

УДК 519.16+519.85

О ПОЛИНОМИАЛЬНОЙ РАЗРЕШИМОСТИ ОДНОЙ КВАДРАТИЧНОЙ ЕВКЛИДОВОЙ ЗАДАЧИ КЛАСТЕРИЗАЦИИ В ОДНОМЕРНОМ СЛУЧАЕ

А. В. Кельманов^{1,2,*}, В. И. Хандеев^{1,2,**}

Представлено академиком РАН Ю.И. Журавлевым 28.02.2019 г.

Поступило 06.03.2019 г.

Рассматривается задача разбиения конечного множества точек евклидова пространства на кластеры. В этой задаче требуется минимизировать сумму по всем кластерам внутрикластерных сумм квадратов расстояний от элементов кластеров до их центров. Центры некоторых кластеров заданы на входе, а центры других кластеров определяются как центроиды (геометрические центры). Известно, что в общем случае задача NP-трудна в сильном смысле. Мы доказываем, что существует точный полиномиальный алгоритм для одномерного случая задачи.

Ключевые слова: кластеризация по критерию минимума суммы квадратов, евклидово пространство, NP-трудная задача, одномерный случай, полиномиальная разрешимость.

DOI: <https://doi.org/10.31857/S0869-56524872126-129>

Предметом исследования этой работы является одна из труднорешаемых задач разбиения конечного множества точек евклидова пространства. Цель исследования — выяснение вопроса о статусе вычислительной сложности одномерного случая задачи. Наше исследование мотивировано открытостью указанного математического вопроса, а также важностью задачи для приложений, в частности, для прикладных проблем анализа данных (Data analysis), интерпретации данных (Data mining), распознавания образов (Pattern recognition), машинного обучения (Machine learning) и обработки больших данных (Big data processing).

1. ФОРМУЛИРОВКА ЗАДАЧИ, ЕЁ ИСТОКИ И СВЯЗАННЫЕ С НЕЙ ПРОБЛЕМЫ

В известной кластеризационной задаче *K-Means* дано N -элементное множество \mathcal{U} точек в евклидовом пространстве размерности d и натуральное число K . Требуется найти разбиение входного множества \mathcal{U} на непустые непересекающиеся кластеры C_1, \dots, C_K , минимизирующее сумму

$$\sum_{k=1}^K \sum_{y \in C_k} \|y - \bar{y}(C_k)\|^2,$$

где $\bar{y}(C_k) = \frac{1}{|C_k|} \sum_{y \in C_k} y$ — центроид k -го кластера.

Другое распространённое название задачи *K-Means* — MSSS (Minimum Sum-of-Squares Clustering), т.е. кластеризация по критерию минимума суммы квадратов. В статистике эта задача известна с прошлого века и связана с именем Фишера (см., например, [1, 2]). На практике (в самых разнообразных приложениях) она возникает в ситуациях, когда имеется гипотеза о том, что имеющаяся совокупность \mathcal{U} выборочных (входных) данных содержит K однородных кластеров (подмножеств) C_1, \dots, C_K , в которых точки разбросаны относительно соответствующих неизвестных средних значений $\bar{y}(C_1), \dots, \bar{y}(C_K)$. Однако соответствие данных и кластеров неизвестно. Очевидно, что в этой ситуации для корректного применения классических статистических методов (проверки гипотез и оценивания) к обработке выборочных данных требуется сначала разбить их на однородные группы (кластеры). Эта ситуация типична, в частности, для отмеченных выше прикладных проблем.

Сильная NP-трудность задачи *K-Means* установлена относительно недавно [3]. Полиномиальная разрешимость этой задачи на числовой прямой доказана ещё в прошлом веке в [4]. В этой работе представлен алгоритм, реализующий схему динамиче-

¹Институт математики им. С.Л. Соболева
Сибирского отделения Российской Академии наук,
Новосибирск

²Новосибирский национальный исследовательский
государственный университет

*E-mail: kelm@math.nsc.ru

**E-mail: khandeev@math.nsc.ru

ского программирования. Он позволяет находить точное решение задачи в одномерном случае за время $O(KN^2)$. Этот алгоритм опирается на точный полиномиальный алгоритм решения известной задачи о ближайшем соседе (Nearest neighbor search problem) [5]. В последние годы для одномерной задачи K -Means построен ряд точных алгоритмов, имеющих улучшенные показатели быстродействия. Обзор ускоренных алгоритмов и их свойств можно найти в [6].

