

# Estimation of Heteroskedastic Semiparametric Regression Curve Using Fourier Series Approach

Rahmawati Pane<sup>1\*</sup>, Sutarman<sup>1</sup>, Andi Tenri Ampa<sup>2</sup>

<sup>1</sup>Department of Mathematics, Universitas Sumatera Utara, Medan, 20155, Indonesia

<sup>2</sup>Department of Statistics, Institut Teknologi Sepuluh Nopember, Surabaya, 60111, Indonesia

**Abstract.** A heteroskedastic semiparametric regression model consists of two main components, i.e. parametric component and nonparametric component. The model assumes that any data  $(\underline{x}_i', t_i, y_i)$  follows  $y_i = \underline{x}_i' \underline{\beta} + f(t_i) + \sigma_i \varepsilon_i$ , where  $i = 1, 2, \dots, n$ ,  $\underline{x}_i' = (1, x_{i1}, x_{i2}, \dots, x_{ir})$  and  $t_i$  is the predictor variable. Parameter vector  $\underline{\beta} = (\beta_1, \beta_2, \dots, \beta_r)' \in \mathbb{R}^r$  is unknown and  $f(t_i)$  is also unknown and is assumed to be in interval of  $C[0, \pi]$ . Random error  $\varepsilon_i$  is independent on zero mean and variance  $\sigma^2$ . Estimation of the heteroskedastic semiparametric regression model was conducted to evaluate the parametric and nonparametric components. The nonparametric component  $f(t_i)$  regression was approximated by Fourier series  $F(t) = bt + \frac{1}{2} \alpha_0 + \sum_{k=1}^K \alpha_k \cos kt$ . The estimation was obtained by means of Weighted Penalized Least Square (WPLS):  $\min_{f \in C(0, \pi)} \left\{ n^{-1} (\underline{y} - X \underline{\beta} - \underline{f})' W^{-1} (\underline{y} - X \underline{\beta} - \underline{f}) + \lambda \int_0^{\pi} \frac{2}{\pi} [f''(t)]^2 dt \right\}$ . The WPLS solution provided nonparametric component  $\hat{f}_\lambda(t) = M(\lambda) \underline{y}^*$  for a matrix  $M(\lambda)$  and parametric component  $\hat{\underline{\beta}} = [X'T(\lambda)X]^{-1} X'T(\lambda) \underline{y}$ .

**Keyword:** Fourier Series, Heteroskedastic Semiparametric Regression, Bandwidth Parameter, WPLS

**Abstrak.** Suatu model regresi semiparametrik heteroskedastis terdiri atas dua komponen utama, yaitu komponen parametrik dan komponen nonparametrik. Model ini mengasumsikan bahwa sebarang data  $(\underline{x}_i', t_i, y_i)$  dengan  $y_i = \underline{x}_i' \underline{\beta} + f(t_i) + \sigma_i \varepsilon_i$ , untuk  $i = 1, 2, \dots, n$ ,  $\underline{x}_i' = (1, x_{i1}, x_{i2}, \dots, x_{ir})$  dan  $t_i$  merupakan variabel prediksi. Vektor parameter  $\underline{\beta} = (\beta_1, \beta_2, \dots, \beta_r)' \in \mathbb{R}^r$  tidak diketahui dan  $f(t_i)$  adalah fungsi yang tidak diketahui bentuknya dan diasumsikan termuat dalam ruang fungsi kontinu  $C[0, \pi]$ . Error random  $\varepsilon_i$  independen dengan mean nol dan variansi  $\sigma^2$ . Estimasi dari model regresi semiparametrik heteroskedastik dilakukan untuk mengevaluasi komponen parametrik dan nonparametrik. Kurva regresi komponen nonparametrik  $f(t_i)$  dihipotesiskan dengan deret Fourier  $F(t) = bt + \frac{1}{2} \alpha_0 + \sum_{k=1}^K \alpha_k \cos kt$ . Estimasi kurva regresi semiparametrik heteroskedastik diperoleh dari menyelesaikan optimasi Weighted Penalized Least Square (WPLS):  $\min_{f \in C(0, \pi)} \left\{ n^{-1} (\underline{y} - X \underline{\beta} - \underline{f})' W^{-1} (\underline{y} - X \underline{\beta} - \underline{f}) + \lambda \int_0^{\pi} \frac{2}{\pi} [f''(t)]^2 dt \right\}$ . Solusi dari WPLS di atas memberikan estimator komponen nonparametrik dalam bentuk  $\hat{f}_\lambda(t) = M(\lambda) \underline{y}^*$  untuk suatu matriks  $M(\lambda)$  dan komponen parametrik  $\hat{\underline{\beta}} = [X'T(\lambda)X]^{-1} X'T(\lambda) \underline{y}$ .

