

Pendeteksian *Outlier* pada Pengamatan dalam Model Linear Multivariat dengan Metode *Likelihood Displacement Statistic-Lagrange*

Outlier Detection in Observation at Multivariate Linear Models with Likelihood Displacement Statistic-Lagrange Method

Makkulau ¹⁾, Susanti Linuwih ²⁾, Purhadi ²⁾ & Muhammad Mashuri ²⁾

¹⁾Jurusan Matematika, FMIPA Universitas Haluoleo, Kendari

²⁾Jurusan Statistika, FMIPA Institut Teknologi Sepuluh Nopember

ABSTRACT

There are two different outliers, i.e outlier in observations and outlier in models. The existing outlier detection method in models is using common Likelihood method. The limitation of this method is the optimal value produced might be not the real optimal values. This research yields a method for outlier detection in multivariate linear models with Likelihood Displacement Statistic-Lagrange method (LDL method). This method uses multiplier Lagrange with constraint the confidence interval of parameter's vector. This parameter's vector is obtained from the data set which is outlier free. This parameter estimation process uses numerical method with Karush-Kuhn Tucker condition in nonlinear programming. This method compares between LDL value and the table F value that follows the distribution of F value to identify the outlier in models.

Keywords : Distribution of F, Likelihood Displacement Statistic-Lagrange, multivariate linear models, nonlinear programming, outlier detection

PENDAHULUAN

Outlier merupakan pengamatan yang menyimpang sedemikian jauh dari pengamatan lain (Hawkins 1980), dapat mempunyai efek bagi pengambilan suatu kesimpulan atau keputusan pada penelitian. *Outlier* dibedakan atas *outlier* pada pengamatan univariat atau multivariat dan *outlier* pada model linear univariat atau multivariat.

Pendeteksian *outlier* pada model linear telah dilakukan antara lain oleh Srivastava & von Rosen (1998), Cook (2000), Adnan *et al.* (2003), dan Diaz-Garcia *et al.* (2007). Xu *et al.* (2005) mengembangkan jarak Cook's univariat untuk mendeteksi *outlier* pada model linear multivariat. Metode yang digunakan adalah Metode *Likelihood Displacement Statistic* (Metode LD), yaitu suatu metode yang menghilangkan pengamatan yang *outlier* pada model; Metode *Likelihood Ratio Statistic for a Mean Shift* (Metode LR), yaitu suatu metode dengan cara pergeseran rata-rata pada model; dan Metode *Multivariate Leverage* yang menggunakan elemen dari *the average diagonal* Q_{A_m} untuk mengukur keekstriman dari m pengukuran pada variabel independen. Makkulau *et al.* (2008) membahas prosedur Metode LD sedangkan Makkulau *et al.* (2009)

membahas aplikasinya di Pabrik Gula Djombang Baru Jombang Provinsi Jawa Timur. Xu *et al.* (2005) dalam mengestimasi parameter dengan Metode LD dari model linear multivariat bersifat umum, sehingga nilai optimal yang diperoleh dari fungsi tujuan dapat saja bukan merupakan nilai yang paling optimal. Oleh karena itu digunakan pengganda *Lagrange*, sehingga nilai optimal yang diperoleh merupakan nilai yang paling optimal pada daerah kepercayaan yang telah ditentukan.

Penelitian ini mengkaji tentang pendeteksian *outlier* pada pengamatan dalam model linear multivariat sebagai pengembangan Metode LD dengan menggunakan pengganda *Lagrange* yang disebut Metode *Likelihood Displacement Statistic-Lagrange* (Metode LDL). Pengganda *Lagrange* yang digunakan adalah daerah kepercayaan dari vektor parameter dimana pengamatan yang diduga *outlier* telah dihilangkan.

Model linear multivariat

Model linear multivariat adalah model linear dengan variabel dependen lebih dari satu (Christensen 1991). Misalkan X_1, X_2, \dots, X_p adalah variabel independen dan Y_1, Y_2, \dots, Y_q

adalah variabel dependen, jika dilakukan n pengamatan yang diambil pada setiap variabel dependen yang ditulis $y_{i1}, y_{i2}, \dots, y_i$ dimana $i = 1, 2, \dots, n$, atau:

$$y_{ih} = \beta_{0h} + \beta_{1h} X_{i1} + \beta_{2h} X_{i2} + \dots + \beta_{ph} X_{ip}$$

dimana $h = 1, 2, \dots, q$.

