

1 Uvod

1.1 Abstraktni uvod

Diskriminantna analiza se že dolga leta uporablja za določevanje lastnosti, ki ohranjajo razlike med razredi. Definirana je kot optimizacijski problem, ki vključuje kovariančne matrike, ki predstavljajo razprešenost podatkov znotraj posameznega razreda in razpršenost oziroma ločenost posameznih razredov. Diskriminantna analiza pa sama po sebi zahteva, da je ena od teh kovariančnih matrik nesingularna, kar omejuje njeno uporabo na matrikah določenih dimenzij. V nadaljevanju tako preučimo več različnih optimizacijskih kriterijev in poskušamo njihovo uporabo razširiti na vse matrike z uporabo posplošenega singularnega razcepa. Na ta način se izognemo pogoju nesingularnosti, ki ga zahteva diskriminantna analiza. Na ta način dobimo posplošeno diskriminantno analizo, ki jo lahko uporabimo tudi kadar je ena matrika nesingularna (v nadaljevanju lahko vidimo, da je matrika nesingularna kadar je velikost vzorca manjša kot pa dimenzija posamezne meritve – **NEJASNOST!**). V nadaljevanju bom testiral učinkovitost posplošene diskriminantne analize in jo, kjer bo to mogoče, primerjal tudi z diskriminantno analizo.

1.2 Matematični uvod

Cilj diskriminantne analize je združevati lastnosti originalnih podatkov na način, ki kar najučinkoviteje ločuje med razredi, v katerih so podatki. Pri takšnem združevanju lastnosti podatkov se dimenzija teh podatkov zmanjša na način, ki najbolj ohranja strukturo določenih razredov.

Tu predpostavimo, da so podatki zloženi v matiko $A \in \mathbb{R}^{m \times n}$, kjer m predstavlja dimenzijo posamezne meritve, n pa predstavlja število meritev oz. podatkov. Denimo, da so podatki v matriki A iz k različnih razredov. Tako so stolpci a_1, a_2, \dots, a_n matrike A združeni v k podmatrik, ki predstavljajo razrede, v katerih so podatki:

$$A = (A_1, A_2, \dots, A_k), \text{ kjer } A_i \in \mathbb{R}^{m \times n_i}.$$

Cilj diskriminantne analize najti preslikavo G^T , ki v novem, manjdimenzionalnem prostoru, kar najbolj ohranja razrede, v katerih so podatki. Za preslikavo G^T torej velja:

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

kjer je $\ell \leq m - 1$. Torej preslikava G^T nek m -dimenzionalen vektor preslika v nov vektor v ℓ -dimezionalnem prostoru (navadno velja $\ell \leq m$), v katerem so razredi podatkov ohranjeni, razpršenost podatkov znotraj razredov je zmanjšana, razlike med razredi pa so povečane.

Tu število n_i predstavlja moč indeksne množice razreda i . To indeksno množico razreda i označujemo z N_i . Očitno velja tudi:

$$\sum_{i=1}^k n_i = n.$$

Za nadaljnje izračune moramo definirati tudi centroid i -tega razreda, ki je izračunan kot povprečje stolpcev v i -tem razredu, torej:

$$c^{(i)} = \frac{1}{n_i} \sum_{j \in N_i} a_j$$

in centroid celotnih podatkov, ki je izračunan kot povprečje vseh stolpcev:

$$c = \frac{1}{n} \sum_{j=1}^n a_j.$$

Razpršenost podatkov v razredih, razpršenost vseh podatkov ter razpršenost oziroma razlike med razredi je smiselno predstaviti s pomočjo matrik. Zato v nadaljevanju definiramo matriko, ki predstavlja matriko razpršenosti podatkov znotraj razredov:

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

matriko, ki predstavlja matriko razpršenosti oz razlik med razred:

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

in matriko celotne razpršenosti podatkov:

$$S_M = \sum_{j=1}^n (a_j - c)(a_j - c)^T$$

Med zgoraj definiranimi matrikami velja tudi enakost: $S_M = S_W + S_B$. – **DOKAŽI**

S pomočjo preslikave G^T preslikamo v ℓ -dimezionalen prostor tudi matrike S_W , S_B in S_M :

$$S_W^\ell = G^T S_W G, \quad S_B^\ell = G^T S_B G, \quad S_M^\ell = G^T S_M G.$$

