

Линейная бинарная классификация данных с интервальной неопределенностью

В. И. Ерохин, А. П. Кадочников, С. В. Сотников

Военно-космическая академия им. А. Ф. Можайского Министерства обороны Российской Федерации, Санкт-Петербург, Россия

Аннотация. Рассмотрена задача линейного бинарного отделения конечных интервальных множеств (классов). С использованием теории интервальных систем линейных неравенств задача сведена к проблеме поиска решения системы линейных неравенств специального вида. В свою очередь проблема поиска квазиоптимального решения указанной системы (или псевдорешения в случае ее несовместности и линейной неразделимости классов) сведена к задачам безусловной минимизации. Приведены иллюстративные численные примеры.

Ключевые слова: классификация, машинное обучение, интервальная неопределенность данных.

DOI

Введение

Линейная бинарная классификация является важной частной задачей машинного обучения (например [1-3]). В теоретическом обосновании большинства известных методов классификации заложены достаточно жесткие предположения: наличие большой (полной, представительной) обучающей выборки, состоящей при этом из точных данных, и (или) наличие известного типа функции распределения вероятностей признаков. Но указанные предположения редко выполняются на практике. При решении практических задач одной из наиболее часто встречающихся ситуаций является та, в которой нам необходимо классифицировать объекты при наличии малого количества информации об их признаках. И при этом, как было отмечено в работе [4, с. 147], в подавляющем большинстве прикладных задач исследователь

имеет дело с неточными исходными данными, неопределенность которых порождается различными факторами (в таблице ниже).

В зависимости от источника неточности и неопределенности данных в настоящее время используются различные модели описания неопределенных данных, включая *вероятностную, нечеткую и интервальную* модели. Каждая из них имеет свою парадигму, опирается на соответствующий теоретический аппарат, имеет свои методы анализа и область применения.

В настоящей работе будет рассмотрена интервальная модель линейной бинарной классификации. В определенном смысле статья продолжает и развивает подход, изложенный в работе [5], в которой рассматриваются задачи линейной бинарной классификации при условиях, что значения признаков находятся в интервалах с конечными границами. Также известны

Данные	Источник неопределенности и неточности
Результаты измерений	Вариабельность, шумы, ошибки измерения (систематические и случайные), ошибки округления
Прогнозные данные	Незнание, неопределенность, неполнота информации, методические ошибки, ошибки округления и дискретизации
Экспертные оценки	Субъективность, незнание

оценки математических ожиданий или средних значениях каждого признака. При этом границы интервалов и средние значения предположительно получены не в результате измерений, а заданы экспертами. Мы будем рассматривать случай, когда параметры даны с интервальной неопределенностью и при этом математические ожидания (средние значения) признаков неизвестны.

1. Постановка задачи

Пусть даны два конечных множества точек $\mathbf{P} = \{\mathbf{p}_i\}_{i=1}^{m_1}$ и $\mathbf{Q} = \{\mathbf{q}_i\}_{i=1}^{m_2}$ в пространстве \mathbb{R}^n . Указанным множествам соответствуют матрицы:

$$P = \begin{bmatrix} \mathbf{p}_1 \\ \vdots \\ \mathbf{p}_{m_1} \end{bmatrix} = (p_{ij}) \in \mathbb{R}^{m_1 \times n}, Q = \begin{bmatrix} \mathbf{q}_1 \\ \vdots \\ \mathbf{q}_{m_2} \end{bmatrix} = (q_{ij}) \in \mathbb{R}^{m_2 \times n}.$$

Рассматриваемые данные имеют следующую очевидную интерпретацию в терминах задач распознавания (классификации): точки \mathbf{p}_i и \mathbf{q}_i являются объектами классификации, их координаты p_{ij} и q_{ij} – признаками, а множества \mathbf{P} и \mathbf{Q} – соответствующими классами. При исследовании задач линейной бинарной классификации оказывается полезной объединенная матрица

$$A = \begin{bmatrix} P \\ -Q \end{bmatrix},$$

которую принято называть информационной матрицей модели линейной бинарной классификации. При этом хорошо известно, что задача линейной бинарной классификации в классической постановке может быть сформулирована как задача отделимости множеств точек \mathbf{P} , \mathbf{Q} гиперплоскостью (в общем случае не проходящей через начало координат), которая, в свою очередь, сводится к задаче решения системы линейных неравенств:

$$Pw \leq 1, Qw \geq 1 \Leftrightarrow Aw \leq b, \tag{1}$$

где $w \in \mathbb{R}^n$ – вектор неизвестных коэффициентов левой части уравнения разделяющей гиперплоскости (правая часть уравнения нормирована к 1); $b \in \mathbb{R}^{m_1+m_2}$ – вектор, первые m_1 элементов которого имеют значение 1, а последующие m_2 элементов – значение -1.

Поиск решения системы линейных неравенств (1) будем называть задачей линейной бинарной классификации с точными данными.

Предположим теперь, что элементы информационной матрицы заданы с интервальной неопределенностью – вместо точных значений матриц P и Q известны их нижние и верхние границы $\underline{P} = (\underline{p}_{ij}) \in \mathbb{R}^{m_1 \times n}$, $\bar{P} = (\bar{p}_{ij}) \in \mathbb{R}^{m_1 \times n}$, $\underline{Q} = (\underline{q}_{ij}) \in \mathbb{R}^{m_2 \times n}$, $\bar{Q} = (\bar{q}_{ij}) \in \mathbb{R}^{m_2 \times n}$ такие, что следующие неравенства выполнены поэлементно:

$$\underline{P} \leq P \leq \bar{P}, \underline{Q} \leq Q \leq \bar{Q}. \tag{2}$$

В соответствии с (2) задача линейной бинарной классификации данных с интервальной неопределенностью может быть сформулирована как задача решения системы линейных неравенств вида (1) для всех информационных матриц A , поэлементно удовлетворяющих интервальным ограничениям:

$$\underline{A} = \begin{bmatrix} \underline{P} \\ -\bar{Q} \end{bmatrix} \leq A \leq \bar{A} = \begin{bmatrix} \bar{P} \\ -\underline{Q} \end{bmatrix}. \tag{3}$$

Заметим, что множество \mathbf{A} всех матриц, удовлетворяющих условиям (3), является бесконечным, а количество его крайних точек (матриц, все элементы которых совпадают с соответствующим верхним или нижним пределом) составляет величину $2^n(m_1 + m_2)$, т.е. экспоненциально растет с ростом размерности пространства объектов. Указанные особенности поиска решения задачи линейной бинарной классификации данных с интервальной неопределенностью вызывают справедливые опасения о ее NP-трудности. К счастью, указанные опасения не оправдываются, что удастся показать, используя результаты теории линейных интервальных неравенств. Соответствующему результату посвящен следующий параграф.

2. Инструментальный результат теории линейных интервальных неравенств

Приводимый ниже результат, известный как теорема о необходимых и достаточных условиях сильной разрешимости интервальной системы линейных неравенств [6], справедлив даже для более общего случая, чем поиск решения системы (1) для всех матриц $A \in \mathbf{A}$.

Теорема 1. Система линейных неравенств $Ax \leq b$ разрешима для всех A, b таких, что $\underline{A} \leq A \leq \bar{A}$, $\underline{b} \leq b \leq \bar{b}$ тогда и только тогда, когда разрешима система:

$$\bar{A}x^+ - \underline{A}x^- \leq \bar{b}, \quad x^+, x^- \geq 0. \tag{4}$$

При этом $x = x^+ - x^-$.

Следующее утверждение является очевидным следствием Теоремы 1.

Утверждение 1. Задача линейной бинарной классификации данных с интервальной неопределенностью, заданная условиями (1) - (3), разрешима тогда и только тогда, когда имеет решение система линейных неравенств:

$$\begin{bmatrix} \bar{P} \\ -\underline{Q} \end{bmatrix} w^+ - \begin{bmatrix} \underline{P} \\ -\bar{Q} \end{bmatrix} w^- \leq \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \quad w^+, w^- \geq 0, \quad (5)$$

где 1 – вектор размерности m_1 , составленный из единиц; -1 – вектор размерности m_2 , составленный из минус единиц. При этом $w = w^+ - w^-$.

