

Nama : Isda Sari

NIM : 0607093

### 3.3 Sampel Acak dan Nilai Harapan dari Rerata Sampel dan Matrik Kovarians

Kita asumsikan bahwa nilai variabel pengamatan berdasarkan himpunan data . Andaikan data bukan hasil pengamatan, tetapi kita berniat mengumpulkan n himpunan dari pengukuran p variabel. Sebelum pengukuran dibuat, nilai-nilai tersebut tidak ada, oleh karena itu diprediksi. Konsekuensinya, diproses sebagai variabel acak.

Didalam konteks ini, misal entri (i,j) di dalam data matrik menjadi variabel acak . Masing-masing himpunan berukuran pada p variabel adalah sebuah vektor acak dan kita mempunyai matrik acak:

$$\underset{(pxn)}{\mathbf{X}} = \begin{bmatrix} \mathbf{X}_{11} & \mathbf{X}_{12} & \cdots & \mathbf{X}_{1n} \\ \mathbf{X}_{21} & \mathbf{X}_{22} & \cdots & \mathbf{X}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}_{p1} & \mathbf{X}_{p2} & \cdots & \mathbf{X}_{pn} \end{bmatrix} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n] \quad (3.8)$$

Sebuah sampel acak didefinisikan:

Jika vektor kolom  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$  pada (3.8) muncul pengamatan independen dari distribusi gabungan dengan fungsi kepadatan  $f(x) = f(x_1, x_2, \dots, x_p)$ , maka  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$  dikatakan bentuk suatu sampel acak dari  $f(x)$ . Secara matematis, hasil  $f(x_1)f(x_2)\dots f(x_n)$ , dimana  $f(x_j) = f(x_{1j}, x_{2j}, \dots, x_{pj})$  adalah fungsi kepadatan dari kolom vektor ke-j.

Perhatikan, dua *point* yang berhubungan dengan definisi sampel acak:

1. Pengukuran p variabel pada sebuah percobaan tunggal, seperti  $\mathbf{X}'_j = [\mathbf{X}_{1j}, \mathbf{X}_{2j}, \dots, \mathbf{X}_{pj}]$ , biasanya berhubungan. Tentu saja, kita harapkan ini menjadi sebuah kasus. Ukuran dari percobaan yang berbeda, bagaimanapun harus menjadi independen.

2. Kebebasan pengukuran dari percobaan ke percobaan tidak mungkin tetap manakala variabel mungkin menyimpang dari waktu ke waktu, seperti himpunan p harga bursa atau p indikator ekonomi.

Jarak Euclid, jika komponen-komponennya suatu vektor bebas dan mempunyai beberapa varians. Misalkan letak baris ke-i  $Y_i' = [X_{i1}, X_{i2}, \dots, X_{in}]$  dari X dipandang sebagai titik pada n dimensi.

Letak dari titik ini ditentukan oleh peluang distribusi gabungan  $f(y_i) = f(x_{i1}, x_{i2}, \dots, x_{in})$  ketika  $x_{i1}, x_{i2}, \dots, x_{in}$  adalah sampel acak,  $f(y_i) = f(x_{i1}, x_{i2}, \dots, x_{in}) = f_i(x_{i1})f_i(x_{i2}) \dots f_i(x_{in})$  dan akibatnya, setiap penambahan koordinat  $x_{ij}$  sama dengan letak distribusi marginal identik  $f_i(x_{ij})$ . Jika n komponen tidak bebas atau distribusi marginal tidak identik, kegagalan dari pengukuran (koordinat) individual pada letak asimetris.

Kesimpulannya, jangkauan distribusi sampling untuk  $\bar{X}$  dan  $S_n$  dengan tidak membuat asumsi mengenai variabel pokok distribusi gabungan. Sehingga kita dapat melihat bagaimana  $\bar{X}$  dan  $S_n$  menyatakan titik estimasi dari populasi korespondensi rerata vektor  $\mu$  dan matrik kovarians  $\Sigma$ .

