

## Analiza discriminantă

Sub numele de analiză discriminantă, sunt reunite o serie de metode explicative, descriptive și predictive, destinate studierii unei populații împărțite în clase. Fiecare individ este caracterizat printr-un ansamblu de variabile independente predictor și o variabilă calitativă identificând clasa din care face parte. Mulțimea de indivizi se împarte în două submulțimi:

- *eșantionul de bază*, pentru care se cunoaște valoarea variabilei calitative, deci încadrarea pe clase a indivizilor;
- *eșantionul neinvestigat*, pentru care nu se cunoaște încadrarea indivizilor în clase, deci nici valoarea variabilei calitative.

Analiza discriminantă urmărește pe de o parte să identifice regulile pe baza cărora toți indivizii din mulțimea investigată să poată fi încadrați în clase, iar pe de altă parte să reducă numărul de variabile necesare realizării împărțirii în clase (discriminării). Primul aspect urmărit scoate în evidență caracterul predictiv, decizional al analizei discriminante pe când al doilea aspect reliefează caracterul ei descriptiv. Analiza discriminantă se aplică în mod frecvent în variate domenii în care recunoașterea formelor este utilă (previzionarea comportamentelor solicitanților de credit, detecția fraudelor, diagnostic medical, prognoze meteo, recunoaștere vocală etc).

### Notatii

$$\text{Matricea de observații: } X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ x_{21} & x_{22} & \dots & x_{2m} \\ \dots & & & \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{bmatrix},$$

unde  $n$  este numărul de observații,  $m$  - numărul de variabile predictor (variabile independente).

$$\text{Variabila discriminantă: } Y = \begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix}. \text{ Este o variabilă calitativă. O valoare } y_i, i=1, n \text{ reprezintă clasa din care face parte observația } i.$$

Vectorii observațiilor sunt notați cu  $w_i, i=1, n$ . Un vector  $X_i$  este linia  $i$  din matricea  $X$ .

Vectorii de variabile sunt notați cu  $X_j, j=1, m$  și sunt constituiți din coloanele matricei  $X$ .

$$\text{Matricea centrilor de grupă: } G = \begin{bmatrix} g_{11} & g_{12} & \dots & g_{1m} \\ g_{21} & g_{22} & \dots & g_{2m} \\ \dots & & & \\ g_{q1} & g_{q2} & \dots & g_{qm} \end{bmatrix}, \text{ unde } q \text{ este numărul de grupe. O valoare } g_{kj} \text{ reprezintă media variabilei}$$

$$\text{predictor } j \text{ pentru grupa } k. \text{ Vectorii centrilor de grupă sunt notați cu } G_k, k=1, q, \quad G_k = \begin{bmatrix} g_{k1} \\ \dots \\ g_{km} \end{bmatrix}. \text{ Media generală: } \bar{X} = \begin{bmatrix} \bar{x}_1 \\ \dots \\ \bar{x}_m \end{bmatrix}.$$

$$\text{Matricea diagonală a frecvențelor grupelor: } D_G = \begin{bmatrix} n_1 & 0 & \dots & 0 \\ 0 & n_2 & \dots & 0 \\ \dots & & & \\ 0 & 0 & \dots & n_q \end{bmatrix}.$$

## Indicatori de variabilitate și împrăștiere

Discriminarea între grupe se realizează cu ajutorul indicatorilor de variabilitate și împrăștiere.

1. **Matricele de împrăștiere** (*sum of square and cross product*) reflectă împrăștierea la nivelul întregii colectivități (*SST*), în interiorul grupelor (*SSW*) și împrăștierea grupelor între ele (*SSB*).

*SST* este matricea de împrăștiere la nivelul întregii colectivități și arată gradul de împrăștiere în jurul mediei generale. Termenul general al matricei este:

$$\begin{aligned} SST_{jl} &= \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{il} - \bar{x}_l) = \sum_{k=1}^q \sum_{i \in k} (x_{ij} - \bar{x}_j)(x_{il} - \bar{x}_l) = \sum_{k=1}^q \sum_{i \in k} (x_{ij} - g_{kj} + g_{kj} - \bar{x}_j)(x_{il} - g_{kl} + g_{kl} - \bar{x}_l) = \\ &= \sum_{k=1}^q \sum_{i \in k} (x_{ij} - g_{kj})(x_{il} - g_{kl}) + \sum_{k=1}^q \sum_{i \in k} (g_{kj} - \bar{x}_j)(g_{kl} - \bar{x}_l) + \sum_{k=1}^q \sum_{i \in k} (x_{ij} - g_{kj})(g_{kl} - \bar{x}_l) + \sum_{k=1}^q \sum_{i \in k} (g_{kj} - \bar{x}_j)(x_{il} - g_{kl}) = \\ &= \sum_{k=1}^q \sum_{i \in k} (x_{ij} - g_{kj})(x_{il} - g_{kl}) + \sum_{k=1}^q n_k (g_{kj} - \bar{x}_j)(g_{kl} - \bar{x}_l) + \sum_{k=1}^q (g_{kl} - \bar{x}_l) \sum_{i \in k} (x_{ij} - g_{kj}) + \sum_{k=1}^q (x_{il} - g_{kl}) \sum_{i \in k} (g_{kj} - \bar{x}_j) \end{aligned}$$

Prima sumă,  $\sum_{k=1}^q \sum_{i \in k} (x_{ij} - g_{kj})(x_{il} - g_{kl})$ , reprezintă termenul general al matricei de împrăștiere în interiorul grupelor,  $SSW$ . A doua sumă,

$\sum_{k=1}^q n_k (g_{kj} - \bar{x}_j)(g_{kl} - \bar{x}_l)$ , este termenul general al matricei împrăstierii între grupe,  $SSB$ . Sumele trei și patru au valoarea 0 deoarece

sumele simple ale abaterilor față de mediile grupelor sunt 0,  $\sum_{i \in k} (x_{ij} - g_{kj}) = 0$ ,  $\sum_{i \in k} (x_{il} - g_{kl}) = 0$ . Deci:

$$SST_{jl} = SSW_{jl} + SSB_{jl}, j=1, m, l=1, m.$$

La nivel matriceal:  $SST = SSW + SSB$ .

