

Обучение без учителя. Кластеризация

Санкт-Петербургский государственный университет
Кафедра статистического моделирования
Семинар «Статистическое и машинное обучение»

Санкт-Петербург, 2025

Обучение без учителя (unsupervised learning) — раздел машинного обучения, изучающий класс задач обработки данных, в которых известны только описания множества объектов (признаки объектов) из обучающей выборки, и требуется обнаружить внутренние зависимости, существующие между объектами. В отличие от обучения с учителем, правильные «ответы» или «метки» для объектов неизвестны.

Задачи обучения без учителя

- **Кластеризация.** Разделение объектов на группы (кластеры) на основании их сходства.
- **Поиск ассоциативных правил.** Выявление связей между объектами. Цель — найти закономерности вида «Если встречается A , то с высокой вероятностью встречается и B », где A, B — некоторые не пересекающиеся наборы признаков.
- **Понижение размерности.** Уменьшение количества признаков при сохранении максимально возможной информации из исходных данных.
- **Заполнение пропусков.** Восстановление отсутствующих данных на основе закономерностей, найденных в имеющихся данных.

Постановка задачи кластеризации

Дано:

X — пространство объектов;

$X^n = \{x_1, \dots, x_n\}$ — выборка из X , где x_i — i -й объект;

$\rho : X \times X \rightarrow [0, \infty)$ — функция расстояния между объектами.

Найти:

Y — множество кластеров;

$a : X \rightarrow Y$ — алгоритм кластеризации.

Причем Y и a такие, что

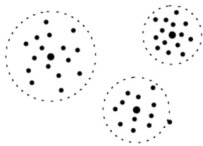
- каждый кластер состоит из близких объектов (относительно ρ);
- объекты разных кластеров существенно различаются.

Некорректность задачи кластеризации

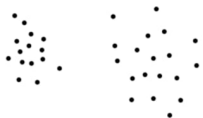
Решение задачи кластеризации неоднозначно:

- точной постановки задачи кластеризации нет;
- существует много критериев качества кластеризации;
- число кластеров, как правило, неизвестно заранее;
- результат кластеризации сильно зависит от метрики ρ , выбор которой также не однозначен.

Типы кластерных структур



кластеры с центрами



внутрикластерные расстояния
меньше межкластерных



ленточные кластеры

Типы кластерных структур



перемычки между кластерами



разреженный фон
из нетипичных объектов



перекрывающиеся кластеры

Типы кластерных структур



кластеры могут вообще отсутствовать



а это вообще не кластеры

- Каждый метод кластеризации имеет свои ограничения и выделяет кластеры лишь некоторых типов.
- Понятие «тип кластерной структуры» зависит от метода и также не имеет формального определения.

В большинстве источников выделяют пять групп алгоритмов:

- **Основанные на центроидах (centroid based):** k-means, k-modes, k-medoids, **Meanshift**, Affinity Propagation
- **Иерархические (hierarchical):** агломеративные (Ward/single/average/complete linkage), BIRCH, на основе теории графов (выделение связных компонент и минимальное остовное дерево), Spectral Clustering
- **Основанные на плотности (density based):** **DBSCAN**, **HDBSCAN**, **OPTICS**
- **Сеточные (grid based):** STING, Wave cluster, **CLIQUE**, OptiGrid, MAFLA
- **Основанные на модели данных (model based):** **Expectation Maximization (EM)**, COBWEB

Model-based clustering — подход с четко поставленной задачей: предполагаем какую-то статистическую модель данных и в ней находим параметры.

Выборка \mathbf{X}^n — случайна, независима и взята из смеси распределений, плотность которой в точке $\mathbf{x} \in \mathbf{X}^n$ представима в виде:

$$p(\mathbf{x}) = \sum_{j=1}^k w_j p_j(\mathbf{x}; \boldsymbol{\theta}_j), \text{ при этом}$$

- $\sum_{j=1}^k w_j = 1, w_j \geq 0$ — априорные вероятности кластеров.
- $p_j(\mathbf{x}; \boldsymbol{\theta}_j)$ — плотность распределения j -го кластера с параметрами $\boldsymbol{\theta}_j$.

