

ESTIMATOR TRUNCATED SPLINE DAN SIFAT LINIERNYA DALAM REGRESI NONPARAMETRIK MULTIVARIABEL

Ni Putu Ayu Mirah Mariati^{a,*}, I Wayan Sudiarsa^b, Gusti Ayu Made Arna Putri^c

^{a,c}Universitas Mahasaraswati Denpasar, Indonesia

^bUniversitas PGRI Mahadewa Indonesia, Denpasar, Indonesia

*email: ayumirahmariati@unmas.ac.id

Abstrak. Analisis regresi digunakan untuk menyelidiki pola hubungan antara variabel dependen dan variabel independen. Hal ini dapat dilakukan dengan dua pendekatan, yaitu pendekatan parametrik dan nonparametrik. Pendekatan parametrik mengasumsikan bentuk model mengikuti suatu pola tertentu. Namun, jika tidak ada informasi tentang bentuk fungsi regresi, maka pendekatan yang digunakan adalah pendekatan regresi nonparametrik. Ada beberapa pendekatan untuk menaksir kurva regresi nonparametrik, salah satunya adalah truncated spline. Keuntungan dari truncated spline adalah dapat menggambarkan perubahan pola perilaku fungsi pada subinterval tertentu. Estimator spline sangat bergantung pada titik knot dan sifat estimatornya adalah linier.

Kata Kunci: Spline, Nonparametric Regression.

PENDAHULUAN

Terdapat tiga pendekatan analisis regresi dalam menaksir kurva regresi, yaitu pendekatan regresi parametrik, regresi nonparametrik, dan regresi semiparametrik (Eubank, 1999). Pendekatan model regresi parametrik memiliki sifat-sifat yang baik dari perspektif Statistika Inferensial seperti sederhana, mudah diinterpretasikan, ekonomis, tidak bias, merupakan estimator linier, efisien, konsisten, dan Best Linear Unbiased Estimator (BLUE). Meskipun sangat reliabel dari perspektif inferensial, pendekatan ini mensyaratkan terpenuhinya asumsi-asumsi yang diperlukan. Tidak semua permasalahan pola hubungan dapat didekati dengan regresi parametrik karena tidak semua permasalahan memiliki informasi mengenai bentuk hubungan atau kurva regresi antara variabel respon dengan variabel prediktor. Jika dipaksakan dengan pendekatan regresi parametrik, maka akan memberikan simpulan yang menyesatkan. Oleh karena itu, terdapat pilihan untuk menggunakan pendekatan nonparametrik yang tidak mensyaratkan informasi mengenai pola hubungan antar variabel. Model regresi nonparametrik yang bentuk kurva regresinya diasumsikan tidak diketahui. Selain itu diasumsikan mulus dimana kurva berada pada ruang fungsi seperti ruang Hilbert, ruang Sobolev, ruang fungsi kontinu dan lain-lain (Hardle, 1994).

Perbedaan antara parametrik dan nonparametrik adalah pada pendekatan parametrik data cenderung dipaksa mengikuti pola tertentu, sedangkan pada pendekatan nonparametrik data diberi kebebasan untuk mencari pola kurva regresinya sendiri sehingga sangat fleksibel dan objektif. Beberapa model regresi nonparametrik yang sering digunakan adalah Histogram, Kernel, kNN, Spline, Neural Network (NN), Local Polynominal, Orthogonal Series, Fourier Series, Wavelet, MARS, dan lain-lain (Wahba, 1985). Semua model memiliki kelebihan dan kekurangan serta motivasinya sendiri dalam memodelkan pola. Dari beberapa model regresi

