



Є.В. Бодянський¹, А.Ю. Шафроненко², І.М. Климова³

¹Доктор технических наук, профессор кафедры Искусственного интеллекта

Харьковский национальный университет радиоэлектроники
yevgeniy.bodyanskiy@nure.ua ORCID 0000-0001-5418-2143

²Кандидат технических наук, доцент кафедры Информатики
Харьковский национальный университет радиоэлектроники
alina.shafronenko@nure.ua ORCID 0000-0002-8040-0279

³Ассистент кафедры Системотехники
Харьковский национальный университет радиоэлектроники
iryana.klymova@nure.ua ORCID 0000-0003-0455-6180

ОНЛАЙН ДОСТОВІРНА НЕЧІТКА КЛАСТЕРІЗАЦІЯ ДАНИХ З ВИКОРИСТАННЯМ ФУНКЦІЇ НАЛЕЖНОСТІ СПЕЦІАЛЬНОГО ТИПУ

Предложен онлайн метод достоверной нечеткой кластеризации, предназначенный для анализа данных, последовательно поступающих на обработку. Особенностью развиваемого подхода является использование функции принадлежности специального вида, описываемой функции плотности распределения Коши. Собственно процедура уточнения центроидов кластеров является по сути правилом самообучения «Победитель получает больше» (WTM), в котором функция соседства порождается введенной функцией принадлежности.

НЕЧЕТКАЯ КЛАСТЕРИЗАЦИЯ, ОБРАБОТКА, ГРАДИЕНТНАЯ ОПТИМИЗАЦИЯ, ФУНКЦИЯ ПРИНАДЛЕЖНОСТИ.

Запропоновано онлайн метод достовірної нечіткої кластеризації, призначений для аналізу даних, послідовно надходять на обробку. Особливістю розвиваючого підходу є використання функції належності спеціального виду, описуваної функції щільності розподілу Коші. Власне процедура уточнення центроїд кластерів є по суті правилом самонавчання «Переможець отримує більше» (WTM), в якому функція сусідства породжується введеною функцією приналежності.

НЕЧІТКА КЛАСТЕРИЗАЦІЯ, ОБРОБКА, ГРАДІЄНТНА ОПТИМІЗАЦІЯ, ФУНКЦІЯ НАЛЕЖНОСТІ.

An online method of reliable fuzzy clustering is proposed, designed to analyze data sequentially received for processing. A feature of the developed approach is the use of the membership function of a special kind described by the density function of the Cauchy distribution. The actual procedure for clarifying the centroids of clusters is essentially a self-learning rule “The Winner Takes More” (WTM), in which the neighborhood function is generated by the introduced membership function.

FUZZY CLUSTERING, PROCESSING, GRADIENT OPTIMIZATION, PROCESSING, ACCESSORIES FUNCTION.

Вступ

Задача кластеризації (класифікації в режимі самонавчання) багатовимірних даних є важливою частиною інтелектуального аналізу даних (Data Mining), в рамках якої склався ряд напрямків і підходів [1, 2]. Один з таких напрямків утворюють методи нечіткої (фаззі-) кластеризації, в основі яких лежить припущення що класи-кластери, які формуються, взаємно перетинаються так, що кожен вектор-спостереження з різними рівнями належності-ймовірності-можливості може належати одночасно до декількох класів одночасно.

Тут найбільш широкого поширення набули алгоритми ймовірнісної нечіткої кластеризації і, перш за все, метод нечітких с-середніх (FCM) [3, 4]. Цей підхід обмежується ймовірнісними обмеженнями на рівні належності так, що «забруднені» збуреннями і викидами спостереження можуть бути віднесені до різних класів з практично однаковими рівнями належності.

У зв'язку з цим в [5] був запропонований можливісний підхід до нечіткої кластеризації (PCM) більш стійкий до шумів і збурень. Разом з тим PCM-алгоритми

страждають від, так званої, проблеми співпадіння, коли в процесі обробки інформації деякі кластери починають зливатися один з одним, що в результаті веде до невірної оцінки кількості сформованих кластерів.

