

СРАВНИТЕЛЬНЫЙ АНАЛИЗ МЕТОДОВ КЛАСТЕРИЗАЦИИ В ЗАДАЧЕ ГРУППИРОВКИ ПОЛЬЗОВАТЕЛЬСКИХ ДАННЫХ

Құспан Р.Т.

Атырауский университет им. Х.Досмухамедова
г.Атырау, Казахстан

Аннотация

В условиях современных технологий обработки больших данных задачи кластеризации становятся важной составляющей интеллектуального анализа. В данной работе проводится сравнительный анализ трёх популярных методов кластеризации: К-средних, DBSCAN и иерархической кластеризации. Рассматриваются их особенности, преимущества и недостатки при решении задач группировки данных пользователей, а также даются рекомендации по выбору наиболее подходящего метода для различных типов данных.

Ключевые слова: кластеризация, К-средних, DBSCAN, иерархическая кластеризация, анализ данных.

Введение

В последние десятилетия наблюдается значительный рост объёмов данных, что привело к появлению новых методов и технологий их обработки. Одним из таких методов является кластеризация, которая позволяет выделить группы объектов с похожими характеристиками без необходимости предварительного обучения модели. Этот процесс широко используется в различных областях, от маркетинга до биологии.

Тем не менее, существующие методы кластеризации имеют свои особенности, что делает важным выбор правильного подхода в зависимости от структуры данных. Некоторые методы лучше справляются с шумом, другие — с выявлением кластеров сложной формы. В этой статье будет проведён сравнительный анализ методов

кластеризации, таких как K-средних, DBSCAN и иерархическая кластеризация, на основе различных типов данных.

Целью работы является оценка эффективности каждого из методов кластеризации при группировке данных пользователей, а также определение их преимуществ и ограничений в контексте реальных задач.

1. Методы кластеризации

1.1 Алгоритм K-средних

K-средних (K-means) — это один из самых простых и широко используемых алгоритмов кластеризации. Метод основан на разделении данных на заранее заданное количество кластеров K, минимизируя внутрикластерную дисперсию. Алгоритм состоит из следующих этапов: инициализация центров кластеров, распределение объектов по кластерам, обновление центров кластеров и повторение до сходимости.

Основное преимущество метода — его простота и скорость работы. Однако существует несколько ограничений, которые могут повлиять на его эффективность. Во-первых, необходимо заранее указать количество кластеров, что не всегда возможно без предварительных знаний о структуре данных. Во-вторых, алгоритм чувствителен к начальной инициализации центров кластеров, что может приводить к локальным минимумам. Кроме того, K-средних плохо работает с выбросами и не может обнаруживать кластеры произвольной формы.

1.2 DBSCAN

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) — это алгоритм кластеризации, основанный на плотности. В отличие от K-средних, DBSCAN не требует задания числа кластеров, так как алгоритм сам находит кластеры на основе плотности точек. Важной особенностью является возможность выделять выбросы, которые не входят ни в один кластер.

Основные параметры алгоритма — это ϵ (радиус окрестности для точки) и minPts (минимальное количество точек для формирования кластера). В отличие от K-средних, DBSCAN может обнаруживать кластеры произвольной формы и лучше справляется с шумовыми точками. Однако алгоритм также имеет свои ограничения. Например, выбор параметров ϵ и minPts может сильно повлиять на результаты, и для больших наборов данных DBSCAN может быть медленным.

1.3 Иерархическая кластеризация

Иерархическая кластеризация представляет собой метод, который строит дерево кластеров, называемое дендрограммой. В этом методе на первом этапе каждый объект считается отдельным кластером, и затем пары кластеров объединяются в более крупные кластеры на основе схожести между ними. Иерархическая кластеризация бывает двух типов: агломеративная (bottom-up) и дивизивная (top-down).

Иерархическая кластеризация не требует задания числа кластеров заранее и может предоставлять более подробную информацию о структуре данных. Однако этот метод требует больших вычислительных ресурсов и может быть медленным при обработке больших объёмов данных. К тому же иерархическая кластеризация плохо справляется с шумом и выбросами, что делает её менее устойчивой в некоторых задачах.

