

DESIGN OF A SAMPLING NETWORK FOR AN ESTUARY IN THE COLOMBIAN CARIBBEAN

por

Ramón Giraldo Henao*, Néstor Méndez** & David Ospina***

Resumen

Giraldo Henao R., Méndez, N., & Ospina D.: Design of a sampling network for an estuary in the Colombian Caribbean. *Rev. Acad. Colomb. Cienc.* **25**(97):509-518, 2001. ISSN 0370-3908.

Se diseñó una red de muestreo para el monitoreo de variables fisicoquímicas y biológicas en el estuario Ciénaga Grande de Santa Marta, ubicado en la costa norte de Colombia. Inicialmente, a través de muestreo sistemático de cuadrículas, se escogieron 115 puntos para medir las variables consideradas. Con base en los datos observados se estimó, para cada variable, la estructura de autocorrelación espacial por medio de la función de semivarianza. Posteriormente, para redes de diferente tamaño, se calcularon las correspondientes varianzas de predicción, tomando como base los modelos de semivarianza ajustados. La comparación de las varianzas de predicción para las diferentes redes y de los costos asociados con cada una de ellas, permitió establecer un conjunto de sitios de muestreo que a un costo razonable disminuyen el error de predicción para las variables de interés.

Palabras clave: Estuario, geoestadística, redes de muestreo.

Abstract

A network for monitoring physical chemistry and biological variables in the Ciénaga Grande de Santa Marta estuary, in the Caribbean coast of Colombia, was designed. Initially, through systematic sampling on a square grid, a set of 115 sampling points was chosen to measure the variables considered. Based on the data provided, a spatial auto-correlation structure for each variable was estimated through the semivariance function. Later, for different size networks, the kriging prediction variances were calculated, taking the adjusted semivariogram models as a basis. The comparison

* Department of Mathematics and Statistics National University of Colombia. Address: Departamento de Matemáticas y Estadística Universidad Nacional de Colombia, Ciudad Universitaria, Bogotá D.C. Colombia, Sudamérica. Phone: 571 3165000 ext. 13164 Fax: 5713165327. E-Mail: rgiraldo@matematicas.unal.edu.co

** Department of Physics, National University of Colombia. Address: Departamento de Física. Universidad Nacional de Colombia. Ciudad Universitaria, Bogotá, D.C. Colombia, Sudamérica. E-Mail: xman@coll.telecom.com.co

*** Department of Mathematics and Statistics National University of Colombia. Bogotá, D.C. Colombia. Address: Departamento de Matemáticas y Estadística. Universidad Nacional de Colombia, Ciudad Universitaria. Bogotá, D.C. Colombia, Sudamérica. Phone: 571 3165000 -13199. Fax: 5713165327. E-Mail: dospina@matematicas.unal.edu.co

among the prediction variances for the different networks and their associated costs allowed establishing a set of sampling sites, that at a reasonable cost, substantially diminishes the prediction error for the variables of interest.

Key Words: Estuary, geostatistics, sampling networks.

Introduction

In most environmental studies, obtaining and analyzing the data is often a slow and costly process that produces skepticism among the decision making entities. Because of this, it is of maximum importance to put sampling systems that provide the best possible information quickly and at low cost.

When a study requires following up the principal factors that control the processes of productivity in estuaries and other aquatic ecosystems, one must frequently recur to biological and ecological type criteria to select the sampling points. This procedure, although it is valid to satisfy very specific objectives, in many cases it is not sufficient and does not allow to have an integrated view of the ecosystem. In the best case, it makes the information that is taken redundant, and in this way increases the cost of the research. Because of this, it is necessary to establish a set of sampling sites that not only identify the conditions of the ecosystem in very strategic regions, but that also supplies general information on the set.

In environmental statistics, there are different approaches for solving the problem of estimating the size of the sample and the location of the sampling sites. Some classical approaches suppose taking independent random samples based on some adaptations of traditional sampling, carrying out the corresponding estimates. This is the case, among others, of the sampling of squares, transects or intercepts. Under these studies, the estimation of the parameters of interest are carried out assuming some probabilistic models (Dale *et al.*, 1991) or through some of the mathematical expressions proper of the sampling design used (Thompson, 1992; Seber 1986). Caselton & Zidek (1984) proposed selecting a monitoring network that is formulated as a decision problem whose solutions may be optimized. This theory assumed a multivariate normal structure in the data.