Цих недоліків позбавлені алгоритми достовірної нечіткої кластеризації [6-8], засновані на апараті теорії достовірності [9]. В рамках цього підходу в процесі розрахунків оцінюються не тільки рівні нечіткої належності, але і рівні довіри, що засновані на мірі належності спеціального виду [10]. Результати експериментів показали [7, 8], що достовірний підхід забезпечує більш високу якість кластеризації в порівнянні з ймовірнісними і можливісними методами.

Вихідною інформацією для рішення задачі нечіткої кластеризації є масив n -вимірних векторів спостережень $X = \{x_1, x_2, \dots, x_N\} \subset R^n$, $x(k) \in X$, $k = 1, 2, \dots, N$, який повинен бути розбитий на m класів-кластерів з деяким рівнем належності-можливості-достовірності $U_q(k)$ k -го вектора x_k - q му кластеру ($1 < m < N, 1 \leq q \leq m$). Необхідно також відзначити, що вихідні дані попередньо передоброблені так, що $-1 \leq x_{ki} \leq 1$ ($1 \leq i \leq n$), де x_{ki} i -та компонента вектора x_k .

Таким чином, завдання кластеризації вирішується в пакетному режимі, коли весь масив даних обробляється багаторазово на основі почергового оцінювання [8]. Якщо ж дані надходять на обробку в вигляді потоку або утворюють надвеликі масиви, пакетний режим не дозволяє ефективно вирішити цю задачу.

У цій ситуації найбільш ефективними є рекурентні процедури нечіткої кластеризації, що дозволяють вирішувати задачу в online режимі, уточнюючи результати в міру надходження кожного нового спостереження. Так, в [11,12] були запропоновані рекурентні варіанти FCM, які є по суті градієнтними процедурами оптимізації прийнятої цільової функції, а в [13,14] були введені рекурентні модифікації РСМ, що призначені для послідовної обробки даних.

У зв'язку з цим є доцільною розробка рекурентної модифікації методу достовірної нечіткої кластеризації, що дозволяє уточнювати характеристики кластерів у міру надходження кожного нового спостереження.

1. Рекурентний метод достовірної нечіткої кластеризації (RCCM)

Найбільш популярний метод ймовірнісної нечіткої кластеризації пов'язаний з мінімізацією цільової функції [4]

$$E(U_q(k), w_q) = \sum_{k=1}^N \sum_{q=1}^m U_q^\beta(k) D^2(x_k, w_q) \quad (1)$$

за обмежень $\sum_{q=1}^m U_q(k) = 1, 0 < U_q(k) < N$. Вирішуючи задачу нелінійного програмування за допомогою методу невизначених множників Лагранжа, приходимо до відомого результату

$$U_q^{(\tau+1)}(k) = \left(D^2(x_k, w_q^{(\tau)}) \right)^{\frac{1}{1-\beta}} \left(\sum_{l=1}^m \left(D^2(x_k, w_l^{(\tau)}) \right)^{\frac{1}{1-\beta}} \right)^{-1}, \quad (2)$$

$$w_q^{(\tau+1)} = \sum_{k=1}^N \left(U_q^{(\tau+1)}(k) \right)^\beta x_k \left(\sum_{k=1}^N \left(U_q^{(\tau+1)}(k) \right)^\beta \right)^{-1}, \quad (3)$$

де $U_q(k)$ -рівень належності векторного спостереження x_k q -му кластеру $Cl_q (1 \leq q \leq m)$, w_q — прототип-центроїд q -го кластера, $\beta > 1$ фаззифікатор, що визначає «розмитість» границь між класами, $D(x_k, w_q)$ — відстань між x і w_q в прийнятій метриці, $\tau = 0, 1, 2, \dots$ — індекс епохи обробки інформації в режимі попереминого оцінювання. При цьому процес обчислень триває до виконання умови

$$w - w \leq \varepsilon \forall 1 \leq q \leq m,$$

де ε — наперед заданий поріг точності обчислень.

