

ROZHODOVÁNÍ S PODPOROU DISKRIMINAČNÍ ANALÝZY

Bohdan Linda

Ústav matematiky, FES, Univerzita Pardubice

Abstract:

The discrimination analysis is a multivariable statistical method with the aim to find the optimal assignment regulation and prediction, into which of k groups a watched element belongs to. Several characteristics are measured at every element. They express its features. It means, that this element is characterised by means of a random variable $\mathbf{X}=(X_1, X_2, \dots, X_m)$. The final result of this procedure is a model, that allows to predict the pertinence of the watched element to the certain group on the basis of measured values (x_1, x_2, \dots, x_m) .

Ve veřejné správě často potřebujeme seskupit nějaké objekty, což mohou být například obce, okresy regiony apod., podle některých znaků, které vzájemně mohou být nesouměřitelné, ale spoluvytvářejí nějakou vlastnost, kterou nelze přímo měřit – například ekonomickou úroveň regionu, výkonnost a další. Podobně tomu je i v ekonomické praxi, kde mnohdy potřebujeme jednotlivé podniky rozdělit do několika skupin podle přímo neměřitelného znaku. A právě zde může diskriminační analýza usnadnit členění, napomoci k přijetí správného rozhodnutí při zařazení sledovaného regionu, podniku či jiného sledovaného objektu.

Objekty, které mají některé podobné vlastnosti, můžeme z důvodu lepšího sledování jejich vlastností sloučovat do skupin. Nebývá ale vždy možné jednoduchým způsobem a jednoznačně o nějakém prvku rozhodnout, do které skupiny patří. Jednou z metod, která může v rozhodování pomoci je diskriminační analýza. Hlavním cílem diskriminační analýzy je nalézt statisticky nejvýhodnější způsob rozlišení skupin prvků a předpovědět, do které z k skupin patří sledovaný prvek. U každého prvku je měřeno několik znaků, vyjadřujících jeho vlastnosti, to znamená, že tento prvek je charakterizován prostřednictvím náhodné veličiny $\mathbf{X}=(X_1, X_2, \dots, X_m)$. Postup začíná analýzou vlastností prvků, u nichž je známá jak příslušnost ke skupině, tak i hodnoty příslušné náhodné veličiny. Konečným výsledkem postupu je model, který umožňuje předpovědět příslušnost sledovaného prvku k určité skupině na základě naměřených hodnot (x_1, x_2, \dots, x_m) .

Diskriminační analýza se s úspěchem aplikuje v řadě oborů. Jedno z prvních použití bylo v archeologii. Při nálezech hrobů s kosterními pozůstatky jsou nalézány také kultovní předměty. Na základě určitých charakteristických vlastností lze pak nález přiřadit k určitému období, kultuře, rase. Ve školské i v personální podnikové praxi se v přijímacím řízení uplatňují soubory testů, jejichž výsledkem je bodové hodnocení uchazeče. Po uplynutí určitého období lze získat informaci, jak úspěšný či neúspěšný je sledovaný jedinec. Disponujeme tedy informativním výběrem a pokud volba testů je taková, že jejich výsledky souvisí s pozorovanou úspěšností, lze přistoupit ke klasifikaci. To znamená, že se snažíme provést s malým rizikem omylu předpověď úspěchu podle výsledku testů. Nejjednodušším případem diskriminační analýzy je předpověď příslušnosti k dichotomicky členěné množině prvků založené na jednorozměrných proměnných.

Jedním ze zásadních problémů metody je přesnost předpovědi. Je-li kritériem zařazení prvku do skupiny jen jednoduchá proměnná, pak postup diskriminační analýzy umožňuje rozlišit zařazení prvku do skupiny někdy téměř dokonale, někdy částečně a někdy to

neumožouje vůbec. Větší rozdíl mezi středními hodnotami sledované veličiny ve skupinách umožňuje lepší rozlišení skupin. V jednoduchém případě dichotomických skupin zpravidla není problémem různá variabilita ve skupinách. V případě většího počtu skupin musíme ale předpokládat rovnost rozptylů ve skupinách.

Obecný princip metody

Úkolem diskriminační analýzy je nalezení optimálního prázovacího pravidla, tzn. pravidla, které minimalizuje pravděpodobnost chybné klasifikace, tedy minimalizuje střední hodnotu chyby rozhodnutí. Může se totiž stát, že prvek, který skutečně pochází z určité skupiny zařadíme do jiné skupiny.

