

Estimation of Spatial Dependence with GEE

YOON DONG LEE¹, HYEMI CHOI²

ABSTRACT

We consider an efficient parametric estimation method of spatial dependence in weak stationary processes. Spatial dependence is modeled through variogram and correlogram. Most of parametric estimation methods of correlogram use two step method; nonparametric estimation and parametric integration. We bind these two steps into one step by using GEE method instead of least squares type optimization. Our one step method is more efficient statistically and gives a clear interpretation of related concepts used in traditional two step methods.

Keywords. Spatial dependence, Variogram, GEE, Quasi-likelihood

1. INTRODUCTION

In spatial models, usually the dependence of stationary spatial process $\{Z(\mathbf{s}) : \mathbf{s} \in \mathbb{R}^d\}$ is described by $2\gamma(\mathbf{h}) = \text{var}(Z(\mathbf{s}+\mathbf{h}) - Z(\mathbf{s}))$ or $\rho(\mathbf{h}) = E(Z(\mathbf{s}+\mathbf{h})Z(\mathbf{s}))$, which are called as *variogram* and *correlogram* respectively. Note that $2\gamma(\mathbf{h}) = 2\sigma^2(1-\rho(\mathbf{h}))$ when $\sigma^2 = \text{var}(Z(\mathbf{s}))$. For variogram estimation, various easy-to-use methods have been proposed.

The methods proposed by Journel and Huijbregts (1978) and Cressie (1985) are divided into two steps, the nonparametric estimation step and the parametric integration step. In the first step, they estimated variogram for a preselected set $H = \{\mathbf{h}_1, \dots, \mathbf{h}_K\}$ of lags by the Matheron's estimator (1962) which is the average of the squared differences $(Z(\mathbf{s}_1) - Z(\mathbf{s}_2))^2$ for all possible pairs of the observed values $Z(\mathbf{s}_1)$ and $Z(\mathbf{s}_2)$ satisfying $\mathbf{s}_1 = \mathbf{s}_2 + \mathbf{h}$ for $\mathbf{s}_1, \mathbf{s}_2 \in \mathbf{D}$, the sampling region. In case of irregularly spaced data, Matheron's method is modified as follows: let $T(\mathbf{h})$ be a predetermined set of lag \mathbf{h}' regarded to be \mathbf{h} approximately and average the squared differences at the pairs satisfying $\mathbf{s}_1 = \mathbf{s}_2 + \mathbf{h}'$, $\mathbf{h}' \in T(\mathbf{h})$. In the

¹Department of Applied Statistics, University of Suwon

²Department of Statistics, Seoul National University

second step, the model variogram $2\gamma(\mathbf{h}; \beta)$ defined through the parameter β is estimated by plugging $\hat{\beta}$, which is obtained by versions of least squares estimator (LSE), based on $2\hat{\gamma}(\mathbf{h})$, $\mathbf{h} \in H$ at the first step .

In this paper, we propose one step estimation method of correlogram $\rho(\mathbf{h}; \beta)$, directly using the squared differences $Y(\mathbf{s}_1, \mathbf{s}_2)$, $(\mathbf{s}_1, \mathbf{s}_2) \in \mathbf{D}^2$. More efficient estimator $\hat{\beta}$ than traditionally used LSEs is proposed, which is based on generalized estimating equation method and its asymptotic properties are shown in this paper.

The generalized estimating equations (GEE) suggested by Liang and Zeger (1986) provides a general tool to handle non-normal dependent data. Albert and McShane (1995) applied GEE to the estimation problem for binary spatial data. They considered the estimation of the mean structure and dependence structure simultaneously by iteratively solving the two estimating equations: one for mean and the other for dependence. They suggested two step method for the dependence estimation.

This paper is only concerned with dependence estimation. The proposed method is expected to be more statistically efficient than the two step estimation methods, since it reduces the loss of the captured information in the first step of the two step method. Moreover, this one step method reduces ambiguity and arbitrariness in selecting the set H of lags of interest and the tolerance region $T(\mathbf{h})$, $\mathbf{h} \in H$. The major difficulties in handling $Y(\mathbf{s}_1, \mathbf{s}_2)$, $(\mathbf{s}_1, \mathbf{s}_2) \in \mathbf{D}^2$ directly, instead of the estimated variogram $2\hat{\gamma}(\mathbf{h})$, $\mathbf{h} \in H$, reside in the non-normality of the quantity $Y(\mathbf{s}_1, \mathbf{s}_2)$ and dependence among the quantities.