nonparametrik, yang paling populer adalah Spline. Penelitian secara teoritis juga banyak berkembang, diantaranya (Dette *et al.*, 2018) Spline merupakan potongan polinomial yang memiliki sifat tersegmentasi dan kontinu (terpotong). Regresi spline terpotong digunakan karena kelebihanannya adalah model ini dapat menemukan estimasi datanya sendiri ke mana pun pola data bergerak (Hidayat, *et al.*, 2019) Dengan titik-titik simpul ini, spline dapat memberikan fleksibilitas yang lebih besar daripada polinomial, sehingga memungkinkan untuk beradaptasi secara efektif dengan karakteristik lokal. Alasan lain untuk memilih regresi spline terpotong adalah karena bersifat objektif dan optimasinya menggunakan metode kuadrat terkecil sehingga secara matematis mudah, sederhana, dan baik dalam membantu inferensi statistik dalam artikelnya menyarankan untuk menggunakan regresi spline polinomial jika plot data tidak jelas, simpangan baku besar, dan untuk penyederhanaan. Pada regresi spline nonparametrik data cross section dan satu respon, bahwa jika nilai parameter penghalusan sangat kecil akan memberikan estimator kurva regresi yang sangat kasar. Sebaliknya, jika nilai parameter penghalusan sangat besar, akan dihasilkan estimator kurva regresi nonparametrik yang sangat halus (Gu, C, 2013). Oleh karena itu, dalam penduga spline untuk data cross section, perlu dipilih parameter penghalusan yang optimal agar diperoleh penduga yang paling tepat untuk data tersebut. Dalam penelitian ini, kurva regresi nonparametrik akan diestimasi menggunakan pendekatan spline terpotong dan sifat-sifat liniernya.

METODE PENELITIAN

Mendapatkan estimasi kurva regresi nonparametrik Spline nonparametrik multivariabel langkah-langkahnya adalah sebagai berikut:

- a. Membentuk model dengan Nonparametrik Spline Multivariabel

$$y_i = \sum_{j=1}^p f_j(x_{ji}) + \varepsilon_i ; i = 1, 2, \dots, n$$

- b. Menghampiri Kurva Regresi Spline *Truncated* derajat m dengan r titik knot

$$f_j(x_{ji}) = \sum_{v=1}^m \beta_{vj} x_{ji}^v + \sum_{k=1}^r \beta_{j(m+k)} (x_{ji} - K_{jk})_+^m$$

- c. Membentuk optimasi *Least Square*
- d. Menyajikan persamaan *Least Square* dalam bentuk matriks

$$\text{Min}_{\beta \in R^{p(m+r)}} \left\{ n^{-1} \|y - X\beta\| \right\}$$

- e. Menyelesaikan optimasi (d) dengan derivatif parsial.
- f. Mencari sifat linier estimatornya.

HASIL DAN PEMBAHASAN

Estimator Spline Truncated

Dalam bagian ini dibahas tentang estimasi Spline multivariabel. Spline merupakan jumlahan dari fungsi polinomial dengan suatu fungsi (*truncated*).

Diberikan model regresi nonparametrik multivariabel:

$$y_i = \mu(x_{1i}, x_{2i}, \dots, x_{pi}) + \varepsilon_i$$

Selanjutnya, kurva regresi f_j dihipotesis dengan fungsi Spline multivariabel maka dapat ditulis menjadi:

$$\mu(x_{1i}, x_{2i}, \dots, x_{pi}) = \sum_{j=1}^p f_j(x_{ji}) + \varepsilon_i; i = 1, 2, \dots, n$$

dengan:

$$f_j(x_{ji}) = \sum_{v=1}^m \beta_{vj} x_{ji}^v + \sum_{k=1}^r \beta_{j(k+m)} (x_{ji} - K_{jk})_+^m$$

dimana $j = 1, 2, \dots, p$; $v = 1, 2, \dots, m$; dan $k = 1, 2, \dots, r$ titik-titik knot yang memperlihatkan perubahan pola perilaku dari fungsi tersebut pada sub-sub interval yang berbeda. Dengan demikian model regresi dapat ditulis menjadi:

$$\begin{aligned} \sum_{j=1}^p f_j(x_{ji}) &= \sum_{j=1}^p \left(\beta_{1j} x_{ji}^1 + \dots + \beta_{mj} x_{ji}^m + \beta_{j(1+m)} (x_{ji} - K_{j1})_+^m + \dots + \beta_{j(r+m)} (x_{ji} - K_{jr})_+^m \right) \\ &= \left(\beta_{11} x_{1i}^1 + \dots + \beta_{m1} x_{1i}^m + \beta_{1(1+m)} (x_{1i} - K_{11})_+^m + \dots + \beta_{1(r+m)} (x_{1i} - K_{1r})_+^m \right) + \dots + \\ &\quad \left(\beta_{1p} x_{pi}^1 + \dots + \beta_{mp} x_{pi}^m + \beta_{p(1+m)} (x_{pi} - K_{p1})_+^m + \dots + \beta_{p(r+m)} (x_{pi} - K_{pr})_+^m \right) \end{aligned} \quad (1)$$

