.R. (1979) Stereological Methods. Vol. 1. Practical Methods for Biological Morphometry. Academic Press, London.

Weibel, E.R. (1980) Stereological Methods. Vol. 2. Theoretical Foundations. Academic Press, London.