Iz danih matrik razpršenosti podatkov bi radi tvorili kriterij kvalitete razredov. Kriterij kvalitete razredov bi imel visoko vrednost, kadar bi bili razredi, v katerih so podatki, strnjeni in dobro ločeni med seboj. Opazimo lahko, da $sled(S_W)$ predstavlja kako skupaj so si podatki v posameznem razredu, saj velja:

$$\begin{aligned} sled(S_W) &= \sum_{t=1}^m \left(\sum_{i=1}^k \sum_{j \in N_i} (a_{jt} - c_t^{(i)})^2 \right) = \sum_{i=1}^k \sum_{j \in N_i} \left(\sum_{t=1}^m (a_{jt} - c_t^{(i)})^2 \right) \\ &= \sum_{i=1}^k \sum_{j \in N_i} \|a_{jt} - c_t^{(i)}\|_2^2. \end{aligned}$$

Podobno $sled(S_B)$ predstavlja ločenost med razredi, saj velja:

$$\begin{aligned} sled(S_B) &= \sum_{t=1}^m \left(\sum_{i=1}^k \sum_{j \in N_i} (c_t^{(i)} - c_t)^2 \right) = \sum_{i=1}^k \sum_{j \in N_i} \left(\sum_{t=1}^m (c_t^{(i)} - c_t)^2 \right) \\ &= \sum_{i=1}^k \sum_{j \in N_i} \left\| c_t^{(i)} - c_t \right\|_2^2. \end{aligned}$$

Optimalna preslikava G^T tako maksimizira $sled(S_B^\ell)$ in minimizira $sled(S_W^\ell)$. Smislen kriterij se tako zdi

$$sled(G^T S_B G) / sled(G^T S_W G),$$

ki pa ga zaradi težke izračunljivosti aproksimiramo kar z

$$sled((S_W^\ell)^{-1} S_B^\ell).$$

– NEJASNOST!

Kljub temu, da je ta optimizacijski kriterij lažje izračunljiv ima svoje pomanjkljivosti. Opazimo lahko, da kriterij lahko uporabimo le v primeru, ko je matrika S_W^ℓ nesingularna oz. da kriterija nemoremo uporabiti, ko je matrika S_W^ℓ singularna (torej kadar je njena determinanta enaka 0). Ker pa za determinanto matrike velja:

$$det(S_W^\ell) = det(G^T S_W G) = det(G^T) \cdot det(S_W) \cdot det(G),$$

je $det(S_W^\ell)$ enaka 0 kadar je $det(S_W)$ enaka 0, torej kadar je matrika S_W singularna. Do te situacije pa lahko pride kar precej pogosto. Matrika $S_W \in \mathbb{R}^{m \times m}$ je singularna namreč v vseh primerih, ko je za matriko $A \in \mathbb{R}^{m \times n}$ velja $m > n$, saj je potem m -dimenzionalna matrika S_W , ki je sestavljena iz n vektorjev (iz $a_j - c^{(i)}$ za $\forall i \in \{1, \dots, k\}$ in $\forall j \in N_j$). Iz n vektorjev pa lahko sestavimo le n dimenzionalen prostor, torej bo $m \times m$ matrika iz teh vektorjev očitno singularna, torej bo njena determinanta enaka 0. Na primer, do tega problema pride v primeru, ko je pridobivanje podatko drago oz. zahtevno in so pridobljeni podatki visokih dimenzij (dimenzija posameznega podatka je večja od števila vseh pridobljenih podatkov).