On the other hand, the problem has been treated assuming that the phenomenon to be studied represents a stochastic process. In this case, it is treated with regionalized variables (variables measured in a region) and it is supposed that they have structures of spatial auto-

correlation. In this respect Russo (1984) and Bresler & Green (1982), proposed procedures that are based in some of the criteria associated with the distance between the pairs of points or with the number of pairs of points by lag, respectively. The fundamental goal of both methods is to find the suitable configuration of points to calculate the semivariogram. The advantage is that it does not require any initial information about the characteristic of interest. Nevertheless, problems can come up when point distributions are proposed that let large zones of the region to be studied with no observation point at all. McBratney *et al.* (1981) and McBratney & Webster (1981), presented a procedure that consists of selecting a sampling network that minimizes the standard prediction of kriging error. This method differs from the two foregoing ones in that it requires initial information of the variable, allowing an estimation of the semivariance function.

In this study, the methodology proposed by McBratney *et al.* (1981) was applied with the goal of designing a sampling network for the estuary called Ciénaga Grande de Santa Marta (CGSM) located on the northern coast of Colombia (Figure 1; IGAC, 1973). The CGSM, because of its large area (450 km²) and its ecological and economic importance (more than 5 lakeside towns living off the ecosystem), has been the center of different kinds of studies. Over the last two decades, it has been showing signs of deterioration, and some civil works trying to recover it, have been carried out. For the monitoring of the changes that have been taking place in the ecosystem, it became necessary to have a set of sampling sites that allow an integrated overview of the behavior of the principal variables that govern its productivity processes.

Method

This study is developed using geostatistical methods. The Geostatistics is a branch of statistics that treats spatial phenomena (Journel & Huijbregts, 1978). Its primordial interest is the estimation, prediction and simulation of such phenomena (Myers, 1987).

Geostatistics provides a way of describing spatial continuity, which is an essential distinctive feature of

many natural phenomena and provides adaptation of classical regression techniques to take advantage of this continuity (Isaaks & Srivastava, 1987). Petiglas (1996), defines it as an application of probability theory to the statistical prediction of regionalized variables. The results of the prediction process may be applied with diverse objectives, among others, in the design of sampling networks (Cressie, 1989). The geostatistical analysis is a two-step procedure. First, the spatial structure of the variable is examined with the semivariance analysis. Once a spatial structure has been identified and accurately described by a suitable model, the kriging procedure provides optimal interpolation of the variable at unsampled sites (Rossi *et al.*, 1995).

The difference between kriging and deterministic methods is that in kriging there is estimation of the prediction variance in each point of prediction and consequently a measure of the prediction error. The prediction variance of each point is calculated by (Cressie, 1991):

$$\sigma_0^2 = \sum_{i=1}^n \lambda_i \gamma_{i0} + \mu \quad (1)$$

γ_{i0} is equal to the semivariance function calculated for the distance between the i^{th} -sampling observation and the point where the prediction is desired. The λ_i are calculated by finding the minimum values of the variance of prediction function subject to the restriction that the predictor will be unbiased, for which the Lagrange multipliers m are used. In Giraldo *et al.* (2000) is summarized both the theory of estimation of the spatial structure and prediction method by Kriging.

From equation (1) it is evident that the prediction variance is not constant as in the classic case, as it depends on the semivariance function, which is a monotone function which increases with the distance between the observation and the points where the prediction is made. McBratney *et al.* (1981) shows that, for any sample density, the maximum distance between one point of observation and an interpolation point is minimum when the configuration of the observed points is made in a triangular grid. Under this point, smaller prediction variances will be obtained. Nevertheless, this same author and Warrick *et al.* (1986) indicate that for logistic reasons referent to the location of the field sites and to minimize travelling from point to point, a square grid may be preferable.

Accordingly, the problem of design sampling networks is limited to establishing for different sampling networks, with either an equilateral triangular grid or a square one, the relationship between the maximum prediction variances (those obtained at the center of the triangle or square) and their associated costs. In this way, the necessary cost to reach a certain degree of security can be deduced immediately or, on the contrary, the prediction variance if the cost is prefixed.

Data and procedure

The information used for the analysis was taken during the intensive sampling campaign carried in March 1997, at the CGSM (Figure 1). The area is particularly arid, with a dry season from December to August (interrupted by a short rainy season from May to June) and a major rainy season from September to November (Wiedemann, 1973). Water samples from the surface of the water column were analyzed for the following variables: temperature ($^{\circ}\text{C}$), salinity, total suspended solids (mg l^{-1}), depth (m), silicates ($\mu\text{mol l}^{-1}$), chlorophyll "a" ($\mu\text{g l}^{-1}$), dissolved oxygen (mg l^{-1}), nitrites ($\mu\text{mol l}^{-1}$) and chlorophyll "c" ($\mu\text{g l}^{-1}$). These variables are considered to be of great influence in the primary productivity processes and in the biodiversity of aquatic ecosystems like the one being considered (Vidal, 1995). Between 103 and 114 observations were obtained for each variable. The data was taken throughout the system by systematic samples of squares of 4 km^2 (Figure 2). The location of each sampling point was carried out using a geo-positioning device.

