

BAB III

ANALISIS KORELASI KANONIK ROBUST DENGAN METODE *MINIMUM COVARIANCE DETERMINAN*

3.1 Deteksi Pencilan Multivariat

Pengidentifikasi pencilan pada kasus multivariat tidaklah mudah untuk dilakukan, hal ini dikarenakan adanya efek *masking* dan *swamping*. *Masking* terjadi pada saat pengamatan pencilan tidak terdeteksi karena adanya pengamatan pencilan lain yang berdekatan, sedangkan *swamping* terjadi pada saat pengamatan baik teridentifikasi sebagai pengamatan pencilan. Ada beberapa cara untuk mendeteksi pencilan dalam data multivariat, yaitu melalui jarak Mahalanobis atau melalui jarak *Robust*.

3.1.1 Jarak Mahalanobis

Pemeriksaan pencilan dalam kasus data multivariat dilakukan berdasarkan jarak kuadrat Mahalanobis dari x_i ke μ , yang didefinisikan sebagai berikut:

$$MD_i^2 = (x_i - \mu)^t \Sigma^{-1} (x_i - \mu), \quad i = 1, 2, \dots, n$$

di mana :

μ =vektor rata-rata

Σ =matriks kovarian

n =banyaknya pengamatan.

Sebuah pengamatan x_i diidentifikasi sebagai pencilan jika jarak Mahalanobisnya lebih besar dari nilai khi-kuadrat tabel dengan taraf signifikansi

$(1 - \alpha)$ dan derajat bebas p ($MD_i^2 > \chi_{p,(1-\alpha)}^2$ atau setara dengan $MD_i > \sqrt{\chi_{p,(1-\alpha)}^2}$).

Penggunaan jarak Mahalanobis dalam mendeteksi pencilan dapat dikatakan kurang baik. Hal ini dikarenakan adanya efek *masking* dan *swamping*. Efek *masking* dapat menurunkan jarak Mahalanobisnya sehingga jarak antar titik terpencil saling berdekatan, sedangkan efek *swamping* dapat meningkatkan jarak Mahalanobisnya sehingga dimungkinkan terjadinya kesalahan berupa pengamatan baik yang teridentifikasi sebagai pengamatan pencilan.

3.1.2 Jarak *Robust*

Jarak *robust* merupakan suatu pendekatan untuk mengidentifikasi pencilan pada data multivariat, yaitu dengan menggunakan penaksir dari μ dan Σ pada metode *robust*. Sehingga metode ini mampu meminimumkan pengaruh dari adanya efek *masking* dan *swamping* dalam pendeteksian pencilan (Rencher, 2002).

Pemeriksaan pencilan dalam kasus data multivariat dilakukan berdasarkan jarak kuadrat *robust* dari x_i ke T , yang didefinisikan sebagai berikut:

$$RD_i^2 = (x_i - T)^t C^{-1} (x_i - T), \quad i = 1, 2, \dots, n$$

di mana :

T =vektor rata-rata dari estimasi *robust*

C =matriks kovarian dari estimasi *robust*

n =banyaknya pengamatan.

Sebuah pengamatan x_i diidentifikasi sebagai pencilan jika jarak *robust*nya lebih besar dari nilai khi-kuadrat tabel dengan taraf signifikansi $(1 - \alpha)$ dan derajat bebas p ($RD_i^2 > \chi_{p,(1-\alpha)}^2$ atau setara dengan $RD_i > \sqrt{\chi_{p,(1-\alpha)}^2}$).

