

## ANALISIS REGRESI KOMPONEN UTAMA ROBUST UNTUK DATA MENGANDUNG PENCILAN

Notiragayu\* dan Khoirin Nisa

Jurusan Matematika, FMIPA Universitas Lampung  
\*Alamat korespondensi e-mai : Notiragayu@unila.ac.id

Diterima 5 Februari 2008, perbaikan 17 Maret 2008, disetujui untuk diterbitkan 19 March 2008

### ABSTRACT

Principal Component Regression (PCR) is one of the widely used statistical techniques for regression analysis with colinearity, and a robust technique on PCR when data contains outlier is an important problem. In this paper we consider the problem of robust PCR based on Minimum Volume Ellipsoid (MVE) estimator and Least Trimmed Square (LTS) regression. We aimed to look at the behavior of the principal component regression coefficient resulted by MVE-LTS and compare them with classical estimator through the bias and the mean square error. The result shows that PCR using MVE-LTS is very robust.

**Keywords** : Principal component regression, colinearity, robust

### 1. PENDAHULUAN

Salah satu masalah yang sering dihadapi dalam analisis regresi berganda adalah multikolinieritas, yaitu adanya hubungan linier (korelasi) antara dua peubah bebas atau lebih dalam suatu persamaan regresi. Hal ini menyebabkan matriks  $X^T X$  menjadi singular sehingga persamaan normalnya tidak lagi mempunyai jawaban yang tunggal. Akibatnya dugaan koefisien regresi kuadrat terkecil  $[b = (X^T X)^{-1} X^T Y]$  memiliki simpangan baku yang besar yang menjadikannya tidak lagi terandalkan<sup>1)</sup>.

Regresi komponen utama (RKU) merupakan salah satu metode yang dapat digunakan untuk mengatasi masalah multikolinieritas. Metode ini mengatasi multikolinieritas dengan dua tahapan, tahap pertama analisis komponen utama terhadap peubah-peubah bebas  $X_i$ , dan tahap kedua analisis regresi terhadap komponen-komponen utama dengan peubah respon  $Y$ . Pada RKU klasik, kedua tahap dilakukan secara tradisional, yaitu komponen utama dibentuk menggunakan vektor eigen dari matriks peragam sampel (S) klasik dan diregresikan terhadap  $Y$  dengan metode kuadrat terkecil. Namun analisis seperti ini sangat sensitif terhadap pencilan (*outliers*) dan akan menghasilkan dugaan parameter yang bias akibat terpengaruh oleh data pencilan<sup>2)</sup>. Cara yang paling sederhana untuk mengatasi pencilan pada RKU adalah dengan menggunakan metode *robust* pada kedua tahap.

Analisis komponen utama *robust* dapat dilakukan dengan beberapa cara, diantaranya dengan algoritma proyeksi<sup>3)</sup>, dengan matriks peragam *robust*<sup>4)</sup>, atau

gabungan dari keduanya<sup>5)</sup>. Sedangkan analisis regresi *robust* dapat dilakukan dengan menggunakan metode *alternating regression*<sup>6)</sup>, metode kuadrat terpangkas terkecil<sup>7)</sup>, metode median kuadrat terkecil<sup>8)</sup>, metode nilai mutlak terkecil<sup>9)</sup>, dsb. Dalam tulisan ini, pada tahapan pertama kami memilih untuk menggunakan matriks peragam *robust* dengan menggunakan metode Volume Ellipsoid Minimum (VEM) yang diperkenalkan oleh Peter J. Rousseeuw pada tahun 1985<sup>10)</sup>, dan pada tahapan kedua menggunakan metode Kuadrat Terpangkas Terkecil (KTT) yang diperkenalkan oleh Rousseeuw dan Leroy pada tahun 1987<sup>11)</sup> dan dikenal sangat tegar serta memiliki sifat-sifat statistik yang baik<sup>12)</sup>. Untuk memperlihatkan ketegaran analisis RKU dengan metode VEM-KTT terhadap data mengandung pencilan, kami melakukan simulasi Monte Carlo dan membandingkannya dengan RKU klasik.