3. Практические задачи линейной бинарной классификации данных с интервальной неопределенностью

Очевидно, что приведенные выше постановки задач линейной бинарной классификации (с точными данными и данными с интервальной неопределенностью), сводящиеся к задачам нахождения решений систем линейных неравенств, можно лишь очень ограниченно рассматривать в качестве практических методов классификации. Указанные методы в настоящее время уже достаточно сложны и многообразны, но, в то же время, продолжают интенсивно развиваться (см., например обзоры [7-10] с внушительными списками цитируемой литературы, и работу [11] в качестве примера очередного нового исследования). В связи с этим обратим внимание на важный практический аспект рассматриваемой проблемы, который сравнительно редко обсуждается в научных публикациях. Речь идет о контексте решаемой задачи классификации, в котором можно выделить два крайних случая:

- задача классификации является одной из частных задач автоматического (без прямого управления человеком) функционирования некоторой сложной технической системы, возможно работающей в режиме реального времени и обрабатывающей большие объемы информации (например [12]);
- задача классификации является «разовой» частной задачей обработки результатов некоторого научного исследования, выполняется чело-

веком «в ручном режиме» для сравнительно небольшого объема данных и не критична к оптимальности полученного решения, времени счета и техническим характеристикам компьютера (объему оперативной памяти и производительности процессора).

Предлагаемые ниже «технические решения» задач линейной бинарной классификации не предназначены для первого случая, но, возможно, окажутся полезными для второго, поскольку их вычислительная реализация возможна с помощью любого стандартного решателя задач безусловной минимизации.

Рассмотрим следующую задачу безусловной минимизации:

$$\left\| \begin{bmatrix} P \\ -Q \end{bmatrix} w - \begin{bmatrix} 1 \\ -1 \end{bmatrix} \right\|_+ \rightarrow \min_w (= \gamma), \quad (6)$$

где $\|\cdot\|_+$ – некоторая векторная норма: $[\cdot]_+$ – операция положительной срезки, применяемая к векторному аргументу поэлементно.

Если *точные* множества \mathbf{P} и \mathbf{Q} являются линейно неразделимыми, то $\gamma > 0$ при любом выборе $\|\cdot\|_+$, а сама задача может быть интерпретирована как проблема поиска оптимального (в смысле минимума $\|\cdot\|_+$ - нормы ошибок классификации) псевдорешения задачи линейной бинарной классификации. Заметим, что при выборе строго выпуклой векторной нормы, например, евклидовой, решение задачи (6) будет единственным.

Перейдем теперь от линейно неразделимых *точных* множеств \mathbf{P} и \mathbf{Q} к соответствующим *интервальным* множествам. Рассмотрим следующую задачу безусловной минимизации:

$$\|d(w)\| = \left\| \begin{bmatrix} P_c \\ -Q_c \end{bmatrix} w + \begin{bmatrix} P_r \\ Q_r \end{bmatrix} |w| + \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right\|_+ \rightarrow \min_w (= \gamma), \quad (7)$$

где $\|\cdot\|_+$ – некоторая абсолютная векторная норма: $P_c = (\underline{P} + \bar{P})/2$, $P_r = (\bar{P} - \underline{P})/2$; $Q_c = (\underline{Q} + \bar{Q})/2$; $Q_r = (\bar{Q} - \underline{Q})/2$.

Заметим, что $P_r, Q_r \geq 0$ в силу условий (2), $\gamma \geq 0$ в силу аксиомы неотрицательности векторной нормы. Покажем, что справедливо следующее утверждение.

Утверждение 2. Задача нахождения псевдорешения системы линейных неравенств (5) эквивалентна задаче безусловной минимизации (7).

