Генерация базы знаний на основе нечеткой кластеризации

Т. А. Моисеева, Т. М. Леденева

Федеральное государственное бюджетное образовательное учреждение высшего образования «Воронежский государственный университет», Воронеж, Россия

Аннотация. В данной статье предложен подход для генерации оптимальной базы правил нечеткой системы, основанный на эллипсоидальной кластеризации наблюдаемых данных. Посылки нечетких правил образуются путем построения проекций эллипсоидов на оси координат, а заключения – либо с использованием осей эллипсоидов, либо также на основе проецирования. Идея оптимизации заключается в использовании эллипсоидов минимального объема, включающего все точки кластера. В статье осуществляется сравнительный анализ различных способов выбора оптимальных параметров для эллипсоидов, покрывающих кластеры. Оценка точности аппроксимации полученной нечеткой системы осуществляется на основе среднеквадратичной ошибки.

Ключевые слова: кластеризация, если-то правила, база знаний, нечеткая система.

DOI 10.14357/20718632230110

Введение

Нечеткие продукционные системы или, более кратко, нечеткие системы (НС), являются обобщением обычных продукционных систем и разрабатываются как ядро экспертных систем различного назначения. НС получили широкое распространение при создании систем управления техническими объектами [1,2] и технологическими процессами, например [3,4], систем диагностики и прогнозирования [5-7]; при решении задач классификации в различных приложениях [8-10]. На вычислительном уровне НС могут рассматриваться как гибкие математические структуры, которые подобно нейронным сетям и сетям радиальных базисных функций могут аппроксимировать большой класс нелинейных систем с необходимой степенью точности. По сравнению с другими способами аппроксимации нелинейности, НС обладают тем преимуществом, что правила в них могут быть сформулированы на есте-

ственном языке, что возможно за счет использования лингвистических моделей представления информации. Использование знаний на естественном языке делает модель прозрачной для интерпретации и анализа, что повышает объяснительное способности НС - свойство, которое является одним их важнейших для интеллектуальных систем. Важнейшим этапом проектирования НС является формирование базы продукционных правил (если-то правил). Данный этап осуществляется либо с участием эксперта – высококвалифицированного специалиста предметной области, либо на основе наблюдаемых данных с использованием специальных процедур. Различают три основных типа нечетких моделей продукционных правил [11]: лингвистическая модель (посылка и заключение являются нечеткими высказываниями), реляционная модель (задается нечеткое отношение между термами посылок и заключений соответствующих лингвистических шкал) и модель Takagi-Sugeno (TS-модель) (посылка является нечетким высказыванием, а заключение представлено функцией, как правило, линейной). База правил НС совместно с базой данных, содержащей информацию о параметрах лингвистических шкал входных и/или выходных переменных, представляет собой базу знаний НС. Уменьшение числа правил в общем случае понижает точность модели, но приводит к ускорению обработки данных, что важно для некоторых приложений, имеющих критическое значение, или функционирующих в режиме реального времени. Актуальной проблемой является определение «необходимого и достаточного» количества правил в базе знаний НС для обеспечения ее качественного функционирования.

1. Постановка задачи

Если для каждой переменной i из n входных переменных построена лингвистическая шкала, включающая m_i термов, то всего суще-

ствует $\prod_{i=1}^{n} m_{i}$ различных вариантов посылок

правил, при этом каждому правилу соответствует свое заключение. Однако при функционировании НС некоторые правила никогда не будут активированы, но система всякий раз будет затрачивать время на их просмотр, проявляя неэффективность. Можно выделить следующие основные направления исследований при формировании базы знаний и, в частности, базы правил НС: 1) разработка методов генерации базы правил как на основе эмпирических гипотез (информации от экспертов), так и путем адаптации к имеющимся экспериментальным данным, образующих обучающую выборку; 2) структурная и параметрическая оптимизация базы правил, которая осуществляется не только за счет тонкой настройки параметров термов лингвистических шкал, но и путем сокращения количества правил на основе анализа их посылок и заключений, а также ликвидации противоречивых правил. Среди существующих подходов к решению перечисленных выше задач выделим следующие. В [12] для построения базы правил используются RBF-сети, но формируемая база получается слишком сложной. Компонентный анализ (РСА) позволяет уменьшить количество входов, но опять же усложня-

ет понимание системы [13]. В [14] предложен подход для структуризации базы правил, что обеспечивает ее иерархическое представление. При наличии экспериментальных данных достаточного объема целесообразно использовать методы кластеризации. В рамках проведенного исследования рассматривался алгоритм нечетких с-средних (fuzzy C-means clustering – FCM) [15]. Главным его недостатком является невозможность корректного разбиения на кластеры в случае, если кластеры имеют различную дисперсию по разным осям. Его преодоление возможно за счет выбора вида нормы, которая определяет степень похожести векторных оценок объектов и позволяет выделить кластеры различной формы. Но требование того, чтобы сумма степеней принадлежности к различным кластерам была равна единице, порождает низкую устойчивость алгоритма к выбросам и сильно перекрывающимся кластерам. При отказе от данного требования сформировался класс возможностных алгоритмов [16], основанных на теории возможностей и интерпретирующих характеристическое значение точки по отношению к кластеру как возможность ему принадлежать, однако результат работы таких алгоритмов может отличаться нестабильностью. В [17] для модификации FCM используется теория интуиционистских нечетких множеств, что позволяет каждому объекту помимо степени принадлежности назначить степень его непринадлежности к кластеру. Использование такого типа множеств, с одной стороны, улучшает точность результатов кластеризации, а с другой - увеличивает вычислительную сложность и предполагает хранение и обработку дополнительной информации. Помимо перечисразновидности ленных, известны FCM. основанные на теории приближенных или «грубых» множеств [18], на теории интервальных нечетких множеств [19], а также в последнее время активно развиваются так называемые ядерные (kernel-based) методы FCM [20], которые позволяют вместо линейных границ между кластерами строить нелинейные границы.