#### 1.1. Regresi Komponen Utama

Regresi komponen utama merupakan regresi dari peubah tak-bebas terhadap komponen-komponen utama yang tidak saling berkorelasi, dimana setiap komponen utama merupakan kombinasi linear dari semua peubah bebas yang telah dispesifikasikan sejak awal. Bentuk persamaan regresi dalam bentuk peubah asal  $X$  dapat ditulis sebagai berikut:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon$$

dengan  $Y$  merupakan peubah tak-bebas,  $X_i$  peubah bebas ke- $i$  ( $i = 1, 2, \dots, p$ ),  $\beta_i$  adalah parameter-parameter regresi, dan  $\varepsilon$  merupakan galat.

Peubah baru sebagai komponen utama (K) adalah hasil transformasi dari peubah asal ( $X$ ) yang modelnya dalam

bentuk matriks adalah  $\mathbf{K} = \mathbf{A} \mathbf{X}$ , dan komponen utama ke- $j$  ditulis :

$$K_j = a_{1j} X_1 + a_{2j} X_2 + \dots + a_{pj} X_p = \mathbf{a}_j' \mathbf{x} \quad (1)$$

dimana vektor pembobot  $\mathbf{a}_j'$  diperoleh dengan memaksimalkan keragaman komponen utama ke- $j$ ,

yaitu  $S_{K_j}^2 = \mathbf{a}_j' \mathbf{S} \mathbf{a}_j$  dengan kendala  $\mathbf{a}_j' \mathbf{a}_j = 1$  serta

$\mathbf{a}_h' \mathbf{a}_j = 0$ , untuk  $h \neq j$ . Vektor pembobot  $\mathbf{a}_j'$  diperoleh dari matriks peragam  $\mathbf{\Sigma}$  yang diduga dengan matriks  $\mathbf{S}$ ,

$$\text{yaitu } \mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})'$$

yang memenuhi kendala di atas adalah vektor eigen dari matriks peragam  $\mathbf{\Sigma}$ .

Model regresi komponen utama dapat ditulis sebagai berikut :

$$Y = \beta_0 + \beta_1 K_1 + \beta_2 K_2 + \dots + \beta_m K_m + \varepsilon,$$

dengan  $m \leq p$

### 1.2. Metode Kuadrat Terpangkas Terkecil (KTT)

Metode KTT menduga koefisien regresi dari data yang mengandung pencilan dengan meminimumkan jumlah kuadrat galat terhadap sebaran data yang sudah terpangkas (*trimmed*) atau sebaran terwinsorkan (*winsorized distribution*). Dengan kata lain, metode KTT menduga koefisien regresi dengan menggunakan metode MKT terhadap subhimpunan data terbaik  $H$ , yaitu

$$\hat{\beta} = \arg \min_{\beta} \left( \sum_{i=1}^h e_i^2 \right) = \arg \min_{\beta} \left( \sum_{i=1}^h (y_i - \hat{y}_i)^2 \right),$$

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

dengan  $h$  merupakan banyaknya anggota dalam subhimpunan  $H$ . Penentuan subhimpunan  $H$  terbaik dilakukan dengan menggunakan algoritma resampling<sup>13)</sup> dari seluruh kemungkinan subhimpunan yang didapat

dibentuk yaitu sebanyak  $\binom{n}{h}$ . Subhimpunan  $H$  yang

diperoleh merupakan sebaran data yang sudah terpangkas (*trimmed distribution*)<sup>14)</sup>.