Доказательство. Поиск решения системы линейных неравенств (5) (в случае ее совместности) или псевдорешения (в случае несовместности) эквивалентен задаче минимизации нормы невязки указанной системы и может быть записана в виде следующей задачи математического программирования:

$$\begin{aligned} & \|g(w^+, w^-)\| = \\ & = \left\| \begin{bmatrix} \bar{P} \\ -\underline{Q} \end{bmatrix} w^+ - \begin{bmatrix} \underline{P} \\ -\bar{Q} \end{bmatrix} w^- - \begin{bmatrix} 1 \\ -1 \end{bmatrix} \right\|_+ \rightarrow \min_{w^+, w^- \geq 0} (= \Gamma). \end{aligned} \quad (8)$$

Очевидно, что минимум в задаче (8) достигается, причем в общем случае $\Gamma \geq 0$ и $\Gamma > 0$ тогда и только тогда, когда система линейных неравенств (5) несовместна. Учитывая соотношения для P_c, P_r, Q_c, Q_r и проводя перегруппировку аргументов w^+, w^- выражение для вектора $g(w^+, w^-)$ можно переписать в виде:

$$\begin{aligned} & g(w^+, w^-) = \\ & = \left[\begin{bmatrix} P_c \\ -Q_c \end{bmatrix} (w^+ - w^-) + \begin{bmatrix} P_r \\ -Q_r \end{bmatrix} (w^+ + w^-) - \begin{bmatrix} 1 \\ -1 \end{bmatrix} \right]_+. \end{aligned}$$

В силу Теоремы 1 выполняется условие $w = w^+ - w^-$. В то же время $|w| \leq w^+ + w^-$. С учетом проделанных выкладок и свойств абсолютной векторной нормы (например [13]), имеем:

$$\begin{aligned} & d(w) = d(w^+ - w^-) \leq g(w^+, w^-) \Rightarrow \\ & \|d(w)\| \leq \|g(w^+, w^-)\| \Rightarrow \gamma = \\ & \min_w \|d(w)\| \leq \min_{w^+, w^- \geq 0} \|g(w^+, w^-)\| = \Gamma. \end{aligned}$$

Для завершения доказательства остается показать, что $\gamma = \Gamma$. При $\Gamma = 0$ указанное утверждение очевидно. Предположим, что выполняются условия $0 < \gamma < \Gamma$ и вектор w^* – решение задачи (7). Рассмотрим векторы $w^+ = [w^*]_+$ и $w^- = [-w^*]_+$. Очевидно, что $w^+ - w^- = w^*$, $w^+ + w^- = |w^*|$, в силу чего $g(w^+, w^-) = d(w^*)$, откуда и следует $\Gamma = \|g(w^+, w^-)\| = \|d(w^*)\| = \gamma$.

Заметим, что условие (7) является задачей выпуклой негладкой безусловной минимизации,

имеющей в случае линейно неразделимых интервальных множеств \mathbf{P} и \mathbf{Q} при использовании строго выпуклой нормы единственное решение, которое можно интерпретировать как оптимальное приближенное линейное разделение двух неразделимых интервальных множеств.

Предположим теперь, что интервальные множества \mathbf{P} и \mathbf{Q} являются линейно разделимыми. В этом случае $\gamma = 0$, но решение задачи (7) в общем случае не единственно. Рассмотрим следующую модификацию указанной задачи:

$$\begin{aligned} & \|d(w, \delta)\| = \\ & = \left\| \begin{bmatrix} P_c \\ -Q_c \end{bmatrix} w + \begin{bmatrix} P_r \\ Q_r \end{bmatrix} |w| + \begin{bmatrix} -1 \\ 1 \end{bmatrix} + \delta \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right\|_+ \rightarrow \min_w (= \gamma_\delta), \end{aligned} \quad (9)$$

где $\delta \geq 0$ – некоторый скалярный параметр.

Очевидно, что при $\delta = 0$ задача (9) совпадает с задачей (7), $\gamma_\delta = 0$ и решение (9) в общем случае не единственно. В то же время можно показать, что начиная с некоторого значения $\delta > 0$ будет выполняться условие $\gamma_\delta = 0$ и задача (9) будет иметь единственное решение. Как показывают расчеты, указанное решение является оптимальным в смысле максимальной ширины зазора между интервальными множествами \mathbf{P} и \mathbf{Q} .