3.2 Analisis Korelasi Kanonik

3.2.1 Pengertian dan Tujuan Analisis Korelasi Kanonik

Analisis korelasi kanonik adalah salah satu teknik analisis statistik, yang digunakan untuk melihat hubungan antara himpunan variabel bebas (X_1, X_2, \dots, X_p) dengan himpunan variabel terikat (Y_1, Y_2, \dots, Y_q) . Analisis ini dapat mengukur tingkat keeratan hubungan antara himpunan variabel bebas dengan himpunan variabel terikat. Analisis korelasi kanonik berfokus pada korelasi antara kombinasi linier dari himpunan variabel terikat dengan kombinasi linier dari himpunan variabel bebas. Ide utama dari analisis ini adalah mencari pasangan dari kombinasi linier ini yang memiliki korelasi terbesar.

3.2.2 Asumsi-asumsi dalam Analisis Korelasi Kanonik

Adapun asumsi-asumsi dalam analisis korelasi kanonik adalah sebagai berikut :

1. Banyaknya Variabel

Variabel bebas dan variabel terikat terdiri dari lebih dari satu variabel dan berskala interval. Jika data berskala ordinal, maka data tersebut harus ditransformasi terlebih dahulu ke skala interval.

2. Uji Multikolinieritas

Menurut Hair (1998), multikolinieritas terjadi ketika dua atau lebih variabel memiliki nilai korelasi yang tinggi. Pengujian multikolinieritas dapat dilakukan dengan melihat besarnya nilai korelasi antar variabel bebasnya dan antar variabel terikatnya. Rumus dari nilai korelasi yang digunakan adalah nilai korelasi r Pearson.

Menurut Hocking (2003) pengujian multikolinieritas dengan menggunakan nilai korelasi antar variabel bebasnya atau variabel terikatnya dapat menggunakan kriteria berikut :

Jika nilai $r_{ij} > 0,95$ maka terdapat kolinieritas yang tinggi

Jika dalam suatu data terdapat kolinieritas yang tinggi, maka menurut Nachrowi (2008) salah cara untuk mengatasinya adalah dengan tidak mengikutsertakan salah satu variabel yang kolinier.

3. Uji Normalitas

Johnson dan Winchern (2007) mengemukakan bahwa untuk menguji apakah suatu himpunan data berdistribusi normal multivariat adalah dengan menggunakan Q-Q plot yang didasarkan pada jarak kuadrat Mahalanobis. Adapun langkah-langkah untuk membuat Q-Q plot adalah sebagai berikut:

- 1) Hitung nilai d_j^2 di mana

$$d_j^2 = (x_j - \bar{x})^t S^{-1} (x_j - \bar{x}) \quad j = 1, 2, \dots, n$$

n = banyaknya pengamatan.

- 2) Urutkan jarak kuadrat Mahalanobis tersebut dari yang terkecil sampai terbesar

$$d_{(1)}^2 \leq d_{(2)}^2 \leq \dots \leq d_{(n)}^2.$$

- 3) Setiap $d_{(j)}^2$ dihitung nilai $\frac{\binom{j-1}{2}}{n}$, di mana (j) adalah indeks bawah yang menunjukkan peringkat ke- j .
- 4) Hitung nilai $q_{j,p} \left(\frac{\binom{j-1}{2}}{n} \right)$, yaitu nilai khi-kuadrat dari $\frac{\binom{j-1}{2}}{n}$ dengan derajat bebas p , di mana p adalah banyaknya variabel
- 5) Gambar plot tersebut dengan koordinat $\left(q_{j,p} \left(\frac{\binom{j-1}{2}}{n} \right), d_{(j)}^2 \right)$

Data dikatakan berdistribusi normal jika plot membentuk garis lurus (linier) atau paling tidak 50% dari nilai d_j^2 lebih kecil dari $\chi_{p,0.05}^2$ ($d_j^2 \leq \chi_{p,0.05}^2$) (Anderson, 1999).