Model linear multivariat yang terdiri dari q model linear secara simultan dapat ditulis sebagai:

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E} \tag{1}$$

dengan

$$\mathbf{Y}_{n \times q} = \begin{bmatrix} y_{11} & y_{12} & \dots & y_{1q} \\ y_{21} & y_{22} & \dots & y_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n1} & y_{n2} & \dots & y_{nq} \end{bmatrix}$$

$$\mathbf{X}_{n \times (p+1)} = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1p} \\ 1 & X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \dots & X_{np} \end{bmatrix}$$

$$\mathbf{B}_{(p+1) \times q} = \begin{bmatrix} \beta_{01} & \beta_{02} & \dots & \beta_{0q} \\ \beta_{11} & \beta_{12} & \dots & \beta_{1q} \\ \beta_{21} & \beta_{22} & \dots & \beta_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_{p1} & \beta_{p2} & \dots & \beta_{pq} \end{bmatrix}$$

$$\text{dan } \mathbf{E}_{n \times q} = \begin{bmatrix} \varepsilon_{11} & \varepsilon_{12} & \dots & \varepsilon_{1q} \\ \varepsilon_{21} & \varepsilon_{22} & \dots & \varepsilon_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ \varepsilon_{n1} & \varepsilon_{n2} & \dots & \varepsilon_{nq} \end{bmatrix}$$

Estimasi parameter dan uji hipotesis model linear multivariat

Pada model linear multivariat matriks *error* $\mathbf{E}_{n \times q} = (\varepsilon_{ih})$ merupakan matriks acak, dimana $i = 1, 2, \dots, n$ dan $h = 1, 2, \dots, q$. Diasumsikan juga bahwa $E(\mathbf{E}) = \mathbf{0}$ dan matriks varian-kovariansi-nya adalah $\text{Var}(\mathbf{E}) = \mathbf{I} \otimes \mathbf{\Sigma}$, dimana \otimes adalah perkalian Kronecker dan $\mathbf{\Sigma} = (\sigma_{hh})$ dengan $h = 1, 2, \dots, q$, $h^* = 1, 2, \dots, q$, sehingga $\mathbf{E} \sim N_p(\mathbf{0}, \mathbf{I} \otimes \mathbf{\Sigma})$.

Dengan mengestimasi parameter \mathbf{B} pada (1), maka diperoleh:

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \tag{2}$$

Estimasi parameter $\mathbf{\Sigma}$ pada (1), yaitu:

$$\hat{\mathbf{\Sigma}} = \frac{1}{n} (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^T (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})$$

$\hat{\mathbf{\Sigma}}$ adalah estimator bias untuk $\mathbf{\Sigma}$. Sedangkan estimasi parameter $\mathbf{\Sigma}$ yang lain adalah:

$$\mathbf{S} = \frac{1}{n - \text{rank}(\mathbf{X})} (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^T (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}) \mathbf{S}$$

merupakan estimator tak bias untuk $\mathbf{\Sigma}$.

Vektorisasi matriks variabel dependen pada (1) ditulis $\text{Vec}(\mathbf{Y})$ (Christensen 1991) adalah:

$$\text{Vec}(\mathbf{Y}) = (\mathbf{I}_q \otimes \mathbf{I}_n) \text{Vec}(\mathbf{B}) + \text{Vec}(\mathbf{E})$$

Dimana

$$\text{Vec}(\mathbf{E}) \sim N_p(\mathbf{0}, \mathbf{\Sigma} \otimes \mathbf{I}_n) \tag{3}$$

dan $\text{Vec}(\mathbf{Y}) \sim N_{nq}((\mathbf{I}_q \otimes \mathbf{I}_n) \text{Vec}(\mathbf{B}), \mathbf{\Sigma} \otimes \mathbf{I}_n)$. Dengan menggunakan sifat hasil kali *kroncker* diperoleh:

$$\text{Vec}(\hat{\mathbf{B}}) = (\mathbf{I}_q \otimes (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) \text{Vec}(\mathbf{Y})$$

dan distribusi $\text{Vec}(\hat{\mathbf{B}})$ adalah:

$$\text{Vec}(\hat{\mathbf{B}}) \sim N_{q(p+1)}(\text{Vec}(\mathbf{B}), \mathbf{\Sigma} \otimes (\mathbf{X}^T \mathbf{X})^{-1})$$

Prosedur uji hipotesis parameter pada model linear multivariat adalah:

$$H_0 : \Lambda^T \mathbf{B} = \mathbf{0} \text{ terhadap } H_1 : \Lambda^T \mathbf{B} \neq \mathbf{0}$$

dimana $\Lambda^T = \mathbf{P}^T \mathbf{X}$ dan \mathbf{P} adalah matriks ortogonal (Christensen 1991). Uji ini didasarkan pada statistik uji:

$$\mathbf{H} = \mathbf{Y}^T \mathbf{M}_{MP} \mathbf{Y} = (\Lambda^T \mathbf{B})^T (\Lambda^T \mathbf{X}^T \mathbf{X})^{-1} (\Lambda^T \mathbf{B})$$

dimana $\mathbf{M}_{MP} = \mathbf{M} \mathbf{P} (\mathbf{P}^T \mathbf{M} \mathbf{P})^{-1} \mathbf{P}^T \mathbf{M}$.

\mathbf{M}_{MP} adalah matriks proyeksi pada \mathbf{H} .

$$\mathbf{M} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

\mathbf{M} adalah matriks proyeksi pada ruang kolom dari \mathbf{X} .

Hipotesis nol ditolak jika nilai maksimum dari *likelihood* di bawah H_0 lebih besar dari nilai maksimum keseluruhan.

