

Ability of Covariance Matrix in Bi-Response Multi-Predictor Penalized Spline Model Through Longitudinal Data Simulation

Anna Islamiyati¹, Fatmawati², Nur Chamidah³

¹Department of Mathematics, Hasanuddin University, Makassar 90245, Indonesia

¹correspondence author: annaislamiyati@unhas.ac.id

^{2,3} Department of Mathematics, Airlangga University, Surabaya 60115, Indonesia

²fatmawati@fst.unair.ac.id, ³nur-c@fst.unair.ac.id

Abstract: Bi-response longitudinal data is assumed to have a correlation between responses and observations on the same subject. This causes a correlation between errors. To overcome this problem, we can use a penalized spline model that involves weighting. In this study, the weight used is the covariance matrix. Based on longitudinal data simulations that contain two responses and three predictors, a small GCV value is obtained from the penalized spline model that involves the covariance matrix. This value is compared with the penalized spline model without the covariance matrix. This shows that we need to involve a covariance matrix in estimating bi-response multi-predictor longitudinal data with a penalized spline model.

Keywords: bi-response, covariance matrix, GCV, penalized spline.

1. INTRODUCTION

One of the problems in statistical modeling is the correlation between errors. In this study, we used a non-parametric penalized spline regression model involving a covariance matrix to overcome this problem. The matrix of covariance has been reviewed by several researchers, including Fessler [1] and Wahba [2] which assume a known covariance matrix. Furthermore, Wang [3] has estimated the matrix of covariance through smoothing spline estimators. In the case of more than one response, the covariance matrix has been used by several researchers. For cross section data, Soo and Bates [4] in multi-response spline regression. Budiantara et., al. [5] used weighted spline estimator. Lestari et., al. [6,7,8] in natural smoothing spline. Chamidah et., al. [9] in local polynomial estimators and kernels. Chamidah and Eridani [10] in the P-spline estimator. Chamidah and Lestari [11] in spline smoothing estimator. Aydin et., al. [12] used modified spline estimators with right-censored data.

For longitudinal data with a penalized spline estimator, Liang and Xiao [13] with a varying coefficient model. Chen and Wang [14] and Heckman et., al. [15] with a mix effect model. In cases of more than one response, Islamiyati, et., al. [16] has estimated the function of the goodness of fit and Islamiyati, et., al. [17] have estimated the covariance matrix through the expectation method. Next, we show the ability of the covariance matrix through simulation studies. Longitudinal data simulations contain two responses and three predictors through quadratic spline functions with several sample numbers and correlations.

2. COVARIANCE MATRIX IN THE BI-RESPONSE MULTI-PREDICTOR PENALIZED SPLINE REGRESSION MODEL

The covariance matrix is obtained from the error of the penalized spline bi-response multi-predictor model without

weighting. Non-parametric bi-response multi-predictor regression model without weighting where $i = 1, 2, \dots, n$, $j = 1, 2, \dots, m_i$ as follows:

$$y_i^{(0)} = f^{(0)}(t_{ij1}) + f^{(0)}(t_{ij2}) + f^{(0)}(t_{ij3}) + \varepsilon_i^{(0)}, \quad (1)$$

where

$$y_i^{(0)} = \begin{pmatrix} y_{i1}^{(0)} \\ y_{i2}^{(0)} \end{pmatrix}, f^{(0)}(t_{ij1}) = \begin{pmatrix} f_1^{(0)}(t_{ij1}) \\ f_2^{(0)}(t_{ij1}) \end{pmatrix}, f^{(0)}(t_{ij2}) = \begin{pmatrix} f_1^{(0)}(t_{ij2}) \\ f_2^{(0)}(t_{ij2}) \end{pmatrix}$$

$$f^{(0)}(t_{ij3}) = \begin{pmatrix} f_1^{(0)}(t_{ij3}) \\ f_2^{(0)}(t_{ij3}) \end{pmatrix} \text{ and } \varepsilon_i^{(0)} = \begin{pmatrix} \varepsilon_{i1}^{(0)} \\ \varepsilon_{i2}^{(0)} \end{pmatrix}.$$