ЕМ-алгоритм для model-based подхода

Предполагается, зная число кластеров k и вид плотностей p_j , оценить параметры w_j, θ_j , максимизируя логарифм функции правдоподобия

$$\ln L(\{\mathbf{x}_i\}; \{w_j\}; \{\theta_j\}) = \sum_{i=1}^n \ln \sum_{j=1}^k w_j p_j(\mathbf{x}_i; \theta_j) \rightarrow \max_{\{w_j\}, \{\theta_j\}}.$$

Для решения данной задачи применяется **ЕМ-алгоритм**.

ЕМ-алгоритм. Шаг E (Expectation)

Для текущих оценок параметров вычисляем вероятность принадлежности каждой точки \mathbf{x}_i к каждой компоненте смеси с параметрами θ_j по формуле Байеса:

$$g_{ij} = \frac{w_j p_j(\mathbf{x}_i; \theta_j)}{\sum_{s=1}^k w_s p_s(\mathbf{x}_i; \theta_s)} - \text{скрытые переменные.}$$

ЕМ-алгоритм. Шаг M (Maximization)

Обновляем параметры, используя скрытые переменные, найденные на предыдущем шаге. Максимизируем логарифм функции правдоподобия методом Лагранжа:

$$\mathcal{L}(\{\mathbf{x}_i\}; \{w_j\}; \{\boldsymbol{\theta}_j\}) = \sum_{i=1}^n \ln \sum_{j=1}^k w_j p_j(\mathbf{x}_i; \boldsymbol{\theta}_j) - \lambda \left(\sum_{j=1}^k w_j - 1 \right).$$

Из равенства нулю производной по w_j следует:

$$w_j = \frac{1}{n} \sum_{i=1}^n g_{ij}, \quad j = 1, \dots, k.$$

Из равенства нулю производной по $\boldsymbol{\theta}_j$ следует:

$$\boldsymbol{\theta}_j = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n g_{ij} \ln p(\mathbf{x}_i; \boldsymbol{\theta}), \quad j = 1, \dots, k.$$

Таким образом, параметры будут уточняться на каждом шаге.

Предположим, что компоненты смеси имеют нормальные распределения со средними $\boldsymbol{\mu}_j$ и матрицами ковариаций Σ_j , тогда имеем следующие оценки параметров:

$$\boldsymbol{\mu}_j = \frac{1}{nw_j} \sum_{i=1}^n g_{ij} \mathbf{x}_i,$$

$$\Sigma_j = \frac{1}{nw_j} \sum_{i=1}^n g_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j)(\mathbf{x}_i - \boldsymbol{\mu}_j)^T.$$

ЕМ-алгоритм для model-based. Плюсы и минусы

Преимущества:

- Имеет формально поставленную задачу;
- Полученное разбиение на кластеры интерпретируемо со статистической точки зрения.

Недостатки:

- Число кластеров k является гиперпараметром;
- Параметры должны быть оценены, что требует большего количества точек данных в каждой компоненте;
- Алгоритм чувствителен к начальным данным.

k-means и его модификации

- k-means — итеративно группирует данные вокруг k-центров кластеров, пересчитывая их положение как среднее точек кластера до сходимости.
- Mini-batch k-means — модификация классического k-means, использующая случайные подвыборки данных на каждой итерации для обучения. Хорошо подходит для больших датасетов.
- k-medoids — вариант k-means, который в качестве центроидов выбирает реальные точки (медоиды) из данных, а не их средние значения, что повышает устойчивость к выбросам.
- k-modes — вариант алгоритма k-means для работы с категориальными данными, который выбирает один из объектов в кластере в качестве моды и минимизирует сумму расстояний Хэмминга между модой и объектами в кластере. Расстояние Хэмминга представляет из себя количество позиций, в которых значения векторов не совпадают.

