

**Pemilihan Titik Knot Optimal Menggunakan Metode GCV Dalam Regresi Nonparametrik Spline *Truncated***

**Fitri Kusunartutik\*, Nur Karomah Dwidayati**

Jurusan Matematika, Fakultas MIPA, Universitas Negeri Semarang, Indonesia  
Gedung D7 Lt. 1, Kampus Sekaran Gunungpati, Semarang, 50229

E-mail: fitrikusunart@gmail.com

Diterima 9 April 2022

Disetujui 11 Agustus 2022

Dipublikasikan 24 Oktober 2022

**Abstrak**

Penelitian ini bertujuan untuk mengkaji metode GCV dalam pemilihan titik knot optimal pada regresi nonparametrik spline *truncated* yang disimulasikan pada data Indeks Pembangunan Manusia (IPM) di Jawa Tengah serta untuk mengetahui faktor apa saja yang mempengaruhi IPM di Jawa Tengah. Dengan menggunakan metode GCV akan digunakan satu, dua, tiga, empat, lima dan kombinasi titik knot yang akan membangun model. Kemudian berdasarkan nilai GCV minimum akan dihasilkan titik knot optimal. Pemilihan titik knot optimal akan memberikan model regresi yang optimal. Hasil penelitian menunjukkan model terbaik yaitu dengan menggunakan kombinasi titik knot 5-5-5-5-5 dengan nilai GCV terkecil sebesar 0,3815048. Hasil uji parameter menunjukkan kelima variable yang diduga mempengaruhi IPM berpengaruh signifikan terhadap IPM. Koefisien determinasi sebesar 99,94251% yang artinya 99,94251% variansi IPM dapat dijelaskan secara signifikan oleh kelima variable yang mempengaruhi yaitu TPAK, Rasio Murid Guru, Kepadatan Penduduk, Angka Kesakitan dan PDRB.

Kata kunci: GCV, spline *truncated*, knot

**Abstract**

*This study aims to examine the GCV method in selecting the optimal knot point in the simulated spline truncated nonparametric regression on the Human Development Index (HDI) data in Central Java and to determine what factors influence HDI in Central Java. Using the GCV method will be used one, two, three, four, five and a combination of knots that will build the model. Then based on the minimum GCV value the optimal knots will be generated. The selection of optimal knots will provide an optimal regression model. The results showed that the best model was using a combination of 5-5-5-5-5 knots with the smallest GCV value was 0.3815048. The results of the parameter test showed that the five variables that were thought to affect HDI had a significant effect on HDI. The coefficient of determination is 99.94251%, which means that 99.94251% of the HDI variance can be explained significantly by the five influencing variables, namely TPAK, Student Teacher Ratio, Population Density, Morbidity and PDRB.*

Keywords: *Generalized Cross Validation (GCV), spline truncated, knot*

**How to cite:**

Kusunartutik, F., Dwidayati, NK. 2020. Pemilihan titik knot optimal menggunakan metode GCV dalam regresi nonparametrik spline *Truncated*. *Indonesian Journal of Mathematics and Natural Science*, 45(2), 69-76.

**PENDAHULUAN**

Analisis regresi merupakan metode yang digunakan untuk mengetahui pola hubungan antara variable respon dan prediktor. Terdapat banyak kasus dimana pola hubungan antara variable tidak diketahui yang artinya pendekatan regresi yang akan digunakan yaitu regresi nonparametrik. Estimasi kurva regresi nonparametrik diharapkan dapat menyesuaikan pola naik turun data. Ada lima estimator dalam regresi nonparametrik yang biasa digunakan antara lain metode Wavelet, Kernel, Fourier, Multivariate Adaptive Regression Splines (MARS) dan Spline (Budiantara, 2009).

