

BAB III
MINIMUM VOLUME ELLIPSOID
PADA ANALISIS KOMPONEN UTAMA ROBUST

Pada bab ini akan dikaji bahasan utama yaitu pencilan dan analisis komponen utama robust sebagai konsep pendukung serta metode *Minimum Volume Ellipsoid* sebagai konsep utama pada tugas akhir ini.

3.1. Analisis Komponen Utama Robust

Analisis komponen utama merupakan metode yang baik untuk menerangkan variasi total dari variabel-variabel yang mewakili suatu data melalui sebagian kecil komponen utama. Namun ketika terdapat pencilan dalam data, analisis komponen utama klasik yang diperoleh dapat menyimpang dari hasil yang diharapkan.

3.1.1. Pencilan

Dalam Hampel et al (1986: 1), pencilan didefinisikan sebagai observasi yang terletak sangat jauh dari bagian terbesar data, dan membahayakan beberapa prosedur statistik klasik. Dengan kata lain, pencilan merupakan penyimpangan beberapa observasi dalam suatu data yang terletak cukup jauh dari pusat observasi. Adapun jenis-jenis pencilan yaitu:

1. *Good leverage* merupakan observasi yang berada di ruang distribusi tetapi sudah tidak berada di daerah mayoritas data.

2. *Bad leverage* merupakan observasi yang tidak berada baik dalam ruang distribusi observasi maupun daerah mayoritas data.
3. Pencilan ortogonal merupakan observasi yang mempunyai jarak yang sangat besar dari daerah mayoritas data sehingga observasi tersebut sudah tidak dapat dilihat dalam ruang distribusinya.

Pendeteksian pencilan sangat penting dilakukan dalam prosedur statistik karena pencilan dapat mempengaruhi informasi yang terdapat di dalam data. Sebagai contoh, suatu data berat badan (dalam kilogram) dari lima orang siswa dalam kelompok A masing-masing 43, 39, 39, 38, dan 41, mean dari data tersebut adalah 40 kg. Jika kemudian ada seorang siswa baru masuk ke dalam kelompok tersebut dan dia memiliki berat badan 100 kg, maka mean berat badan dalam kelompok A menjadi 50 kg. Hasil mean akhir tersebut tidak dapat mewakili informasi sebenarnya, karena mayoritas siswa dalam kelompok A memiliki berat badan di bawah 50 kg. Hal ini menunjukkan bahwa pencilan dapat mempengaruhi keseimbangan data.

Pencilan yang terdeteksi dari suatu sampel mungkin saja dihilangkan dari data observasi agar dapat dilakukan analisis lebih lanjut. Tindakan tersebut dapat dilakukan apabila hanya terdapat satu buah pencilan, namun hal ini tidak mungkin dilakukan jika pencilan yang terdeteksi lebih dari satu, karena pencilan dapat mengandung informasi yang penting dari suatu data. Oleh karena itu, dalam mengatasi pencilan ini dibutuhkan suatu metode penaksir robust yang tangguh terhadap pencilan sehingga analisis komponen utama yang dilakukan tidak lagi dipengaruhi pencilan.

3.1.2. Jarak Mahalanobis

Umumnya, metode yang digunakan untuk mendeteksi pencilan dalam analisis multivariat berdasarkan pada jarak Mahalanobis. Bentuk dari persamaan jarak Mahalanobis adalah

$$d_i = \{(\mathbf{x}_i - \mathbf{a})\mathbf{S}^{*-1}(\mathbf{x}_i - \mathbf{a})^t\}^{1/2} \quad (3.1)$$

dengan \mathbf{x}_i adalah vektor observasi, \mathbf{a} adalah estimasi mean, dan \mathbf{S}^* adalah estimasi dari matriks varians kovarians. Bentuk persamaan ini sama dengan T^2 Hotelling untuk vektor observasi individual, perbedaannya yaitu bahwa \mathbf{a} dan \mathbf{S}^* bukan lagi penaksir standar tetapi merupakan bagian dari proses robust. Kebanyakan teknik yang dipakai pada proses robust merupakan iterasi sehingga \mathbf{a} dan \mathbf{S}^* dapat berubah pada setiap iterasi. Jarak Mahalanobis ini dipakai untuk memberi indeks $w(d_i)$ yang digunakan untuk menghasilkan \mathbf{a} dan \mathbf{S}^*

