

Comparison of Performances of Plug-in Spatial Classification Rules based on Bayesian and ML Estimators

Kestutis Ducinkas, Egle Zikariene and Lina Dreiziene

Department of Mathematics and Statistics, Klaipeda University, Herkaus Manto str. 84, Klaipeda, Lithuania

Keywords: Bayes Rule, Spatial Discriminant Function, Gaussian Random Field, Actual Risk, Training Labels Configuration.

Abstract: The problem of classifying a scalar Gaussian random field observation into one of two populations specified by a different parametric drifts and common covariance model is considered. The unknown drift and scale parameters are estimated using given a spatial training sample. This paper concerns classification procedures associated to a parametric plug-in Bayes Rule obtained by substituting the unknown parameters in the Bayes rule by their estimators. The Bayesian estimators are used for the particular prior distributions of the unknown parameters. A closed-form expression is derived for the actual risk associated to the aforementioned classification rule. An estimator of the expected risk based on the derived actual risk is used as a performance measure for the classifier incurred by the plug-in Bayes rule. A stationary Gaussian random field with an exponential covariance function sampled on a regular 2-dimensional lattice is used for the simulation experiment. A critical performance comparison between the plug-in Bayes Rule defined above and a one based on ML estimators is performed.

1 INTRODUCTION

It is known that for completely specified populations and loss function, Bayes Rule (BR) is an optimal classification procedure in the sense of minimum risk (expected losses) (Anderson, 2003). When this is not the case, the missing information is usually provided by a training sample. For parametrically specified populations, the training sample is used to obtain the estimators of statistical parameters and plugging them into the BR. The resulting classification rule is usually called a plug-in Bayes rule (PBR). Actual risk (ACR) or conditional risk is usually used as a performance measure for the PBR. Performance comparison of the PBR based on the different types of estimators can easily be done by the closed-form expression of the ACR.

Many authors have investigated the performance of the PBR when the parameters are estimated from training samples consisting of dependent observations by using the frequentist approach for the estimation (Kharin, 1996; Saltyte-Benth and Ducinkas, 2005; Batsidis and Zografos, 2011). A closed-form expression for the ACR in supervised classification of Gaussian random field observations is derived by Ducinkas (2009). Only the ML

estimators of the drift parameters and the scale parameter of covariance function are considered.

In the present paper we use Bayesian estimators instead of the ML estimators for the classification problem described above. Proposed methodology is useful for classification of images corrupted by the Gaussian spatial correlated noise (Ducinkas, Stabingiene and Stabingis, 2011).

The closed-form expression for the ACR associated with the PBR is derived. The estimator of expected risk is based on the derived actual risk and is used as the performance of the PBR which is measured by the average of the ACR usually called an empirical estimator of expected risk.

This estimator of expected risk is a case of the stationary Gaussian random field on 2-dimensional regular lattice with an exponential covariance function. The dependence of the ACR values on the statistical hyperparameters is investigated. Numerical comparison for the case of ML estimators is implemented.

2 THE MAIN CONCEPTS AND DEFINITIONS

This paper concerns classifying a Gaussian random field (GRF) $\{Z(s) : s \in D \subset R^p\}$ observations into one of two populations denoted by Ω_1, Ω_2 .

The model of an observation $Z(s)$ in a population Ω_j is

$$Z(s) = \mu_j(s) + \varepsilon(s) \tag{1}$$

where $\mu_j(s)$, is a drift (a deterministic function of locations) and $\varepsilon(s)$ is the stochastic error term, $j=1,2$.

Suppose the drift $\mu_j(s)$ can be represented as $\mu_j(s) = x' \beta_j$, where x is a $q \times 1$ vector of non-random regressors and β_j is a $q \times 1$ vector of the parameters, $j = 1, 2$.

