

# Reproducing Kernel dan Estimator Spline dalam Regresi Nonparametrik

Ruliana\*

## Abstrak

Spline merupakan polinomial yang memiliki sifat tersegmen yang memberikan fleksibel lebih dari polinomial biasa, yang bisa menyesuaikan diri secara efektif terhadap karakteristik suatu data. *Estimator Spline* didapatkan dari meminimalkan kriteria *penalized least square* yaitu kriteria kesesuaian kurva terhadap data (*goodness of fit*) dengan kekasaran kurva (*roughness penalty*), dimana keseimbangan antara kesesuaian kurva ini dikontrol oleh suatu parameter penghalus  $\lambda$ , dimana  $\lambda > 0$ . Jika suatu data berpasangan  $\{(t_i, y_i), i = 1, 2, \dots, n\}$ ,  $t_i \in \mathbb{R}$ ,  $y_i \in \mathbb{R}$  dan hubungan antara  $t_i$  dan  $y_i$ , diasumsikan mengikuti model  $y_i = L_i f + \varepsilon_i$ ,  $i = 1, 2, \dots, n$ , dengan  $f \in H_R$ , dan  $L_i$ ,  $i = 1, 2, \dots, n$  adalah fungsional linear terbatas pada  $H_R$  dimana  $H_R = H_0 \oplus H_1$ , maka estimasi dari  $f$  adalah  $\hat{f} = \sum_{v=1}^m \alpha_v \varphi_v + \sum_{i=1}^n \beta_i \xi_i \in \Pi^{2m-1}$ . yang merupakan *polynomial natural spline*. Estimator ini dipengaruhi oleh pemilihan  $\lambda$ . Dengan menggunakan data simulasi dalam kasus ini diperoleh bahwa  $I_{GCV} = \frac{L(\lambda_{GCV})}{L(\lambda_{opt})}$ , lebih efisien dari  $I_{CV} = \frac{L(\lambda_{CV})}{L(\lambda_{opt})}$  dalam mendapatkan parameter penghalus optimal.

**Kata Kunci:** Spline, Penalized Least Square, reproducing Kernel, parameter penghalus, GCV, CV, asimtotik efisien.

## 1. Pendahuluan

Analisis regresi merupakan salah satu alat statistik yang banyak digunakan untuk mengetahui hubungan antara sepasang variabel atau lebih. Misalkan diberikan data  $\{(t_i, y_i), i = 1, 2, \dots, n\}$ ,  $t_i \in \mathbb{R}$ ,  $y_i \in \mathbb{R}$  dan hubungan antara  $t_i$  dan  $y_i$  diasumsikan mengikuti model regresi

$$y_i = f(t_i) + \varepsilon_i, t_i \in [0, 1], i = 1, 2, \dots, n, \quad (1)$$

dimana  $f(t_i)$  adalah fungsi regresi dan  $\varepsilon_i$  adalah sesatan random yang diasumsikan berdistribusi normal independen dengan mean nol dan variansi  $\sigma^2$ .

Ada dua metode yang dapat digunakan untuk menaksir fungsi  $f$ , yaitu metode regresi parametrik dan metode regresi nonparametrik. Pertimbangan untuk memilih metode mana yang akan digunakan ada kaitannya dengan asumsi dari fungsi  $f$  tersebut. Metode regresi parametrik akan sesuai jika bentuk fungsi  $f$  diketahui atau ada sumber lain yang bisa digunakan untuk

\* Jurusan Matematika, Fakultas Matematika dan Ilmu Pengetahuan Alam, Universitas Negeri Makassar, email: ruliana@yahoo.co.id

menentukan bentuk fungsi  $f$  tersebut. Tetapi jika fungsi  $f$  tersebut tidak diketahui bentuknya, maka metode regresi nonparametrik lebih sesuai digunakan (Eubank, 1988). Dalam hal ini fungsi  $f$  hanya diasumsikan termuat dalam suatu ruang fungsi tertentu, dimana pemilihan ruang fungsi tersebut biasanya dimotivasi oleh sifat kemulusan (*smoothness*) yang dimiliki oleh fungsi  $f$  tersebut.

Beberapa pendekatan nonparametrik yang cukup populer dalam mengestimasi fungsi  $f$  antara lain *Estimator Deret Orthogonal* (Eubank 1988), *Spline* (Wahba, 1990), *Estimator Kernel* (Härdle, 1990), *Histogram* (Green, 1994). Khusus untuk *Spline*, pendekatan ini merupakan *piecewise polynomial*, yaitu polinomial yang memiliki sifat tersegmen. Sifat inilah yang memberikan fleksibel lebih daripada polinomial biasa, sehingga memungkinkan untuk menyesuaikan diri secara efektif terhadap karakteristik lokal dari fungsi atau data.

