

Міністерство Освіти і Науки України
Київський Національний Університет імені Тараса Шевченка
Факультет Інформаційних Технологій
Кафедра Інформаційних систем та технологій

Звіт з лабораторної роботи № 3
з дисципліни «Data Science та машинне навчання»
Тема: « **МЕТОДИ РОЗДІЛЯЮЧИХ ФУНКЦІЙ** »

Виконав студент 1-го курсу
магістратури
групи ІРма-12
Гаврасієнко Є.О.

Київ – 2025
Мета роботи:

1. Опанувати непараметричні методи «навчання з учителем», засновані на

лінійних розділяючих функціях і методику побудови шматково-лінійних вирішальних правил;

2. Отримати навички статистичного оцінювання показників якості класифікації з використанням системи MathCAD для моделювання та подання об'єктів у вигляді даних спостережень.

Завдання:

1. За заданими (згідно з варіантом) двовимірними даними спостережень $\xi_i = (x_i, y_i)$ двох класів об'єктів a_1 і a_2 за правилом найближчого сусіда провести межі між класами:

a. за вибірковими значеннями – межу $g_1(x, y) = 0$;

b. за вибірковим середнім – межу $g_2(x, y) = 0$.

2. Побудувати вирішальні правила g_1 та g_2 .

3. Згенерувати масиви N даних спостережень ($N = 100$) класів a_1 і a_2 у припущенні, що спостерігається двовимірний випадковий вектор, компоненти якого – некорельовані нормально розподілені величини. В якості параметрів розподілу класів відповідно) взяти їх статистичні оцінки, отримані за заданими вихідними даними.

4. Змоделювати процеси розпізнавання спостережень згідно вирішуючих правил g_1 і g_2 і порівняти ефективності класифікаторів за емпіричними оцінками ймовірностей правильних рішень.

5. Оформити звіт про лабораторну роботу, який повинен містити короткі теоретичні відомості, алгоритми моделювання даних та прийняття рішень, графічні подання реалізацій спостережень та меж між класами, висновки.

Хід роботи

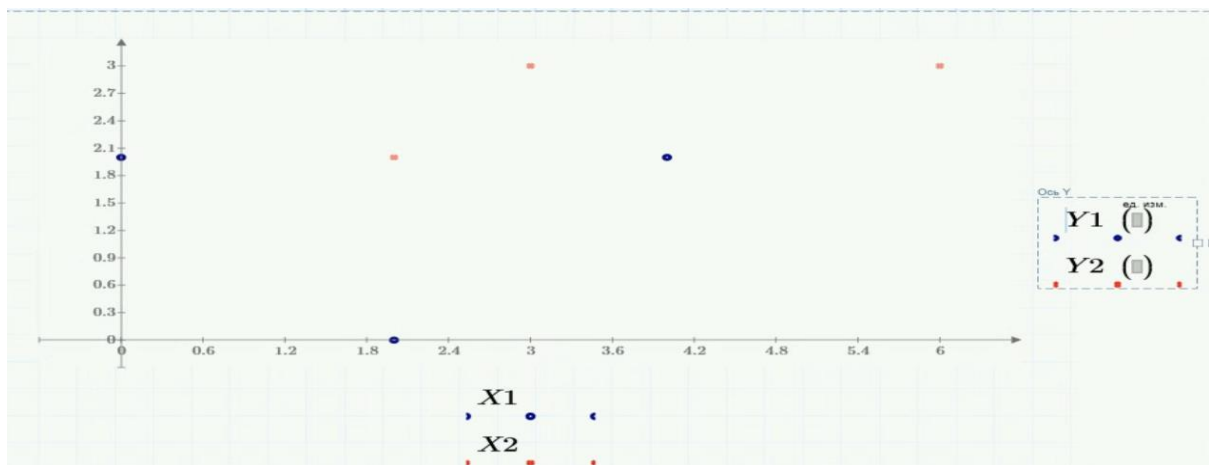
Варіант для виконня згідно таблиці - 2

Параметри для варіанту:

Вхідні дані			
Клас a_1		Клас a_2	
$x_{11} := 0$	$y_{11} := 2$	$x_{24} := 2$	$y_{24} := 2$
$x_{12} := 2$	$y_{12} := 0$	$x_{25} := 3$	$y_{25} := 3$
$x_{13} := 4$	$y_{13} := 2$	$x_{26} := 6$	$y_{26} := 3$

Завдання 1:

Нанесемо точки на площину для більш чіткого розуміння їх положення:

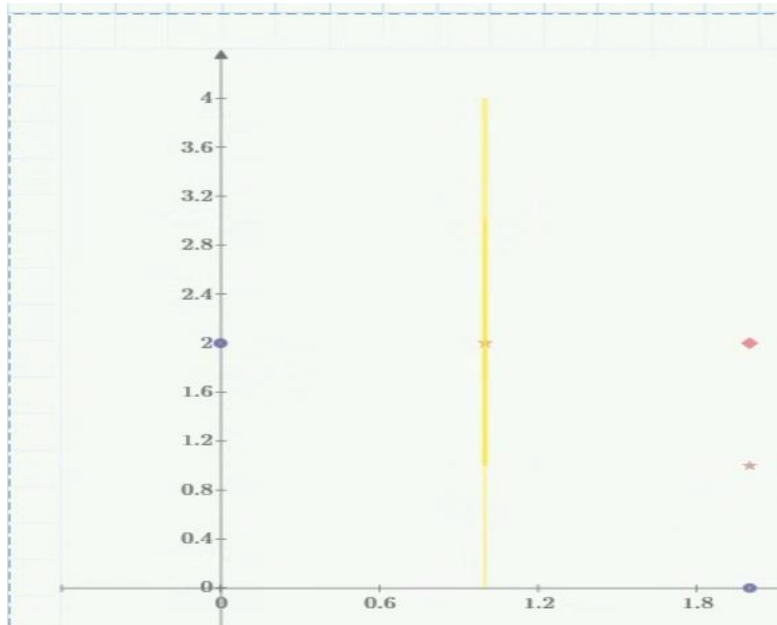


Далі будемо шукати координати середини відрізків, що сполучають об'єкти різних класів для проведення нормалі (перпендикулярних ліній) до цих відрізків, які і будуть слугувати межами між об'єктами цих класів. Шукати ці точки будемо за правилом **найближчого сусіда**.

Спочатку знайдемо координати середини відрізка між точками $\xi_1 \in a_1$ та $\xi_4 \in a_2$ та розрахуємо функцію нормалі для цієї пари точок:

$$norm3(n) := \frac{(y_{13} + y_{25})}{2} + \left(\frac{(y_{25} - y_{13})}{-(x_{25} - x_{13})} \right) \cdot \left(n - \frac{(x_{13} + x_{25})}{2} \right)$$

Зобразимо нормаль на площині:



Схожим чином виконаємо розрахунок нормалей для інших точок і нанесемо їх на площину

$$norm1 := 1 \dots 4 \quad norm5 := 0,25 \dots 2,6$$

$$norm3(n) := \frac{(y_{13} + y_{25})}{2} + \left(\frac{(y_{25} - y_{13})}{-(x_{25} - x_{13})} \right) \cdot \left(n - \frac{(x_{13} + x_{25})}{2} \right)$$