Mean Shift — это алгоритм кластеризации, который реализует локальный градиентный подъём по оценке плотности данных (KDE). Сдвиг — это направление наибольшего возрастания плотности, а итоговые «центры» кластеров — локальные максимумы плотности (моды). Точки, которые сходятся к одному максимуму, считаются принадлежащими одному кластеру.

В отличие от популярного алгоритма k-means, Mean Shift не требует предварительного указания количества кластеров, но требует задать параметр окна для KDE.

Mean Shift. Определения

- Ядерная оценка плотности в \mathbb{R}^p (p — кол-во признаков) с окном ширины h (bandwidth) и радиально-симметричным ядром:

$$f_h(\mathbf{x}) = \frac{1}{nh^p} \sum_{i=1}^n k \left(\left\| \frac{\mathbf{x} - \mathbf{x}_i}{h} \right\|^2 \right), \text{ где } k : [0, \infty] \rightarrow \mathbb{R} \text{ — профиль ядра.}$$

- Средневзвешенное значение плотности в окне:

$$m(\mathbf{x}) = \frac{\sum_{i=1}^n \mathbf{x}_i g \left(\left\| \frac{\mathbf{x} - \mathbf{x}_i}{h} \right\|^2 \right)}{\sum_{i=1}^n g \left(\left\| \frac{\mathbf{x} - \mathbf{x}_i}{h} \right\|^2 \right)}, \text{ где } g(u) = -k'(u).$$

$m(\mathbf{x})$ пропорционален градиенту функции $f_h(\mathbf{x})$.

- Вектор сдвига (mean shift) равен $m(\mathbf{x}) - \mathbf{x}$ и направлен в сторону максимального увеличения плотности.

Mean Shift. Алгоритм

- ❶ Инициализация — каждая точка данных становится потенциальным центром кластера.
- ❷ Для каждой точки создается скользящее окно с фиксированным радиусом (bandwidth).
- ❸ Вычисляется центроид $m(\mathbf{x})$ всех точек в пределах окна как взвешенное среднее.
- ❹ Центр окна перемещается к вычисленному центроиду.
- ❺ Процесс повторяется до тех пор, пока окна не перестанут существенно смещаться (достигнута сходимость).
- ❻ Пересекающиеся окна объединяются, выбирается окно с наибольшим количеством точек.
- ❼ Точки данных назначаются кластерам в соответствии с окном, в котором они находятся.

Mean Shift. Плюсы и минусы

Преимущества:

- Не требует предварительного задания количества кластеров k ;
- Может находить кластеры произвольной формы;
- Теоретически обоснован как метод поиска мод плотности.

Недостатки:

- Чувствителен к параметру bandwidth h : малое h приводит к большому количеству мелких кластеров, может создать отдельный кластер для каждой точки; большое h приводит к малому количеству крупных кластеров, может объединить все точки в один кластер;
- Высокая вычислительная сложность для больших наборов данных;
- Нормализация/стандартизация данных может существенно изменить расклад.

DBSCAN (Density-based spatial clustering of applications with noise) — это эвристический алгоритм кластеризации, основанный на плотности. Алгоритм группирует вместе те объекты, которые тесно расположены, помечая как выбросы объекты, которые находятся в областях с малой плотностью.

Помимо того, что DBSCAN может обнаруживать кластеры произвольной формы и выбросы в данных, его главная особенность заключается в самостоятельном определении необходимого количества кластеров, что избавляет от необходимости в их подборе.

Алгоритм достаточно прост, наряду с k-means один из самых популярных

DBSCAN. Типы объектов

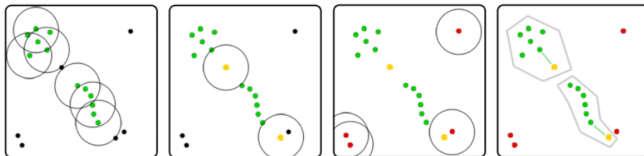
Для $x \in X$ его ε -окрестность $U_\varepsilon(x) = \{u \in X : \rho(x, u) \leq \varepsilon\}$.