В разі $\beta = 2$ і евклідової метрики

$$D^2(x_k, w_q) = \|x_k - w_q\|_2^2,$$

приходимо до популярного алгоритму нечітких с-середніх (FCM) [13] виду

$$U_q^{(\tau+1)}(k) = \left\| x_k - w_q^{(\tau)} \right\|^{-2} \left(\sum_{l=1}^m \left\| x_k - w_l^{(\tau)} \right\|^{-2} \right)^{-1}, \quad (4)$$

$$w_q^{(\tau+1)} = \sum_{k=1}^N \left(U_q^{(\tau+1)}(k) \right)^2 x_k \left(\sum_{k=1}^N \left(U_q^{(\tau+1)}(k) \right)^2 \right)^{-1}. \quad (5)$$

Якщо дані надходять на обробку послідовно в онлайн режимі, задача нелінійного програмування може бути вирішена за допомогою алгоритму Ерроу-Гурвіца-Удзави, що є за суттю градієнтною процедурою пошуку сідлової точки функції Лагранжа на основі критерія (1) з обмеженнями на суму належностей.

При цьому співвідношення (2), (3) можуть бути переписані у формі

$$\begin{cases} U_q(k+1) = \left(D^2(x_{k+1}, w_q^{(k)}) \right)^{\frac{1}{1-\beta}} \left(\sum_{l=1}^m \left(x_{k+1}, w_l^{(k)} \right)^{\frac{1}{1-\beta}} \right)^{-1}, \\ w_q(k+1) = w(k) + \eta(k+1) U_q^\beta(k+1) (x_{k+1} - w_q(k)), \end{cases} \quad (6)$$

(тут $\eta(k)$ — параметр кроку навчання, а (4), (5)

$$\begin{cases} U(k+1) = \|x - w(k)\| \left(\|x - w(k)\| \right), \\ w(k+1) = w + \eta(k+1) U(k+1) (x - w(k)), \end{cases} \quad (7)$$

що є узагальненням рекурентних процедур Парка-Деггера [11] і Чанга-Лі [12].

Можливісні алгоритми нечіткої кластеризації засновані на мінімізації цільової функції [5]

$$E(U(k), w, \mu) = U(k) D(x, w) + \mu (1 - U(k)), \quad (8)$$

де $\mu_q \geq 0$ визначає відстань, на якій рівень належності приймає значення 0,5, тобто. $U_q(k) = 0$, якщо $D^2(x_k, w_q) = \mu$.

Мінімізація критерія (8) дозволяє отримати аналітичний розв'язок у вигляді

$$U(k) = \left(1 + \left(\frac{D(x, w)}{\mu} \right) \right)^{-1}, \quad (9)$$

$$w_q^{(\tau+1)} = \sum_{k=1}^N \left(U_q^{(\tau+1)}(k) \right)^\beta x_k \left(\sum_{k=1}^N \left(U_q^{(\tau+1)}(k) \right)^\beta \right)^{-1}, \quad (10)$$

$$\mu_q^{(\tau+1)} = \sum_{k=1}^N \left(U_q^{(\tau+1)}(k) \right)^\beta D^2(x_k, w_q^{(\tau+1)}) \left(\sum_{k=1}^N \left(U_q^{(\tau+1)}(k) \right)^\beta \right)^{-1}, \quad (11)$$

який у квадратичному випадку набуває форму

$$U_q^{(\tau+1)}(k) = \left(1 + \frac{\|x_k - w_q^{(\tau)}\|^2}{\mu_q^{(\tau)}} \right)^{-1}, \quad (12)$$

$$w_q^{(\tau+1)} = \sum_{k=1}^N \left(U_q^{(\tau+1)}(k) \right)^2 x_k \left(\sum_{k=1}^N \left(U_q^{(\tau+1)}(k) \right)^2 \right)^{-1}, \quad (13)$$

$$\mu_q^{(\tau+1)} = \sum_{k=1}^N \left(U_q^{(\tau+1)}(k) \right)^2 \|x_k - w_q^{(\tau+1)}\|^2 \left(\sum_{k=1}^N \left(U_q^{(\tau+1)}(k) \right)^2 \right)^{-1}. \quad (14)$$

Онлайн версії (9)-(14) при цьому мають вигляд [13,14]

$$\begin{cases} U_q(k+1) = \left(1 + \left(\frac{D^2(x_{k+1}, w_q(k))}{\mu_q(k)} \right)^{\frac{1}{\beta-1}} \right)^{-1}, \\ w_q(k+1) = w_q(k) + \eta(k+1) U_q^\beta(k+1) (x_{k+1} - w_q(k)), \\ \mu_q(k+1) = \sum_{p=1}^{k+1} U_q^\beta(p) D^2(x_p, w_q(k+1)) \left(\sum_{p=1}^{k+1} U_q^\beta(p) \right)^{-1} \end{cases} \quad (15)$$