4. Uji Linieritas

Linearitas dapat dikatakan penting untuk analisis korelasi kanonik dan itu mempengaruhi dua aspek hasil korelasi kanonik. Pertama, koefisien korelasi kanonik antara sepasang variabel kanonik adalah berdasarkan hubungan linier. Jika variabel kanonik berhubungan secara nonlinier, maka koefisien korelasi kanonik tidak akan menangkap hubungan tersebut. Kedua, analisis korelasi kanonik memaksimalkan hubungan linier antara variabel kanonik. Jadi, meskipun analisis korelasi kanonik adalah metode multivariat yang paling umum, masih dibatasi untuk mengidentifikasi hubungan linier. Jika hubungan tidak linier, maka satu atau kedua variabel kanonik harus diubah, itupun jika memungkinkan (Hair, 1988).

Pengujian linieritas dilakukan antara sepasang variabel kanoniknya dan dapat dilihat dari nilai korelasi kanoniknya. Jika nilai tersebut tergolong signifikan secara statistik maka dapat dipastikan bahwa asumsi linieritas telah dipenuhi untuk pasangan variabel kanonik tersebut.

3.2.3 Penentuan Korelasi Kanonik dan Koefisien Variabel Kanonik

Analisis korelasi kanonik adalah suatu teknik yang digunakan untuk menentukan tingkatan asosiasi linier antara dua himpunan variabel, di mana masing-masing himpunan terdiri dari beberapa variabel. Kelompok pertama dari p variabel diwakili oleh $(p \times 1)$ vektor acak X . Kelompok kedua dari q variabel diwakili oleh $(q \times 1)$ vektor acak Y . Asumsikan, dalam pengembangan teoritis, bahwa X mewakili himpunan yang lebih kecil, sehingga $p \leq q$.

Misalkan untuk vektor acak X dan Y :

$$\begin{aligned} E(X) &= \mu_X ; & Cov(X) &= \Sigma_{XX} \\ E(Y) &= \mu_Y ; & Cov(Y) &= \Sigma_{YY} \\ & & Cov(X, Y) &= \Sigma_{XY} = \Sigma_{YX}. \end{aligned}$$

Vektor acaknya :

$$W = [X_1, X_2, \dots, X_p, Y_1, Y_2, \dots, Y_q]^t,$$

vektor rata-ratanya :

$$\mu = \begin{bmatrix} E(X) \\ E(Y) \end{bmatrix} = \begin{bmatrix} \mu_X \\ \mu_Y \end{bmatrix}$$

dan matriks kovariannya:

$$\Sigma = E(W - \mu)(W - \mu)^t = \begin{bmatrix} \Sigma_{XX} & \Sigma_{YX} \\ \Sigma_{XY} & \Sigma_{YY} \end{bmatrix}$$

Tugas pokok dari analisis korelasi kanonik adalah meringkaskan kumpulan antara himpunan X dan Y . Kombinasi linear menyediakan ringkasan sederhana mengukur suatu himpunan dari variabel. Himpunan $U = \alpha^t X$ dan $V = \beta^t Y$ dengan α dan β merupakan koefisien kanonik dan

$$Var(U) = \alpha^t \Sigma_{XX} \alpha$$

$$Var(V) = \beta^t \Sigma_{YY} \beta$$

$$Cov(U, V) = \alpha^t \Sigma_{XY} \beta.$$

Kemudian dapat dicari koefisien vektor α dan β sedemikian sehingga,

$$Corr(U, V) = \frac{\alpha^t \Sigma_{XY} \beta}{\sqrt{(\alpha^t \Sigma_{XX} \alpha)(\beta^t \Sigma_{YY} \beta)}} \quad (3.1)$$

sebisa mungkin bernilai besar. Kombinasi linier yang dapat dibentuk adalah sebanyak min (p, q) pasang.