Once the data were obtained, descriptive measures were calculated to summarize the information. Through dispersion charts (that are not included in the text) the spatial stationarity assumption (Giraldo *et al.*, 2000) was evaluated. For each variable, the spatial auto-correlation structure was estimated through the experimental semivariogram. In the same way, the theoretical models were adjusted using the GeoEAS software (Englund & Sparks, 1988) and sampling networks were simulated with square grids of 4 (the observed), 9, 16, 25 and 36 km^2 , respectively (Figure 3). The corresponding prediction variances of each variable were estimated taking as the basis the estimated spatial correlation models.

The prediction variances that were obtained were related to the associated costs of measuring each variable in each sampling density. The final decision on the proposed sampling network was based on practical criteria founded on the prediction variance-cost relationship.

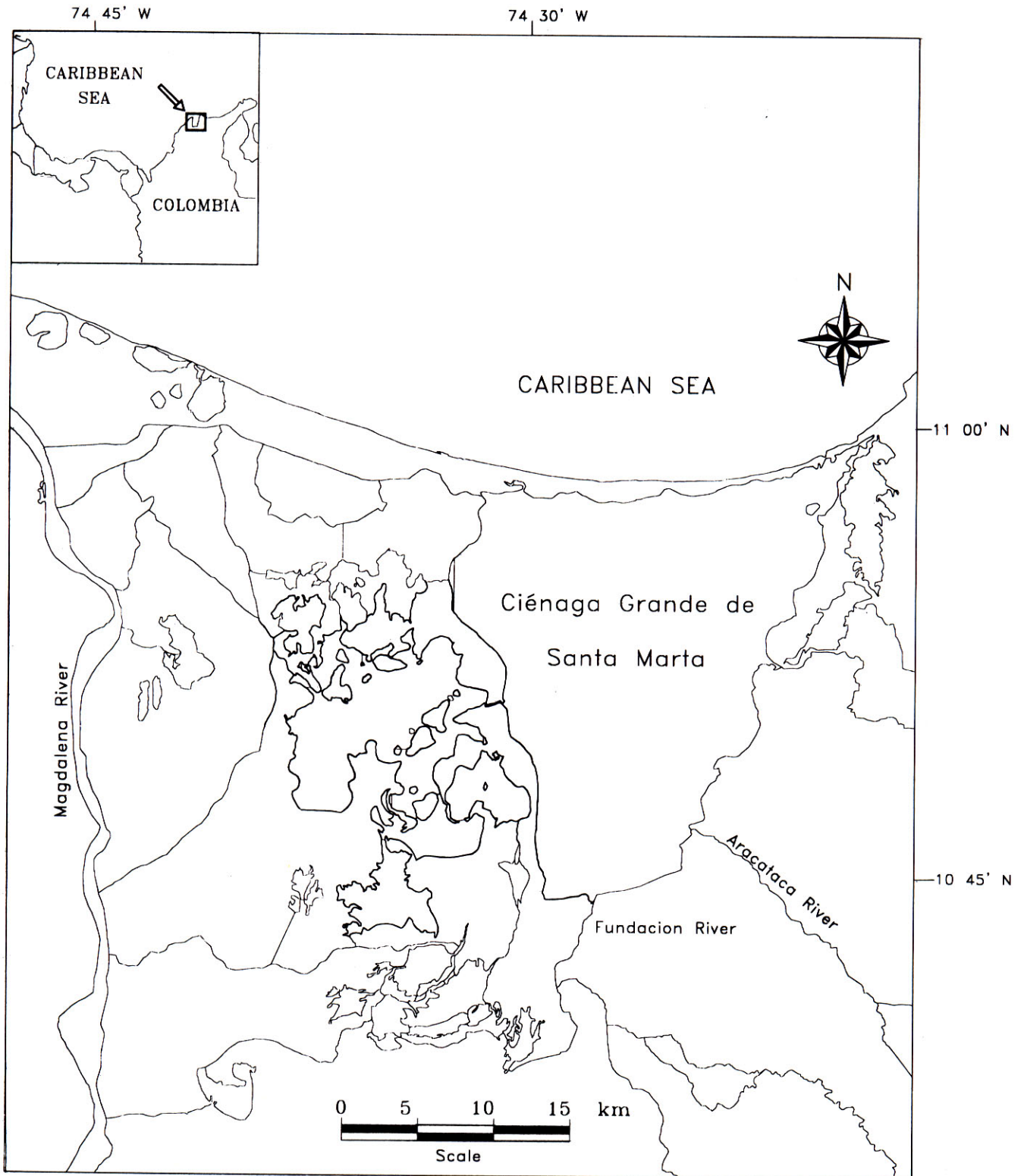


Figure 1. Geographical location of the Ciénaga Grande de Santa Marta estuary.