Объектом нашего внимания является близкая по постановке к задаче K -Means слабоизученная

Задача 1 (K -Means and Given J -Centers). Дано: N -элементное множество \mathcal{U} точек в евклидовом пространстве размерности d , натуральное число K и набор $\{c_1, \dots, c_J\}$ точек. Найти: разбиение множества \mathcal{U} на $K+J$ непустых кластеров $C_1, \dots, C_K, D_1, \dots, D_J$ такое, что

$$F = \sum_{k=1}^K \sum_{y \in C_k} \|y - \bar{y}(C_k)\|^2 + \sum_{j=1}^J \sum_{y \in D_j} \|y - c_j\|^2 \rightarrow \min,$$

где $\bar{y}(C_k)$ — центроид k -го кластера.

Эту задачу можно рассматривать как модификацию задачи K -Means. С другой стороны, введённые обозначения позволяют назвать задачу 1 как K -Means and Given J -Centers.

В отличие от задачи K -Means, задача 1 моделирует прикладную проблему, в которой у части кластеров, а именно у D_1, \dots, D_J , соответствующие центры c_1, \dots, c_J квадратичного разброса данных известны заранее (заданы). Эта прикладная проблема также типична для анализа и интерпретации данных, распознавания образов и машинного обучения. В частности, двухкластерный вариант задачи — 1 -Mean and Given 1 -Center — связан с решением прикладной проблемы обработки сигналов, а именно с проблемой совместного обнаружения квазипериодически повторяющегося импульса неизвестной формы и оценивания этой формы в условиях гауссовского шума с нулевым средним [7–9]. В этом двухкластерном варианте задачи нулевое среднее соответствует кластеру с заданным в начале координат центром. Первое сообщение об этом двухкластерном варианте задачи 1, по-видимому, было сделано в [7]. Заметим, что более простые экстремальные задачи, которые индуцируются прикладными проблемами помехоустойчивого обнаружения и различения импульсов заданных форм, характерны, в частности, для радиолокации, электронной разведки, гидроакустики, геофизики, технической и медицинской диагностики, космического мониторинга (см., например, [10–12]).

Сильная NP-трудность задачи 1 установлена в [13–15]. Вопрос о разрешимости этой задачи на числовой прямой до последнего времени оставался открытым.

Основной результат настоящей работы — доказательство полиномиальной разрешимости задачи 1 в одномерном случае.

2. ВСПОМОГАТЕЛЬНЫЕ УТВЕРЖДЕНИЯ

Наше доказательство опирается на несколько приведённых ниже вспомогательных утверждений, которые раскрывают структуру оптимального решения задачи 1. Для краткости мы приводим эти утверждения без доказательства, ограничиваясь изложением их идей.

Обозначим через $C_1^*, \dots, C_K^*, D_1^*, \dots, D_J^*$ оптимальные кластеры в задаче 1.

Лемма 1. Пусть $d = 1$ в задаче 1. Тогда если $c_m < c_\ell$, где $1 \leq m \leq J, 1 \leq \ell \leq J$, то для любых $x \in D_m^*$ и $z \in D_\ell^*$ выполнено неравенство $x \leq z$.

Лемма 2. Пусть $d = 1$ в задаче 1. Тогда:

- 1) если $\bar{y}(C_m^*) < c_\ell$, где $1 \leq m \leq K, 1 \leq \ell \leq J$, то для любых $x \in C_m^*$ и $z \in D_\ell^*$ справедливо неравенство $x \leq z$;
- 2) если $\bar{y}(C_m^*) > c_\ell$, где $1 \leq m \leq K, 1 \leq \ell \leq J$, то для любых $x \in C_m^*$ и $z \in D_\ell^*$ выполнено неравенство $x \geq z$.