**Kata Kunci:** Deret Fourier, Regresi Semiparametrik Heteroskedastik, Parameter Bandwidth, WPLS

\*Corresponding author at: Department of Mathematics, Universitas Sumatera Utara, Medan, 20155, Indonesia

E-mail address: rahmawatipane@usu.ac.id

Received 5 December 2019 | Revised 14 January 2020 | Accepted 17 February 2020

## 1. Introduction

Regression analysis is an important method in statistics. The objectives are to investigate the relationship between the predictor variables and the response variables and to evaluate the contribution of the predictor variable  $x'_i = (1, x_{i1}, x_{i2}, \dots, x_{ir})$  and  $t_i$  to the response variable  $y_i$ . When the relationship of the variables is known, the regression is called parametric regression [1]. Otherwise, the regression is identified as nonparametric regression.

In many cases, the response variable has linear relationship with one of the predictor variables, but with other predictor variables, it remains unknown. If this happens, Wahba [2] suggested the use of semiparametric regression approach. Semiparametric regression is a combination of parametric with nonparametric regressions. According to [3] and [4], semiparametric regression model assumed that data  $(x'_i, t_i, y_i)$  follows model  $y_i = x'_i\beta + f(t_i) + \sigma_i\varepsilon_i$  where  $i = 1, 2, \dots, n$ , and  $x'_i$  and  $t_i$  are the predictor variables.  $x'_i\beta$  is the parametric component and  $f(t_i)$  is the nonparametric component. The semiparametric regression has been widely used with various smoothing approaches. Some researches applied Spline function [5], [2], Wavelet [6], Kernel [7], [8], [9], [10], or local polynomial on the nonparametric component.

In fact, some data have periodic properties which are in a good agreement with Fourier series function. The Fourier series is known as a trigonometric polynomial that is very “flexible” to effectively approximate complicated functions. The Fourier series is best applied to describe any function whose terms is sines and cosines. The Fourier series estimator, in general, is used when the data in a question is periodic and unknown [11], [12]. A nonparametric regression by means of Fourier series has been developed by Bilodeau [12] to represent periodic functions with the same trend and variance (homoskedastic). However, a serious problem arises when it is applied to a heteroskedastic function; a periodic function with different random error variance (heteroskedastic). Therefore, it is crucial to develop an estimator based on Fourier series which includes elements of the trend, and the variance of the error in heteroskedastic semiparametric regression model.

In this report, a Fourier series estimator was derived to estimate the parametric and nonparametric components of a heteroskedastic function by means of semiparametric regression. The result could be an alternative model to the data pattern recurring/ periodic with trends in semiparametric regression heteroskedastic.

## 2. Semiparametric Regression Model

A semiparametric regression model is a combination of parametric with nonparametric regression models. The model is employed when the relationship between the predictor variable and the response variable can be described as certain curves. Given that paired data  $(\underline{t}_i, \underline{x}'_i, y_i)$  and  $\underline{x}'_i = (x_{i1}, x_{i2}, \dots, x_{ir})$  is the predictor variables that the relationship with the response variable  $y_i$  is known, and  $\underline{t}_i = (t_{i1}, t_{i2}, \dots, t_{pi})'$  is the predictor variable that the relationship with the response variable is unknown. The relationship of the variables  $(\underline{t}_i, \underline{x}'_i, y_i)$  is assumed to follow a heteroskedastic semiparametric regression model:

$$y_i = \underline{x}'_i \underline{\beta} + f(t_i) + \sigma_i \varepsilon_i, \quad i = 1, 2, \dots, n \quad (1)$$

The parameter  $\underline{\beta} = (\beta_1, \beta_2, \dots, \beta_r)' \in \mathbb{R}^r$  is a vector parameter of the unknown parametric component with size of  $r \times 1$ . The random error  $\varepsilon_i, i = 1, 2, \dots, n$  is mutually independent with zero mean and variance  $\sigma^2$ .  $f(t_i)$  is the nonparametric component with unknown function and is assumed to be defined in  $C[0, \pi]$ , where  $C[0, \pi] = \{f, f \text{ is a continuous function in the interval of } [0, \pi]\}$ .