Алгоритм имеет 2 гиперпараметра:

- ✱ величина окрестности ε ;
- ✱ минимальное количество объектов в окрестности m .

Каждый объект может быть одного из трёх типов:

- **корневой** (core): в ε -окрестности не менее m точек;
- **граничный** (border): в ε -окрестности меньше m точек, но среди них есть как минимум одна корневая;
- **шумовой** (noise): не корневой и не граничный.



DBSCAN. Алгоритм

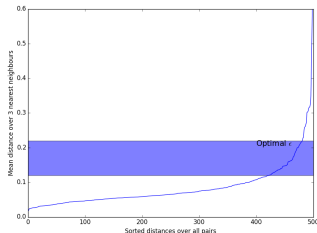
Вход: выборка $X^n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, параметры ε и m ;

Выход: разбиение выборки на кластеры и шумовые выбросы;

- 1 $U = X^n$ — мн-во еще не обработанных объектов, $a = 0$;
- 2 **Пока** есть некластеризованные точки, т.е. $U \neq \emptyset$;
- 3 взять случайную точку $\mathbf{x} \in U$;
- 4 если $|U_\varepsilon(\mathbf{x})| < m$, **то**
- 5 позначить \mathbf{x} как шумовой;
- 6 **иначе**
- 7 создать новый кластер: $K = U_\varepsilon(\mathbf{x})$; $a = a + 1$;
- 8 **для всех** $\mathbf{x}' \in K$
- 9 **если** $|U_\varepsilon(\mathbf{x}')| \geq m$ **то** $K = K \cup U_\varepsilon(\mathbf{x}')$;
- 10 **иначе** позначить \mathbf{x}' как граничный элемент K ;
- 11 соотнести объект классу a **для всех** $\mathbf{x}' \in K$;
- 12 $U = U \setminus K$

DBSCAN. Подбор параметров

- Значение параметра m предлагается выбирать как $m \geq p + 1$, где p — количество признаков (размерность данных). Также встречается $m = 2 \cdot p$.
- Чтобы подобрать ε , используют следующий алгоритм:
 - ✳ Строится график, где по оси y для каждой точки будет среднее расстояние по m ближайшим соседям, а по оси x — точки, отсортированные в порядке возрастания этого расстояния.
 - ✳ Следует взять ε где-нибудь в полосе, где происходит самый сильный перегиб.



Преимущества:

- Относительно быстро работает на больших данных;
- Устойчив к выбросам;
- Не требуется заранее указывать количество кластеров;
- Способен находить кластеры произвольной формы, а также шумовые точки в данных;
- Хорошо поддаётся модифицированию.

Недостатки:

- Чувствителен к выбору параметров ϵ и m ;
- Проблемы с высокоразмерными данными (проклятие размерности);
- Плохо работает с кластерами разной плотности. Не способен соединять кластеры через проёмы, и, наоборот, способен связывать явно различные кластеры через плотно населённые перемычки.

OPTICS (Ordering points to identify the clustering structure) — модификация DBSCAN для решения его ключевой проблемы: неспособности эффективно обрабатывать данные с кластерами различной плотности.

Основная идея OPTICS заключается в линейном упорядочивании точек таким образом, чтобы пространственно близкие точки становились соседними в этом упорядочивании. Этот порядок затем можно использовать для извлечения кластеров с любыми параметрами плотности.

- **Основное расстояние** (Core Distance) d_{core} — расстояние от точки до ее m -го ближайшего соседа. То есть это такое минимальное расстояние $\varepsilon' \leq \varepsilon$, при котором точка все еще остается корневой (основной).

Если такого расстояния $\leq \varepsilon$ не существует, d_{core} считается неопределённой.