Throughout the paper, we assume that process $Z(\cdot)$ is normally or nearly normally distributed with the mean 0 and the finite variance σ^2 . We describe the proposed method in section 2. The statistical and computational methods to reduce the computational burden are also outlined.

2. GEE FOR CORRELOGRAM ESTIMATION

When the process $Z(\mathbf{s})$ is Gaussian with mean 0 and variance σ^2 , the squared term $Y(\mathbf{s}_1, \mathbf{s}_2) = (Z(\mathbf{s}_1) - Z(\mathbf{s}_2))^2$ has a chi-squared distribution with one degree of freedom and a scale parameter $2\gamma(\mathbf{s}_1 - \mathbf{s}_2; \beta)$, $\beta \in \mathbb{R}^q$, *i.e.* $Y(\mathbf{s}_1, \mathbf{s}_2) \sim 2\gamma(\mathbf{s}_1 - \mathbf{s}_2; \beta) \cdot \chi^2(1)$. The mean and the variance of $Y(\mathbf{s}_1, \mathbf{s}_2)$ are $2\gamma(\mathbf{s}_1 - \mathbf{s}_2; \beta)$

and $2\{2\gamma(\mathbf{s}_1 - \mathbf{s}_2; \boldsymbol{\beta})\}^2$, respectively. The correlation $r((\mathbf{s}_1, \mathbf{s}_2), (\mathbf{s}_3, \mathbf{s}_4))$ between $Y(\mathbf{s}_1, \mathbf{s}_2)$ and $Y(\mathbf{s}_3, \mathbf{s}_4)$ is expressed in terms of the correlogram $\rho(\mathbf{h})$ as follows:

$$r((\mathbf{s}_1, \mathbf{s}_2), (\mathbf{s}_3, \mathbf{s}_4)) = \frac{\{\rho(\mathbf{s}_1 - \mathbf{s}_3) + \rho(\mathbf{s}_2 - \mathbf{s}_4) - \rho(\mathbf{s}_1 - \mathbf{s}_4) - \rho(\mathbf{s}_2 - \mathbf{s}_3)\}^2}{4\{1 - \rho(\mathbf{s}_1 - \mathbf{s}_2)\}\{1 - \rho(\mathbf{s}_3 - \mathbf{s}_4)\}}. \quad (2.1)$$

By modeling $\rho(\mathbf{h})$ with respect to a parameter $\boldsymbol{\beta}$, $r((\mathbf{s}_1, \mathbf{s}_2), (\mathbf{s}_3, \mathbf{s}_4))$ depends on $\boldsymbol{\beta}$, say $r((\mathbf{s}_1, \mathbf{s}_2), (\mathbf{s}_3, \mathbf{s}_4); \boldsymbol{\beta})$. For notational convenience, let $[\mathbf{s}_1, \mathbf{s}_2]$ be $(\mathbf{s}_1, \mathbf{s}_2)$ or $(\mathbf{s}_2, \mathbf{s}_1)$, and $Y[\mathbf{s}_1, \mathbf{s}_2]$ denote $Y(\mathbf{s}_1, \mathbf{s}_2)$ or $Y(\mathbf{s}_2, \mathbf{s}_1)$, noting that $Y(\mathbf{s}_1, \mathbf{s}_2) = Y(\mathbf{s}_2, \mathbf{s}_1)$. Now, consider $Y[\mathbf{s}_1, \mathbf{s}_2]$ only for the paired sampling sites between which the distance is less than a fixed constant λ ,

$$\mathcal{D}_n^\lambda = \{[\mathbf{s}_i, \mathbf{s}_j] : 0 < d(\mathbf{s}_i, \mathbf{s}_j) < \lambda, \mathbf{s}_i, \mathbf{s}_j \in \mathbf{D}_n\}$$