Рассмотрим особенности реализации метода FCM. В общем случае норму можно задать через симметрическую положительно определенную матрицу A в виде

 $||x-c||_A^2 = \langle (x-c)A, x-c \rangle$, где x – векторная оценка объекта, с – центр кластера, при этом тип матрицы определяет форму кластера. Матрица A = I, где I – единичная матрица, определяет евклидову норму, которая позволяет выделять кластеры в виде сфер. Диагональная матрица A, в которой элементы главной диагонали интерпретируются как веса координат центров, позволяет выделять кластеры в виде эллипсоидов, ориентированных вдоль координатных осей. Если в качестве матрицы A использовать ковариационную матрицу, то получим норму Махаланобиса, которая выделяет кластеры в виде эллипсоидов, оси которых могут быть ориентированы в произвольных направлениях. Если выборка данных нечастая или шумная, то ковариация кластеров велика, в результате чего получаются большие эллипсоиды. Если данные плотные, то ковариация кластеров мала, в результате чего, получаются мелкие эллипсоиды. В базовом алгоритме FCM расстояние между объектом и центром кластера определяется на основе евклидовой нормы [15]. Алгоритм Густавсона-Кесселя позволяет найти кластеры различной геометрической формы, используя адаптивную норму для каждого кластера [21]. После разбиения данных на кластеры, каждому кластеру можно поставить в соотпродукционное правило. позволяет сгенерировать минимально необходимое количество правил, при этом за счет оптимизации формы кластеров можно улучшить точность аппроксимации, а, следовательно, и качество обработки данных. На Рис. 1 база знаний НС как бы покрывает график неизвестной функции «заплатками», усредняя результат в местах их пересечения. Если получены такие «заплатки», то их можно объединять, сокращая тем самым количество правил в базе знаний. Функции принадлежности термов лингвистических шкал входных и/или выходной переменных получаются поточечным проецированием матрицы разбиения на оси переменных. Точечно-определенные нечёткие множества затем аппроксимируются подходящей функцией.

Другой подход [21] заключается в том, что эллипсоид вписывается в n-мерный параллелепипед, который проецируется на оси переменных. Функция принадлежности задается

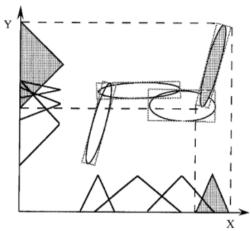


Рис. 1. Формирование продукционных правил на основе эллипсоидальной кластеризации

аналитически с использованием длины проекции и координат центра кластера. Цель статьи заключается в представлении подхода для формирования базы знаний НС на основе эллипсоидальной кластеризации наблюдаемых данных и его экспериментального исследования.

2. Основы эллипсоидальной кластеризации

2.1. Теоретическое обоснование метода

Пусть $x \in \mathbb{R}^n$, A — симметрическая положительно определенная матрица, c — n-мерный вектор центра эллипсоида, тогда определим эллипсоид выражением

$$E(c,A) = \{x \in \mathbb{R}^n : (A^{-1}(x-c),(x-c)) \le 1\}.$$

Заметим, что данный класс эллипсоидов характеризуется n+n(n+1)/2 параметрами, при этом n(n+1)/2 параметров задают элементы симметрической матрицы A, а n параметров определяют вектор c. Объем эллипсоида определяется выражением $v_E = \pi^{n/2} \left(\det A \right)^{1/2} \bigg/ \Gamma \left(\frac{n}{2} + 1 \right)$, (где $\Gamma(\cdot)$ — гаммафункция Эйлера) и пропорционален величине

функция Эилера) и пропорционален величине $\det A$.

Известно, что эллипсоиды используются для аппроксимации произвольных областей. Имеют место следующие утверждения.

Теорема 1. [22]. Для любого ограниченного множества S в R^n существует единственный

эллипсоид E^* наименьшего объема, содержащий S .

Формально задача построения эллипсоида E^* наименьшего объема формулируется как оптимизационная задача отыскания вектора c^* и матрицы A^* для заданного множества S в виде

$$\begin{cases} \det A \to \min \\ S \subset E^* = E(c^*, A^*). \end{cases}$$

Пусть в R^n заданы эллипсоиды $E\left(c_1,A_1\right)$ и $E\left(c_2,A_2\right)$, их сумма $E\left(c_1,A_1\right)+E\left(c_2,A_2\right)$ есть совокупность точек $x=x^1+x^2$, таких что $x^1\in E\left(c_1,A_1\right)$ и $x^2\in E\left(c_2,A_2\right)$.

Теорема 2. [23]. Параметры эллипсоида $E\left(c^*,A^*\right)$ наименьшего объема, содержащего сумму эллипсоидов $E\left(c_1,A_1\right)+E\left(c_2,A_2\right)$, определяются формулами

$$c^* = c_1 + c_2$$
, $A^* = (p^{-1} + 1)A_1 + (p + 1)A_2$,

где p — единственный положительный корень алгебраического уравнения $\sum_{i=1}^{n} \frac{1}{p+\lambda_{i}} = \frac{n}{p(p+1)}$,

 $\lambda_i \geq 0 \left(i = \overline{1,n}\right)$ — корни характеристического уравнения $\det\left(A_1 - \lambda A_2\right) = 0$, причем каждый корень считается столько раз, какова его кратность.

Таким образом, можно сформулировать следующие свойства эллипсоидов, которые обосновывают их выбор для аппроксимации кластеров: эллипсоид в R^n определяется сравнительно небольшим числом параметров; класс эллипсоидов инвариантен относительно аффинных преобразований; в общем случае эллипсоиды позволяют получать двустороннюю аппроксимацию произвольных выпуклых множеств [22]; существует эллипсоид минимального объема, который «покрывает» заданное множество точек в R^n ; для суммы эллипсоидов существует эллипсоид минимального объема.

2.2. Построение эллипсоидов минимального объема

Алгоритмы кластеризации позволяют построить эллипсоиды в пространстве наблюдаемых данных. При наложении эллипсоида на

точки кластера наблюдается расхождение между пространством, заключенным внутри эллипсоида, и реальным геометрическим местом точек кластера. Таким образом, возникает задача определения оптимальных значений параметров эллипсоида с целью сокращения этого расхождения. По сути, необходимо определить эллипсоид минимального объема, который содержит заданный кластер. Задача нахождения эллипсоида минимального объема подробно рассматривалась в [24-27].

Рассмотрим данную задачу для одного кластера, содержащего некоторое множество точек из R^n . Пусть k — количество точек, принадлежащих кластеру; P — симметрическая положительно определенная матрица, определяющая эллипсоид E с центром в точке c объема v_E . Формальная постановка задачи имеет вид:

ная постановка задачи имеет вид.
$$\begin{cases} (\det P)^{\frac{1}{2}} \to \min \\ (P^{-1}(x_i - c), (x_i - c)) \le 1 & (i = \overline{1, k}), \\ P > 0, \end{cases}$$
 (1)

где P > 0 — положительно определенная матрица, вектор c и матрица P — неизвестные.