The error term is generated by a zero mean stationary GRF $\{\varepsilon(s) : s \in D\}$ with the covariance function defined by the model for all $s, u \in D$

$$\text{cov}\{\varepsilon(s), \varepsilon(u)\} = \sigma^2 C(s-u), \tag{2}$$

where σ^2 is the variance or the scale parameter and $C(\cdot)$ is the spatial correlation function.

In the case when the covariance function parameters are known, the model (1), (2) is called a universal kriging model (Cressie, 1993).

For the given training sample, consider the classification problem of the $Z_0 = Z(s_0)$ into one of two populations when $x'(s_0)\beta_1 \neq x'(s_0)\beta_2, s_0 \in D$.

Denote by $S_n = \{s_i \in D; i = 1, \dots, n\}$ the set of locations where the training sample $T' = (Z(s_1), \dots, Z(s_n))$ is taken. It specifies the spatial sampling design or the spatial framework for the training sample (Shekhar et al., 2002).

We shall assume the deterministic spatial sampling design and all analyses are carried out conditional on S_n .

Assume that each training sample realization $T=t$ and S_n are arranged in the following way. The first n_1 components are the observations of $Z(s)$ from Ω_1 and the remaining $n_2 = n - n_1$ components are the observations of $Z(s)$ from Ω_2 . So S_n is partitioned into a union of two disjoint subsets, i. e.

$S_n = S^{(1)} \cup S^{(2)}$, where $S^{(j)}$ is a subset of S_n that contains n_j locations of the feature observations from $\Omega_j, j = 1, 2$. So each partition $\xi(S_n) = \{S^{(1)}, S^{(2)}\}$ with marked labels determines the training labels configuration (TLC).

For TLC $\xi(S_n)$, define the variable $d = |D^{(1)} - D^{(2)}|$, where $D^{(j)}$ is a sum of distances between the location s_0 and the locations in $S^{(j)}, j = 1, 2$.

The $n \times 2q$ design matrix of the training sample T denoted by X is specified by $X = X_1 \oplus X_2$, where the symbol \oplus denotes a direct sum of the matrices and X_j is a $n_j \times q$ matrix of the regressors for the observations from $\Omega_j, j = 1, 2$.

As it follows, we assume that S_n and TLC $\xi(S_n)$ are fixed. This is the case, when spatial classified training data are collected at the fixed locations (stations).

So the model of the training sample is

$$T = X\beta + E, \tag{3}$$

where $\beta = (\beta_1', \beta_2')$ is a $2q \times 1$ vector of the regression parameters and E is the n - vector of the random errors that has a multivariate Gaussian distribution $N_n(0, \sigma^2 C_n)$.

Here C_n denotes a spatial correlation matrix among the observations forming the training sample T .

Denote by c_0 the vector of correlations among Z_0 and the components of T . Let t denote the realization of T .

Since Z_0 follows the model specified in (1), the conditional distribution of Z_0 given $T = t, \Omega_j, j = 1, 2$ is Gaussian with the mean

$$\mu_t^0 = E(Z_0 | T = t; \Omega_j) = x'_0 \beta_j + \alpha'_0 (t - X\beta) \tag{4}$$

and the variance

$$\sigma_t^2 = \text{var}(Z_0 | T = t; \Omega_j) = \sigma^2 \rho_0, \tag{5}$$

where $x'_0 = x'(s_0), \alpha'_0 = c'_0 C_n^{-1}, \rho_0 = (1 - c'_0 C_n^{-1} c_0)$.

Under the assumption of complete parametric certainty of populations and for known finite non-negative losses $\{L(i, j), i, j = 1, 2\}$, BR minimizing the risk of classification is associated with the spatial

discriminant function (SDF) formed by a log ratio of conditional likelihoods (McLachlan, 2004).

$$W_t(Z_0, \Psi) = \left(Z_0 - \frac{1}{2}(\mu_{1t}^0 + \mu_{2t}^0) \right) \times \left(\mu_{1t}^0 - \mu_{2t}^0 \right) / \sigma_t^2 + \gamma, \quad (6)$$

where $\gamma = \ln(\pi_1^*/\pi_2^*)$ and $\Psi = (\beta', \sigma^2)'$.