### Akibat 3.1

Misalkan  $X_1, X_2, \dots, X_n$  adalah sampel acak dari distribusi gabungan dengan rerata vektor  $\mu$  dan matrik kovarians  $\Sigma$ . Maka  $\bar{X}$  adalah estimasi tak bias pada  $\mu$  dan matrik kovariansnya adalah  $\frac{1}{n}\Sigma$ , dimana

$$E(\bar{X}) = \mu \quad (\text{populasi rerata vektor})$$

$$Cov(\bar{X}) = \frac{1}{n}\Sigma \quad (\text{populasi matrik varians-kovarians dibagi oleh ukuran sampel})$$

untuk matrik kovarians  $S_n$ ,  $E(S_n) = \frac{n-1}{n}\Sigma = \Sigma - \frac{1}{n}\Sigma$  maka  $E\left(\frac{n}{n-1}S_n\right) = \Sigma$

Jadi  $\left[\frac{n}{(n-1)}\right]S_n$  merupakan estimasi tak bias pada  $\Sigma$ . Sedangkan  $S_n$

merupakan estimasi bias dengan  $(bias) = E(S_n) - \Sigma = -\left(\frac{1}{n}\right)\Sigma$

Bukti:

Diketahui  $\bar{X} = (X_1 + X_2 + \dots + X_n)/n$ , dengan menggunakan pernyataan (2-24)

$$E(X + Y) = E(X) + E(Y)$$

$$E(AXB) = AE(X)B$$

Diperoleh:

$$\begin{aligned}
 E(\bar{X}) &= E\left(\frac{1}{n}X_1 + \frac{1}{n}X_2 + \dots + \frac{1}{n}X_n\right) \\
 &= E\left(\frac{1}{n}X_1\right) + E\left(\frac{1}{n}X_2\right) + \dots + E\left(\frac{1}{n}X_n\right) \\
 &= \frac{1}{n}E(X_1) + \frac{1}{n}E(X_2) + \dots + \frac{1}{n}E(X_n) \\
 &= \frac{1}{n}\mu + \frac{1}{n}\mu + \dots + \frac{1}{n}\mu \\
 &= \mu
 \end{aligned}$$

Kemudian 
$$\begin{aligned}
 (\bar{X} - \mu)(\bar{X} - \mu)' &= \left(\frac{1}{n}\sum_{j=1}^n (X_j - \mu)\right)\left(\frac{1}{n}\sum_{l=1}^n (X_l - \mu)\right)' \\
 &= \frac{1}{n^2}\sum_{j=1}^n \sum_{l=1}^n (X_j - \mu)(X_l - \mu)'
 \end{aligned}$$

maka 
$$Cov(\bar{X}) = E(\bar{X} - \mu)(\bar{X} - \mu)' = \frac{1}{n^2}\left(\sum_{j=1}^n \sum_{l=1}^n E(X_j - \mu)(X_l - \mu)'\right)$$

untuk  $j \neq l$ , setiap entri di  $E(X_j - \mu)(X_l - \mu)'$  adalah nol.

Sebab entri-entrinya adalah kovarians antara komponen dari  $X_j$  dan komponen  $X_l$  dan independen.

Pada pernyataan (2-29) yaitu:  $Cov(X_i, X_k) = 0$ , jika  $X_i$  dan  $X_k$  independen

Karena itu, 
$$Cov(\bar{X}) = \frac{1}{n^2}\left(\sum_{j=1}^n E(X_j - \mu)(X_j - \mu)'\right).$$

Karena  $\Sigma = E(X_j - \mu)(X_j - \mu)'$  adalah populasi umum matrik kovarians untuk setiap  $X_j$ , kita peroleh

$$Cov(\bar{X}) = \frac{1}{n^2}\left(\sum_{j=1}^n E(X_j - \mu)(X_j - \mu)'\right) = \frac{1}{n^2}\underbrace{(\Sigma + \Sigma + \dots + \Sigma)}_{n \text{ suku}} = \frac{1}{n^2}(n\Sigma) = \left(\frac{1}{n}\right)\Sigma$$

Berlaku pada nilai harapan  $S_n$  dengan catatan  $(X_{ij} - \bar{X}_i)(X_{kj} - \bar{X}_k)$

