

Tinjauan Singkat Tentang Regresi Parametrik dan Regresi non Parametrik

Muhammad Abdy

Universitas Negeri Makassar
e-mail: muh.abdy@unm.ac.id

Abstrak

Dalam analisis regresi, terdapat dua pendekatan yang digunakan untuk mengestimasi fungsi regresi, yaitu pendekatan parametrik dan pendekatan nonparametrik. Metode regresi parametrik akan sesuai jika ada teori, pengalaman masa lalu atau sumber lain yang dapat digunakan untuk menentukan bentuk fungsi regresi sehingga bentuk dari fungsi regresi diasumsikan diketahui kecuali untuk berhingga parameter yang tidak diketahui. Inferensi tentang fungsi regresi ekuivalen dengan inferensi tentang parameter. Dalam regresi nonparametrik, tidak ada asumsi tentang bentuk fungsi regresi sehingga memberikan fleksibilitas di dalam bentuk yang mungkin dari fungsi regresi. Terdapat banyak tehnik untuk mengestimasi fungsi regresi nonparametrik, antara lain estimator kernel, spline, wavelet, k-NN dan lain-lain.

Kata kunci: regresi parametrik, regresi nonparametrik, fungsi regresi, tehnik penghalusan.

1. PENDAHULUAN

Analisis Regresi merupakan salah satu alat statistik yang banyak digunakan untuk mengetahui hubungan antara sepasang variabel atau lebih. Misalnya y adalah variabel respon dan x adalah variabel predictor, maka untuk n pengamatan, secara umum hubungan variabel itu dapat dinyatakan sebagai :

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

dengan ε adalah variabel acak yang diasumsikan independen, mean nol, dan varians σ^2 . Fungsi $r(\cdot)$ merupakan fungsi yang tidak diketahui yang disebut fungsi regresi.

Ada dua pendekatan yang dapat digunakan untuk mengestimasi fungsi $r(\cdot)$, yaitu pendekatan parametrik dan pendekatan nonparametrik. Fungsi $r(\cdot)$ yang diperoleh dengan pendekatan parametrik merupakan regresi parametrik dan dengan pendekatan nonparametrik merupakan regresi nonparametrik. Sedangkan dengan pendekatan keduanya dinamakan regresi semiparametrik.

Pada regresi parametrik, asumsi yang paling mendasar adalah bahwa bentuk fungsi regresi $r(\cdot)$ diketahui kecuali untuk sejumlah berhingga parameter yang tidak diketahui, yakni ada vektor parameter $\beta = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$ dan ada fungsi $r(\cdot; \beta)$ sehingga $r(\cdot) = r(\cdot; \beta)$. Fungsi regresi $r(\cdot; \beta)$ dapat berbentuk linear atau non linear dalam parameter. Inferensi tentang fungsi regresi $r(\cdot; \beta)$ ekuivalen dengan inferensi tentang parameter β .

Pada umumnya bentuk fungsi regresi $r(\cdot)$ tidak diketahui bentuknya. Tetapi mungkin kita menggunakan sifat kualitatif saja, seperti kontinu, differensiabel atau mulus, dalam arti fungsi tersebut termuat dalam suatu ruang fungsi, yaitu ruang Sobolev. Apabila terjadi hal yang demikian, maka kita menggunakan regresi nonparametrik. Dalam regresi nonparametrik, tidak ada asumsi tentang bentuk fungsi regresi $r(\cdot)$, sehingga memberikan fleksibilitas dalam bentuk yang mungkin dari fungsi regresi. Ada beberapa teknik untuk menduga fungsi regresi $r(\cdot)$ dalam regresi nonparametrik, antara lain estimator

spline, estimator wavelet, k-NN, estimator histogram, estimator kernel, estimator deret orthogonal, estimator deret Fourier, dan lain-lain.

2. REGRESI PARAMETRIK

Pandang kembali model (1). Dalam tulisan ini kita hanya meninjau fungsi regresi $r(\cdot; \beta)$ yang berbentuk linier, yaitu apabila terdapat fungsi yang diketahui f_1, f_2, \dots, f_p sedemikian hingga

$$r(x) = \sum_{j=1}^p \beta_j f_j(x) \quad (2)$$

Berdasarkan bentuk (2) diatas maka model (1) berbentuk:

$$y_i = \sum_{j=1}^p \beta_j f_j(x_i) + \varepsilon_i, i = 1, 2, \dots, n. \quad (3)$$