$$\mathbf{a} = \sum_{i=1}^n w_a(d_i) \frac{\mathbf{x}_i}{\sum_{i=1}^n w_a(d_i)} \quad (3.2)$$

$$\mathbf{S}^* = \frac{\sum_{i=1}^n w_s(d_i) [\mathbf{x}_i - \mathbf{a}]^t [\mathbf{x}_i - \mathbf{a}]}{f\{w_s(d_i)\}} \quad (3.3)$$

Definisi 3.2.1 Jarak Mahalanobis (Chen Y, Chen X, Xu, 2008: 222)

Untuk matriks data $\mathbf{X}_{n \times p}$ dengan n observasi dan p variabel, jarak Mahalanobis dihitung sebagai

$$MD^2(\mathbf{x}_i, \mathbf{X}) = (\mathbf{x}_i - \mathbf{T}(\mathbf{X}))\mathbf{S}(\mathbf{X})^{-1}(\mathbf{x}_i - \mathbf{T}(\mathbf{X}))^t$$

untuk masing-masing observasi, dengan $\mathbf{T}(\mathbf{X})$ adalah penaksir lokasi multivariat (dalam kasus ini merupakan *arithmetic mean*) dan $\mathbf{S}(\mathbf{X})$ adalah penaksir varians kovarians klasik. Titik-titik dengan $MD^2(\mathbf{x}_i, \mathbf{X})$ yang besar diidentifikasi sebagai

titik-titik tak beraturan (pencilan) dan ditaksir menggunakan distribusi χ^2 dengan derajat kebebasan yang sesuai.

3.1.3. Estimasi Robust dengan Penaksir *Affine Equivariant*

Dalam Rousseeuw dan Leroy (1987: 248), sebagai penaksir lokasi multivariat, persamaan untuk T yang merupakan *translation equivariant* dapat dituliskan sebagai

$$T(\mathbf{x}_1 + \mathbf{b}, \dots, \mathbf{x}_n + \mathbf{b}) = T(\mathbf{x}_1, \dots, \mathbf{x}_n) \quad (3.4)$$

untuk beberapa vektor \mathbf{b} berdimensi p . Persamaan tersebut juga disebut sebagai *location equivariance*. Sedangkan penaksir lokasi multivariat adalah mean aritmetik

$$T(\mathbf{X}) = \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad (3.5)$$

yang merupakan penaksir *least square* dalam kerangka ini karena persamaan tersebut akan meminimalkan $\sum_{i=1}^n \|\mathbf{x}_i - \mathbf{T}\|^2$ dengan $\|\dots\|$ adalah norm Euclid.

T disebut *affine equivariant* jika dan hanya jika

$$T(\mathbf{x}_1 \mathbf{A} + \mathbf{b}, \dots, \mathbf{x}_n \mathbf{A} + \mathbf{b}) = T(\mathbf{x}_1, \dots, \mathbf{x}_n) \mathbf{A} + \mathbf{b} \quad (3.6)$$

untuk semua vektor \mathbf{b} dan sebarang matriks nonsingular \mathbf{A} .

Untuk penaksir varians kovarians, *affine equivariance* berarti bahwa

$$\mathbf{C}(\{\mathbf{x}_1 \mathbf{A} + \mathbf{b}, \dots, \mathbf{x}_n \mathbf{A} + \mathbf{b}\}) = \mathbf{A}^t \mathbf{C}(\{\mathbf{x}_1, \dots, \mathbf{x}_n\}) \mathbf{A} \quad (3.7)$$

dengan \mathbf{A} adalah matriks nonsingular berukuran $p \times p$ dan \mathbf{b} merupakan vektor.