Here $\pi_j^* = \pi_j(L(j, 3-j) - L(j, j))$, $j = 1, 2$, where π_1, π_2 ($\pi_1 + \pi_2 = 1$) are prior probabilities of the populations Ω_1 and Ω_2 , respectively.

Note, that in the present paper we implemented the following values of the prior probabilities $\hat{\pi}_j = n_j / n$, $j = 1, 2$.

So the classification procedure based on the SDF allocates the observation in the following way:

Classify the observation Z_0 given $T = t$ to the population Ω_1 if $W_t(Z_0, \Psi) \geq 0$, and to the population Ω_2 , otherwise.

Definition 1. The risk for the classification rule based on the SDF $W_t(Z_0; \Psi)$ is defined as

$$R_0 = \sum_{i=1}^2 \sum_{j=1}^2 \pi_i L(i, j) P_{ij}, \quad (7)$$

where, for $i, j = 1, 2$, $P_{ij} = P_u((-1)^j W_t(Z_0; \Psi) < 0)$.

Here, for $i = 1, 2$, the probability measure P_u is based on the conditional distribution of Z_0 given $T = t$, Ω_i specified in (4), (5).

Note that under the condition (3), the squared Mahalanobis distance between Ω_1 and Ω_2 in the location s_0 based on marginal distributions of Z_0 and the same squared Mahalanobis distance given $T = t$ are specified by $\Delta^2 = (\mu_1^0 - \mu_2^0)^2 / \sigma^2$ and $\Delta_0^2 = (\mu_{1t}^0 - \mu_{2t}^0)^2 / \sigma_t^2$, respectively.

From (4), (5) it follows that $\Delta_0^2 = \Delta^2 / \rho_0$. In the population Ω_j , the conditional distribution of $W_t(Z_0; \Psi)$ given $T = t$ is the normal distribution with the mean $E_j(W_t(Z_0; \Psi)) = (-1)^{j+1} \Delta_0^2 / 2 + \gamma$ and the variance $Var_j(W_t(Z_0; \Psi)) = \Delta_0^2$, $j = 1, 2$.

By using the properties of normal distributions we obtain

$$R_0 = \sum_{j=1}^2 \left(\pi_j^* \Phi(-\Delta_0/2 + (-1)^j \gamma/\Delta_0 + \pi_j L(j, j)) \right),$$

where $\Phi(\cdot)$ is the standard normal distribution function.

In practical applications not all statistical parameters of populations are known. In such cases the estimators of the unknown parameters can be found from the training sample. When the estimators of the unknown parameters are plugged into the SDF, the plug-in SDF (PSDF) is obtained. In this paper we assume that the true values of the parameters β and σ^2 are unknown (complete parametric uncertainty).

Let $\hat{\beta}, \hat{\sigma}^2$ be the estimators of the corresponding parameters from the training sample.

Put $\hat{\Psi} = (\hat{\beta}', \hat{\sigma}^2)'$. After replacing the parameters with their estimates in (6) the PSDF gets the following form

$$W_t(Z_0; \hat{\Psi}) = \left(Z_0 - \alpha'_0(t - X\hat{\beta}) - \frac{1}{2}x'_0 H \hat{\beta} \right) \times \left(x'_0 G \hat{\beta} \right) / (\hat{\sigma}^2 \rho_0) + \gamma, \quad (8)$$

with $H = (I_q, I_q)$ and $G = (I_q, -I_q)$, where I_q denotes the identity matrix of the order q .

Definition 2. The expectation of the ACR with respect to the distribution of T is called the expected risk (ER).

Recall that the actual risk incurred by the PSDF is obtained by replacing Ψ by the ML estimator $\hat{\Psi}$ in (6) (Ducinkas and Dreiziene, 2011).