Spline merupakan potongan polinom yang memiliki sifat tersegmen pada titik knot (Eubank, 1999). Spline bersifat fleksible, artinya model regresi akan cenderung mencari sendiri estimasi data mengikuti pola yang bergerak. Spline dapat memodelkan data dengan pola berubah-ubah pada sub-sub interval tertentu, dikarenakan spline merupakan jenis potongan polynomial yang bersifat tersegmen. Dalam membentuk model regresi spline terdapat dua hal yang harus diperhatikan yaitu menentukan orde model dan menentukan banyak dan lokasi titik knot (Montoya *et al.*, 2014). Berdasarkan hasil penelitian Budiantara (2006) menjelaskan bahwa penggunaan spline truncated menghasilkan perhitungan matematik yang relative lebih mudah dan sederhana. Selain itu optimasinya bisa dilakukan tanpa melibatkan *penalty* yaitu dengan menggunakan optimasi kuadrat terkecil (*least square*). Ada empat metode yang bisa digunakan untuk memilih titik knot optimal dalam estimator spline, antara lain metode *Generalized Cross Validation* (GCV), *Unbiased Risk* (UBR), *Generalized Maximum Likelihood* (GML) dan *Cross Validation* (CV) (Hardle, 1990).

Dalam penelitian ini menghasilkan titik knot optimal beserta dengan model regresi spline dengan metode GCV pada data Indeks Pembangunan Manusia (IPM) Jawa Tengah. Hal ini sejalan dengan penelitian terdahulu yang telah dilakukan oleh Demu *et al.* (2017) bahwa data IPM di Indonesia dapat dianalisis menggunakan model regresi nonparametrik spline *truncated*. Scholica *et al.* (2018) mengkaji regresi spline *truncated* dalam pemodelan presentase *unmet need* menggunakan dua titik knot dengan metode GCV. Pada penelitian ini akan dicari faktor yang mempengaruhi IPM sehingga dapat memberikan masukan kepada pemerintah bidang apa saja yang harus ditingkatkan agar nilai IPM dapat ditingkatkan.

### Regresi Nonparametrik

Diberikan  $n$  pengamatan  $(x_i, y_i)$  dengan  $i = 1, 2, 3, \dots, n$  untuk  $x_i$  dan  $y_i$  dalam  $R$ . Variabel  $x_i$  merupakan variabel prediktor pada pengamatan ke- $i$ , variabel  $y_i$  merupakan variabel respon pada pengamatan ke- $i$ , dan  $R$  merupakan bilangan riil. Hubungan antara  $x_i$  dan  $y_i$  diasumsikan mengikuti model regresi berikut.

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

dengan  $f(x_i)$  merupakan fungsi  $f$  yang tidak diketahui dan  $\varepsilon_i$  adalah residual yang diasumsikan berdistribusi normal independen dengan mean nol dan variansi  $\sigma^2$  (Fox, 2000).

### Regresi Nonparametrik Spline *Truncated*

Wahba (1990) mendefinisikan fungsi *spline truncated* derajat  $m$  dengan titik knot  $(K_1, K_2, \dots, K_r)$  dalam fungsi  $f(x)$  dengan persamaan sebagai berikut:

$$f(x) = \sum_{l=0}^m \beta_l x^l + \sum_{k=1}^r \beta_{k+m} (x - K_k)_+^m \quad (2)$$

Fungsi *truncated* yaitu

$$(x - K_k)_+^m = \begin{cases} (x - K_k)^m & , x \geq K_k \\ 0 & , x < K_k \end{cases} \quad (3)$$

### Metode GCV

Menurut Eubank (1988), GCV merupakan modifikasi dari *cross validation* (CV). CV merupakan metode untuk memilih model berdasarkan pada kemampuan prediksi dari model tersebut. Bentuk model dari CV yaitu sebagai berikut:

$$CV = n^{-1} \sum_{i=1}^n \left( \frac{y_i - f(x_i)}{1 - g_{ii}} \right)^2 \quad (4)$$

Dengan  $g_{ii}$  adalah elemen diagonal ke- $i$  dari matriks  $G$ .

Persamaan GCV diperoleh dengan mengganti  $g_{ii}$  pada persamaan dengan  $\sum_{i=1}^n g_{ii} = n^{-1} Tr(I - G)$ . Nilai dari  $Tr(I - G)$  adalah penjumlahan elemen diagonal matriks  $(I - G)$ . Fungsi GCV didefinisikan sebagai :

$$GCV = n^{-1} \frac{\sum_{i=1}^n (y_i - f(x_i))^2}{(n^{-1} \text{Tr}(I - G))^2}$$

$$GCV = \frac{MSE}{(n^{-1} \text{Tr}(I - G))^2} \quad (5)$$

dengan  $n^{-1} \text{Tr}(I - G) < n$  dan  $G = W(W^T W)^{-1} W^T$ .