Outlier pada pengamatan dalam model linear multivariat

Outlier pada pengamatan dalam model linear multivariat dapat dibagi atas 3 kategori, yaitu *outlier* terhadap nilai X ; *outlier* terhadap nilai Y ; dan *outlier* terhadap nilai X dan Y .

Outlier dapat diklasifikasikan ke dalam empat kelompok berdasarkan penyebabnya, yaitu observasi umum; titik *leverage* baik; *outlier* vertikal; dan titik *leverage* jelek (Rousseeuw & Hubert 1997).

Metode pendeteksian Outlier pada pengamatan dalam model linear multivariat

Fungsi *likelihood* dalam model linear multivariat ditulis sebagai berikut (Christensen 1991 serta Rencher & Schaalje 2008):

$$L(\mathbf{B}, \mathbf{\Sigma}) = (2\pi)^{-\frac{nq}{2}} |\mathbf{\Sigma}|^{-\frac{n}{2}} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X}\mathbf{B})^T (\mathbf{Y} - \mathbf{X}\mathbf{B}))\right) \tag{4}$$

dimana $\Delta = \left(-\frac{1}{2} \text{tr}(\mathbf{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X}\mathbf{B})^T (\mathbf{Y} - \mathbf{X}\mathbf{B}))\right)$

Pendeteksian *outlier* pada pengamatan dalam model linear multivariat oleh Xu et al. (2005) mengembangkan jarak Cook's univariat

dengan menggunakan Metode LD, Metode LR, dan Metode *Multivariate Leverage*. Pendeteksian *outlier* dengan Metode LD dilakukan dengan cara menghilangkan pengamatan yang *outlier* pada model. Misalkan ada m pengamatan dikumpulkan pada himpunan tertentu, dengan m pengamatan diduga *outlier*. Indeks A_m adalah kumpulan dari m pengamatan yang diduga *outlier*. Dengan kata lain, indeks A_m artinya ada *outlier*, sehingga:

Y_{A_m} : himpunan dari variabel dependen dengan pengamatan yang ada *outlier*.

$Y_{A_m}^c$: himpunan dari variabel dependen dengan pengamatan tanpa *outlier*.

Berikut adalah definisi dari fungsi *Likelihood Displacement Statistic* (LD) untuk pengamatan yang ada *outlier*.

Definisi 1 (Christensen 1991)

LD dengan pengamatan yang ada *outlier* untuk B dengan diberikan Σ adalah:

$$LD_{A_m}(B|\Sigma) = \left(2 \ln(\hat{B}, \hat{\Sigma}) - \ln(\hat{B}_{A_m}, \hat{\Sigma}_{A_m}^c) \right) \quad (5)$$

dimana $\hat{\Sigma}(\hat{B}_{A_m}^c)$ adalah estimasi maksimum *likelihood* dari Σ ketika B diestimasi oleh $\hat{B}_{A_m}^c$.

Definisi 2 (Christensen 1991)

LD dengan pengamatan yang ada *outlier* dan bersyarat adalah:

$$LD_{A_m}(B_1, \Sigma | B_2, \Sigma) = \left(2 \ln(\hat{B}_1, \hat{\Sigma}_1) - \ln(\hat{B}_2, \hat{\Sigma}_2) \right) \quad (6)$$

dimana:

$$\hat{\Sigma}_1 = \left(\hat{B}_{1, A_m}^c, \hat{\Sigma}(\hat{B}_{1, A_m}^c) \right)$$

$$\hat{\Sigma}_2 = \left\{ \left(\hat{B}_{2, A_m}^c, \hat{\Sigma}(\hat{B}_{2, A_m}^c) \right) \right\} \left(\hat{B}_{1, A_m}^c, \hat{\Sigma}(\hat{B}_{1, A_m}^c) \right)$$

{ } menotasikan suatu fungsi yang bentuknya $\hat{\theta}_2(\hat{\theta}_{1, A_m}^c)$

Estimator dari B dengan pengamatan tanpa *outlier* adalah:

$$\hat{B}_{A_m}^c = \hat{B} - \left(X^T X \right)^{-1} K_{A_m}^T \left(Q - \right) E_{A_m}^{-1} \hat{E}$$

dimana

$$Q_{A_m} = X_{A_m}^T X_{A_m}^{-1} X; \hat{E}_{A_m} = Y_{A_m}^T - X_{A_m}$$

$$I + \left(I - Q_{A_m} \right)^{-1} Q_{A_m} = \left(I - Q_{A_m} \right)^{-1}$$

dan estimator dari Σ dengan pengamatan tanpa *outlier* yaitu $\hat{\Sigma}(\hat{B}_{A_m}^c)$ adalah:

$$\hat{\Sigma}(\hat{B}_{A_m}^c) = \frac{n}{n-m} \hat{\Sigma} \frac{1}{n-m} \hat{E}_{A_m}^T I Q_{A_m}^{-1}$$

sehingga fungsi *likelihood* dalam model linear multivariat dengan pengamatan yang ada *outlier* adalah:

$$LD_{A_m}(B|\Sigma) = \left(2 \ln(\hat{B}, \hat{\Sigma}) - \ln(\hat{B}_{A_m}^c, \hat{\Sigma}_{A_m}^c) \right)$$

Optimasi nonlinear

Masalah optimasi ditentukan oleh karakteristik fungsi tujuan dan fungsi kendala. Permasalahan optimasi disebut nonlinear jika fungsi tujuan dan fungsi kendalanya mempunyai bentuk nonlinear pada salah satu atau keduanya (Bazaara et al.1993).