Лемма 3. Пусть $d = 1$ в задаче 1. Тогда если $\bar{y}(C_m^*) < \bar{y}(C_\ell^*)$, где $1 \leq m \leq K, 1 \leq \ell \leq K$, то для любых $x \in C_m^*$ и $z \in C_\ell^*$ справедливо неравенство $x \leq z$.

Доказательство лемм 1–3 проводится методом от противного с применением равенства

$$\begin{aligned} (x - c_m)^2 + (z - c_\ell)^2 &= \\ &= 2(x - z)(c_\ell - c_m) + (z - c_m)^2 + (x - c_\ell)^2, \end{aligned}$$

справедливость которого следует из известной формулы для суммы квадратов диагоналей трапеции.

Лемма 4. Если $d = 1$ в задаче 1, то для каждого $k \in \{1, 2, \dots, K\}$ и каждого $j \in \{1, 2, \dots, J\}$ справедливо $\bar{y}(C_k^*) \neq c_j$.

Лемма 5. Если $d = 1$ в задаче 1, то для каждой $k, j \in \{1, 2, \dots, K\}$ таких, что $k \neq j$, справедливо $\bar{y}(C_k^*) \neq \bar{y}(C_j^*)$.

Доказательство этих лемм также проводится методом от противного.

Леммы 1–5 устанавливают взаимное расположение на числовой прямой оптимальных кластеров D_1^*, \dots, D_J^* с заданными центрами и кластеров

C_1^*, \dots, C_K^* с неизвестными центрами. Эти леммы служат базой для доказательства следующего утверждения.

Теорема 1. Пусть $d = 1$ в задаче 1 и, кроме того, точки y_1, \dots, y_N входного множества \mathcal{Y} , а также точки c_1, \dots, c_J упорядочены так, что

$$y_1 < \dots < y_N, \\ c_1 < \dots < c_J.$$

Тогда оптимальному разбиению \mathcal{Y} на кластеры $C_1^*, \dots, C_K^*, D_1^*, \dots, D_J^*$ соответствует разбиение последовательности $1, 2, \dots, N$ натуральных чисел на непересекающиеся целочисленные отрезки.

3. ПОЛИНОМИАЛЬНАЯ РАЗРЕШИМОСТЬ ЗАДАЧИ В ОДНОМЕРНОМ СЛУЧАЕ

Следующая теорема является основным результатом работы.

Теорема 2. Существует полиномиальный алгоритм, который при $d = 1$ находит оптимальное решение задачи 1 за время $\mathcal{O}(KJN^2)$.

Доказательство теоремы проводится конструктивно, а именно: мы даём обоснование алгоритма, который реализует схему динамического программирования и позволяет находить точное решение задачи за полиномиальное время.

Идея доказательства состоит в следующем. Без ограничения общности мы полагаем, что точки y_1, \dots, y_N входного множества \mathcal{Y} , а также точки c_1, \dots, c_J упорядочены, как в теореме 1.

Пусть $\mathcal{Y}_{s,t} = \{y_s, \dots, y_t\}$, где $1 \leq s \leq t \leq N$, — подмножество из $t - s + 1$ точек входного множества \mathcal{Y} с номерами от s до t .

Положим

$$f_{s,t}^j = \sum_{i=s}^t (y_i - c_j)^2, \quad j = 1, 2, \dots, J,$$

$$f_{s,t} = \sum_{i=s}^t (y_i - \bar{y}(\mathcal{Y}_{s,t}))^2,$$

где $\bar{y}(\mathcal{Y}_{s,t})$ — центроид подмножества $\mathcal{Y}_{s,t}$.

