

Analisis Regresi Komponen Utama Robust dengan Metode *Minimum Covariance Determinant – Least Trimmed Square* (MCD-LTS)

Siska Diah Ayu Larasati¹, Khoirin Nisa^{1,*} dan Eri Setiawan¹

¹Jurusan Matematika, Fakultas MIPA, Universitas Lampung
Jl. Soemantri Brojonegoro 1 Bandar Lampung

*Email korespondensi: khoirin.nisa@fmipa.unila.ac.id

Dikirim: 27-02-2020, Diterima: 17-03-2020, Diterbitkan: 31-03-2020

Abstrak

Regresi Komponen Utama (RKU) merupakan metode yang digunakan untuk mengatasi masalah multikolinearitas dengan mereduksi dimensi variabel bebas sehingga diperoleh variabel baru yang lebih sederhana tanpa kehilangan sebagian besar informasi yang terkandung pada variabel bebasnya. Apabila pada data pengamatan terindikasi adanya pencilan, maka diperlukan metode yang tegar terhadap pencilan yaitu metode RKU *robust*. Dalam paper ini kami menggunakan metode *robust* yang merupakan kombinasi antara Analisis Komponen Utama *Robust* dengan metode *Minimum Covariance Determinant* (MCD) dan Analisis Regresi *Robust* dengan metode *Least Trimmed Square* (LTS). Tujuan dari penelitian ini adalah mengkaji analisis RKU *robust* dengan metode MCD-LTS serta mengetahui ketegaran RKU *robust* dengan melihat kepekaannya terhadap pencilan. Hasil yang diperoleh dibandingkan dengan RKU klasik berdasarkan nilai bias dan *Mean Square Error* (MSE) pada beberapa ukuran sampel dan persentase pencilan yang berbeda. Hasil dari penelitian ini menunjukkan bahwa RKU *robust* menggunakan MCD-LTS efektif dan efisien dalam mengatasi masalah multikolinearitas dan pencilan pada analisis regresi.

Kata kunci : multikolinearitas, pencilan, regresi komponen utama, *robust*.

Abstract

Principal Component Regression (PCR) is a method used to overcome multi collinearity problems by reducing the dimensions of independent variables to obtain new simpler variables without losing most of the information contained in the variables. If the data analyzed contain outliers, a robust method on PCR is required. In this paper we use a robust method which is a combination of Robust Principal Component Analysis using the Minimum Covariance Determinant (MCD) method and Robust Regression Analysis using Least Trimmed Square (LTS) method. The purpose of this study is to examine the robust PCR analysis using the MCD-LTS method and to know the robustness of the method by looking at its sensitivity to outliers. Results for this purpose we compared the MCD-LTS PCR to the classic PCR based on the bias and Mean Square Error (MSE) values on several different sample sizes and percentages of outliers. The results of this study indicate that robust PCR using MCD-LTS is effective and efficient in overcoming the problem of multicollinearity and outliers in regression analysis.

Keywords: multicollinearity, outliers, principal component regression, *robust*.

1. Pendahuluan

Analisis regresi merupakan salah satu metode analisis data yang digunakan untuk menyelidiki hubungan antara variabel respon dengan satu atau beberapa variabel bebas. Apabila hubungan yang diselidiki terdiri dari satu variabel respon dan satu variabel bebas maka disebut dengan analisis regresi linear sederhana. Sedangkan analisis regresi linear berganda terdiri dari satu variabel respon dan lebih dari satu variabel bebas. Bentuk persamaan regresi dalam bentuk variabel asal X dapat ditulis sebagai berikut:

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon, \quad (1)$$

dengan Y merupakan variabel tak bebas, X_j variabel bebas ke- j ($j = 1, 2, \dots, p$), β_0 dan β_j adalah parameter-parameter regresi, dan ε merupakan galat.

Permasalahan yang sering dihadapi dalam analisis regresi linear berganda yaitu adanya multikolinieritas. Masalah ini terjadi ketika adanya korelasi yang kuat antara variabel bebas. Hal ini dapat menyebabkan $\mathbf{X}^T\mathbf{X}$ memiliki kondisi buruk (*ill condition*) atau hampir singular yang pada akhirnya akan menyebabkan nilai penduga ragam bagi parameter regresi menjadi lebih besar [1].

Salah satu metode statistik yang digunakan untuk mengatasi masalah multikolinieritas adalah regresi komponen utama (RKU). Pada RKU klasik terdapat dua tahap yaitu komponen utama dibentuk menggunakan vektor eigen dari matriks kovarian sampel (\mathbf{S}) klasik dan diregresikan terhadap Y dengan metode kuadrat terkecil [2]. Variabel baru sebagai komponen utama (Q) adalah hasil transformasi dari variabel asal (X) yang modelnya dalam bentuk matriks yaitu $\mathbf{Q} = \mathbf{A}\mathbf{X}$, dan komponen utama ke- j ditulis:

$$\begin{aligned} Q_j &= a_{1j}X_1 + a_{2j}X_2 + \dots + a_{jp}X_p \\ &= \mathbf{a}_j^T \mathbf{X} \end{aligned} \quad (2)$$

dimana vektor pembobot \mathbf{a}_j^T diperoleh dengan memaksimalkan keragaman komponen utama ke- j , yaitu $\mathbf{S}_{Q_j}^2 = \mathbf{a}_j^T \mathbf{S} \mathbf{a}_j$ dengan kendala $\mathbf{a}_j^T \mathbf{a}_j = 1$ serta $\mathbf{a}_l^T \mathbf{a}_j = 0$, untuk $l \neq j$. Vektor pembobot \mathbf{a}_j^T diperoleh dari matriks kovarian $\mathbf{\Sigma}$ yang diduga dengan matriks \mathbf{S} , yaitu $\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$. Vektor \mathbf{a}_j^T yang memenuhi kendala diatas adalah vektor eigen dari matriks kovarian $\mathbf{\Sigma}$. Model regresi komponen utama dapat ditulis sebagai berikut:

$$Y = \beta_0 + \beta_1 Q_1 + \dots + \beta_k Q_k + \varepsilon \quad (3)$$

dengan $k \leq p$.