Apabila persamaan (1) disajikan dalam bentuk matrik, maka didapat:

$$\begin{aligned} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} &= \begin{bmatrix} x_{11}^1 & \dots & x_{11}^m & (x_{11} - K_{11})_+^m & \dots & (x_{11} - K_{1r})_+^m \\ x_{12}^1 & \dots & x_{12}^m & (x_{12} - K_{11})_+^m & \dots & (x_{12} - K_{1r})_+^m \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ x_{1n}^1 & \dots & x_{1n}^m & (x_{1n} - K_{11})_+^m & \dots & (x_{1n} - K_{1r})_+^m \end{bmatrix} \begin{bmatrix} \beta_{11} \\ \vdots \\ \beta_{m1} \\ \beta_{1(1+m)} \\ \vdots \\ \beta_{1(r+m)} \end{bmatrix} + \dots + \\ &\quad \begin{bmatrix} x_{p1}^1 & \dots & x_{p1}^m & (x_{p1} - K_{p1})_+^m & \dots & (x_{p1} - K_{pr})_+^m \\ x_{p2}^1 & \dots & x_{p2}^m & (x_{p2} - K_{p1})_+^m & \dots & (x_{p2} - K_{pr})_+^m \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ x_{pn}^1 & \dots & x_{pn}^m & (x_{pn} - K_{p1})_+^m & \dots & (x_{pn} - K_{pr})_+^m \end{bmatrix} \begin{bmatrix} \beta_{1p} \\ \vdots \\ \beta_{mp} \\ \beta_{p(1+m)} \\ \vdots \\ \beta_{p(r+m)} \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix} \end{aligned} \quad (2)$$

Persamaan (2) dapat ditulis menjadi:

$$y = X(K)\beta + \varepsilon$$

dengan,

$$y = [y_1, \dots, y_n]$$

$$K = [K_{11}, \dots, K_{1r} \dots K_{p1}, \dots, K_{pr}]$$

$$X(K) = [A_1 \quad \dots \quad A_p]$$

$$\underline{\beta} = (\underline{\beta}'_1, \dots, \underline{\beta}'_p)'$$

$$\underline{\beta}'_1 = (\beta_{11}, \dots, \beta_{m1}, \beta_{1(1+m)}, \dots, \beta_{1(r+m)})', \dots, \underline{\beta}'_p = (\beta_{1p}, \dots, \beta_{mp}, \beta_{p(1+m)}, \dots, \beta_{p(r+m)})'$$

dan $\underline{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$

Estimator parameter $\underline{\beta}$ didapat dari menyelesaikan optimasi:

$$\underset{\underline{\beta} \in \mathbb{R}^{p(m+r)}}{\text{Min}} \left\{ n^{-1} \|y - X \underline{\beta}\|^2 \right\} = \left\{ (y - X(K) \underline{\beta})' (y - X(K) \underline{\beta}) \right\}$$

Untuk menyelesaikan optimasi dengan menggunakan derivatif parsial, misalkan:

$$\square(\underline{\beta}) = \left\{ (y - X(K) \underline{\beta})' (y - X(K) \underline{\beta}) \right\} \tag{3}$$

Persamaan (3) diturunkan terhadap $\underline{\beta}$:

$$\frac{\partial \square(\underline{\beta})}{\partial \underline{\beta}} = -2X'(K)y + 2X'(K)X(K)\underline{\beta} \tag{4}$$

Setelah persamaan (4) diturunkan hasilnya disamakan dengan 0, maka diperoleh persamaan:

$$0 = -2X'(K)y + 2X'(K)X(K)\hat{\underline{\beta}} \tag{5}$$

Hasil yang diperoleh dari persamaan (5) adalah

$$X'(K)X(K)\hat{\underline{\beta}} = X'(K)y$$

Sehingga estimator $\hat{\underline{\beta}}$ diberikan oleh:

$$\hat{\underline{\beta}} = (X'(K)X(K))^{-1} X'(K)y$$

dengan $\hat{\underline{\beta}} = (\hat{\beta}'_1, \dots, \hat{\beta}'_p)'$.