Dalam tulisan ini pendekatan regresi yang digunakan adalah regresi nonparametrik *spline*. Dari bentuk (1), apabila fungsi  $f \in W_2^m[0,1] = \{f : f^{(k)}, k = 0, 1, \dots, m-1, f^{(m)} \in L_2[0,1]\}$ , dimana  $L_2[0,1]$  menyatakan himpunan fungsi-fungsi kuadrat terintegral pada interval  $[0,1]$ , dengan kesesuaian kurva terhadap data adalah  $n^{-1} \sum_{i=1}^n (y_i - f(t_i))^2$  dan kekasaran kurva adalah

$\int_0^1 (f''(t))^2 dt$  maka estimasi  $f$  diperoleh dengan meminimumkan *Penalized Least Square*,

$$n^{-1} \sum_{i=1}^n \{y_i - f(t_i)\}^2 + \lambda \int_0^1 \{f''(t)\}^2 dx, \lambda > 0. \quad (2)$$

Permasalahan optimasi pada (2) dapat diselesaikan dengan pendekatan *Reproducing Kernel* (Kemildrorf & Wahba, 1971). Sifat dari *Reproducing Kernel* adalah dapat ditentukan representasi dari suatu fungsional linear terbatas, sehingga kurva regresi  $f \in W_2^m[0,1]$  yang merupakan penyelesaian yang optimal dari persamaan (2) bisa diperoleh (Wahba, 2000).

Dalam menghasilkan estimasi kurva regresi yang baik, pemilihan  $\lambda$  yang optimal merupakan hal yang penting. Dengan menggunakan *Reproducing Kernel* untuk mengestimasi kurva regresi pada (1), selanjutnya akan dipilih suatu nilai  $\lambda$  yang optimal. Beberapa metode untuk memilih  $\lambda$  yaitu metode *Cross Validation* (CV), *Generalized Cross Validation* (GCV), metode *Generalized Maximum Likelihood* atau GML (Kou, 2003).

Dalam tulisan ini dibahas penyelesaian optimal dari (2) dengan menggunakan *Reproducing Kernel*, dan pemilihan parameter penghalus  $\lambda$  dengan menggunakan metode GCV dan CV serta melihat sifat efisien dari kedua metode tersebut dengan menggunakan data simulasi.

## 2. Spline dalam Regresi Nonparametrik

Ada beberapa pendekatan nonparametrik yang cukup populer dalam mengestimasi  $f$ , salah satunya adalah *spline*. *Spline* adalah potongan polinomial order  $r$ . Titik bersama dari potongan-potongan tersebut disebut dengan *knots*. *Spline* order  $r$  dengan *knots* pada  $\zeta_1, \dots, \zeta_k$  diberikan dalam fungsi  $S$  dengan bentuk:

$$S(t) = \sum_{i=0}^{r-1} \theta_i t^i + \sum_{j=1}^k \delta_j (t - \zeta_j)_+^{r-1} \quad (3)$$

dengan  $(t - \zeta_j)_+^{r-1} = \begin{cases} (t - \zeta_j) & , \text{jika}(t - \zeta_j) \geq 0 \\ 0 & , \text{jika}(t - \zeta_j) < 0 \end{cases}$  (Eubank, 1988).

Spline (3) mempunyai sifat sebagai berikut:

1.  $S$  merupakan potongan polinomial derajat  $r-1$  pada setiap subinterval  $[\zeta_j, \zeta_{j+1}]$ .
2.  $S$  mempunyai turunan ke  $(r-2)$  yang kontinu.
3.  $S$  mempunyai turunan ke  $(r-1)$  yang merupakan fungsi tangga dengan titik-titik lompatan pada  $(\zeta_1, \dots, \zeta_k)$ .

### 3. Reproducing Kernel

**Definisi 1.** (Wahba, 2000)

*Reproducing Kernel* dari  $H$  adalah suatu fungsi  $R$  yang didefinisikan pada  $[0,1] \times [0,1]$  sedemikian hingga untuk setiap titik tetap  $t \in [0,1]$  berlaku  $R_t \in H$  dengan  $R_t(s) = R(s, t)$  dan

$$L_t f = \langle R_t, f \rangle = f(t). \quad (4)$$

Didefinisikan ruang  $H_0$  memuat koleksi fungsi  $H_0 = \{ \mathbf{h} : \mathbf{h}^{(i-1)}(0) = (0) , i = 1, \dots, m, \mathbf{h}^{(i-1)}$  kontinu absolut pada  $[0,1]$  dan  $\mathbf{h}^m \in [0,1] \}$ .  $H_0$  adalah ruang inner product yang lengkap yang disebut *Hilbert space*.  $H_0$  merupakan *Reproducing Kernel Hilbert Space* karena untuk setiap  $\mathbf{h} \in H_0$  dan pemetaan  $L_t : H_0 \rightarrow R, \mathbf{h} \rightarrow L_t(\mathbf{h}) = 0$   $L_t$  memenuhi sifat linearitas, dengan domain  $[0,1]$  didalam ruang  $H_0$ , dengan *reproducing kernel* dari  $H_0$  adalah  $R_t^0 \in H_0$  yaitu  $R_t(s) = \sum_{i=1}^m \phi_i(s) \phi_i(t)$ , dan untuk titik tetap  $t$  berlaku  $L_t \mathbf{h} = \langle R_t^0, \mathbf{h} \rangle = h(t), \mathbf{h} \in H_0$ .