Penambahan titik knot tidak selalu menghasilkan nilai GCV yang semakin kecil. Nilai GCV yang optimal dihasilkan dari pemilihan titik knot yang optimal (Putri & Dwidayati, 2018)

## METODE

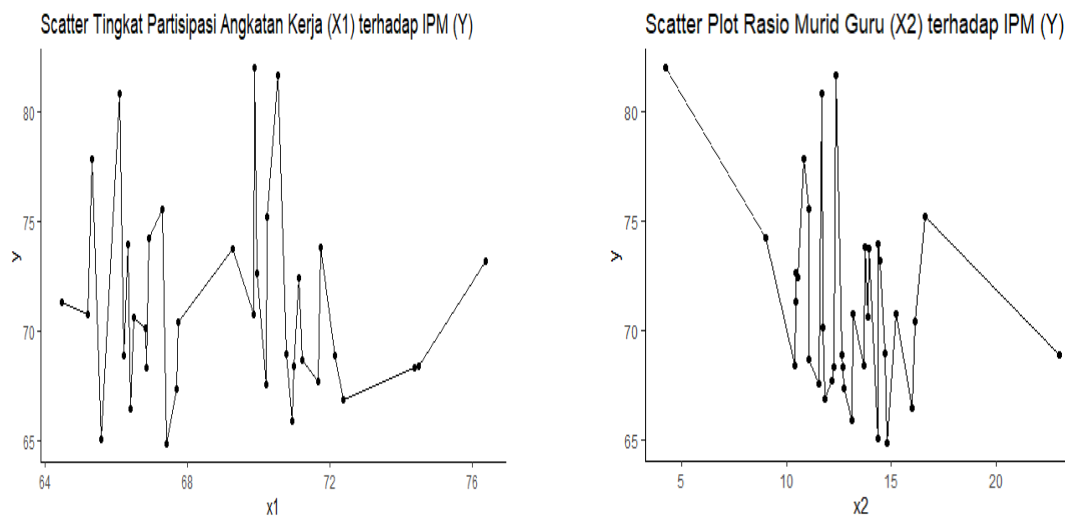
Data yang digunakan bersumber dari data sekunder yang diperoleh dari Badan Pusat Statistika Provinsi Jawa Tengah. Unit observasi sebanyak 35 Kabupaten/Kota. Adapun variable respon yang digunakan dalam penelitian ini adalah Indeks Pembangunan Manusia. Variabel prediktor antara lain Tingkat Partisipasi Angkatan Kerja (TPAK), Rasio Murid Guru, Kepadatan Penduduk, Angka Kesakitan, dan PDRB. Dalam penelitian ini akan digunakan kombinasi lima titik knot dengan orde linear.