Variabel x dalam model (3) dapat berupa variabel acak maupun non-acak. Kita akan meninjau hanya untuk variabel x non-acak. Dalam notasi matriks, model (3) ditulis sebagai:

$$\mathbf{y} = \mathbf{x}\beta + \varepsilon \quad (4)$$

dimana $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_n]^T$ adalah vektor respon ukuran $n \times 1$,

$\beta = [\beta_1 \ \beta_2 \ \dots \ \beta_p]$ adalah vektor parameter ukuran $p \times 1$,

$\varepsilon = [\varepsilon_1 \ \varepsilon_2 \ \dots \ \varepsilon_p]$ adalah vektor error yang diasumsikan berdistribusi normal dengan mean 0 dan varians $\sigma^2 I_n$.

dan $\mathbf{x} = \begin{bmatrix} f_1(x_1) & f_2(x_1) & \dots & f_p(x_1) \\ f_1(x_2) & f_2(x_2) & \dots & f_p(x_2) \\ \dots & \dots & \dots & \dots \\ f_1(x_n) & f_2(x_n) & \dots & f_p(x_n) \end{bmatrix}$ adalah matriks data berukuran $n \times p$:

Matriks \mathbf{x} tersebut dapat mempunyai rank penuh, yaitu jika $p \leq n$, atau tidak mempunyai rank penuh, yaitu rank $(\mathbf{x}) < p$. dalam tulisan ini, hanya akan ditinjau matriks dengan rank penuh.

Untuk mengestimasi vektor parameter β , terdapat banyak metode pengestimasi yang tersedia, antara lain metode kuadrat terkecil biasa, metode kuadrat kecil terboboti, metode kemungkinan maksimum dan lain-lain. Untuk asumsi vektor error diatas, metode kuadrat terkecil biasa paling sering digunakan. Prinsip metode ini adalah meminimumkan jumlah kuadrat error $(\varepsilon\varepsilon^T)$ terhadap β . Perhatikan kembali bentuk (4), akan diperoleh :

$$\begin{aligned} \varepsilon^T \varepsilon &= (\mathbf{y} - \mathbf{x}\beta)^T (\mathbf{y} - \mathbf{x}\beta) \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{x}\beta + \beta^T \mathbf{x}^T \mathbf{x}\beta, \\ \frac{\partial}{\partial \beta} \varepsilon^T \varepsilon &= -2\mathbf{x}^T \mathbf{y} + 2\mathbf{x}^T \mathbf{x}\beta \\ \text{jika } \frac{\partial}{\partial \beta} \varepsilon^T \varepsilon &= 0, \text{ maka } \mathbf{x}^T \mathbf{y} = \mathbf{x}^T \mathbf{x}\beta \end{aligned} \quad (5)$$

Persamaan terakhir diatas disebut persamaan normal. Dengan asumsi matrix \mathbf{x} mempunyai rank penuh maka persamaan normal tersebut mempunyai penyelesaian tunggal $\hat{\beta} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$. Kemudian karena $\frac{\partial^2}{\partial \beta \partial \beta^T} (\varepsilon^T \varepsilon) = 2\mathbf{x}^T \mathbf{x}$ merupakan matriks definit positif maka $\hat{\beta}$ akan meminimumkan jumlah kuadrat error. Berikut ini teorema menyangkut sifat estimator kuadrat terkecil $\hat{\beta}$ dan estimasi σ^2 .

Teorema 1.

Estimator kuadrat terkecil $\hat{\beta}$ merupakan estimator UMVUE (Uniform Minimum Variance Unbiased Estimator) untuk β .

Bukti:

1. Sifat tak bias:

$$\begin{aligned} E(\hat{\beta}) &= E[(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}] \\ &= (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T E(\mathbf{y}) \\ &= (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{x}\beta = \beta \end{aligned}$$

2. Sifat varians minimum:

Pertama, akan dicari $\text{var}(\hat{\beta})$.

$$\begin{aligned}\text{var}(\hat{\beta}) &= \text{var}((\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}) \\ &= (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \text{var}(\mathbf{y}) \mathbf{x} (\mathbf{x}^T \mathbf{x})^{-1} \\ &= (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \sigma^2 \mathbf{I}_n \mathbf{x} (\mathbf{x}^T \mathbf{x})^{-1} \\ &= \sigma^2 (\mathbf{x}^T \mathbf{x})^{-1}\end{aligned}$$

selanjutnya, ambil $\beta^* = [(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T + \mathbf{c}] \mathbf{y}$ dengan \mathbf{c} adalah fungsi dari \mathbf{x} , yang merupakan sebarang estimator linear untuk β .