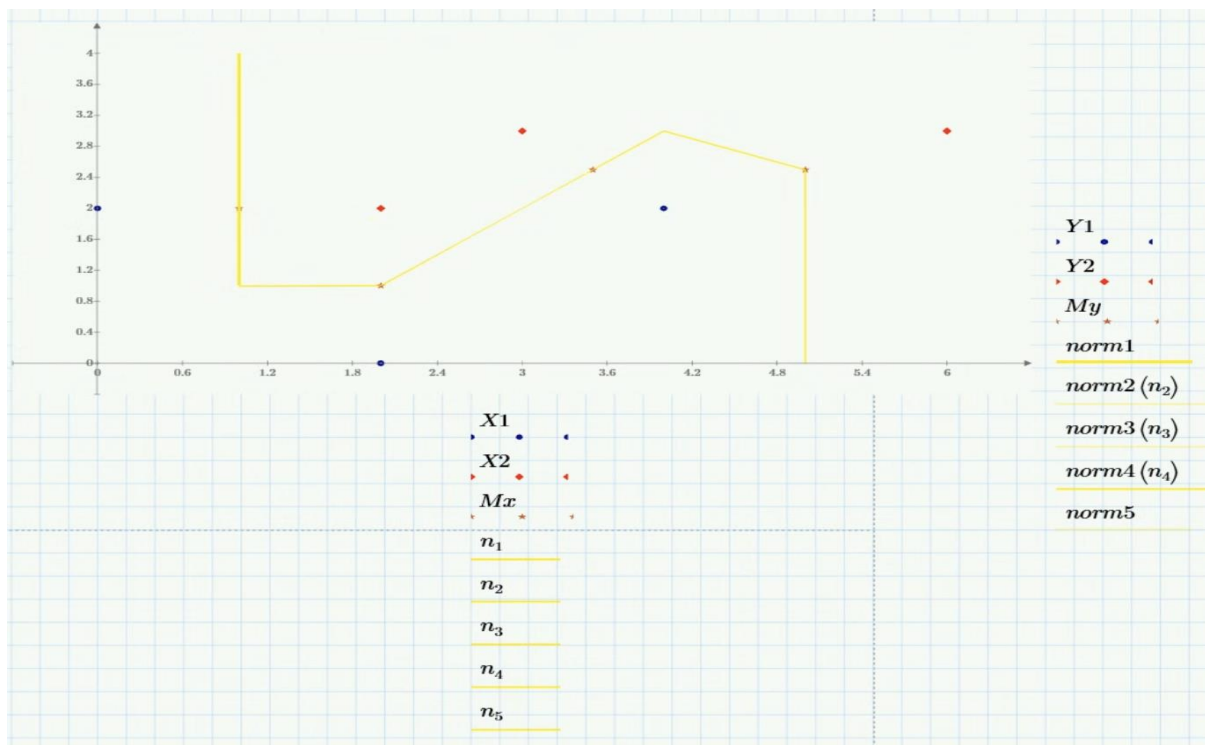
$$norm4(n) := \frac{(y_{13} + y_{26})}{2} + \left(\frac{(y_{26} - y_{13})}{-(x_{26} - x_{13})} \right) \cdot \left(n - \frac{(x_{13} + x_{26})}{2} \right)$$

$$norm2(n) := \frac{(y_{12} + y_{24})}{2} + \left(\frac{(y_{24} - y_{12})}{300} \right) \cdot \left(n - \frac{(x_{12} + x_{24})}{2} \right) +$$

Прорахуємо середні значення для всіх комбінацій

$$\begin{aligned}
 Mxn &:= \left[\frac{(x_{11} + x_{24})}{2} \right] \\
 Myn &:= \left[\frac{(y_{11} + y_{24})}{2} \right] \\
 Mx &:= \begin{bmatrix} \frac{(x_{11} + x_{24})}{2} \\ \frac{(x_{12} + x_{24})}{2} \\ \frac{(x_{13} + x_{25})}{2} \\ \frac{(x_{13} + x_{26})}{2} \end{bmatrix} \\
 My &:= \begin{bmatrix} \frac{(y_{11} + y_{24})}{2} \\ \frac{(y_{12} + y_{24})}{2} \\ \frac{(y_{13} + y_{25})}{2} \\ \frac{(y_{13} + y_{26})}{2} \end{bmatrix}
 \end{aligned}$$

Нанесемо нормалі на графік для зображення в подальшому межі між класами



Проміжки, які ми використали для побудови неперервної ламаної, що буде відділяти класи один від одного

$$\begin{aligned}
 n_1 &:= 1 & n_3 &:= 2..4 & n_5 &:= 5 \\
 n_2 &:= 1..2 & n_4 &:= 4..5
 \end{aligned}$$