where $[\mathbf{s}_1, \mathbf{s}_2]$ denotes the pair of two points $\mathbf{s}_1, \mathbf{s}_2$ satisfying $\mathbf{s}_1 \prec \mathbf{s}_2$ with a complete ordering \prec on \mathbf{D} , and $d(\cdot, \cdot)$ is a metric defined on $\mathbb{R}^d \times \mathbb{R}^d$. A set \mathbf{D}_n of sampling sites denotes \mathbf{D} of which number of elements is n . Let n_λ be the number of elements contained in \mathcal{D}_n^λ . When we define \mathbf{Y}_n^λ to be the vector of Y s corresponding to all paired sampling sites in \mathcal{D}_n^λ in an order defined on \mathcal{D}_n^λ . Note that \mathbf{Y}_n^λ is a n_λ -dimensional vector. The fixed constant λ plays a role of the maximum lag. As in time series analysis, the maximum lag can be an appropriately defined bounded value regardless of sample size n (*cf.* Box and Jenkins, 1976, p33). Thus we assume that $n_\lambda = O(n)$.

The expectation $\boldsymbol{\mu}_n(\boldsymbol{\beta})$ and the (scaled) covariance $V_n(\boldsymbol{\beta})$ of \mathbf{Y}_n^λ are given by

$$\begin{aligned} \boldsymbol{\mu}_n(\boldsymbol{\beta}) &= 2\sigma^2\{1 - \rho(\mathbf{s}_i - \mathbf{s}_j; \boldsymbol{\beta})\}, \quad \text{an } n_\lambda \times 1 \text{ vector,} \\ V_n(\boldsymbol{\beta}) &= A_n(\boldsymbol{\beta})R_n(\boldsymbol{\beta})A_n(\boldsymbol{\beta}), \quad \text{an } n_\lambda \times n_\lambda \text{ matrix,} \end{aligned}$$

where $\boldsymbol{\beta}$ is a q -dimensional vector of unknown parameters, $R_n(\boldsymbol{\beta})$ is an $n_\lambda \times n_\lambda$ matrix of which elements are $r([\mathbf{s}_i, \mathbf{s}_j], [\mathbf{s}_k, \mathbf{s}_l]; \boldsymbol{\beta})$, $[\mathbf{s}_i, \mathbf{s}_j], [\mathbf{s}_k, \mathbf{s}_l] \in \mathcal{D}_n^\lambda$, and $A_n(\boldsymbol{\beta})$ a diagonal matrix with diagonal elements $(1 - \rho(\mathbf{s}_i - \mathbf{s}_j; \boldsymbol{\beta}))$ and all off-diagonal elements zero, defined along the same order of index used in defining \mathbf{Y}_n^λ .

Now, the generalized estimating equation for our model is defined as

$$U_n(\boldsymbol{\beta}) \equiv \Gamma_n^T(\boldsymbol{\beta})V_n^{-1}(\boldsymbol{\beta})(\mathbf{Y}_n^\lambda - \boldsymbol{\mu}_n(\boldsymbol{\beta})) = \mathbf{0} \quad (2.2)$$

where $\Gamma_n(\boldsymbol{\beta})$ is a $n_\lambda \times q$ matrix of which rows are defined to be $(\partial/\partial\boldsymbol{\beta})\{\rho(\mathbf{s}_1 - \mathbf{s}_2; \boldsymbol{\beta})\}$ for all $[\mathbf{s}_1, \mathbf{s}_2] \in \mathcal{D}_n^\lambda$.

The solution $\hat{\beta}_n$ of (2.2) is obtained by the Gauss-Newton type iteration, obtained from Fisher's scoring approximation,

$$\beta^{new} = \beta^{old} - (1/2\sigma^2) \cdot M_n^{-1}(\beta^{old})U_n(\beta^{old}) ,$$

where $M_n(\beta) = \Gamma_n^T(\beta)V_n^{-1}(\beta)\Gamma_n(\beta)$, as given in McCullagh and Nelder (1989). The asymptotic covariance matrix of $\hat{\beta}_n$ is $2M_n^{-1}(\beta)$. The consistency and asymptotic normality of GEE estimator $\hat{\beta}_n$ of the true value β_* are obtained as follows.

