

УДК 519.237.8

М. А. Іванчук, асистент

Буковинський державний медичний університет, м. Чернівці

КЛАСИФІКАЦІЯ МНОЖИН МЕТОДОМ ЛІНІЙНОГО ВІДОКРЕМЛЕННЯ ЇХ ОПУКЛИХ ОБОЛОНОК

У статті представлений метод лінійного відокремлення опуклих оболонок (ЛВОО) для класифікації двох множин в евклідовому просторі R^n . Наводяться приклади для порівняння результату класифікації методом ЛВОО, а також дискримінантним аналізом, наївним байєсівським класифікатором та методом опорних векторів.

Ключові слова: лінійна класифікація, опуклі оболонки, відокремлююча гіперплощина.

Вступ. Розглянемо наступну задачу класифікації. Нехай задано дві групи об'єктів A та B . Множина A складається з m_A об'єктів, множина B — з m_B об'єктів. Для кожного з об'єктів цих множин відомо n кількісних ознак, впорядкованих за інформативністю, наприклад, за інформативною мірою Кульбака [1, с. 93–120]. Задача полягає в знаходженні такого алгоритму класифікації, який би на заданому рівні значущості вірно розподіляв об'єкти за відповідними групами, використовуючи при цьому найменшу кількість ознак. Прикладом подібної задачі на практиці є задача медичного прогнозування: існує дві групи хворих (A — хворі з ускладненнями та B — хворі без ускладнень). Для кожного з хворих відомо n ознак — факторів, що впливають на розвиток ускладнень. Необхідно, використовуючи найменшу кількість ознак, побудувати експертну систему, що з заданим рівнем значущості прогнозуватиме у хворих виникнення ускладнень.

1. Постановка задачі. Розглянемо математичну модель даної задачі. Будемо представляти об'єкти з множин A та B як точки у евклідовому просторі R^n . Кожна координата відповідатиме за одну з n ознак, що описують дані об'єкти. Тоді множини A та B можна записати як дві множини точок $A = \{a_i = (a_i^1, a_i^2, \dots, a_i^n), i = \overline{1, m_A}\}$ та $B = \{b_i = (b_i^1, b_i^2, \dots, b_i^n), i = \overline{1, m_B}\}$. Виділимо з них випадковим чином множину точок контрольної групи $Z = \{z_i = (z_i^1, z_i^2, \dots, z_i^n), i = \overline{1, m_Z}\}$ для перевірки якості класифікації. У випадку малих значень m_A та m_B можна використати метод ковзаючого контролю [2, с. 266–268].

Позначимо через α наперед заданий допустимий рівень помилок при розподіленні елементів множин A та B за відповідними класами, θ — наперед заданий допустимий рівень помилок при перевірці класифікаційної моделі на контрольній групі Z .

Оскільки задача полягає в знаходженні класифікатора, що використовує найменшу кількість ознак, будемо проводити класифікацію множин в просторах розмірностей $k = \overline{1, n}$. Серед класифікаторів, що задовольняють наперед задані допустимі рівні помилок α та θ , вибиратимемо той, що відповідає простору меншої розмірності.

Означення. Гіперплощину

$$L_{i^k} = \{x \in R^k : \langle p, x \rangle = \gamma\}, \quad p \neq 0, \quad k = \overline{1, n}$$

будемо називати *відокремлюючою гіперплощиною* для множин A та B в просторі R^k , якщо не менше ніж $(1 - \alpha)(m_A + m_B)$ точок множин A та B можна помістити в різні півпростори $L_{i^k A} = \{x \in R^k : \langle p, x \rangle > \gamma\}$ та $L_{i^k B} = \{x \in R^k : \langle p, x \rangle < \gamma\}$.