Dacă ținem cont de gradele de libertate, matricele de împrăștiere sunt calculate ca valori medii astfel:  $MST = \frac{SST}{n-1}$ ,  $MSW = \frac{SSW}{n-q}$ ,

$$MSB = \frac{SSB}{q-1}.$$

2. *Matricele de covarianță*. Termenii generali ai matricelor de covarianță diferă de cei ai matricelor de împrăștiere prin faptul că se calculează ca valori medii.

$$T_{jl} = \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{il} - \bar{x}_l), j=1, m, l=1, m - \text{covarianța totală}$$

$$W_{jl} = \sum_{k=1}^q \frac{n_k}{n} \cdot \frac{1}{n_k} \sum_{i \in k} (x_{ij} - g_{kj})(x_{il} - g_{kl}), j=1, m, l=1, m - \text{covarianța intra-grupă (în interiorul grupelor)}$$

$$B_{jl} = \sum_{k=1}^q \frac{n_k}{n} (g_{kj} - \bar{x}_j)(g_{kl} - \bar{x}_l), j=1, m, l=1, m - \text{covarianța inter-grupă (între grupe)}.$$

Relația între termeni este aceeași ca și la matricele de împrăștiere:

$$T_{jl} = W_{jl} + B_{jl}, j=1, m, l=1, m.$$

La nivel matriceal:

$$T = W + B.$$

3. *Varianța totală*. Varianța totală este reflectată de valorile de pe diagonalele principale ale matricei de covarianță și se calculează ca suma varianțelor celor  $m$  variabile:

$VT = \text{Trace}(T)$  - varianța totală generală,

$VW = \text{Trace}(W)$  - varianța totală intra-grupă,

$VB = \text{Trace}(B)$  - varianța totală inter-grupă,

Varianța totală se poate calcula și pe baza matricelor de împrăștiere:

$$VT = \text{Trace}(SST), VW = \text{Trace}(SSW), VB = \text{Trace}(SSB).$$

4. *Varianța generalizată*. Varianța generalizată se calculează ca determinant al matricelor de covarianță sau de împrăștiere:

$$VGT = |T|, VGW = |W|, VGB = |B|, \text{ sau } VGT = |SST|, VGW = |SSW|, VGB = |SSB|.$$

## Semnificația modelului. Teste statistice

Testarea modelului se face în două etape:

- un test F (Fisher) bazat pe statistica Wilks care arată dacă ansamblul de variabile predictor poate face discriminarea pe grupe a instanțelor

- teste statistice individuale pentru fiecare variabilă predictor prin care se decide dacă o variabilă poate fi sau nu un bun predictor

### Testul F global

$$H_0: g_1 = g_2 = \dots = g_q$$

$$H_1: \text{există două grupe } i, k, \text{ astfel încât } g_i \neq g_k$$

Se calculează un indicator  $\Lambda = \frac{|SSB|}{|SSB + SSW|}$ . Cu cât valoarea lui  $\Lambda$  este mai mare cu atât șansele ca ipoteza  $H_0$  să fie respinsă

sunt mai mari. Statistica testului este o valoare F calculată astfel:  $F = \frac{1 - \Lambda^{1/b}}{\Lambda^{1/b}} \frac{ab - c}{m(q-1)}$ , unde  $a$ ,  $b$  și  $c$  se calculează astfel:

$$a = n - q - \frac{m - q + 2}{2}, b = \begin{cases} \sqrt{\frac{m^2(q-1) - 4}{m^2 + (q-1)^2 - 5}} & \text{dacă } m^2 + (q-1)^2 - 5 > 0 \\ 1 & \text{dacă } m^2 + (q-1)^2 - 5 \leq 0 \end{cases}, c = \frac{m(q-1) - 2}{2}.$$

Dacă  $F < F_{m(q-1), ab-c; \alpha}^{Critic}$  ipoteza nulă este respinsă cu un nivel de încredere  $1-\alpha$ .

### Teste statistice la nivel de variabile predictor

O variabilă predictor poate fi considerată bun predictor dacă poate separa cât mai clar grupele. Deci raportul dintre varianța medie intergrupe și varianța medie intragrupe este cât mai mare. Varianțele sunt preluate din matricele de împrăștiere/covarianță de pe diagonala principală. Fiind vorba de un raport între două varianțe, testul aplicat este testul  $F$ . Astfel pentru variabila predictor  $X_j$ , ipotezele nulă și alternativă sunt:

$$H_0: g_{1j} = g_{2j} = \dots = g_{qj}$$

$$H_1: \exists k, i \text{ două grupe astfel încât } g_{kj} \neq g_{ij}$$

$$\text{Statistica testului: } F_j = \frac{MSB_{jj}}{MSW_{jj}}.$$

Valoarea critică pentru  $q-1$  și  $n-q$  grade de libertate și un prag de semnificație  $\alpha$  este  $F_{q-1; n-q; \alpha}^{Critic}$ .

Dacă  $F_j > F_{q-1; n-q; \alpha}^{Critic}$  ipoteza nulă este respinsă cu un nivel de încredere  $1-\alpha$ .