

Kajian Metode Estimasi Parameter dalam Regresi Semiparametrik Spline

Wahyu Wibowo¹, Sri Haryatmi¹, I Nyoman Budiantara²

wahyu.stk@gmail.com

¹Jurusan Matematika, Universitas Gadjah Mada Yogyakarta

²Jurusan Statistika, Institut Teknologi Sepuluh Nopember Surabaya

Intisari

Pada regresi semiparametrik spline, estimasi kurva regresi dapat diselesaikan dengan metode kuadrat terkecil, kuadrat terkecil terpenalti, dan reproducing kernel Hilbert space. Masing-masing pendekatan memiliki karakteristik yang berbeda. Pada pendekatan kuadrat terkecil, masalah utama adalah memilih banyak knot dan lokasinya. Pada kuadrat terkecil terpenalti dan reproducing kernel mempunyai permasalahan yang sama, yaitu menentukan nilai parameter penghalus yang optimal. Namun, secara matematis, metode reproducing kernel memberi hasil yang lebih umum dibandingkan dengan kuadrat terkecil terpenalti karena berlaku untuk semua derajat polinomial spline yang akan dipergunakan. Makalah ini akan menjelaskan penggunaan metode kuadrat terkecil, kuadrat terkecil terpenalti, dan reproducing kernel hilbert space serta perbedaan masing-masing dalam estimasi kurva regresi semiparametrik spline.

Kata kunci : *Regresi semiparametrik, least square, penalized least square, reproducing kernel Hilbert space, spline*

Abstract

Curve estimation for spline semiparametric regression can be done by using least square, penalized least square and reproducing kernel Hilbert space method. Each methods has different characteristics. Least square method has problem about how to choose the number of knots and their location. Penalized least square and reproducing kernel has the same problem about how to choose the optimal smoothing parameter. However, reproducing kernel approach is more general mathematically than penalized least square due to be valid for any degree of polynomial spline that to be used. This paper will explaine about least square, penalized least square, reproducing kernel Hilbert space as weel as difference of each method in curve estimation for spline semiparametric regression.

Keyword : *semiparametric regression, least square, penalized least square, reproducing kernel Hilbert space, spline*

1. Pendahuluan

Regresi semiparametrik merupakan kombinasi antara regresi parametrik dan regresi nonparametrik. Kombinasi dalam hal ini dimaksudkan bahwa dalam regresi semiparametrik memuat sekaligus model regresi parametrik dan model regresi nonparametrik. Regresi semiparametrik ini muncul karena adanya kasus-kasus pemodelan yang hubungan antar variabelnya selain ada yang linear juga ada yang tidak diketahui bentuknya. Keberadaan dua komponen yang berbeda dalam regresi semiparametrik ini menjadikan pemakaian model ini menjadi luas dan secara teori berkembang sangat pesat. Perkembangan ini selain karena aplikasi juga karena berkembangnya perangkat keras teknologi komputasi yang mempermudah dan mempercepat komputasi.

Estimasi kurva regresi semiparametrik spline dapat diselesaikan dengan metode kuadrat terkecil, kuadrat terkecil terpenalti, dan reproducing kernel hilbert space. Masing-masing pendekatan memiliki karakteristik yang berbeda. Pada pendekatan kuadrat terkecil, masalah utama adalah memilih banyak knot dan lokasinya. Pada kuadrat terkecil terpenalti dan reproducing kernel mempunyai permasalahan sama, yaitu menentukan nilai parameter penghalus yang optimal. Namun, secara matematis, metode reproducing kernel memberi hasil yang lebih umum dibandingkan dengan kuadrat terkecil terpenalti karena berlaku untuk semua derajat polinomial spline yang akan dipergunakan

Makalah ini akan menjelaskan penggunaan metode kuadrat terkecil, kuadrat terkecil terpenalti, dan reproducing kernel hilbert space serta perbedaan masing-masing dalam estimasi kurva regresi semiparametrik spline. Penjelasan akan dimulai dengan konstruksi model regresi semiparametrik spline, metode kuadrat terkecil, kuadrat terkecil terpenalti, dan reproducing kernel Hilbert space. Sebagai penutup akan diberikan kesimpulan berkaitan dengan metode-metode tersebut.