Отже, за отриманою межею, визначимо вирішальне правило (3.3):

```

 $a_i :=$ 
  if  $x_i \leq 1$ 
    if  $y_i \geq 2$ 
      1
    else
      2
  else if  $1 \leq x_i \leq 2$ 
    if  $y_i \geq 1$ 
      1
    else
      2
  else if  $2 \leq x_i \leq 3.5$ 
    if  $y_i \geq 2.5$ 
      1
    else
      2
  else if  $3.5 \leq x_i \leq 5$ 
    if  $y_i \geq 2.5$ 
      1
    else
      2
  else if  $x_i > 5$ 
    2

```

- Якщо $x < 1$:
 - Якщо $y \geq 2$, то клас **1**, інакше клас **2**.
- Якщо x знаходиться в діапазоні $[1, 2]$:
 - Якщо $y \geq 1$, то клас **1**, інакше клас **2**.
- Якщо x знаходиться в діапазоні $[2, 3.5]$:
 - Якщо $y \geq 2.5$, то клас **1**, інакше клас **2**.
- Якщо x знаходиться в діапазоні $[3.5, 5]$:
 - Якщо $y \geq 2.5$, то клас **1**, інакше клас **2**.
- Якщо $x > 5$, то клас **2**.

Оскільки дані спостережень – випадкові величини ξ_i , можна провести межу між класами після статистичної обробки даних – знаходження середніх вибірових значень. Визначимо компоненти векторів - статистичних оцінок математичних очікувань (МО) класів a_1 та a_2 :

$$x_i := \begin{bmatrix} x_{11} \\ x_{12} \\ x_{13} \\ x_{24} \\ x_{25} \\ x_{26} \end{bmatrix} \quad y_i := \begin{bmatrix} y_{11} \\ y_{12} \\ y_{13} \\ y_{24} \\ y_{25} \\ y_{26} \end{bmatrix}$$

$$\begin{aligned}
 m1_0 &:= \frac{1}{3} \cdot \sum_{i=1}^3 x_i & m1_1 &:= \frac{1}{3} \sum_{i=1}^3 y_i \\
 m2_0 &:= \frac{1}{3} \cdot \sum_{i=4}^6 x_i & m2_1 &:= \frac{1}{3} \sum_{i=4}^6 y_i \\
 m1_0 &:= \frac{(x_{11} + x_{12} + x_{13})}{2} = 3 & m1_1 &:= \frac{(y_{11} + y_{12} + y_{13})}{2} = 2 \\
 m2_0 &:= \frac{(x_{24} + x_{25} + x_{26})}{2} = 5.5 & m2_1 &:= \frac{(y_{24} + y_{25} + y_{26})}{2} = 4
 \end{aligned}$$