Pada $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, maximum likelihood menghasilkan penaksir equivariant

$$T(\mathbf{X}) = \bar{\mathbf{x}}$$

$$\mathbf{S}(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{T}(\mathbf{X}))^t (\mathbf{x}_i - \mathbf{T}(\mathbf{X})). \quad (3.8)$$

Untuk menghasilkan penaksir tak bias dari Σ , penyebut pada persamaan $\mathbf{S}(\mathbf{X})$ dari n diganti dengan $n - 1$, dan \mathbf{x}_i menjadi vektor sampel ke- i yang dihitung sebagai

$$\mathbf{T}(\mathbf{X}) = \bar{\mathbf{x}}_i$$

$$\mathbf{S}(\mathbf{X}) = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{T}(\mathbf{X}))^t (\mathbf{x}_i - \mathbf{T}(\mathbf{X})). \quad (3.9)$$

Untuk menggeneralisasi pendekatan maksimum likelihood, digunakan langkah iteratif *affine equivariant* untuk menaksir μ dan Σ . Definisi solusi simultan dari sistem persamaan

$$\left\{ \begin{array}{l} \frac{1}{n} \sum_{i=1}^n u_1 \left([(\mathbf{x}_i - \mathbf{T}) \mathbf{S}^{-1} (\mathbf{x}_i - \mathbf{T})^t]^{\frac{1}{2}} \right) (\mathbf{x}_i - \mathbf{T}) = \mathbf{0} \\ \frac{1}{n} \sum_{i=1}^n u_2 ((\mathbf{x}_i - \mathbf{T}) \mathbf{S}^{-1} (\mathbf{x}_i - \mathbf{T})^t) (\mathbf{x}_i - \mathbf{T})^t (\mathbf{x}_i - \mathbf{T}) = \mathbf{S}. \end{array} \right. \quad (3.10)$$

Sebagai ukuran tingkat kerobustan suatu metode terhadap keberadaan pencilan digunakan *breakdown point*. *Breakdown point* untuk mean adalah 1 yang artinya dengan hanya menggantikan 1 nilai ekstrim pada data asal, maka akan didapati perubahan mean yang sangat besar (Suryana, 2008). *Breakdown point* merupakan persentase terkecil pada penyimpangan yang dapat memiliki efek besar pada penaksir. *Breakdown point* pada semua *affine equivariant* yang banyak dipakai adalah $1/(p+1)$ atau sebesar 50% *breakdown point* (Rousseeuw dan Leroy, 1987: 253).

3.2. Pendeteksian Pencilan Menggunakan Metode MVE

Penggunaan jarak Mahalanobis klasik pada data multivariat berdasarkan mean sampel dan matriks varians kovarians sampel terkadang tidak efektif ketika berhadapan dengan sekelompok pencilan, dengan kata lain data masih dipengaruhi oleh pencilan. Akibatnya, jika terdapat pencilan dalam data, pencilan tersebut dapat mempengaruhi mean sampel dan matriks varians kovarians sampel sedemikian sehingga pencilan-pencilan tersebut memiliki jarak Mahalanobis yang kecil. Karena itu, pencilan tetap tidak terdeteksi. Hal ini disebut dengan efek *masking*. Rousseeuw dan Leroy (1987: 256) mengusulkan untuk menghitung jarak robust berdasarkan penaksir lokasi dan *scatter* dengan *high breakdown point*.

3.2.1. MVE dan Resampling

MVE yang diperkenalkan oleh Rousseeuw merupakan penaksir robust dengan *high breakdown* pertama pada lokasi dan *scatter* multivariat yang sering digunakan. MVE terkenal karena kepekaannya terhadap pencilan yang membuatnya dapat diandalkan untuk mendeteksi pencilan dan tersedia secara luas, serta implementasi MVE mudah dioperasikan pada algoritma perhitungannya (Rousseeuw dan Van Aelst, 2009: 1).