Analisis RKU seperti di atas sangat sensitif terhadap pencilan (*outliers*) dan akan menghasilkan dugaan parameter yang bias akibat terpengaruh oleh data pencilan. Untuk itu maka diperlukan metode yang tegar terhadap pencilan yang disebut sebagai metode *robust*. Penelitian yang berkaitan dengan RKU *robust* telah dilakukan oleh para peneliti, diantaranya dapat dilihat pada [3-5].

Salah satu metode RKU *robust* untuk mengatasi masalah multikolinieritas dan sekaligus pencilan yaitu dengan menggunakan matriks kovarian *robust* pada analisis komponen utama dan metode regresi *robust* pada tahap analisis regresi. Salah satu metode *robust* bagi matriks kovarian sampel (\mathbf{S}) yaitu metode MCD [6] dan kemudian hasil komponen-komponen utama *robust* yang terbentuk akan diregresikan dengan variabel respon menggunakan metode regresi *robust Least Trimmed Square* (LTS).

Penduga MCD adalah pasangan $(\bar{\mathbf{X}}, \mathbf{S})$, dimana $\bar{\mathbf{X}}$ adalah vektor rata-rata dan \mathbf{S} adalah matriks kovariansi yang meminimumkan nilai determinan \mathbf{S} pada subsampel yang berisikan tepat sebanyak h anggota dari n pengamatan, dimana nilai standar dari $h = [(n + k + 1)/2]$. Pada populasi dengan jumlah pengamatan yang kecil, penduga MCD dapat dengan cepat dihitung dan ditemukan. Tetapi jika jumlah pengamatan besar, maka akan banyak sekali kombinasi subsampel dari H yang harus ditemukan dan penghitungan pun akan cukup memakan waktu. Dalam mengatasi keterbatasan ini, maka digunakan suatu algoritma baru untuk metode MCD yang dinamakan dengan metode *fast-MCD* [7-8].

Langkah-langkah penduga MCD dalam menduga nilai koefisien regresi dengan menggunakan *fast-MCD* dimulai dengan mengambil subsampel pertama dari data pengamatan secara acak, misalkan subsampel tersebut H_1 dengan jumlah elemen sebanyak h . Selanjutnya menghitung vektor rata-rata $\bar{\mathbf{X}}_{MCD}$ dan matriks kovarians \mathbf{S}_{MCD} dari H_1 dengan memisalkan $\bar{\mathbf{X}}_1$ dan \mathbf{S}_1 menggunakan persamaan berikut

$$\bar{\mathbf{X}}_{MCD} = \frac{1}{h} \sum_{i \in H} \mathbf{x}_i, \quad \mathbf{x}_i \text{ merupakan vektor pengamatan ke-}i, \quad (4)$$

$$\mathbf{S}_{MCD} = \frac{1}{h} \sum_{i \in H} [\mathbf{x}_i - \bar{\mathbf{X}}_{MCD}][\mathbf{x}_i - \bar{\mathbf{X}}_{MCD}]^T. \quad (5)$$

Jika $\det(\mathbf{S}_1) = 0$, maka berhenti. Namun apabila $\det(\mathbf{S}_1) \neq 0$, maka dilanjutkan dengan menghitung jarak *robust* (*robust distance* = RD) untuk setiap pengamatan yang kemudian diurutkan dari yang terkecil hingga terbesar menggunakan persamaan berikut

$$RD_i = \sqrt{(\mathbf{x}_i - \bar{\mathbf{X}}_{MCD})^T \mathbf{S}_{MCD}^{-1} (\mathbf{x}_i - \bar{\mathbf{X}}_{MCD})} \quad (6)$$

Pada kasus subsampel selanjutnya, yaitu H_2 akan diambil sebanyak h pengamatan dengan jarak RD terkecil. Demikian seterusnya hingga mencapai konvergen $(S_{i+1}) = (S_1)$. Selanjutnya memilih himpunan H yang memiliki determinan \mathbf{S}_{MCD} terkecil, serta mencari $\bar{\mathbf{X}}_{MCD}$ dan \mathbf{S}_{MCD} dari himpunan H terpilih.

Metode LTS menduga koefisien regresi dari data yang mengandung pencilan dengan meminimumkan jumlah kuadrat galat terhadap sebaran data yang sudah terpankaskan (*trimmed*) yang sering juga disebut dengan istilah “sebaran terwinsorkan” (*winsorized distribution*) [9]. Metode ini menduga koefisien regresi dengan menggunakan metode *Ordinary Least Square* (OLS) terhadap subhimpunan data berukuran h , yaitu:

$$\begin{aligned} \hat{\beta} &= \arg \min_{\beta} \left(\sum_{i=1}^h e_i^2 \right) \\ &= \arg \min_{\beta} \left(\sum_{i=1}^h (y_i - \hat{y}_i)^2 \right) \end{aligned} \tag{7}$$

dengan $\frac{(3n+k+1)}{4} \leq h \leq n$.