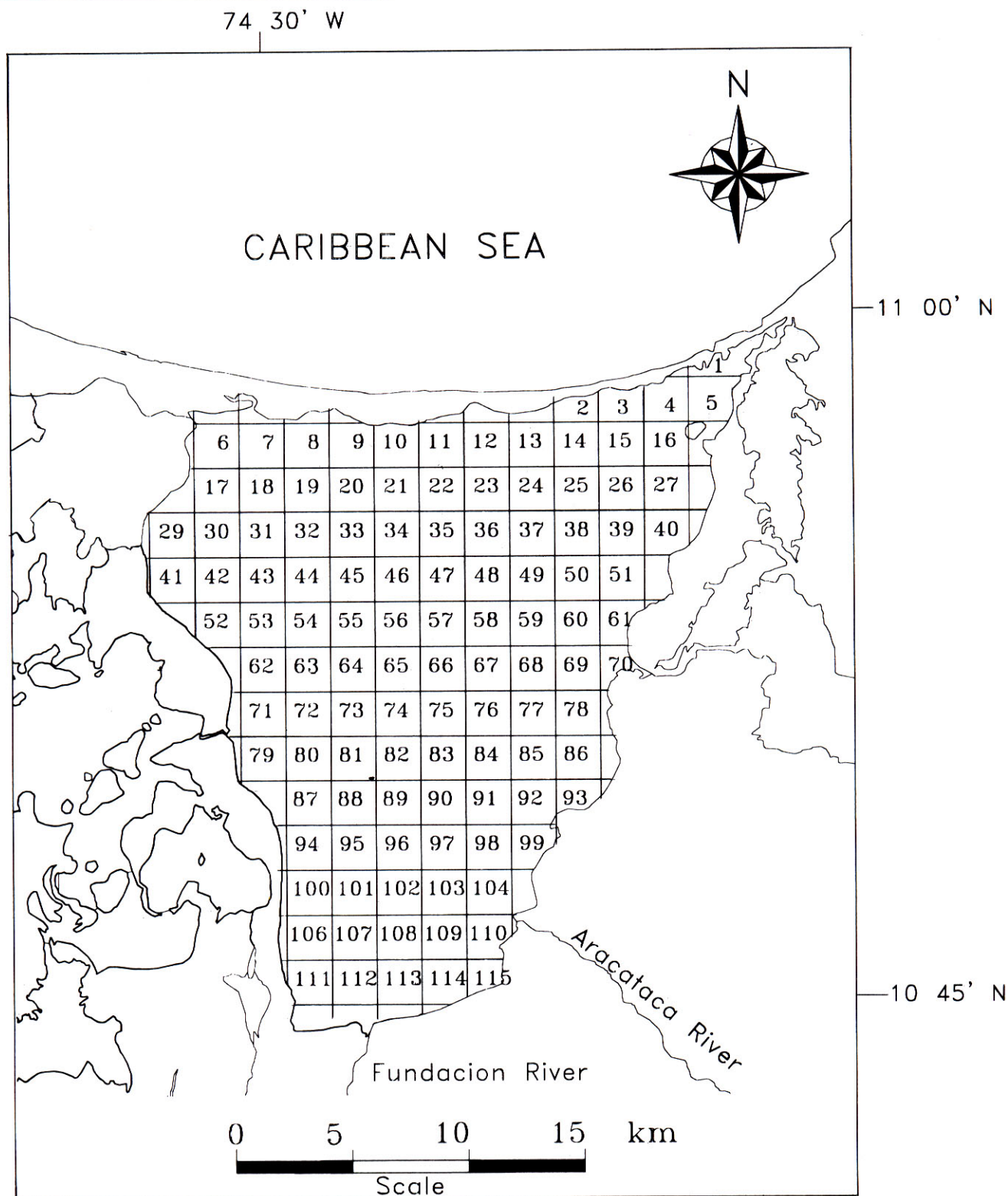


Figure 2. Initial sampling network (squares of 4 km²) used to capture data. Sampling points were at center of each square.