Позначимо Z_A^k — точки контрольної множини, що відносяться до групи A^k , відповідно Z_B^k — точки контрольної множини, що відносяться до групи B^k , $Z^k = Z_A^k \cup Z_B^k$. Для кожної з знайдених відокремлюючих гіперплощин $L_{i^k}, i^k = \overline{1, l^k}$ будемо розпізнавати точки контрольної множини Z^k відповідно до їх розташування відносно цієї гіперплощини. Позначимо $Z_{i^k A}^+ = \{z \in Z_A^k : z \in L_{i^k A}\}$ — множина вірно розпізнаних точок Z_A^k , $Z_{i^k B}^+ = \{z \in Z_B^k : z \in L_{i^k B}\}$ — множина вірно розпізнаних точок Z_B^k , $Z_{i^k A}^- = \{z \in Z_A^k : z \in L_{i^k B}\}$ — множина невірно розпізнаних точок Z_A^k , $Z_{i^k B}^- = \{z \in Z_B^k : z \in L_{i^k A}\}$ — множина невірно розпізнаних точок Z_B^k .

Тоді відносна помилка відокремлюючої гіперплощини L_{i^k} буде

$$\theta_{i^k} = \frac{m_{Z_{i^k A}^-} + m_{Z_{i^k B}^-}}{m_{Z^k}}.$$

Означення. Нехай $\exists i^k : \theta_{i^k} < \theta$, тоді *оптимальною відокремлюючою гіперплощиною в просторі R^k* називатимемо відокремлюючу гіперплощину $L_{i^k_{\min}}$ таку, що $i^k_{\min} = \arg \min_{i^k} \{\theta_{i^k} < \theta\}$.

Означення. *Оптимальним класифікатором* для множин A та B на заданому рівні помилок θ називатимемо оптимальну відокремлюючу гіперплощину, що відповідає простору найменшої розмірності.

2. Пошук відокремлюючих гіперплощин. Для знаходження відокремлюючих гіперплощин $L_{i,k}$ в просторі R^k , $k = \overline{1, n}$ пропонуємо наступну методику.

Побудуємо методом Джарвіса [3, с. 114–183] однозв'язні опуклі оболонки conv_{A^k} та conv_{B^k} множин A^k та B^k . При цьому можливі наступні варіанти взаємного розташування множин та їх опуклих оболонок:

1 випадок. $\text{conv}_{A^k} \cap \text{conv}_{B^k} = \emptyset$. Тобто опуклі оболонки множин не перетинаються (рис. 1).

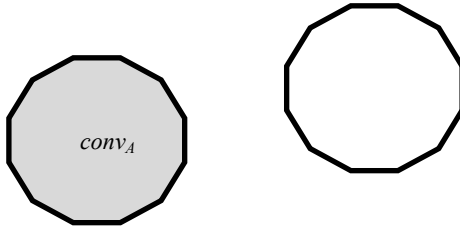


Рис. 1. 1 випадок взаємного розташування множин

Виходячи з наслідку з теореми Хана-Банаха [4, с. 139–150] та теореми про відокремлюючу вісь [5] сформулюємо наступне твердження [6]:

Твердження. Нехай задано два опуклих багатогранники conv_A та conv_B , що не перетинаються та їх найближчі точки $a_{\min} \in \text{conv}_A$ та $b_{\min} \in \text{conv}_B$ такі, що

$$|\overline{a_{\min} b_{\min}}| = \min \left\{ |\overline{a_i b_j}| : a_i \in \text{conv}_A, b_j \in \text{conv}_B, i = \overline{1, m_A}, j = \overline{1, m_B} \right\}.$$

Серед гіперграней, що містять ці точки знайдеться хоча б одна така, паралельно якій через точку на відрізку $\overline{a_{\min} b_{\min}}$ можна провести відокремлюючу гіперплощину L , що розділяє conv_A та conv_B .

Нехай g_k — одна з гіперграней, паралельно якій, згідно твердження, можна провести відокремлюючу гіперплощину. Серед множини гіперплощин, паралельних гіперграні g_k в якості відокремлюючої гіперплощини для множин A та B прийматимемо ту, що знаходиться на однаковій відстані від найближчих точок опуклих оболонок a_{\min} та b_{\min} .