Teorema Korelasi Kanonik

Misalkan $p \leq q$ dan vektor acak X dan Y mempunyai, $Cov(X) = \Sigma_{XX}$, $Cov(Y) = \Sigma_{YY}$ dan $Cov(X, Y) = \Sigma_{XY} = \Sigma_{YX}$ di mana Σ mempunyai rank lengkap. Untuk koefisien vektor α dan β , bentuk kombinasi linear $U = \alpha^t X$ dan $V = \beta^t Y$. Maka

$$\max_{\alpha, \beta} Corr(U, V) = \rho_1$$

diperoleh dengan kombinasi linear (variabel kanonik bagian pertama).

$$U_1 = e_1^t \Sigma_{XX}^{-\frac{1}{2}} X \quad \text{dan} \quad V_1 = f_1^t \Sigma_{YY}^{-\frac{1}{2}} Y.$$

Bagian ke- k dari variabel kanonik, $k = 2, 3, \dots, p$,

$$U_k = e_k^t \Sigma_{XX}^{-\frac{1}{2}} X \quad \text{dan} \quad V_k = f_k^t \Sigma_{YY}^{-\frac{1}{2}} Y$$

memaksimumkan

$$\text{Corr}(\mathbf{U}_k, \mathbf{V}_k) = \rho_k^*$$

di antara kombinasi linear yang tidak berkorelasi dengan variabel kanonik

sebelumnya. $\rho_1^{*2} \geq \rho_2^{*2} \geq \dots \geq \rho_p^{*2}$ adalah nilai eigen dari $\Sigma_{XX}^{-\frac{1}{2}} \Sigma_{XY} \Sigma_{YY}^{-1} \Sigma_{YX} \Sigma_{XX}^{-\frac{1}{2}}$

dan $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p$ adalah vektor eigen ($p \times 1$). (banyaknya $\rho_1^{*2}, \rho_2^{*2}, \dots, \rho_p^{*2}$

juga nilai eigen p paling besar dari matriks $\Sigma_{YY}^{-\frac{1}{2}} \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY} \Sigma_{YY}^{-\frac{1}{2}}$ yang

berkorespondensi dengan ($q \times 1$) vektor eigen $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_p$. Tiap \mathbf{f}_i adalah

proporsi untuk $\Sigma_{YY}^{-\frac{1}{2}} \Sigma_{YX} \Sigma_{XX}^{-\frac{1}{2}} \mathbf{e}_1$.

Variabel kanonik mempunyai sifat sebagai berikut:

$$\text{Var}(\mathbf{U}_k) = \text{Var}(\mathbf{V}_k) = 1$$

$$\text{Cov}(\mathbf{U}_k, \mathbf{U}_l) = \text{Corr}(\mathbf{U}_k, \mathbf{U}_l) = 0 \quad k \neq l$$

$$\text{Cov}(\mathbf{V}_k, \mathbf{V}_l) = \text{Corr}(\mathbf{V}_k, \mathbf{V}_l) = 0 \quad k \neq l$$

$$\text{Cov}(\mathbf{U}_k, \mathbf{V}_l) = \text{Corr}(\mathbf{U}_k, \mathbf{V}_l) = 0 \quad k \neq l$$

untuk $k, l = 1, 2, \dots, p$

Dengan bahasa yang lebih sederhana Dillon dan Goldstein dalam Kumaat (2001) mengemukakan bahwa, untuk mendapatkan fungsi kanonik, langkah-langkah yang harus ditempuh adalah :

1. Untuk memperoleh koefisien korelasi kanonik langkah-langkahnya adalah dengan menyusun matriks kovarian (\mathbf{S}) atau matriks korelasi (\mathbf{R}). Matriks \mathbf{S}

dipakai apabila data yang diolah memiliki satuan yang sama, sedangkan matriks \mathbf{R} dipakai bila data tersebut tidak memiliki satuan yang sama.