Введем вспомогательную задачу: вложим множество точек x_i в гиперплоскость $x_{n+1} = 1$

пространства
$$R^{n+1}$$
. Обозначим $q_i = \begin{pmatrix} x_i \\ 1 \end{pmatrix}$. В но-

вом пространстве будем искать матрицу S, определяющую эллипсоид с фиксированным центром в нуле. В [26] утверждается, что сечение такого эллипсоида плоскостью $x_{n+1} = 1$ и будет решением исходной задачи. Формальная постановка вспомогательной задачи имеет вид:

$$\begin{cases} \left(\det S\right)^{1/2} \to \min \\ \left(S^{-1}q_i, q_i\right) \le 1 \quad \left(i = \overline{1, k}\right), \\ S > 0. \end{cases}$$

Обозначая $M = S^{-1}$ и преобразуя целевую функцию к виду $-\ln(\det M) \to \min$, получим:

$$\begin{cases}
-\ln\left(\det M\right) \to \min \\
1 - \left(Mq_i, q_i\right) \ge 0, i = \overline{1, k}, \\
M > 0.
\end{cases}$$
(2)

Задача (2) является задачей выпуклого программирования с линейными ограничениями. Функция Лагранжа для данной задачи имеет вид $L(M,\lambda) = -\ln\left(\det M\right) - \sum_{i=1}^k \lambda_i \left(1-q_i M q_i^T\right)$, где λ — вектор множителей Лагранжа, $\lambda_i \geq 0$ $\left(i=\overline{1,k}\right)$, при этом ограничение M>0 будем считать нефункциональным, включив его в область определения, т.е. $dom L = \left\{(M,\lambda): M>0, \lambda \in R^k\right\}$.

Используя определение скалярного произведения Фробениуса [28] и определение производной для матрицы [29], получим следующие вспомогательные формулы для производных функций $\ln(\det(X))$ и aXa^T , где X — симметрическая положительно определенная матрица, и a — фиксированный вектор:

$$\frac{\partial \left(\ln\left(\det\left(X\right)\right)\right)}{\partial X} = X^{-1}, \ \frac{\partial \left(aXa^{T}\right)}{\partial X} = aa^{T}.$$

Теперь выпишем условия Куна-Таккера для задачи (2), используя полученные формулы. Пусть $Q = (q_1, \ldots, q_k), \Lambda = diag(\lambda),$ откуда $M^{-1} = Q\Lambda Q^T$, тогда:

$$\begin{cases}
\frac{\partial L(M,\lambda)}{\partial M} = \frac{\partial \left(-\ln\left(\det M\right)\right)}{\partial M} - \sum_{i=1}^{k} \lambda_{i} \frac{\partial \left(1 - q_{i} M q_{i}^{T}\right)}{\partial M} = \\
-M^{-1} + Q \Lambda Q^{T} = 0 \\
1 - q_{i} M q_{i}^{T} \ge 0, \quad i = \overline{1, k}, \\
\lambda_{i} \left(1 - q_{i} M q_{i}^{T}\right) = 0, \quad i = \overline{1, k}, \\
\lambda_{i} \ge 0, \quad i = \overline{1, k}.
\end{cases} \tag{3}$$

Так как задача (2) является задачей выпуклого программирования с линейными ограничениями, то условия Куна-Таккера являются необходимыми и достаточными. Можно пока-

зать, что
$$\sum_{i=1}^k \lambda_i \left(1 - q_i^T M q_i\right) = \sum_{i=1}^k \lambda_i - (n+1) = 0$$
, откуда $\sum_{i=1}^k \lambda_i = n+1$ или $\mathbf{1}^T \lambda = n+1$.

Двойственная задача Лагранжа для задачи (2) имеет следующий вид:

$$\begin{cases} \log \det \left(M^{-1} \right) \to \max \\ \mathbf{1}^T \lambda = n+1, \ \lambda \ge 0. \end{cases} \tag{4}$$

После замены переменных $u = \frac{\lambda}{n+1}$, U = diag(u) в задаче (4) получим:

$$\begin{cases} \log \det \left(QUQ^T \right) \to \max \\ \mathbf{1}^T u = 1, u \ge 0. \end{cases}$$
 (5)

На основе анализа матрицы $M^{-1} = Q\Lambda Q^T$ можно показать, что условие M>0 в задаче (2) выполняется.

Пусть u — решение задачи (5), тогда параметры минимального покрывающего эллипсоида для задачи (1) вычисляются по формулам c = Xu, $P = n \left(XUX^T - Xu\left(Xu\right)^T\right)$, где X — матрица, в столбцах которой стоят векторы наблюдаемых данных.

Для решения задачи (5) воспользуемся алгоритмом Хачияна [26], который представлен ниже. В данном алгоритме используются следующие дополнительные обозначения: $V(u) = QUQ^T$, e_j — нулевой вектор с единичной j-й координатой, EPS — требуемая точность.

Алгоритм Хачияна Inputs: q_i , $i = \overline{1,k}$; EPS Outputs: uInitialization: $u := \frac{1}{n} \cdot 1$ Do $g_i(u) := q_i V(u)^{-1} q_i$, $i = \overline{1,k}$ $j := \arg\max_i g_i(u)$, $i = \overline{1,k}$ $\Delta u := e_j - u$; $\varepsilon := \frac{g_j(u) - (n+1)}{n+1}$ $\alpha := \frac{g_j(u) - (n+1)}{(n+1)(g_j(u)-1)}$ $u := u + \alpha \Delta u$ While $\varepsilon > EPS$ Return u

3. Механизм формирования продукционных правил

Функции принадлежности термов переменных можно получить либо на основе матрицы разбиения, либо, если построен эллипсоид, используя длину проекции на оси и координаты центра кластера. Рассмотрим механизм формирования нечетких продукционных правил. Обозначим $K^{-1} = A$ и будем считать, что матрица A является положительно определенной. Если n — количество входов нечеткой системы, а p — количество выходов, то положим q = n + p.

Для формирования функций принадлежности термов, соответствующих эллипсоиду, он вписывается в параллелепипед, который затем проецируется на оси пространства состояний.