Lemma 1. The actual risk for $W_t(Z_0; \hat{\Psi})$ specified in (8) is

$$R_A(\hat{\Psi}) = \sum_{j=1}^2 \left(\pi_j^* \Phi(\hat{Q}_j) + \pi_j L(j, j) \right), \quad (9)$$

where

$$\hat{Q}_j = (-1)^j \left((a_j - \hat{b}) \operatorname{sgn}(x'_0 G \hat{\beta}) / \sigma_0 + \hat{\sigma}_0^2 \gamma / \left(\sigma_0 |x'_0 G \hat{\beta}| \right) \right)$$

and $a_j = x'_0 \beta_j + \alpha'_0(t - X\beta)$,

$\hat{b} = \alpha'_0(t - X\hat{\beta} + x'_0 H \hat{\beta} / 2)$ for $j = 1, 2$.

The proof of the lemma is presented in the appendix.

The ACR is useful in providing a guide to the performance of the plug-in classification rule when it actually formed from the training sample. The ER is the performance measure to the PBDF similar as the mean squared prediction error (MSPE) is the performance measure to the plug-in kriging predictor (Diggle, Ribeiro and Christensen, 2002). The estimators of the MSPE are used for the spatial sampling design criterion for the prediction

(Zimmerman, 2006; Zhu and Zhang, 2006). These facts strengthen the motivation for the deriving of the estimators of the ER associated with the PSDF.

In this paper we propose an empirical estimator of the ER incurred by the rule based on the proposed PSDF. The following steps are performed to construct this estimator:

1. Simulate M training sample T realizations according to the model specified in (3).
2. For each simulated realization of $T = t_i$, $i = \overline{1, M}$ compute the appropriate estimates $\hat{\beta}_{(i)}$ and $\hat{\sigma}_{(i)}^2$, respectively. Set $\hat{\Psi}'_{(i)} = (\hat{\beta}'_{(i)}, \hat{\sigma}_{(i)}^2)$
3. By using (9) compute the empirical estimator of the ER

$$\hat{R}^{(\kappa)} = \sum_{i=1}^M R_A(\hat{\Psi}_{(i)}) / M, \quad (10)$$

where κ - denotes the abbreviation of the estimator type, i.e. takes the values BA or ML.

3 THE ESTIMATORS OF PARAMETERS

It is known, that the ML estimators of β and σ^2 based on T are $\hat{\beta}_{ML} = (X^T R^{-1} X)^{-1} X^T R^{-1} T$, $\hat{\sigma}_{ML}^2 = (T - X \hat{\beta}_{ML})' R^{-1} (T - X \hat{\beta}_{ML}) / n$.

Using the properties of the multivariate Gaussian distribution is easy to prove that $\hat{\beta}_{ML} \sim N_{2q}(\beta, \Sigma_\beta)$, $\Sigma_\beta = \sigma^2 (X^T R^{-1} X)^{-1}$, $\hat{\sigma}_{ML}^2 \sim \sigma^2 \chi_{n-2q}^2 / (n-2q)$.

The ML estimator of β and the bias adjusted ML estimator of σ^2 are used in the PLDF, i.e. $\hat{\beta} = \hat{\beta}_{ML}$, $\hat{\sigma}^2 = \hat{\sigma}_{ML}^2 n / (n-2q)$ (Ducinkas, 2009).

In the Bayesian approach the likelihood is given by $T | \beta, \sigma^2 \sim N_n(\beta, \sigma^2 R)$. The conjugate prior are chosen for the parameters so $p(\beta, \sigma^2) = p(\beta | \sigma^2) p(\sigma^2)$, where $p(\beta | \sigma^2) \sim N_{2q}(\beta^{(0)}, \sigma^2 \Sigma^{(0)})$ is the Gaussian prior for β conditional on σ^2 , $p(\sigma^2) \sim IG(u^{(0)}, v^{(0)})$ - the prior density for σ^2 , where $u^{(0)}, v^{(0)} > 0$. So the conjugate prior is the Normal-Inverse Gamma (NIG) and denotes as