The error of the penalized spline model is different for each subject of observation from longitudinal data. Therefore, the vector $y_i^{(0)}$ in equation (1) can be formed into:

$$\left. \begin{aligned} y_{i1}^{(0)} &= f^{(0)}(t_{i11}) + f^{(0)}(t_{i12}) + f^{(0)}(t_{i13}) + \varepsilon_{i1}^{(0)} \\ y_{i2}^{(0)} &= f^{(0)}(t_{i21}) + f^{(0)}(t_{i22}) + f^{(0)}(t_{i23}) + \varepsilon_{i2}^{(0)} \\ &\vdots \\ y_{in}^{(0)} &= f^{(0)}(t_{in1}) + f^{(0)}(t_{in2}) + f^{(0)}(t_{in3}) + \varepsilon_{in}^{(0)} \end{aligned} \right\}. \quad (2)$$

The response function in equation (2) can be expressed in the vector form from the 1st response and the 2nd response is as follows:

$$y_i^{(0)} = f^{(0)}(t_{ij1}) + f^{(0)}(t_{ij2}) + f^{(0)}(t_{ij3}) + \varepsilon_i^{(0)}, \quad (3)$$

where $y_i^{(0)} = (y_{i1}^{(0)}, y_{i2}^{(0)})^T$ and $\varepsilon_i^{(0)} = (\varepsilon_{i1}^{(0)}, \varepsilon_{i2}^{(0)})^T$.

Model (3) is estimated by penalized least square, that is:

$$PLS = \mathbf{y}_i^{(0)T} \mathbf{y}_i^{(0)} - 2\beta^{(0)T} \mathbf{X}_i^{(0)T} \mathbf{y}_i^{(0)} + \beta^{(0)T} \mathbf{X}_i^{(0)T} \mathbf{X}_i^{(0)} \beta^{(0)} + \beta^{(0)T} \mathbf{D}_{\hat{\lambda}_i^{(0)}} \beta^{(0)}. \quad (4)$$

$$\text{where } \mathbf{y}_i^{(0)} = \begin{pmatrix} y_{1,i}^{(0)} \\ y_{2,i}^{(0)} \end{pmatrix}, \beta^{(0)} = \begin{pmatrix} \beta_1^{(0)} \\ \beta_2^{(0)} \end{pmatrix}, \mathbf{X}_i^{(0)} = \begin{pmatrix} \mathbf{X}_{1,i}^{(0)} & \mathbf{0} \\ \mathbf{0} & \mathbf{X}_{2,i}^{(0)} \end{pmatrix},$$

$$\mathbf{D}_{\hat{\lambda}_i^{(0)}} = \begin{pmatrix} \mathbf{D}_{\hat{\lambda}_1^{(0)}} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_{\hat{\lambda}_2^{(0)}} \end{pmatrix}, \mathbf{y}_{1,i}^{(0)} = (y_{1,1}^{(0)}, y_{1,2}^{(0)}, \dots, y_{1,n}^{(0)})^T \text{ and}$$

$\mathbf{y}_{2,i}^{(0)} = (y_{2,1}^{(0)}, y_{2,2}^{(0)}, \dots, y_{2,n}^{(0)})^T$. The vector $\beta^{(0)}$ is a vector of nonparametric regression coefficients in the 1st response and the 2nd response.

Equation (4) is derived from $\beta^{(0)}$ and obtained:

$$\hat{\beta}^{(0)} = \begin{bmatrix} \hat{\beta}_1^{(0)} \\ \hat{\beta}_2^{(0)} \end{bmatrix} = \begin{bmatrix} \left(\mathbf{X}_{1,i}^{(0)T} \mathbf{X}_{1,i}^{(0)} + \mathbf{D}_{\hat{\lambda}_1^{(0)}} \right)^{-1} \mathbf{X}_{1,i}^{(0)T} \mathbf{y}_{1,i}^{(0)} \\ \left(\mathbf{X}_{2,i}^{(0)T} \mathbf{X}_{2,i}^{(0)} + \mathbf{D}_{\hat{\lambda}_2^{(0)}} \right)^{-1} \mathbf{X}_{2,i}^{(0)T} \mathbf{y}_{2,i}^{(0)} \end{bmatrix}, \quad (5)$$