Параллелепипед имеет 2^q вершин $\left(\pm\alpha/\sqrt{\lambda_1},...,\pm\alpha/\sqrt{\lambda_q}\right)$, где α — параметр эллипсоида, заданного формулой $\left(A(x_i-c),(x_i-c)\right) \leq \alpha$ $\left(i=\overline{1,k}\right)$, в повернутой плоскости координат. Единичное собственное значение матрицы A определяет направление косинуса для каждой оси эллипсоида. Направление косинуса $\cos\gamma_{ij}$ — это угол между j—ым собственным вектором и i-ой осью эллипсоида. Проекция параллелепипеда на i-ую ось отцентрирована на величину c_i по i-ой оси и

имеет длину
$$\rho = 2\alpha \sum_{j=1}^{q} \frac{\left|\cos \gamma_{ij}\right|}{\sqrt{\lambda_{i}}}$$
. Параметры эл-

липсоида определяют положение, форму и размер нечеткого правила, значение ρ определяет носитель треугольного нечеткого числа на i -ой оси. Ориентация собственных векторов определяет размер проекции, собственные вектора являются единичными и ортогональными. Для q=2 только одно число ориентирует эллипс, и

матрица поворота имеет вид
$$P = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

, тогда проекции определяются следующим образом:

$$\rho_1 = 2\alpha \left(\frac{\left| \cos \theta \right|}{\sqrt{\lambda_1}} + \frac{\left| \sin \theta \right|}{\sqrt{\lambda_2}} \right), \, \rho_2 = 2\alpha \left(\frac{\left| \sin \theta \right|}{\sqrt{\lambda_1}} + \frac{\left| \cos \theta \right|}{\sqrt{\lambda_2}} \right).$$

Заметим, что проекции не используют всю информацию, содержащуюся в эллипсоидальных нечетких «заплатках», такую, как размер и ориентацию эллипсоидов. Улучшение характеристик можно получить за счет непосредственного использования этих характеристик.

4. Вычислительный эксперимент

Для анализа подходов к формированию базы знаний, включающей TS-правила, которая обеспечивает высокую точность аппроксимации точечно-заданной функции, был разработан программный комплекс на языке Java. Для визуализации результатов расчетов и лингвистических шкал использовались библиотека Python matplotlib и Matlab. В данном программном комплексе реализован следующий функционал: кластеризация методами нечетких ссредних и Густавсона-Кесселя; построение эллипсоидов на основе матриц ковариаций кластеров и эллипсоидов минимального объема; построение лингвистических шкал входных и/или выходной переменных с термами в форме треугольных функций принадлежности на основе эллипсоидов; построение заключений TSправил с использованием осей эллипсоидов.

Для сравнительного анализа методов построения базы знаний с базой TS-правил был проведен вычислительный эксперимент. Поскольку одно из требований заключалось в необходимости визуализации лингвистических шкал, то особое внимание было уделено обработке двумерных данных, когда кластеры строились в форме эллипсов. Анализ точности осуществлялся на основе величины среднеквадратической ошибки. Для эмпирического оценивания применялась процедура кросс-валидации по пяти блокам. В экспериментах использовалось 2000 точек, сгенерированных помощью многоэкстремальных тестовых функций путем добавления шума (всего исследовалось 8 функций).

Рассмотрим более подробно один из случаев. На Рис. 2 и 3 представлены результаты кластеризации методом Густавсона-Кесселя точек, сгенерированных на основе тестовой функции $f_4(x) = x + x \sin x, x \in [-7,5]$. Эллипсы строились с использованием матриц ковариаций кластеров и минимальных эллипсов соответственно. Заметим, что расположение кластеров такое, что

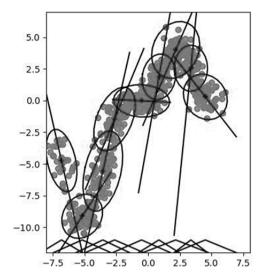


Рис. 2. База правил, построенная по матрицам ковариаций, для 9 кластеров

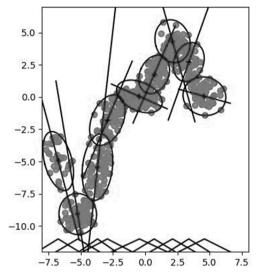


Рис. 3. База правил, построенная по минимальным эллипсам, для 9 кластеров

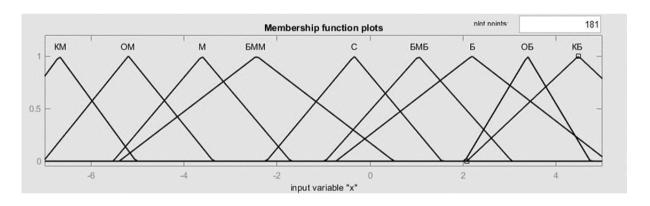


Рис. 4. Лингвистическая шкала входной переменной, сгенерированная по матрицам ковариаций

целесообразно использовать TS-правила, при этом для входной переменной нужно построить лингвистическую шкалу, термы которой используются в посылках правил, а в заключениях правил будут стоять функции, определяемые на основе осей соответствующих эллипсов.

На Рис. 4 представлена построенная на основе матрицы ковариации лингвистическая шкала.

Здесь термы лингвистической шкалы задаются в форме треугольных нечетких чисел (a_i – модальное значение; l_i , r_i – коэффициенты неопределенности, определяющие левую и правую границу соответственно). Важнейшим вопросом является интерпретация полученных

термов для входной переменной величина. В данном случае считалось, что базовое терм-множество содержит три терма $T = \left\{a_1 = \text{малая}, a_2 = \text{средняя}, a_3 = \text{большая}\right\}.$ Для порождения новых термов использовались следующие модификаторы: более или менее (БМ), очень (O), крайне (K), что позволило сформировать лингвистическую шкалу для входной переменной в виде:

 $\{KM, OM, M, EMM, C, EME, E, OE, KE\}$, где, например, терм KM означает крайне малая, терм EMM — более или менее малая и т.д. Функции принадлежности термов в форме треугольных нечетких чисел имеют следующий вид:

$$KM = \begin{pmatrix} -8.33301 \\ -6.68395 \\ -5.03489 \end{pmatrix} OM = \begin{pmatrix} -7.02779 \\ -5.19951 \\ -3.37123 \end{pmatrix} M = \begin{pmatrix} -5.52715 \\ -3.61744 \\ -1.70773 \end{pmatrix} EMM = \begin{pmatrix} -5.39600 \\ -2.43992 \\ 0.51615 \end{pmatrix}$$

$$EMB = \begin{pmatrix} -0.96315 \\ 1.04415 \\ 3.05145 \end{pmatrix} E = \begin{pmatrix} -0.73421 \\ 2.19217 \\ 5.11857 \end{pmatrix} OE = \begin{pmatrix} 2.03946 \\ 3.38994 \\ 4.74043 \end{pmatrix} KE = \begin{pmatrix} 2.08290 \\ 4.48389 \\ 6.88487 \end{pmatrix}$$