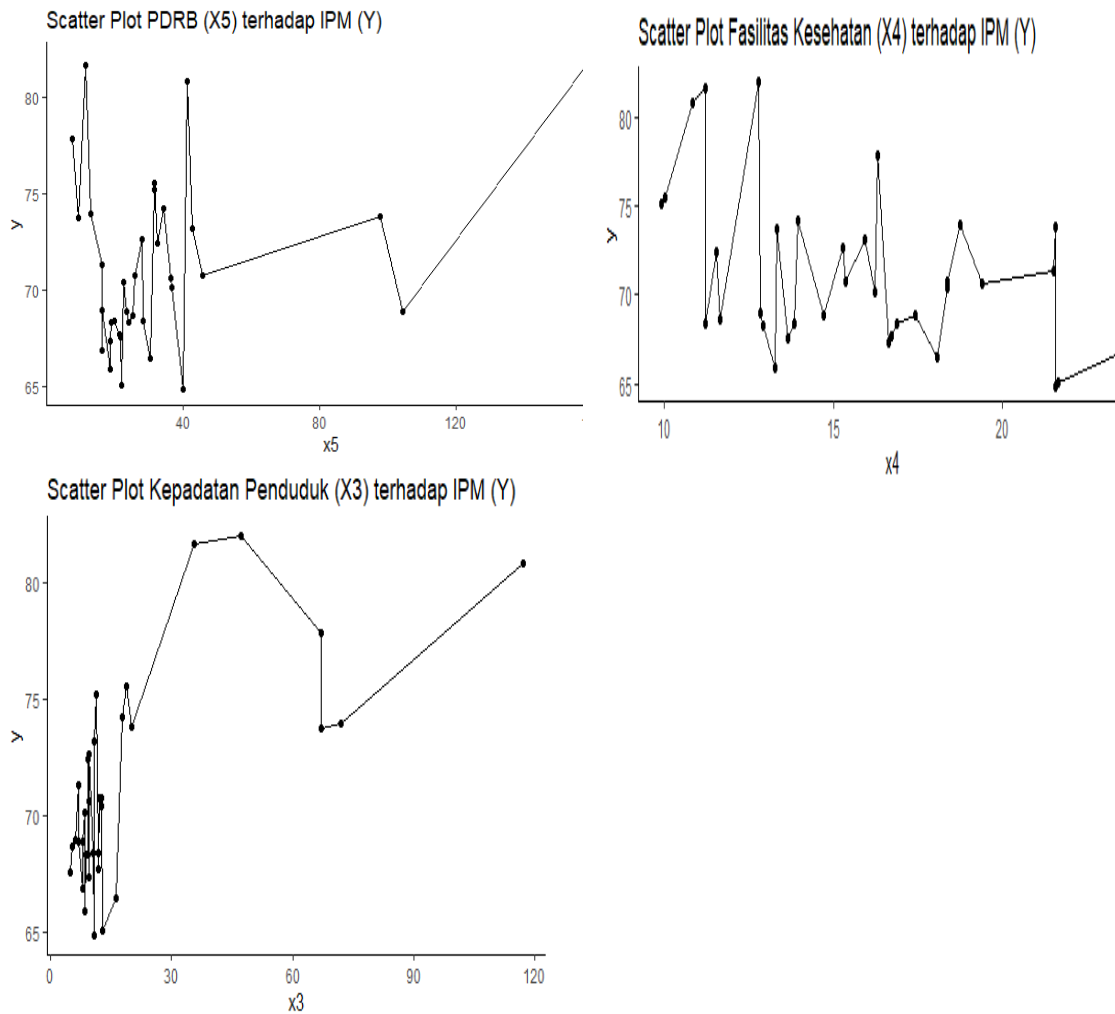
Adapun tahapan penelitian yaitu membuat *scatter plot* variable prediktor (X) terhadap variable respon (Y); memodelkan data dengan model regresi nonparametrik spline dengan satu, dua, tiga, empat, lima dan kombinasi titik knot; menentukan titik knot optimal berdasarkan nilai GCV minimum; memilih estimator regresi spline terbaik; melakukan uji signifikansi parameter; melakukan uji asumsi residual; membuat kesimpulan.

## HASIL DAN PEMBAHASAN

### Scatter Plot

Gambar 1 merupakan *scatter plot* yang menunjukkan pola hubungan variable prediktor terhadap variable respon. Dari Gambar 1 data mengalami pola naik turun yang tajam pada interval tertentu. Dapat dilihat bahwa pola hubungan antara variable IPM terhadap masing-masing variable predkctor yaitu TPAK, Rasio Murid Guru, Kepadatan Penduduk, Angka Kesakitan dan PDRB/1 jt tidak mengikuti pola hubungan tertentu. Selanjutnya dilakukan uji asumsi parametrik. Dengan bantuan R Studio didapatkan bawa model tidak mengikuti asumsi distribusi normal sehingga digunakan regresi nonparametrik. Langkah pertama yang dilakukan dalam membentuk model regresi spline truncated yaitu menentukan titik knot optimal dengan satu, dua, tiga, empat, dan lima titik knot yang digunakan. Tabel 1 merupakan hasil GCV minimum.