Solusi $\hat{\beta}$ pada Persamaan (7) dapat diperoleh dengan menggunakan konsep turunan seperti pada penyelesaian pendugaan metode OLS. Lain halnya pada LTS, persamaan tersebut dihitung pada subhimpunan H terbaik yang dilakukan dengan menggunakan algoritma resampling dari seluruh kemungkinan subhimpunan yang dapat dibentuk yaitu sebanyak $\binom{n}{h}$. Subhimpunan H yang diperoleh merupakan sebaran data yang sudah terpankaskan [10].

Dalam paper ini akan disajikan analisis RCU *Robust* dengan metode MCD-LTS serta mengkaji ketegaran metode tersebut terhadap pencilan melalui simulasi data.

2. Metodologi Penelitian

Data yang digunakan dalam penelitian ini merupakan data simulasi. Pada simulasi dibangkitkan data menggunakan *software* SAS 9.4 dengan variabel bebas X sebanyak enam variabel ($k = 6$) dan galat berdistribusi normal $N(0,1)$ sehingga Y merupakan kombinasi linear dari k variabel bebas ditambah galat. Masing- masing data dibangkitkan dengan ukuran sampel ($n = 20, 60, 100, 200$) serta 5 jenis persentase pencilan (5%, 10% 15%, 20%, 25%) dan diulang sebanyak 1000 kali. Kemudian pencilan dibangkitkan dari distribusi normal dengan mean 50 dan simpangan baku 0,05 yaitu $\varepsilon^* \sim N(50, (0.05)^2)$.

Untuk mendapatkan data kolinearitas pada setiap himpunan data, X_{ik} dibangkitkan menggunakan simulasi Monte Carlo berdasarkan McDonald & Galarneau [11] dengan persamaan sebagai berikut:

$$X_{ij} = (1 - \rho^2)^{\frac{1}{2}} x_{ij} + \rho x_{i(p+1)} \tag{8}$$

dimana : $i = 1, 2, \dots, n$ dan $j = 1, 2, \dots, p$.

Di sini $x_{i1}, x_{i2}, \dots, x_{i(p+1)}$ merupakan data yang dibangkitkan berdistribusi normal $N(0,1)$ dan ρ ditentukan sehingga korelasi antarvariabel bebas diberikan oleh ρ^2 . Dua himpunan dari variabel yang saling berkorelasi dalam artikel ini dibuat berdasarkan nilai $\rho = 0.99$.

Adapun tahapan simulasi data yang dilakukan adalah sebagai berikut:

1. Membangkitkan matriks data.
2. Melakukan regresi komponen utama klasik dengan langkah sebagai berikut:
 - a. Menghitung matriks kovarian dari variabel asal (X).
 - b. Menghitung nilai eigen λ_i dan vektor eigen a_i dari matriks kovarian.
 - c. Menghitung skor komponen utama Q .
 - d. Memilih komponen utama yang berpadanan dengan nilai eigen terbesar.
 - e. Menghitung nilai duga koefisien regresi komponen utama berdasarkan komponen yang terpilih dengan metode OLS, simpan sebagai $\hat{\beta}_j^{(0)}$, untuk $j = 0, 1$.
3. Membangkitkan sebuah matriks *noise* dari distribusi $N(0, (0.01)^2)$.
4. Membangkitkan matriks pencilan dari distribusi $N(50, (0.05)^2)$ sehingga diperoleh matriks kontaminasi. Elemen dari matriks kontaminasi adalah nol kecuali untuk beberapa elemen yang dijadikan pencilan.
5. Menambahkan matriks *noise* dan matriks kontaminasi pada data simulasi. Sehingga diperoleh matriks data yang telah terkontaminasi pencilan.
6. Melakukan regresi komponen utama klasik seperti langkah 2 pada data yang telah terkontaminasi pencilan.
7. Menyimpan nilai $\hat{\beta}_j^{(s)}$ untuk $j = 0, 1$ pada regresi komponen utama klasik di langkah 6.
8. Melakukan regresi komponen utama *robust* dengan langkah sebagai berikut:
9. Menghitung matriks kovarian MCD dari data yang terkontaminasi pencilan.
10. Menghitung nilai eigen λ_i dan vektor eigen a dari matriks kovarian MCD.
 - a. Menghitung skor komponen utama Q .
 - b. Memilih komponen utama yang memiliki nilai eigen lebih dari 1.

- c. Menghitung nilai duga koefisien regresi komponen utama berdasarkan komponen yang terpilih dengan *Least Trimmed Square* (LTS).
11. Menyimpan nilai $\hat{\beta}_j^{(s)}$ untuk $j = 0, 1$ pada regresi komponen utama *robust*.
 12. Ulangi langkah 4 sampai 10 sebanyak 1000 kali untuk seluruh ukuran data.
 13. Menghitung nilai bias dan *Mean Square Error* (MSE) dari RKU klasik dan RKU *robust* dengan menggunakan rumus sebagai berikut:

$$\text{Bias}(\hat{\beta}_j) = \frac{1}{m} \sum_{i=1}^m |\hat{\beta}_{ij}^{(s)} - \hat{\beta}_j^{(0)}|, \quad \text{MSE}(\hat{\beta}_j) = \frac{1}{m} \sum_{i=1}^m (\hat{\beta}_{ij}^{(s)} - \hat{\beta}_j^{(0)})^2$$

dimana, $j = 0, 1$ dan $m = 1000$.