2. Model

Pandang n sampel random dengan variabel pada masing-masing sampel adalah (y_i, x_i, t_i) , $i=1,2,\dots,n$. Dalam hal ini diasumsikan y_i sebagai variabel respon, x_i sebagai variabel prediktor yang diketahui berpengaruh linear, dan t_i sebagai variabel prediktor yang tidak diketahui bentuk pengaruhnya terhadap respon. Selanjutnya dibentuk model regresi semiparametrik

$$y_i = \beta_0 + \beta_1 x_i + f(t_i) + \varepsilon_i, \quad i=1,2,\dots,n \quad (1)$$

(β_0, β_1) parameter untuk komponen parametrik, dan f komponen nonparametrik, dalam hal ini merupakan fungsi yang tidak diketahui. $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ adalah error random yang saling independen dengan mean nol dan variansi σ^2 .

Model (1) dapat dinyatakan dalam notasi matrik menjadi

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{f} + \boldsymbol{\varepsilon} \quad (2)$$

$$\text{dimana } \mathbf{Y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} \\ \vdots & \vdots \\ 1 & x_{1n} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}, \quad \mathbf{f} = \begin{pmatrix} f(t_1) \\ \vdots \\ f(t_n) \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

Masalah estimasi pada regresi semiparametrik muncul karena adanya komponen nonparametrik berupa fungsi f yang tidak diketahui bentuknya. Oleh karena itu, hampiran terhadap bentuk fungsi tersebut dapat dilakukan dengan lebih dari satu bentuk fungsi. Beberapa diantaranya adalah spline, kernel, fourier, wavelet, dan polinomial lokal. Secara aplikasi, hampiran-hampiran ini memiliki kelebihan yang berbeda. Sebagai contoh, hampiran spline sangat cocok untuk data yang memiliki perilaku berubah-ubah dalam sub-sub interval tertentu. Hampiran fourier sangat cocok untuk data yang memiliki pola berulang atau musiman. Penjelasan mengenai macam-macam hampiran ini dapat dilihat pada Eggermont dan Lariccia (2009).

3. Metode Kuadrat Terkecil

Metode kuadrat terkecil merupakan metode yang sangat lazim dipergunakan dalam regresi linear. Prinsip metode ini adalah meminimumkan kuadrat residual. Metode ini juga bisa dipergunakan dalam regresi semiparametrik spline. Wibowo, dkk (2009, 2010) telah menggunakan metode ini untuk estimasi parameter pada regresi semiparametrik dan juga

sifat-sifat statistik estimator yang diperoleh. Penggunaan metode ini mensyaratkan bentuk spesifik fungsi $f(t)$ dalam model (1). Bentuk fungsi spline yang biasa dipergunakan adalah fungsi basis spline polinomial *truncated*.

Fungsi spline polinomial *truncated* derajat p dengan k titik knots $\kappa_1, \kappa_2, \dots, \kappa_k$ disajikan dalam bentuk,

$$f(t_i) = \alpha_0 + \alpha_1 t_i + \alpha_1 t_i^2 + \dots + \alpha_p t_i^p + \alpha_{p1} (t_i - \kappa_1)_+^p + \dots + \alpha_{pk} (t_i - \kappa_k)_+^p \quad (3)$$

dengan $(t_i - \kappa)_+^p = \begin{cases} (t_i - \kappa)^p, & t_i \geq \kappa \\ 0, & t_i < \kappa \end{cases}$

Sehingga model (1) dapat dinyatakan menjadi

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_i + \alpha_0 + \alpha_1 t_i + \alpha_2 t_i^2 + \dots + \alpha_p t_i^p + \\ &\quad \alpha_{p1} (t_i - \kappa_1)_+^p + \dots + \alpha_{pk} (t_i - \kappa_k)_+^p + \varepsilon_i, \quad (4) \\ i &= 1, 2, \dots, n \end{aligned}$$

Apabila dinyatakan dalam matriks, diperoleh

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\alpha} + \boldsymbol{\varepsilon} \quad (5)$$

dengan

$$\mathbf{Z} = \begin{bmatrix} 1 & x & \dots & x^p & (t_i - \kappa_1)_+^p & \dots & (t_i - \kappa_k)_+^p \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 1 & x & \dots & x^p & (t_i - \kappa_1)_+^p & \dots & (t_i - \kappa_k)_+^p \end{bmatrix}_{n \times (1+p+k)} \quad \boldsymbol{\alpha} = \begin{bmatrix} \alpha_0 & \alpha_1 & \dots & \alpha_p & \alpha_{p1} & \dots & \alpha_{pk} \end{bmatrix}_{1 \times (1+p+k)}^T$$