Menurut Rousseeuw, Croux, dan Haesbroeck (2002: 192), MVE dapat didefinisikan sebagai elipsoid terkecil yang mencakup paling sedikit h elemen pada himpunan $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \in \mathbb{R}^p$, dengan penaksir lokasi MVE merupakan pusat elipsoid dan penaksir *scatter* MVE berkoresponden menjadi matriks bentukan (*shape matrix*).

Definisi 3.2.1.1 (Rousseeuw dan Van Aelst, 2009: 72)

Pasangan (\mathbf{T}, \mathbf{S}) dengan penaksir lokasi $\mathbf{T}(\mathbf{X})$ dan penaksir *scatter* $\mathbf{S}(\mathbf{X})$ meminimalkan determinan \mathbf{S} pada kondisi

$$\{i | (\mathbf{X}_i - \mathbf{T})\mathbf{S}^{-1}(\mathbf{X}_i - \mathbf{T})^t \leq c^2\} \geq h$$

dengan n adalah jumlah observasi, p adalah jumlah variabel, $h = [(n + p + 1)/2]$, c adalah konstanta dan \mathbf{X} adalah himpunan data observasi. Proses meminimalkan tersebut berada pada setiap $\mathbf{T} \in \mathbb{R}^p$ dan $\mathbf{S} \in PDS(p)$. $PDS(p)$ adalah himpunan matriks simetri definit positif berukuran $p \times p$.

Nilai c adalah konstanta tetap yang dipilih untuk menentukan besar dari \mathbf{S}_n . Umumnya, c dipilih sedemikian sehingga \mathbf{S}_n merupakan penaksir konsisten pada matriks varians kovarians untuk data dari distribusi normal multivariat, dengan kata lain $c = \sqrt{\chi_{p,\alpha}^2}$ dengan $\alpha = h/n$. Dari definisi ini, jelas bahwa MVE menaksir pusat dan *scatter* untuk h lebih memusatkan observasi dalam data. Nilai h dapat dipilih dan menentukan ke-robust-an untuk menghasilkan estimasi MVE.

Persamaan standar untuk nilai h yang sering dipakai dinyatakan dengan $h = \left\lceil \frac{n+p+1}{2} \right\rceil \approx \frac{n}{2}$ karena ini dapat menghasilkan *breakdown point* yang maksimal. Tetapi jika diketahui pecahan pencilon sebagian besar berada pada $0 < \alpha < \frac{1}{2}$, maka dapat digunakan estimator MVE (α) dengan $h = [n(1 - \alpha)]$. Untuk mendapatkan hasil yang baik, nilai $\alpha = \frac{1}{4}$ adalah pilihan yang tepat.

MVE memiliki *breakdown point* mendekati 50%, yang berarti bahwa $\mathbf{T}(\mathbf{X})$ dapat tetap terbatas dan besar nilai eigen dari $\mathbf{S}(\mathbf{X})$ akan tetap jauh dari nol

dan tidak terhingga ketika kurang dari separuh dari jumlah data digantikan oleh nilai-nilai sebarang (Rousseeuw dan Van Aelst, 1991: 196).

Definisi 3.2.1.2 (Rousseeuw dan Van Aelst, 2009: 73)

Jarak robust $RD(x_i)$ yang berhubungan dengan MVE didefinisikan

$$RD(x_i) = \sqrt{(x_i - T(X))S(X)^{-1}(x_i - T(X))^t}, \quad i = 1, \dots, n.$$

Karena dalam mencari solusi eksak untuk menyelesaikan masalah pada MVE agak sulit, sebagai alternatifnya akan dicari pendekatan MVE menggunakan algoritma resampling yang disebut juga algoritma $(p + 1)$ -subset. Algoritma ini memiliki keuntungan dalam penaksirannya karena masih bersifat *affine equivariant* dan masih mempertahankan *high breakdown point*. Dalam tugas akhir ini digunakan adaptasi dari penaksir $(p + 1)$ -subset untuk MVE dengan *high breakdown point*. Gambaran utama dari algoritma resampling ini yaitu dengan merata-ratakan lebih dari beberapa nilai percobaan dibandingkan dengan hanya mengambil salah satu yang optimal.