4. Иллюстративные численные примеры

На Рис. 1 и 2 представлены результаты решения модельных задач линейной бинарной классификации линейно разделимых и неразделимых интервальных множеств при $n = 2$, полученные с помощью численного решения задач (7) и (9), соответственно. Расчеты проводились в среде Mathcad® 15.0. Матрицы P_c, P_r, Q_c, Q_r вычислялись с помощью генератора псевдослучайных равномерно распределенных чисел. Численное решение задач (7) и (9) (определение вектора w^*) осуществлялось с помощью встроенного в среду Mathcad решателя Minimize с использованием метода сопряженных градиентов и конечно-разностной центральной аппроксимацией производных. Минимальное значение $\delta > 0$, гарантирующее выполнение условия $\gamma_\delta > 0$ в задаче (9) выбиралось с помощью одномерного

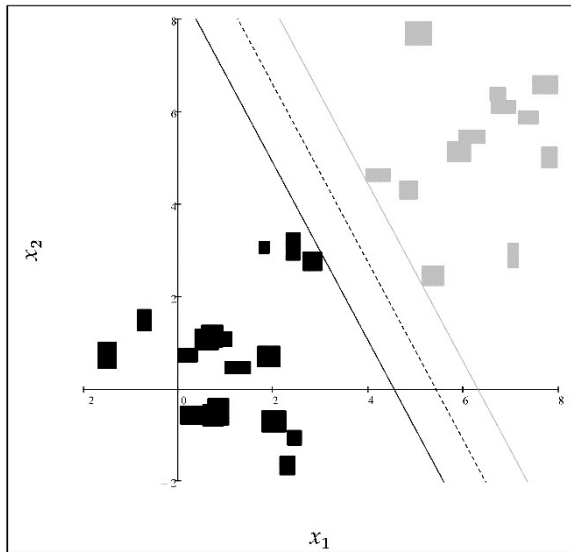


Рис. 1. Строгое отделение двух множеств с интервальной неопределенностью

- ■ – объекты интервального множества **P**;
- — — прямая $xw^* = 1 - \alpha$ (граница множества **P**);
- ■ – объекты интервального множества **Q**;
- — — прямая $xw^* = 1 + \beta$ (граница множества **Q**);
- — — прямая $xw^* = 1 + (\beta - \alpha) / 2$ (квазиоптимальная разделяющая прямая)

поиска. Уравнение разделяющей прямой подвергалось дополнительной коррекции, в результате которой принимало вид:

$$xw^* = 1 + \frac{\beta - \alpha}{2},$$

где $x = [x_1 \ x_2]$ – произвольный вектор (вектор-строка) из пространства объектов;

$$\alpha = \min_i (1 - P_c w^* - P_r |w^*|)_i,$$

$$\beta = \min_i (-1 + Q_c w^* - Q_r |w^*|)_i.$$

Абсолютные значения параметров α и β пропорциональны расстояниям от разделяющей прямой до интервальной границы ближайшего к указанной прямой объекта множеств **P** и **Q**, соответственно. При этом $\alpha, \beta > 0$ в случае линейной отделимости и $\alpha, \beta < 0$ – в случае линейной неотделимости. Учитывая указанную информацию, оказывается возможным построить две дополнительные прямые, заданные уравнениями $xw^* = 1 - \alpha$ (граница множества **P**) и $xw^* = 1 + \beta$ (граница множества **Q**). В случае линейно разделимых множеств указанные прямые образуют

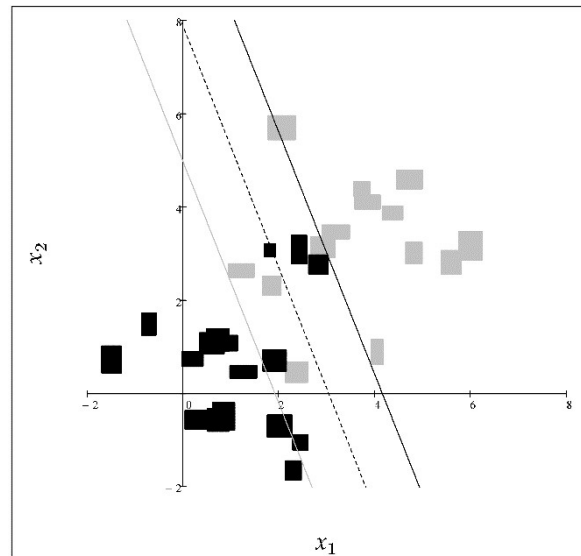


Рис. 2. Псевдоотделение двух линейно неразделимых множеств с интервальной неопределенностью