Ruang  $H_1$  memuat koleksi fungsi  $H_1 = \{ g : g^{(k)}(0) = (0), k = 0, 1, \dots, m-1, \mathbf{g}^{(k)}$  kontinu absolut pada  $[0,1]$  dan  $\mathbf{g}^m \in L_2[0,1] \}$ .  $H_1$  adalah *Reproducing Kernel Hilbert Space* dengan *Inner Product*

$$\langle \mathbf{g}_1, \mathbf{g}_2 \rangle_1 = \int_0^1 \mathbf{g}_1^{(m)}(t) \mathbf{g}_2^{(m)}(t) dt \quad (5)$$

dan *Reproducing Kernel* dari  $H_1$  adalah

$$R_1(s, t) = \int_0^1 \frac{(s-u)_+^{m-1} (t-u)_+^{m-1}}{((m-1)!)^2} du \quad (6)$$

$H$  adalah ruang Hilbert dan terdapat dengan tunggal *reproducing kernel*  $R_t$  untuk titik  $t \in [0,1]$  dalam  $H$  dan  $L_t f = \langle R_t, f \rangle = f(t), f \in H, L_t$  adalah fungsional linear terbatas pada  $H$  yang memetakan fungsi  $f$  didalam ruang  $H$  kebilangan riil,  $L_t : f \rightarrow f(t)$ .

Dalam masalah regresi, untuk data berpasangan  $\{(t_i, y_i), i = 1, 2, \dots, n\}, t_i \in R, y_i \in R$  dan hubungan antara  $t_i$  dan  $y_i$  diasumsikan mengikuti model regresi secara umum yaitu  $y_i = L_t f + \varepsilon_i, i = 1, 2, \dots, n$ , dengan  $f$  adalah fungsi regresi dan  $\varepsilon_i$  adalah sesatan random yang

diasumsikan berdistribusi independen dengan mean nol dan varians  $\sigma^2$ . Menaksir fungsi  $f$  adalah mencari  $f$  yang ada didalam ruang Hilbert yang meminimumkan  $n^{-1} \sum_{i=1}^n (y_i - L_i f)^2 + \|P_1 f\|_R^2$ .

#### 4. Estimator Spline

Misalkan diberikan model data

$$y_i = L_i f + \varepsilon_i \quad i = 1, 2, \dots, n, \quad (7)$$

dengan  $f \in H_R$  dan  $L_i \quad i = 1, 2, \dots, n$  adalah fungsional linear terbatas pada  $H_R$  serta  $H_R$  mempunyai dekomposisi  $H_R = H_0 \oplus H_1$ , estimasi dari  $f$  adalah  $f \in H_R$  yang meminimumkan

$$\frac{1}{n} \sum_{i=1}^n (y_i - L_i f)^2 + \lambda \|P_1 f\|_R^2 \quad (\text{Wahba, 1990}). \quad (8)$$

**Teorema 1.** (Wahba, 1990)

Apabila  $H_R = H_0 \oplus H_1$  dan  $\phi_1, \dots, \phi_m$  basis dari ruang  $H_0$  serta  $T_{nxm}$  adalah matriks full rank kolom berorder  $nxm$  yang diberikan oleh:

$$T_{nxm} = \{L_i \phi_v\}, \quad i = 1, 2, \dots, n; \text{ dan } v = 1, 2, \dots, m$$

maka  $\hat{f}$  yang meminimumkan  $\frac{1}{n} \sum_{i=1}^n (y_i - \langle \eta_i, f \rangle)^2 + \lambda \|P_1 f\|_R^2$  (9)

$$\text{adalah } \hat{f} = \sum_{v=1}^m \alpha_v \phi_v + \sum_{i=1}^n \beta_i \xi_i$$

$$\text{dengan } \alpha = (\alpha_1, \dots, \alpha_m)' = (T' M^{-1} T)^{-1} T' M^{-1} y$$

$$\beta = (\beta_1, \dots, \beta_n)' = M^{-1} (I - T (T' M^{-1} T)^{-1} T' M^{-1}) y$$

$$\xi_i = P_1 \eta_i$$

$$M = \Sigma + n \lambda I$$

$$\Sigma = \{ \langle \xi_i, \xi_j \rangle \}, \quad i, j = 1, 2, 3, \dots, n$$

**Bukti:**

$\hat{f}$  dapat ditulis dalam bentuk

$$\hat{f} = \sum_{v=1}^m \alpha_v \phi_v + \sum_{i=1}^n \beta_i \xi_i = \phi' \alpha + \xi' \beta + \psi, \quad \psi \in H_R$$

yang tegak lurus dengan  $\phi_1, \dots, \phi_m, \xi_1, \dots, \xi_n$ .