$$\begin{aligned}E(\beta^*) &= E[(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T + \mathbf{c}] \mathbf{y} \\ &= [(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T + \mathbf{c}] E(\mathbf{y}) \\ &= [(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T + \mathbf{c}] \mathbf{x} \beta\end{aligned}$$

Agar β^* tak bias untuk β , maka haruslah $\mathbf{c} \mathbf{x} = 0$.

$$\begin{aligned}\text{Selanjutnya: } \text{var}(\beta^*) &= [(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T + \mathbf{c}] \text{var}(\mathbf{y}) [\mathbf{c}^T + \mathbf{x} (\mathbf{x}^T \mathbf{x})^{-1}] \\ &= \sigma^2 [(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{c} + \mathbf{c} \mathbf{c}^T + (\mathbf{x} \mathbf{x}^T)^{-1} + \mathbf{c} \mathbf{x} (\mathbf{x}^T \mathbf{x})^{-1}] \\ &= \sigma^2 (\mathbf{x}^T \mathbf{x} + \mathbf{c} \mathbf{c}^T)\end{aligned}$$

sehingga $\text{var}(\beta^*) - \text{var}(\hat{\beta}) = \sigma^2 \mathbf{c} \mathbf{c}^T$ merupakan matriks definit non-negatif. Akibatnya $\text{var}(\beta^*) \geq \text{var}(\hat{\beta})$.

Teorema 2

Jika $E(\mathbf{y}) = \mathbf{x} \beta$ dengan \mathbf{x} matriks berukuran $n \times p$ dengan rank $p \leq n$, $\text{var}(\mathbf{y}) = \sigma^2 \mathbf{I}_n$, maka $s^2 = \frac{(\mathbf{y} - \mathbf{x} \hat{\beta})^T (\mathbf{y} - \mathbf{x} \hat{\beta})}{n-p}$ merupakan estimator tak bias untuk σ^2 .

Sebelum membuktikan teorema tersebut terlebih dahulu akan dibuktikan lemma berikut:

Lemma:

Misalkan $\mathbf{z} = \{z_i\}$ merupakan vektor variabel acak berukuran $n \times 1$ dan \mathbf{A} matriks simetris berukuran $n \times n$. Jika $E(\mathbf{z}) = \theta$ dan $\text{var}(\mathbf{z}) = \{\sigma^2\} = \Sigma$ maka $E(\mathbf{z}^T \mathbf{A} \mathbf{z}) = \text{tr}(\mathbf{A} \Sigma) + \theta^T \mathbf{A} \theta$.

Bukti:

$$\begin{aligned}E(\mathbf{z}^T \mathbf{A} \mathbf{z}) &= E[(\mathbf{z} - \theta)^T \mathbf{A} (\mathbf{z} - \theta) + \theta^T \mathbf{A} \mathbf{z} + \mathbf{z}^T \mathbf{A} \theta - \theta^T \mathbf{A} \theta] \\ &= E[(\mathbf{z} - \theta)^T \mathbf{A} (\mathbf{z} - \theta)] + \theta^T \mathbf{A} \theta \\ &= \sum_i \sum_j a_{ij} E(z_i - \theta_i)(z_j - \theta_j) + \theta^T \mathbf{A} \theta \\ &= \sum_i \sum_j a_{ij} \sigma_{ij}^2 + \theta^T \mathbf{A} \theta \\ &= \text{tr}(\mathbf{A} \Sigma) + \theta^T \mathbf{A} \theta\end{aligned}$$

selanjutnya teorema 2 dapat dibuktikan sebagai berikut:

$$\begin{aligned}(n-p)s^2 &= (\mathbf{y} - \mathbf{x} \hat{\beta})^T (\mathbf{y} - \mathbf{x} \hat{\beta}) \\ &= [(\mathbf{I} - \mathbf{P}) \mathbf{y}]^T (\mathbf{I} - \mathbf{P}) \mathbf{y} \\ &= \mathbf{y}^T (\mathbf{I} - \mathbf{P}) \mathbf{y}\end{aligned}$$

$$\begin{aligned}E[\mathbf{y}^T (\mathbf{I} - \mathbf{P}) \mathbf{y}] &= \text{tr}(\mathbf{I} - \mathbf{P}) \sigma^2 \mathbf{I}_n + \beta^T \mathbf{x}^T (\mathbf{I} - \mathbf{P}) \mathbf{x} \beta \\ &= \sigma^2 \text{tr}(\mathbf{I} - \mathbf{P}) + \beta^T \mathbf{x}^T \mathbf{x} \beta - \beta^T \mathbf{x}^T \mathbf{P} \mathbf{x} \beta \\ &= \sigma^2 (n-p) + \beta^T \mathbf{x}^T \mathbf{x} \beta - \beta^T \mathbf{x}^T \mathbf{x} \beta \\ &= \sigma^2 (n-p)\end{aligned}$$