Pada optimasi dengan kendala, persamaan yang akan dioptimasi dapat dituliskan sebagai berikut:

Maksimumkan (minimumkan):

$$g(Z) = g(Z_1, Z_2, \dots, Z_p)$$

Kendala : $k_p(Z) = b_p$ dan $Z \geq 0$.

Optimasi dengan kendala ini dapat diselesaikan dengan menggunakan pengganda *Lagrange* seperti berikut ini:

$$L = g(Z) - \sum_{p=1}^n \lambda_p (k_p(Z) - b_p)$$

METODE

Pembahasan dalam makalah ini dilakukan dengan pendekatan teoretik. Secara garis besar, langkah pertama yang dilakukan dalam penelitian ini adalah mencari solusi persamaan Schrodinger PDM sistem osilator harmonik secara analitik, yaitu berupa nilai eigen energi (En) dan fungsi eigen (Ψn), dengan menggunakan metode transformasi. Langkah kedua adalah melakukan aproksimasi $m(x)=1$ dan $\mu(x)=x$ untuk melihat apakah hasilnya mereduksi menjadi hasil untuk sistem osilator harmonik dengan massa konstan. Langkah kedua ini sekaligus sebagai verifikasi apakah hasil yang didapat konsisten atau tidak. Pada bagian akhir akan ditinjau pula beberapa bentuk fungsi massa yang bergantung posisi dan potensial yang dibangkitkannya, serta interpretasinya secara fisis.

Pendeteksian *outlier* dengan Metode LDL berdasarkan langkah-langkah:

- a. Mengumpulkan m pengamatan yang diduga *outlier*.
- b. Mendeteksi *outlier* pada model dengan asumsi $Ved(E) \sim N_p(\Sigma_0, I)$, dimulai dengan membuat

$L(\mathbf{B}, \Sigma)$, lalu menentukan $\hat{\mathbf{B}}$ dan $\hat{\Sigma}$ untuk mendapatkan $L(\hat{\mathbf{B}}, \hat{\Sigma})$.

- c. Memaksimumkan $L(\mathbf{B}, \Sigma)$ dengan kendala jika ada m buah pengamatan adalah outlier menggunakan pengganda Lagrange, lalu membuat $\ln L(\mathbf{B}, \Sigma)$ dan menentukan $\text{Vec}(\hat{\mathbf{B}})$ dengan program nonlinear.

HASIL DAN PEMBAHASAN

Nilai eigen energi E_n dan fungsi eigen $\Psi(x)$ sistem osilator harmonik PDM

Penelitian ini dibatasi hanya pada pendeteksian outlier pada pengamatan dalam model linear multivariat dengan Metode LDL. Pendeteksian outlier dalam model linear multivariat dimulai dengan memisalkan ada m pengamatan yang diduga outlier (A_m) dari $Y_{o/n}, Y_{o/n}^L, Y_{o/n}^L, Y_{o/n}^L$, sehingga \mathbf{Y}_{A_m} adalah himpunan dari variabel dependen dengan pengamatan yang ada outlier dan $\mathbf{Y}_{A_m}^C$ adalah himpunan dari variabel dependen dengan pengamatan tanpa outlier. Sebelumnya, jika dipunyai variabel independen sebanyak p dan variabel dependen sebanyak q , maka model linear multivariat (1) secara simultan dapat ditulis sebagai:

$$\mathbf{Y}_{n \times q} = \begin{pmatrix} J_{o/n} & \mathbf{M} & \mathbf{X}_{1, n \times (p+1)} \end{pmatrix} \mathbf{B}_{(p+1) \times q} + \mathbf{E}_{n \times q}$$

$$= \mathbf{X}_{n \times (p+1)} \mathbf{B}_{(p+1) \times q} + \mathbf{E}_{n \times q} \tag{7}$$

Dari (7) dapat ditulis dalam bentuk vektor:

$$\text{Vec}(\mathbf{Y}) = (\mathbf{I}_q \otimes \mathbf{J}_{o/n}) \text{Vec}(\mathbf{B}) + \text{Vec}(\mathbf{E})$$

Pendeteksian outlier pada dalam model linear multivariat dengan asumsi seperti pada (3) dimulai dengan membuat fungsi likelihood untuk populasi $L(\mathbf{B}, \Sigma)$, lalu menentukan $\hat{\mathbf{B}}$ dan $\hat{\Sigma}$. Estimasi parameter \mathbf{B} pada (7) dengan fungsi likelihood seperti pada (4) dengan Metode MLE dimulai dengan melogoritmanaturalkan (4), sehingga:

$$\ln L(\mathbf{B}, \Sigma) = \frac{nq}{2} \ln |\Sigma| - \frac{n}{2} \text{tr} \left[\Sigma^{-1} (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^T (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}) \right] \tag{8}$$