- **Достижимое расстояние** (Reachability Distance) — мера того, насколько легко точка может быть достигнута из других точек в наборе данных:

$$d_{reach}(\mathbf{a}, \mathbf{b}) = \max(d_{core}(\mathbf{b}), \rho(\mathbf{a}, \mathbf{b})), \text{ где } \rho(\mathbf{a}, \mathbf{b}) = \text{dist}(\mathbf{a}, \mathbf{b}).$$

- Как основное, так и достижимое расстояния не определены, если нет достаточно плотного кластера (применительно к ε).

Алгоритм строит упорядочивание точек $\pi = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ такое, что последовательность достижимых расстояний $r_i = d_{reach}(\mathbf{x}_{i-1}, \mathbf{x}_i)$ отражает переходы между плотными областями данных.

- ❶ $d_{reach}(\mathbf{x}) = \infty, \quad \forall \mathbf{x} \in \mathbf{X}^n.$
- ❷ Для каждой непосещённой точки $\mathbf{x} \in \mathbf{X}^n$:
 - Вычислить $d_{core}(\mathbf{x})$.
 - Добавить \mathbf{x} в упорядочение π .
 - Если $d_{core}(\mathbf{x})$ определена, то для всех $\mathbf{y} \in U_\varepsilon(\mathbf{x})$, которые ещё не включены в упорядочение, пересчитывается их достижимость:

$$d_{reach}(\mathbf{y}) = \min(d_{reach}(\mathbf{y}), \max(d_{core}(\mathbf{x}), \rho(\mathbf{x}, \mathbf{y}))).$$

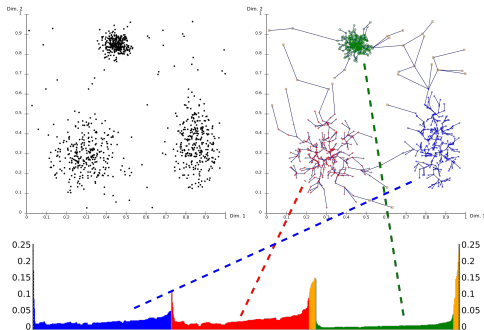
- Следующей выбирается точка

$$p' = \arg \min_{\mathbf{y} \notin \pi} d_{reach}(\mathbf{y})$$

и процесс повторяется до тех пор, пока все точки не будут упорядочены.

OPTICS. Извлечение кластеров

- Результатом работы алгоритма является **график достижимости** — двумерный график, где по оси x откладываются точки в порядке их обработки алгоритмом, а по оси y — достижимое расстояние.
- Поскольку точки, принадлежащие кластеру, имеют небольшое достижимое расстояние до ближайшего соседа, **кластеры** выглядят как **долины** на графике достижимости. Чем глубже долина, тем плотнее кластер.
- **Пики** на графике представляют расстояния между кластерами или переходы от кластера к шуму.



HDBSCAN (Hierarchical DBSCAN) — модификации DBSCAN, которая автоматически находит подходящее значение ϵ для каждого кластера, используя иерархический подход, что позволяет определять кластеры с разной плотностью.

По сравнению с DBSCAN для него требуется большее количество вычислений, что увеличивает время работы алгоритма.

- **Основное расстояние** (Core Distance) d_{core} — расстояние от точки до ее m -го ближайшего соседа. Это величина, показывающая, насколько плотна область вокруг точки.
- **Взаимное достижимое расстояние** (Mutual Reachability Distance) — специальная метрика:

$$d_{mreach}(\mathbf{a}, \mathbf{b}) = \max(d_{core}(\mathbf{a}), d_{core}(\mathbf{b}), \rho(\mathbf{a}, \mathbf{b})), \text{ где } \rho(\mathbf{a}, \mathbf{b}) = \text{dist}(\mathbf{a}, \mathbf{b}).$$

Расстояния между точками в областях высокой плотности (с малым d_{core}) остаются неизменными, а расстояния между точками в областях низкой плотности (с большим d_{core}) увеличиваются. Таким образом, общий эффект использования MRD в качестве метрики расстояния заключается в том, что точки в областях с низкой плотностью отдаляются от точек в областях с высокой плотностью.