$$\hat{f} = \sum_{v=1}^m \alpha_v \phi_v + \sum_{i=1}^n \beta_i \xi_i = \phi' \alpha + \xi' \beta$$

$$\{ \langle L_i f \rangle \} = \{ \langle \eta_i, \hat{f} \rangle \}; \quad \eta_i \in H_R = H_0 + H_1$$

$$= \{ \langle \eta_i, \phi' \alpha + \xi' \beta \rangle \} = \{ \langle \eta_i, \phi' \alpha \rangle \} + \{ \langle \eta_i, \xi' \beta \rangle \}, \quad i = 1, \dots, n$$

$$= T\alpha + \Sigma\beta \quad (10)$$

$$\begin{aligned} \|Pf\|_R^2 &= \langle Pf, Pf \rangle \\ &= \langle P(\phi' \alpha + \xi' \beta), P(\phi' \alpha + \xi' \beta) \rangle \\ &= \beta' \Sigma \beta \end{aligned} \quad (11)$$

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (y_i - L_i f)^2 &= \frac{1}{n} (y - Lf)' (y - Lf), \text{ dari persamaan (10) selanjutnya} \\ &= \|(y - T\alpha - \Sigma\beta)\|^2 \end{aligned} \quad (12)$$

Dari persamaan (11) dan (12), maka persamaan (9) dapat ditulis sebagai

$$\mathfrak{R}(\alpha, \beta) = \frac{1}{n} \|(y - T\alpha - \Sigma\beta)\|^2 + \lambda \beta' \Sigma \beta \quad (13)$$

Selanjutnya akan dicari  $\alpha$  dan  $\beta$  yang meminimumkan persamaan (13)

$$\begin{aligned} \mathfrak{R}(\alpha, \beta) &= \frac{1}{n} \|(y - T\alpha - \Sigma\beta)\|^2 + \lambda \beta' \Sigma \beta \\ \mathfrak{R}(\alpha, \beta) &= \|(y - T\alpha - \Sigma\beta)\|^2 + n\lambda \beta' \Sigma \beta \end{aligned} \quad (14)$$

Dengan menurunkan secara parsial persamaan (14) terhadap  $\beta$  dan hasilnya disamakan dengan nol

$$\begin{aligned} \frac{\partial \mathfrak{R}}{\partial \beta} &= -2\Sigma' y + 2\Sigma' T\alpha + 2(\Sigma' \Sigma + n\lambda \Sigma) \beta = -\Sigma' y + \Sigma' T\alpha + (\Sigma' \Sigma + n\lambda \Sigma) \beta = 0 \\ \beta &= M^{-1} (y - T\alpha) \end{aligned} \quad (15)$$

Dengan cara yang sama persamaan (14) diturunkan terhadap  $\alpha$  di dapatkan

$$\begin{aligned} \frac{\partial \mathfrak{R}}{\partial \alpha} &= -2T' y + 2T' T\alpha + 2T' \Sigma \beta = 0 \\ -T' y + T' T\alpha + T' \Sigma \beta &= 0. \end{aligned} \quad (16)$$

Sehingga

$$\begin{aligned} \alpha &= (T' M^{-1} T)^{-1} T' M^{-1} y \\ M &= \Sigma + n\lambda I. \end{aligned} \quad (17)$$

Dari persamaan (17) disubstitusi pada persamaan (15)

$$\begin{aligned} \beta &= M^{-1} (y - T\alpha) \\ \beta &= M^{-1} (I - T(T' M^{-1} T)^{-1} T' M^{-1}) y. \end{aligned} \quad (18)$$

Sehingga terbukti bahwa  $\hat{f}$  yang meminimumkan

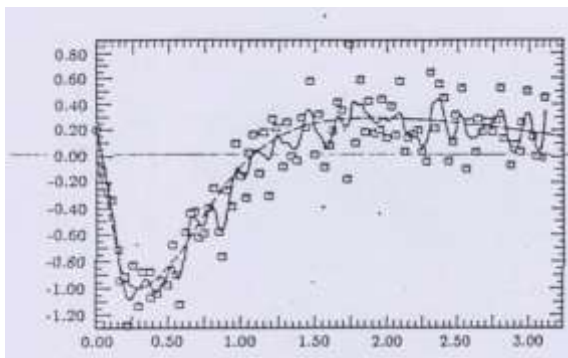
$$\frac{1}{n} \sum_{i=1}^n (y_i - \langle \eta_i, f \rangle)^2 + \lambda \|P_1 f\|_R^2 \text{ adalah } f_\lambda = \sum_{v=1}^m \alpha_v \varphi_v + \sum_{i=1}^n \beta_i \xi_i,$$