Kondisi optimal (maksimum atau minimum) dicapai bila memenuhi kondisi berikut ini:

Jika logaritma natural dari fungsi likelihood (8) diturunkan terhadap \mathbf{B} dan disamakan dengan nol, maka:

$$\frac{\partial \ln L(\mathbf{B}, \Sigma)}{\partial \mathbf{B}} = -\frac{1}{2} \text{tr} \left(-2\Sigma^{-1} (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}) \mathbf{X}^T \right) \mathbf{0}$$

diperoleh:

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

Kemudian jika (8) diturunkan terhadap Σ dan disamakan dengan nol, maka:

$$\frac{\partial \ln L(\mathbf{B}, \Sigma)}{\partial \Sigma} = -\frac{1}{2} \text{tr} \left(n\hat{\Sigma}^{-1} - \hat{\Sigma}^{-1} \hat{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^T (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}) \right) \mathbf{0}$$

diperoleh:

$$\hat{\Sigma} = \frac{1}{n} (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^T (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})$$

Untuk memaksimumkan fungsi likelihood $L(\mathbf{B}, \Sigma)$ dengan kendala $(\text{Vec}(\hat{\mathbf{B}}) - \text{Vec}(\mathbf{B}))^T (\text{Var}(\text{Vec}(\hat{\mathbf{B}})))^{-1} (\text{Vec}(\hat{\mathbf{B}}) - \text{Vec}(\mathbf{B}))$ jika ada m buah variabel dependen adalah outlier dengan menggunakan pengganda Lagrange, dimulai dengan membuat $\ln L(\mathbf{B}, \Sigma)$ dan menentukan $\text{Vec}(\hat{\mathbf{B}})$ dengan program nonlinear. Metode ini disebut dengan Metode LDL yaitu suatu metode yang menghilangkan pengamatan yang outlier pada model dengan kendala yang menggunakan pengganda Lagrange.

Fungsi likelihood untuk \mathbf{B} yaitu (4), sehingga fungsi likelihood untuk $\hat{\mathbf{B}}_{A_m}^C$ adalah:

$$L(\mathbf{B}_{A_m}^C, \Sigma_{A_m}^C) = (2\pi)^{-\frac{(n-m)p}{2}} \left| \Sigma_{A_m}^C \right|^{-\frac{n-m}{2}} \exp \left[-\frac{1}{2} \text{tr} \left((\Sigma_{A_m}^C)^{-1} (\mathbf{Y}_{A_m}^C - \mathbf{X}_{A_m}^C \hat{\mathbf{B}}_{A_m}^C)^T (\mathbf{Y}_{A_m}^C - \mathbf{X}_{A_m}^C \hat{\mathbf{B}}_{A_m}^C) \right) \right]$$

dimana

$$\square = \frac{1}{2} \text{tr} \left((\Sigma_{A_m}^C)^{-1} (\mathbf{Y}_{A_m}^C - \mathbf{X}_{A_m}^C \hat{\mathbf{B}}_{A_m}^C)^T (\mathbf{Y}_{A_m}^C - \mathbf{X}_{A_m}^C \hat{\mathbf{B}}_{A_m}^C) \right)$$

Berdasarkan (7), diperoleh:

$$\mathbf{Y}_{A_m}^C = \mathbf{X}_{A_m}^C \mathbf{B}_{A_m}^C; \text{ dan}$$

$$\text{Vec}(\mathbf{Y}_{A_m}^C) \sim N_p(\mathbf{0}, \square \mathbf{I}_{-n})_m,$$

dan berdasarkan (2) diperoleh:

$$\hat{\mathbf{B}}_{A_m}^C = \left(\left(\mathbf{X}_{A_m}^C \right)^T \right)^{-1} \mathbf{X}_{A_m}^C \mathbf{Y}_{A_m}^C$$

dimana

$$\left(\mathbf{X}_{A_m}^C \right)^T \mathbf{X}_{A_m}^C = \mathbf{X}^T \mathbf{X} - \dots$$

$$\left(\mathbf{X}_{A_m}^C \right)^T \mathbf{Y}_{A_m}^C = \mathbf{X}^T \mathbf{Y} - \mathbf{X}_{A_m}^T \mathbf{Y}_{A_m}$$