### 2.1. Estimation of Fourier Series-Based Heteroskedastic Semiparametric Regression Model

Estimating the semiparametric regression model by Fourier series is to define the parametric component  $\underline{x}'_i \underline{\beta}$  as well as the nonparametric component  $f(t_i)$ . The nonparametric component of the heteroskedastic semiparametric regression curve can be approached using Fourier series as follows:

$$F(t) = bt + \frac{1}{2} \alpha_0 + \sum_{k=1}^K \alpha_k \cos kt \quad (2)$$

The estimators for regression curve  $f$  and parameter  $\underline{\beta}$  can be obtained by solving Weighted Penalized Least Square (WPLS) optimization:

$$\min_{f \in C(0, \pi)} \left\{ n^{-1} (\underline{y} - X \underline{\beta} - \underline{f})' W^{-1} (\underline{y} - X \underline{\beta} - \underline{f}) + \lambda \int_0^\pi \frac{2}{\pi} [f''(t)]^2 dt \right\} \quad (3)$$

where the weighted matrix  $W = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)$ ,  $\underline{y} = (y_1, y_2, \dots, y_n)'$  and  $\underline{f}$  have size of  $(n \times 1)$  and  $\lambda$  is the bandwidth parameter. The WPLS optimization was executed by taking goodness of fit  $G(f)$  and penalty  $P(f)$  of the Eq. 3.

$$\text{Where} \quad G(\underline{f}) = n^{-1} (\underline{y} - X \underline{\beta} - \underline{f})' W^{-1} (\underline{y} - X \underline{\beta} - \underline{f}) \quad (4)$$

$$\text{and} \quad P(f) = \int_0^\pi \frac{2}{\pi} [f''(t)]^2 dt \quad (5)$$

the heteroskedastic semiparametric regression model from Eq. 1 can be rewritten in the form of:

$$\underline{y} = X \underline{\beta} + B(\underline{t}) \underline{\delta} + W^{1/2} \underline{\varepsilon}$$

Where

$$X = \begin{pmatrix} 1 & x_{11} & x_{21} & L & x_{r1} \\ 1 & x_{12} & x_{22} & L & x_{r2} \\ M & M & M & M & M \\ 1 & x_{1n} & x_{2n} & L & x_{rn} \end{pmatrix}, B(t) = \begin{pmatrix} t_1 & 1 & \cos t_1 & \cos 2t_1 & L & \cos Kt_1 \\ t_2 & 1 & \cos t_2 & \cos 2t_2 & L & \cos Kt_2 \\ M & M & M & M & M & M \\ t_n & 1 & \cos t_n & \cos 2t_n & L & \cos Kt_n \end{pmatrix},$$

$$W^{1/2} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n) \text{ and } \underline{\delta} = \left(b, \frac{1}{2}\alpha_0, \alpha_1, \alpha_2, \dots, \alpha_K\right)' \in \mathfrak{R}^{(K+2)}.$$

The goodness of fit from Eq. 4 can be defined as:

$$G(\underline{\delta}) = n^{-1}(\underline{y} - X\underline{\beta} - B(t)\underline{\delta})' W^{-1}(\underline{y} - X\underline{\beta} - B(t)\underline{\delta}) \quad (6)$$

Since  $F(t) = bt + \frac{1}{2}\alpha_0 + \sum_{k=1}^K \alpha_k \cos kt$ , then:

$$\int_0^\pi \frac{2}{\pi} [f''(t)]^2 dt = \int_0^\pi \frac{2}{\pi} \left[ bt + \frac{1}{2}\alpha_0 + \sum_{k=1}^K \alpha_k \cos kt \right]^2 dt.$$