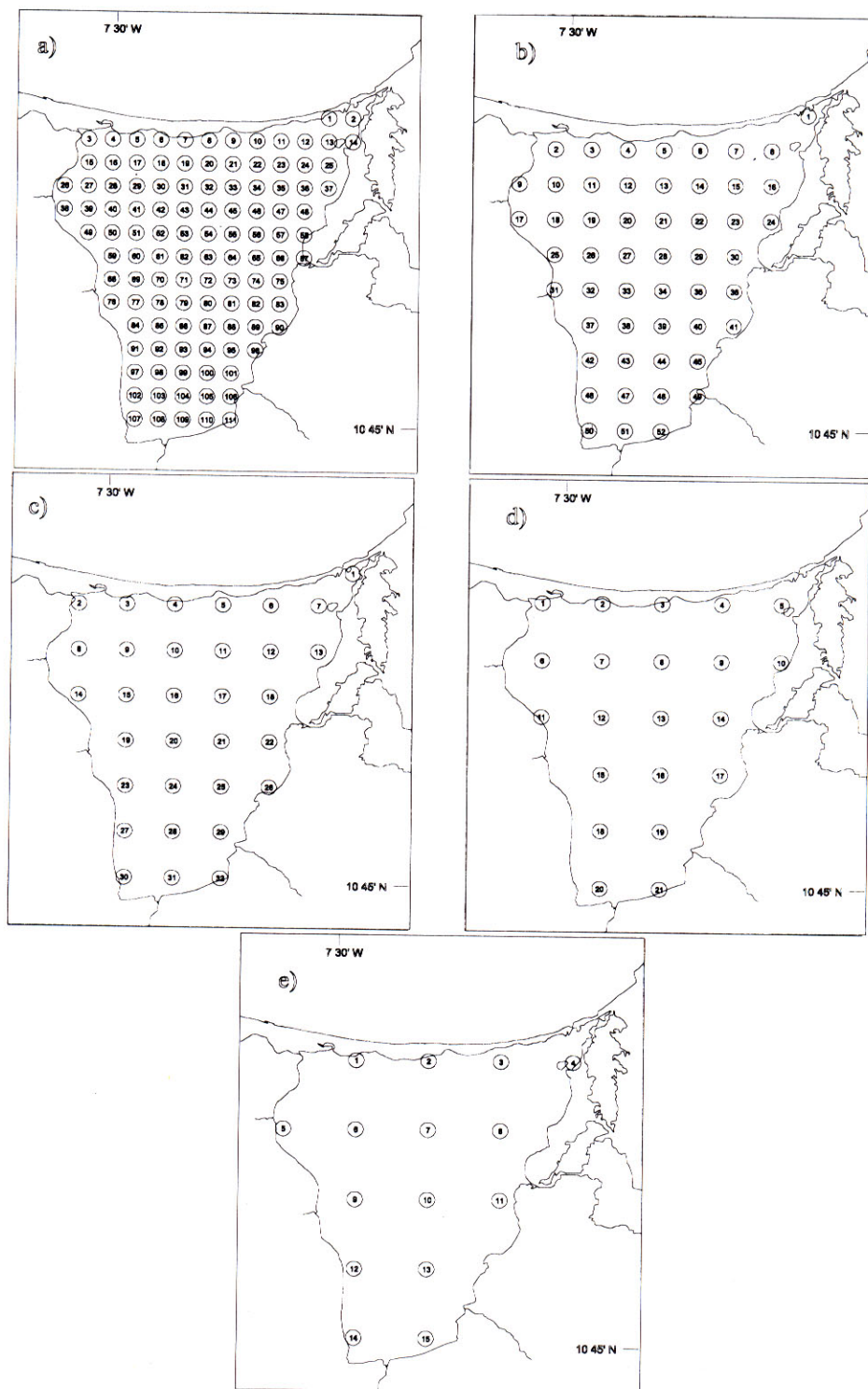


Figure 3. The sampling networks under which the estimation of the prediction variances were made for each one of the variables considered, assuming the estimated semivariance models. The distances between the sampling points: a) 2000 m; b) 3000 m; c) 4000 m; d) 5000 m and e) 6000 m.

Results and discussion

In general, according to the mean, minimum, and maximum values calculated and observed (Table 1), the variables considered have atypical magnitudes compared to those historically observed for the same time of the year (Giraldo *et al.*, 1995). The mean value for salinity was "low" compared to that registered for the same time of the year in other studies (about 25). The opposite happened in variables like nitrites, silicates, and chlorophyll, in which the observed magnitudes were similar to those reported in other studies (Vidal, 1995) for the rainiest months of the year, in which it is reasonable to find them in high concentrations, due to the larger supply of water from the rivers that run into the ecosystem. The above may be due to a possible increase in the flow of the rivers that run into the estuary during the rainy month preceding the sampling, as a consequence of the *el niño* phenomenon in the area at the end of 1996.