і (при $\beta = 2$)

$$\begin{cases} U_q(k+1) = \left(1 + \frac{\|x_{k+1} - w_q(k)\|^2}{\mu_q(k)} \right)^{-1}, \\ w_q(k+1) = w_q(k) + \eta(k+1) U_q^2(k+1) (x_{k+1} - w_q(k)), \\ \mu_q(k+1) = \sum_{p=1}^{k+1} U_q^2(p) \|x_p - w_q(k+1)\|^2 \left(\sum_{p=1}^{k+1} U_q^2(p) \right)^{-1}. \end{cases} \quad (16)$$

Достовірність нечітка кластеризація пов'язана з мінімізацією цільової функції

$$E(Cr_q(k), w_q) = \sum_{k=1}^N \sum_{q=1}^m Cr_q^\beta(k) D^2(x_k, w_q) \quad (17)$$

за обмежень

$$\begin{aligned} 0 \leq Cr_q(k) \leq 1 \forall q, k; \sup Cr_q(k) \geq 0,5 \forall k; \\ Cr_q(k) + \sup Cr_l(k) = 1 \end{aligned}$$

для будь-яких q і k , для яких $Cr_q(k) \geq 0,5$. Тут $Cr_q(k)$ – достовірність того, що спостереження x_k належить кластеру Cl_q . При цьому рівень достовірності розраховується на основі функції приналежності [15]

$$U(k) = \varphi(D(x, w)), \quad (18)$$

що задовольняє умовам:

- $\varphi_q(\cdot)$ монотонно зменшується на інтервалі $[0, \infty]$,
- $\varphi_q(0) = 1$,
- $\varphi_q(\infty) \rightarrow 0$.

Нескладно помітити, що функція (18) є за суттю мірою подібності, заснований на відстані [16].

В якості такої функції в [15] було запропоновано використовувати вираз

$$U_q(k) = \left(1 + D^2(x_k, w_q) \right)^{-1}, \quad (19)$$

що є звичайною дзвонуватою функцією належності, яка використовується в системах нечіткого висновування.

Цікаво зауважити, що вираз (2) може бути перепи-саний у формі

$$\begin{aligned} U_q(k) &= \left(D^2(x_k, w_q(k)) \right)^{\frac{1}{1-\beta}} \left(\sum_{l=1}^m \left(D^2(x_k, w_l(k)) \right)^{\frac{1}{1-\beta}} \right)^{-1} = \\ &= \left(D^2(x_k, w_q(k)) \right)^{\frac{1}{1-\beta}} \left(D^2(x_k, w_q(k)) \right)^{\frac{1}{1-\beta}} + \sum_{l \neq q}^m \left(D^2(x_k, w_l(k)) \right)^{\frac{1}{1-\beta}} \right)^{-1} = (20) \\ &= \left(1 + \left(D^2(x_k, w_q(k)) \right)^{\frac{1}{1-\beta}} \sum_{l \neq q}^m \left(D^2(x_k, w_l(k)) \right)^{\frac{1}{1-\beta}} \right)^{-1}, \end{aligned}$$

а для евклідової метрики і $\beta = 2$ приймає вид функції щільності розподілу Коші з параметром ширини σ_q^2 [16]

$$U_q(k) = \left(1 + \frac{\|x_k - w_q(k)\|^2}{\sigma_q^2} \right)^{-1}, \quad (21)$$

$$\sigma_q^2 = \left(\sum_{l=1}^m \|x_k - w_l(k)\|^2 \right)^{-1}. \quad (22)$$

Нескладно бачити, що функція належності (19) є окремим випадком (21) при $\sigma_q^2 = 1$.