Это делает алгоритм устойчивым к разной плотности.

- ❶ Для каждой точки рассчитывается d_{core} .
- ❷ Вычисляем MRD между всеми парами точек, строим полный взвешенный граф.
- ❸ На графе взаимной достижимости строим минимальное остовное дерево (MST).
- ❹ Строим иерархию кластеров (преобразовываем MST в дендрограмму):
 - Сортируем рёбра MST по весу в порядке возрастания;
 - Начиная с самого маленького расстояния, рёбра последовательно добавляются, соединяя точки и формируя кластеры.

HDBSCAN. Алгоритм

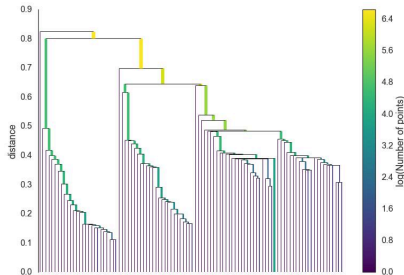
- 5 Введём параметр \hat{m} — минимальное число точек в кластере и новую шкалу для дендрограммы $\lambda = 1/d_{mreach}$.
- 6 Будем рассматривать дендрограмму снизу вверх по мере убывания λ (возрастания d_{mreach}). Мы скажем, что множество узлов C_i , получающееся рассмотрением связного поддерева в дендрограмме на высоте λ , является кластером, если оно содержит хотя бы \hat{m} вершин. Максимальное такое λ назовем $\lambda_{i,death}$.
- 7 Минимальную величину λ , на которой поддерево C_t распадается на два кластера C_i и C_j , назовём $\lambda_{i,birth} = \lambda_{j,birth} = \lambda_{t,death}$.
- 8 Разбиваем точки на кластеры так, чтобы

$$\max \sum_{i \in \mathcal{C}} \int_{\lambda_{i,birth}}^{\lambda_{i,death}} p_{i,\lambda} d\lambda,$$

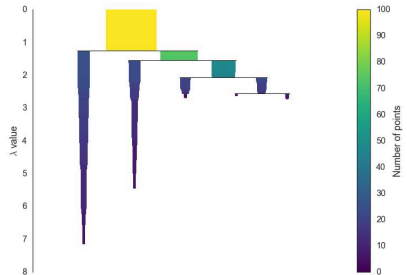
где $p_{i,\lambda}$ — это число точек в кластере C_i на уровне λ .

- 9 Не попавшие в разбиение точки объявляем выбросами.

HDBSCAN. Иллюстрация



(a) Дендрограмма



(b) Сжатая дендрограмма

CLIQUE (CLustering In QUest) — это алгоритм, который комбинирует в себе идеи сеточного подхода к кластеризации и подхода, основанного на подпространствах, для обнаружения плотных кластеров в подпространствах максимальной размерности.

Мотивация: В данных с большим количеством признаков «проклятие размерности» приводит к тому, что данные становятся очень разреженными. Кластеры часто существуют только в пределах небольшого подмножества признаков, а в других измерениях точки могут быть распределены случайным образом. CLIQUE создан для решения именно этой проблемы.

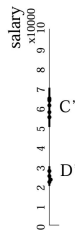
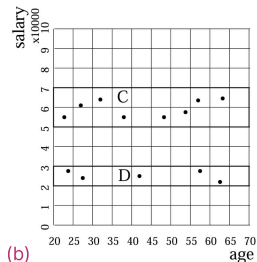
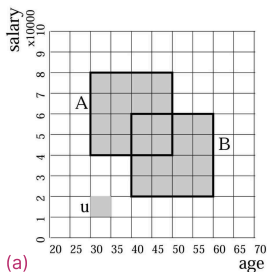
У алгоритма два входных параметра:

- * ξ — параметр сетки;
- * τ — пороговое значение плотности.