Gambar 1. Scatter plot  $x_1$ ,  $x_2$ ,  $x_3$ ,  $x_4$  dan  $x_5$

Model	GCV Minimum
Satu Titik Knot	8,43
Dua Titik Knot	7,80
Tiga Titik Knot	6,86
Empat Titik Knot	4,95
Lima Titik Knot	0,38
Kombinasi Titik Knot (5-5-5-5-5)	0,38

Berdasarkan Tabel 1 diperoleh nilai GCV minimum sebesar 0,38 dengan menggunakan kombinasi lima titik knot (5-5-5-5-5) dengan lokasi titik knot pada variable  $x_1$  terdapat pada titik 65,3; 65,71; 69,81; 70,22 dan 75,55. Untuk variable  $x_2$  pada titik 5,52; titik 6,17; titik 12,63; titik 13,28 dan titik 15,87, Untuk variable  $x_3$  pada titik 12,54; 16,42; 55,19; 59,07 dan 109,47. Untuk variable  $x_4$  pada titik 10,87; 11,34; 16,07; 16,54 dan 22,68. Sementara untuk variable  $x_5$  pada 18,09; 23,33; 75,75; 80,99 dan 149,14. Dengan demikian diperoleh model regresi spline *truncated* optimal menggunakan kombinasi lima titik knot, yang selanjutnya dilakukan pengujian model. Uji signifikansi parameter dilakukan untuk melihat apakah parameter pada variable prediktor berpengaruh signifikan terhadap variable respon atau tidak. Tabel 2 merupakan hasil uji serentak.

Tabel 2. Uji serentak signifikansi parameter

Sumber	Df	SS	MS	$F_{hitung}$	$p-value$
Regresi	30	682,119	22,74	231,77	0,0000394
Error	4	0,392	0,0981		
Total	34	682,511			

Dari Tabel 2 diperoleh bahwa  $p-value = 0,0000394 < 0,05$ . Jadi dapat diputuskan bahwa  $H_0$  ditolak, artinya minimal terdapat satu parameter pada model regresi yang berpengaruh signifikan. Selanjutnya dilakukan uji signifikansi parameter secara parsial yang disajikan dalam Tabel 3.

Berdasarkan Tabel 3 diperoleh bahwa semua parameter yang dihasilkan berpengaruh signifikan terhadap model regresi. Artinya masing-masing variable dapat dikatakan berpengaruh secara signifikan terhadap IPM. Langkah berikutnya yaitu uji asumsi residual yang dilakukan untuk mengetahui apakah residual memenuhi asumsi identik, independen, dan berdistribusi normal (IIDN) atau tidak. Hasil uji asumsi identik disajikan dalam Tabel 4

Tabel 4. Hasil uji Glejser

Sumber	Df	SS	MS	$F_{hitung}$	$p-value$
Regresi	30	0,18	0,0062	0,8609	0,653
Error	4	0,03	0,0072		
Total	34	0,22			

Berdasarkan Tabel 4 diperoleh nilai  $p-value$  sebesar 0,653 dimana  $p-value > 0,05$ . Jadi  $H_0$  gagal ditolak. Jadi dapat disimpulkan tidak terjadi heteroskedastisitas. Artinya asumsi residual identik terpenuhi.

Uji Durbin Watson menghasilkan nilai  $d_{hitung}$  sebesar 2,429467. Langkah selanjutnya yaitu menghitung nilai  $d_L$ . Dari tabel Durbin Watson dengan  $n = 35$  dan  $k = 6$  didapatkan  $d_L = 1,0974$ . Diperoleh bahwa  $0 < d_{hitung} = 2,429467 > d_L = 1,0974$  atau  $4 - d_L = 2,9026 < d_{hitung} = 2,429467 > 4$ . Sehingga  $H_0$  diterima. Artinya tidak terjadi autokorelasi. Jadi residual telah memenuhi asumsi independen. Dengan melakukan Uji Kolmogorov-Smirnov diperoleh nilai  $p-value = 0,17539 > 0,05$ , maka  $H_0$  diterima. Artinya residual telah mengikuti pola distribusi normal.