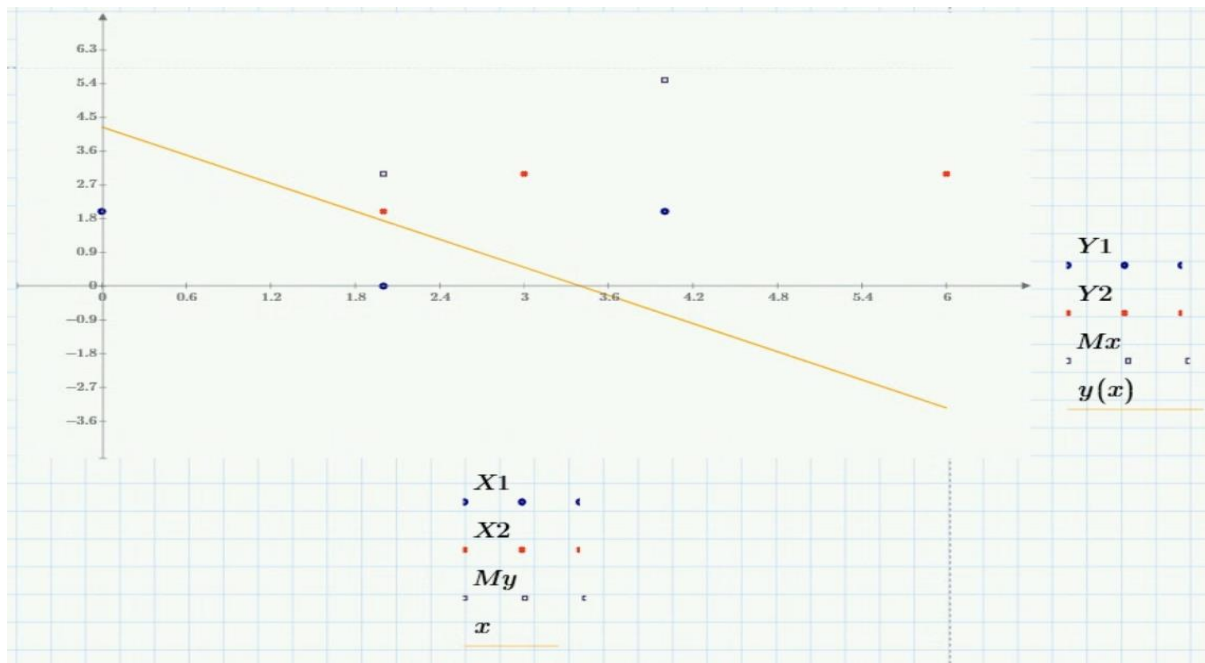
Знайдемо рівняння межі g_2 між класами a_1 та a_2 після усереднення даних спостережень. Ця межа проходить через точку з координатами:

$$\begin{aligned}
 My &:= \begin{bmatrix} m1_1 \\ m2_1 \end{bmatrix} & Mx &:= \begin{bmatrix} m1_0 \\ m2_0 \end{bmatrix} \\
 M &\left(\frac{(m1_0 + m2_0)}{2}, \frac{(m1_1 + m2_1)}{2} \right) \rightarrow M(4.25, 3.0) \\
 x_0 &:= 4.25 & y_0 &:= 3.0 \\
 k &:= \frac{(m2_1 - m1_1)}{(m2_0 - m1_0)} = 0.8 \\
 kt &:= \frac{1}{k} = 1.25 \\
 g_2(x, y) &:= -y - (x - x_0) \cdot kt + y_0
 \end{aligned}$$

Отже, за отриманою межею, визначимо вирішальне правило (3.4):

$$\begin{aligned}
 y(x) &:= -kt \cdot x + x_0 & x &:= m1_1 & x &:= 0, 1..6 \\
 \{ \text{Якщо } y &\leq y(x), \text{ тоді } \gamma_1, \text{ інакше } \gamma_2 \}
 \end{aligned}$$

Зобразимо лінію на площині разом з точками:



Завдання 3

Для експериментальної перевірки якості роботи класифікаторів за правилами (3.3) і (3.4) змодельуємо результати спостережень – масиви $(\{x1_i\}; \{y1_i\})$ та $(\{x2_i\}; \{y2_i\})$, відповідні класам $a1$ і $a2$. Вважаємо, що $x1$ та $y1$ – некорельовані компоненти двовимірної випадкової величини, що підпорядковується нормальному закону розподілу з МО і середньоквадратичним відхиленням (СКВ) В якості значень параметрів розподілу прийmemo їх статистичні оцінки:

$$i := 1, 2 \dots 6$$

$$D1_0 := \frac{1}{2} \cdot \sum_{i=1}^3 (x_i - m1_0)^2 \quad \sigma1_0 := \sqrt{m1_0} = 1.732$$

$$D1_1 := \frac{1}{2} \cdot \sum_{i=1}^3 (y_i - m1_1)^2 \quad \sigma1_1 := \sqrt{m2_0} = 2.345$$

$$D2_0 := \frac{1}{2} \cdot \sum_{i=3}^6 (x_i - m2_0)^2 \quad \sigma2_0 := \sqrt{m1_1} = 1.414$$

$$D2_1 := \frac{1}{2} \cdot \sum_{i=3}^6 (y_i - m2_1)^2 \quad \sigma2_1 := \sqrt{m2_1} = 2$$