In other expression, the last equation can be rewritten as follows:

$$\int_0^\pi \frac{2}{\pi} [f''(t)]^2 dt = \sum_{k=1}^K k^4 \alpha_k^2.$$

Therefore:

$$\lambda \int_0^\pi \frac{2}{\pi} [f''(t)]^2 dt = \lambda \underline{\delta}' \begin{pmatrix} 0 & 0 & 0 & L & 0 \\ 0 & 0 & 0 & L & 0 \\ 0 & 0 & 1^4 & L & 0 \\ M & M & M & 0 & M \\ 0 & 0 & 0 & L & K^4 \end{pmatrix} \underline{\delta} = \lambda \underline{\delta}' D \underline{\delta} \quad (7)$$

with matrix  $D = \text{diag}(0, 0, 1^4, 2^4, \dots, K^4)$ . If the goodness of fit (4) and penalty (5) are combined, then the WPLS optimization (3) can be expressed as:

$$\begin{aligned} \min_{\underline{\delta} \in \mathfrak{R}^{(K+2)}} \{G(\underline{\delta}) + \lambda P(\underline{\delta})\} \\ = \min_{\underline{\delta} \in \mathfrak{R}^{(K+2)}} \{n^{-1}(\underline{y} - X\underline{\beta} - B(t)\underline{\delta})' W^{-1}(\underline{y} - X\underline{\beta} - B(t)\underline{\delta}) + \lambda \underline{\delta}' D \underline{\delta}\} \\ = \min_{\underline{\delta} \in \mathfrak{R}^{(K+2)}} \{Q(\underline{\delta})\} \end{aligned} \quad (8)$$

Taking the partial derivative of the Eq. (8) yields:

$$\begin{aligned} \frac{\partial \{Q(\underline{\delta})\}}{\partial \underline{\delta}} &= \frac{\partial}{\partial \underline{\delta}} \{n^{-1} \underline{y}^* W^{-1} \underline{y}^* - 2n^{-1} \underline{\delta}' B'(t) W^{-1} \underline{y}^* \\ &\quad + \underline{\delta}' (n^{-1} B'(t) W^{-1} B(t) + \lambda D) \underline{\delta}\} \end{aligned} \quad (9)$$

The normal solution of the Eq. (9) is

$$\begin{aligned} -2n^{-1} B'(t) W^{-1} \underline{y}^* + 2(n^{-1} B'(t) W^{-1} B(t) + \lambda D) \underline{\delta} &= 0 \\ -2(n^{-1} B'(t) W^{-1} B(t) + \lambda D) \underline{\delta} &= -2n^{-1} B'(t) W^{-1} \underline{y}^* \end{aligned}$$

The estimator for  $\underline{\delta}$  is given as:

$$\hat{\underline{\delta}}_\lambda = (n^{-1} B'(t) W^{-1} B(t) + \lambda D)^{-1} n^{-1} B'(t) W^{-1} (\underline{y} - X\hat{\underline{\beta}}).$$

The Fourier series estimator for the estimation curve of the nonparametric component in the heteroskedastic semiparametric regression is define as:

$$\begin{aligned}\hat{f}_\lambda(t) &= B(t)\hat{\delta}_\lambda \\ &= B(t)(n^{-1}B'(t)W^{-1}B(t) + \lambda D)^{-1}n^{-1}B'(t)W^{-1}(\underline{y} - X\hat{\beta}) \\ &= H(\lambda)\underline{y}^*\end{aligned}\quad (10)$$

where  $\underline{y}^* = \underline{y} - X\hat{\beta}$  and  $H(\lambda) = B(t)(n^{-1}B'(t)W^{-1}B(t) + \lambda D)^{-1}n^{-1}B'(t)W^{-1}$ .

The nonparametric component  $\hat{f}_\lambda(t)$  is somewhat depending on the parametric component  $\hat{\beta}$ .