Předpokládejme, že je dáno k skupin prvků ($k = 2, 3, \dots$) a každý zprvků je charakterizován pomocí náhodné veličiny $X = (X_1, X_2, X_3, \dots, X_m)$ se známým typem rozdělení pravděpodobností. To znamená, že prvky patřící do i -té skupiny můžeme považovat za náhodný výběr ze základního souboru s rozdělením $f_i(x)$ se střední hodnotou m_i a variancí maticí S_i .

Cílem diskriminační analýzy je zjistit, do které z k skupin patří sledovaný prvek.

Předpokládejme, že náhodná veličina X nabývá hodnot z R_m (m -rozměrného reálného vektorového prostoru). Utvoříme rozklad prostoru R_m na k množin M_1, \dots, M_k (tj. pro množiny M_i musí platit: $R_m = \bigcup_{i=1}^k M_i$; $M_i \cap M_j = \emptyset$; $i \neq j$). Padne-li hodnota veličiny X do množiny M_i ,

tvrdíme, že sledovaný prvek náleží do i -té skupiny. Problémem je nalézt takový rozklad, aby rozhodnutí o příslušnosti k dané třídě bylo optimální. Abychom mohli rozhodnout o tom, zdali je rozklad optimální, potřebujeme mít k dispozici nějakou kritériální funkci. Takovým kritériem může být například veličina Z , představující ztrátu, která vznikne chybným zařazením prvků v důsledku nesprávného rozkladu R_m na podmnožiny M_i . Protože zařazení prvku chceme provádět na základě konkrétní realizace náhodného vektoru X , je logické, že veličiny Z by měla být kromě M_i i funkcí X , jinými slovy řečeno, Z bude náhodnou veličinou. Proto při určování optimálního rozkladu budeme pracovat s její střední hodnotou. Za optimální budeme považovat takový rozklad, který minimalizuje střední hodnotu ztráty Z .

Základem kritéria Z bude ztráta z_{ij} , která vznikne chybným zařazením prvku i -té skupiny do j -té a kterou ve většině konkrétních situací dokážeme určit. Nejdříve určíme střední hodnotu náhodné veličiny Z_i , představující ztrátu, když prvek i -té skupiny zařadíme nesprávně (tj. do kterékoliv skupiny $j = 1, 2, \dots, m, j \neq i$).

Platí pro ni:

$$EZ_i = z_{i1} \int_{M_1} f_i(x_1, \dots, x_m) dx_1 \dots dx_m + \dots + z_{ik} \int_{M_k} f_i(x_1, \dots, x_m) dx_1 \dots dx_m \quad i = 1, 2, \dots, k.$$

Funkce EZ_i je podmíněná střední hodnota ztráty, podmínkou je, že objekt přísluší i -té skupině. Nepodmíněná střední hodnota ztráty EZ (tj. střední hodnota ztráty bez ohledu na to, do které skupiny prvek patří) je pak

$$EZ = p_1 Z_1 + p_2 Z_2 + \dots + p_k Z_k,$$

kde p_i je pravděpodobnost, že objekt patří do i -té skupiny.

$$\text{Položíme-li } h_j(x_1, \dots, x_m) = \sum_{i=1}^k p_i z_{ij} f_i(x_1, \dots, x_m),$$

pak střední hodnotu ztráty můžeme přepsat ve tvaru

$$EZ = \sum_{i=1}^k \int_{M_i} h_i(x_1, \dots, x_m) dx_1 \dots dx_m$$

K nalezení optimálního rozkladu můžeme použít následující vtu:

VĚ TA: Jestliže existuje rozklad $M_1^0, M_2^0, \dots, M_k^0$ takový, že pro libovolné $(x_1, x_2, \dots, x_m) \in M_s^0$ platí $h_s(x_1, x_2, \dots, x_m) \leq h_j(x_1, x_2, \dots, x_m)$, pak tomuto rozkladu odpovídá nejmenší střední hodnota ztráty EZ^0 .