### 1.3. Metode Volume Ellipsoid Minimum (VEM)

Metode Volume Ellipsoid Minimum (VEM) merupakan metode penduga *robust* untuk vektor nilai tengah dan matriks peragam. Pada prinsipnya metode ini adalah mencari ellipsoid dengan volume paling minimum yang melingkupi suatu subhimpunan dari minimal  $h$  pengamatan. Subhimpunan berukuran  $h$  ini disebut *halfset* karena  $h$  sering dipilih lebih dari setengah  $n$  pengamatan. Penduga nilai tengah adalah pusat ellipsoid secara geometris dan penduga matriks peragam adalah matriks pembentuk ellipsoid.

Mencari penduga VEM secara esensial dilakukan dengan dua proses. Bagian pertama mencari *halfset* terbaik yang memuat  $h$  pengamatan. Bagian kedua mencari volume paling minimum dari ellipsoid yang melingkupi *halfset*. Untuk sebuah *halfset* terdiri dari banyak ellipsoid yang dapat melingkupinya. Kedua proses ini dilakukan secara iterative menggunakan algoritma resampling seperti dalam pendugaan matriks peragam *robust* menggunakan metode determinan peragam minimum<sup>15)</sup>.

Penduga VEM didefinisikan sebagai pasangan  $(\bar{\mathbf{X}}, \mathbf{S})$

di mana  $\bar{\mathbf{X}}$  adalah vektor- $p$  dan  $\mathbf{S}$  adalah matriks semi-definit positif  $p \times p$  yang memenuhi :

$$\# \{i \mid (\mathbf{x}_i - \bar{\mathbf{X}})' \mathbf{S}^{-1} (\mathbf{x}_i - \bar{\mathbf{X}}) \leq c^2\} \geq h$$

dengan  $n$  adalah jumlah pengamatan,  $p$  adalah jumlah peubah,  $h = [(n+p+1)/2]$ ,  $c$  suatu konstanta dan  $\mathbf{x}_i$  adalah data pengamatan<sup>16)</sup>.

## 2. METODE PENELITIAN

Dalam penelitian ini kami menggunakan data yang dibangkitkan dengan menggunakan perangkat lunak SAS/IML versi 8. Data terdiri dari 5 peubah bebas  $X_i$  dan sebuah peubah tak bebas  $Y$ . Untuk menciptakan multikolinieritas dalam data kami membangkitkan  $X_i$  dari sebaran normal peubah ganda dengan vektor nilai tengah  $\boldsymbol{\mu} = \mathbf{0}$  dan matriks peragam  $\mathbf{\Sigma}$  nondiagonal yang dibangkitkan secara acak namun dirancang agar memiliki nilai eigen pertama yang proporsinya melebihi 75%. Sedangkan  $Y$  dibangkitkan sebagai kombinasi liner dari  $X_i$  ditambah dengan unsur galat. Model regresi komponen utama yang digunakan adalah :  $\mathbf{Y} = \boldsymbol{\beta} \mathbf{K} + \varepsilon$ , dimana  $\mathbf{K}$  merupakan komponen utama pertama dengan proporsi keragaman di atas 75%. Data pencilan dibangkitkan dari sebaran normal dengan  $\mu = 10$  dan  $\sigma = 1$ .

Simulasi Monte Carlo dilakukan sebanyak 100 replikasi untuk ukuran sampel besar (mengingat metode VEM dan KTT bersifat tak bias asimtotik) yaitu  $n = 100, 300,$  dan  $500$ . Pada setiap sampel diberikan kontaminasi pencilan sebesar  $\theta = 5\%, 10\%, 15\%$  dan  $20\%$ . Nilai bias dan KTG dihitung atas 100 himpunan sampel sebagai berikut :

$$\text{Bias}(\beta) = \left| \left( \frac{1}{100} \sum_{s=1}^{100} \hat{\beta}^{(s)} \right) - \hat{\beta}^{(0)} \right|$$

$$\text{MSE}(\beta) = \frac{1}{100} \sum_{s=1}^{100} (\hat{\beta}^{(0)} - \hat{\beta}^{(s)})^2$$

Dengan  $\hat{\beta}^{(0)}$  adalah koefisien regresi untuk data tanpa pencilan, dan  $\hat{\beta}^{(s)}$  adalah koefisien regresi untuk data yang telah diberi kontaminasi pencilan.