База правил, сгенерированная на основе данных, порожденных тестовой функцией f_4 , с использованием матриц ковариаций кластеров (Рис. 3), имеет следующий вид:

- [1] ecnu x = KM, mo y = -10x 71.60132;
- [2] ecnu x = OM, mo y = 10x + 42.84763;
- [3] $ec\pi u x = M$, mo y = 6.05212x + 16.35037;
- [4] ecnu x = EMM, mo y = 2.08470x + 3.92466;
- [5] $ec\pi u \ x=C, mo \ v = 0.08148x + 0.07207$;
- [6] ecnu x = EME, mo y = 3.13768x 1.19029;
- [7] ecnu x = E, mo v = 1.34315x + 1.01897;
- [8] ecnu x = OE, mo y = -10x + 36.37498;
- [9] $ec\pi u x = KB$, mo y = -1.41547x + 6.62371.

На Рис. 3 кластеры представлены эллипсами минимальной площади, которые получены с помощью алгоритма Хачияна. Здесь также девять кластеров, их проецирование на ось абсцисс позволяет получить соответственно девять термов в форме треугольных нечетких чисел, которые представлены ниже. Их интерпретация будет той же, что и в первом случае, но параметры определяются иначе. Выходная переменная в заключениях правил задается линейной функцией, определяемой на основе соответствующих осей эллипса.

База правил включает девять TS-правил следующего вида:

[1]
$$ecnu \ x=KM, mo \ y=-10x-71.70912$$
;

[2]
$$ecnu \ x=OM$$
, $mo \ y=-0.52878x-11.83678$;

[3]
$$ec\pi u \ x=M, mo \ y=10x+31.90432$$
;

[4]
$$ecnu \ x=BMM$$
, $mo \ y=2.13652x+4.64637$;

[5]
$$ec\pi u \ x=C$$
, $mo \ y=-0.63758x-0.13278$;

[6]
$$ec\pi u \ x = EME$$
, mo $v = 0.92283x + 1.16436$;

[7]
$$ec\pi u \ x=B$$
, mo $y=2.43915x+9.62634$;

[8]
$$ec\pi u \ x = OE$$
, mo $y = 4.12675x + -12.03818$;

[9]
$$ecnu \ x=KB$$
, $mo \ y=-0.31180x+1.49631$.

В Табл. 1 приведены величины среднеквадратической ошибки Err для каждой из тестовых функций для разбиения на 5, 9 и 17 кластеров.

На основе результатов проведенных экспериментов можно сделать следующие выводы:

- 1. С увеличением количества кластеров ошибка уменьшается и описание данных становится все более подробным, но, возможно, что за небольшое приращение точности придется заплатить увеличением количества правил в базе, что повлияет на быстродействие.
- 2. При увеличении количества кластеров использование эллипсов минимальной площади в основном дает небольшое улучшение точности

$$KM = \begin{pmatrix} -8.11680 \\ -6.69428 \\ -5.27176 \end{pmatrix} OM = \begin{pmatrix} -7.46619 \\ -5.30839 \\ -3.15059 \end{pmatrix} M = \begin{pmatrix} -5.08854 \\ -3.72533 \\ -2.36212 \end{pmatrix} EMM = \begin{pmatrix} -5.15348 \\ -2.89479 \\ -0.63610 \end{pmatrix}$$

$$EMB = \begin{pmatrix} -1.03436 \\ 0.89224 \\ 2.81884 \end{pmatrix} E = \begin{pmatrix} 0.33632 \\ 2.14627 \\ 3.95621 \end{pmatrix} OE = \begin{pmatrix} 2.16972 \\ 3.54588 \\ 4.92204 \end{pmatrix} KE = \begin{pmatrix} 2.55783 \\ 4.59508 \\ 6.63234 \end{pmatrix}$$

	5 кластеров		9 кластеров		17 кластеров	
Тестовая	На основе	Эллипсы	На основе	Эллипсы	На основе	Эллипсы
функция	матриц	минимальной	матриц	минимальной	матриц	минимальной
	ковариаций	площади	ковариаций	площади	ковариаций	площади
f_1	13.13689	12.55516	11.84933	11.28037	12.87200	12.01481
f_2	3.96905	3.88835	3.36676	3.55838	3.30281	3.08225
f_3	2.72498	2.62667	2.31930	2.30666	2.240366	2.16684
f_4	5.27414	4.41781	1.80160	1.69607	1.72709	1.56693
f_5	3.78032	2.44583	2.93046	2.37752	2.52319	1.87073
f_6	15.76963	15.49591	15.22970	15.02886	15.23434	14.95771
Среднее значение ошибки	7.442001	6.904955	6.249525	6.04131	6.316633	5.943212

Табл. 1. Оценка величины среднеквадратичной ошибки

при генерации правил типа Такаги-Сугено, но при малом количестве кластеров выигрыш оказывается значительным.

- 3. Одним из основных свойств базы правил является интерпретируемость, которая обеспечивается подходящим количеством правил. Поскольку каждое правило связано с определенным кластером и имеется стратегия на уменьшение количества правил, то необходим обоснованный выбор процедуры кластеризации. Вычислительный эксперимент продемонстрировал преимущества метода Густавсона-Кесселя по сравнению с методом с-средних. В обоих методах изначально задается количество кластеров. Если данные визуализированы (например, исходные данные это временные ряды), то эта проблема решаема.
- 4. Если у полученных эллипсов большая ось значительно превосходит меньшую, то целесообразно использовать правила в виде ТЅмодели, иначе лингвистическую модель. Этим выбором обеспечивается улучшение качества аппроксимации.
- 5. Если получены эллипсоиды минимального объема, то можно сократить количество правил за счет сложения соответствующих эллипсоидов.
- 6. Установлено, что функции принадлежности термов, полученные при проецировании эллипсов на ось абсцисс для некоторых тестовых функций, сильно перекрывают друг друга,

что ощутимо сказывается на точности аппроксимации. Дальнейшая оптимизация может быть связана с объединением сильно перекрывающихся термов, что требует использования специальных операций над нечеткими множествами, например, параметрических треугольных норм и конорм [30, 31].