Obvykle se v diskriminační analýze volí ztráty z_{ij} zjednodušeným způsobem takto:

$$\begin{aligned} z_{ij} &= 1 \text{ pro } i, j = 1, 2, \dots, m \text{ a } i \neq j, \\ z_{ii} &= 0 \text{ pro } i = 1, 2, \dots, m \end{aligned}$$

V takovém pápadě minimalizovat EZ znamená minimalizovat počet chybně zařazených objektů. Výraz $g_j(x_1, \dots, x_m)$ můžeme vyjádřit jako

$$g_j(x_1, \dots, x_m) = \sum_{i=1}^k p_i f_i(x_1, \dots, x_m) - p_j f_j(x_1, \dots, x_m),$$

Pro danou hodnotu $\mathbf{x} = (x_1, \dots, x_m)$ je potom $g_s(x_1, \dots, x_m) \leq g_j(x_1, \dots, x_m)$, $j = 1, 2, \dots, k$, právě když je

$$p_s f_s(x_1, \dots, x_m) \geq p_j f_j(x_1, \dots, x_m).$$

Je-li u zkoumaného prvku zjištěna hodnota $\mathbf{x} = (x_1, \dots, x_m)$, je prvek zařazen do skupiny t , pro niž je splněna nerovnost

$$p_t f_t(x_1, \dots, x_m) > p_j f_j(x_1, \dots, x_m), \quad (j = 1, 2, \dots, k; j \neq t) \quad (1)$$

V praxi obvykle předpokládáme, že náhodná veličina \mathbf{X} má m -rozměrné normální rozložení pravděpodobností se známým vektorem středních hodnot a se známou varianční maticí. To znamená, že v jednotlivých skupinách M_j , $j = 1, 2, \dots, k$ bude mít vektor \mathbf{X} hustotu pravděpodobnosti

$$f_j(x_1, \dots, x_m) = (2\pi)^{-\frac{1}{2}m} \cdot [\det(\mathbf{\Sigma}_j)]^{-\frac{1}{2}} \cdot e^{-\frac{1}{2}(x_1 - \mu_j)' \mathbf{\Sigma}_j^{-1} (x_1 - \mu_j)}$$

Nerovnost (1) je ekvivalentní s nerovností

$$\ln p_s + \ln f_s(x_1, \dots, x_m) > \ln p_j + \ln f_j(x_1, \dots, x_m) \quad ; \quad j \neq s ; j = 1, 2, \dots, m$$

Označme-li

$$D_j = -0,5 \ln [\det(\mathbf{S}_j)] - 0,5 (\mathbf{x} - \mathbf{m}_j)' \mathbf{S}_j^{-1} (\mathbf{x} - \mathbf{m}_j) + \ln p_j,$$

bude vztah (1) platit, právě tehdy, když $D_t > D_j$, tzn.

$$-\frac{1}{2} \ln[\det(\mathbf{S}_j)] - \frac{1}{2} (\mathbf{x} - \mathbf{m}_t)' \mathbf{S}_t^{-1} (\mathbf{x} - \mathbf{m}_t) + \ln p_t > -\frac{1}{2} \ln[\det(\mathbf{S}_j)] - \frac{1}{2} (\mathbf{x} - \mathbf{m}_j)' \mathbf{S}_j^{-1} (\mathbf{x} - \mathbf{m}_j) + \ln p_j$$

Při praktickém provádění diskriminační analýzy vypočítáme pro všechny hodnoty náhodné veličiny \mathbf{X} hodnoty D_1, \dots, D_k . Zkoumaný prvek přísluší té skupině, která odpovídá největší z hodnot D_j . Hodnoty \mathbf{m}_j a \mathbf{S}_j zpravidla neznáme, ale použijeme jejich odhady.

Pravděpodobnost p_j se zpravidla volí úměrná rozsahu j -té skupiny, pokud tyto rozsahy nejsou známy, volí se $p_j = 1/k$.

Diskriminační analýza zatím v ekonomické a veřejno-správní praxi u nás nenalezla širší uplatnění. Je to dáno jednak tím, že její teoretické pozadí je vzhledem k jiným statistickým disciplinám poměrně složité a také tím, že výpočty jsou náročné. Díky výpočetní technice však v dnešní době tyto překážky ztrácejí na významu.

Literatura:

1. Andìl,J.: *Matematická statistika*. SNTL, Praha 1978
2. Antoch,J.,Vorlířková,D.: *Vybrané metody statistické analýzy dat*. Akademia Praha 1992
3. Bock,H.H.:*Automatische Klassifikation*. In:Walter E.a kol.: *Statistische Methoden II*. Berlin, Springer 1970
4. Bowerman,B.,O'Connell,R.,T.: *Applied Statistics*. Irwin, USA 1997

Kontaktní adresa:

doc.RNDr. Bohdan Linda, CSc.
Ústav matematiky, FES, Univerzita Pardubice
Studentská 84
532 10 Pardubice
tel.: 466 036 020
e-mail: Bohdan.Linda@upce.cz