where $\hat{\lambda}_1^{(0)}, \hat{\lambda}_2^{(0)}$ is the smoothing parameter for the 1st and 2nd response. Based on equation (5), the estimation of nonparametric bi-response regression function in longitudinal data through an unweighted penalized spline estimator is:

$$\hat{f}^{(0)} = \begin{pmatrix} \mathbf{X}_{1,i}^{(0)} \hat{\beta}_1^{(0)} \\ \mathbf{X}_{2,i}^{(0)} \hat{\beta}_2^{(0)} \end{pmatrix} = \begin{pmatrix} \mathbf{X}_{1,i}^{(0)} (\mathbf{X}_{1,i}^{(0)T} \mathbf{X}_{1,i}^{(0)} + \mathbf{D}_{\hat{\lambda}_1^{(0)}})^{-1} \mathbf{X}_{1,i}^{(0)T} \mathbf{y}_{1,i}^{(0)} \\ \mathbf{X}_{2,i}^{(0)} (\mathbf{X}_{2,i}^{(0)T} \mathbf{X}_{2,i}^{(0)} + \mathbf{D}_{\hat{\lambda}_2^{(0)}})^{-1} \mathbf{X}_{2,i}^{(0)T} \mathbf{y}_{2,i}^{(0)} \end{pmatrix}. \quad (6)$$

If the matrix $\mathbf{A}_1^{(0)} = \mathbf{X}_{1,i}^{(0)} (\mathbf{X}_{1,i}^{(0)T} \mathbf{X}_{1,i}^{(0)} + \mathbf{D}_{\hat{\lambda}_1^{(0)}})^{-1} \mathbf{X}_{1,i}^{(0)T}$ and $\mathbf{A}_2^{(0)} = \mathbf{X}_{2,i}^{(0)} (\mathbf{X}_{2,i}^{(0)T} \mathbf{X}_{2,i}^{(0)} + \mathbf{D}_{\hat{\lambda}_2^{(0)}})^{-1} \mathbf{X}_{2,i}^{(0)T}$, then the equation (6) become:

$$\hat{f}^{(0)} = \begin{pmatrix} \mathbf{A}_1^{(0)} \mathbf{y}_{1,i}^{(0)} \\ \mathbf{A}_2^{(0)} \mathbf{y}_{2,i}^{(0)} \end{pmatrix} = \begin{pmatrix} E(\mathbf{y}_{1,i}^{(0)}) \\ E(\mathbf{y}_{2,i}^{(0)}) \end{pmatrix}. \quad (7)$$

Furthermore, the error from the model is obtained as follows:

$$\varepsilon_i^{(0)} = \mathbf{y}_i^{(0)} - \hat{f}^{(0)} \quad (8)$$

Based on the error value obtained from equation (8), the variance covariance matrix can be estimated as follows:

$$\hat{\Omega} = \begin{bmatrix} \hat{\Sigma}_{1,1} & \mathbf{0} & \dots & \mathbf{0} & \hat{\Sigma}_{12,1} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \hat{\Sigma}_{1,2} & \dots & \mathbf{0} & \mathbf{0} & \hat{\Sigma}_{12,2} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \hat{\Sigma}_{1,n} & \mathbf{0} & \mathbf{0} & \dots & \hat{\Sigma}_{12,n} \\ \hat{\Sigma}_{21,1} & \mathbf{0} & \dots & \mathbf{0} & \hat{\Sigma}_{2,1} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \hat{\Sigma}_{2,2} & \dots & \mathbf{0} & \mathbf{0} & \hat{\Sigma}_{2,2} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \hat{\Sigma}_{2,n} & \mathbf{0} & \mathbf{0} & \dots & \hat{\Sigma}_{2,n} \end{bmatrix}, \quad (9)$$

