

Klasická diskriminací analýza

data:

$$\begin{pmatrix} x_{11} & \dots & x_{1p} & y_1 \\ \vdots & & \vdots & \vdots \\ x_{m1} & \dots & x_{mp} & y_m \end{pmatrix}$$

$x_i = (x_{i1}, \dots, x_{ip})^T$... vektor regresorů (prediktorů) pro i -té pozorování

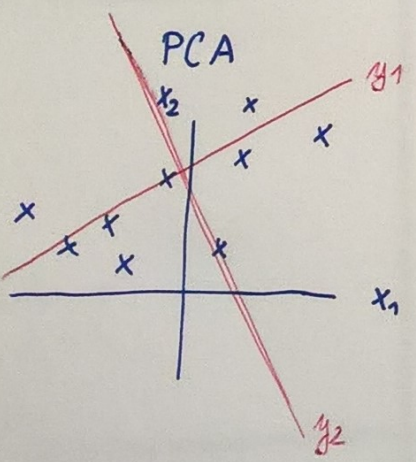
y_i - udává příslušnost i -lého pozorování k dané skupině
 - kategoriální proměnná \rightarrow kategoriemi $\{1, 2, \dots, J\}$

Cíl: Vybrat rozhodovací pravidlo, které bude co nejlépe klasifikovat nové pozorování $x^* = (x_1^*, \dots, x_p^*)^T$ do správné skupiny.

Budeme předpokládat, že naše data jsou generována spojilým náhodným vektorem $X = (X_1, \dots, X_p)^T$.

1) Kanonická diskriminací analýza

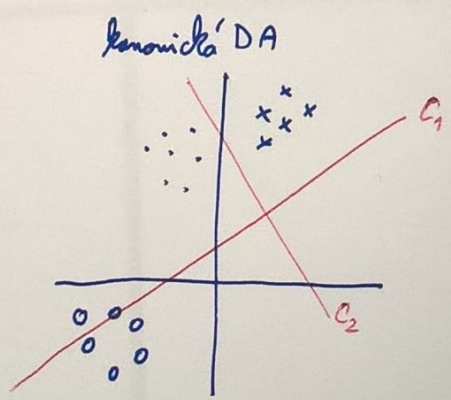
- kombinace PCA a MANOVA



$$y_1 = c_1^T X$$

$$y_2 = c_2^T X$$

$$\vdots$$



$$c_1 = \mu_1^T X$$

$$c_2 = \mu_2^T X$$

$$\vdots$$

maximalizuj se podíl mezi skupinové a vnitroskupinové variability

$$\lambda = \frac{\mu_1^T B \mu_1}{\mu_1^T \Sigma \mu_1}$$

$$B = \frac{1}{J} \sum_{i=1}^J (\mu_i - \mu) (\mu_i - \mu)^T$$

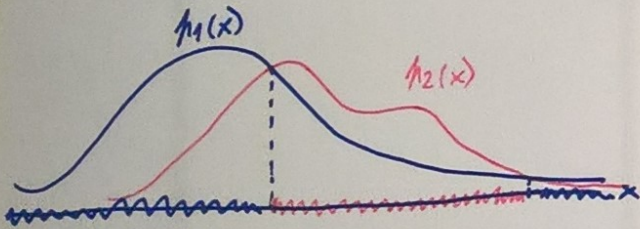
Σ je společný rozptyl skupin

2) Fisherova diskriminační analýza

- měří nat. vel. X má v_j -lé skupině sdruženou hustotu $f_j(x)$, j .

$X | Y=j$ má hustotu $f_j(x)$ pro $j=1, \dots, J$
 \uparrow
 známá p -rozměrná hustota

Pr. / $p=1, 2$ skupiny ($J=2$)



x^* klasifikuj do "pravděpodobnější" skupiny

a co když skupiny nejsou zhruba stejně rozměrné?

oznámě $\pi_j = P(Y_i=j) = P(\text{dané pozorování patří do třídy } j) \text{ pro } j=1, \dots, J$... a priori pravděpodobnost

rovněž $\pi_1 + \pi_2 + \dots + \pi_J = 1$.

$$\text{Bayesova věta: } P(Y=j | X=x) = \frac{P(X=x | Y=j) \cdot P(Y=j)}{P(X=x | Y=1) \cdot P(Y=1) + \dots + P(X=x | Y=J) \cdot P(Y=J)} = \frac{f_j(x) \cdot \pi_j}{\underbrace{f_1(x) \cdot \pi_1 + \dots + f_J(x) \cdot \pi_J}_{\text{konstanta (nezávislá na } j)}} \quad (j=1, \dots, J)$$

aposteriori pravděpodobnost

Bayesovo rozhodovací pravidlo - pozorování x řadí do třídy $\operatorname{argmax}_{j=1, \dots, J} P(Y=j | X=x) = \operatorname{argmax}_{j=1, \dots, J} f_j(x) \cdot \pi_j = \operatorname{argmax}_{j=1, \dots, J} \{ \log f_j(x) + \log \pi_j \}$

- minimalizuj pravděpodobnost špatné klasifikace

a) Lineární diskriminační analýza (LDA)

$$X|Y=j \sim N_p(\mu_j, \Sigma)$$

$$f_j(x) = \frac{1}{\sqrt{(2\pi)^p |\Sigma|}} \cdot \exp \left\{ -\frac{1}{2} (x - \mu_j)^T \Sigma^{-1} (x - \mu_j) \right\} \quad x \in \mathbb{R}^p$$

Baysovské rozhodovací pravidlo:

$$\log \pi_j + \log f_j(x) = \log \pi_j + \underbrace{\frac{p}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma| - \frac{1}{2} x^T \Sigma^{-1} x}_{\text{nezávislé na } j} + \frac{1}{2} x^T \Sigma^{-1} \mu_j + \frac{1}{2} \mu_j^T \Sigma^{-1} x - \frac{1}{2} \mu_j^T \Sigma^{-1} \mu_j = \underbrace{\log \pi_j - \frac{1}{2} \mu_j^T \Sigma^{-1} \mu_j + \mu_j^T \Sigma^{-1} x}_{L_j(x)} + K$$

$L_j(x)$... lineární diskriminační funkce

neznamé parametry π_j , μ_j a Σ_j třeba odhadnout z dat

- $\hat{\pi}_j = \frac{n_j}{n}$... pokud víme kategorický skupin v datech

- $= \frac{1}{J}$... pokud jsou všechny skupiny zastoupeny rovnoměrně

- $\hat{\mu}_j$ odhadneme pomocí vzájemného průměru v j -lé skupině

- S_j ... vzájemná kovarianční matice v j -lé skupině

$$\hat{\Sigma} = \frac{1}{n-J} \left((n_1-1) S_1 + \dots + (n_J-1) S_J \right)$$

B) Kvadratická diskriminační analýza (QDA)

$$X|Y=j \sim N_p(\mu_j, \Sigma_j)$$

$$f_j(x) = \frac{1}{\sqrt{(2\pi)^p |\Sigma_j|}} \cdot \exp \left\{ -\frac{1}{2} (x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j) \right\}, x \in \mathbb{R}^p$$

Bayesovo rozhodovací pravidlo:

$$\log \pi_j + \log f_j(x) = \log \pi_j - \frac{p}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_j| - \frac{1}{2} \mu_j^T \Sigma_j^{-1} \mu_j + \mu_j^T \Sigma_j^{-1} x - \frac{1}{2} x^T \Sigma_j^{-1} x$$

metainna j

$Q_j(x)$... kvadratická diskriminační funkce