Selanjutnya, ditentukan $\mathbf{C} = [\mathbf{x} \ \mathbf{z}]$, $\boldsymbol{\omega} = \begin{bmatrix} \boldsymbol{\beta} \\ \boldsymbol{\alpha} \end{bmatrix}$, sehingga (5) dapat dinyatakan dengan

$$\mathbf{Y} = \mathbf{C}\boldsymbol{\omega} + \boldsymbol{\varepsilon}$$

Selanjutnya didefinisikan kuadrat residual sebagai berikut,

$$\begin{aligned} \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} &= (\mathbf{Y} - \mathbf{C}\boldsymbol{\omega})^T (\mathbf{Y} - \mathbf{C}\boldsymbol{\omega}) \\ &= \mathbf{Y}^T \mathbf{Y} - 2\boldsymbol{\omega}^T \mathbf{C}^T \mathbf{Y} + \boldsymbol{\omega}^T \mathbf{C}^T \mathbf{C} \boldsymbol{\omega} \quad (6) \end{aligned}$$

Estimator $\boldsymbol{\omega}$ diperoleh dengan cara meminimumkan (6) terhadap $\boldsymbol{\omega}$. Selanjutnya (6) diturunkan terhadap $\boldsymbol{\omega}$ dan disamadengangkan nol, sehingga diperoleh :

$$\frac{\partial \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}}{\partial \boldsymbol{\omega}} = -2\mathbf{C}^T \mathbf{Y} + 2\mathbf{C}^T \mathbf{C} \boldsymbol{\omega} = 0 \quad (7)$$

$$\mathbf{C}^T \mathbf{C} \boldsymbol{\omega} = \mathbf{C}^T \mathbf{Y}$$

yang merupakan persamaan normal. Penyelesaian persamaan ini akan merupakan estimator $\boldsymbol{\omega}$. Sesuai aljabar matrik, karena matrik \mathbf{C} mempunyai rank $(3 + p + k)$ dan $\mathbf{C}^T \mathbf{C}$ matrik *positive-definite*, maka $\mathbf{C}^T \mathbf{C}$ akan merupakan matrik nonsingular. Sehingga persamaan (6) akan mempunyai penyelesaian tunggal,

$$\hat{\boldsymbol{\omega}} = (\mathbf{C}^T \mathbf{C})^{-1} \mathbf{C}^T \mathbf{Y} \quad (8)$$

Dalam hal ini, $\hat{\boldsymbol{\omega}}$ merupakan estimator kuadrat terkecil $\boldsymbol{\omega}$. Mengingat persamaan (2), estimator (8) berlaku hanya untuk derajat polinomial p dan banyak knots k yang tertentu. Sehingga lebih tepat kalau estimator ini dinyatakan dengan

$$\hat{\boldsymbol{\omega}}_{p,k} = (\mathbf{C}^T \mathbf{C})^{-1} \mathbf{C}^T \mathbf{Y}$$

$$\hat{\mathbf{Y}} = \mathbf{C}\hat{\boldsymbol{\omega}}_{p,\kappa} = \mathbf{C}(\mathbf{C}^T\mathbf{C})^{-1}\mathbf{C}^T\mathbf{Y}$$

$$\text{dengan } \mathbf{H}(p; \kappa_1, \dots, \kappa_k) = \mathbf{C}(\mathbf{C}^T\mathbf{C})^{-1}\mathbf{C}^T$$

Permasalahan selanjutnya adalah bagaimana menentukan derajat polinomial p dan banyak knots k yang akan digunakan dalam estimator. Untuk keperluan ini akan dipergunakan kriteria *Generalized Cross Validation* (GCV), yang didefinisikan :

$$GCV(p; \kappa_1, \kappa_2, \dots, \kappa_k) = \frac{n^{-1} \|(\mathbf{I} - \mathbf{H}(p; \kappa_1, \kappa_2, \dots, \kappa_k))\mathbf{y}\|^2}{[n^{-1} \text{trace}(\mathbf{I} - \mathbf{H}(p; \kappa_1, \kappa_2, \dots, \kappa_k))]^2}$$