### 3. HASIL DAN PEMBAHASAN

Pengaruh pencilan pada analisis RKU pada setiap tahapan dijelaskan sebagai berikut. Pada tahap pertama, pencilan mempengaruhi matriks peragam dari data X yang secara tak langsung mempengaruhi nilai eigennya. Hal ini mengakibatkan buruknya representasi komponen utama terhadap sebaran data awal karena komponen utama dibentuk berdasarkan vektor eigen matriks peragam seperti pada persamaan (1). Dengan menggunakan metode VEM, nilai eigen yang diperoleh mendekati nilai eigen dari data awal, sehingga skor

komponen utama yang dihasilkannya jauh lebih baik untuk menjelaskan sebaran data awal dibandingkan metode klasik. Pada Tabel 1 disajikan nilai eigen untuk matriks peragam dari X pada n=100.

Pada Tabel 1 di bawah terlihat bahwa metode *robust* jauh lebih baik dibandingkan dengan metode klasik. Semakin besar prosentase pencilan mengakibatkan penyimpangan yang semakin besar terhadap nilai eigen metode klasik, sementara nilai eigen metode *robust* hanya sedikit terpengaruh. Untuk ukuran sampel dan prosentase pencilan lainnya memberikan hasil yang serupa seperti di atas sehingga tidak ditampilkan di sini.

Tabel 1. Nilai eigen dari matriks peragam dari X dengan n=100

Nilai Eigen	data awal	data dengan 5% pencilan		data dengan 10% pencilan	
		Klasik	VEM	Klasik	VEM
$\lambda_1$	14,588879	31,143393	14,054304	49,868027	13,877666
$\lambda_2$	2,7551894	10,041588	2,6013773	11,675904	2,6451922
$\lambda_3$	7,06E-16	18,183967	0,0032795	2,0697765	0,0030162
$\lambda_4$	-2,72E-16	0,0937708	0,002237	0,0853384	0,0028041
$\lambda_5$	-2,18E-15	0,0476747	0,0021113	0,058925	0,0017368

Tabel 1 (lanjutan)

Data dengan 15% pencilan		Data dengan 20% pencilan	
klasik	VEM	klasik	VEM
70,452.787	13,452839	90,000848	12,033517
11,58036	2,6848096	12,210814	2,7617919
2,0126182	0,0028809	21,895263	0,002751
0,1756584	0,0023371	0,2636281	0,0023845
0,0713062	0,0022188	0,1998619	0,0021649

Tabel 2. Nilai dugaan koefisien RKU dengan jumlah pencilan 5%

n	Koefisien RKU		
	$\hat{\beta}^{(0)}$	$\hat{\beta}_{\text{klasik}}$	$\hat{\beta}_{\text{VEM-KTT}}$
100	0,86566	2,10824	1,06507
300	0,87215	2,15147	0,87497
500	0,90962	2,18290	0,89653

Tabel 3. Nilai bias dan KTG koefisien RKU

N	Prosentase pencilan	Bias		KTG	
		Klasik	Robust	Klasik	Robust
100	0%	0,000405	0,087775	0,000009	0,007850
	5%	1,212908	0,033743	1,471467	0,001265
	10%	1,327582	0,079398	1,762496	0,006457
	15%	1,346974	0,085516	1,814347	0,001663
	20%	1,357721	0,109189	1,843410	0,012100
300	0%	0,000132	0,161447	0,000003	0,029577
	5%	1,282254	0,039634	1,644248	0,002745
	10%	1,358500	0,058209	1,805530	0,003438
	15%	1,346712	0,123704	1,813635	0,015334
	20%	1,342722	0,132010	1,842902	0,001529
500	0%	0,000058	0,128654	0,000002	0,022731
	5%	1,256842	0,063360	1,579686	0,004417
	10%	1,319958	0,048786	1,712291	0,002529
	15%	1,309704	0,059196	1,745324	0,004278
	20%	1,327861	0,078258	1,763215	0,006942

Pada tahap kedua, pencilan pada peubah respon Y mempengaruhi dugaan koefisien regresi. Sebagai contoh pada Tabel 2 di atas disajikan dugaan koefisien RKU pada setiap ukuran sampel dengan jumlah pencilan sebesar 5%.