Estimasi dari \mathbf{B} setelah outlier dikeluarkan ($\hat{\mathbf{B}}_{A_m}^C$) adalah:

$$\hat{\mathbf{B}}_i^c = \left(\mathbf{X}'\mathbf{X} \quad -\mathbf{X}'_i \right) \mathbf{X}_i^{-1} \left(\mathbf{X}'\mathbf{Y} \right)$$

$$= \hat{\mathbf{B}} - \left(-\mathbf{X}'\mathbf{X}^{-1} \mathbf{I}_i \left(\mathbf{I} - \mathbf{Q} \right)^{-1} \mathbf{Q} \right) \mathbf{Y}_i \left(\mathbf{I} - \mathbf{Q} \right)^{-1}$$

sehingga:

$$\hat{\mathbf{B}}_{A_m}^c = \hat{\mathbf{B}} - \left(\mathbf{X}'\mathbf{X} \right)^{-1} \mathbf{X}'_{A_m} \left(\mathbf{Q} - \mathbf{E}_{A_m}^{-1} \right)^{-1} \mathbf{E}_{A_m}$$

dan diperoleh pula:

$$\hat{\mathbf{B}}_{A_m}^c \sim N \left(\hat{\mathbf{B}}_i, \mathbf{X}'\mathbf{X}^{-1} \mathbf{I}_i \mathbf{Q} + \square \right)$$

dimana

$$\square = \left(\left(\mathbf{X}'\mathbf{X} \right)^{-1} \mathbf{X}'_{A_m} \left(\mathbf{Q} - \mathbf{E}_{A_m}^{-1} \right)^{-1} \text{Var} \left(\mathbf{X}'\mathbf{X} \right)^{-1} \left(\mathbf{I}_i - \mathbf{Q} \right) \right)^T$$

Permasalahan di atas bersifat umum, sehingga nilai optimal yang diperoleh bisa saja bukan nilai yang paling optimal. Oleh karena itu digunakan pengganda *Lagrange*, sehingga nilai optimal yang diperoleh diharapkan merupakan nilai yang paling optimal pada daerah kepercayaan yang telah ditentukan. Secara umum daerah kepercayaan $(1-\alpha)$ 100% untuk $\text{Vec}(\mathbf{B})$ adalah:

$$\left(\text{Vec}(\hat{\mathbf{B}}_i^c) - \text{Vec}(\hat{\mathbf{B}}_i^c) \right)^T \left(\hat{\text{Var}}(\text{Vec}(\hat{\mathbf{B}}_i^c)) \right)^{-1} \left(\text{Vec}(\hat{\mathbf{B}}_i^c) - \text{Vec}(\hat{\mathbf{B}}_i^c) \right) \leq K$$

sehingga daerah kepercayaan $(1-\alpha)$ 100 untuk model dimana *outlier*-nya telah dihilangkan $\left(\text{Vec}(\hat{\mathbf{B}}_{A_m}^c) \right)$ adalah:

$$\left(\text{Vec}(\hat{\mathbf{B}}_{A_m}^c) - \text{Vec}(\hat{\mathbf{B}}_{A_m}^c) \right)^T \left(\hat{\text{Var}}(\text{Vec}(\hat{\mathbf{B}}_{A_m}^c)) \right)^{-1} \left(\text{Vec}(\hat{\mathbf{B}}_{A_m}^c) - \text{Vec}(\hat{\mathbf{B}}_{A_m}^c) \right) \leq K \tag{9}$$

untuk $F_{v_1, v_2, \alpha} = K$

dimana $v_1 = p$ dan $v_2 = n - p - 1$.

Untuk menyelesaikan permasalahan nonlinear di atas digunakan pengganda *Lagrange*. Misalkan fungsi tujuan (8) dengan kendala (9), maka fungsi *Lagrange*-nya dapat ditulis:

$$F(\mathbf{B}, \Sigma, \lambda) = \ln L(\mathbf{B}, \Sigma) - \lambda \{ F^T \cdot \text{S.F.} - K \} \tag{10}$$

dimana:

$$F = \left(\text{Vec}(\hat{\mathbf{B}}_{A_m}^c) - \text{Vec}(\hat{\mathbf{B}}_{A_m}^c) \right)^T$$

$$\text{dan } S = \hat{\text{Var}} \left(\text{Vec}(\hat{\mathbf{B}}_{A_m}^c) \right)^{-1}$$

Kondisi optimal dicapai bila memenuhi kondisi berikut ini:

Fungsi *Lagrange* (10) diturunkan terhadap \mathbf{B} dan disamakan dengan nol, maka:

$$\frac{\partial F(\mathbf{B}, \Sigma, \lambda)}{\partial \mathbf{B}} = \frac{\partial \ln L(\mathbf{B}, \Sigma)}{\partial \mathbf{B}} = 0$$

Berdasarkan (8), maka diperoleh:

$$\hat{\mathbf{B}} = \left(\mathbf{X}'\mathbf{X} \right)^{-1} \mathbf{X}'\mathbf{Y}$$

Dengan menggunakan sifat hasil kali *kroncker* diperoleh:

$$\text{Vec}(\hat{\mathbf{B}}) = \left(\mathbf{I}_q \otimes \left(\mathbf{X}'\mathbf{X} \right)^{-1} \mathbf{X}' \right) \text{Vec}(\mathbf{Y})$$

$$\text{dengan } E(\text{Vec}(\hat{\mathbf{B}})) = \text{Vec}(\mathbf{B})$$

$$\text{dan } \hat{\text{Var}}(\text{Vec}(\hat{\mathbf{B}})) = \hat{\Sigma} \otimes \left(\mathbf{X}'\mathbf{X} \right)^{-1}$$