Остаточний пакетний алгоритм достовірної нечіткої кластеризації може бути записаний у формі [7,8]:

$$U_q^{(\tau+1)}(k) = \left(1 + D^2(x_k, w_q^{(\tau)}) \right)^{-1}, \quad (23)$$

$$U_q^{*(\tau+1)}(k) = U_q^{(\tau+1)}(k) \left(\sup U_l^{(\tau+1)}(k) \right)^{-1}, \quad (24)$$

$$Cr_q^{(\tau+1)}(k) = \frac{1}{2} \left(U_q^{*(\tau+1)}(k) + 1 - \sup_{l \neq q} U_l^*(k) \right), \quad (25)$$

$$w_q^{(\tau+1)} = \sum_{k=1}^N \left(Cr_q^{(\tau+1)}(k) \right)^\beta x_k \left(\sum_{k=1}^N \left(Cr_q^{(\tau+1)}(k) \right)^\beta \right)^{-1} \quad (26)$$

На підставі (17), (21)-(26) можна записати онлайн версію алгоритма достовірної нечіткої кластеризації у вигляді

$$\begin{cases} \sigma_q^2(k+1) = \left(\sum_{l \neq q}^m \|x_{k+1} - w_l(k)\|^2 \right)^{-1}, \\ U_q(k+1) = \left(1 + \frac{\|x_{k+1} - w_q(k)\|^2}{\sigma_q^2(k+1)} \right)^{-1}, \\ U_{(k+1)}^* = U_q(k+1) \left(\sup U_l(k+1) \right)^{-1}, \\ Cr_q(k+1) = \frac{1}{2} \left(U_q^*(k+1) + 1 - \sup_{l \neq q} U_l^*(k+1) \right), \\ w_q(k+1) = w_q(k) + \eta(k+1) Cr_q^\beta(k+1) (x_{k+1} - w_q(k)). \end{cases} \quad (27)$$

Як видно, з обчислювальної точки зору online алгоритм достовірної нечіткої кластеризації не складніше рекурентних версій FCM і PCM, зберігаючи при цьому переваги достовірного підходу.

2. Результати обчислювального експерименту

Щоб перевірити ефективність розроблених методів, а також довести їх переваги перед аналогами, експериментальні дослідження були проведені за допомогою двох різних баз даних. Також проведено порівняльний аналіз якості кластеризації даних щодо основних характеристик оцінок якості, таких як: коефіцієнт розподілу (ПК), що визначає «перекриття» між групами точок; індекс розділу (SC), що кількісно визначає співвідношення суми компактності та розділеності кластерів; індекс C_i і Бені (XB), що вимірюють співвідношення загальної мінливості всередині кластерів та їх поділу.

Таблиця 1

Оцінка якості методів нечіткої кластеризації за допомогою першого набору даних

Методи кластеризації даних	Перший набір даних		
	PC	SC	XB
Fuzzy C-means (FCM)	0.50	1.62	0.19
Gustafson-Kessel	0.27	1.66	1.62
Gath-Geva	0.25	1.54	1.35
Адаптивна ймовірнісна нечітка кластеризація	0.25	1.44	0.01
Адаптивна нечітка можлива кластеризація даних	0.26	1.22	0.01
Адаптивна нечітка достовірнісна кластеризація даних	0.21	1.13	0.01

Таблиця 2

Оцінка якості методів нечіткої кластеризації за допомогою другого набору даних

Методи кластеризації даних	Другий набір даних		
	PC	SC	XB
Fuzzy C-means (FCM)	0.48	1.60	0.19
Gustafson-Kessel	0.26	1.64	1.62
Gath-Geva	0.26	1.50	1.35
Адаптивна ймовірнісна нечітка кластеризація	0.25	1.42	0.01
Адаптивна нечітка можлива кластеризація даних	0.37	1.11	0.18
Адаптивна нечітка достовірнісна кластеризація даних	0.23	1.22	0.01

Порівняно з відомими методами запропонований підхід до онлайн кластеризації даних демонструє дещо надійніші результати.

Висновок

Розглянуто задачу нечіткої кластеризації на основі ймовірнісного, можливісного і достовірного підходів на основі пакетного і online режимів надходження і обробки інформації. Введена рекурентна версія достовірного алгоритму, що є за суттю процедурою градієнтної оптимізації прийнятого критерія нечіткої достовірної кластеризації.

Список літератури

[1] Xu R., Wunsch D.C. II. Clustering— Hoboken, N.J.: John Wiley & Sons, Inc., 2009.
 [2] Aggarwal C.C. Data Mining: Text Book. Springer, 2015.