Estimator kurva regresi $\hat{f}(x)$ diperoleh dari:

$$\begin{aligned} \hat{f}(x) &= [\hat{f}(x_1)', \hat{f}(x_2)', \dots, \hat{f}(x_p)']' \\ &= X(K)\hat{\underline{\beta}} \\ &= (X'(K)X(K))^{-1} X'(K)y \end{aligned}$$

Sehingga diperoleh,

$$\hat{f}(x) = B(K)y.$$

SIFAT LINIER ESTIMATOR

Model regresi nonparametrik spline dapat ditulis dalam bentuk matriks berikut. $\underline{y} = X(K)\underline{\beta} + \underline{\varepsilon}$ dengan $K = [K_{11}, \dots, K_{1r}, \dots, K_{p1}, \dots, K_{pr}]$. Jika dituliskan $\underline{f} = X(K)\underline{\beta}$ dengan $X(K)$ merupakan matriks fungsi K, maka diperoleh: $\underline{y} = \underline{f} + \underline{\varepsilon}$.

$$\begin{aligned} \hat{\underline{f}} &= X(K)\hat{\underline{\beta}} \\ &= (X'(K)X(K))^{-1} X'(K)\underline{y} \\ &= B(K)\underline{y} \end{aligned}$$

Berdasarkan persamaan di atas, terlihat bahwa estimator spline $\hat{\underline{f}}$ merupakan estimator yang linier. Kelinieran ini dapat memberikan kemudahan bagi peneliti dalam membentuk statistik inferensi untuk pendekatan spline.

SIMPULAN DAN SARAN

Simpulan

Jika diketahui model regresi nonparametrik: $y_i = \mu(x_{1i}, x_{2i}, \dots, x_{pi}) + \varepsilon_i$ kurva regresi f_j dihampiri dengan fungsi Spline multivariabel maka dapat ditulis menjadi: $\mu(x_{1i}, x_{2i}, \dots, x_{pi}) = \sum_{j=1}^p f_j(x_{ji}) + \varepsilon_i; i = 1, 2, \dots, n$ sehingga diperoleh, $\hat{f}(x) = B(K)\underline{y}$. Sifat estimator yang diperoleh adalah linier.

Saran

Dapat dikembangkan secara teoritis dalam mencari sifat estimator yang lainnya dan diaplikasikan pada kasus riil.

DAFTAR PUSTAKA

- Dette, H., Möllenhoff, K., Volgushev, S., and Bretz, F. (2018). Equivalence of Regression Curves. *Journal of the American Statistical Association*. 113 (522), 711–729.
- Draper, N. R and Smith, H. (1998). *Applied Regression Analysis*. USA: John Wiley&Sons.
- Eubank, R. L. (1999). *Nonparametric Regression and Spline Smoothing*. New York: Marcel Dekker.
- Fitriana, D.m Budiantara, I.N.m and Ratnasari, V. (2017). Semiparametric Spline Truncated Regression on Modelling AHH in Indonesia. *Proceedings The 3rd International Seminar on Science and Technology*, 2, 26–31.
- Fithriasari, K., Hariastuti, I., and Wening, K. S. (2020). Handling Imbalance Data in Classification Model with Nominal Predictors, *International Journal Of Computing Science And Applied Mathematics*. 6(1), 33-37.
- Gu, C. (2013). *Smoothing Spline ANOVA Models*. New York: Springer
- Hardle, W. K. (1994). *Applied Nonparametric Regression*. New York: Cambridge University Press.

- Hidayat, R., Budiantara, I.N., Otok, B.W., and Ratnasari, V. (2019). A reproducing kernel hilbert space approach and smoothing parameters selection in spline-kernel regression, *Journal of Theoretical and Applied Information Technology*, 97 (2), 465–475.
- Liu, X., and Preve, D. (2016). Measure of location-based estimators in simple linear regression, *Journal of Statistical Computation and Simulation*. 86(9), 1771–1784.
- Mozumder, S. I., Rutherford, M., and Lambert, P. (2017). Direct Likelihood Inference On The Cause-Specific Cumulative Incidence Function: A Flexible Parametric Regression Modelling Approach. *Statistics in Medicine*. 37, 1–16.
- Takezawa, K. (2006). *Introduction to Nonparametric Regression*. USA: John Wiley&Sons.
- Wahba, G. (1985). A Comparison of GCV and GML for Choosing The Smoothing Parameter in The Generalized Spline Smoothing Problem, *The Annals of Statistics*. 13 (4), 1378-1402.
- Wahba, G. (1990). *Spline Models For Observation Data*. Pennsylvania: SIAM.
- Wang, X., Shen, J., and Ruppert, D. (2011). On the asymptotics of penalized spline smoothing, *Electronic Journal of Statistics*. 5, 1–17.
- Wang, Y. (2011). *Smoothing Splines. Methods and Applications*. New York,: Chapman & Hall CRC
- Wahba, G. (1990). *Spline Models for Observational Data*. Pennsylvania: SIAM.