Penaksir dasar dalam penentuan lokasi dan *scatter* multivariat adalah *affine equivariance*, yang berarti bahwa penaksir menunjukkan dengan baik di bawah transformasi *affine* pada data. Sehingga, penaksir T dan S pada lokasi dan *scatter* multivariat yang *affine equivariant* berikut pada beberapa matriks data X

$$\begin{aligned} T(XA + \mathbf{1}_n v^t) &= A^t T(X) + v \\ S(XA + \mathbf{1}_n v^t) &= A^t S(X) A \end{aligned} \quad (3.11)$$

untuk semua matriks nonsingular $p \times p$ A dan $v \in \mathbb{R}^p$. Vektor $\mathbf{1}_n = (1, 1, \dots, 1)^t \in \mathbb{R}^n$.

Affine equivariance dalam penaksir cukup penting karena dapat membuat analisis

menjadi independen pada skala pengukuran variabel baik dalam translasi maupun rotasi pada data.

3.2.2. Metode MVE

Dari definisi 3.2.1.1 dapat diketahui bahwa perhitungan MVE yang tepat untuk himpunan \mathbf{X}_n akan menuntut pengujian semua $\binom{n}{h}$ elipsoid yang mengandung h observasi pada \mathbf{X}_n untuk mencari elipsoid dengan volume terkecil. Penyelesaian masalah kombinatorial ini mungkin dilakukan jika himpunan data kecil dan berdimensi rendah. Karena jumlah elipsoid umumnya sangat besar, maka untuk menyelesaikan persoalan pada data yang berdimensi tinggi dilakukan pendekatan algoritma.

Patokan algoritma MVE membatasi pencarian elipsoid yang ditentukan oleh subset yang memuat $p + 1$ (p merupakan jumlah variabel pada matriks \mathbf{X}_n) observasi yang berbeda diberikan indeks $J = \{i_1, i_2, \dots, i_{p+1}\} \subset \{1, \dots, n\}$, kemudian berdasarkan persamaan (3.10), hitung mean dan matriks varians kovarians sampel sehingga persamaan menjadi

$$\mathbf{T}_{\mathcal{J}} = \bar{\mathbf{x}}_J = \frac{1}{p+1} \sum_{i \in \mathcal{J}}^{p+1} \mathbf{x}_i \quad (3.12)$$

$$\mathbf{S}_J = \frac{1}{p} \sum_{i \in J}^{p+1} (\mathbf{x}_i - \bar{\mathbf{x}}_J)^t (\mathbf{x}_i - \bar{\mathbf{x}}_J) \quad (3.13)$$

dengan matriks varians kovarians \mathbf{S}_J non singular untuk $\mathbf{x}_{i_1}, \mathbf{x}_{i_2}, \dots, \mathbf{x}_{i_{p+1}}$. Jika $p + 1$ bukan pada posisi umum, maka observasi pada \mathbf{X}_n ditambahkan hingga subset dengan matriks varians kovarians sampel nonsingular diperoleh (atau

sebuah subsampel singular berukuran h diperoleh). Elipsoid ditentukan oleh T_J dan S_J yang kemudian dinaikkan atau diturunkan hingga mengandung tepat h titik. Sedangkan faktor penskalaan diberikan oleh persamaan D_J^2/c^2 dengan $c = \sqrt{\chi_{p,\alpha}^2}$ dan

$$D_J^2 = [(x_i - \bar{x}_J)C_J^{-1}(x_i - \bar{x}_J)^t]_{h:n} \quad (3.14)$$

dengan $h = (n + p + 1)/2$ dan $h:n$ menunjukkan jarak kuadrat terkecil ke- h di antara jarak kuadrat pada n observasi dalam X_n . Elipsoid yang dihasilkan dapat memenuhi definisi 3.4.1.1 dan volumenya proporsional dengan menggunakan persamaan