where

$$\hat{\Sigma}_{1,i} = \text{diag}(\hat{\sigma}_{1,i}^2, \hat{\sigma}_{1,i}^2, \dots, \hat{\sigma}_{1,i}^2), \hat{\Sigma}_{2,i} = \text{diag}(\hat{\sigma}_{2,i}^2, \hat{\sigma}_{2,i}^2, \dots, \hat{\sigma}_{2,i}^2), \\ \hat{\Sigma}_{12,i} = \text{diag}(\hat{\sigma}_{12,i}, \hat{\sigma}_{12,i}, \dots, \hat{\sigma}_{12,i}), \hat{\Sigma}_{21,i} = \text{diag}(\hat{\sigma}_{21,i}, \hat{\sigma}_{21,i}, \dots, \hat{\sigma}_{21,i}) \\ \hat{\Sigma}_{12,i} = \hat{\Sigma}_{21,i}.$$

Furthermore, the variance in the first response is

$$\hat{\sigma}_{1,i}^2 = \frac{\mathbf{y}_{1,i}^{(0)T} \mathbf{P}_{11} \mathbf{y}_{1,i}^{(0)} - \mathbf{y}_{1,i}^{(0)T} \mathbf{A}_1^{(0)T} \mathbf{P}_{11} \mathbf{A}_1^{(0)} \mathbf{y}_{1,i}^{(0)}}{\text{tr}(\mathbf{P}_{11})},$$

$\mathbf{P}_{11} = [\mathbf{I} - \mathbf{A}_1^{(0)}]^T [\mathbf{I} - \mathbf{A}_1^{(0)}]$, $\mathbf{A}_1^{(0)}$ is a smoothing matrix in the first response, $\mathbf{y}_{1,i}^{(0)}$ is the first response vector in the i subject. For variance in the second response with the i subject is $\hat{\sigma}_{2,i}^2 = \frac{\mathbf{y}_{2,i}^{(0)T} \mathbf{P}_{22} \mathbf{y}_{2,i}^{(0)} - \mathbf{y}_{2,i}^{(0)T} \mathbf{A}_2^{(0)T} \mathbf{P}_{22} \mathbf{A}_2^{(0)} \mathbf{y}_{2,i}^{(0)}}{\text{tr}(\mathbf{P}_{22})}$, where

$\mathbf{P}_{22} = [\mathbf{I} - \mathbf{A}_2^{(0)}]^T [\mathbf{I} - \mathbf{A}_2^{(0)}]$, $\mathbf{A}_2^{(0)}$ is a smoothing matrix in the second response, $\mathbf{y}_{2,i}^{(0)}$ is the second response vector in the i subject. The variance in the second response is $\hat{\sigma}_{12,i} = \frac{\mathbf{y}_{1,i}^{(0)T} \mathbf{P}_{12} \mathbf{y}_{2,i}^{(0)} - \mathbf{y}_{1,i}^{(0)T} \mathbf{A}_1^{(0)T} \mathbf{P}_{12} \mathbf{A}_2^{(0)} \mathbf{y}_{2,i}^{(0)}}{\text{tr}(\mathbf{P}_{12})} = \hat{\sigma}_{21,i}$.

The equation (9) can be simplified to be:

$$\hat{\Omega} = \begin{pmatrix} \hat{\Sigma}_1 & \hat{\Sigma}_{12} \\ \hat{\Sigma}_{21} & \hat{\Sigma}_2 \end{pmatrix}, \quad (5)$$

where

$$\hat{\Sigma}_{1,i} = \text{diag}(\hat{\Sigma}_{1,1}, \hat{\Sigma}_{1,2}, \dots, \hat{\Sigma}_{1,n}), \hat{\Sigma}_{2,i} = \text{diag}(\hat{\Sigma}_{2,1}, \hat{\Sigma}_{2,2}, \dots, \hat{\Sigma}_{2,n}) \text{ and} \\ \hat{\Sigma}_{12,i} = \text{diag}(\hat{\Sigma}_{12,1}, \hat{\Sigma}_{12,2}, \dots, \hat{\Sigma}_{12,n}), \hat{\Sigma}_{21,i} = \hat{\Sigma}_{12,i}.$$