Merupakan elemen  $(i, k)$  pada  $(X_i - \bar{X})(X_j - \bar{X})'$

Mewakili penjumlahan dari matrik persegi dan perkalian dapat dituliskan sebagai berikut:

$$\sum_{j=1}^n (\mathbf{X}_j - \bar{\mathbf{X}})(\mathbf{X}_j - \bar{\mathbf{X}})' = \sum_{j=1}^n (\mathbf{X}_j - \bar{\mathbf{X}})\mathbf{X}_j' + \left( \sum_{j=1}^n (\mathbf{X}_j - \bar{\mathbf{X}}) \right) (-\bar{\mathbf{X}})' = \sum_{j=1}^n \mathbf{X}_j \mathbf{X}_j' - n\bar{\mathbf{X}}\bar{\mathbf{X}}'$$

karena  $\sum_{j=1}^n (\mathbf{X}_j - \bar{\mathbf{X}}) = 0$  dan  $n\bar{\mathbf{X}}' = \sum_{j=1}^n \mathbf{X}_j'$ . Sehingga nilai harapannya adalah

$$E\left(\sum_{j=1}^n \mathbf{X}_j \mathbf{X}_j' - n\bar{\mathbf{X}}\bar{\mathbf{X}}'\right) = \sum_{j=1}^n E(\mathbf{X}_j \mathbf{X}_j') - nE(\bar{\mathbf{X}}\bar{\mathbf{X}}')$$

Untuk beberapa vektor  $V$  dengan  $E(V) = \mu_v$  dan  $Cov(V) = \Sigma_v$ , kita mempunyai

$$E(VV') = \Sigma_v + \mu_v \mu_v'$$

Konsekuensinya  $E(\mathbf{X}_j \mathbf{X}_j') = \Sigma + \mu \mu'$  dan  $E(\bar{\mathbf{X}}\bar{\mathbf{X}}') = \frac{1}{n} \Sigma + \mu \mu'$

Gunakan hasil ini,  $\sum_{j=1}^n E(\mathbf{X}_j \mathbf{X}_j') - nE(\bar{\mathbf{X}}\bar{\mathbf{X}}') = n\Sigma + n\mu\mu' - n\left(\frac{1}{n}\Sigma + \mu\mu'\right) = (n-1)\Sigma$

karena  $S_n = \left(\frac{1}{n}\right)\left(\sum_{j=1}^n \mathbf{X}_j \mathbf{X}_j' - n\bar{\mathbf{X}}\bar{\mathbf{X}}'\right)$  sedemikian sehingga  $E(S_n) = \frac{(n-1)}{n} \Sigma$ .

### Akibat 3.1

Ditunjukkan bahwa entri  $(i, k)$  dari  $\left[\frac{n}{(n-1)}\right]S_n$  adalah estimasi takbias pada  $\sigma_{ik}$ . Bagaimanapun, sampel individu standar deviasi  $\sqrt{S_{ii}}$  yang bukan estimasi takbias dari koresponden banyak populasi  $\sqrt{\sigma_{ii}}$ , dihitung dengan salah satu  $n$  atau  $n-1$  sebagai pembagi. Selain itu, koefisien korelasi  $r_{ik}$  yang bukan estimasi takbias dari banyak populasi  $\rho_{ik}$ . Meskipun, bias  $E(\sqrt{S_{ii}}) - \sqrt{\sigma_{ii}}$  atau  $E(r_{ik}) - \rho_{ik}$ , biasanya dapat diabaikan jika ukuran  $n$  sampel cukup besar.

Pertimbangan pada bias mendorong suatu perubahan definisi dari matrik varians-kovarians sampel. Akibat 3.1 menyediakan kepada kita dengan suatu estimasi takbias  $S$  untuk  $\Sigma$ . (Takbias) Matrik varians-kovarians sample

$$S = \left(\frac{n}{n-1}\right)S_n = \frac{1}{n-1} \sum_{j=1}^n (\mathbf{X}_j - \bar{\mathbf{X}})(\mathbf{X}_j - \bar{\mathbf{X}})' \dots (3-11)$$

Disini  $S$  mempunyai entri  $(i, k)$ ,  $(n-1)^{-1} \sum (\mathbf{X}_{ij} - \bar{X}_i)(\mathbf{X}_{kj} - \bar{X}_k)$ . Definisi dari sampel kovarians ini biasa dipakai dalam tes statistik multivariat. Oleh karena itu, akan digantikan  $S_n$  sebagai matrik kovarians sampel dalam berbagai materi.