Мы доказываем, что оптимальное значение F^* целевой функции задачи 1 находится по формуле

$$F^* = F_{K,J}(N),$$

а значения функции

$$F_{k,j}(n), \quad k = -1, 0, 1, \dots, K, \quad j = -1, 0, 1, \dots, J, \\ n = 0, 1, \dots, N,$$

вычисляются по нижеприведённым рекуррентным формулам. При этом формула

$$F_{k,j}(n) = \begin{cases} 0, & n = k = j = 0; \\ +\infty, & n = 0; k = 0, 1, \dots, K; \\ & j = 0, 1, \dots, J; k + j \neq 0; \\ +\infty, & k = -1; j = -1, 0, \dots, J; \\ & n = 0, 1, \dots, N; \\ +\infty, & j = -1; k = -1, 0, \dots, K; \\ & n = 0, 1, \dots, N; \end{cases} \quad (1)$$

задаёт начальные и граничные условия для последующих вычислений. Основная формула

$$F_{k,j}(n) = \min \left\{ \min_{i=1}^n \{F_{k-1,j}(i-1) + f_{i,n}\}, \right. \\ \left. \min_{i=1}^n \{F_{k,j-1}(i-1) + f_{i,n}^j\} \right\},$$

$$k = 0, 1, \dots, K; j = 0, 1, \dots, J; n = 1, 2, \dots, N, \quad (2)$$

определяет рекурсию. В целом формулы (1), (2) реализуют прямой ход алгоритма.

Далее мы доказываем, что оптимальные кластеры $C_1^*, \dots, C_K^*, D_1^*, \dots, D_J^*$ находятся с помощью следующего рекуррентного правила, которое реализуется на обратном ходе алгоритма. Пошаговая запись правила:

Шаг 0. $k := K, j := J, n := N$.

Шаг 1. Если

$$\min_{i=1}^n (F_{k-1,j}(i-1) + f_{i,n}) \leq \min_{i=1}^n (F_{k,j-1}(i-1) + f_{i,n}^j),$$

то

$$C_k^* = \{y_{i^*}, y_{i^*+1}, \dots, y_n\},$$

где

$$i^* = \arg \min_{i=1}^n (F_{k-1,j}(i-1) + f_{i,n}),$$

выполнить $k := k - 1; n := i^* - 1$.

Если же

$$\min_{i=1}^n (F_{k-1,j}(i-1) + f_{i,n}) > \min_{i=1}^n (F_{k,j-1}(i-1) + f_{i,n}^j),$$

то

$$D_j^* = \{y_{i^*}, y_{i^*+1}, \dots, y_n\},$$

где

$$i^* = \arg \min_{i=1}^n (F_{k,j-1}(i-1) + f_{i,n}^j),$$

выполнить $j := j - 1; n := i^* - 1$.

Шаг 2. Если $k > 0$ или $j > 0$, то переход на шаг 1, иначе — конец вычислений.

Справедливость этого правила доказывается по индукции.

Наконец мы доказываем, что время работы алгоритма есть величина $\mathcal{O}(KJN^2)$, т.е. алгоритм полиномиален. Справедливость этой оценки времени работы алгоритма следует из того, что вычисления по формуле (2) производятся $\mathcal{O}(KJN)$ раз, причём для каждого вычисления требуется $\mathcal{O}(N)$ операций.

Таким образом, NP-трудная в сильном смысле задача 1 в одномерном случае разрешима за полиномиальное время. Построение приближённых эффективных алгоритмов с гарантированными оценками точности для общего случая представляется делом ближайшей перспективы.

Источники финансирования. Разделы 2, 3 выполнены при финансовой поддержке РФФИ, проекты 19–01–00308 и 18–31–00398, остальные разделы — при поддержке программы ФНИ РАН, проект 0314-2019-0015.