However, for the purpose of this study these differences are not an obstacle, since in fact it is assumed that the establishment of an optimum set of sampling points does not depend on the magnitude of the variables, but on the spatial correlation structure.

The values calculated for the coefficient of variation (Table 1) indicate that, with the exception of nitrite and chlorophyll "c", the variables are in general homogeneous (variation coefficients less than 40%). This would imply, from at least one classical point of view, that relatively small sampling sizes can be established for the follow-up of the variables.

The adjusted semivariance models (Table 2) show strong spatial association patterns of the variables in the area. The ranges that are reached, up to 30 km, are relatively high if one takes into account that the distance between the extreme north and south of the system (the

Table 1. Descriptive statistics of the physical-chemical and biological variables measured in Ciénaga Grande de Santa Marta, Colombia. Sampling carried out in March, 1997. Calculations based on data of n sites. S.D. = Standard Deviation; C.V. = Coefficient of Variation

Variable	n	Mean	S.D	Minimum	Maximun	V. C (%)
Depht	114	1.47	0.35	0.25	2.50	24.1
Temperature	114	29.43	2.12	26.00	33.20	7.2
Salinity	114	17.62	2.85	13.02	34.95	16.1
Dissolved oxygen	114	8.80	3.25	3.03	16.29	36.9
Total Suspended Solids	103	218.20	41.18	103.00	318.00	18.8
Nitrites	112	0.43	0.30	0.01	1.61	70.8
Silicates	112	245.29	61.51	10.99	358.74	25.1
Chlorophyll a	107	132.44	31.58	2.91	198.35	23.8
Chlorophyll c	107	8.94	7.76	0.00	31.41	86.8

Table 2. Theoretical semivariance models adjusted to experimental semivariograms calculated from information on physical-chemical and biological variables measured in two sampling expeditions carried out in March and in October 1997 in the Ciénaga Grande de Santa Marta, Colombia

Variable	Type	Nugget	Sill	Range (m)	r ²
Depht	Gaussian	0.07	0.12	24850	0.993
Temperature	Gaussian	0.32	6.45	15230	0.999
Salinity	Lineal	0.18	12.30	20000	0.897
Dissolved oxygen	Gaussian	1.83	14.32	12940	0.998
Total suspended solids	Linear	1087.10	1138.20	22000	0.907
Nitrites	Linear	0.07	0.04	22000	0.876
Silicates	Gaussian	1810.00	2089.00	7240	0.999
Chlorophyll a	Spherical	116.68	597.40	6940	0.980
Chlorophyll c	Linear	29.95	81.68	30000	0.944

longest distance) is not more than 30 km. The above gives rise to a reduction in the kriging prediction variance suggesting a smaller number of sampling sites. This underlines the fact that, with respect to the other two parameters, in no case the value of the nugget was greater than 50% of the value of the sill (Table 2). This, according to **Diaz-Francés** (1993), is recommendable for the spatial correlation model to adequately describe reality.

Given that the goal of the study is not to make a description of the distribution of the variables in the system, there is no summary of the results that were obtained in the prediction process in the text. Nevertheless, it is valid to mention that the prediction errors (Table 3) in almost all cases were less than 5% of

the predicted values, which indicates that if any specific probability distribution is assumed, low confidence intervals would be obtained.

As it was expected, the estimated standard prediction errors grow as a function of the distance between the sampling points (Table 3). Salinity is the variable with which the greatest gain in precision was attained (35%) when changing to the less dense network (Figure 3 (e)), to the densest (Figure 3(a)), (Table 4). Other variables such as temperature, dissolved oxygen, silicates, and chlorophyll "a" had precision increases that varied between 15.9% and 23.8% (Table 4). Finally, for depth, nitrites, total suspended solids and chlorophyll "c", the increase in precision was only in percentages between

Table 3. Maximum standard prediction error (square root of the variance) of each variable for sampling networks with grids of 4, 9, 16, 25, and 36 km²

Variables	Network Size (Distance in meters between sampling points)				
	2000	3000	4000	5000	6000
Depth	0.2825	0.2874	0.2930	0.3002	0.3070
Temperature	0.6380	0.6690	0.7046	0.7632	0.8373
Salinity	0.9096	1.0511	1.1676	1.2965	1.4075
Dissolved oxygen	1.5145	1.5917	1.6752	1.7977	1.9431
Total suspended solids	35.6363	36.4021	37.0459	37.8076	38.5197
Nitrites	0.2832	0.2875	0.2913	0.2958	0.3003
Silicates	47.6524	50.2070	52.3806	54.6797	56.6932
Chlorophyll a	19.4634	21.2041	22.5233	23.5582	24.2163
Chlorophyll c	6.1071	6.2977	6.4536	6.6336	6.7967