### 3.4 Varians Umum

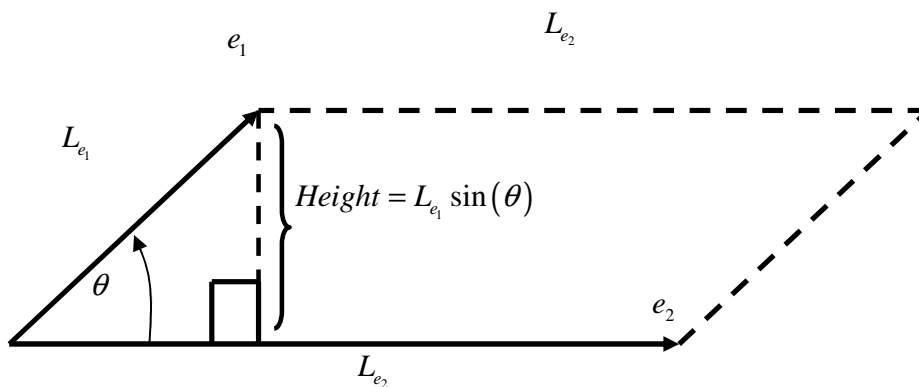
Untuk variabel tunggal biasanya varians sampel digunakan untuk menggambarkan banyaknya variasi pengukuran pada variabel itu. Ketika p variabel merupakan pengamatan pada masing-masing unit, variasi digambarkan oleh matrik varians-kovarians sampel

$$S = \begin{bmatrix} s_{11} & s_{12} & \cdots & s_{1p} \\ s_{21} & s_{22} & \cdots & s_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ s_{1p} & s_{2p} & \cdots & s_{pp} \end{bmatrix} = \left\{ s_{ik} = \frac{1}{n-1} \sum_{j=1}^n (x_{ij} - \bar{x}_i)(x_{kj} - \bar{x}_k) \right\}$$

Matrik kovarians sampel memuat p varians dan  $\frac{1}{2} p(p-1)$  kemungkinan kovarians yang lain. Kadang-kadang itu yang diperlukan untuk menentukan suatu nilai numerik tunggal pada variasi yang dinyatakan oleh S. Satu pilihan untuk nilai determinan dari S, yang mereduksi varians sampel biasa dari suatu karakteristik tunggal  $p = 1$ . Determinan ini disebut varians sampel umum.

$$\text{Varians sampel umum} = |S| \quad \dots (3-12)$$

Ketika  $p > 1$ , di dalam prosesnya beberapa informasi sampel hilang. Penjelasan secara geometrik dari  $|S|$  akan membantu kita mendeskripsikannya. Anggap dua vektor deviasi dalam bidang  $e_1 = y_1 - \bar{x}_1$  dan  $e_2 = y_2 - \bar{x}_2$ . Misalkan  $L_{e_1}$  merupakan panjang dari  $e_1$  dan  $L_{e_2}$  merupakan panjang dari  $e_2$ . Geometrinya,



dan luas trapesium di gambar ini adalah  $(L_{e_1} \sin(\theta)) L_{e_2}$ .

Karena  $\cos^2(\theta) + \sin^2(\theta) = 1$  kita dapat menggambarkan luas ini sebagai

$$\text{Luas} = L_{e_1} L_{e_2} \sqrt{1 - \cos^2(\theta)}$$

Dari (3-5) dan (3-7)

$$L_{e_i}^2 = e_i' e_i = \sum_{j=1}^n (x_{ij} - \bar{x}_i)^2 \quad (3-5)$$

$$r_{ik} = \frac{s_{ik}}{\sqrt{s_{ii}} \sqrt{s_{kk}}} = \cos(\theta_{ik}) \quad (3-7)$$

Sehingga  $L_{e_1} = \sqrt{\sum_{j=1}^n (x_{1j} - \bar{x}_1)^2} = \sqrt{(n-1)s_{11}}$ ,  $L_{e_2} = \sqrt{\sum_{j=1}^n (x_{2j} - \bar{x}_2)^2} = \sqrt{(n-1)s_{22}}$  dan  $\cos(\theta) = r_{12}$

Oleh karena itu,  $Luas = (n-1)\sqrt{s_{11}}\sqrt{s_{22}}\sqrt{1-r_{12}^2} = (n-1)\sqrt{s_{11}s_{22}(1-r_{12}^2)}$  ... (3-13)