Tabel 3. Uji Parsial Signifikansi Parameter

Var	Parameter	$p-value$
Konstanta	$\beta_0$	0,000013
$x_1$	$\beta_1$	0,00
	$\beta_2$	0,00
	$\beta_3$	0,0000003
	$\beta_4$	0,01288
	$\beta_5$	0,000043
	$\beta_6$	0,00
$x_2$	$\beta_7$	0,000031
	$\beta_8$	0,000006
	$\beta_9$	0,000012
	$\beta_{10}$	0,00
	$\beta_{11}$	0,00
$x_3$	$\beta_{12}$	0,00
	$\beta_{13}$	0,00
	$\beta_{14}$	0,00
	$\beta_{15}$	0,00
	$\beta_{16}$	0,0001
	$\beta_{17}$	0,00
$x_4$	$\beta_{18}$	0,00
	$\beta_{19}$	0,00
	$\beta_{20}$	0,00
	$\beta_{21}$	0,00

Var	Parameter	<i>p</i> -value
$x_5$	$\beta_{22}$	0,00
	$\beta_{23}$	0,00
	$\beta_{24}$	0,00
	$\beta_{25}$	0,00
	$\beta_{26}$	0,00
	$\beta_{27}$	0,00
	$\beta_{28}$	0,00
	$\beta_{29}$	0,00
	$\beta_{30}$	0,00

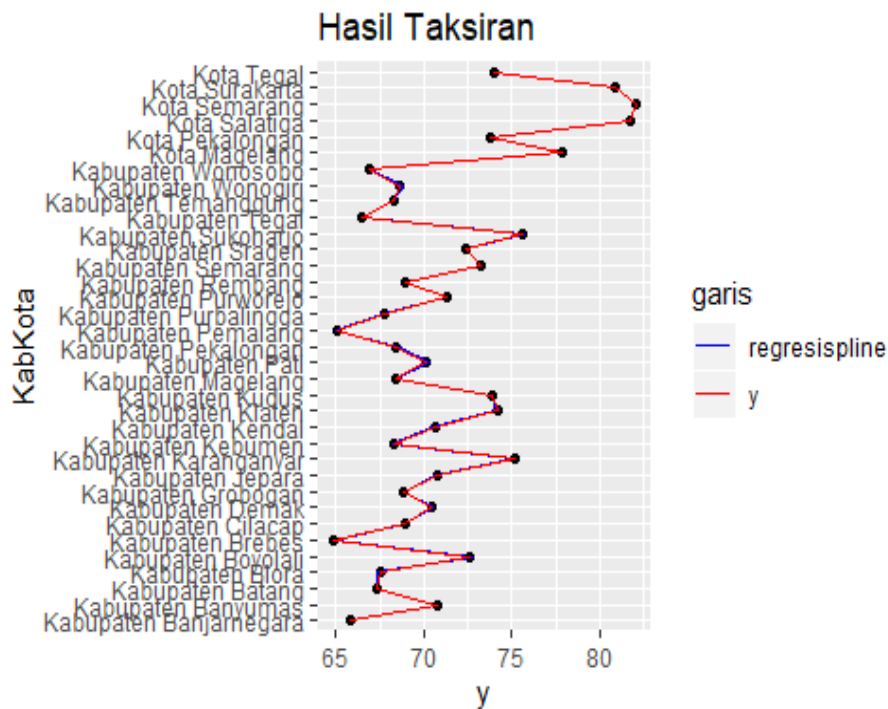
Model terbaik regresi nonparametrik spline *truncated* dihasilkan melalui pemilihan titik knot optimal (5-5-5-5) pada data Indeks Pembangunan Manusia Jawa Tengah dapat dituliskan sebagai berikut. Tanda (+) pada model menunjukkan bahwa terdapat fungsi sepotong (*truncated*) sebagai berikut.