2. Mencari nilai eigen (λ) berdasarkan matriks \mathbf{S} atau \mathbf{R} pada langkah 1 dengan menggunakan rumus :

$$\left| \mathbf{R}_{XX}^{-\frac{1}{2}} \mathbf{R}_{XY} \mathbf{R}_{YY}^{-1} \mathbf{R}_{YX} \mathbf{R}_{XX}^{-\frac{1}{2}} - \lambda \mathbf{I} \right| = 0$$

Nilai eigen tersebut digunakan untuk memperoleh vektor eigen, di mana vektor eigen merupakan koefisien variabel kanonik

3. Mencari vektor eigen berdasarkan nilai eigen yang telah diperoleh pada langkah 2 dengan persamaan berikut :

$$\left(\mathbf{R}_{XX}^{-\frac{1}{2}} \mathbf{R}_{XY} \mathbf{R}_{YY}^{-1} \mathbf{R}_{YX} \mathbf{R}_{XX}^{-\frac{1}{2}} - \lambda \mathbf{I} \right) \mathbf{a} = 0$$

$$\left(\mathbf{R}_{YY}^{-\frac{1}{2}} \mathbf{R}_{YX} \mathbf{R}_{XX}^{-1} \mathbf{R}_{XY} \mathbf{R}_{YY}^{-\frac{1}{2}} - \lambda \mathbf{I} \right) \mathbf{b} = 0$$

Vektor eigen tersebut dinotasikan dengan \mathbf{a} dan \mathbf{b} , yang merupakan nilai koefisien variabel kanonik atau disebut juga sebagai pembobot kanonik. Variabel kanonik yang dapat dibentuk berdasarkan vektor eigen tersebut ada sebanyak min (p, q) pasang

4. Setelah memperoleh vektor eigen, selanjutnya dicari korelasi kanonik yang dapat dihitung dengan menggunakan rumus :

$$\text{Corr}(\mathbf{U}, \mathbf{V}) = \frac{\mathbf{a}^t \mathbf{R}_{XY} \mathbf{b}}{\sqrt{(\mathbf{a}^t \mathbf{R}_{XX} \mathbf{a})(\mathbf{b}^t \mathbf{R}_{YY} \mathbf{b})}}$$

5. Mencari proporsi atau keragaman data yang dijelaskan oleh setiap pasangan variabel kanonik, dengan menggunakan rumus :

$$e_i = \frac{\lambda_i}{1 - \lambda_i}, i = 1, 2, \dots, \min(p, q)$$

$$\text{Keragaman data} = \frac{e_i}{\sum e_i}$$

Keragaman data ini digunakan untuk memilih pasangan variabel kanonik mana yang akan dianalisis lebih lanjut. Batasan minimum keragaman kumulatif yang dikemukakan oleh Dillon dan Goldstein (1984) adalah 80%.

6. Melakukan pengujian hipotesis untuk setiap korelasi kanonik. Berdasarkan Johnson dan Winchern (2007), pengujian korelasi kanonik secara individu dilakukan melalui pendekatan distribusi khi-kuadrat dengan rumusan hipotesisnya adalah :

$H_0 : \rho_i = 0$, artinya tidak ada hubungan yang signifikan antar pasangan variabel kanonik ke-i

$H_1 : \rho_i \neq 0$, artinya ada hubungan yang signifikan antar pasangan variabel kanonik ke-i

Kriteria yang digunakan adalah tolak H_0 pada tingkat signifikansi α , jika

$$-\left(n - 1 - \frac{1}{2}(p + q + 1)\right) \ln \prod_{k=i}^p (1 - \rho_k^2) > \chi_{(p-i-1)(q-i-1)}^2(\alpha).$$

Guna memudahkan dalam pencarian korelasi kanonik, berikut algoritma untuk analisis korelasi kanonik klasik dan analisis korelasi kanonik *robust*.