14. Membandingkan RKU *robust* dengan RKU klasik berdasarkan nilai rata-rata bias dan rata-rata MSE dari seluruh dugaan koefisien regresi yang dihasilkan.

3. Hasil dan Pembahasan

3.1. Analisis Regresi Komponen Utama *Robust*

RKU klasik membentuk komponen utama menggunakan vektor eigen dari matriks kovarian klasik yang kemudian diregresikan menggunakan metode OLS. Metode RKU klasik tanpa pencilan menjadi acuan sebagai nilai koefisien RKU sebenarnya, sedangkan RKU klasik dengan pencilan menjadi acuan untuk mengetahui tingkat ketegaran RKU *robust*.

Variabel Q_1, Q_2, \dots, Q_j disebut komponen utama dari X . Jika matriks kovarian dari variabel asal $X_j, j = 1, 2, \dots, p$ dilambangkan dengan Σ , maka diperoleh varian komponen utama yaitu:

$$\begin{aligned} \text{var}(Q_j) &= E[(Q_j - E(Q_j))(Q_j - E(Q_j))^T] \\ &= E[(\mathbf{a}_j^T \mathbf{X} - E(\mathbf{a}_j^T \mathbf{X}))(\mathbf{a}_j^T \mathbf{X} - E(\mathbf{a}_j^T \mathbf{X}))^T] \\ &= E[(\mathbf{a}_j^T \mathbf{X} - \mathbf{a}_j^T E(\mathbf{X}))(\mathbf{a}_j^T \mathbf{X} - \mathbf{a}_j^T E(\mathbf{X}))^T] \\ &= E[(\mathbf{a}_j^T \mathbf{X} - \mathbf{a}_j^T \boldsymbol{\mu})(\mathbf{a}_j^T \mathbf{X} - \mathbf{a}_j^T \boldsymbol{\mu})^T] \\ &= E[(\mathbf{a}_j^T \mathbf{X} - \mathbf{a}_j^T \boldsymbol{\mu})(\mathbf{X}^T \mathbf{a}_j - \boldsymbol{\mu}^T \mathbf{a}_j)] \\ &= \mathbf{a}_j^T E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] \mathbf{a}_j \\ &= \mathbf{a}_j^T \Sigma \mathbf{a}_j. \end{aligned} \tag{9}$$

Komponen utama ke- j yang merupakan kombinasi linear dari variabel asal yang memaksimalkan $\text{var}(Q_j)$ dan tidak berkorelasi dengan komponen utama yang lain. Oleh karena itu, Q_j harus memenuhi batasan $\mathbf{a}_j^T \mathbf{a}_j = 1, \mathbf{a}_i^T \mathbf{a}_j = 0; i < j$ [12]. Untuk mencari komponen utama pertama, pilih vektor \mathbf{a}_1 yang memaksimumkan varian komponen utama pertama $\mathbf{a}_1^T \mathbf{X}$, yaitu memaksimumkan $\text{var}(Q_1) = \mathbf{a}_1^T \Sigma \mathbf{a}_1$ dengan kendala $\mathbf{a}_1^T \mathbf{a}_1 = 1$. Pendekatan standarnya adalah dengan teknik pengganda Lagrange atau dalam hal ini akan memaksimumkan fungsi $f(\mathbf{a}_1) = \mathbf{a}_1^T \Sigma \mathbf{a}_1$ dengan kendala $\mathbf{a}_1^T \mathbf{a}_1 = 1$. Maka fungsi objektif dari permasalahan ini adalah:

$$\max f(\mathbf{a}_1) = \max\{\mathbf{a}_1^T \Sigma \mathbf{a}_1 - \lambda(\mathbf{a}_1^T \mathbf{a}_1 - 1)\},$$

dimana λ adalah pengganda Lagrange. Maka dari penyelesaian pemaksimuman fungsi objektif di atas dapat diperoleh bahwa $\text{var}(Q_1) = \lambda$; yaitu nilai eigen dari matriks kovarian Σ ; dan \mathbf{a}_1 merupakan vektor eigen dari matriks kovarian Σ yang berpadanan dengan nilai eigen λ . Karena permasalahan dalam pembentukan komponen utama pertama adalah memaksimumkan variannya, maka vektor eigen \mathbf{a}_1 dari Σ yang dipilih adalah yang berpadanan dengan nilai eigen terbesar pertama.

Seperti yang telah dinyatakan di atas, dapat ditunjukkan bahwa komponen utama kedua (Q_2), ketiga (Q_3), keempat (Q_4),... (Q_j), maka vektor eigen $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_p$ adalah vektor eigen dari matriks kovarian Σ yang berpadanan dengan nilai eigen terbesar selanjutnya $\lambda_2, \lambda_3, \lambda_4, \dots, \lambda_p$ dan $\text{var}(\mathbf{a}_j^T \mathbf{X}) = \lambda_j$, untuk $j = 1, 2, \dots, p$. Pada analisis komponen utama (AKU) *robust*, matriks kovarian *robust* MCD digunakan untuk pembentukan komponen-komponen utama. Pada prinsipnya metode MCD adalah mencari subsampel yang anggotanya sebanyak h elemen dari matriks data dengan h merupakan bilangan bulat terkecil dari $\frac{n+k+1}{2}$. Misalkan subsampel itu adalah H_1 maka terdapat sebanyak C_h^n kombinasi yang harus ditemukan untuk mendapatkan dugaan vektor rata-rata dan matriks kovarian. Apabila n yang digunakan kecil, pendugaan MCD mudah dan relatif lebih cepat untuk ditemukan. Tetapi, apabila n besar, maka akan banyak kombinasi subsampel yang diperlukan

untuk mendapatkan pendugaan MCD tersebut. Untuk mengatasi keterbatasan ini, maka digunakan pendekatan FAST-MCD yang dikembangkan oleh Rousseeuw dan Van Driessen [6].