$$\text{jadi } E[(n-p)s^2] = (n-p)E(s^2) = \sigma^2 (n-p)$$

$$\text{Akibatnya } E(s^2) = \sigma^2$$

3. REGRESI NONPARAMETRIK

Metode regresi nonparametrik mulai dikenal sejak abad XIX. Engel, seorang ekonom telah mengkonstruksi suatu kurva yang dikenal sebagai regressogram. Perkembangan regresi nonparametrik tidak terlalu pesat, yang lebih cepat berkembang adalah regresi dengan pendekatan parametrik. Namun sejak beberapa dekade terakhir, regresi nonparametrik begitu pesat berkembang seiring dengan perkembangan teknologi komputer.

Pandang kembali bentuk (1). Bentuk tersebut dengan tidak ada informasi sebelumnya mengenai bentuk $r(x)$ akan menghasilkan regresi nonparametrik. Jika diamati nilai-nilai variabel respon y dan variabel-variabel prediktor x ditentukan, diperoleh hasil pengamatan bivariate $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ yang mengikuti model (1) dengan $r(\cdot)$ adalah kurva regresi yang tidak

diketahui dan akan diestimasi serta ε_i adalah error acak independent dengan mean nol dan varians sama σ^2 . Tidak ada asumsi tentang bentuk fungsi regresi $r(\cdot)$. fungsi regresi $r(\cdot)$ umumnya hanya diasumsikan termuat dalam suatu ruang fungsi yang berdimensi tak hingga. Untuk mengkonstruksi model regresi nonparametrik, terlebih dahulu dipilih ruang fungsi yang sesuai, dimana fungsi regresi $r(\cdot)$ diyakini termasuk didalamnya. Pemilihan ruang fungsi ini biasanya dimotivasi oleh kemulusan (smoothness) yang diasumsikan dimiliki oleh fungsi regresi, sehingga fungsi regresi tersebut diasumsikan termuat dalam suatu ruang fungsi, yang dinamakan ruang Sobolev : $W_2^m[a, b] = \{f|f, f', \dots, f^{(m-1)} \text{ kontinu absolut pada } [a, b], f^{(m)} \in L_2[a, b]\}$ dimana $L_2[a, b]$ adalah himpunan fungsi yang kuadratnya terintegral pada $[a, b]$ atau dapat dinyatakan dengan $L_2[a, b] = \{f| \int_a^b [f^{(m)}(t)]^2 dt < \infty\}$. Data pengamatan kemudian digunakan untuk mengestimasi fungsi $r(\cdot)$ dengan teknik penghalusan tertentu.

4. TEKNIK PENGHALUSAN PADA REGRESI NONPARAMETRIK

Jika fungsi regresi $r(\cdot)$ diyakini licin maka pengamatan pada titik-titik dekat x akan memuat informasi tentang nilai $r(\cdot)$ pada x . Dengan demikian adalah mungkin untuk menggunakan rata-rata lokal dari pengamatan dekat x untuk mengkonstruksi estimator dari $r(\cdot)$. prosedur rata-rata lokal ini dapat dipandang sebagai dasar pemikiran dari tehnik-tehnik penghalusan. Secara formal prosedur ini didefinisikan sebagai:

$$\hat{r}(x) = \frac{1}{n} \sum_{i=1}^n W_{ni}(x) y_i \quad (6)$$

dengan $\{W_{ni}(x); i = 1, 2, \dots, n\}$ adalah barisan pembobot dan $\hat{r}(x)$ merupakan estimator dari $r(x)$. Dari (6), besarnya rata-rata dikontrol oleh $\{W_{ni}(x); i = 1, 2, \dots, n\}$ yang diukur oleh parameter penghalus. Terdapat banyak Teknik penghalusan untuk mengestimasi fungsi regresi $r(\cdot)$, antara lain estimator histogram, estimator kernel, estimator k-NN, estimator spline, wavelet dan lain-lain. Pada tulisan ini hanya kan dibahas sepiantas estimator kernel.