To search for the parametric component estimator  $\hat{\beta}$ , a least mean square method was used:

$$\underline{y} = X\beta + B\alpha + \varepsilon$$

$$\underline{y} = X\beta + \hat{f}_\lambda(t) + \varepsilon$$

$$\varepsilon = \underline{y} - X\beta - \hat{f}_\lambda(t)$$

$$\varepsilon_i = y_i - x_i'\beta - \hat{f}_\lambda(t)$$

$$\sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - x_i'\beta - \hat{f}_\lambda(t))^2$$

$$\begin{aligned}\varepsilon'\varepsilon &= (y_i - x_i'\beta - \hat{f}_\lambda(t))' (y_i - x_i'\beta - \hat{f}_\lambda(t)) \\ &= [\underline{y} - X\beta - (H(\lambda)(\underline{y} - X\beta))]' [\underline{y} - X\beta - (H(\lambda)(\underline{y} - X\beta))] \\ &= \underline{y}'(I - H(\lambda))'(I - H(\lambda))\underline{y} - 2\beta'X'(I - H(\lambda))'(I - H(\lambda))\underline{y} \\ &\quad + \beta'X'(I - H(\lambda))'(I - H(\lambda))X\beta\end{aligned}$$

Assume that

$$\begin{aligned}\underline{y}'(I - H(\lambda))'(I - H(\lambda))\underline{y} - 2\beta'X'(I - H(\lambda))'(I - H(\lambda))\underline{y} \\ + \beta'X'(I - H(\lambda))'(I - H(\lambda))X\beta = Q(\beta)\end{aligned}$$

Deriving  $Q(\beta)$  on  $\beta$  and taking its zero value will yield:

$$\begin{aligned}\frac{\partial\{Q(\beta)\}}{\partial\beta} &= \frac{\partial}{\partial\beta} \{ \underline{y}'(I - H(\lambda))'(I - H(\lambda))\underline{y} - 2\beta'X'(I - H(\lambda))'(I - H(\lambda))\underline{y} \\ &\quad + \beta'X'(I - H(\lambda))'(I - H(\lambda))X\beta \} \\ &= -2\beta'X'(I - H(\lambda))'(I - H(\lambda))\underline{y} + 2X'(I - H(\lambda))'(I - H(\lambda))X\beta \\ &\quad - 2\beta'X'(I - H(\lambda))'(I - H(\lambda))\underline{y} + 2X'(I - H(\lambda))'(I - H(\lambda))X\beta = 0 \\ &\quad - 2X'(I - H(\lambda))'(I - H(\lambda))X\beta = -2\beta'X'(I - H(\lambda))'(I - H(\lambda))\underline{y} \\ \hat{\beta} &= [X'(I - H(\lambda))'(I - H(\lambda))X]^{-1} X'(I - H(\lambda))'(I - H(\lambda))\underline{y} \\ \hat{\beta} &= [X'T(\lambda)X]^{-1} X'T(\lambda)\underline{y}\end{aligned}$$

which is the parametric component estimator, where  $T(\lambda) = (I - H(\lambda))'(I - H(\lambda))$ . Furthermore

$$\begin{aligned}\hat{f}_{\lambda}(t) &= H(\lambda)(y - X\beta) \\ &= H(\lambda)y - H(\lambda)X[X'T(\lambda)X]^{-1}X'T(\lambda)y \\ &= H(\lambda)y\{I - X[X'T(\lambda)X]^{-1}X'T(\lambda)\} \\ &= M(\lambda)y\end{aligned}\quad (11)$$

is the nonparametric component estimator with  $M(\lambda) = H(\lambda)\{I - X[X'T(\lambda)X]^{-1}X'T(\lambda)\}$ . Hence, the heteroskedastic semiparametric regression model for the Eq. (1) is:

$$\begin{aligned}\hat{y} &= X\hat{\beta} + \hat{f}_{\lambda}(t) \\ &= X[X'T(\lambda)X]^{-1}X'T(\lambda)y + M(\lambda)y \\ &= \{X[X'T(\lambda)X]^{-1}X'T(\lambda) + M(\lambda)\}y \\ &= N(\lambda)y\end{aligned}\quad (12)$$

where:  $N(\lambda) = X[X'T(\lambda)X]^{-1}X'T(\lambda) + M(\lambda)$ .