$NIG(\beta^{(0)}, \Sigma^{(0)}, u^{(0)}, v^{(0)})$. Combining the prior with the likelihood gives a joint Normal-Inverse Gamma posterior (Diggle, Ribeiro and Christensen, 2002):

$$p(\beta, \sigma^2 | T) = \frac{p(\beta, \sigma^2) p(T | \beta, \sigma^2)}{p(T)} = NIG(\beta^{(0)}, \Sigma^{(0)}, u^{(0)}, v^{(0)})$$

where

$$\beta^{(1)} = \left(X^T R^{-1} X + (\Sigma^{(0)})^{-1} \right)^{-1} \left(X^T R^{-1} T + (\Sigma^{(0)})^{-1} \beta^{(0)} \right)$$

$$\Sigma^{(1)} = \left(X^T R^{-1} X + (\Sigma^{(0)})^{-1} \right)^{-1}, \quad u^{(1)} = u^{(0)} + \frac{n}{2},$$

$$v^{(1)} = v^{(0)} + \frac{1}{2} \left((\beta^{(0)})' (\Sigma^{(0)})^{-1} \beta^{(0)} + T' R^{-1} T - (\beta^{(0)})' (\Sigma^{(1)})^{-1} \beta^{(0)} \right).$$

The marginal posterior for β on integrating out σ^2 is a multivariate Student distribution $p(\beta | T) \sim t_p(\beta^{(1)}, \Sigma^*)$, where $p = 2u^{(1)}$, $\Sigma^* = (v^{(1)} / u^{(1)}) \Sigma^{(1)}$.

The marginal posterior for σ^2 is $p(\sigma^2 | T) \sim IG(u^{(1)}, v^{(1)})$, where $IG(\cdot, \cdot)$.

So the BA of β and σ^2 are $\hat{\beta}_{BA} = \beta^{(1)}$, $\hat{\sigma}_{BA}^2 = v^{(1)} / (u^{(1)} - 1)$, respectively.

4 EXAMPLE AND DISCUSSIONS

A numerical example is considered to investigate the influence of the parameter estimation methods to the proposed empirical estimator of the ER in the finite (even small) training sample case. With an insignificant loss of generality a case with $L(i, j) = 1 - \delta_{ij}$, $i, j = 1, 2$ is considered.

In this example, the observations are assumed to arise from the stationary Gaussian random field with the constant mean and the exponential correlation function specified by $c(h) = \exp\{-|h|/\phi\}$ is considered. Here ϕ denotes the range parameter. Assume that D is a regular 2-dimensional lattice with unit spacing. Consider the case of $s_0 = (1, 1)$ and the fixed set of training locations S_8 containing 8 second-order neighbours of s_0 . This example also illustrates the comparison of two different TLC.

Consider two TLC ξ_1 , ξ_2 for S_8 specified by $\xi_1 = \{S^{(1)} = \{(0, 2), (1, 2), (2, 2), (2, 1)\}, S^{(2)} = \{(2, 0), (1, 0), (0, 0), (0, 1)\}\}$.

$\xi_2 = \{S^{(1)} = \{(0,2), (1,2), (2,2), (2,1), (2,0)\}, S^{(2)} = \{(1,0), (0,0), (0,1)\}\}$. These are illustrated in Figure 1.

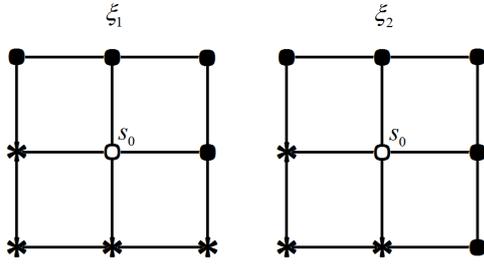


Figure 1: Two different TLC with $S^{(1)}$ and $S^{(2)}$ points marked as dots and asterisks.