2 випадок. $conv_{A^k} \cap conv_{B^k} \neq \emptyset$, але кількість точок множини $O_{A^k} = \{o^k : o^k \in A^k \cap conv_{B^k}\}$ задовольняє нерівність $m_{O_{A^k}} \leq \alpha \cdot m_{A^k}$ (рис. 2)

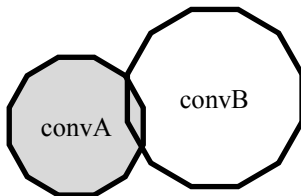


Рис. 2. 2 випадок взаємного розташування множин

Означення. Точку a_i^k називатимемо *промахом* множини A^k , якщо $a_i^k \in conv_{B^k}$.

Означення. Множину $O_{A^k} = \{o^k \in A^k : o^k \in conv_{B^k}\}$ будемо називати *множиною промахів* A^k .

Відкинемо промахи з тієї множини, у якої їх менше та побудуємо для неї нову опуклу оболонку. Будемо нові опуклі оболонки та відкидаємо з них промахи до тих пір, поки не виконується умова

$$conv_{A^k} \cap conv_{B^k} = \emptyset$$

Коли ця умова виконується, $conv_{A^k}$ та $conv_{B^k}$ — це дві замкнені обмежені опуклі множини, що не перетинаються. Отримуємо випадок, аналогічний попередньому.

3 випадок. $A^k \notin conv_B$ (рис.3а), але $conv_{A^k} \cap conv_{B^k} \neq \emptyset$ (рис. 3б)

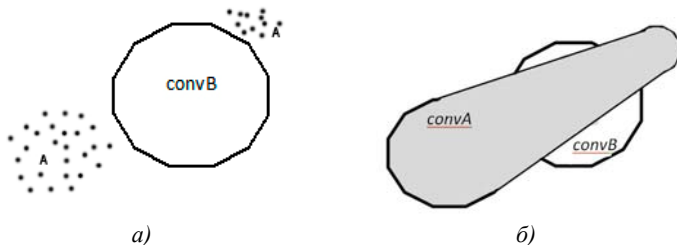


Рис. 3. 3 випадок взаємного розташування множин

Тобто більшість точок множини A^k розташовані по один бік множини B^k . При цьому деякі точки множини A^k розташовані по інший бік множини B^k . Якщо їх кількість знаходиться в межах заданого рівня помилок α , вважатимемо їх неявними промахами та відкидатимемо.

Означення. Пару точок a_i та a_j називатимемо *підозрілою на неявний промах* для множини A^k , якщо $a_i, a_j \notin \text{conv}_{B^k}$, але існує гіпергрань $q \in \text{conv}_{B^k}$ така, що точки a_i та a_j лежать в різних півпросторах, утворених гіперплощиною q (рис. 4).

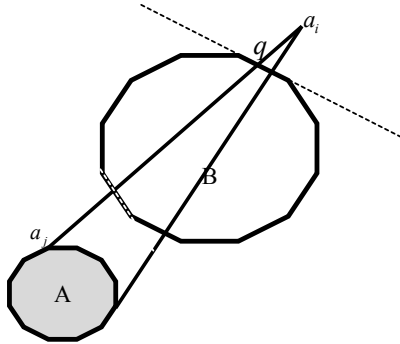


Рис. 4. Пара точок, підозрілих на неявний промах

Означення. Точка a_i є *неявним промахом* для множини A^k , якщо вона з не менше ніж $(1 - \alpha) \cdot m_A$ точками множини A^k складає пару, підозрілу на неявний промах.

Множину неявних промахів позначатимемо P_{A^k} .

Після відкидання неявних промахів та промахів множин приходимо до першого випадку.

3. Складність алгоритму методу ЛВОО. В усіх описаних випадках взаємного розташування множин найбільшу складність має процедура побудови опуклих оболонок. Якщо опуклі оболонки будувати методом Джарвіса, складність алгоритму для кожного k в найгіршому випадку дорівнює $O(m^2)$ [3, с. 135]. Отже, загальну складність алгоритму методу ЛВОО можна оцінити як $O(nm^2)$, якщо використані всі ознаки. Якщо оптимальний класифікатор отримується при використанні невеликої кількості ознак ($k \ll m$), складність алгоритму методу ЛВОО можна оцінити як $O(m^2)$.