dengan  $\alpha$  dan  $\beta$  yaitu

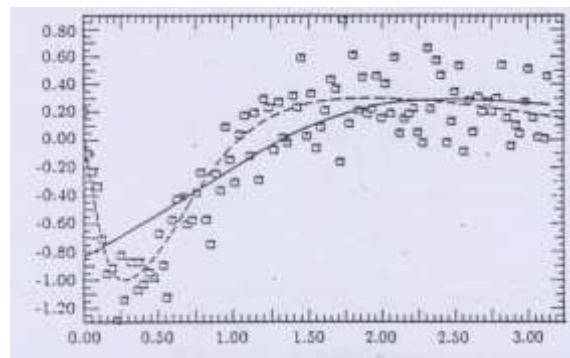
$$\alpha = (T' M^{-1} T)^{-1} T' M^{-1} y$$

$$\beta = M^{-1} (I - T(T' M^{-1} T)^{-1} T' M^{-1}) y.$$

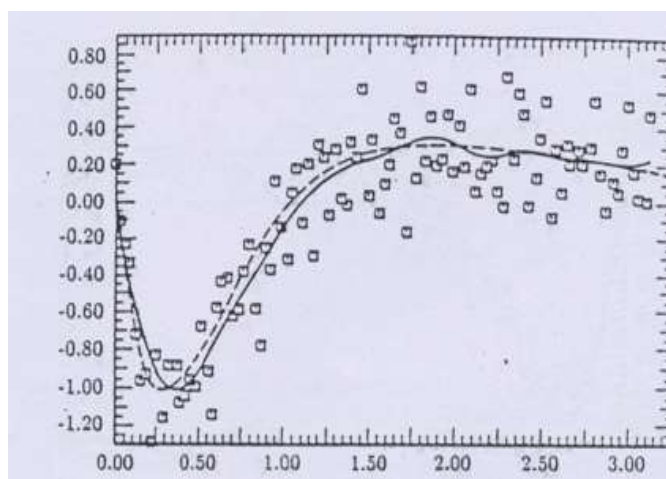
Bentuk estimator spline, dipengaruhi oleh parameter penghalus  $\lambda$ . nilai  $\lambda$  yang sangat kecil atau besar akan memberikan bentuk fungsi penyelesaian yang sangat kasar atau sangat mulus (Wahba 1990; Eubank, 1988).



Gambar 1. Estimasi Data dengan  $\lambda$  yang Kecil.



Gambar 2. Estimasi Data dengan  $\lambda$  yang Besar



Gambar 3. Estimasi Data dengan  $\lambda$  Optimal.

## 5. Pemilihan Parameter Penghalus

Dalam regresi nonparametrik, bentuk fungsi  $f$  tidak diketahui, fungsi  $f$  diasumsikan *smooth* dalam arti  $f$  merupakan anggota ruang Sobolev  $W_2^m[0,1] = \{f : f^{(k)}, k = 0,1,\dots,m-1, \text{ kontinu absolute pada } [0,1] \text{ dan } f^{(m)} \in L_2[0,1]\}$  dengan  $L_2[0,1]$  menyatakan himpunan fungsi-fungsi kuadrat terintegral pada interval  $[0,1]$ .

Idealnya akan dipilih suatu nilai  $\lambda$  yang meminimumkan fungsi kerugian  $L(\lambda)$ , akan tetapi dalam regresi nonparametrik tidak dapat dilakukan secara nyata sebab  $L(\lambda)$  masih memuat

fungsi  $f$  yang tidak diketahui. Sehingga perlu mengestimasi data dan kemudian estimatornya diminimumkan terhadap  $\lambda$  untuk mendapat estimator  $f$  yang paling baik (Eubank, 1988). Salah satu pemilihan parameter penghalus  $\lambda$  adalah menggunakan metode *Generalized Cross Validation* (GCV). Ide dasar dari GCV adalah memodifikasi *Cross Validation* atau CV (Green & Silverman 1994).

**Definisi 2.** (Eubank, 1988)

Jika  $f_\lambda = (f_{\lambda 1}, \dots, f_{\lambda n})^T$  adalah estimator untuk  $f$  maka fungsi kerugian (loss function)

didefinisikan sebagai  $L(\lambda) = n^{-1} \sum_{i=1}^n (f_i - f_{\lambda i})^2$ .

**Teorema 2.** (Craven dan Wahba, 1979)

Jika  $f \in W_2^m[0,1]$  dan  $\lambda_{GCV}$ ,  $\lambda_{res}$  menyatakan nilai  $\lambda$  yang meminimumkan  $EGCV(\lambda)$  dan  $EL(\lambda)$ , maka terdapat suatu barisan  $\lambda_{GCV}$  sedemikian sehingga

$$\frac{EL(\lambda_{GCV})}{EL(\lambda)} \rightarrow 1, n \rightarrow \infty.$$

Jika  $\sigma^2$  diketahui maka  $\lambda$  optimal bisa didapat langsung dari kriteria prediksi *mean square error* atau fungsi kerugian yang didefinisikan oleh

$$L(\lambda) = n^{-1} \sum_{i=1}^n (f_i - f_{\lambda i})^2 \quad (19)$$

fungsi dari persamaan (19) memberikan penaksiran dari kebaikan suatu estimator. Dalam hal  $\sigma^2$  tidak diketahui maka dapat digunakan metode *Cross Validation* (CV) dan *Generalized Cross Validation* (GCV), untuk mendapatkan nilai  $\lambda$  optimal. *Cross validation* (CV) adalah metode untuk memilih  $\lambda$  yang meminimumkan