Analisis korelasi kanonik klasik :

1. Uji Asumsi
2. Deteksi Pencilan
3. Menentukan matriks kovarians
4. Menentukan nilai eigen dan vektor eigen

5. Menentukan korelasi kanonik dan pembobot kanonik
6. Menentukan proporsi keragaman

Analisis korelasi kanonik *robust*:

1. Uji Asumsi
2. Deteksi Pencilan
3. Menentukan matriks kovarians dengan metode MCD
4. Menentukan nilai eigen dan vektor eigen
5. Menentukan korelasi kanonik dan pembobot kanonik
6. Menentukan proporsi keragaman

3.2.4 *Canonical Loadings dan Cross Loadings*

Canonical loadings merupakan korelasi sederhana antara variabel asal dengan masing-masing variabel kanoniknya. Semakin besar nilai *canonical loadings* menunjukkan semakin dekat hubungan antara variabel asal dengan variabel kanoniknya. Menurut Hair (1998) *canonical loadings* variabel terikat diperoleh dengan rumus

$$R_{YV} = R_{YY}b$$

R_{YY} merupakan korelasi sederhana antar variabel Y dan b merupakan vektor koefisien kanonik variabel V . Sedangkan *canonical loadings* untuk variabel bebas diperoleh dengan rumus

$$R_{XU} = R_{XX}a$$

R_{XX} merupakan korelasi sederhana antar variabel X dan a merupakan vektor koefisien kanonik variabel U .

Canonical Cross loadings merupakan korelasi sederhana antara variabel asal dengan masing-masing variabel kanonik lawannya. Semakin besar nilai *canonical cross loadings* menunjukkan semakin kuat hubungan variabel asal dengan variabel kanonik lawannya. *Canonical Cross loadings* diperoleh dengan cara

$$R_{XV} = R_{XU}\rho_k$$

$$R_{YU} = R_{YV}\rho_k$$

ρ_k adalah nilai korelasi kanonik dari variabel kanonik ke- k .

3.2.5 Redudansi

Redudansi merupakan sebuah indeks yang menghitung proporsi keragaman yang dapat dijelaskan oleh variabel kanonik yang dipilih baik variabel kanonik terikat maupun variabel kanonik bebas.

Proporsi keragaman variabel asal yang diterangkan oleh variabel kanoniknya diperoleh dari perhitungan rata-rata *canonical loadings* yang dikuadratkan.

$$R_{X|U_i}^2 = \frac{\sum_{j=1}^k R_{X_j U_i}^2}{k}$$

$$R_{Y|V_i}^2 = \frac{\sum_{j=1}^k R_{Y_j V_i}^2}{k}$$

Proporsi keragaman variabel asal yang diterangkan oleh variabel kanonik lawannya diperoleh melalui perkalian kuadrat korelasi kanonik dengan rata-rata *canonical loadings* yang dikuadratkan, atau dapat dituliskan

$$R_{X|V_i}^2 = \rho_k^2 R_{X|U_i}^2$$

$$R_{Y|U_i}^2 = \rho_k^2 R_{Y|V_i}^2$$

ρ_k adalah nilai korelasi kanonik dari variabel kanonik ke- k .

3.3 MINIMUM COVARIANCE DETERMINANT

3.3.1 Definisi MCD

Misalkan $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n)^t$ merupakan kumpulan data sejumlah n pengamatan terdiri dari p -variabel di mana $n \geq p + 1$. Penaksir MCD merupakan pasangan $\mathbf{T} \in \mathbb{R}^p$ dan \mathbf{C} adalah matriks definit positif simetris berdimensi $p \times p$ dari suatu sub sampel berukuran h pengamatan di mana $\frac{(n+p+1)}{2} \leq h \leq n$ dengan

$$\mathbf{T} = \frac{1}{h} \sum_{i=1}^h \mathbf{w}_i$$

dan

$$\mathbf{C} = \frac{1}{h} \sum_{i=1}^h (\mathbf{w}_i - \mathbf{T}_i)(\mathbf{w}_i - \mathbf{T}_i)^t$$

yang meminimumkan determinan \mathbf{C} (Buttler, dkk. 1993).