Зауважимо, що складність алгоритму лінійного дискримінантного аналізу оцінюється як $O(mnt + t^3)$, де $t = \min(m, n)$ [7]; найвигіднішого байєсівського класифікатора — $O(mn)$ [8]; методу опорних векторів — $O(m^3)$ [9].

Приклади застосування алгоритму ЛВОО. Описаний вище алгоритм ЛВОО був реалізований нами в програмі MATLAB R2013b. Наведемо приклади його застосування.

Приклад 1. Розглянемо відому базу даних іриси Фішера [10]. Є дані про 150 рослин трьох видів іриси (*setosa*, *versicolor*, *virginica*). Для кожної рослини відомо 4 ознаки — довжина і ширина чашолистика та довжина і ширина пелюстки. Необхідно за цими даними побудувати алгоритм класифікації. Ця база даних є цікавою тим, що перший клас (*setosa*) лінійно відокремлюється від двох інших, а класи *versicolor* та *virginica* між собою лінійно не відокремлюються.

Класифікацію проводили за допомогою дискримінантного аналізу [10], наївного байєсівського класифікатора [11, с. 47-48], методом опорних векторів [12, с.156-163] та за допомогою методу ЛВОО (рис. 5). Оскільки об'єми груп невеликі — кожна група складається з 50 рослин, контроль за якістю класифікації проводили методом ковзаючого контролю [2, с. 266-268].

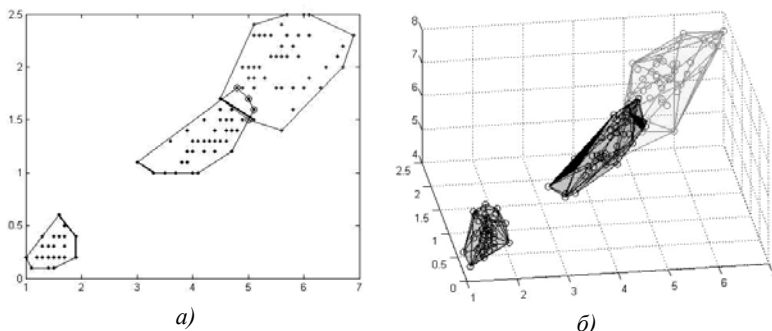


Рис. 5. Застосування методу ЛВОО для класифікації ірисів Фішера а) $n = 2$, б) $n = 3$

Сумарна кількість помилок в кожному класі, отриманих при застосування різних методів класифікації представлені в таблиці 1.

Таблиця 1

Сумарна кількість помилок класифікації

Метод	Setosa	Versicolor	Virginica
Дискримінантний аналіз	0	2	1
Наївний байєсівський класифікатор	0	3	3
Метод опорних векторів (лінійне ядро)	0	4	0
Метод опорних векторів (ядро — гаусівська радіальна базисна функція)	0	1	0
Метод ЛВОО	0	1	0

Приклад 2. Для хворих на гострий деструктивний панкреатит проводили класифікацію між групами хворих з післяопераційними

ускладненнями та без них. Для класифікації використовували дані про 89 хворих (група A — 28 хворих з ускладненнями, група B — 61 хворий без ускладнень).

При класифікації методом ЛВОО контроль за якістю класифікації проводили методом ковзаючого контролю. Оскільки практична мета даної класифікаційної моделі — не пропустити хворих, у яких можуть виникнути післяопераційні ускладнення, тобто зменшити кількість помилок гіподіагностики, то за допустимий рівень помилок прийняли відсоток помилок класифікації множини A : $\theta_A = 0,05$. Заданий рівень помилок був досягнений при використанні $n = 3$ ознак (рис.6). Класифікацію іншими методами відповідно проводили також за трьома ознаками.