Nilai p dan k dipilih dengan menyelesaikan optimasi

$$GCV(p_{opt}; \kappa_{1_{opt}}, \kappa_{2_{opt}}, \dots, \kappa_{k_{opt}}) = \min_{\substack{p \in \mathbb{N}^+ \\ \kappa_1 \in \mathbb{N}^+ \\ \vdots \\ \kappa_k \in \mathbb{N}^+}} (GCV(p; \kappa_1, \kappa_2, \dots, \kappa_k)) \quad \text{Metode estimasi kuadrat terkecil secara matematis dapat diselesaikan dengan langkah-langkah yang sederhana dan menghasilkan model statistik yang mudah diinterpretasikan. Kesulitan metode ini adalah menentukan derajat polinomial, banyak knot dan lokasi knot-knot tersebut.}$$

4. Metode Kuadrat Terkecil Terpenalty

Metode kuadrat terkecil terpenalty merupakan perluasan metode kuadrat terkecil dengan menambahkan parameter penghalus dan penalti pada fungsi yang akan dipergunakan. Fungsi yang akan dipergunakan merupakan keluarga fungsi yang terdifferensial pada interval $[a,b]$ dan kontinu absolut pada turunan pertama. Bentuknya diberikan sebagai berikut,

$$S(\beta, f) = \frac{1}{n} \sum_{i=1}^n \{y_i - (\beta_0 + x_i \beta_1) - f(t_i)\}^2 + \lambda \int_a^b \{f^{(m)}(t)\}^2 dt \quad (9)$$

Dalam hal ini λ merupakan parameter penghalus, sedangkan penalti diberikan oleh

$$\int_a^b \{f^{(m)}(t)\}^2 dt$$

Jika $\lambda \rightarrow 0$, maka hasil estimasi mendekati hasil metode kuadrat terkecil. Sebaliknya, jika $\lambda \rightarrow \infty$, maka estimasi akan menginterpolasi titik-titik data. Estimator terbaik merupakan kompromi antara nilai jumlah kuadrat residual dan parameter penghalus λ yang bisa didapatkan dengan meminimumkan nilai GCV.

Dalam Grenn dan Silverman (1994) dinyatakan bahwa fungsi yang meminimumkan (9) merupakan fungsi *natural cubic spline*, yang diberikan sebagai berikut :

$$f(t) = a_i + b_i(t - t_i) + c_i(t - t_i)^2 + d_i(t - t_i)^3, \quad t_i \leq t \leq t_{i+1}, \quad i = 1, 2, \dots, n$$

t_1, t_2, \dots, t_n merupakan knot yang memenuhi $a < t_1 < t_2 < \dots < t_n < b$. Selanjutnya, penalti dalam (9) akan dinyatakan dalam bentuk nilai turunan kedua dengan langkah sebagai berikut.

1. Didefinisikan $f_i = f(t_i)$ dan $\gamma_i = f''(t_i)$, $i = 1, 2, \dots, n$.
2. Pandang $\mathbf{f} = (f_1, \dots, f_n)^T$, $\gamma = (\gamma_2, \dots, \gamma_{n-1})^T$, $h_i = t_{i+1} - t_i$, $i = 1, 2, \dots, n-1$
3. Susun matrik, namakan \mathbf{Q} dan \mathbf{R} , dengan ketentuan sebagai berikut.

\mathbf{Q} adalah matrik berukuran $n \times (n-2)$ dengan element q_{ij} , $i=1,2,\dots,n$ dan $j=2,3,\dots,n-1$.

$$q_{j-1,j} = h_{j-1}^{-1}, q_{ij} = -h_{j-1}^{-1} - h_j^{-1}, q_{j+1,j} = h_j^{-1} \quad j=2,3,\dots,n-1, \quad q_{ij} = 0, \text{ jika } |i-j| \geq 2.$$

Matriks \mathbf{Q} is diindeks mulai $j=2$, sehingga elemen teratas \mathbf{Q} adalah q_{12} .

Selanjutnya \mathbf{R} adalah matriks symmetric berukuran $(n-2) \times (n-2)$ dengan elements r_{ij} , $i,j=2,3,\dots,(n-1)$ sebagai berikut

$$r_{ij} = \frac{1}{3}(h_{i-1} + h_i), \quad r_{i,j+1} = r_{i+1,i} = \frac{1}{6}h_i, \text{ dan } r_{ij} = 0 \text{ untuk } |i-j| \geq 2$$

4. Definiskan matrik \mathbf{K} dengan

$$\mathbf{K} = \mathbf{Q}\mathbf{R}^{-1}\mathbf{Q}^T$$

Selanjutnya hubungan antara $\mathbf{Q}, \mathbf{R}, \mathbf{f}, \gamma$ and \mathbf{K} terdapat dalam pada Green dan Silverman (1994), yang dinyatakan dalam teorema berikut.