3. A SIMULATION STUDY

The ability of the covariance matrix is shown through simulation studies with the following functions:

$$y_1 = f_1(t_{i1}) + f_1(t_{i2}) + f_1(t_{i3}) + \varepsilon_1, \varepsilon_1 \sim N(0, \Sigma_1),$$

$$y_2 = f_2(t_{i1}) + f_2(t_{i2}) + f_2(t_{i3}) + \varepsilon_2, \varepsilon_2 \sim N(0, \Sigma_2),$$

where $f_1(t_{ij1}) = 0,5 + 2t_{ij1} - 2t_{ij1}^2$, $f_1(t_{ij2}) = -5 - 2t_{ij2} + 4,5t_{ij2}^2$,
 $f_1(t_{ij3}) = -2 + 0,5t_{ij3} + 2t_{ij3}^2$, $f_2(t_{ij1}) = -5 - 2t_{ij1} + 4,5t_{ij1}^2$,
 $f_2(t_{ij2}) = 0,5 + 2t_{ij2} - 2t_{ij2}^2$, and $f_2(t_{ij3}) = -2 - 5t_{ij3} - 2t_{ij3}^2$.

Simulations carried out on several correlation coefficient values with error variance is Σ . The error value of the penalized spline bi-response multi-predictor model at $r = -0,9; -0,6; 0,7; 0,8$ are shown in Figure 1.

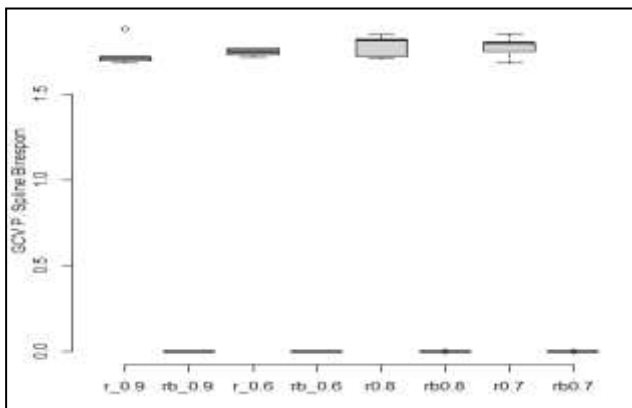


Fig 1. GCV value of the model with the covariance matrix (lower) and without the covariance matrix (upper)

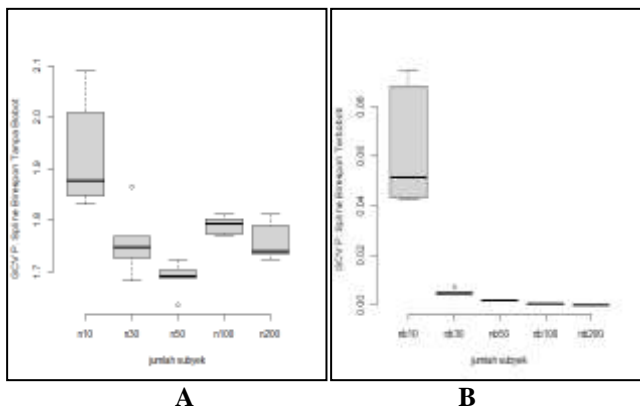


Fig 2. GCV values on a number of subjects $n = 10, 30, 50, 100, 200$, A: without weighting; B: by weighting the covariance matrix.

Figure 1 and Figure 2 shows the value of GCV involving the covariance matrix in the bi-response multi-predictor penalized spline regression model approaching 0, smaller than the GCV value of the un-weighted model. This means that we can use the covariance matrix as the weighting of the error without weighting penalized spline regression models.

4. CONCLUSION

The matrix of covariance as weighting matrix is obtained from an error from the non-parametric regression model of the penalized spline bi-response without weighting. It is quite accurate in estimating the model of the penalized spline bi-response multi-predictors in longitudinal data. The accuracy of the model is shown through GCV values from simulation studies.

