

BAB IV DASAR-DASAR DISKRIMINAN

Misalkan bahwa kita mempunyai k populasi, yaitu $\pi_1, \pi_2, \dots, \pi_k$. Kita akan mengklasifikasikan suatu individu dengan observasi $\underline{x} = (x_1, x_2, \dots, x_p)'$ atau suatu group dari n individu dengan observasi $\underline{x}_j = (x_{j1}, \dots, x_{jp})'$, dimana $j = 1, \dots, n$, dalam p karakter yang berbeda, kedalam salah satu dari $\pi_1, \pi_2, \dots, \pi_k$.

Misal k populasi dituliskan dengan :

$$\pi_1, \pi_2, \dots, \pi_k$$

$$(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_p)', (\underline{y}_1, \underline{y}_2, \dots, \underline{y}_p)', \dots, (\underline{z}_1, \underline{z}_2, \dots, \underline{z}_p)'$$

atau :

$$\begin{pmatrix} x_{11} & x_{21} & \dots & x_{p1} \\ x_{12} & x_{22} & \dots & x_{p2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1n} & x_{2n} & \dots & x_{pn} \end{pmatrix}, \begin{pmatrix} y_{11} & y_{21} & \dots & y_{p1} \\ y_{12} & y_{22} & \dots & y_{p2} \\ \vdots & \vdots & \ddots & \vdots \\ y_{1n} & y_{2n} & \dots & y_{pn} \end{pmatrix}, \dots, \begin{pmatrix} z_{11} & z_{21} & \dots & z_{p1} \\ z_{12} & z_{22} & \dots & z_{p2} \\ \vdots & \vdots & \ddots & \vdots \\ z_{1n} & z_{2n} & \dots & z_{pn} \end{pmatrix}$$

Selanjutnya kita asumsikan bahwa setiap π_i mempunyai fungsi distribusi F_i dari vektor random $\underline{x} = (x_1, \dots, x_p)'$. Kita asumsikan bahwa bentuk fungsi F_i , untuk setiap i diketahui dan F_i berbeda untuk nilai i yang berlainan, walaupun parameter dari F_i diketahui maupun

Misalkan E^p adalah ruang euclid: p dimensi, kita akan membagi E^p kedalam daerah yang saling asing yaitu R_1, R_2, \dots, R_k , sehingga jika x atau $(x_j, j = 1, \dots, p)$ berada dalam R_i , kita dapat menunjukkan individu atau suatu group dalam π_i . Tetapi, dalam menggunakan aturan klasifikasi, kita mungkin membuat suatu kesalahan, dengan suatu kesalahan klasifikasi individu untuk π_i apabila sebenarnya individu itu untuk π_j ($i \neq j$).

Misalkan bahwa harga kesalahan klasifikasi individu untuk π_j apabila sebenarnya individu tersebut untuk π_i adalah $C(j/i)$. Umumnya $C(j/i)$ tidak semuanya sama, tergantung pada kesalahan relatif yang penting. Selanjutnya diasumsikan bahwa tidak ada penyimpangan untuk klasifikasi yang benar. Dengan kata lain $C(i/i) = 0$, untuk semua i .

Kita akan mengklasifikasikan suatu individu tunggal dengan observasi x kedalam salah satu π_i ($i = 1, 2, \dots, k$). Misalkan $\zeta = (R_1, \dots, R_k)$. Kita akan menunjukkan pembagian aturan klasifikasi untuk ruang E dalam keadaan terpisah, menjadi R_1, \dots, R_k dalam ζ . Probabilitas kesalahan klasifikasi individu dengan observasi x pada π_i tetapi dimasukkan dalam π_j adalah :

$$P(j/i, \zeta) = \int_{R_j} f_i(x) dx, \dots \dots \dots (4.1)$$

dimana :

$$dx = \prod_{i=1}^p dx_i$$

Harapan harga kesalahan klasifikasi pada observasi π_i diatas adalah diberikan dengan :

$$r_i(\underline{r}) = \sum_{j=1, j \neq i}^k C(j/i) P(j/i, \underline{r}) \quad , \quad i = 1, 2, \dots, k$$

..... (4.2)

4.1. ATURAN KLASIFIKASI

DEFINISI 4.1.1.