Theorem 1:

Vektor \mathbf{f} dan γ menentukan *natural cubic spline f* jika dan hanya jika memenuhi

$$\mathbf{Q}^T \mathbf{f} = \mathbf{R} \gamma$$

Selanjutnya penalti akan dapat dinyatakan menjadi

$$\int_a^b \{f''(t)\}^2 dt = \gamma^T \mathbf{R} \gamma = \mathbf{f}^T \mathbf{K} \mathbf{f}$$

Bukti dapat dilihat pada Green dan Silverman (1994). Sehingga (9) dapat dinyatakan dalam notasi matrik menjadi ; $S(\beta, \mathbf{f}) = (\mathbf{Y} - \mathbf{X}\beta - \mathbf{f})^T (\mathbf{Y} - \mathbf{X}\beta - \mathbf{f}) + \lambda \mathbf{f}^T \mathbf{K} \mathbf{f}$

Sebagai hasil teorema ini, jumlah kuadrat terpenalti (9) dapat dinyatakan

$$S(\beta, \mathbf{f}) = n^{-1} (\mathbf{Y} - \mathbf{X}\beta - \mathbf{f})^T (\mathbf{Y} - \mathbf{X}\beta - \mathbf{f}) + \lambda \mathbf{f}^T \mathbf{K} \mathbf{f} \quad (10)$$

Memminimumkan (9) equivalent dengan memminimumkan,

$$S(\beta, \mathbf{f}) = (\mathbf{Y} - \mathbf{X}\beta - \mathbf{f})^T (\mathbf{Y} - \mathbf{X}\beta - \mathbf{f}) + \lambda \mathbf{f}^T \mathbf{K} \mathbf{f} \quad (11)$$

Selanjutnya, dengan menurunkan (10) terhadap β and \mathbf{f} , dan menyamadengangkan nol, diperoleh

$$\begin{aligned} \hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{Y} - \hat{\mathbf{f}}) \\ \hat{\mathbf{f}} &= (\mathbf{I} + \lambda \mathbf{K})^{-1} (\mathbf{Y} - \mathbf{X}\hat{\beta}) \end{aligned} \quad (12)$$

Dalam hal ini terhadap dua smoother, $\mathbf{S}_\beta = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ menghasilkan dugaan $\mathbf{X}\hat{\beta}$ dan smoother lainnya adalah $\mathbf{S}_f = (\mathbf{I} + \lambda \mathbf{K})^{-1}$ menghasilkan dugaan $\hat{\mathbf{f}}$. Selanjutnya dugaan parametrik diberikan oleh $\mathbf{S}_\beta(\mathbf{Y} - \mathbf{f}) = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{Y} - \mathbf{f}) \equiv \mathbf{X}\hat{\beta}$ sedangkan dugaan untuk nonparametrik adalah $\hat{\mathbf{f}} = \mathbf{S}_f(\mathbf{Y} - \mathbf{X}\hat{\beta})$. Substitusi dugaan nonparametrik ke bagian pertama (11), menghasilkan

$$\mathbf{X}^T \mathbf{X} \beta = \mathbf{X}^T (\mathbf{Y} - \mathbf{S}_f(\mathbf{Y} - \mathbf{X}\hat{\beta})) \Leftrightarrow \mathbf{X}^T (\mathbf{I} - \mathbf{S}_f) \mathbf{X} \beta = \mathbf{X}^T (\mathbf{I} - \mathbf{S}_f) \mathbf{Y}$$

Persamaan ini adalah persamaan normal untuk *generalized least square* normal, dengan elemen non-diagonal berupa matriks pembobot $(\mathbf{I} - \mathbf{S}_f)$. Sehingga, parameter β dan \mathbf{f} dapat diselesaikan dengan,

$$\hat{\beta} = (\mathbf{X}^T(\mathbf{I} - \mathbf{S}_f)\mathbf{X})^{-1}\mathbf{X}^T(\mathbf{I} - \mathbf{S}_f)\mathbf{Y}$$

$$\hat{\mathbf{f}} = \mathbf{S}_f(\mathbf{Y} - \mathbf{X}\hat{\beta}) = -\mathbf{S}_f\mathbf{X}\hat{\beta} + \mathbf{S}_f\mathbf{Y}$$