By the definition, variable d represents the asymmetry of the TLC with respect to the location s_0 . It is easy to obtain that $d = 0$ and $n_1 = 4$ for ξ_1 and $d = 2\sqrt{2}$ with $n_1 = 5$ for ξ_2 .

So we can conclude that ξ_1 is the symmetric TLC and ξ_2 is the asymmetric one. Denote by $\hat{R}^{(ML)}$ and $\hat{R}^{(BA)}$ the empirical estimators of the ER given in (10) with the implemented ML and BA parameter estimators, respectively.

The values of the empirical estimators of $\hat{R}^{(*)}$ are presented for ξ_1 with $\pi_1 = 0,5$, and for ξ_2 with $\pi_1 = \frac{5}{8}$ in Table 1. Two cases of the simulated realizations of the training sample are selected here, i.e. $M = 10^2$ and $M = 10^4$. Analyzing the figures of the number of simulated realizations of the training sample in Table 1 we see that for all

$\Delta = (0.5, 0.7, \dots, 1.9)$ values $\hat{R}^{(ML)} \geq \hat{R}^{(BA)}$. So we can conclude that the BA case has an advantage against the ML case by the ER minimum criterion. The quantitative comparison of the two cases of the parameter estimators is also done by the values of the index $\eta = \hat{R}^{(ML)} / \hat{R}^{(BA)}$. The values of this index are shown in Figure 2 for $M = 10^4$ and $\xi_l, l = 1, 2$.

It is easy to see from Figure 2 that for both the TLC $\eta \geq 1$ and the values of this index increases when Δ increases. The same situation is for both TLC considered.

5 CONCLUSIONS

In this paper, the comparison of two approaches to parameter estimation is done based on the values of the ER incurred by the classification rule based on the PSDF.

The proposed optimality criterion is based on the derived formula of the actual risk.

The simulation experiment shows the advantage of Bayesian estimation approach against the frequentist (ML) approach. This advantage is greater for strongly separated populations (larger values of Δ) than for the close populations. These conclusions are valid to the symmetric training labels configuration as well to the asymmetric one. Hence the results of this paper give us strong arguments to expect that Bayesian estimators of spatial population parameters could be effectively used in spatial Gaussian data classification incurred by plug-in Bayes rules.

Table 1: Values of $R_1^{(*)}$ for the different estimators and the TLC ξ_1 and ξ_2 .

$\Delta \setminus M$	ξ_1				ξ_2			
	ML	BA	ML	BA	ML	BA	ML	BA
	10^2		10^4		10^2		10^4	
0.5	0.352	0.335	0.365	0.349	0.359	0.340	0.328	0.315
0.7	0.275	0.262	0.277	0.256	0.310	0.280	0.257	0.239
0.9	0.195	0.177	0.201	0.180	0.196	0.182	0.189	0.171
1.1	0.143	0.123	0.140	0.122	0.143	0.119	0.136	0.118
1.3	0.094	0.083	0.098	0.083	0.097	0.085	0.097	0.081
1.5	0.067	0.053	0.067	0.055	0.061	0.051	0.066	0.053
1.7	0.049	0.035	0.047	0.036	0.044	0.034	0.047	0.035
1.9	0.034	0.024	0.033	0.024	0.033	0.024	0.032	0.023

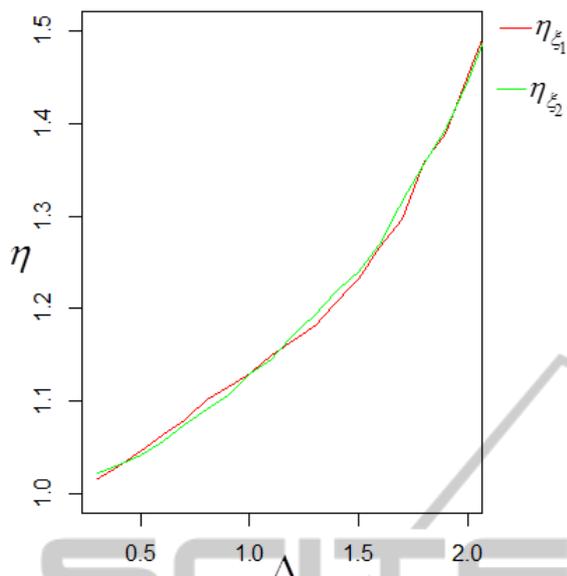


Figure 2: Values of η for different TLC.