Sehingga untuk pengamatan dalam model linear multivariat dengan m *outlier* dihilangkan, diperoleh:

$$\hat{\mathbf{B}}_{A_m}^c = \hat{\mathbf{B}} - \left(\mathbf{X}'\mathbf{X} \right)^{-1} \mathbf{X}'_{A_m} \left(\mathbf{Q} - \mathbf{E}_{A_m}^{-1} \right)^{-1} \mathbf{E}_{A_m}$$

$$\text{dengan } E(\text{Vec}(\hat{\mathbf{B}}_{A_m}^c)) = \text{Vec}(\mathbf{B}_{A_m}^c)$$

Fungsi *Lagrange* (10) diturunkan terhadap Σ dan disamakan dengan nol, maka:

$$\frac{\partial F(\mathbf{B}, \Sigma, \lambda)}{\partial \Sigma} = \frac{\partial \ln L(\mathbf{B}, \Sigma)}{\partial \Sigma} = 0$$

Berdasarkan (8), maka diperoleh:

$$\hat{\Sigma} = \frac{1}{n} (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^T (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})$$

Fungsi *Lagrange* (10) diturunkan terhadap λ dan disamakan dengan nol, maka diperoleh:

$$\left(\text{Vec}(\hat{\mathbf{B}}_i^c) - \text{Vec}(\hat{\mathbf{B}}_i^c) \right)^T \left(\hat{\text{Var}}(\text{Vec}(\hat{\mathbf{B}}_i^c)) \right)^{-1} \left(\text{Vec}(\hat{\mathbf{B}}_i^c) - \text{Vec}(\hat{\mathbf{B}}_i^c) \right) = K$$

Selanjutnya menentukan nilai estimasi dari parameter yang optimal dengan metode numerik yang menggunakan kondisi Karush-Kuhn-Tucker (KKT) pada program nonlinear.

Program nonlinear Metode LDL dengan kondisi KKT adalah:

Maksimumkan:

$$F(\mathbf{B}, \Sigma, \lambda) = \ln L(\mathbf{B}, \Sigma) - \lambda \{ F^T \cdot \text{S.F.} - K \}$$

Kendala:

$$F^T \cdot \text{S.F.} \leq K$$

$$\text{Vec}(\hat{\mathbf{B}}) = \left(\mathbf{I}_q \otimes \left(\mathbf{X}'\mathbf{X} \right)^{-1} \mathbf{X}' \right) \text{Vec}(\mathbf{Y})$$

$$\hat{\Sigma} = \frac{1}{n} (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^T (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})$$

$$\lambda \{ F^T \cdot \text{S.F.} - K \} = 0$$

Solusi terakhir akan mendapatkan X_1, X_2, \dots, X_p yang optimal.

Untuk kasus khusus θ_1 dari θ , maka LD dapat dimodifikasi sebagai:

$$\text{LDL}_{A_m}(\theta_1 | \theta_2) = \frac{1}{2} \ln \left(\hat{\Phi} - \ln \left(\hat{\Phi}_{1, A_m} \hat{\Phi}_2 \hat{\Phi}_4 \right) \right)$$

dimana $\theta = (\mathbf{B}, \Sigma)$, $\theta_1 = \mathbf{B}$, $\theta_2 = \Sigma$, sehingga $\hat{\theta} = (\hat{\mathbf{B}}, \hat{\Sigma})$, dan diperoleh:

$LDL_{A_m}(\mathbf{B}\hat{\Sigma}) = 2(\ln L(\hat{\mathbf{B}}\hat{\Sigma}) - \ln L(\hat{\mathbf{B}}^c\hat{\Sigma}(\hat{\mathbf{B}}^c)))$
 Fungsi *likelihood* dengan kendala sebanyak m pengamatan yang diduga *outlier* adalah:

$$L(\hat{\mathbf{B}}\hat{\Sigma}, \hat{\mathbf{B}}^c) = (2\pi)^{\frac{mn}{2}} |\hat{\Sigma}|^{-\frac{n}{2}} \exp(\mathbf{A})$$

dimana:

$$\mathbf{A} = -\frac{1}{2} \text{tr} \left\{ (\hat{\Sigma}(\hat{\mathbf{B}}^c))^{-1} (\mathbf{Y} - \mathbf{X}_A^c \hat{\mathbf{B}}_A^c)^T (\mathbf{Y} - \mathbf{X}_A^c \hat{\mathbf{B}}_A^c) \right\}$$

$$\hat{\mathbf{B}}_A^c = \hat{\mathbf{B}} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_A^T (\mathbf{Q} - \hat{\Sigma}_A^{-1} \hat{\mathbf{E}}_A)$$

dan

$$\hat{\Sigma}(\hat{\mathbf{B}}) = \frac{1}{n} \left(\hat{\mathbf{E}}_A^T \mathbf{I} \hat{\mathbf{E}}_A + \mathbf{Q}_A \mathbf{Q}_A^T \right)$$