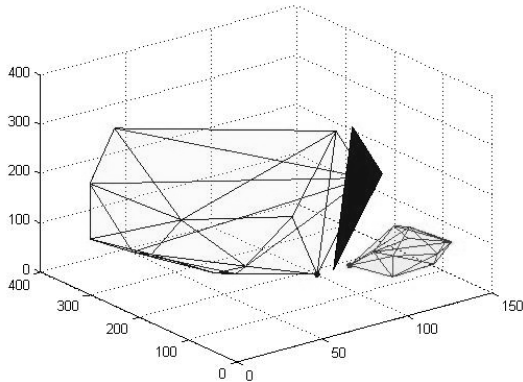


Рис. 6. Застосування методу ЛВОО для прогнозування виникнення ускладнень у хворих на гострий деструктивний панкреатит ($n = 3$)

Сумарна кількість помилок в кожному класі, отриманих при застосування різних методів класифікації представлені в таблиці 2.

Таблиця 2

Сумарна кількість помилок класифікації

Метод	Хворі з ускладненнями	Хворі без ускладнень
Дискримінантний аналіз	1	4
Наївний байєсівський класифікатор	2	2
Метод опорних векторів (лінійне ядро)	1	4
Метод опорних векторів (ядро — гаусівська радіальна базисна функція)	1	3
Метод ЛВОО	1	4

Висновок. В роботі запропонований метод лінійної класифікації двох однозв'язних множин. Суть методу полягає в знаходженні відокре-

млюючої гіперплощини, що паралельна одній з гіперграней опуклих оболонок множин, що класифікуються. За оптимальну відокремлюючу гіперплощину вибирається гіперплощина з найменшою кількістю помилок в контрольній групі. В якості класифікатора вибирається оптимальна відокремлююча гіперплощина, що відповідає простору меншої розмірності. Складність методу ЛВОО оцінюється як $O(nm^2)$.

Список використаних джерел:

1. Кульбак С. Теория информации и статистика / С. Кульбак. — М. : Наука, 1967. — 408 с.
2. Вапник В. Н. Восстановление зависимостей по эмпирическим данным / В. Н. Вапник. — М. : Наука, 1979. — 448 с.
3. Препарата Ф. Вычислительная геометрия: Введение / Ф. Препарата, М. Шеймос. — М. : Мир, 1989. — С. 478.
4. Колмогоров А. Н. Элементы теории функций и функционального анализа / А. Н. Колмогоров, С. В. Фомин. — 7-е изд. — М. : ФИЗМАТЛИТ, 2004. — 572 с.
5. SAT (Separating Axis Theorem) [Електронний ресурс]. — Режим доступу: <http://www.codezealot.org/archives/55>.
6. Ivanchuk M. Mathematical Modeling of the Expert System Predicting the Severity of Acute Pancreatitis / M. Ivanchuk, V. Maksimyuk, I. Malyk // Journal of Computational Medicine. — Vol. 2014. — Article ID 532453.
7. Deng Cai, Xiaofei He, Jiawei Han Training Linear Discriminant Analysis in Linear Time [Електронний ресурс]. — Режим доступу: http://researchweb.iit.ac.in/~nataraj.j/poseSearchReports/icde08_dengcai.pdf
8. Chris Fleizach, Satoru Fukushima A naive Bayes classifier on 1998 KDD Cup [Електронний ресурс]. — Режим доступу: <http://sysnet.ucsd.edu/~cfleizac/cse250b/project1.pdf>.
9. Vector Machines: Fast SVM Training on Very Large Data Sets / I. W. Tsang, J. T. Kwok, Pak-Ming Cheung // Journal of Machine Learning. — 2005. — Research 6. — P. 363–392. Submitted 12/04; Published 4/05
10. Fisher R. A. The Use of Multiple Measurements in Taxonomic Problems / R. A. Fisher // Annals of Eugenics 7. — 1936. — P. 179–188.
11. Прикладная статистика: классификация и снижение размерности / С. А. Айвазян, В. М. Бухштабер, И. С. Енюков, Л. Д. Мешалкин. — М. : Финансы и статистика, 1989. — 608 с.
12. Vapnik V. N. The nature of statistical learning theory / V. N. Vapnik. — 2nd ed. — New York, 2000. — 314 p.

The method of convex hulls linear separation for classification of two sets in Euclidian space in R^n is proposed. Examples for comparing the results of using the method of convex hulls linear separation and discriminant analysis, naive Bayes classifier and SVM are described in the manuscript.

Key words: *linear classification, convex hulls, separating hyperplane.*

Отримано: 14.04.2015