Estimator penalized least square masih tergantung pada parameter smoothing λ , oleh karena itu harus dipilih yang optimum dengan meminimumkan GCV. GCV didefinisikan dalam metode ini diformulakan dengan

$$GCV(\lambda) = \frac{\|\mathbf{Y} - \hat{\mathbf{Y}}\|^2}{(1 - n^{-1}tr(\mathbf{A}))^2}$$

\mathbf{A} merupakan *hat* matrik yang memenuhi

$$\mathbf{AY} = \mathbf{X}\hat{\beta} + \hat{\mathbf{f}} = [\mathbf{X} \quad \mathbf{I}] \begin{bmatrix} \hat{\beta} \\ \hat{\mathbf{f}} \end{bmatrix} = \hat{\mathbf{Y}}$$

dengan

$$\mathbf{A} = [\mathbf{X} \quad \mathbf{I}] \begin{bmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \\ \mathbf{X} & (\mathbf{I} + \lambda \mathbf{K}) \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{X}^T \\ \mathbf{I} \end{bmatrix}$$

Bentuk kurva regresi hasil metode kuadrat terkecil terpenalti kuadrat terkecil bergantung pada parameter penghalus λ . Namun metode ini menjadi lebih kompleks jika terdapat lebih dari satu prediktor nonparametrik. Selain itu, interpretasi pengaruh prediktor nonparametrik terhadap respon tidak diberikan secara eksplisit melalui model statistik, akan tetapi melalui kurva regresi.

5. Reproducing Kernel Hilbert Space

Penggunaan metode reproducing kernel Hilbert space dalam regresi semiparametrik pada dasarnya merupakan perluasan metode kuadrat terkecil terpenalti dengan menggunakan fungsi yang terdiferensial pada interval $[a,b]$ dan turunannya yang ke- m kontinu absolut pada interval tersebut. Dengan kata lain, fungsi tersebut termuat di dalam ruang Sobolev $W_2^m[a,b]$ dengan $W_2^m([a,b]) = \{f : f, \dots, f^{(m-1)} \text{ kont. abs., } \int_a^b (f^{(m)}(t))^2 dt < \infty\}$ Reproducing Kernel Hilbert Space (RKHS) H_R adalah suatu ruang Hilbert dari fungsi bernilai real pada $[0,1]$ dengan sifat bahwa untuk setiap $t \in [0,1]$, fungsional $L_t(f) = f(t)$ merupakan fungsional linear terbatas, dalam arti bahwa terdapat M sedemikian hingga berlaku

$$|L_t f| = |f(t)| \leq M \|f\|$$

H_R dapat didekomposisi menjadi $H_R = H_0 \oplus H_1$ dengan H_0 ruang Null, dan H_1 adalah ruang yang tegak lurus dengan ruang Null.

Reproducing kernel dari H_R adalah fungsi R yang didefinisikan pada $[0,1] \times [0,1]$, sedemikian hingga untuk setiap titik $t \in [0,1]$ berlaku $R_t \in H_R$ dengan $R_t(s) = R(t,s)$ dan $L_t f = \langle R_t, f \rangle, f \in H_R$

Untuk menggunakan metode reproducing kernel dalam estimasi parameter regresi semiparametrik, model (1) dinyatakan menjadi

$$y_i = \beta_0 + \beta_1 x_i + L_{t_i} f + \varepsilon_i, \quad i = 1, 2, \dots, n \quad (13)$$

dengan $L_{t_i} f = \langle R_{t_i}, f \rangle = f(t_i)$

Selanjutnya, estimasi parameter diperoleh dengan meminimumkan penalized least square

$$S(\beta, f) = \frac{1}{n} \sum_{i=1}^n \{y_i - (\beta_0 + x_i \beta_1) - L_{t_i} f\}^2 + \lambda \|P_1 f\|_R^2 \quad (14)$$

dengan $\int_a^b \{f^{(m)}(t)\}^2 dt = \|P_1 f\|_{W_m}^2$, P_1 proyeksi orthogonal f pada H_1 dalam H_R .

Estimator yang akan diperoleh terdiri dari estimator spline parsial sebagai komponen nonparametrik (\mathbf{f}), dan estimator parametrik (β). Estimator spline parsial diperoleh berdasarkan teorema berikut.