Misalkan terdapat $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$ merupakan himpunan vektor-vektor pengamatan sejumlah n dengan k variabel. Misalkan $H_1 \subset \{1, 2, \dots, n\}$ dengan $|H_1| = h$, maka berdasarkan Persamaan (4)-(5) vektor rata-rata dan matriks kovarian MCD secara berturut-turut diberikan oleh:

$$\bar{\mathbf{X}}_{MCD1} = \frac{1}{h} \sum_{i \in H} \mathbf{x}_i$$

$$\mathbf{S}_{MCD1} = \frac{1}{h} \sum_{i \in H} [\mathbf{x}_i - \bar{\mathbf{X}}_{MCD1}][\mathbf{x}_i - \bar{\mathbf{X}}_{MCD1}]^T$$

Jika $\det(\mathbf{S}_1) \neq 0$, maka definisikan jarak *robust* seperti pada Persamaan (6):

$$RD_i = \sqrt{(\mathbf{x}_i - \bar{\mathbf{X}}_{MCD1})^T \mathbf{S}_{MCD1}^{-1} (\mathbf{x}_i - \bar{\mathbf{X}}_{MCD1})}$$

dengan $i = 1, 2, \dots, n$. Kemudian mengurutkan nilai RD mulai dari yang terkecil hingga terbesar. Selanjutnya ambil H_2 sejumlah h pengamatan dengan jarak terkecil kemudian hitung $\bar{\mathbf{X}}_{MCD2}$ dan \mathbf{S}_{MCD2} dari himpunan H_2 .

Penjelasan MCD mensyaratkan $\det(\mathbf{S}_{MCD1}) \neq 0$, karena jika $\det(\mathbf{S}_{MCD1}) = 0$ maka nilai objektif minimum untuk mendapatkan determinan terkecil telah ditemukan. Proses pada metode MCD akan berhenti jika ditemukan himpunan bagian yang konvergen yaitu $d(\mathbf{S}_{i+1}) = d(\mathbf{S}_i)$. Dengan demikian pilih subsampel H yang memiliki nilai determinan matriks kovarians terkecil. Kemudian dari subsampel yang terpilih akan dicari nilai $\bar{\mathbf{X}}_{MCD}$ dan \mathbf{S}_{MCD} . Berdasarkan hal tersebut, maka untuk AKU *robust* akan diperoleh vektor eigen dari matriks kovarian MCD yang berpadanan dengan nilai eigen terbesar pertama dari matriks kovarian MCD.

Pada paper ini hanya menggunakan satu komponen utama, sehingga pada langkah selanjutnya akan diregresikan antara komponen utama tersebut dengan variabel tak bebas Y . Dengan menggunakan metode LTS diperoleh nilai koefisien regresi $\hat{\beta}_0$ dan $\hat{\beta}_1$.

Penduga LTS diperoleh dengan mempertimbangkan jumlah kuadrat galat:

$$JKG = \sum_{i=1}^h \varepsilon_i^2 = \sum_{i=1}^h (y_i - \beta_0 - \beta_1 Q_i)^2 \tag{10}$$

Untuk mendapatkan nilai β_0 dan β_1 pada Persamaan (10) maka JKG diturunkan terhadap β_0 dan β_1 kemudian menyamakannya dengan nol. Jika JKG diturunkan terhadap β_0 maka diperoleh:

$$\begin{aligned} \frac{\partial JKG}{\partial \beta_0} &= -2 \sum_{i=1}^h (y_i - \beta_0 - \beta_1 Q_i) = 0 \\ \sum_{i=1}^h y_i - h\beta_0 - \beta_1 \sum_{i=1}^h Q_i &= 0 \\ \beta_0 &= \frac{\sum_{i=1}^h y_i - \beta_1 \sum_{i=1}^h Q_i}{h} \end{aligned}$$

Tulis $\bar{y} = \frac{\sum_{i=1}^h y_i}{h}$ dan $\bar{Q} = \frac{\sum_{i=1}^h Q_i}{h}$ maka diperoleh:

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{Q} \tag{11}$$

Selanjutnya jika JKG diturunkan terhadap β_1 dan disamakan dengan nol maka diperoleh

$$\frac{\partial JKG}{\partial \beta_1} = -2 \sum_{i=1}^h (y_i - \beta_0 - \beta_1 Q_i) Q_i = 0.$$

Dengan cara yang sama dapat diperoleh:

$$\hat{\beta}_1 = \frac{h \sum_{i=1}^h y_i Q_i - \sum_{i=1}^h y_i \sum_{i=1}^h Q_i}{h \sum_{i=1}^h Q_i^2 - \sum_{i=1}^h (Q_i)^2} \tag{12}$$

Jadi, $\hat{\beta}_0$ dan $\hat{\beta}_1$ pada Persamaan (11) dan Persamaan (12) merupakan penduga bagi nilai koefisien RKU *robust* yaitu menggunakan metode MCD-LTS.

