

Posplošena diskriminantna analiza z uporabo posplošenega singularnega razcepa

Jernej Banevec

Mentor: izred. prof. dr. Marjeta Knez

Fakulteta za matematiko in fiziko

kratka predstavitev diplomskega dela

8.11.2017

Pregled vsebine

- 1 Posplošena diskriminantna analiza
- 2 Posplošena diskriminantna analiza kot optimizacijski problem
- 3 Posplošen singularni razcep
- 4 Zgled
- 5 Viri

Posplošena diskriminantna analiza

- Posplošitev linearne diskriminantne analize
- Ena zelo uporabljenih statističnih metod
- Oblike podatkov:
 - Združeni v matriki $A \in \mathbb{R}^{m \times n}$
 - $m \dots$ dimenzija posamezne meritve
 - $n \dots$ število meritev oz. podatkov
 - Podatki grupirani v k razredov oz. gruč

Posplošena diskriminantna analiza

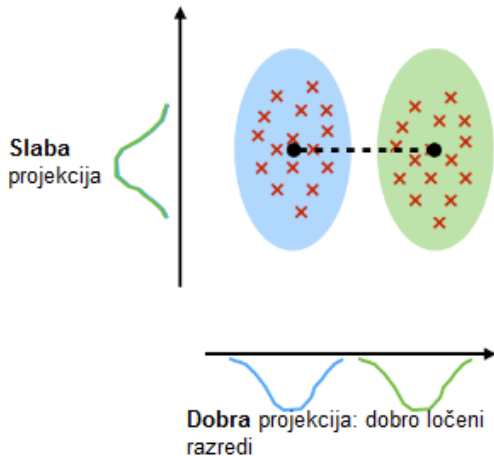
- Iščemo preslikavo

$$G : \mathbb{R}^m \rightarrow \mathbb{R}^\ell,$$

kjer je $\ell \leq m - 1$

- Cilj:
 - Ohraniti razporejenost razredov
 - Zmanjšati razpršenost podatkov znotraj razredov
 - Povečati razlike med razredi

Primerjava dobra proti slabi preslikavi



Definicije

- Centroid i -tega razreda: $c^{(i)} = \frac{1}{n_i} \sum_{j \in N_i} a_j$
- Centroid celotnih podatkov: $c = \frac{1}{n} \sum_{j=1}^n a_j$
- Matrika razpršenosti podatkov znotraj razreda:

$$S_W = \sum_{i=1}^k \sum_{j \in N_i} (a_j - c^{(i)})(a_j - c^{(i)})^T$$

- Matrika razlik med razredi:

$$S_B = \sum_{i=1}^k n_i (c^{(i)} - c)(c^{(i)} - c)^T$$

- Matrika celotne razpršenosti:

$$S_M = S_W + S_B$$

- Tu so vse matrike elementi $\mathbb{R}^{m \times m}$

Definicije

- Preslikava S_W na prostor dimenzije ℓ :

$$S_W^\ell = GS_W G^T$$

- Preslikava S_B na prostor dimenzije ℓ :

$$S_B^\ell = GS_B G^T$$

- Preslikava S_M na prostor dimenzije ℓ :

$$S_M^\ell = GS_M G^T$$

- Tu so vse matrike elementi $\mathbb{R}^{\ell \times \ell}$

- $A^\ell = GA$

Posplošena diskriminantna analiza kot optimizacijski problem

- $sled(S_W) = \sum_{i=1}^k \sum_{j \in N_i} \|a_j - c^{(i)}\|_2^2$
- $sled(S_B) = \sum_{i=1}^k n_i \|c^{(i)} - c\|_2^2$
- Želimo:
 - povečati $sled(S_B^\ell)$
 - zmanjšati $sled(S_W^\ell)$
- Dobimo optimizacijski problem, pri katerem iščemo takšno preslikavo G , ki maksimizira

$$sled(GS_B G^T) / sled(GS_W G^T) \approx sled((S_W^\ell)^{-1} S_B^\ell)$$

- Uporabno le ko je S_W^ℓ nesingularna oz. obrnljiva

Posplošeni singularni razcep

- Originalna definicija posplošenega singularnega razcepa (Van Loan)