Заключение

База знаний, основу которой составляет база правил, является настраиваемым компонентом любой нечеткой системы. От ее качества зависит успех решения конкретной прикладной задачи. По сути, можно выделить два основных подхода к формированию базы правил – с привлечением эксперта, опыт которого тиражируется базой знаний, и на основе обучающей выборки. Эллипсоидальная кластеризация позволяет выделить группы однородных данных, которые можно описать либо с помощью лингвистической модели, либо ТS-модели. В качестве перспективных направлений проведенного исследования выделим следующие:

- формирование базы правил, включающей правила, относящиеся к различным типам, и построение механизмов логического вывода, которые работают с такой базой;
- исследование влияния минимизации площади эллипса на точность аппроксимации с помощью правил лингвистического типа.

Литература

- Sunardi. Tsukamoto Fuzzy Inference System on Internet of Things-Based for Room Temperature and Humidity Control / Sunardi, A. Yudhana, Furizal // IEEE Access. – 2023. – Vol. 11. – Pp. 6209-6227.
- Lin, J. Fuzzy PID Control for Multi-joint Robotic Arm / Lin, X. Liu, Z. Ren // 2022 IEEE 20th International Conference on Industrial Informatics (INDIN), Perth, Australia. – 2022. – Pp. 723-728.
- Санаева Г. Н. Иерархическая система нечеткого регулирования процесса получения ацетилена окислительным пиролизом природного газа / Г. Н. Санаева, А. Е. Пророков, В. Н. Богатиков, Д. П. Вент // Вестник Астраханского государственного технического университета. Серия: Управление, вычислительная техника и информатика. 2020. № 1. С. 7-17.
- Рябчиков М. Ю. Система управления температурой пара после пароперегревательной установки с применением нечеткой логики для упреждающей компенсации возмущений / М. Ю. Рябчиков, Е. С. Рябчикова, С. А. Филиппов // Мехатроника, автоматизация, управление. – 2021. – Т. 22. – № 4. – С. 181-190.
- Нифантов В. М. Диагностика, оценка и прогнозирование технического состояния технологического оборудования при помощи нечеткой экспертной системы в централизованных системах технического обслуживания и ремонта / В. М. Нифантов // Труды Кольского научного центра РАН. – 2019. – Т. 10. – № 9. – С. 187-197.
- 6. Леденева Т. М. Нечеткое моделирование медицинских экспертных систем / Т. М. Леденева, С. Л. Подвальный, Р. К. Стрюков, Дегтярев С. В. // Биомедицинская радиоэлектроника, 2016. № 9. С. 16-24.
- Дубенко Ю. В. Нечеткая система определения оптимальных методов для прогнозирования параметров сложных технических систем / Ю. В. Дубенко, Е. Е. Дышкант // Известия высших учебных заведений. Поволжский регион. Технические науки. 2018. Т. 47. № 3. С. 58-69.
- Hoang, T. -M. An efficient IDS using FIS to detect DDoS in IoT networks / T. -M. Hoang, N. -H. Tran, V. -L. Thai, D. -L. Nguyen, N. -H. Nguyen // 2022 9th NAFOSTED Conference on Information and Computer Science (NICS), Ho Chi Minh City, Vietnam. – 2022. – Pp. 193-198.
- Дикарев П. В. Система распознавания аварийных режимов воздушных линий электропередачи с использованием нечеткой логики / П. В. Дикарев, А. А. Шилин, С. Ю. Юдин // Энерго- и ресурсосбережение: промышленность и транспорт. 2022. Т. 38. № 1. С. 6-12.
- 10. Тутыгин В. С. Система распознавания болезней растений по изображениям листьев на основе нечеткой логики и нейронных сетей / В. С. Тутыгин, Б. Х. М. А. Аль Винди, И. А. Рябцев // Современная наука: актуальные проблемы теории и практики. Серия: Естественные и технические науки. 2019. № 3. С. 107-115.
- 11. Пегат А. Нечеткое моделирование и управление: Пер. с англ. / А. Пегат. М.: БИНОМ. Лаборатория знаний, 2009. 798 с. (Piegat, A. Fuzzy Modeling and Control. Physica-Verlag Heidelberg, 2001. 728 р.)

- Preuss H. P. Neuro-fuzzy / H. P. Preuss, V. Tresp // Automatisierungstechnische Praxis. 1994. Vol. 36. № 5. Pp. 10-24.
- 13. Luukka P. A new non-linear fuzzy robust PCA Algorithm and similarity classifier in Classification of medical data sets / P. Luukka // International Journal of Fuzzy Systems. 2011. Vol. 13. Pp. 153-162.
- 14. Сергиенко М. А. Методы проектирования нечеткой базы знаний / М. А. Сергиенко // Вестник Воронежского государственного университета. Серия: системный анализ и информационные технологии. 2008. № 2. С. 67-71.
- Bezdek J. C. FCM: the Fuzzy c-means clustering algorithm / J. Bezdek, R. Ehrlich, W. Full // Computers and Geosciences. – Dec. 1984. – Vol. 10. – Pp. 191-203.
- Krishnapuram R. A possibilistic approach to clustering / R. Krishnapuram, J. M. Keller // Fuzzy Systems. – 1993. – Vol. 1. – No 2. – Pp. 98-110.
- Xu Z. Clustering algorithm for intuitionistic fuzzy sets / Z.
 Xu, J. Chen, J. Wu. // Inf. Sci. 2008. Vol. 178. –
 Pp. 3775-3790.
- 18. Ji Z. Generalized rough fuzzy c-means algorithm for brain MR image segmentation / Z. Ji, Q. Sun, Y. Xia, Q. Chen, D. Xia, D. D. Feng // Computer methods and programs in biomedicine. – 2012. – Vol. 108. – Issue 2. – Pp. 644-655.
- Hwang C. Uncertain Fuzzy Clustering: Interval Type-2 Fuzzy Approach to C -Means / C. Hwang and F. C.-H. Rhee // IEEE Transactions on Fuzzy Systems. – 2007. – Vol. 15. – No. 1. – Pp. 107-120.
- 20. Das S. A new kernelized fuzzy C-means clustering algorithm with enhanced performance / S. Das, H. K. Baruah // International Journal of Research in Advent Technology. 2014. Vol. 2. № 6. Pp. 43-51.
- Dickerson J. A. Fuzzy function approximation with ellipsoidal rules / J. A. Dickerson, B. Kosko // IEEE Transaction on fuzzy systems. August 2004. Vol. 26. № 4. Pp. 542-560.
- 22. Грюнбаум Б. Этюды по комбинаторной геометрии и теории выпуклых тел: Пер. с англ. / Б. Грюнбаум. М.: Наука, 1971. 349 с. (Grünbaum B. Measures of symmetry for convex sets. Proc. Sympos. Pure Math., Providence, USA, 1963. Vol. 7. Pp. 233-270.; Grünbaum B. Borsuk's problem and related questions. Proc. Sympos. Pure Math., Providence, USA, 1963. Vol. 7. Pp. 271-284.)
- Черноусько Ф. Л. Оценивание фазового состояния динамических систем. Метод эллипсоидов / Ф. Л. Черноусько. М.: Наука: Главная редакция физикоматематической литературы, 1988. 320 с.
- 24. Raposo A. A. M. A construction of the minimum volume ellipsoid containing a set of points using BRKGA metaheuristic / A. A. M. Raposo, V. M. de Souza, L. R. A. G. Filho // Proceeding Series of the Brazilian Society of Computational and Applied Mathematics. − 2021. − Vol. 8. − №. 1. − Pp. 010339.
- 25. Sun P. Computation of minimum-volume enclosing ellipsoids / P. Sun, R. Freund // Oper. Res. 2004. Vol. 52. No. 5. Pp. 690-706.