- ■ – объекты интервального множества **P**;
- — — прямая $xw^* = 1 - \alpha$ (граница множества **P**);
- ■ – объекты интервального множества **Q**;
- — — прямая $xw^* = 1 + \beta$ (граница множества **Q**);
- — — прямая $xw^* = 1 + (\beta - \alpha) / 2$ (квазиоптимальная разделяющая прямая)

«разделительную полосу» (возможно, максимальной ширины), а в случае неразделимых множеств – полосу (возможно, минимальной ширины), в которой находятся объекты обеих множеств с учетом их интервальных границ.

5. Некоторые замечания и комментарии

Заметим, что построение гиперплоскости, разделяющей два интервальных множества, сводится к задаче поиска решения системы линейных неравенств, формально отличающейся от соответствующей системы в задаче линейного бинарного разделения точечных множеств только более высокой размерностью и требованием неотрицательности переменных. По этой причине можно утверждать:

- не существует теоретических препятствий для построения «интервальных» методов линейной бинарной классификации, являющихся аналогами метода опорных векторов и его модификаций;
- для устойчивой линейной бинарной классификации интервальных множеств можно, также

как и в случае множеств точечных, использовать методы матричной коррекции систем линейных неравенств, предложенные в работах [14-16].

Заключение

В работе рассмотрена задача линейной бинарной классификации интервальных множеств. Она формулируется как проблема поиска решения интервальной системы линейных неравенств. Необходимые и достаточные условия существования решения данной проблемы и его вид устанавливает теорема, являющаяся известным результатом теории интервальных систем линейных неравенств. С использованием указанной теоремы задача линейной бинарной классификации интервальных множеств сведена к проблеме поиска решения системы линейных неравенств специального вида. Для построения решения (или псевдорешения в случае линейной неразделимости классов) предложены соответствующие задачи безусловной минимизации. Приведены иллюстративные численные примеры.

В заключение можно обратить внимание на потенциальную возможность доработки (подстраивания) широкого спектра известных и вновь разрабатываемых методов линейной бинарной классификации (в том числе – методов псевдорешения задачи при отсутствии линейного разделения) к случаю интервальных множеств. Указанная возможность обусловлена практически совпадающей формой систем линейных неравенств, описывающей задачу бинарной линейной классификации, как точных так и интервальных множеств.

Литература

1. Вапник В.Н., Червоненкис А.Я. Теория распознавания образов (статистические проблемы обучения). М.: Наука. 1974. 416 с.
2. Воронцов К. В. Лекции по линейным алгоритмам классификации // URL: [http://www.](http://www.machinelearning.ru/wiki/images/6/68/voron-ML-Lin.pdf)

[machinelearning.ru/wiki/images/6/68/voron-ML-Lin.pdf](http://www.machinelearning.ru/wiki/images/6/68/voron-ML-Lin.pdf). (дата обращения: 20.01.2023). 2009.