### 3. Conclusions

Given that a heteroskedastic semiparametric regression model  $y_i = x_i'\beta + f(t_i) + \sigma_i\varepsilon_i$ ,  $i = 1, 2, \dots, n$ . The nonparametric regression component  $f$  form is unknown and is assumed to be smooth; it is defined in a continuous function space  $C[0, \pi]$ . The random error  $\varepsilon_i$  is mutually independent on the zero mean and variance  $\sigma^2$ . Fourier series from Eq. 2 is used to approach the nonparametric component regression curve  $f(t)$ . The estimation curve of the nonparametric component of the heteroskedastic semiparametric regression was obtained by solving WPLS optimization:

$$\min_{f \in C(0, \pi)} \left\{ n^{-1}(y - X\beta - f)'W^{-1}(y - X\beta - f) + \lambda \int_0^{\pi} \frac{2}{\pi} [f''(t)]^2 dt \right\}$$

- The solution of the WPLS optimization gave the Fourier series-based nonparametric component estimator, which was described in  $\hat{f}_{\lambda}(t) = M(\lambda)y$ , for matrix  $M(\lambda)$  and parametric component  $\hat{\beta} = [X'T(\lambda)X]^{-1}X'T(\lambda)y$ .
- The finest estimator for nonparametric component by Fourier series in the heteroskedasticity semiparametric regression model is strongly dependent on the optimum values of bandwidth parameter  $\lambda$  and  $K$  that can be obtained by various methods, for instance by means of generalized cross validation (GCV).

### REFERENCES

- [1] I. N. Budiantara, "Spline Model with Optimal Knot", *Journal of Basic Science FMIPA Jember University* vol. 7, no. 1, pp. 77-85, 2006.

- [2] G. Wahba, *Spline Models for Observational Data*, Philadelphia: Society for Industrial and Applied Mathematics, 1990.
- [3] W. Hardle, H. Liang and J. Gao, *Partially Linear Models*, Physica-Verlag Heidelberg, 2003.
- [4] W. Hardle, Y. Mori and P. Vieu, *Statistical Methods for Biostatistics and Related Fields*, New York: Springer-Verlag Berlin Heidelberg, 2007.
- [5] H. Becher, G. Kauermann, P. Khomski and B. Kouyaté, “Using Penalized Splines to Model Age- and Season-of-birth Dependent Effects of Childhood Mortality Risk Factors in Rural Burkina Faso”, *Biometrical Journal*, vol. 51, no. 1, pp. 110-122, 2009.
- [6] A. Antoniadis, J. Bigot and T. Sapatinas, “Wavelet Estimators in Nonparametric Regression: A Comparative Simulation Study”, *Journal of Statistical Software*, vol. 6, pp. 1-83, 2001.
- [7] M. Y. Cheng, R. L. Paige, S. Sun and K. Yan, “Variance Reduction for Kernel Estimators in Clustered/Longitudinal Data Analysis”, *Journal of Statistical Planning and Inference*, vol. 140, no. 6, pp. 1389-1397, 2010.
- [8] N. Wang, “Marginal Nonparametric Kernel Regression Accounting for Within-Subject Correlation”, *Biometrika*, vol. 90, no. 1, pp. 43-52, 2003.
- [9] M. Kayri and G. Zirhlioğlu, “Kernel Smoothing Function and Choosing Bandwidth for Non-parametric Regression Methods”, *Ozean Journal of Applied Sciences.*, vol. 2, no. 1, pp. 49-60, 2009.
- [10] C. O. Wu and C. T. Chiang, “Kernel Smoothing on Varying Coefficient Models with Longitudinal Dependent Variable”, *Statistica Sinica*, vol. 10, no. 2, pp. 433-456, 2000.
- [11] A. Tripena and I. N. Budiantara, *Fourier Estimator in Nonparametric Regression*, International Conference on Natural Sciences and Applied Natural Sciences, Ahmad Dahlan University, Yogyakarta, 2007.
- [12] M. Bilodeau, “Fourier Smoother and Additive Models”, *The Canadian Journal of Statistics*, vol. 20, no. 3, pp. 257-269, 1992.
- [13] I. N. Budiantara, “Inferensi Statistik untuk Model Spline”, *Jurnal Mat-Stat*, vol. 7, no. 1, pp. 1-14, 2007.