Diberikan dua aturan klasifikasi \underline{r} dan \underline{r}^* , kita katakan bahwa \underline{r} lebih baik atau sama dengan \underline{r}^* , jika $r_i(\underline{r}) \leq r_i(\underline{r}^*)$ untuk semua i . Dan \underline{r} lebih baik dari \underline{r}^* , apabila $r_i(\underline{r}) < r_i(\underline{r}^*)$.

DEFINISI 4.1.2.

Aturan klasifikasi \underline{r} dikatakan Admissible jika tidak ada aturan klasifikasi \underline{r}^* yang lebih baik daripada \underline{r} .

DEFINISI 4.1.3.

Suatu kelas dari aturan klasifikasi dikatakan komplit, jika untuk semua aturan \underline{r}^* diluar kelas ini, kita dapat menentukan aturan \underline{r} didalam kelas yang lebih baik dari \underline{r}^* .

DEFINISI 4.1.4.

Aturan klasifikasi \underline{r}^* dikatakan minimax diantara kelas dari semua aturan \underline{r} jika :

$$\max_i r_i(\underline{r}^*) = \min_{\underline{r}} \max_i r_i(\underline{r})$$

Misalkan bahwa p_i menunjukkan proporsi dari π_i dalam populasi, $i = 1, 2, \dots, k$. Jika p_i diketahui, kita dapat menentukan mean dari kesalahan klasifikasi individu dengan menggunakan aturan \underline{r} . Karena probabilitas pengambilan observasi dari π_i adalah p_i , maka probabilitas

pengambilan observasi dari π_i dan secara tepat klasifikasi

itu kedalam π_i dengan pertolongan aturan ζ , diberikan dengan :

$$p_i P(i/i, \zeta), \quad i = 1, 2, \dots, k$$

Dengan cara sama, probabilitas pengambilan observasi π_i dari kesalahan klasifikasi ke π_j dimana $i \neq j$ adalah :

$$p_i P(j/i, \zeta)$$

Jadi harga :

$$\sum_{i=1}^k p_i \sum_{j=1, j \neq i}^k C(j/i) P(j/i, \zeta) \dots \dots \dots (4.4)$$

adalah mean harga kesalahan klasifikasi untuk aturan ζ terhadap prior probabilitas $g = (p_1, p_2, \dots, p_k)$.

DEFINISI 4.1.5.

Diberikan g , aturan klasifikasi ζ yang meminimumkan mean harga kesalahan klasifikasi disebut Aturan Bayes terhadap g .

Dapat disimpulkan bahwa aturan Bayes dapat menghasilkan probabilitas kesalahan klasifikasi yang besar, disini ada beberapa usaha untuk mengatasi kesulitan ini. Dalam kasus dimana prior probabilitas p_i diketahui, aturan Bayes adalah optimum, dalam arti aturan ini meminimumkan mean harga harapan.

Kita sekarang akan mengevaluasi bentuk eksplisit dari aturan Bayes ini, dalam kasus dimana setiap π_i diambil dari fungsi density probabilitas f_i , $i = 1, \dots, k$. Kita asumsikan bahwa semua prosedur klasifikasi dianggap sama, jika hal itu hanya berbeda pada kumpulan dari ukuran probabilitas berharga nol.

4.2. CONTOH

Misalkan bahwa :

$$f_i(x) = \begin{cases} \beta_i^{-1} \exp(-x/\beta_i) & , & 0 < x < \infty \\ 0 & \text{untuk yang lain} \end{cases}$$