5. REFERENCES

- [1] Fessler, J.A. (1991). Nonparametric fixed-interval smoothing with vector splines IEEE Trans. Signal Processing, 39:852-859.
- [2] Wahba, G. (1992). Multivariate functional and operator estimation, based on smoothing splines and reproducing kernels Nonlinear Modelling and Forecasting SFI Studies in the Science of Complexity (Edited by M. Casdagli and S. Eubank), Vol **XII**:95-112, Addison Wesley Reading.
- [3] Wang, Y. (1998). Smoothing spline models with correlated random errors J. Amer. Statist. Assoc. 93:341-348.
- [4] Soo, Y.W., & Bates, D.M (1996). Multi-response Spline Regression *Computational Statistics & Data Analysis*, 22:619-631
- [5] Budiantara, I. N., Lestari, B., & Islamiyati, A. (2009). Weighted spline estimator in heteroscedastic multi-response nonparametric regression for longitudinal data, Proceeding IndoMS International Conference on Mathematics and Its Application (IICMA) 2009, Yogyakarta, Indonesia, 12-13 October 2009, pp. 921-934.
- [6] Lestari, B., Budiantara, I.N., Sunaryo, S., & Mashuri, M. (2012). Spline smoothing for multi-response nonparametric regression model in case of heteroscedasticity of variance. *Journal of Mathematics and Statistics*, 8(3):337-384.
- [7] Lestari, B., Anggraeni, D., & Saifuddin, T. (2017). Konstruksi dan estimasi matriks kovariansi dalam model regresi nonparametric multirespon berdasarkan estimator smoothing spline untuk beberapa kasus ukuran sampel. Proceeding of National Seminar on Mathematics and Its Application (SNMA 2017), Airlangga University, Surabaya-Indonesia, 238-242.
- [8] Lestari, B., Anggraeni, D., and Saifuddin, T. (2018). Estimation of covariance matrix based on spline estimator in homoscedastic multi-response nonparametric regression model in case of unbalance number of observation. Far East Journal of Mathematical Science (FJMS), in press.
- [9] Chamidah, N., Budiantara, I.N., Sunaryo, S., & Ismaini, Z. (2012). Designing of child growth chart based on multi response local polynomial modeling. *Journal of Mathematics and Statistics*. 8:342-347.

- [10] Chamidah, N., & Eridani. (2015). Designing of growth reference chart by using bi-response semi-parametric regression approach based on P-spline estimator. *International Journal of Applied Mathematics and Statistics*, 53(3):150-158.
- [11] Chamidah, N., & Lestari, B. (2016). Spline estimator in homoscedastic multi-response nonparametric regression model in case of unbalanced number of observations. *Far East Journal of Mathematical Sciences (FJMS)*, 100(9):1433-1453.
- [12] Aydin, D., Aydin, D., & Yilmaz, E. (2018). Modified spline regression based on randomly right-censored data: A comparative study, Volume 47, Issue 9, Communications in Statistics-Simulation and Computation.
- [13] Liang, H., & Xiao, Y. 2006. Penalized for Longitudinal Data with an Application in AIDS Studies, *Journal of Modern Applied Statistical Methods*. 73:13 – 22.
- [14] Chen, H., & Wang, Y. 2011. A Penalized Spline Approach to Functional Mixed Effects Model Analysis. *Biometrics*. 67:861-870.
- [15] Heckman, N., Lockhart R., & Nielsen J.D. 2013. Penalized Regression, Mixed Effect Models and Appropriate Modelling. *Electronic Journal of Statistics*, 7:1517-1552.
- [16] Islamiyati A., Fatmawati & Chamidah N. (2017). Fungsi Goodness of Fit dalam Kriteria Penalized Spline pada Estimasi Regresi Nonparametrik Birespon untuk Data Longitudinal. Proceeding of National Seminar on Mathematics and Its Applications (SNMA 2017), Airlangga University, Surabaya-Indonesia, 216-221.
- [17] Islamiyati A., Fatmawati and Chamidah N. (2018). Estimation of Covariance Matrix on Bi-Response Longitudinal Data Analysis with Penalized Spline Regression. *Journal of Physics: Conf. Series*, 979 012093, IOP Publishing. doi :10.1088/1742-6596/979/1/012093.