REFERENCES

- Anderson, T. W., 2003. *An Introduction to Multivariate Statistical Analysis*, Wiley. New York.
- Batsidis, A. and Zografos, K., 2011. Errors of misclassification in discrimination of dimensional coherent elliptic random field observations. *Statistica Neerlandica*, 65, p. 446-461.
- Cressie, N. A. C., 1993. *Statistics for spatial data*, Wiley. New York.
- Diggle, P. J., Ribeiro, P. J. and Christensen, O. F., 2002. An introduction to model-based geostatistics. *Lecture notes in statistics*, 173, p. 43-86.
- Ducinkas, K., 2009. Approximation of the expected error rate in classification of the Gaussian random field observations. *Statistics and Probability Letters*, 79, p. 138-144.
- Ducinkas, K., Dreiziene, L., 2011. Supervised classification of the scalar Gaussian random field observations under a deterministic spatial sampling design. *Austrian Journal of Statistics*. 40, No. 1, 2, p. 25-36.
- Ducinkas, K. Stabingiene, L., Stabingis, G., 2011. Image classification based on Bayes discriminant functions. *Procedia Environmental Sciences*, 7, p. 218-223.
- Kharin, Y., 1996. *Robustness in Statistical Pattern Recognition*, Kluwer Academic Publishers. Dordrecht.
- McLachlan, G. J., 2004. *Discriminant analysis and statistical pattern recognition*, Wiley. New York.
- Saltyte-Benth, J. and Ducinkas, K., 2005. Linear discriminant analysis of multivariate spatial-temporal regressions. *Scandinavian Journal of Statistics*, 32, p. 281 – 294.

- Shekhar, S., Schrater, P. R., Vatsavai, R. R., Wu, W. and Chawla, S., 2002. Spatial Contextual Classification and Prediction Models for Mining Geospatial Data. *IEEE Transactions on Multimedia*, 4, p. 174-188.
- Zhu, Z. and Zhang, H., 2006. Spatial sampling design under infill asymptotic framework. *Environmetrics*, 17, p. 323-337.
- Zimmerman, D. L., 2006. Optimal network design for spatial prediction, covariance parameter estimation, and empirical prediction. *Environmetrics*, 17, p. 635-652.

APPENDIX

Proof of Lemma. Recall that the actual risk (ACR) for PSDF $W_t(Z_0; \hat{\Psi})$ (Ducinkas and Dreiziene, 2011) is defined as

$$R(\hat{\Psi}) = \sum_{i=1}^2 \sum_{j=1}^2 \pi_i L(i, j) \hat{P}_{ij} \quad \text{where for } i, j = 1, 2,$$

$$P_{ij} = P_{it} \left((-1)^j W_t(Z_0; \hat{\Psi}) < 0 \right).$$

In the population Ω_j , it is easy to derive that the conditional distribution of $W_t(Z_0; \hat{\Psi})$ given $T = t$ is normal distribution with the mean $E_j(W_t(Z_0; \hat{\Psi})) = (a_j - \hat{b})' x_0' G \hat{\beta} / \hat{\sigma}_0^2 + \gamma$ (14) and the variance $Var_j(W_t(Z_0; \hat{\Psi})) = (x_0' G \hat{\beta})^2 \sigma_0^2 / \hat{\sigma}_0^4$, $j = 1, 2$.

Then by using the properties of the normal distribution and we complete the proof of lemma 1.