PROPOSITION 2.1. *Under mild regularity conditions, $\hat{\beta}_n \rightarrow \beta_*$ in probability as n goes to ∞ , and $\sqrt{n}(\hat{\beta}_n - \beta_*)$ is asymptotically multivariate Gaussian with zero mean and covariance matrix*

$$\Sigma(\beta_*) = \lim_{n \rightarrow \infty} 2 n \cdot M^{-1}(\beta_*) .$$

The proof is based on the relations

$$\begin{aligned} \text{cov}(\mathbf{Y}_n^\lambda) &= 2(2\sigma^2)^2 V(\beta_*), E(U_n U_n^T) = 2(2\sigma^2)^2 M(\beta_*) \\ E((\partial/\partial\beta)U_n) &= -(2\sigma^2)M(\beta_*) \\ \text{cov}(\sqrt{n}\hat{\beta}_n) &\simeq \{E((\partial/\partial\beta)U_n)\}^{-1}E(U_n U_n^T)\{E((\partial/\partial\beta)U_n)\}^{-1}, \end{aligned}$$

which follows from the view point of the standardized estimating function in Heyde (1997) and covariance of quasi-likelihood estimators for dependent observations in McCullagh and Nelder (1989, p332). The result needs the mixing property defined on the process $Z(\cdot)$ and so does on $Y(\cdot)$. Otherwise, it is shown analogously to Liang and Zeger (1986) under m -dependence condition. The proposition was stated in the view point of *increasing-domain asymptotics*, in which the region where the process is observed is increased as the number of observation increases, as usual in asymptotics of time series. In spatial setting, other asymptotics such as *infill-domain asymptotics* and *mixed-domain asymptotics* are possible to consider, of which details are found in Lahiri, Lee and Cressie (2002) and Lahiri (200x).

The main difficulty in estimation of $\hat{\beta}_n$ is the computation of a large matrix $V_n(\beta)$ and its inversion when n is quite large. To reduce the dimension n_λ of \mathbf{Y}_n^λ , taking smaller λ , or more generally performing the proposed method on a smaller subset $\tilde{\mathcal{D}}_n^\lambda$ of \mathcal{D}_n^λ seems possible without any difference in related results.

The elements of the set $\tilde{\mathcal{D}}_n^\lambda$ are selected by one's preference as for the selection of the set H of lags and the tolerance region $T(\cdot)$ in LSEs. However, for the purpose of formalization we use the set \mathcal{D}_n^λ and \mathbf{Y}_n^λ themselves.

REFERENCES

- Albert, P. and McShane, M. (1995). A Generalized Estimating Equations Approach for Spatially Correlated Binary Data: Applications to the Analysis of Neuroimaging Data. *Biometrics*, **51**, 627-638.
- Box, G.E.P. and Jenkins, G.M. (1976). *Time Series Analysis: Forecasting and Control*, 2nd edition, Holden-Day, San Francisco, CA.
- Cressie, N. (1985). Fitting variogram models by weighted least squares. *Journal of the International Association for Mathematical Geology*, **17**, 563-586.
- Heyde, C.C. (1997). *Quasi-Likelihood and Its Application: A General Approach to Optimal Parameter Estimation*, Springer-Verlag, New York.
- Journel, A.G. and Huijbregts, C.J. (1978). *Mining Geostatistics*. Academic Press, London.
- Lahiri, S.N. (200x). Central Limit Theorems for Weighted Sums of a Spatial Process Under a Class of Stochastic and Fixed Design. (???)
- Lahiri, S., Lee, Y.D. and Cressie, N. (2002) On asymptotic distribution and asymptotic efficiency of least squares estimators of spatial variogram parameters. *Journal of Statistical Planning and Inference*, **103**, 65-85.
- Liang, K. and Zeger, S (1986). Longitudinal data analysis using generalized linear models. *Biometrika*, **73**, 13-22.
- Matheron, G. (1962). Traite de Geostatistique Appliquee, Tom I. *Memoires du Bureau de Recherches Geologiques et Minières*, No. 14. Editions Technip, Paris.
- McCullagh, P. and Nelder, J. (1989). *Generalized Linear Model*, 2nd edition, Chapman and Hall.