Poznamo več načinov, kako aplicirati diskriminantno analizo na matriki $A \in \mathbb{R}^{m \times n}$ z $m > n$. V grobem jih ločimo na tiste, kjer dimenzijo zmanjšamo v dveh korakih in na tiste, kjer dimenzijo podatkov zmanjšamo v enem koraku. V prvem načinu se faza diskriminante analize nadaljuje v fazo, v kateri zanemarimo oblike posameznih razredov. Najpopularnejša metoda za prvi del tega procesa je zmanjšanje rang s pomočjo singularnega razcepa. To je tudi glavno orodje metode imenovane principalna komponentna analiza. Kakorkoli, celotna predstava dvostopenjskih načinov je precej občutljiva na zmanjšanje dimenzije v prvi fazi. Sam se bom bolj osredotočil na način, ki posploši diskriminantno analizo tako, da teoretično optimalno zmanjša dimenzijo, brez na bi uvedel dodaten korak. Zato obravnavamo kriterij

$$sled((S_2^Y)^{-1} S_1^Y),$$

kjer sta matriki S_2 in S_1 izbrani iz matrik S_W , S_B in S_M . Klasična diskriminantna analiza predstavi svojo rešitev s pomočjo posplošenega problema lastnih vrednosti, kadar je matrika S_2 nesingularna. Z prestrukturiranjem problema tako, da uporabimo posplošeni singularni razcep, pa razširimo uporabnost diskriminantne analize tudi na primer, ko je matrika S_2 singularna.

2 Matematična priprava - posplošeni singularni razcep

Originalna definicija posplošenega singularnega razcepa (Van Loan)

Izrek 1 (Posplošeni singularni izrek (Van Loan):). *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 $1 \leq i \leq q$.

Dokaz. Iz matrik K_A in K_B tvorimo zdrženo matriko $K = \begin{pmatrix} K_A \\ K_B \end{pmatrix}$, na kateri naredimo singularni razcep. Iz singularnega razcepa dobimo matriki $Q \in \mathbb{R}^{(p+n) \times (p+n)}$ in matriko $Z_1 \in \mathbb{R}^{m \times m}$, tako da velja

$$Q^T \begin{pmatrix} K_A \\ K_B \end{pmatrix} Z_1 = \text{diag}(\gamma_1, \dots, \gamma_m) \quad \text{kjer za velja } \gamma_1 \geq \dots \geq \gamma_k > \gamma_{k+1} = \dots = \gamma_m. \quad (1)$$

– **DOPOLNI! Ali je samo diag za napisat ali celo razširjeno matriko?**

V kolikor matriko Z_1 razdelimo na dve matriki: $Z_{11} \in \mathbb{R}^{m \times k}$, ki je sestavljena iz prvih k stolpcev matrike Z_1 in $Z_{12} \in \mathbb{R}^{(m-k) \times n}$, ki je sestavljena iz preostalih $m - k$ stolpcev matrike Z_1 , lahko vidimo da velja:

$$Q^T K \begin{pmatrix} Z_{11} \\ Z_{12} \end{pmatrix} = \begin{pmatrix} Z_{11} \\ Z_{12} \end{pmatrix}.$$

Po predpostavki velja $p \geq m$ in ker je očitno tudi $m \geq k$ sledi: $p \geq n \geq k$. Sedaj definirajmo matriko

$$D := \text{diag}(\gamma_1, \dots, \gamma_k \in \mathbb{R}^{k \times k}).$$

Tako iz zgornje enačbe (??) dobimo:

$$\begin{pmatrix} AZ_{11} & AZ_{12} \\ BZ_{11} & BZ_{12} \end{pmatrix} = Q \begin{pmatrix} D & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix},$$

iz česar sledi:

$$\begin{pmatrix} AZ_{11} \\ BZ_{11} \end{pmatrix} = Q \begin{pmatrix} D \\ 0 \end{pmatrix}$$

in v kolikor še matriko Q razdelimo na podmatrike na naslednji način:

$$Q = \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{pmatrix},$$

kjer je matrika $Q_{11} \in \mathbb{R}^{k \times k}$, matrika $Q_{12} \in \mathbb{R}^{k \times (p+n-k)}$, matrika $Q_{21} \in \mathbb{R}^{(p+n-k) \times k}$ in matrika matrika $Q_{22} \in \mathbb{R}^{(p+n-k) \times (p+n-k)}$, ugotovimo, da je:

$$Q \begin{pmatrix} D \\ 0 \end{pmatrix} = \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{pmatrix} \begin{pmatrix} D \\ 0 \end{pmatrix} = \begin{pmatrix} Q_{11}D \\ Q_{21}D \end{pmatrix}.$$

Iz tega neposredno sledi enakost:

$$AZ_{11} = Q_{11}D \implies AZ_{11}D^{-1} = Q_{11} =: A_1,$$

kjer dodatno definiramo matriko A_1 in enakost:

$$BZ_{11} = Q_{21}D \implies BZ_{11}D^{-1} = Q_{21} =: B_1,$$

kjer dodatno definiramo matriko B_1 .