[3] Bezdek J.C. Pattern Recognition with Fuzzy Objective Function Algorithms. — N.Y.: — Plenum Press, 1981.
 [4] Höppner F., Klawonn F., Kruse R., Runkler T. Fuzzy Clustering Analysis: Methods for Classification, Data Analysis and Image Recognition.-Chichester: John Wiley & Sons, 1999. — 289 p.
 [5] R. Krishnapuram, J.M. Keller. A possibilistic approach to clustering. Fuzzy Systems, 1993, 1, № 2, P. 98–110.
 [6] Chintalapudi K. K. and M. kam, “A noise resistant fuzzy c-means algorithm for clustering,” IEEE conference on Fuzzy Systems Proceedings, vol. 2, May 1998, pp. 1458–1463.
 [7] Zhou J., Wang Q., Hung C.-C., Yi X. Credibilistic clustering: the model and algorithms. Int.J. of Uncertainty, Fuzziness and Knowledge-Based Systems- 2015-23-№4 — P. 545–564.
 [8] Zhou, J., Wang, Q., Hung, C. C. Credibilistic clustering algorithms via alternating cluster estimation.- J. Intell. Manuf.-2017-28 — P. 727–738.
 [9] Liu, B., & Liu, Y. Expected value of fuzzy variable and fuzzy expected value models. IEEE Transactions on Fuzzy Systems, - 2002-10-№4 — P. 445–450.
 [10] Liu, B. A survey of credibility theory. Fuzzy Optimization and Decision Making-2006-5-№4 — P. 387–408.
 [11] D.C. Park, I. Dagher. Gradient based fuzzy c-means (GB-FCM) algorithm. Proc. IEEE Int. Conf. on Neural Networks, 1984, P. 1626–1631.
 [12] F.L. Chung, T. Lee. Fuzzy competitive learning. Neural Networks, 1994, 7, №3, P. 539–552.
 [13] Bodyanskiy Ye, Kolodyazhnyi V., Stephan A. Recurcive fuzzy clustering algorithms. —Proc 10th East West Fuzzy Coll. 2002, -Zittau-Görlitz, HS, 2002 — P. 276–283.
 [14] Bodyanskiy, Ye. Computational intelligence techniques for data analysis / Ye. Bodyanskiy // Lecture Notes in Informatics.-Bonn: V.p. 72GI, 2005. — P. 15–36.
 [15] Zhou, J., & Hung, C.-C. (2007). A generalized approach to possibilistic clustering algorithms. Int. J. of Uncertainty, Fuzziness and Knowledge-Based Systems. — 2007 — 15. — P.117–138.
 [16] Young F.W., Hamer R.M. Theory and Applications of Multi-dimensional Scaling-Hillsdale, N.J.: Erlbaum, 1994.
 [17] Hu Zh., Bodyanskiy Ye, Tyshchenko O., Shafronenko A. Fuzzy clustering of incomplete data by means of similarity measures- Proc.2019 IEEE 2nd Ukr. Conf. on Electrical and Computer Engineering (UKRCON), Track 6.-Lviv, Ukraine, 2019. — P. 149–152.
 [18] Bezdek J.C. A convergence theorem for the fuzzy ISODATA clustering algorithms. — IEEE Trans. Pattern Anal. Mach. Intell. - 1980 - 2. — P. 1–8.
 Bodyanskiy Ye, Gorshkov Ye, Kokshenev I., Kolodyazhnyi V. Outlier resistant recursive fuzzy clustering algorithms. Ed. By B. Reusch «Computational Intelligence Theory and Applications» - Advances in Soft Computing-Vol.38.-Berlia Heidelberg, Springer Verlag, 2006 — P. 647–652.
 [19] Bodyanskiy Ye, Gorshkov Ye, Kokshenev I., Kolodyazhnyi V. Robust recursive fuzzy clustering algorithms- Proc. 12th East West Fuzzy Coll 2005 - Zittau-Görlitz, FH,2005 — P. 301–308.
 [20] Bodyanskiy Ye, Shafronenko A., Mashtalir S., Online robust fuzzy clustering of data with omissions using similarity measure of special type - Lecture Notes in Computational Intelligence and Decision Making-Cham: Springer, 2020 — P. 637–646

Надійшла до редколегії 4.09.2019