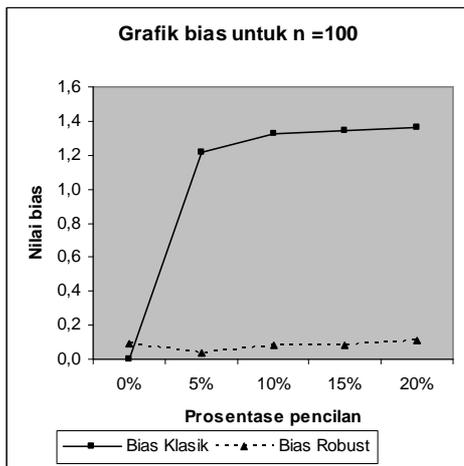
Pada Tabel 2 terlihat bahwa dugaan metode *robust* lebih mendekati dugaan data awal dibandingkan dengan metode klasik. Penyimpangan terhadap dugaan koefisien regresi akan menghasilkan model yang salah sehingga berdampak buruk pada peramalan.

Untuk membandingkan nilai dugaan koefisien RKU secara keseluruhan kami melakukan replikasi sebanyak 100 kali untuk setiap ukuran sampel dan prosentase pencilan yang telah ditentukan. Berdasarkan nilai-nilai dugaan dari 100 replikasi diperoleh nilai bias dan MSE seperti pada Tabel 3.

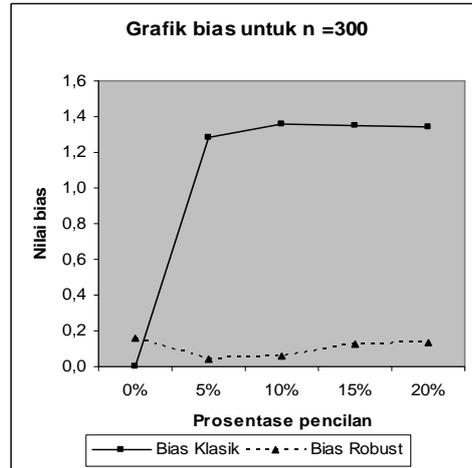
Untuk data tanpa pencilan (0%), nilai bias dan KTG yang dihasilkan oleh metode *robust* VEM-KTT lebih besar dari metode klasik yang berarti bahwa metode klasik lebih baik dari metode *robust*. Namun nilai bias dan KTG metode VEM-KTT pada pencilan 0% masih relatif kecil yaitu bias < 0,2 dan KTG < 0,03 sehingga masih cukup baik sebagai penduga RKU.

Untuk data mengandung pencilan 5% - 20%, nilai bias dan KTG metode klasik lebih besar dari bias dan KTG metode *robust*. Dan perhatikan bahwa penambahan pencilan dalam data meningkatkan nilai bias dan KTG metode klasik yang berarti bahwa nilai dugaannya semakin buruk. Sedangkan nilai bias dan KTG metode *robust* VEM-KTT dan nilainya stabil < 0,2 untuk bias dan < 0,03 untuk KTG. Semakin kecil nilai bias suatu metode pendugaan menunjukkan bahwa nilai dugaan yang dihasilkan mendekati nilai parameter sebenarnya. Dan semakin kecil nilai KTG suatu metode pendugaan menunjukkan bahwa nilai dugaan yang dihasilkan semakin stabil. Untuk mempermudah perbandingan, data bias dan KTG pada Tabel 3 disajikan dalam bentuk grafik pada Gambar 1- 6 di bawah.