- Пространство признаков разбивается на ξ равных частей по каждому признаку. Пересечение одного интервала из каждого измерения называется **единицей** (unit).
- **Селективность** (selectivity) единицы определяется как доля общего количества точек данных, содержащихся в этой единице.
- **Плотная единица** (dense unit) — единица считается плотной, если ее селективность превышает пороговое значение плотности τ .
- **Кластер** — максимальное множество связанных (т.е. имеющих общую грань) плотных единиц в одном и том же подпространстве.

- CLIQUE генерирует описания кластеров в виде выражений в **дизъюнктивной нормальной форме (ДНФ)**, покрывая кластер минимальным количеством максимальных, возможно перекрывающихся, прямоугольников и описывая кластер как объединение этих прямоугольников.
- Кроме того, последний шаг CLIQUE принимает в качестве входных данных покрытие для каждого кластера и находит минимальное покрытие, определяемое с точки зрения количества максимальных областей (прямоугольников), необходимых для покрытия кластера.

CLIQUE. Иллюстрация определений



- ★ График (a): сетка 10×10 , единица — u , кластер $A \cup B$, его мин. описание в виде ДНФ

$$((30 \leq \text{age} < 50) \wedge (4 \leq \text{salary} < 8)) \vee ((40 \leq \text{age} < 60) \wedge (2 \leq \text{salary} < 6)).$$

- ★ График (b): $\tau = 20\%$, ни одна двумерная единица не является плотной, и в исходном пространстве данных нет кластеров. Однако если спроецировать точки на измерение зарплаты, то получим три одномерные плотные единицы. Две из них соединены, поэтому в одномерном подпространстве зарплаты есть два кластера.

CLIQUE. 3 основных этапа

1 Поиск плотных подпространств

- Представляем данные как ξ -мерную сетку;
- Находим **плотные ячейки** (где точек $>$ порога τ);
- Поиск плотных единиц в подпр-ах высокой размерности ведется снизу вверх, начиная с 1-мерных подпр-в. Используем принцип **монотонности**: «Если k -мерная ячейка плотная, то все её $(k-1)$ -мерные проекции тоже плотные».

2 Формирование кластеров

- После того как найдены все плотные k -мерные единицы в определенном подпространстве, алгоритм объединяет их в кластеры;
- Строится граф, где вершины — это плотные единицы, а ребра соединяют соседние единицы;
- Все связанные плотные единицы объединяются в один кластер. Для этого используется алгоритм поиска в глубину (DFS).

3 Генерация описаний

- Покрываем кластер **минимальным количеством максимальных прямоугольников**;
- Строим **минимальное покрытие**: удаляются те прямоугольники, которые полностью покрываются другими, более крупными прямоугольниками из покрытия;
- Записываем результат в виде **ДНФ-формул**.

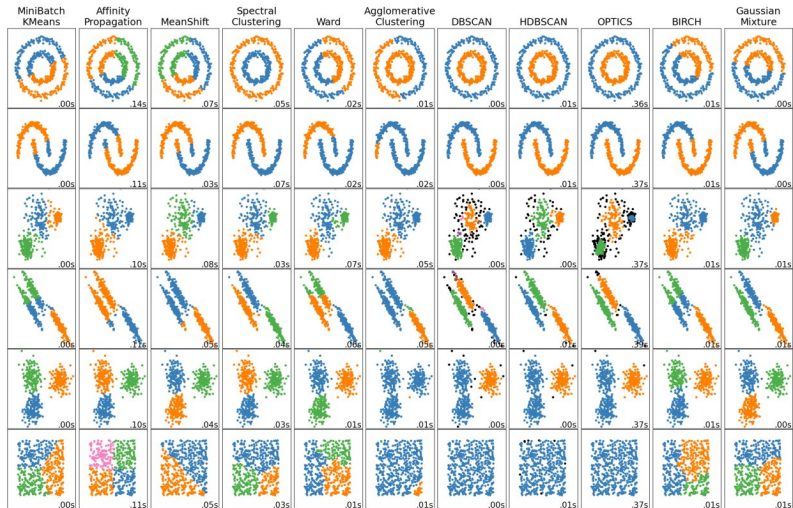
Преимущества:

- Хорошо работает для данных с большой размерностью;
- Не требует знания числа кластеров заранее, они формируются автоматически из плотных ячеек;
- Каждый кластер можно описать в логической форме;
- Алгоритм устойчив к пропущенным значениям во входных данных.