Teorema 2 :

Apabila $H_R = H_0 \oplus H_1$ dan $\phi_1, \phi_2, \dots, \phi_m$ merupakan basis di ruang H_0 , serta $\mathbf{T}_{n \times m}$ merupakan matrik *full rank* berorde $n \times m$ yang diberikan oleh

$$\mathbf{T}_{n \times m} = \{L_i \phi_v\}, \quad i=1, 2, \dots, n; \quad v=1, 2, \dots, m$$

maka fungsi f yang meminimumkan

$$S(\beta, f) = \frac{1}{n} \sum_{i=1}^n \{y_i - (\beta_0 + x_i \beta_1) - L_i f\}^2 + \lambda \|P_1 f\|_R^2$$

adalah

$$\hat{\mathbf{f}}_\lambda = \sum_{v=1}^m d_v \phi_v + \sum_{i=1}^n c_i \varphi_i \quad (15)$$

dengan

$$\varphi_i = P_1 \eta_i$$

$$\begin{aligned} \mathbf{d} &= (d_1, d_2, \dots, d_m)' = (\mathbf{T}' \mathbf{M}^{-1} \mathbf{T})^{-1} \mathbf{T} \mathbf{M}^{-1} (\mathbf{Y} - \mathbf{X} \beta) \\ \mathbf{c} &= (c_1, c_2, \dots, c_n)' = \mathbf{M}^{-1} (\mathbf{I} - \mathbf{T} (\mathbf{T}' \mathbf{M}^{-1} \mathbf{T})^{-1} \mathbf{T}' \mathbf{M}^{-1} (\mathbf{Y} - \mathbf{X} \beta)) \end{aligned}$$

$$\mathbf{M} = \Sigma + n\lambda, \quad \Sigma = \left\{ \langle \varphi_i, \varphi_j \rangle \right\}; \quad i, j = 1, 2, \dots, n$$

Sedangkan estimator parametrik yang bersesuaian adalah

$$\hat{\beta}_\lambda = (\mathbf{X}' (\mathbf{I} - \mathbf{A}(\lambda)) \mathbf{X})^{-1} \mathbf{X}' (\mathbf{I} - \mathbf{A}(\lambda))' \mathbf{Y} \quad (16)$$

Bukti diberikan pada akhir makalah. Sampai di sini, telah diperoleh estimator untuk regresi semiparametrik spline dengan metode (RKHS). Nilai taksiran untuk variabel respon dapat dinyatakan menjadi

$$\hat{\mathbf{Y}} = \mathbf{X} \hat{\beta}_\lambda + \hat{\mathbf{f}}_\lambda \quad (17)$$

Kelebihan estimator yang diperoleh dengan metode RKHS adalah keumumannya yang berlaku untuk sebarang polinomial spline sampai derajat- m . Interpretasi pengaruh prediktor nonparametrik terhadap respon dapat dengan mudah dilakukan melalui kurva regresi. Jika terdapat lebih dari satu prediktor nonparametrik, maka akan muncul kesulitan secara matematis dan grafis.

6. Kesimpulan dan Saran

Masalah utama pada estimasi parameter regresi semiparametrik adalah adanya komponen nonparametrik berupa fungsi yang tidak diketahui bentuknya. Penggunaan metode kuadrat terkecil mengasumsikan bentuk fungsi spline polinomial truncated dan memberikan kemudahan interpretasi melalui model statistik. Penggunaan kuadrat terkecil terpenalti menghasilkan fungsi *natural cubic spline* sebagai komponen nonparametrik dan bentuk kurva regresinya tergantung pada parameter penghalus. Estimator yang diperoleh dengan reproducing kernel Hilbert space juga tergantung pda parameter penghalus, namun

estimatornya berlaku untuk sebarang derajat polinomial spline. Sehingga estimator yang diperoleh dengan pendekatan reproducing kernel Hilbert space bersifat lebih umum dibandingkan dengan estimator yang diperoleh dengan metode kuadrat terkecil maupun kuadrat terkecil terpenalti.