Ker je matrika Q ortogonalna dodatno velja $A_1^T A_1 + B_1^T B_1 = I_k$, kjer je I_k identična matrika dimzije $k \times k$. To enačbo lahko dobimo tako, da razpišemo spodnjo enačbo:

$$\begin{aligned} Q^T Q &= \begin{pmatrix} Q_{11}^T & Q_{21}^T \\ Q_{12}^T & Q_{22}^T \end{pmatrix} \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{pmatrix} = \begin{pmatrix} Q_{11}^T Q_{11} + Q_{21}^T Q_{21} & Q_{11}^T Q_{12} + Q_{21}^T Q_{22} \\ Q_{12}^T Q_{11} + Q_{22}^T Q_{21} & Q_{12}^T Q_{12} + Q_{22}^T Q_{22} \end{pmatrix} \\ &= \begin{pmatrix} A_1^T A_1 + B_1^T B_1 & Q_{11}^T Q_{12} + Q_{21}^T Q_{22} \\ Q_{12}^T Q_{11} + Q_{22}^T Q_{21} & Q_{12}^T Q_{12} + Q_{22}^T Q_{22} \end{pmatrix} = I = \begin{pmatrix} I_k & 0 \\ 0 & I_{p+n-k} \end{pmatrix}. \end{aligned}$$

□

Problem tega izreka ja, da se ga ne da uporabiti, kadar dimenzije matrike K_A niso ustrezne. Zaradi tega pretirano zavezujočega pogoja se odločita C.C. Paige in M.A. Saunders ta posplošeni singularni izrek še dodatno posplošiti. Tako dobimo naslednji izrek:

Izrek 2 (Posplošeni singularni izrek (Paige in Saunders):). *Naj bosta dani matriki $K_A \in \mathbb{R}^{p \times m}$ in $K_B \in \mathbb{R}^{n \times m}$. Potem za $K = \begin{pmatrix} K_A \\ K_B \end{pmatrix}$ in $t = \text{rang}(K)$ obstajajo ortogonalne matrike $U \in \mathbb{R}^{p \times p}$, $V \in \mathbb{R}^{n \times n}$, $W \in \mathbb{R}^{t \times t}$ in $Q \in \mathbb{R}^{m \times m}$, da velja:*

$$U^T K_A Q = \Sigma_A \begin{pmatrix} W^T R & 0 \end{pmatrix} \quad \text{in} \quad V^T K_B Q = \Sigma_B \begin{pmatrix} W^T R & 0 \end{pmatrix},$$

kjer je

$$\Sigma_A = \begin{pmatrix} I_A & & \\ & D_A & \\ & & 0_A \end{pmatrix} \quad \text{in} \quad \Sigma_B = \begin{pmatrix} 0_B & & \\ & D_B & \\ & & I_B \end{pmatrix}.$$

$R \in \mathbb{R}^{t \times t}$ je nesingularna matrika, matriki $I_A \in \mathbb{R}^{r \times r}$ in $I_B \in \mathbb{R}^{(t-r-s) \times (t-r-s)}$ identični matriki, kjer je

$$r = \text{rang}(K) - \text{rang}(K_B) \quad \text{in} \quad s = \text{rang}(K_A) + \text{rang}(K_B) - \text{rang}(K),$$

$0_A \in \mathbb{R}^{(p-r-s) \times (t-r-s)}$ in $0_B \in \mathbb{R}^{(n-t+r) \times r}$ ničelni matriki, ki imata lahko tudi ničelno število vrstic ali stolpcev, matriki $D_A = \text{diag}(\alpha_{r+1}, \dots, \alpha_{r+s})$ in $D_B = \text{diag}(\beta_{r+1}, \dots, \beta_{r+s})$ pa diagonalni matriki, ki zadoščata pogoju:

$$1 > \alpha_{r+1} \geq \dots \geq \alpha_{r+s} > 0 \quad \text{in} \quad 0 < \beta_{r+1} \leq \dots \leq \beta_{r+s} < 1$$