$$P_J = (\det(D_J^2 S_J))^{0.5}. \quad (3.15)$$

Ulangi langkah-langkah tersebut sebanyak subsampel J , dan pilih salah satu dengan P_J terendah. Kemudian untuk subsampel J yang telah diperoleh, hitung

$$T(X) = T_J$$

dan

$$S(X) = c^2(n, p)(\chi_{p,\alpha}^2)^{-1} D_J^2 S_J \quad (3.16)$$

dengan $c^2(n, p)$ adalah *correction term* sampel kecil yang dihitung sebagai $\left[1 + \frac{15}{n-p}\right]^2$ dan $\chi_{p,\alpha}^2$ merupakan median dari distribusi χ^2 dengan derajat kebebasan p . Hal ini menunjukkan bahwa sampling intensif dan perhitungan membutuhkan pencarian solusi dari analisis MVE. Total jumlah subsampel bergantung pada n dan p . Berdasarkan estimasi MVE mean $T(X)$ dan varians kovarians $S(X)$, pemberian indeks dapat dihitung menggunakan

$$W_i = (\mathbf{x}_i - \mathbf{T}(\mathbf{X}))\mathbf{S}(\mathbf{X})^{-1}(\mathbf{x}_i - \mathbf{T}(\mathbf{X}))^t \quad (3.17)$$

untuk observasi i , dengan $W_i > \chi_{p;\alpha}^2$ didefinisikan sebagai pencilan. PCA yang digambarkan dengan pencilan yang didefinisikan melalui MVE memiliki indeks 0 dan data normal diberikan indeks 1.

Dari uraian tersebut, algoritma pendeteksian pencilan pada metode MVE adalah:

1. Bentuk subsampel yang mengandung $p+1$ observasi disimbolkan dengan indeks J sebanyak $\binom{n}{h}$.
2. Hitung mean dan matriks varians kovarians koresponden dengan \mathbf{S}_J nonsingular.
3. Hitung $D_J^2 = [(\mathbf{x}_i - \bar{\mathbf{x}}_J)\mathbf{S}_J^{-1}(\mathbf{x}_i - \bar{\mathbf{x}}_J)^t]_{h:n}$
4. Hitung $P_J = (\det(D_J^2 \mathbf{S}_J))^{0,5}$ untuk menghasilkan elipsoid.
5. Ulangi langkah 1-4 sebanyak subsampel J , kemudian pilih P_J dengan nilai terendah.
6. Hitung $\mathbf{T}(\mathbf{X})$ dan $\mathbf{S}(\mathbf{X}) = c^2(n, p)(\chi_{p;\alpha}^2)^{-1} D_J^2 \mathbf{S}_J$ dari P_J yang akan digunakan sebagai estimasi mean dan matriks varians kovarians.
7. Observasi dengan $W_i > \chi_{p;\alpha}^2$ diidentifikasi sebagai pencilan.

3.3. *Weighted PCA*

Untuk menangani masalah pencilan, selanjutnya digunakan *Weighted Principal Component Analysis* (WPCA) yang diperkenalkan oleh Joaquim F. Pinto da Costa et al (2011) dalam jurnal “*A Weighted Principal Component Analysis and Its Application to Gene Expression Data*”. PCA klasik tidak cukup tepat untuk menangani masalah pencilan. Dalam tugas akhir ini akan digunakan koefisien korelasi baru sebagai alternatif korelasi Pearson yang disebut *Weighted PCA* (WPCA).

Pada PCA klasik, vektor-vektor eigen dari matriks varians kovarians atau matriks korelasi Pearson memuat koefisien-koefisien dari kombinasi linier dari matriks asli yang berkorespondensi dengan variabel-variabel baru (komponen-komponen utama). Seperti telah diketahui koefisien korelasi Pearson sangat sensitif terhadap kehadiran pencilan dan gangguan. Untuk menangani hal ini, akan digunakan *rank* (rangking) dari setiap observasi, yang dimulai dengan merangking observasi untuk setiap variabel dari satu (rangking tertinggi) sampai n (rangking terendah).