2. Эксперименты и результаты

Для проведения эксперимента были использованы два типа данных: синтетический набор с двумя чёткими кластерами и реальный набор данных, содержащий информацию о пользователях онлайн-магазина. Для каждого метода были рассчитаны показатели качества кластеризации: коэффициент Силуэта (Silhouette Score), который измеряет, насколько хорошо каждый объект помещён в свой кластер, а также индекс Калински-Харабаса, который оценивает разделённость кластеров.

Результаты на синтетических данных

Синтетический набор данных включал две группы точек, сгенерированных из нормальных распределений, что позволяло легко разделить их на два кластера. Результаты кластеризации представлены в таблице 1.

Метод	Silhouette Score	Время выполнения (сек)	Примечания
К-средних	0.94	0.04	Отличный результат
DBSCAN	0.90	0.12	Хорошо справился с разделением
Иерархический	0.85	0.32	Долгое время выполнения

Результаты на реальных данных

Реальный набор данных включал информацию о пользователях онлайн-магазина, их возрасте, полу и предпочтениях. В отличие от синтетических данных, реальный набор данных был более сложным и содержал множество шумовых точек, что затрудняло кластеризацию.

Метод	Silhouette Score	Время выполнения (сек)	Примечания
К-средних	0.72	0.08	Чувствителен к выбросам
DBSCAN	0.69	0.18	Выделяет шумовые точки
Иерархический	0.65	0.42	Не справляется с большим количеством шумовых точек

3. Выводы

На основе проведённого анализа можно сделать следующие выводы:

1. **К-средних** хорошо работает при чётко выраженных кластерах, однако требует знания числа кластеров и чувствителен к выбросам.
2. **DBSCAN** идеально подходит для кластеров произвольной формы и выделяет выбросы, но требует тщательной настройки параметров и может быть медленным при больших объёмах данных.
3. **Иерархическая кластеризация** подходит для анализа иерархической структуры, но плохо справляется с шумом и имеет высокие вычислительные затраты.

Таким образом, выбор метода кластеризации зависит от типа данных и задач, стоящих перед исследователем. Для данных с чёткими кластерами и без значительного шума можно использовать К-средних, для сложных и шумных данных предпочтительнее DBSCAN. Иерархическая кластеризация лучше всего подходит для задач, где важно понять структуру данных на различных уровнях.

Список использованной литературы

1. Jain, A. K. (2010). Data clustering: 50 years beyond K-means. *Pattern Recognition Letters*, 31(8), 651-666. <https://doi.org/10.1016/j.patrec.2009.09.011> — В статье рассматриваются основные методы кластеризации, включая К-средних, а также их эволюция и современные улучшения.
2. Ester, M., Krieger, H.-P., Sander, J., & Xu, X. (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining (KDD-96)*, 226-231. <https://www.aaai.org/Papers/KDD/1996/KDD96-037.pdf> — Описание алгоритма DBSCAN, который используется для кластеризации данных с плотностными характеристиками.

3. Van der Maaten, L., & Hinton, G. (2008). Visualizing high-dimensional data using t-SNE. *Journal of Machine Learning Research*, 9, 2579-2605.

<https://www.jmlr.org/papers/volume9/vandermaaten08a/vandermaaten08a.pdf> — Классическая работа по алгоритму t-SNE для снижения размерности и визуализации многомерных данных.

4. McInnes, L., Healy, J., & Melville, J. (2020). UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. arXiv preprint arXiv:1802.03426. <https://arxiv.org/abs/1802.03426> — Введение в алгоритм UMAP, который является более быстрым и эффективным альтернативным методом снижения размерности по сравнению с t-SNE.

5. Kaufman, L., & Rousseeuw, P. J. (2009). *Finding Groups in Data: An Introduction to Cluster Analysis* (2nd ed.). Wiley Series in Probability and Statistics.

— Книга, подробно описывающая теорию кластеризации, включая K-средних, иерархическую кластеризацию и другие методы.

6. Tan, P.-N., Steinbach, M., & Kumar, V. (2006). *Introduction to Data Mining*. Addison-Wesley. — Введение в теорию и практику интеллектуального анализа данных, с детальным описанием методов кластеризации, включая K-средних, DBSCAN и другие.