Dalam menentukan penaksir MCD, jika jumlah pengamatan (n kecil) maka penaksir MCD dapat segera ditentukan. Namun jika jumlah pengamatan (n besar) maka akan membutuhkan waktu yang cukup lama untuk menentukan penaksir MCD. Karena keterbatasan ini maka Rousseeuw dan Van Driessen (1999) mengembangkan suatu algoritma FAST-MCD yaitu Teorema C-Steps berikut.

Teorema C-Steps:

Misalkan $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n)^t$ merupakan himpunan sejumlah n pengamatan terdiri dari p -variabel. Misal $H_1 \subset \{1, 2, \dots, n\}$ dengan sejumlah elemen H_1 , jumlah $|H_1| = h$, tetapkan

$$\mathbf{T}_1 = \frac{1}{h} \sum_{i=1}^h \mathbf{w}_i$$

dan

$$\mathbf{C}_1 = \frac{1}{h} \sum_{i=1}^h (\mathbf{w}_i - \mathbf{T}_1)(\mathbf{w}_i - \mathbf{T}_1)^t$$

Jika $\det(\mathbf{C}_1) \neq 0$ maka jarak relatif,

$$d_1(i) = \sqrt{(\mathbf{w}_i - \mathbf{T}_1)^t \mathbf{C}_1^{-1} (\mathbf{w}_i - \mathbf{T}_1)}, i = 1, 2, \dots, n.$$

Selanjutnya ambil H_2 sedemikian sehingga $\{d_1(i); i \in H_2\} := \{(d_1)_{1:n}, (d_1)_{2:n}, \dots, (d_1)_{h:n}\}$, di mana $(d_1)_{1:n} \leq (d_1)_{2:n} \leq \dots \leq (d_1)_{n:n}$ menyatakan urutan jarak, dan hitung \mathbf{T}_2 dan \mathbf{C}_2 berdasarkan H_2 , maka

$$\det(\mathbf{C}_2) \leq \det(\mathbf{C}_1)$$

dan akan sama jika dan hanya jika $\mathbf{T}_1 = \mathbf{T}_2$ dan $\mathbf{C}_1 = \mathbf{C}_2$ (Driessen, 1999).

Untuk pengamatan $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n\} \in \mathbb{R}^p$, MCD ditentukan dengan memilih subset dari observasi $\{\mathbf{w}_{i_1}, \mathbf{w}_{i_2}, \dots, \mathbf{w}_{i_h}\}$ dari ukuran h , dengan $1 \leq h \leq n$ di mana memiliki determinan matriks kovarian yang minimum.

3.3.2 *Affine Equivariance*

Affine Equivariance memiliki implikasi bahwa suatu estimator dapat bertransformasi dengan baik dalam suatu transformasi linier nonsingular. Sehingga meskipun data dirotasi atau ditranslasi tidak akan memiliki pengaruh pada pendeteksian pencilan (Hubert, 2009). Transformasi linier yang dimaksud adalah pada analisis diskriminan, analisis faktor, analisis korelasi kanonik (Yohai, 2006).

Estimator MCD dari rata-rata dan kovarian merupakan *affine equivariance*. Maksud dari *affine equivariance* adalah untuk suatu matriks nonsingular A dan vektor konstan $\mathbf{b} \in \mathbb{R}^d$

$$\mathbf{T}_{MCD}(\mathbf{AX} + \mathbf{b}) = \mathbf{AT}_{MCD}(\mathbf{X}) + \mathbf{b}$$

$$\mathbf{C}_{MCD}(\mathbf{AX} + \mathbf{b}) = \mathbf{AC}_{MCD}(\mathbf{X})\mathbf{A}^t.$$

3.3.3 *Breakdown Point*

Breakdown point adalah jumlah pengamatan minimal yang dapat menggantikan sejumlah pengamatan awal yang berakibat pada nilai taksiran yang dihasilkan sangat berbeda dari taksiran sebenarnya (Lopuhaa dan Rousseeuw, 1991). *Breakdown point* juga merupakan alat untuk mengukur kerobustan dari suatu penaksir.