$$CV(\lambda) = n^{-1} \sum_{j=1}^n (y_j - f_\lambda^{[j]}(t))^2 \quad (20)$$

$f_\lambda^{[j]}(t)$  adalah  $f \in W_2^m[0,1]$  yang meminimumkan

$$n^{-1} \sum_{i=1}^n \{y_i - f(t_i)\}^2 + \lambda \int_0^1 \{f''(t)\}^2 dx.$$

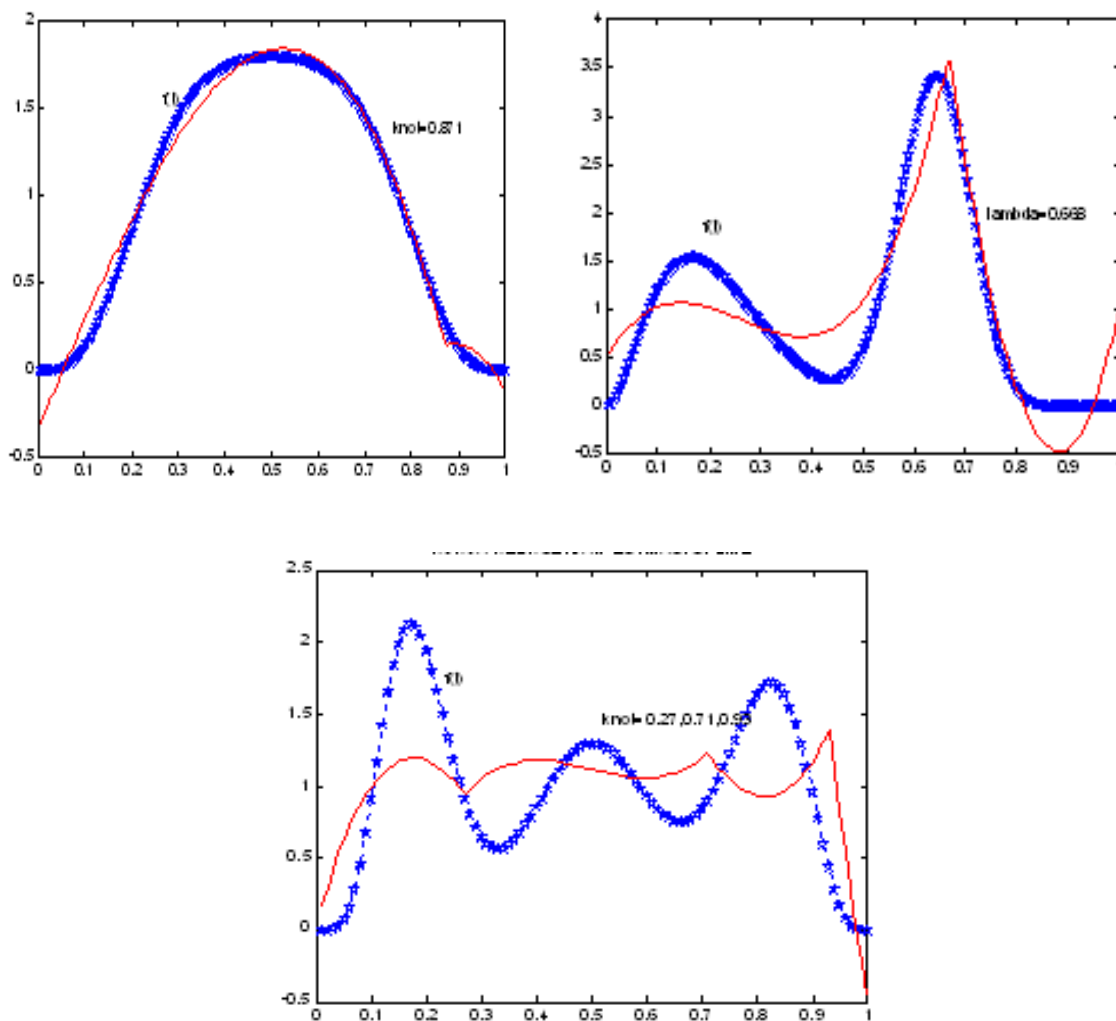
Sedangkan *Generalized Cross Validation* (GCV), adalah metode untuk memilih  $\lambda$  yang meminimumkan

$$GCV(\lambda) = n^{-1} \sum_{j=1}^n \frac{(y_j - f_\lambda(t_j))^2}{\left(1 - n^{-1} \sum_{i=1}^n a_{ij}(\lambda)\right)^2} \quad (5.3)$$

## 6. Studi Simulasi Metode GCV dan CV dalam Mendapatkan $\lambda$ Optimal

Dalam penelitian ini, diselidiki perilaku optimal estimator spline berdasarkan data simulasi, dengan metode GCV dan CV. Data dibangkitkan dari distribusi Beta dengan berbagai nilai  $a$  dan  $b$ :

$$B(a,b) = \frac{\tau(a+b)}{\tau(a)\tau(b)} x^{a-1} (1-x)^{b-1} I_{(0,1)}(x).$$



Gambar 4. Plot Tiga Fungsi Beta dengan  $a$  dan  $b$  yang Berbeda.