СПИСОК ЛИТЕРАТУРЫ

1. Fisher R.A. Statistical Methods and Scientific Inference. N.Y.: Hafner, 1956.
2. MacQueen J.B. Some Methods for Classification and Analysis of Multivariate Observations. In: Proc. 5th Berkeley Symp. on Mathematical Statistics and Probability. Berkeley: Univ. of California Press, 1967. V. 1. P. 281–297.
3. Aloise D., Deshpande A., Hansen P., Popat P. NP-hardness of Euclidean Sum-of-Squares Clustering // Mach. Learn. 2009. V. 75. № 2. P. 245–248.
4. Rao M. Cluster Analysis and Mathematical Programming // J. Amer. Stat. Assoc. 1971. V. 66. P. 622–626.
5. Bellman R. Dynamic Programming. Princeton: Princeton Univ. Press, 1957.
6. Grönlund A., Larsen K.G., Mathiasen A., Nielsen J.S., Schneider S., Song M. Fast Exact k -Means, k -Medians and Bregman Divergence Clustering in 1D. CoRR arXiv:1701.07204 (2017).
7. Кельманов А.В., Хамидуллин С.А., Кельманова М.А. Совместное обнаружение и оценивание повторяющегося фрагмента в зашумленной числовой последовательности при заданном числе повторов // Тез. докл. Рос. конф. “Дискретный анализ и исследование операций” (ДАИО-4). Новосибирск: Изд-во ИМ СО РАН, 2004. С. 185.
8. Gimadi E.Kh., Kel'manov A.V., Kel'manova M.A., Khamidullin S.A. A Posteriori Detection of a Quasi-Periodic Fragment in Numerical Sequences with Given Number of Recurrences // Sib. J. Industrial Math. 2006. V. 9. № 1 (25). P. 55–74.
9. Gimadi E.Kh., Kel'manov A.V., Kel'manova M.A., Khamidullin S.A. A Posteriori Detecting a Quasiperiodic Fragment in a Numerical Sequence // Pattern Recogn. and Image Anal. 2008. V. 18. № 1. P. 30–42.
10. Kel'manov A.V., Khamidullin S.A. Posterior Detection of a Given Number of Identical Subsequences in a Quasi-Periodic Sequence // Comput. Math. Math. Phys. 2001. V. 41. № 5. P. 762–774.
11. Kel'manov A.V., Jeon B. A Posteriori Joint Detection and Discrimination of Pulses in a Quasiperiodic Pulse Train // IEEE Trans. Signal Processing. 2004. V. 52. № 3. P. 645–656.
12. Carter J.A., Agol E., et al. Kepler-36: A Pair of Planets with Neighboring Orbits and Dissimilar Densities // Science. 2012. V. 337. № 6094. P. 556–559.
13. Кельманов А.В., Пяткин И.В. // ДАН. 2009. Т. 471. № 5. С. 590–592.
14. Kel'manov A.V., Pyatkin A.V. On a Version of the Problem of Choosing a Vector Subset // J. Appl. Ind. Math. 2009. V. 3. № 4. P. 447–455.
15. Kel'manov A.V., Pyatkin A.V. Complexity of Certain Problems of Searching for Subsets of Vectors and Cluster Analysis // Comput. Math. Math. Phys. 2009. V. 49. № 11. P. 1966–1971.

ON POLYNOMIAL SOLVABILITY OF ONE QUADRATIC EUCLIDEAN CLUSTERING PROBLEM ON A LINE

A. V. Kel'manov^{1,2}, V. I. Khandeev^{1,2}

¹Sobolev Institute of Mathematics, Novosibirsk, Russian Federation

²Novosibirsk State University, Novosibirsk, Russian Federation

Presented by Academician of the RAS Yu.I. Zhuravlev February 28, 2019

Received March 6, 2019

We consider the problem of partitioning a finite set of points in Euclidean space into clusters so as to minimize the sum over all clusters of the intracenter sums of the squared distances between clusters elements and their centers. The centers of some clusters are given as an input, while the other centers are defined as centroids (geometrical centers). It is known that the general case of the problem is strongly NP-hard. We show that there exists an exact polynomial algorithm for the one-dimensional case of the problem.

Keywords: minimum sum-of-squares clustering, Euclidean space, NP-hard problem, one-dimensional case, polynomial solvability.