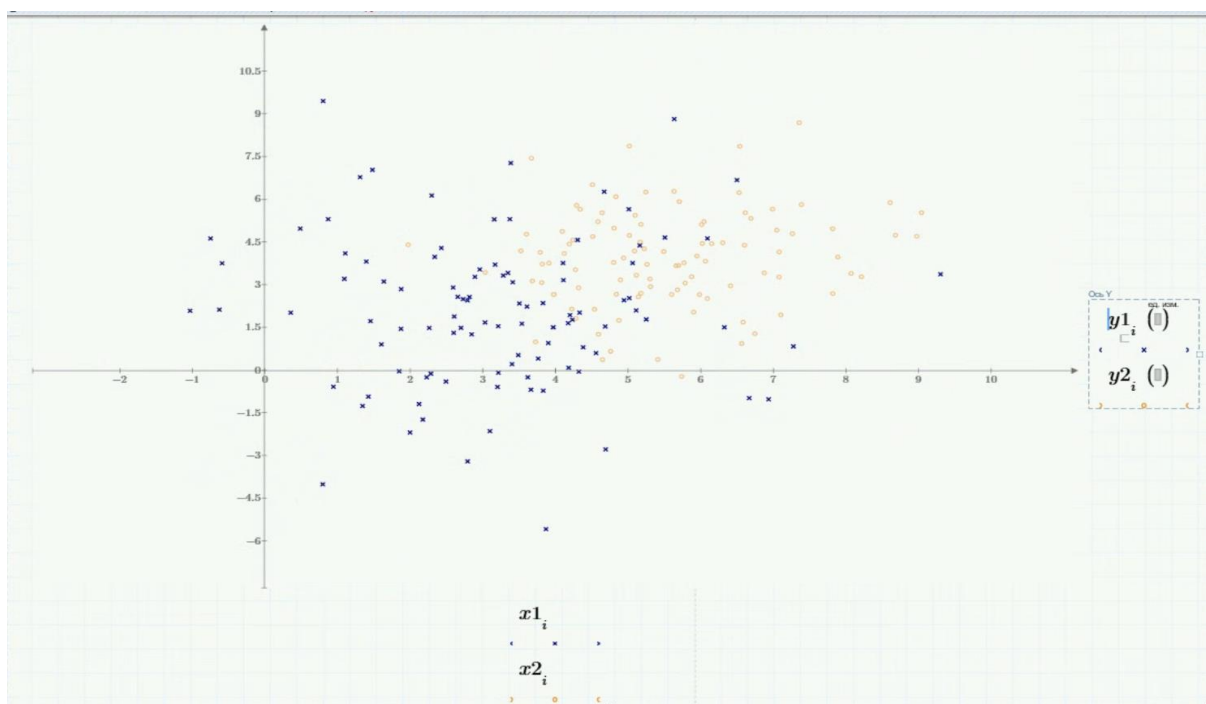
Визначимо функцію користувача, яка здійснює алгоритм генерації масиву реалізацій нормально розподіленої випадкової величини:

$$\begin{aligned}
 n &:= 48 & k &:= 1 \dots n \\
 \text{Norm}(z, m, \sigma) &:= \sqrt{\frac{12}{n}} \cdot \sigma \cdot \left(\sum_k \text{rnd}(1) - \frac{n}{2} \right) + m \\
 N &:= 100 & i &:= 0 \dots N-1
 \end{aligned}$$