Izrek (Posplošeni singularni razcep)

Za matriki $K_A \in \mathbb{R}^{p \times m}$ z $p \geq m$ in $K_B \in \mathbb{R}^{n \times m}$ obstajata ortogonalni matriki $U \in \mathbb{R}^{p \times p}$ in $V \in \mathbb{R}^{n \times n}$ ter nesingularna matrika $X \in \mathbb{R}^{m \times m}$, da velja

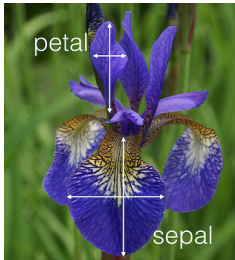
$$U^T K_A X = \text{diag}(\alpha_1, \dots, \alpha_m) \text{ in } V^T K_B X = \text{diag}(\beta_1, \dots, \beta_q)$$

kjer $q = \min(n, m)$, $\alpha_i \geq 0$ za $1 \leq i \leq m$ in $\beta_i \geq 0$ za $1 \leq i \leq q$.

- Pozneje dodatno posplošimo singularni razcep

Zgled - roža Iris

- Precej poznan zgled
- Trije razredi:
 - 1 Iris - setosa
 - 2 Iris - versicolor
 - 3 Iris - virginica



Zgled - roža Iris

- Podatki:

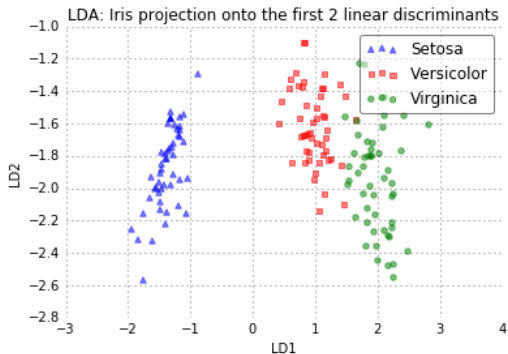
$$A = \begin{bmatrix} a_{1_{\text{sepal length}}} & a_{2_{\text{sepal length}}} & \cdots & a_{150_{\text{sepal length}}} \\ a_{1_{\text{sepal width}}} & a_{2_{\text{sepal width}}} & \cdots & a_{150_{\text{sepal width}}} \\ a_{1_{\text{petal length}}} & a_{2_{\text{petal length}}} & \cdots & a_{150_{\text{petal length}}} \\ a_{1_{\text{petal width}}} & a_{2_{\text{petal width}}} & \cdots & a_{150_{\text{petal width}}} \end{bmatrix}, Y = \begin{bmatrix} \omega_{\text{setosa}} \\ \omega_{\text{setosa}} \\ \vdots \\ \omega_{\text{versicolor}} \\ \vdots \\ \omega_{\text{virginica}} \end{bmatrix}$$

- Na teh podatkih uporabimo posplošeno diskriminantno analizo

- 1 Računanje centroidov
- 2 Računanje matrik razpršenosti podatkov
- 3 Lastni razcep $(S_W^\ell)^{-1} S_B^\ell$
- 4 Urejanje lastnih vrednosti (po parih z vektorji)
- 5 Transformiranje na nov podprostor

Zgled - roža Iris

- Slikamo na 2-dimenzionalen (pod)prostor
- Dobimo sledečo sliko:



- Howland P. in Park H., Generalizing discriminant analysis using the generalized singular value decomposition. TPAMI. 2004, 26, 8, 995-1006
- Jieping Ye, Characterization of a Family of Algorithms for Generalized Discriminant Analysis on Undersampled Problems. Journal of Machine Learning Research. 2005, 6, 483-502