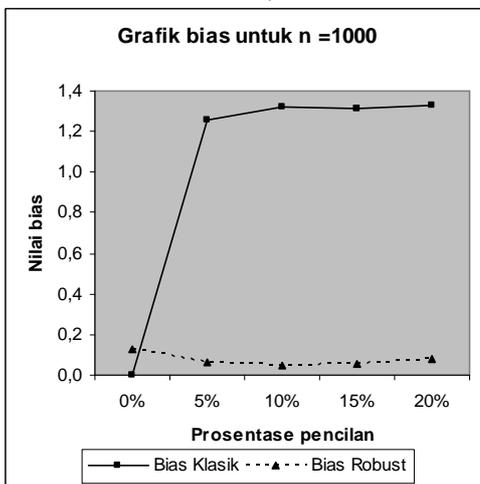
Pada Gambar 1- Gambar 3 dapat dilihat grafik nilai bias dari metode klasik sangat baik untuk pencilan 0% namun meningkat pesat pada saat data diberi pencilan 5% dan seterusnya. Sedangkan bias VEM-KTT pada pencilan 0% lebih besar dari nilai bias pada pencilan 5%-20%. Demikian pula nilai KTG pada Gambar 4 – Gambar 6, nilai KTG metode klasik memiliki pola yang sama dengan nilai biasnya, sedangkan nilai KTG metode *robust* VEM-KTT cenderung stabil < 0,03.