$$\hat{y} = 11,88 + 8,97x_1 + 17,65(x_1 - 65,3)_+^1 - 9,43(x_1 - 65,71)_+^1 + 4,33(x_1 - 69,81)_+^1 - 5,79(x_1 - 70,22)_+^1 + 19,11(x_1 - 75,55)_+^1 + 42,95x_2 - 20(x_2 - 5,52)_+^1 - 26,36(x_2 - 6,17)_+^1 + 13(x_2 - 12,64)_+^1 - 10,49(x_2 - 13,28)_+^1 + 7,88(x_2 - 21,69)_+^1 + 0,45x_3 - 2,14(x_3 - 12,54)_+^1 + 1,52(x_3 - 16,42)_+^1 - 2,65(x_3 - 55,19)_+^1 + 4,33(x_3 - 59,07)_+^1 - 12,78(x_3 - 109,47)_+^1 + 44,42x_4 - 141,07(x_4 - 10,87)_+^1 + 96,31(x_4 - 11,34)_+^1 + 8,71(x_4 - 16,07)_+^1 - 10,41(x_4 - 16,54)_+^1 + 15,39(x_4 - 22,68)_+^1 - 1,97x_5 + 2,77(x_5 - 18,09)_+^1 - 0,85(x_5 - 23,33)_+^1 + 12,38(x_5 - 75,75)_+^1 - 15,03(x_5 - 80,99)_+^1 + 21,34(x_5 - 149,14)_+^1.$$

Fungsi *truncated* ditampilkan sebagai berikut

$$\hat{y}(x) = \begin{cases} 11,88+8,97x_1 & x_1 < 65,3 \\ -1152,74+8,68x_1 & 65,3 \leq x_1 < 65,71 \\ -532,85+0,75x_1 & 65,71 \leq x_1 < 69,81 \\ -835,42+3,58x_1 & 69,81 \leq x_1 < 70,22 \\ -428,53+2,21x_1 & 70,22 \leq x_1 < 75,55 \\ -1872,49+16,9x_1 & 75,55 \leq x_1 \\ -42,95x_2 & x_2 < 5,52 \\ 110,48-22,95x_2 & 5,52 \leq x_2 < 6,17 \\ 273,14-3,41x_2 & 6,17 \leq x_2 < 12,64 \\ 108,84-9,59x_2 & 12,64 \leq x_2 < 13,28 \\ 248,11-0,9x_2 & 13,28 \leq x_2 < 21,69 \\ 77,31-6,98x_2 & 21,69 \leq x_2 \\ -44,42x_3 & x_3 < 12,54 \\ 26,9-1,7x_3 & 12,54 \leq x_3 < 16,42 \\ 1,94-0,18x_3 & 16,42 \leq x_3 < 55,19 \\ 147,93-2,82x_3 & 55,19 \leq x_3 < 59,06 \\ -107,68-1,5x_3 & 59,07 \leq x_3 < 109,47 \\ 1291,6-11,28x_3 & 109,47 \leq x_3 \\ -44,42x_4 & x_4 < 10,87 \\ 1532,82-96,65x_4 & 10,87 \leq x_4 < 11,34 \\ 440,89-0,35x_4 & 11,34 \leq x_4 < 16,07 \\ 300,88-8,37x_4 & 16,07 \leq x_4 < 16,54 \\ 473,08-2,05x_4 & 16,54 \leq x_4 < 22,68 \\ 123,92-13,35x_4 & 22,68 \leq x_4 \\ 1,97x_5 & x_5 < 18,09 \\ -50,1+0,80x_5 & 18,09 \leq x_5 < 23,33 \\ -30,38+0,04x_5 & 23,33 \leq x_5 < 75,75 \\ -968,23+12,34x_5 & 75,75 \leq x_5 < 80,99 \\ 248,75+2,69x_5 & 80,99 \leq x_5 < 149,14 \\ -2933,9+18,65x_5 & 149,14 \leq x_5 \end{cases}$$

Kurva hasil estimasi dari regresi spline *truncated* menggunakan kombinasi titik knot (5-5-5-5) pada data Indeks Pembangunan Manusia Jawa Tengah dapat dilihat pada Gambar 2.



Gambar 2. Kurva regresi spline truncated

Berdasarkan Gambar 2 dapat dilihat bahwa secara visual, kurva regresi spline *truncated* mendekati nilai yang sebenarnya.