Table 4. Increase in precision (percentage reduction of the standard prediction error) of each sampling network (observed and simulated) with respect to the 6000 meter network (the least dense)

Variables	Network Size (Distance in meters between sampling points)				
	2000	3000	4000	5000	6000
Depth	8.0	6.4	4.6	2.2	0
Temperature	23.8	20.1	15.8	8.8	0
Salinity	35.4	25.3	17.0	7.9	0
Dissolved oxygen	22.1	18.1	13.8	7.5	0
Total suspended solids	7.5	5.5	3.8	1.8	0
Nitrites	5.7	4.3	3.0	1.5	0
Silicates	15.9	11.4	7.6	3.6	0
Chlorophyll a	19.6	12.4	7.0	2.7	0
Chlorophyll c	10.1	7.3	5.0	2.4	0

5.7% and 10.1% (Table 4). Obviously, when comparing the intermediate networks, those with grid distances between 3000, 4000, and 5000 m. (Figure 3 (b), 3(c), and 3(d)), with the 6000 m network (Figure 3(e)), the relative increase in precision was much less (Table 4).

The sampling costs associated with each variable under each sampling density were different, with the exception of the variables of temperature, depth, and salinity, whose cost for the 2000 m network was much lower than the other variables (Figure 4). For some of the variables (dissolved oxygen, silicates, and chlorophyll) going from a 3000 m network to a 2000 m the sampling cost increased in more than Colombian \$500.000 (about US \$240).

Hence, for temperature and salinity, it would be much more convenient to make an intense sampling (the densest network) as this would increase the efficiency in a considerable percentage (23% and 35%, respectively, with net costs increased in only about Colombian \$190.000 (about US \$90) (Figure 4). For depth, even if the sampling costs are not significantly increased (Figure 4), is more recommendable to sample it in the less dense network, given that the efficiency is increased by a maximum of 7% in comparison with the other networks (Table 4). For nitrite, total suspended solids and chlorophyll "c", there is only a little increase in the efficiency with increasing network density (Table 3); on the contrary, the costs, especially in the 2000 m network, increase considerably. Hence the less dense networks (5000 m and 6000 m between sampling points) are the most adequate for the follow-up of these variables. In the remaining variables (dissolved oxygen, silicates and chlorophyll "a"), the decision is complex given that there are considerable increases in the costs (Figure 4) and the efficiencies with increasing density (Table 4).

A global analysis of the increases in cost and in efficiency (Table 5) clearly show that the 2000 m network is the least recommendable given that, compared to the

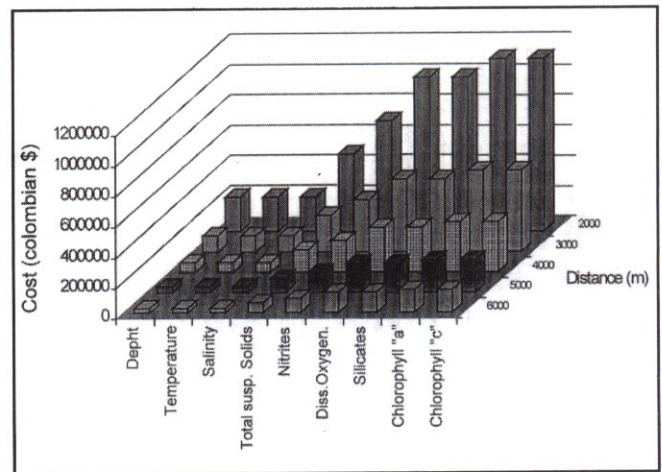


Figure 4. The estimates sampling costs for each variable under five sampling networks, in which the prediction variances were calculated.

3000 m one, there is a high increase in costs (more than 200%) but the relative efficiency increases in only 4.9%. While in relative terms, the change in efficiency and the costs going from one network to another with a greater number of points is similar (with the exception of the 2000 m one), the networks with distances between sampling points of 4000 m and 5000 m should be considered [(Figure 3(b) and 3 (c))] to be the most advisable, given that they produce a greater efficiency that the one obtained in the 6000 m, with slightly higher costs (Figure 4).

The suggestion given in the foregoing paragraph about the optimum sampling arrangement to monitor the variables considered in the ecosystem under study, are not in any way absolute. In the final analysis, while comparing the functions of cost and of statistical efficiency, many purely empirical criteria have been used. Nevertheless it is considered that the agencies that make the final decision should have a tool that allows them to plan the most adequate monitoring strategy for the future.