3. Deisenroth M.P., Faisal A.A., Ong C.S. Mathematics for machine learning. Cambridge University Press. 2020. URL: <https://mml-book.com> (accessed January 20, 2023).
4. Воцинин А.П. Задачи анализа с неопределенными данными – интервальность и/или случайность? // Труды Международной конференции по вычислительной математике МКВМ-2004. Рабочие совещания. Ред.: Ю.И. Шокин, А.М. Федотов, С.П. Ковалев и др. Новосибирск: Изд. ИВМиМГ СО РАН. 2004. С. 147-158.
5. Уткин Л.В., Жук Ю.А., Селиховкин И.А. Модель классификации на основе неполной информации о признаках в виде их средних значений // Искусственный интеллект и принятие решений. 2012. № 2. С. 16-26.
6. Фидлер М., Недома Й., Рамик Я. и др. Задачи линейной оптимизации с неточными данными / Пер. с англ. М.-Ижевск: НИЦ «Регулярная и хаотическая динамика», Институт компьютерных исследований, 2008. 288 с.
7. Bennett K. P., Campbell C. Support vector machines: hype or hallelujah? // ACM SIGKDD explorations newsletter. 2000. V. 2. No 2. P. 1-13.
8. Moguerza J. M., Muñoz A. Support vector machines with applications // Statistical Science. 2006. V. 21. No 3. P. 322-336.
9. Carrizosa E., Morales D. R. Supervised classification and mathematical optimization // Computers & Operations Research. 2013. V. 40. No 1. P. 150-165.
10. Silva A. P. D. Optimization approaches to supervised classification // European Journal of Operational Research. 2017. V. 261. No 2. P. 772-788.
11. Nueda M.J., Gandía C., Molina M.D. LPDA: A new classification method based on linear programming // PLoS ONE. 2022. V. 17. No 7. P. 1-13.
12. Sevakula R.K., Verma N.K. Improving Classifier Generalization: Real-Time Machine Learning based Applications // Studies in Computational Intelligence. Springer Nature. 2022. V. 989. 166 p.
13. Ланкастер П. Теория матриц / Пер. с англ. М.: Наука. 1982. 272 с.
14. Горелик В.А., Муравьева О.В. Построение разделяющей гиперплоскости, устойчивой к коррекции данных // Моделирование, декомпозиция и оптимизация сложных динамических процессов. 2013. Т. 28. № 1. С. 42-49.
15. Муравьева О.В. Исследование параметрической устойчивости решений систем линейных неравенств и построение разделяющей гиперплоскости // Дискретный анализ и исследование операций. 2014. Т. 21. № 3. С.53-63.
16. Муравьева О.В. Параметрическая устойчивость систем линейных неравенств // Таврический вестник информатики и математики. 2015. № 2 (27). С. 101-109.

Ерохин Владимир Иванович. Доктор физико-математических наук, профессор. Старший научный сотрудник. Военно-космическая академия имени А.Ф. Можайского Министерства обороны РФ. Области исследований: вычислительная математика, математическое программирование, моделирование и оптимизация сложных технических систем. E-mail: erohin_v_i@mail.ru (ответственный за переписку).

Кадочников Андрей Павлович. Кандидат технических наук. Заведующий лабораторией. Военно-космическая академия имени А.Ф. Можайского Министерства обороны РФ. Области исследований: моделирование сложных технических систем. E-mail: kado162@mail.ru

Сотников Сергей Владимирович. Кандидат технических наук. Старший научный сотрудник. Военно-космическая академия имени А.Ф. Можайского Министерства обороны РФ. Области исследований: автоматизированные системы управления сложными техническими системами. E-mail: svsotnikov@gmail.com

Linear Binary Classification under Interval Uncertainty of Data

V. I. Erokhin, A. P. Kadochnikov, S. V. Sotnikov

A. F. Mozhaisky Military-Space Academy of Ministry of Defence of the Russian Federation, St. Petersburg, Russia

Abstract. The problem of linear binary classification of interval sets is considered. This problem is formulated as a problem of finding a solution to an interval system of linear inequalities. Necessary and sufficient conditions for the existence of a solution to this problem and its form are established by a theorem, which is a well-known result of the theory of interval systems of linear inequalities. The problem of linear binary classification of interval sets is reduced to the problem of finding a solution to a system of linear inequalities of a special form. To construct a solution (or a pseudo-solution in the case of linear inseparability of classes), the corresponding problems of unconditional minimization are proposed. Illustrative numerical examples are given. The article notes the potential possibility of adapting a very wide range of known and newly developed methods of linear binary classification to the case of interval sets, due to the practically coinciding form of systems of linear inequalities that describe the problem of binary linear classification of both exact and interval sets.

Keywords: classification, machine learning, interval data uncertainty.