Формальними аргументами цієї функції є номери елементів масиву реалізацій і параметри нормального розподілу МО m і СКО σ випадкової величини, яка моделюється. Отримаємо 100 даних спостережень класу a_1 та a_2

$$\begin{aligned}
 x1_i &:= \text{Norm}(i, m1_0, \sigma1_0) & y1_i &:= \text{Norm}(i, m1_1, \sigma1_1) \\
 x2_i &:= \text{Norm}(i, m2_0, \sigma2_0) & y2_i &:= \text{Norm}(i, m2_1, \sigma2_1)
 \end{aligned}$$

Зобразимо їх на площині:



Завдання 4

Виконаємо розпізнавання контрольної вибірки ($\{x_i\}; \{y_i\}$) за вирішальним правилом (3.3). Нехай контрольна вибірка належить класу a_1 , тоді

$$N := 100 \quad i := 0 \dots N - 1$$

$$x_i := x1_i \quad y_i := y1_i$$

Формалізуємо опис процедури ухвалення рішення (3.3):

$$\begin{array}{l}
 N := 100 \quad i := 0 \dots N - 1 \\
 x_i := x1_i \quad y_i := y1_i \\
 a_i := \begin{array}{l} \text{if } x_i \leq 1 \\ \quad \parallel \\ \quad \text{if } y_i \geq 2 \\ \quad \quad \parallel \\ \quad \quad 1 \\ \quad \text{else} \\ \quad \quad \parallel \\ \quad \quad 2 \\ \text{else if } 1 \leq x_i \leq 2 \\ \quad \parallel \\ \quad \text{if } y_i \geq 1 \\ \quad \quad \parallel \\ \quad \quad 1 \\ \quad \text{else} \\ \quad \quad \parallel \\ \quad \quad 2 \\ \text{else if } 2 \leq x_i \leq 3.5 \\ \quad \parallel \\ \quad \text{if } y_i \geq 2.5 \\ \quad \quad \parallel \\ \quad \quad 1 \\ \quad \text{else} \\ \quad \quad \parallel \\ \quad \quad 2 \\ \text{else if } 3.5 \leq x_i \leq 5 \\ \quad \parallel \\ \quad \text{if } y_i \geq 2.5 \\ \quad \quad \parallel \\ \quad \quad 1 \\ \quad \text{else} \\ \quad \quad \parallel \\ \quad \quad 2 \\ \text{else if } x_i > 5 \\ \quad \parallel \\ \quad 2 \end{array} = \begin{bmatrix} 2 \\ 1 \\ 1 \\ 2 \\ 1 \\ 2 \\ 2 \\ 1 \\ 1 \\ 2 \\ 1 \\ 1 \\ 1 \\ \vdots \end{bmatrix}
 \end{array}$$

Проведемо розпізнавання контрольної вибірки ($\{x_i\}$; $\{y_i\}$) за вирішальним правилом (3.4):

$$b_i := \text{if}(a_i \leq 1, 1, 0)$$

$$P22 := \frac{1}{N} \cdot \sum_i b_i = 0.48$$

$$P12 := 1 - P22 = 0.52$$

Проведемо дослідження для правил 3.3 та 3.4 для класу а2:

$$N := 100 \quad i := 0..N-1$$

$$x_i := x2_i \quad y_i := y2_i$$

$$b_i := \text{if}(a_i \leq 1, 1, 0)$$

$$P22 := \frac{1}{N} \cdot \sum_i b_i = 0.32$$

$$P12 := 1 - P22 = 0.68$$

За проведеними дослідженнями побудуємо таблицю результатів:

	Правило 3.3	Правило 3.4
P11	0.48	0.46
P21	0.52	0.54
P22	0.32	0.30
P12	0.68	0.70

Висновки

Виходячи з результатів досліджень, побудовані мною вирішальні правила за вибірковими значеннями та за узагальненим значенням дали дуже схожі

результати , а отже побудовано мною межі розділення класів є вірними. При цьому , метод вибірових значень статистично був точніший у задачі правильної класифікації об'єкта і менше припустився помилок. Отже можна зробити висновок, що узагальнене значення може негативно вплинути на вірну класифікацію об'єктів.