Dalam menentukan  $\lambda$  optimal dengan metode *Generalized Cross Validation* (GCV) dan *Cross Validation CV*, dimana  $I_{GCV} = \frac{L(\lambda_{GCV})}{L(\lambda_{opt})}$ ,  $I_{CV} = \frac{L(\lambda_{CV})}{L(\lambda_{opt})}$  didapatkan hasil yang dapat dilihat dalam Tabel 1.



**Tabel 1.** Hasil Simulasi Lambda Optimal.

$\sigma^2$	$\lambda$ Optimal	n=25		n=50		n=100		n=200	
		$\lambda$	nilai	$\lambda$	Nilai	$\lambda$	Nilai	$\lambda$	nilai
$\sigma^2=0,0125$	$\lambda$ (GCV)	0,715	0,0077	0,720	0,0074	0,714	0,0069	0,715	0,0068
	$\lambda$ (CV)	0,701	0,0106	0,702	0,0081	0,704	0,0071	0,712	0,0069
	$\lambda$ (LOSS)	0,730	0,0047	0,723	0,0057	0,715	0,0061	0,714	0,0063
$\sigma^2=0,025$	$\lambda$ (GCV)	0,734	0,0078	0,721	0,0079	0,718	0,0075	0,712	0,0072
	$\lambda$ (CV)	0,718	0,0105	0,709	0,0086	0,712	0,0078	0,709	0,0073
	$\lambda$ (LOSS)	0,73	0,0048	0,718	0,0058	0,721	0,0062	0,715	0,0064
$\sigma^2=0,05$	$\lambda$ (GCV)	0,723	0,0101	0,732	0,0097	0,714	0,0091	0,717	0,0093
	$\lambda$ (CV)	0,696	0,0127	0,714	0,0107	0,702	0,0091	0,708	0,0094
	$\lambda$ (LOSS)	0,729	0,0052	0,717	0,0061	0,72	0,0062	0,714	0,0064
$\sigma^2=0,1$	$\lambda$ (GCV)	0,687	0,0201	0,696	0,0194	0,702	0,0172	0,717	0,0163
	$\lambda$ (CV)	0,675	0,0213	0,678	0,0208	0,693	0,0173	0,714	0,0164
	$\lambda$ (LOSS)	0,729	0,0074	0,72	0,0063	0,717	0,0064	0,72	0,0066

**Tabel 2.** Ratio GCV dan Ratio CV.

$\sigma^2$	Ratio	n=25	n=50	n=100	n=200
0,0125	$W_r$ (GCV)	1,562904	1,236415	1,10849220	1,05192572
	$W_r$ (CV)	1,648624	1,220889	1,05407279	1,04367263
0,025	$W_r$ (GCV)	1,564278	1,235004	1,10814558	1,05192572
	$W_r$ (CV)	1,649797	1,219301	1,09358752	1,04367263
0,05	$W_r$ (GCV)	1,572061	1,234827	1,10814558	1,05570839
	$W_r$ (CV)	1,655471	1,219125	1,09341421	1,04745530
0,1	$W_r$ (GCV)	1,587077	1,253831	1,12623377	1,06497100
	$W_r$ (CV)	1,678571	1,239211	1,11168831	1,05672100

Pada Tabel 2 terlihat bahwa, dengan bertambahnya ukuran sampel, nilai ratio dari GCV semakin mendekati 1, begitu juga untuk nilai ratio dari CV, namun dalam hal ini ratio GCV selalu lebih kecil dari ratio CV untuk semua ukuran sampel  $n$  dan  $\sigma^2$  yang diberikan.

**Tabel 3.** Nilai  $I_{GCV}$ , dan  $I_{CV}$ .

$\sigma^2$	efisiensi	n=25	n=50	n=100	n=200
<b>0,0125</b>	$I_{GCV}$	1,021277	1,000000	1,000000	1,000000
	$I_{CV}$	1,042553	1,017544	1,016393	1,000000
<b>0,025</b>	$I_{GCV}$	1,000000	1,000000	1,000000	1,000000
	$I_{CV}$	1,020833	1,000000	1,000000	1,000000
<b>0,005</b>	$I_{GCV}$	1,019231	1,000000	1,000000	1,000000
	$I_{CV}$	1,000000	1,000000	1,016129	1,016129
<b>0,1</b>	$I_{GCV}$	1,040541	1,000000	1,000000	1,000000
	$I_{CV}$	1,054054	1,01587	1,000000	1,000000

Dari Tabel 3 terlihat bahwa nilai  $I_{GCV}$  untuk ukuran sampel yang besar lebih banyak yang mendekati 1 dibandingkan dengan nilai  $I_{CV}$ .