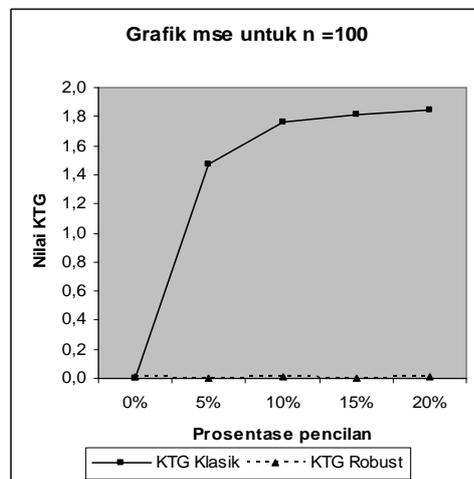
Gambar 1. Grafik bias pada n=100



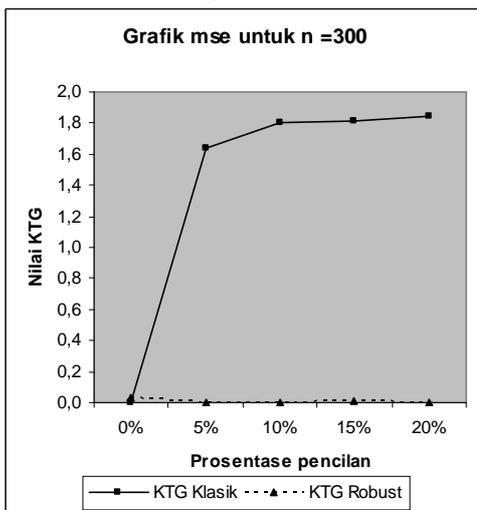
Gambar 2. Grafik bias pada n=300



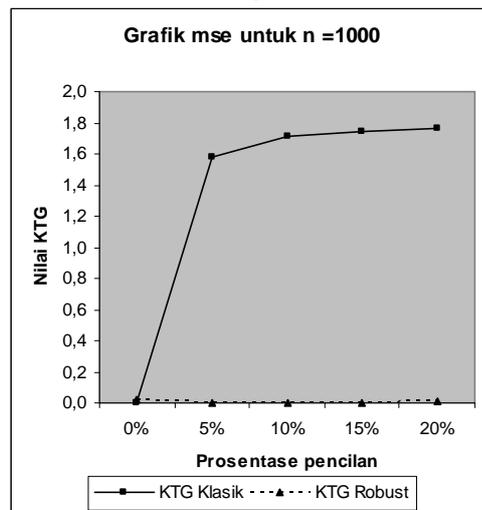
Gambar 3. Grafik bias pada n=500



Gambar 4. Grafik KTG pada n=100



Gambar 5. Grafik KTG pada n=300



Gambar 6. Grafik KTG pada n=1000

#### 4. KESIMPULAN

Analisis Regresi Komponen Utama terhadap data mengandung pencilan memerlukan metode *robust* agar menghasilkan dugaan yang tak bias dan stabil. Berdasarkan hasil simulasi data diperoleh bahwa metode VEM-KTT menghasilkan dugaan RKU yang sangat baik dan stabil untuk data dengan pencilan 5% - 20 %.

#### DAFTAR PUSTAKA

1. Myers, R.H. 1990. *Classical and Modern Regression with Applications*. PWS-KENT Publishing Company, Boston.
2. Filzmoser, P. 1999. Robust Principal Component and Factor Analysis in the Geostatistical Treatment of Environmental Data. *Environmetrics*, **10**: 363-375.
3. Croux, C., and Ruiz-Gazen, A. 2005. High breakdown estimators for principal components: the projection-pursuit approach revisited. *Journal of Multivariate Analysis*, **95**: 206-226.
4. Croux, C., Haesbroeck, G. 2000. Principal component analysis based on robust estimators of the covariance or correlation matrix: influence functions and efficiencies, *Biometrika*, **87** : 603-618
5. Hubert, M., Rousseeuw, P. J., Vanden Branden, K. 2005. ROBPCA: a new approach to robust principal components analysis. *Technometrics*, **47**:64–79.
6. Croux. C , P. Filzmoser, G. Pison & P.J. Rousseeuw. 2004. Fitting Multiplicative Models by Robust Alternating Regressions. *Technical Report* No. 350.
7. Rousseeuw, P.J., Van Aelst. S., Van Driessen, K., Agullio, J. 2004. Robust multivariate regression. *Technometrics* **46**: 293-305.
8. Rousseeuw, P.J. 1984. Least median of squares regression. *Journal of the American Statistical Association*. **79** (388): 871-880.
9. Frome, E. 2003. Least Absolute Values (LAV) Regression. [http://www.epm.ornl.gov/~frome/E\\_L\\_Frome LAV Regression.html](http://www.epm.ornl.gov/~frome/E_L_Frome_LAV_Regression.html)
10. Hubert, M., Rousseeuw, P.J., Van Aelst, S. 2004. Robustness, *Encyclopedia of Actuarial Sciences*, edited by Sundt, B. and Teugels, J., Wiley, New York, pp. 1515-1529.
11. Rousseeuw, R. J. & A. M. Leroy. 1987. *Robust Regression and Outlier Detection*. New York: Wiley.
12. Filzmoser P. 2005. Identification of Multivariate Outliers: A Performance Study. *Austrian Journal Of Statistics*. **34** (2): 127–138
13. Nisa, K. 2006. Analisis Regresi Robust Menggunakan Metode *Least Trimmed Square* untuk Data Mengandung Pencilan. *Jurnal Ilmiah MIPA*. **IX** (2) : 93-100.
14. Chen, C. 2002. Robust Regression and Outlier Detection with the ROBUSTREG procedure. *SUGI Paper* 265-27. SAS Institute: Cary, NC. SAS OnLineDoc. SAS Institute, Cary, NC: IML Robust Regression, <http://v8doc.sas.com/sashtml>
15. Rousseeuw, P.J., Van Zomeren, B.C. 1990. Unmasking multivariate outliers and leverage points, *Journal of the American Statistical Association*, **85**:633-639.
16. Nisa, K., Herawati, N., Setiawan, E., Nusyirwan. 2006. Robust Principal Component Analysis Using Minimum Covariance Determinant Estimator. *Proceedings of International Conference on Mathematics and Natural Sciences*. pp 789-792.