3.2. Hasil Simulasi

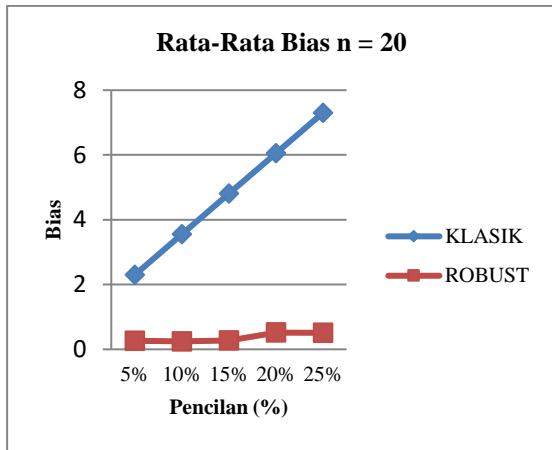
Berikut ini kami sajikan nilai rata-rata bias dan MSE hasil simulasi dengan ulangan sebanyak 1000 kali.

Tabel 1. Rata-Rata Bias dan MSE Pendugaan Koefisien Regresi Komponen Utama

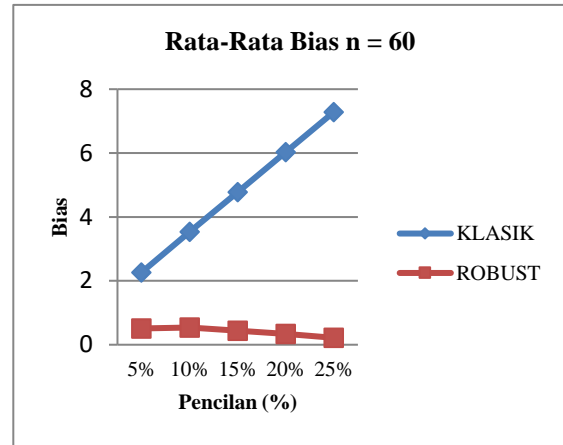
n	Pencilan (%)	Rata-Rata Bias		Rata-Rata MSE	
		RKU Klasik	RKU Robust	RKU Klasik	RKU Robust
20	5	2,2985	0,2673	5,3236	0,2301
	10	3,557	0,2459	14,7344	0,148
	15	4,8047	0,2688	30,3489	0,1727
	20	6,0513	0,5145	52,2118	0,4211
	25	7,3025	0,5111	80,3393	0,4713
60	5	2,2633	0,5075	5,1785	0,4866
	10	3,5315	0,5388	14,6277	0,4946
	15	4,7819	0,4392	30,2556	0,3221
	20	6,0328	0,3416	52,1328	0,1981
	25	7,2864	0,2159	80,2717	0,0882
100	5	2,2817	0,1531	5,254	0,0474
	10	3,5458	0,1026	14,6875	0,0284
	15	4,7951	0,0435	30,3097	0,02
	20	6,0449	0,2353	52,1833	0,0893
	25	7,2975	0,238	80,3188	0,0925
200	5	2,2725	0,1806	5,2161	0,0641
	10	3,522	0,0525	14,5889	0,0121
	15	4,7766	0,0837	30,2328	0,0124
	20	6,0278	0,1374	52,1125	0,0287
	25	7,2747	0,1429	80,2252	0,0224

Pada Tabel 1 terlihat bahwa nilai rata-rata bias dan MSE dari data mengandung pencilan 5% - 25% yang dihasilkan dari metode RKU klasik lebih besar dibandingkan RKU *robust*. Pada RKU klasik, terlihat juga bahwa setiap penambahan persentase dalam satu ukuran sampel maka selalu terjadi peningkatan nilai rata-rata bias dan MSE, bahkan diperoleh nilai rata-rata bias dan MSE > 1 di masing-masing persentase pencilan dan ukuran sampel tersebut. Hal ini dapat diartikan nilai dugaan koefisien RKU klasik semakin buruk. Sedangkan pada RKU *robust* rata-rata bias dan MSE yang dihasilkan selalu lebih kecil atau di bawah nilai rata-rata bias dan MSE RKU klasik yaitu nilai rata-rata bias dan MSE < 1 di masing-masing persentase pencilan dan ukuran sampel tersebut.

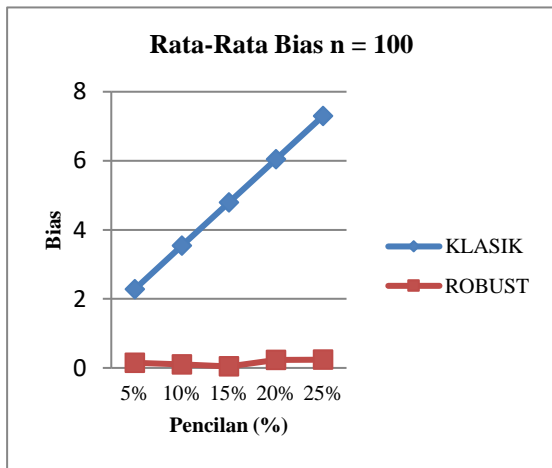
Berdasarkan nilai rata-rata bias dan MSE dengan ukuran sampel yaitu 20, 60, 100, dan 200 serta persentase pencilan yaitu 5%, 10%, 15%, 20%, dan 25% yang disajikan pada Tabel 1, maka akan diperoleh grafik perbandingan dari metode RKU klasik dan *robust* sebagai berikut:



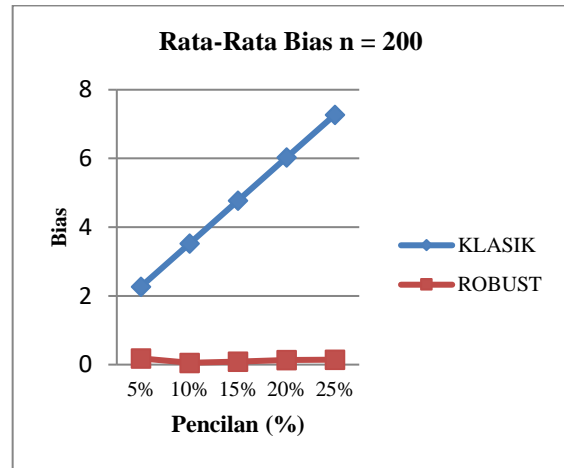
Gambar 1. Grafik Rata-Rata Bias untuk n = 20.



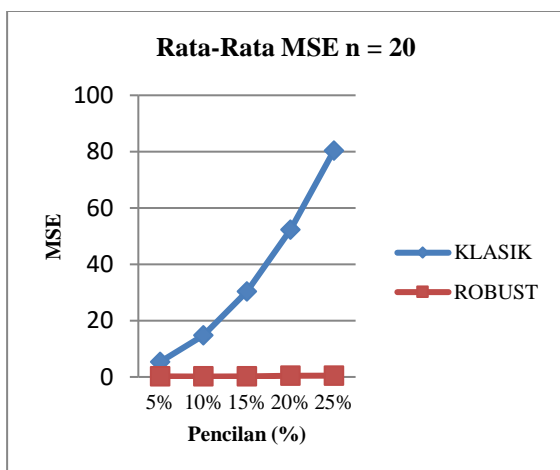
Gambar 2. Grafik Rata-Rata Bias untuk n = 60.



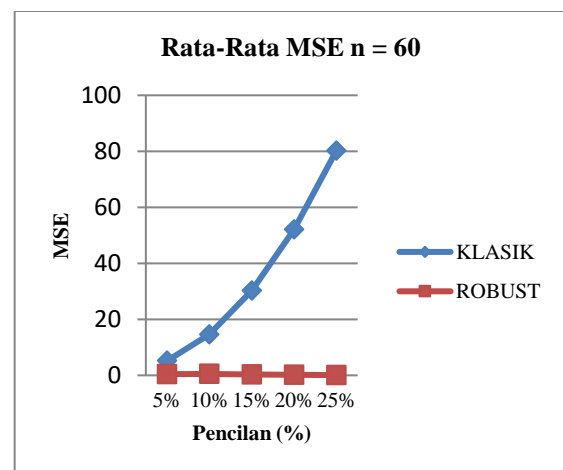
Gambar 3. Grafik Rata-Rata Bias untuk n = 100.



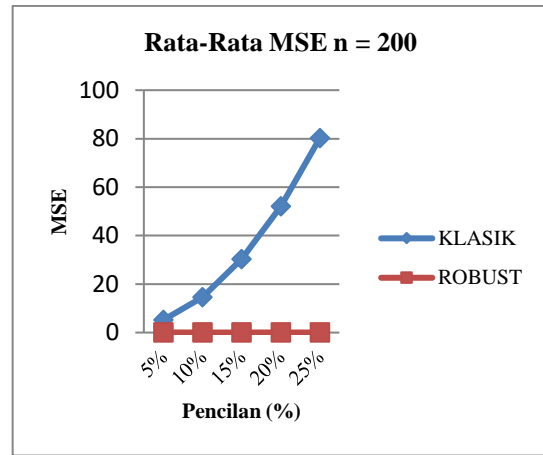
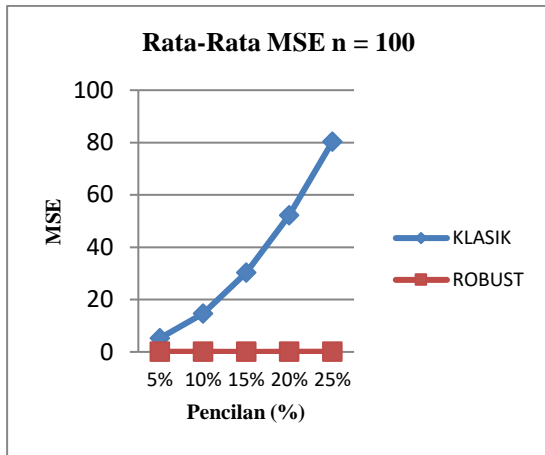
Gambar 4. Grafik Rata-Rata Bias untuk n = 200.



Gambar 5. Grafik Rata-Rata MSE untuk n = 20.



Gambar 6. Grafik Nilai Rata-Rata MSE untuk n = 60.



Gambar 7. Grafik Rata-Rata MSE untuk n = 100. Gambar 8. Grafik Rata-Rata MSE untuk n = 200.