Juga,

$$|S| = \begin{vmatrix} s_{11} & s_{21} \\ s_{12} & s_{22} \end{vmatrix} = \begin{vmatrix} s_{11} & \sqrt{s_{11}}\sqrt{s_{22}}r_{12} \\ \sqrt{s_{11}}\sqrt{s_{22}}r_{12} & s_{22} \end{vmatrix} = s_{11}s_{22} - s_{11}s_{22}r_{12}^2 = s_{11}s_{22}(1-r_{12}^2) \quad \dots(3-14)$$

Jika kita dibanding (3-14) dan (3-13), terlihat bahwa  $|S| = \frac{(luas)^2}{(n-1)^2}$

Asumsikan bahwa  $|S| = (n-1)^{-(p-1)} (volume)^2$  menyatakan volume dengan n space oleh p-1 vektor deviasi  $e_1, e_2, \dots, e_{p-1}$  kita dapat menentukan hasil untuk p vektor deviasi dengan induksi

$$Varians sampel umum = |S| = (n-1)^{-p} (volume)^2 \quad \dots\dots (3-15)$$

dengan vektor deviasi  $e_1 = y_1 - \bar{x}_1 \mathbf{1}, e_2 = y_2 - \bar{x}_2 \mathbf{1}, \dots, e_p = y_p - \bar{x}_p \mathbf{1}$

Varians umum digambarkan ke dalam p-space titik sebaran yang mewaliki data.

Seringnya menggambarkan secara intuisi menyangkut sebaran titik rerata sample

$\bar{x}' = [\bar{x}_1, \bar{x}_2, \dots, \bar{x}_p]$ . Anggap pengukuran jarak

$$0 < (jarak)^2 = x'Ax \quad \text{jika } x \neq 0 \quad \dots\dots (2-19)$$

dengan  $\bar{x}$  berperan pada titik tertentu  $\mu$  dan  $S^{-1}$  berperan pada A. Dengan pilihan ini,

koordinat  $x' = [x_1, x_2, \dots, x_p]$  pada titik suatu jarak konstan c dari  $\bar{x}$  yang memenuhi

$$(x - \bar{x})' S^{-1} (x - \bar{x}) = c^2 \quad \dots\dots (3-16)$$

Ketika p=1  $(x - \bar{x})' S^{-1} (x - \bar{x}) = \frac{(x_1 - \bar{x}_1)^2}{s_{11}}$  merupakan jarak persegi dari  $x_1$  ke  $\bar{x}_1$

dalam satuan standar deviasi.

Persamaan (3-16) menetapkan pusat hiperelipsoid di  $\bar{x}$  (suatu elips jika p=2). Ini dapat

ditunjukkan menggunakan kalkulus integral bahwa volume pada hiperelipsoid

berhubungan dengan  $|S|$ . Yaitu,

$$\left\{x : (x - \bar{x})' S^{-1} (x - \bar{x}) \leq c^2\right\} = k_p |S|^{1/2} c^p$$

atau

$$(\text{volume elipsoid})^2 = (\text{kons tan})(\text{variansi sampel umum})$$

dimana  $k_p = \frac{2\pi^{p/2}}{p\Gamma(p/2)}$  ; dengan  $\Gamma(z)$  merupakan notasi fungsi gamma.

Walaupun variansi umum mempunyai beberapa intuisi yang memuaskan penafsiran geometris, inilah kelemahannya dalam menggambarkan matrik kovarians sampel S.

Ilustrasi:

Misalkan tiga matrik kovarians sampel dan diperoleh koefisien kovarians sebagai berikut:

$$S = \begin{bmatrix} 5 & 4 \\ 4 & 5 \end{bmatrix}$$

$$S = \begin{bmatrix} 5 & -4 \\ -4 & 5 \end{bmatrix}$$

$$S = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}$$

$$r_{12} = \frac{s_{12}}{\sqrt{s_{11}}\sqrt{s_{22}}} = \frac{4}{\sqrt{5}\sqrt{5}} = 8$$

$$r_{12} = \frac{-4}{\sqrt{5}\sqrt{5}} = -8$$

$$r_{12} = \frac{0}{\sqrt{3}\sqrt{3}} = 0$$

Setiap matrik kovarians mempunyai beberapa variansi umum yaitu  $|S| = 9$ , namun jelas berbeda struktur korelasi (kovarians). Perbedaan struktur korelasi tidak ditemukan oleh  $|S|$ .