## SIMPULAN

Berdasarkan nilai GCV minimum didapatkan model regresi nonparametrik spline *truncated* terbaik pada data IPM Jawa Tengah beserta lima variable yang mempengaruhi yaitu menggunakan kombinasi lima titik knot optimal (5-5-5-5-5). Seluruh variable berpengaruh signifikan terhadap model

$$\hat{y} = 11,88 - 8,97x_1 + 17,65(x_1 - 65,3)_+^1 - 9,43(x_1 - 65,71)_+^1 + 4,33(x_1 - 69,81)_+^1 - 5,79(x_1 - 70,22)_+^1 + 19,11(x_1 - 75,55)_+^1 + 42,95x_2 - 20(x_2 - 5,52)_+^1 - 26,36(x_2 - 6,17)_+^1 + 13(x_2 - 12,64)_+^1 - 10,49(x_2 - 13,28)_+^1 + 7,88(x_2 - 21,69)_+^1 + 0,45x_3 - 2,14(x_3 - 12,54)_+^1 + 1,52(x_3 - 16,42)_+^1 - 2,65(x_3 - 55,19)_+^1 + 4,33(x_3 - 59,07)_+^1 - 12,78(x_3 - 109,47)_+^1 + 44,42x_4 - 141,07(x_4 - 10,87)_+^1 + 96,31(x_4 - 11,34)_+^1 + 8,71(x_4 - 16,07)_+^1 - 10,41(x_4 - 16,54)_+^1 + 15,39(x_4 - 22,68)_+^1 - 1,97x_5 + 2,77(x_5 - 18,09)_+^1 - 0,85(x_5 - 23,33)_+^1 + 12,38(x_5 - 75,75)_+^1 - 15,03(x_5 - 80,99)_+^1 + 21,34(x_5 - 149,14)_+^1$$

Pada penelitian ini hanya terbatas pada orde linier, untuk penelitian selanjutnya dapat menggunakan orde lain seperti orde kuadrat dan kubik pada model regresi nonparametrik spline dengan kombinasi titik knot sehingga memungkinkan mendapatkan hasil yang lebih baik.

## DAFTAR PUSTAKA

- Budiantara, I. N. 2006. Model spline dengan knots optimal. *Jurnal Ilmu Dasar*, 7(6), 77-85.
- Budiantara, I.N. 2009. Spline dalam regresi nonparametrik dan semiparametrik sebuah pemodelan statistika masa kini dan masa mendatang, Pidato Pengukuhan untuk Jabatan Guru Besar dalam Bidang Ilmu Matematika, ITS Press, Institut Teknologi Sepuluh Nopember Surabaya.
- Demu, K.R., Saputro, D. R. S., & Widyaningsih, P. 2017. Model regresi nonoparametrik sline truncated pada data indeks pembangunan manusia (IPM) di Indonesia. Universitas Sebelas Maret.
- Eubank, R.L. 1999. *Nonparametric Regression and Spline Smoothing*. Second Edition. New York: Marcel Dekker, Inc.
- Eubank, R.L. 1988. *Spline Smoothing and Nonparametric Regression*. Marcel Dekker.Inc. New York.

- Fox, J. 2000. Nonparametric simple regression; smoothing scatterplots. International Educational and Professional Publisher. New Delhi, London.
- Hardle, W. 1990. *Applied Nonparametric Regression*. New York: Cambridge University Press.
- Montoya, E.L, Ulloa, N, & Miller, V. 2014. A simulation study comparing knot selection methods with equally spaced knots in a penalized regression spline. *International Journal of Statistics and Probability*, 3(3), 96-110. <http://doi.org/10.5539/ijsp.v3n3p96>
- Putri, TDP., & Dwidayati, NK. 2018. Pemodelan regresi nonparametrik spline multivariable. *Unnes Journal of Mathematics*, 7(2).
- Scholica, C. N., Budiantara, I.N., & Ratna, M. 2018. Regresi nonparametrik spline truncated untuk memodelkan presentase unmet need di Kabupaten Gresik. *Jurnal Sains dan Seni ITS*, 7(2), D61-D68. <http://doi.org/10.12962/j23373520.v7i2.35259>
- Wahba, G. 1990. Spline models for observational data. SIAM, CBMS-NSF Regional Conference Series and Applied Mathematics.