DOI

References

1. Vapnik V.N., and Chervonenkis A.Ya. 1974. *Teoriya raspoznavaniya obrazov (statisticheskie problemy obucheniya)* [pattern recognition theory (statistical learning problems)]. Moscow: Nauka. 416 p.
2. Vorontsov K.V. *Lektsii po linejnym algoritmam klassifikatsii* [Lectures on linear classification algorithms]. 2009. URL: <http://www.machinelearning.ru/wiki/images/6/68/voron-ML-Lin.pdf> (accessed January 20, 2023).
3. Deisenroth M. P., Faisal A. A., Ong C. S. 2020. *Mathematics for machine learning*. Cambridge University Press. URL: <https://mml-book.com> (accessed January 20, 2023).
4. Voshchinin A.P. *Zadachi analiza s neopredelennymi dannymi – interval'nost' i/ili sluchajnost'?* [Problems of analysis with uncertain data – interval and/or randomness?] // *Proceedings of International Conference on Computational Mathematics IICM-2004. Workshops* / Eds.: Yu. I. Shokin, A.M. Fedotov, S.P. Kovalyov, et al. – Novosibirsk: ICM&MG Publ., 2004.
5. Utkin L.V., Zhuk Yu.A., Selikhovkin I.A. *Model' klassifikatsii na osnove nepolnoj informatsii o priznakakh v vide ikh srednikh znachenij* [Classification model based on incomplete information about signs in the form of their average values] // *Iskusstvennyj intellekt i prinyatie reshenij* [Artificial intelligence and decision making]. 2012. V. 2. P. 16-26.
6. Fiedler M., Nedoma J., Ramnik J. et al. *Linear Optimization Problems with Inexact Data*. N.Y.: Springer Science+Business Media, Inc. 2006.
7. Bennett K.P., Campbell C. *Support vector machines: hype or hallelujah?* ACM SIGKDD explorations newsletter. 2000. V. 2. No 2. P. 1-13.
8. Moguerza J.M., Muñoz A. *Support vector machines with applications*. Statistical Science. 2006. V. 21. No 3. P. 322-336.
9. Carrizosa E., Morales D.R. *Supervised classification and mathematical optimization*. Computers & Operations Research. 2013. V. 40. No 1. P. 150-165.
10. Silva A.P.D. *Optimization approaches to supervised classification*. European Journal of Operational Research. 2017. V. 261. No 2. P. 772-788.
11. Nueda M.J., Gandí, C, Molina M.D. *LPDA: A new classification method based on linear programming* // PLoS ONE. 2022. V. 17. No 7. P. 1-13.
12. Sevakula R.K., Verma N.K. *Improving Classifier Generalization: Real-Time Machine Learning based Applications* // *Studies in Computational Intelligence*. Springer Nature. 2022. V. 989. 166 p.
13. Lankaster P. *Theory of matrices*. New York-London: Academic Press. 1969.
14. Gorelik V.A., Murav'eva O.V. *Postroenie razdelyayushchej giperploskosti, ustojchivoj k korrektsii dannykh* [Construction of a separating hyperplane, resistant to data correction] // *Modelirovanie, dekompozitsiya i optimizatsiya slozhnykh dinamicheskikh protsessov* [Modeling, Decomposition and Optimization of Complex Dynamic Processes] 2013. V.28. No 1. P. 42-49.
15. Murav'eva O.V. *Studying the stability of solutions to systems of linear inequalities and // constructing separating hyperplanes*. Journal of applied and industrial mathematics. 2014. V. 8. No 3. P. 349-356.
16. Murav'eva O.V. *Parametricheskaya ustojchivost' sistem linejnykh neravenstv* [Parametric stability of systems of linear inequalities] // *Tavrisheskij vestnik informatiki i matematiki* [Tauride Bulletin of Informatics and Mathematics] 2015. V.2. No 27. P. 101-109.

Erokhin Vladimir I. Doctor of physical and mathematical sciences, professor. Senior researcher. A.F. Mozhaisky Military-Space Academy, Ministry of Defense of the Russian Federation. Research areas: computational mathematics, mathematical programming, modeling and optimization of complex technical systems. E-mail: erohin_v_i@mail.ru

Kadochnikov Andrey P. Candidate of technical sciences. Head of laboratory. A.F. Mozhaisky Military- Space Academy, Ministry of Defense of the Russian Federation. Research areas: modeling of complex technical systems. E-mail: kado162@mail.ru

Sotnikov Sergey V. Candidate of technical sciences. Senior researcher. A.F. Mozhaisky Military-Space Academy, Ministry of Defense of the Russian Federation. Research areas: automated control systems for complex technical systems. E-mail: svstotnikov@gmail.com