Daftar Pustaka

- Green, P.J. and Silverman, B.W., 1994, *Nonparametric Regression and Generalized Linear Model*, Chapman & Hall, London
- Eggermont, P.P.B., and Lariccia, V.N., 2009, *Maximum Penalized Likelihood Estimation, Volume II : Regression*, Springer Series in Statistics
- Wahba, G., 1990, *Spline Model for Observational Data*, SIAM, XII, Philadelphia
- Wibowo, W., Haryatmi, S., Budiantara, I.N., 2009, Least Square Methods for Estimating Curve of Spline Semiparametric Regression, *Proceeding of National Seminar on Mathematic and Mathematic Education*, Yogyakarta State University, December 5th 2009, p. 633-645, ISBN : 978-979-16353-3-2
- Wibowo, W., Haryatmi, S., Budiantara, I.N., 2010, Inference And Confidence Interval For Regression Curve In Spline Semiparametric Model, *Proceeding of National Seminar on Mathematic and Mathematic Education*, University of Muhammadiyah Malang, January 30th, 2010

Bukti Teorema 2 :

Untuk bukti estimator nonparametrik, dapat dilihat pada Wahba (1990), sedangkan untuk estimator parametrik diberikan secara singkat.

Untuk mendapatkan estimator parametrik, persamaan (15) dinyatakan dalam bentuk matrik menjadi

$$\mathbf{f}_\lambda = \mathbf{T}\mathbf{d} + \sum \mathbf{c}$$

Selanjutnya, dengan mengingat hubungannya dengan model (2), maka (15) dapat dinyatakan juga menjadi

$$\mathbf{f}_\lambda = \mathbf{Y} - \mathbf{X}\beta - \mathbf{M}\mathbf{c} + \sum \mathbf{c}$$

Dalam hal ini, \mathbf{f}_λ dapat dilihat sebagai hasil smoothing terhadap $\mathbf{Y} - \mathbf{X}\beta$, sehingga secara umum menjadi

$$\mathbf{f}_\lambda = \mathbf{A}(\lambda)(\mathbf{Y} - \mathbf{X}\beta)$$

$$\begin{aligned} \mathbf{f}_\lambda &= \mathbf{Y} - \mathbf{X}\beta - (\sum + n\lambda - \sum) \mathbf{c} \\ &= \mathbf{Y} - \mathbf{X}\beta - n\lambda \mathbf{c} \end{aligned}$$

$$\begin{aligned} n\lambda \mathbf{c} &= \mathbf{Y} - \mathbf{X}\beta - \mathbf{f}_\lambda \\ &= \mathbf{Y} - \mathbf{X}\beta - \mathbf{A}(\lambda)(\mathbf{Y} - \mathbf{X}\beta) \\ &= (\mathbf{I} - \mathbf{A}(\lambda))(\mathbf{Y} - \mathbf{X}\beta) \end{aligned}$$

Di lain pihak,

$$\begin{aligned} \sum \mathbf{c} &= \mathbf{f}_\lambda - \mathbf{T}\mathbf{d} \\ &= \mathbf{A}(\lambda)(\mathbf{Y} - \mathbf{X}\beta) - \mathbf{T}\mathbf{d} \\ n\lambda \mathbf{c}' \sum \mathbf{c} &= n\lambda \mathbf{c}' \mathbf{A}(\lambda)(\mathbf{Y} - \mathbf{X}\beta) \\ &= (\mathbf{I} - \mathbf{A}(\lambda))(\mathbf{Y} - \mathbf{X}\beta) \mathbf{A}(\lambda)(\mathbf{Y} - \mathbf{X}\beta) \end{aligned}$$

$$S(\beta, \mathbf{f}_\lambda) = (\mathbf{Y} - \mathbf{X}\beta - \mathbf{f}_\lambda)^T (\mathbf{Y} - \mathbf{X}\beta - \mathbf{f}_\lambda) + n\lambda \mathbf{c}' \sum \mathbf{c}$$

$$= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{I} - \mathbf{A}(\lambda))^T (\mathbf{I} - \mathbf{A}(\lambda))(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + n\lambda \mathbf{c}' \Sigma \mathbf{c} = (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{I} - \mathbf{A}(\lambda))^T (\mathbf{I} - \mathbf{A}(\lambda))(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \\ (\mathbf{I} - \mathbf{A}(\lambda))(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})\mathbf{A}(\lambda)(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})$$

$$\frac{S(\boldsymbol{\beta}, \mathbf{f}_\lambda)}{\partial \boldsymbol{\beta}} = 0$$

$$\boldsymbol{\beta}_\lambda = (\mathbf{X}'(\mathbf{I} - \mathbf{A}(\lambda))\mathbf{X})^{-1} \mathbf{X}'(\mathbf{I} - \mathbf{A}(\lambda))' \mathbf{Y}$$