- 26. Kumar P. Minimum-volume enclosing ellipsoids and core sets / P. Kumar, E. Yildirim // Journal of Optimization Theory and Applications. 2005. Vol. 126. No. 1. Pp. 1-21.
- Todd M. J. On Khachiyan's algorithm for the computation of minimum volume enclosing ellipsoids / M. J. Todd, E. A. Yıldırım // Discrete Appl. Math. – 2007. – Vol. 155. – Pp. 1731-1744.
- 28. Foucart S. A Mathematical Introduction to Compressive Sensing / S. Foucart, H. Rauhut. NY: Birkhäuser New York, 2013. 625 p.
- 29. Magnus J. R. Matrix differential calculus with applications in statistics and econometrics / J. R. Magnus, H. Neudecker. John Wiley & Sons, 1988. 504 p.
- 30. Ledeneva T. M. New Family of Triangular Norms for Decreasing Generators in the Form of a Logarithm of a Linear Fractional Function / T. M. Ledeneva // Fuzzy sets and systems. 2022. № 427. Pp. 37-54.
- 31. Ledeneva T. M. Additive generators of fuzzy operations in the form of a Linear Fractional Function / T. M. Ledeneva // Fuzzy sets and systems. 2020. № 386. Pp. 1-24.

Моисеева Татьяна Александровна. Федеральное государственное бюджетное образовательное учреждение высшего образования «Воронежский государственный университет» (ФГБОУ ВО ВГУ), г. Воронеж. Преподаватель кафедры математического обеспечения ЭВМ, аспирант 3-го года обучения кафедры вычислительной математики и прикладных информационных технологий. Количество печатных работ: 11. Область научных интересов: нечеткие системы, кластеризация. E-mail: tatiana.vsu@gmail.com

Леденева Татьяна Михайловна. Федеральное государственное бюджетное образовательное учреждение высшего образования «Воронежский государственный университет» (ФГБОУ ВО ВГУ), г. Воронеж. Зав. кафедрой вычислительной математики и прикладных информационных технологий, д-р. техн. наук, профессор. Количество печатных работ: 220 (в т. ч. 3 монографии). Область научных интересов: моделирование нечетких логических связок, нечеткие системы, моделирование целенаправленных систем. E-mail: ledeneva-tm@yandex.ru

Knowledge Base Generation Based on Fuzzy Clustering

T. A. Moiseeva, T. M. Ledeneva

Voronezh State University, Voronezh, Russia

Abstract. The article states fuzzy Takagi-Sugeno rule base generation problem based on ellipsoidal clustering. After obtaining clusters of ellipsoidal shape the problem of building minimal volume ellipsoids, enclosing all clusters points, appears. The premises of the generated fuzzy rules are formed by constructing projections of ellipsoids on the coordinate axis, and conclusions — either using ellipsoid axes, or based on the projection. In the article, the authors suggest to use Khachiyan's algorithm for building minimal volume enclosing ellipsoid in order to increase the accuracy of approximation and they compare two approaches of choosing optimal parameters of ellipsoids which enclose all clusters points.

Keywords: clustering algorithms, "if-then" rules, knowledge base, fuzzy systems.

DOI 10.14357/20718632230110

References

- Sunardi, A. Yudhana, Furizal. 2023. Tsukamoto Fuzzy Inference System on Internet of Things-Based for Room Temperature and Humidity Control. IEEE Access. 11:6209-6227.
- Lin, J., X. Liu, Z. Ren. 2022. Fuzzy PID Control for Multi-joint Robotic Arm. 2022 IEEE 20th International Conference on Industrial Informatics (INDIN), Perth, Australia. 723-728.
- Sanaeva, G. N., A. E. Prorokov, V. N. Bogatikov, D. P. Vent. 2020. Ierarhicheskaya sistema nechetkogo regulirovaniya processa polucheniya acetilena okislitelnym pirolizom prirodnogo gaza [Hierarchical system of fuzzy regulation of acetylene production process by oxidative pyrolysis of natural gas]. Vestnik Astrahanskogo gosudarstvennogo tekhnicheskogo universiteta. Seriya: Uprav-
- lenie, vychislitelnaya tekhnika i informatika [Proceedings of Astrakhan State Technical University. Series: Management, computer science and informatics]. 1:7-17.
- Ryabchikov, M. Yu., E. S. Ryabchikova, S. A. Filippov. 2021. Sistema upravleniya temperaturoj para posle paroperegrevatelnoj ustanovki s primeneniem nechetkoj logiki dlya uprezhdayushchej kompensacii vozmushchenij [A fuzzy logic-based system for controlling the temperature of steam exiting a superheater for the purpose of preemptive perturbation compensation]. Mekhatronika, avtomatizaciya, upravlenie. 22(4):181-190
- 5. Nifantov, V. M. 2019. Diagnostika, ocenka i prognozirovanie tekhnicheskogo sostoyaniya tekhnologicheskogo oborudovaniya pri pomoshchi nechetkoj ekspertnoj sistemy v centralizovannyh sistemah tekhnicheskogo obsluzhivaniya i remonta [Diagnostics, assessment and forecasting of technical condition of technological equipment