Недостатки:

- Результат сильно зависит от входных параметров: размера сетки и пороговой плотности;
- Размер сетки растёт экспоненциально с размерностью.

Сравнение методов



Функционалы качества кластеризации

Задачу кластеризации можно ставить как задачу *дискретной оптимизации*: необходимо так приписать номера кластеров y_i объектам \mathbf{x}_i , чтобы значение выбранного функционала качества приняло наилучшее значение.

- **Среднее внутрикластерное расстояние:**

$$F_0 = \frac{\sum_{i < j} \mathbf{I}_{\{y_i = y_j\}} \rho(\mathbf{x}_i, \mathbf{x}_j)}{\sum_{i < j} \mathbf{I}_{\{y_i = y_j\}}} \rightarrow \min.$$

- **Среднее межкластерное расстояние:**

$$F_1 = \frac{\sum_{i < j} \mathbf{I}_{\{y_i \neq y_j\}} \rho(\mathbf{x}_i, \mathbf{x}_j)}{\sum_{i < j} \mathbf{I}_{\{y_i \neq y_j\}}} \rightarrow \max.$$

На практике вычисляют отношение пары функционалов, чтобы учесть как межкластерные, так и внутрикластерные расстояния:

$$F_0 / F_1 \rightarrow \min$$

Функционалы качества кластеризации. Silhouette

- Пусть a и b есть среднее расстояние между наблюдением и всеми другими точками в том же кластере/в следующем ближайшем кластере
- Коэффициент силуэта для наблюдения есть

$$s = \frac{b - a}{\max(a, b)}$$

- Для выборки коэффициент силуэта задается средним значением коэффициентов каждого наблюдения
- Значения в интервале $[-1; 1]$. Чем больше, тем лучше. Если коэффициент близок к 0, то это свидетельство в сторону того, что кластеры "накладываются" друг на друга

Функционалы качества кластеризации. Calinski-Harabasz Index (Variance Ratio Criterion)

- Индекс Калинского–Харабэса определяется как отношение между средней межкластерной дисперсией и средней дисперсией внутри кластеров

$$s = \frac{\text{tr}(B)}{k-1} / \frac{\text{tr}(W)}{n-k}$$

- Тут k – число кластеров, а матрицы B и W имеют вид

$$W = \sum_{q=1}^k \sum_{x \in C_q} (x - c_q)(x - c_q)^T, \quad B = \sum_{q=1}^k n_q (c_q - c_x)(c_q - c_x)^T$$

где C_q – наблюдения из кластера q , c_q – центр кластера q , c_x – центр всего набора данных X , а n_q – число наблюдений в кластере q .

- Значение индекса тем больше, чем разделеннее кластеры и чем более сгруппированы в кластерах наблюдения. Ну и применять для случая сферических/эллиптических кластеров

Функционалы качества кластеризации. Davies-Bouldin Index

- Индекс Дэвиса–Болдина оценивает среднее "сходство" между кластерами, где сходство есть мера, которая сравнивает расстояние между кластерами с размером самих кластеров
- Пусть s_i есть среднее расстояние между каждой точкой кластера i и центроидом этого кластера, а d_{ij} есть расстояние между центроидами кластеров. Определим схожесть между кластерами следующим образом

$$R_{ij} = \frac{s_i + s_j}{d_{ij}}$$

- Тогда индекс имеет вид

$$DB = \frac{1}{k} \sum_{i=1}^k \max_{i \neq j} R_{ij}$$

- Чем индекс меньше, тем лучше. Близок к 0 – отличное разделение. Недостаток тот же, что и у индекса Калинского–Харабэса