Pada tahun 1956 Murray Rosenblatt mengusulkan suatu estimator untuk fungsi kepadatan $f(x)$ yaitu:

$$\hat{f}_h(x) = \frac{1}{n} \sum_{i=1}^n K_h[x - x_i] = \frac{1}{nh} \sum_{i=1}^n K \left[\frac{x - x_i}{h} \right] \quad (7)$$

Estimator ini dinamakan estimator kepadatan kernel untuk fungsi kepadatan $f(x)$. Fungsi K merupakan fungsi pembobot yang dinamakan fungsi kernel dan h dinamakan parameter penghalus atau sering disebut bandwidth.

Selanjutnya estimator kernel untuk fungsi regresi $r(\cdot)$ dikonstruksi sebagai berikut:

$$\begin{aligned} \text{Perhatikan bahwa } r(x) &= E[Y|X = x] = \int_{-\infty}^{\infty} yf(y|x)dy \\ &= \frac{\int_{-\infty}^{\infty} yf(y,x)dy}{f(x)} \end{aligned} \quad (8)$$

Estimator untuk $r(x)$ adalah dengan mengganti pembilang dan penyebut pada (8) dengan (7) sehingga diperoleh :

$$\hat{r}_h(x) = \frac{\int_{-\infty}^{\infty} y \hat{f}_h h_1, h_2(x, y) dy}{\hat{f}_h(x)} = \frac{n^{-1} \sum_1^n K_h(x - x_i) y_i}{n^{-1} \sum_1^n K_h(x - x_j)}$$

Estimator tersebut diatas diusulkan oleh Nadaraya dan Watson pada tahun 1964 sehingga biasa dinamakan estimator Nadaraya-Watson. Dibawah asumsi tertentu distribusi estimator tersebut konvergen kedistribusi normal sehingga dapat dibentuk interval kepercayaan titik demi titik untuk fungsi regresi $r(\cdot)$. misalkan estimator kernel hanya dihitung pada pengamatan $\{x_i; i = 1, 2, \dots, n\}$, maka jika $h \rightarrow 0$ maka

$$\hat{r}_h(x) \rightarrow \frac{K(0)y_i}{K(0)} = y_i$$

Jadi bandwidth yang kecil akan menghasilkan data y_i sebagai estimator. Sebaliknya, jika $h \rightarrow \infty$ maka

$$K\left(\frac{x-x_i}{h}\right) \rightarrow K(0), \text{ akibatnya } \hat{r}_h(x_i) \rightarrow n^{-1} \sum_1^n y_i$$

Jadi bandwidth yang besar akan menghasilkan estimator yang sangat mulus dan menuju rata-rata variabel respon. Terdapat suatu metode untuk memilih bandwidth yang paling baik (optimal), tapi tidak dibahas dalam tulisan ini.

5. KESIMPULAN

Perbedaan mendasar pada asumsi tentang fungsi regresi $r(\cdot)$ menjadi salah satu pertimbangan untuk memilih metode mana yang akan dipakai. Metode regresi parametrik akan sesuai jika ada teori, pengalaman masa lalu atau sumber lain yang dapat digunakan untuk menentukan bentuk fungsi regresi. Jika hanya sedikit informasi tentang fungsi regresi, maka data lebih banyak mengandung tentang fungsi regresi, sehingga metode regresi nonparametrik lebih sesuai digunakan.

Sebagai pertimbangan lain, jika diperhatikan sifat asimptotik dari estimator, ternyata jika digunakan kriteria risk kuadrat sebagai kebaikan estimator, umumnya estimator nonparametrik akan konvergen dalam probabilitas ke fungsi regresi yang sebenarnya dengan laju $n^{-\delta}$ untuk suatu $0 < \delta < 1$. Sedangkan untuk estimator parametrik kecepatan kekonvergenya dapat mencapai n^{-1} . Dengan demikian estimator parametrik apabila sesuai akan lebih efisien dibanding estimator nonparametrik. Tetapi penggunaannya sangat tergantung pada masalah yang dihadapi.

DAFTAR PUSTAKA

- Budiantara, I. N dan Subanar. (1996), *Regresi Spline dan Permasalahannya*, Naskah Publikasi UGM
- Draper, N and Smith, H. (1981), *Applied Regression Analysis*, John Wiley & Sons, Inc, New York.
- Eubank, R. (1998), *Spline Smoothing and Nonparametrik Regression*, Marcel Dekker, New York.
- Hart, J. D, and Wehrly, T. E. (1986), *Kernel Regression Estimation Using Repeated Measurement Data*, Journal of the American Statistical Association, 81, 1080-1088.
- Suyono (1997), *Perbandingan Regresi Parametrik dan Regresi Nonparametrik*, Tesis S2 UGM.