Misalkan diberikan data bivariat (X_i, Y_i) , $i = 1, 2, \dots, n$, beri nama R_i untuk rangking variabel X_i dan Q_i untuk rangking variabel Y_i . Jika akan menghitung koefisien korelasi Pearson untuk rangking data, maka yang diperoleh adalah koefisien korelasi rank Spearman (r_s), yaitu :

$$r_s = \frac{\sum_{i=1}^n (R_i - \bar{R})(Q_i - \bar{Q})}{\sqrt{\sum_{i=1}^n (R_i - \bar{R})^2 \sum_{i=1}^n (Q_i - \bar{Q})^2}} \quad (3.18)$$

dengan \bar{R} dan \bar{Q} adalah rata-rata *rank*. (Rachmatin, 2011: 2)

R_i dan Q_i untuk merepresentasikan *rank* dari dua variabel yang berkorespondensi dengan observasi (sampel) i . Untuk tujuan komputasi, persamaan dalam mencari koefisien korelasi *rank* adalah

$$r_s = 1 - \frac{6 \sum_{i=1}^n (R_i - Q_i)^2}{n^3 - n} \quad (3.19).$$

Definisi 3.2.3.1 (da Costa, 2011: 247)

$$W_2 D_i^2 = (R_i - Q_i)^2 (2n + 2 - R_i - Q_i)^2$$

yang mewakili lebih dari $W D_i^2$ penyesuaian yang lebih penting pada *rank-rank* atas. Biasanya koefisien korelasi *rank* didefinisikan seperti koefisien Spearman, sebagai fungsi linier dari jarak antara dua vektor *rank*.

Da Costa mengusulkan bahwa transformasi mengandung pensubstitusian nilai observasi i pada variabel pertama dengan nilai $R'_i = R_i(2n + 2 - R_i)$, dengan R_i merupakan *rank* dari observasi. Begitu pula untuk variabel lainnya, diperoleh dari $Q'_i = Q_i(2n + 2 - Q_i)$. Untuk menghitung rata-ratanya digunakan persamaan berikut

$$\begin{aligned} \overline{R'} &= \frac{1}{n} \sum_{i=1}^n R_i(2n + 2 - R_i) = \frac{(n+1)(4n+5)}{6} \\ \overline{Q'} &= \frac{1}{n} \sum_{i=1}^n Q_i(2n + 2 - Q_i) = \frac{(n+1)(4n+5)}{6} \end{aligned} \quad (3.20)$$

Algoritma WPCA yang diperkenalkan da Costa (Rachmatin, 2011) adalah sebagai berikut:

Misalkan diberikan data matriks X dengan ukuran $n \times p$ (n observasi dan p variabel).

1. Transformasi X menjadi R ($n \times p$) \rightarrow matriks *rank* untuk data.

$$\mathbf{X} = [x_{ij}] \quad i = 1, 2, \dots, n$$

$$j = 1, 2, \dots, p$$

$\mathbf{R} = [R_{ij}]$ dengan $R_{ij} = \text{rank}$ observasi ke- i kolom ke- j .

2. Transformasi \mathbf{R} menjadi \mathbf{R}' dengan $\mathbf{R}' = [R_{ij}']$

$$R_{ij}' = R_{ij}(2n + 2 - R_{ij})$$

3. Standardisasi R_{ij}' dengan transformasi

$$R_{ij}'' = \frac{R_{ij}' - \bar{R}_j'}{S_{R_j} \sqrt{n}}$$

$$\text{dengan } \bar{R}_j' = \frac{1}{n} \sum_{i=1}^n R_{ij}'$$

$$S_{R_j} = \sqrt{\frac{1}{n} \sum_{i=1}^n (R_{ij}' - \bar{R}_j')^2}$$

4. Jika \mathbf{X}'' menyatakan matriks data hasil transformasi langkah ke-3, hitung $(\mathbf{X}'')^t \cdot \mathbf{X}'' = r_{w2}$.

Tentukan nilai eigen dan vektor eigen dari $(\mathbf{X}'')^t \cdot \mathbf{X}''$.