Table 5. Costs increases (net and relative) and efficiency (average of the nine variables) with increasing network density

Networks From - To	Increase (Colombian \$)	Relative increase in cost (%)	Relative increase in efficiency (%)
6000 m - 5000 m	332.900	140	4.27
5000 m - 4000 m	611.150	152	4.63
4000 m - 3000 m	1'112.000	162	4.15
3000 m - 2000 m	3'282.350	213	4.90

Literature cited

- Bresler, E. & Green, R. E.** 1982. Soil parameters and sampling scheme for characterizing soil hydraulic properties of a watershed. Technical report Number 148, University of Hawaii, Honolulu.
- Casleton, W. F. & Zidek, J. V.** 1984. Optimal monitoring networks designs. *Statistics & Probability Letters*, 2: 223-227.
- Cressie, N.** 1989. Geostatistics. *The American Statistician*, 43(4): 611-623.
- . 1991. *Statistical for spatial data*, John Wiley & Sons, New York.
- Díaz-Francés, E.** 1993. Introducción a los conceptos básicos de geoestadística. Memorias del seminario en Estadística y Medio Ambiente, CIMAT, Guanajuato, México.
- Dale, V. H., Franklin, R. L. A., Post, W. M. & Gardner, R. H.** 1991. Sampling ecological information: Choice of sample size. *Ecological Modeling*, 57: 1-10.
- Englund, E. and Sparks, A.** 1988. *GeoEAS User's guide*. E. P. A. Las Vegas, Nevada.
- Giraldo, R., J. Martínez, S. Zea, H. Hurtado & E. Madera.** 1995. Análisis de clasificación de series temporales: El caso de la salinidad en la Ciénaga Grande de Santa Marta, Colombia. *An. Inst. Invest. Mar. Punta Betín*, 24: 123-134.
- Giraldo, R., W. Troncoso, J. E. Mancera & N. Méndez.** 2000. Geoestadística. Una herramienta para la modelación en estuarios. *Rev. Acad. Col. Ciencias*: 24(90), 59-72.
- IGAC.** 1973. Monografía del Departamento del Magdalena. Instituto Geográfico Agustín Codazzi, Bogotá.
- Isaaks, E. & Srivastava, R. M.** 1987. *Applied Geostatistics*. Oxford University Press, New York.
- Journel, A.G. & Huijbregts, C. J.** 1978. *Mining Geostatistics*, Academic Press, New York.
- McBratney, A. B., Webster, R. & Burgess, T. M.** 1981. The design of optimal sampling schemes for local estimation and mapping of regionalized variables I. *Computers and Geosciences*, 7(4): 331-334.
- McBratney, A. B. & Webster, R.** 1981. The design of optimal sampling schemes for local estimation and mapping of regionalized variables II. *Computers and Geosciences*, 7(4) : 335-365.
- Myers, D. E.** 1987. Optimization of sampling locations for variogram calculations. *Water Resources Research*, 23(3): 283-293.
- Petitgas, P.** 1996. Geostatistics and their applications to fisheries survey data. in Megrey, B. A. and Moksness, E. (eds), *Computers in Fisheries Research*, pp 114-142, Chapman-Hall, London.
- Rossi, J. P., Lavelle, P. & J. E. Tondoh.** 1995. Statistical tool for soil biology. *Geostatistical analysis. Eur. J. Soil Biol.* 31(4), 173-181.
- Russo, D.** 1984. Design of an optimal sampling network for estimating the variogram. *Soil Sci. Soc. Am. J.*, 48: 708-716.
- Seber, G. A. F.** 1986. A review of estimating animal abundance. *Biometrics*, 42: 267- 269.
- Thompson, S. K.** 1992. *Sampling*, John Wiley & Sons, New York.
- Vidal, L. A.** 1995. Estudio del fitoplancton en el sistema lagunar estuarino tropical Ciénaga Grande de Santa Marta, Colombia, durante el año 1987. Tesis M. Sc. Biol. Mar., Universidad Nacional de Colombia, Santa Fe de Bogotá.
- Warrick, A. W., D. E. Myers and Nielsen. D. R.** 1986. Geostatistical methods applied to soil science. *Methods of soil analysis. Part 1, Physical and mineralogical methods*, Agronomy Monograph Number 9: 53 – 81,
- Wiedemann, H. V.** 1973. Reconnaissance of the C. G. S. M., Colombia: Physical parameters and geologic history. *Milt. Inst. Colombo-Aleman Inv. Cientif.*, 7: 85-119.