Setelah mendapatkan $\hat{\mathbf{B}}_A^c$ dan $\hat{\Sigma}(\hat{\mathbf{B}}_A^c)$, selanjutnya ditentukan fungsi *Likelihood* untuk pengamatan yang ada *outlier* yaitu:

$$LDL_{A_m} = LDL_{A_m}(\hat{\mathbf{B}}\hat{\Sigma}) = 2(\ln L(\hat{\mathbf{B}}\hat{\Sigma}) - \ln L(\hat{\mathbf{B}}^c\hat{\Sigma}(\hat{\mathbf{B}}^c))) \quad (11)$$

Dengan melogaritmaturalkan dan membuang $\hat{\mathbf{B}}_A^c$ pada (11) di atas, diperoleh:

$$LDL_{A_m} = n \left(\frac{1}{n} \hat{\mathbf{E}}_A^T \mathbf{C}_{A_m} \hat{\mathbf{E}}_A + \frac{1}{n} |\hat{\Sigma}| \right) \quad (12)$$

dimana $\mathbf{C}_{A_m} = (\mathbf{Q} - \mathbf{Q}_A)^T \mathbf{I} (\mathbf{Q} - \mathbf{Q}_A)^{-1}$.

LDL_{A_m} didekati dengan $LDL_{A_m} = \sum_{i=1}^m \frac{v_i}{v_1} \frac{v_2}{v_1} Z_i^2$

sehingga $LDL_{A_m} \sim F_{v_1, v_2}$,

dimana $v_1 = p$ dan $v_2 = n - p - 1$.

Penentuan *outlier* dilakukan dengan membandingkan $LDL_{A_{hitung}}$ pada (12) dan F_{tabel} dengan uji hipotesis adalah:

H_0 : A_m bukan *outlier* dan

H_1 : A_m adalah *outlier*.

Jika $LDL_{A_{hitung}} > \lambda.F_{tabel}$, maka tolak H_0 , artinya pengamatan tersebut adalah *outlier*.

KESIMPULAN

Penentuan dan pendeteksian *outlier* pada model linear multivariat menggunakan pengganda *Lagrange* yang disebut Metode

LDL. Pengganda *Lagrange* yang digunakan adalah daerah kepercayaan dari vektor parameter dimana pengamatan yang *outlier* telah dihilangkan dari model. Penentuan nilai estimasi dari parameter yang optimal dilakukan melalui metode numerik dengan menggunakan kondisi KKT pada program nonlinear. Metode ini membandingkan nilai LDL hitung dengan nilai tabel F yang akan mengikuti distribusi F .

DAFTAR PUSTAKA

Adnan R, Mohamad MN & Setan H. 2003. Multiple Outliers Detection Procedures in Linear Regression. *Matematika*. 1:29-45.
 Bazaara MS, Sherali HD & Shetty CM. 1993. *Nonlinear Programming: Theory and Algorithms*. 2nd edition. John Wiley & Sons. New York.
 Christensen R. 1991. *Linear Model for Multivariate, Time Series, and Spatial Data*, Springer-Verlag. New York.
 Cook RD. 2000. Detection of Influential Observation in Linear Regression. *Technometrics*. 42(1):65-68.
 Diaz-Garcia JA, Gonzalez-Farias G & Alvarado-Castro V. 2007. Exact Distributions for Sensitivity Analysis in Linear Regression. *Applied Mathematical Sciences*. 22:1083-1100.
 Hawkins DM. 1980. *Identifications of Outliers*. Chapman and Hall. New York.
 Makkulau, Linuwih S, Purhadi & Mashuri M. 2008. Prosedur Pendeteksian Outlier pada Model Linear Multivariat dengan Metode Likelihood Displacement Statistic. *Prosiding Seminar Nasional Matematika IV*. Jurusan Matematika FMIPA ITS. Desember 2008. Surabaya.
 Makkulau, Linuwih S, Purhadi & Mashuri M. 2009. Pendeteksian Outlier Model Linear Multivariat pada Produksi Gula dan Tetes Tebu. *Prosiding Seminar Nasional Matematika*. Jurusan Matematika FMIPA Universitas Jember, Februari 2009. Jember.
 Rencher AC & Schaalje GB. 2008. *Linear Models in Statistics*. 2nd edition. John Wiley & Sons. New York.
 Rousseeuw PJ & Hubert M. 1997. Recent Developments in PROG-RESS, dalam L1-Statistical Procedure and Related Topics, edited by Y. Dodge. *Institute of Mathematical Statistics Lecture Notes and Monograph Series*. Hayward, California. 31:201-214.
 Srivastava. MS & Von Rosen D. 1998. Outliers in Multivariate Regression Models. *Journal of Multivariate Analysis*. 65:195-208.
 Xu J, Abraham B & Steiner SH. 2005. *Outlier Detection Methods in Multivariate Regression Models*. <http://www.bisrg.uwaterloo.ca/archive/R-R-06-07.pdf> [4 April 2007].