- by using fuzzy expert system in centralized maintenance and repair systems]. Trudy Kolskogo nauchnogo centra RAN [Transactions Kola Science Centre]. 10(9):187-197.
- Ledeneva, T.M., S.L. Podvalny, R.K. Stryukov, S.V. Degtyarev. 2016. Nechetkoe modelirovanie medicinskih ekspertnyh sistem [Fuzzy modelling medical expert systems]. Biomedicinskaya radioelektronika [Biomedicine radioengineering]. 9:16-24.
- Dubenko, Yu. V., E. E. Dyshkant. 2018. Nechetkaya sistema opredeleniya optimalnyh metodov dlya prognozirovaniya parametrov slozhnyh tekhnicheskih sistem [Fuzzy system for determining optimal methods to forecast parameters of complex technical systems]. Izvestiya vysshih uchebnyh zavedenij. Povolzhskij region. Tekhnicheskie nauki [University proceedings. Volga region. Technical sciences]. 47(3):58-69.
- Hoang, T. -M., N. -H. Tran, V. -L. Thai, D. -L. Nguyen and N. -H. Nguyen. 2022. An efficient IDS using FIS to detect DDoS in IoT networks. 9th NAFOSTED Conference on Information and Computer Science (NICS), Ho Chi Minh City, Vietnam. 193-198.
- Dikarev, P. V., A. A. Shilin, S. Yu. Yudin. 2022. Sistema raspoznavaniya avarijnyh rezhimov vozdushnyh linij elektroperedachi s ispolzovaniem nechetkoj logiki [System for recognition of emergency modes of overhead power lines using fuzzy logic]. Energo- i resursosberezhenie: promyshlennost i transport [Energy and resource saving: industry and transport]. 38(1):6-12.
- 10. Tutygin, V. S., B. H. M. A. Al Vindi, I. A. Ryabcev. Sistema raspoznavaniya boleznej rastenij po izobrazheniyam listev na osnove nechetkoj logiki i nejronnyh setej [System of recognition of plant diseases on leaves images on the basis of fuzzy logic and neural network]. Sovremennaya nauka: aktualnye problemy teorii i praktiki. Seriya: Estestvennye i tekhnicheskie nauki [Modern Science: actual problems of theory and practice. Series "Natural and Technical Sciences]. 3:107-115.
- Piegat, 2001. A. Fuzzy Modeling and Control. Physica-Verlag Heidelberg. 728 p.
- 12. Preuss, H. P., Tresp, V. 1994. Neuro-fuzzy. Automatisier-ungstechnische Praxis. 36(5):10-24.
- Luukka, P. 2011. A new non-linear fuzzy robust PCA Algorithm and similarity classifier in Classification of medical data sets. International Journal of Fuzzy Systems. 13:153-162.
- 14. Sergienko M.A. 2008. Metody proektirovaniya nechetkoj bazy znanij [Designing methods of hierarchical knowledge base]. Vestnik Voronezhskogo gosudarstvennogo universiteta. Seriya: sistemnyj analiz i informacionnye tekhnologii [Proceedings of Voronezh State University. Series: Systems Analysis and Information Technologies]. 2:67-71.
- Bezdek, J. C., Ehrlich, R., Full, W. 1984. FCM: the Fuzzy c-means clustering algorithm. Computers and Geosciences. 10:191-203.

- Krishnapuram, R., J. M. Keller. 1993. A possibilistic approach to clustering. Fuzzy Systems. 1(2):98-110.
- Xu, Z., J. Chen, J. Wu. 2008. Clustering algorithm for intuitionistic fuzzy sets. Inf. Sci. 178:3775-3790.
- Ji, Z., Q. Sun, Y. Xia, Q. Chen, D. Xia, D. D. Feng. 2012. Generalized rough fuzzy c-means algorithm for brain MR image segmentation. Computer methods and programs in biomedicine. 108(2):644-655.
- Hwang, C., F. C.-H. Rhee. 2007. Uncertain Fuzzy Clustering: Interval Type-2 Fuzzy Approach to C -Means. IEEE Transactions on Fuzzy Systems. 15(1):107-120.
- Das, S., H. K. Baruah. 2014. A new kernelized fuzzy C-means clustering algorithm with enhanced performance. International Journal of Research in Advent Technology. 2(6):43-51.
- Dickerson, J. A., B. Kosko. 2004. Fuzzy function approximation with ellipsoidal rules. IEEE Transaction on fuzzy systems. August 2004. 26(4):542-560.
- Grünbaum, B. 1963. Measures of symmetry for convex sets. Proc. Sympos. Pure Math., Providence, USA. 7:233-270.; Grünbaum, B. 1963. Borsuk's problem and related questions. Proc. Sympos. Pure Math., Providence, USA. 7:271-284.
- 23. Chernoysko, F. L. 1988. Ocenivanie fasovogo sostoyaniya dinamicheskih system [Estimation of the phase state of dynamic systems]. Moscow: Science: The main editorial office of the physical and mathematical literature. 320 p.
- 24. Raposo, A. A. M., V. M. de Souza, L. R. A. G. Filho. 2021. A construction of the minimum volume ellipsoid containing a set of points using BRKGA metaheuristic. Proceeding Series of the Brazilian Society of Computational and Applied Mathematics. 8(1):010339.
- 25. Sun, P., R. Freund. 2004. Computation of minimum-volume enclosing ellipsoids .Oper. Res. 52(5):690–706.
- 26. Kumar, P., E. Yildirim. 2005. Minimum-volume enclosing ellipsoids and core sets. Journal of Optimization Theory and Applications. 126(1):1–21.
- 27. Todd, M. J., E. A. Yıldırım. 2007. On Khachiyan's algorithm for the computation of minimum volume enclosing ellipsoids. Discrete Appl. Math. 155:1731–1744.
- 28. Foucart, S., H. Rauhut. 2013. A Mathematical Introduction to Compressive Sensing. Birkhäuser New York, NY. 625 p.
- 29. Magnus, J. R., H. Neudecker. 1988. Matrix differential calculus with applications in statistics and econometrics. John Wiley & Sons. 504 p.
- Ledeneva, T.M. 2022. New Family of Triangular Norms for Decreasing Generators in the Form of a Logarithm of a Linear Fractional Function. Fuzzy sets and systems. 427:37-54.
- 31. Ledeneva, T.M. 2020. Additive generators of fuzzy operations in the form of a Linear Fractional Function. Fuzzy sets and systems. 386:1-24.

Moiseeva T. A. Voronezh State University, 1 Universitetskaya pl., Voronezh, 394018, Russia. E-mail: tatiana.vsu@gmail.com

Ledeneva T. M. Doctor of Technical Sciences, Professor, Voronezh State University, 1 Universitetskaya pl., Voronezh, 394018, Russia. E-mail: ledeneva-tm@yandex.ru