Berdasarkan Gambar 1 – Gambar 4, terlihat bahwa letak garis nilai rata-rata bias RKU klasik dengan ukuran sampel 20, 60, 100, 200 berada di atas garis nilai rata-rata bias RKU *robust*. Gambar 1 – Gambar 4 juga terlihat bahwa titik-titik yang disambungkan garis untuk nilai rata-rata bias pada RKU klasik di setiap persentase pencilan selalu naik secara konstan sedangkan pada RKU *robust* titik-titik yang disambungkan garis untuk nilai rata-rata bias di setiap persentase pencilan tidak secara konstan naik atau turun. Selanjutnya, berdasarkan Gambar 5 – Gambar 8 terlihat bahwa letak garis nilai rata-rata MSE RKU klasik dengan masing-masing ukuran sampel tersebut berada di atas garis nilai rata-rata MSE RKU *robust*. Gambar 5 – Gambar 8 juga terlihat bahwa titik-titik yang disambungkan garis untuk nilai rata-rata MSE pada RKU klasik di setiap persentase pencilan selalu naik secara konstan yaitu berkisar dari angka > 5 hingga angka ≥ 80 sedangkan pada RKU *robust* titik-titik yang disambungkan garis untuk nilai rata-rata MSE di setiap persentase pencilan tidak secara konstan naik atau turun, dapat dikatakan hanya di sekitar angka < 0,5.

4. Kesimpulan

Berdasarkan hasil dan pembahasan, maka diperoleh kesimpulan sebagai berikut:

1. Analisis Komponen Utama *Robust* dapat dilakukan menggunakan matriks kovarian *robust MCD* dan metode regresi *robust LTS* dengan penduga nilai koefisien RKU *robust* sebagai berikut:

$$\hat{\beta}_0 = \frac{\sum_{i=1}^h y_i - \beta_1 \sum_{i=1}^h Q_i}{h}$$

$$\hat{\beta}_1 = \frac{h \sum_{i=1}^h y_i Q_i - \sum_{i=1}^h y_i \sum_{i=1}^h Q_i}{h \sum_{i=1}^h Q_i^2 - \sum_{i=1}^h (Q_i)^2}$$

dimana,

Q = Komponen utama.

h = Banyaknya anggota dalam subsampel terbaik.

2. Pada hasil simulasi diperoleh bahwa pada suatu ukuran sampel dengan presentase pencilan 5%-25%, nilai rata-rata bias dan MSE metode RKU klasik selalu naik secara konstan disetiap penambahan presentase pencilan sedangkan metode RKU *robust* diperoleh nilai rata-rata bias yang lebih kecil dibandingkan RKU klasik. Hal ini menunjukkan bahwa metode RKU *robust* merupakan metode yang kekar terhadap pencilan, sedangkan metode RKU klasik sangat sensitif terhadap danya pencilan. Dengan demikian maka RKU *robust* menggunakan MCD-LTS memberikan hasil yang baik dan merupakan metode efektif dan efisien dalam mengatasi masalah multikolinearitas dan pencilan.

Daftar Pustaka:

- [1] Draper, N.R. dan Smith, H. 2011. *Applied Regression Analysis*, 3rd Edition, India: Wiley.
- [2] Notiragayu dan Nisa, K. 2008. Analisis Regresi Komponen Utama Robust untuk Data mengandung Pencilan. *Jurnal Sains MIPA*. **14** : 1 45-50.
- [3] Abdel Bary, M.N. 2017. Robust Regression Diagnostic for Detecting and Solving Multicollinearity and Outlier Problems: Applied Study by Using Financial Data. *Applied Mathematical Sciences*. **11** : **13** 601 – 622.
- [4] Agarwal, A., Shah, D., Shen, D. dan Song, D. 2020. On Robustness of Principal Component Regression. ArXiv: 1902.10920v7 1–57.
- [5] Oh, T-H., Matsushita, Y., Kweon, I.S. dan Wipf, D. 2016. A Pseudo-Bayesian Algorithm for Robust PCA. *Proceeding of 30th Conference on Neural Information Processing Systems (NIPS)*. Barcelona, Spain.
- [6] Nisa, K., Herawati, N., Setiawan, E., Nusyirwan. (2006). Robust Principal Component Analysis using the Minimum Covariance Determinant Estimator. *Proceeding of International Conference on Mathematics and Natural Sciences (ICMNS)* : 789-792.
- [7] Rousseeuw, P.J. dan Van Driessen, K. 1999. A Fast Algorithm for the Minimum Covariance Determinant Estimator. *Technometrics*. **41** : **3** 212-223.
- [8] Dayanti, N. P., Suciptawati, N.L. dan Susilawati, M. 2016. Penerapan Bootstrap dalam Metode Minimum Covariance Determinant (MCD) dan Least Median Square (LMS) pada Analisis Regresi Linear Berganda., *E-Jurnal Matematika*. **5** : **1** 22-26.
- [9] Nisa, K. 2006. Analisis Regresi Robust Menggunakan Metode Least Trimmed Square untuk Data Mengandung Pencilan. *Jurnal Ilmiah MIPA*. **IX** : **2** 1-9.
- [10] Chen, C. 2002. Robust Regression and Outlier Detection with the ROBUSTREG Procedure. *Proceeding of the 27th SAS User Group International (SUGI) Conference*. Cary, NC: SAS Institute Inc. 1-13.
- [11] McDonald, G.C. dan Galarneau, D.I. 1975. A Monte Carlo Evaluation of Some Ridge-Type Estimators. *Journal of American Statistical Association*. **70** 407-416.
- [12] Jolliffe, I.T. 2002. *Principal Component Analysis second edition*. Springer..