pri $\alpha_i^2 + \beta_i^2 = 1$ za $i = r+1, \dots, r+s$

Dokaz. Dovolj je, če ta izrek dokažemo za vsa kompleksna števila. Iz dejstva, da je množica realnih števil (\mathbb{R}) podmnožica množice kompleksnih števil (\mathbb{C}), sledi, da potem ta izrek velja tudi za vsa realna števila. Definirajmo matriko K , ki je sestavljena kot matrika sestavljena iz K_A in K_B , torej

$$K := \begin{pmatrix} K_A \\ K_B \end{pmatrix}.$$

Na zgoraj definirani matriki K lahko sedaj naredimo singularni razcep. Tako vemo, da za matriko K obstajata unitarni matriki $P \in \mathbb{C}^{(m+p) \times (m+p)}$ in $Q \in \mathbb{C}^{n \times n}$, da velja

$$P^H K Q = \begin{pmatrix} R & 0 \\ 0 & 0 \end{pmatrix}$$

, kjer ima matrika R enak rang kot matrika K . Matriki Q in P sedaj ločimo na sledeče podmatrike:

$$Q = \begin{pmatrix} Q_1 & Q_2 \end{pmatrix} \quad \text{in} \quad P = \begin{pmatrix} P_1 & P_2 \end{pmatrix} = \begin{pmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{pmatrix},$$

kjer je matrika $Q_1 \in \mathbb{C}^{m \times t}$ sestavljena iz prvih k stolpcev matrike Q , matrika $P_1 \in \mathbb{C}^{(p+n) \times t}$ pa iz prvih t stolpcev matrike P in njena podmatrika $P_{11} \in \mathbb{C}^{p \times t}$ pa iz prvih m vrstic matrike P_1 . Vemo, da ker je matrika P unitarna matrika, velja $\|P\|_2 \leq 1$ in posledično velja še $\|P_{11}\|_2 \leq \|P_1\|_2 \leq \|P\|_2 \leq 1$. Iz izreka iz numeričnih metod velja, da posledično nobena lastna vrednost matrike P_{11} ni večja od 1.

Singularni razcep podobno kot na matriki K naredimo tudi na matriki P_{11} . Tako dobimo takšni matriki $U \in \mathbb{C}^{p \times p}$ in $W \in \mathbb{C}^{t \times t}$, da velja

$$U^H P_{11} W = \Sigma_A,$$

kjer je

$$\Sigma_A = \begin{pmatrix} 0_B & & \\ & D_B & \\ & & I_B \end{pmatrix},$$

kjer je matrika D_B diagonalna matrika z diagonalnimi vrednostmi $\alpha_{r+1}, \dots, \alpha_{r+s}$, za katere velja $1 > \alpha_{r+1} \geq \dots \geq \alpha_{r+s} > 0$.

(Spodnji del je vprašljiv? – **DOPOLNI!**)

Na matriki P_{21} uporabimo čudežni razcep (s Householderjevimi zrcaljenji) in tako dobimo matriko $V \in \mathbb{C}^{n \times n}$, da velja

$$V^H P_{21} W = L = (\ell_{ij})_{i,j} = \begin{pmatrix} 0 & \\ & L_1 \end{pmatrix},$$

kjer je matrika L_1 spodnjetrokotna z diagonalnimi elementi večjimi od 0. Opazimo lahko, da velja spodnja enakost

$$\begin{pmatrix} U^T & 0 \\ 0 & V^T \end{pmatrix} \begin{bmatrix} P_{11} \\ P_{21} \end{bmatrix} W = \begin{bmatrix} U^H P_{11} W \\ V^H P_{21} W \end{bmatrix} = \begin{pmatrix} \Sigma_A \\ L \end{pmatrix}.$$

Zgornja matrika $\begin{pmatrix} \Sigma_A \\ L \end{pmatrix}$ je unitarna, saj je produkt unitarnih matrik. Posledično so njeni stolpci ortonormirani. \square

3 Matematična rešitev problema

3.1 Posplošitev linearne diskriminantne analize

3.2 Posplošitev maksimizacijskega kriterija $\text{sled}((S_W^Y)^{-1} S_B^Y)$

4 Algoritem

5 Zaključek

6 Viri