Hasil simulasi ini adalah dengan melihat ratio dan efisiensi maka menggunakan metode *Generalized Cross Validation* (GCV) didalam mencari  $\lambda$  yang paling optimal adalah lebih efisien dari metode *Cross Validation* (CV).

## 6. Kesimpulan

Kesimpulan yang diperoleh dari hasil penelitian ini, yaitu bahwa apabila  $H_R = H_0 \oplus H_1$  dan  $\phi_1, \dots, \phi_m$  span dari ruang  $H_0$  serta  $T_{nxm}$  adalah matriks full rank kolom berorder  $nxm$  yang diberikan oleh  $T_{nxm} = \{L_i \phi_v\}$ ,  $i = 1, 2, \dots, n$  dan  $v = 1, 2, \dots, m$ , maka  $\hat{f}$  yang meminimumkan  $\frac{1}{n} \sum_{i=1}^n (y_i - \langle \eta_i, f \rangle)^2 + \lambda \|P_1 f\|_R^2$  adalah  $\hat{f} = \sum_{v=1}^m \alpha_v \phi_v + \sum_{i=1}^n \beta_i \xi_i$ , dengan  $\alpha = (\alpha_1, \dots, \alpha_m)'$  dan  $\beta = (\beta_1, \dots, \beta_n)'$   $= (T' M^{-1} T)^{-1} T' M^{-1} y$  dan  $\beta = M^{-1} (I - T(T' M^{-1} T)^{-1} T' M^{-1}) y$ . Dengan melihat ratio dan efisiensi maka menggunakan metode *Generalized Cross Validation* (GCV) didalam mencari  $\lambda$  yang paling optimal adalah lebih efisien dari metode *Cross Validation* (CV).

## Daftar Pustaka

- [1] Berberian, K. S., 1961, *Introduction to Hilbert Space*, New York.
- [2] Budiantara, I. N. dan Subanar, 1997, Pemilihan parameter penghalus dalam regresi spline terbobot, *Majalah Ilmiah Matematika dan IPA-UGM*, Yogyakarta.
- [3] Cox, D., 1981, Asymptotics for M-Type smoothing splines, *Technical Report*, 654, University of Wisconsin.
- [4] Craven, P. dan Wahba, G., 1979, Smoothing noisy data with spline functions estimating the correct degree of smoothing by method of Generalized Cross-Validation, *Journal Numerische Mathematics*, 31, 377-403.
- [5] Eubank, R. L., 1988, *Spline smoothing and Nonparametrik Regression*, Marcel Dekker, New York.
- [6] Green, P. J. dan Silverman, B.W., 1994, *Nonparametrik Regression and Generalized Linear Models (a Roughness Penalty Approach)*, Chapman & Hall, New York.
- [7] Hardle, G., 1990, *Applied Nonparametrik Regression*, Cambridge University Press, New York.
- [8] Kimeldorf, G. dan Wahba, G., 1971, Some result on Tchebycheffian spline function, *Journal of Mathematical Analysis and Application*, 33, 82-95.
- [9] Kou, S.C., 2003, On efficiency of selection criteria in spline regression, *Probab, Theory Relat. Fields*, 127, 153-176.
- [10] Kreyszig, E., 1978, *Introductory Functional Analysis with Applications*, University of Windsor, John Wiley & Sons, New York.
- [11] Mallows, C., 1973, Some comments on Cp, *Journal Technometrics*, 15, 661-675.

- [12] Oeherlet, G. W., 1992, Relaxed boundary smoothing splines, *Journal the Annals of Statistics*, 20, 146-160.
- [13] Rohatgi, V. K., 1976, *An Introduction to Probability Theory and Mathematical Statistics*, John Wiley & Sons, New York.
- [14] Silverman, B., 1984, A fast and efficient cross-validation method for smoothing parameter choice in spline regression, *Journal of the American Statistical Association*, 79, 584-589.
- [15] Silverman, B., 1985, some aspects of the spline smoothing approach to non parametric regression curve fitting (With discussion), *Journal Royal Statistical Society*, 47, 1-52.
- [16] Speckman, P., 1980, Minimax estimates of linear functionals in a Hilbert space, *Unpublished Manuscript*.
- [17] Speckman, P., 1985, Spline smoothing and optimal rates of convergence in nonparametric regression models, *Journal Ann. Statist*, 13, 970-983.
- [18] Wahba, G., 1975, Smoothing noisy data by spline functions, *Journal Numer. Math*, 24, 383-393.
- [19] Wahba, G., 1985, A comparison of GCV and GML for choosing the smoothing parameter in the generalized spline smoothing problem, *Journal the Annals of Statistics*, 13, 1378-1402.
- [20] Wahba, G., 1990, Spline model for observational data, *Society For Industrial and Applied Mathematics*, Philadelphia.
- [21] Wahba, G., 2000, An introduction to model building with reproducing kernel Hilbert spaces, *Technical Report* , 1020, University of Wisconsin.