$i = 1, 2, \dots, k$ dan $\beta_1 < \dots < \beta_k$ adalah parameter yang tidak diketahui, dan misalkan bahwa $p_i = i/k$, $i = 1, 2, \dots, k$. Jika x adalah observasi, aturan Bayes dengan $C(i/j)$ sama, menghendaki kita untuk mengklasifikasikan x pada π_j apabila :

$$p_i f_i(x) \geq \max_{(j \neq i)} p_j f_j(x) \quad (4.5)$$

dengan perkataan lain, untuk $i < j$ apabila :

$$\beta_i^{-1} \exp(-x/\beta_i) \geq \beta_j^{-1} \exp(-x/\beta_j) \text{ dipenuhi,}$$

jika dan hanya jika :

$$x \leq \frac{\beta_i \beta_j}{\beta_j - \beta_i} \log(\beta_j - \beta_i)$$

Maka sangatlah mudah untuk menunjukkan bahwa ini adalah fungsi naik dari β_j untuk $\beta_i < \beta_j$ tertentu dan fungsi naik dari β_i untuk $\beta_j < \beta_i$ tertentu. Karena $f_i(x)$ menurun dalam x untuk $x > 0$, berarti bahwa kita mengklasifikasikan x pada π_i apabila :

$$x_{j-1} \leq x \leq x_i,$$

dimana : $x_0 = 0$, $x_k = \infty$ dan

$$x_i = \frac{\beta_i \beta_{i+1}}{\beta_{i+1} - \beta_i} (\log \beta_{i+1} - \log \beta_i)$$

β_i , maka aturan Bayes tidak membuat observasi pada individu dan akan selalu mengklasifikasikan pada π_k .

Sebagai contoh yang lain adalah : misalkan bahwa

$$f_1(X) = \begin{cases} \frac{1}{(2\pi)^{1/2}} \exp \{-1/2(X-\mu_1)^2\} & -\infty < X < \infty \\ 0 & \text{untuk yang lain} \end{cases}$$

dimana μ_1 adalah parameter yang tidak diketahui, dan misalkan $p_j = 1/k$, $i = 1, 2, \dots, k$. Aturan Bayes dengan $C(i/j)$ sama, menghendaki kita untuk mengklasifikasikan observasi X pada π_j apabila :

$$(X - \mu_j)^2 \geq \max_{i, i \neq j} \{(X - \mu_i)^2\}, \dots \dots \dots (4.6)$$

Untuk kasus khusus $k = 2$, aturan klasifikasi Bayes lawan prior (p_1, p_2) diberikan dengan menunjukkan X pada :

$$\left\{ \begin{array}{l} \pi_1, \text{ apabila } \frac{f_1(X)}{f_2(X)} \geq \frac{C(1/2)p_2}{C(2/1)p_1} \\ \pi_2, \text{ apabila } \frac{f_1(X)}{f_2(X)} < \frac{C(1/2)p_2}{C(2/1)p_1} \\ \pi_1 \text{ atau } \pi_2, \text{ apabila } \frac{f_1(X)}{f_2(X)} = \frac{C(1/2)p_2}{C(2/1)p_1} \end{array} \right. \dots \dots (4.7)$$

Meskipun demikian, jika berada dalam π_i , $i = 1, 2$

$$P \left\{ \frac{f_1(X)}{f_2(X)} = \frac{C(1/2)p_2}{C(2/1)p_1} \mid \pi_i \right\} = 0 \dots \dots \dots (4.8)$$

maka aturan Bayes adalah tunggal kecuali untuk kumpulan dari ukuran probabilitas berharga nol.

4.3. ATURAN LIKELIHOOD

Rasio likelihood aturan klasifikasi $\underline{r} = (R_1, R_2,$

$$R_j = C_j f_j(x) > \max_{i, i \neq j} C_i f_i(x) \dots\dots\dots (4.9)$$

untuk konstanta positif C_1, \dots, C_k . Khusus, jika C_i semua sama, maka aturan klasifikasinya disebut Aturan Likelihood Maximum.

Jika distribusi F_i tidak lengkap diketahui (mean, variansinya tidak diketahui), pelengkap informasi pada F_i atau parameter didalamnya diperoleh langsung dari sampel dalam populasi yang bersesuaian. Maka asumsi lengkap diketahui dari F_i , aturan klasifikasi yang baik apabila $C = (R_1, R_2, \dots, R_k)$ antara lain memenuhi Bayes, Minimax dan Aturan Rasio Likelihood.

