



Национальный
исследовательский

Томский
государственный
университет



Кластеризация данных

Сергей Аксёнов, к.т.н., доцент кафедры
Теоретических основ информатики ИПМКН ТГУ

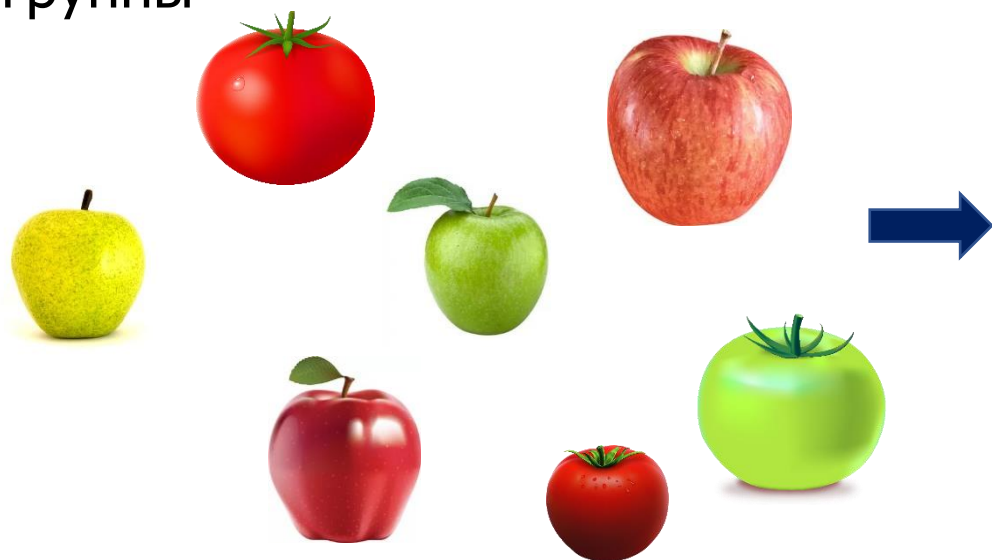
Томск-2023

Обучение без учителя

Меток класса нет. Метод используется для изучения данных.

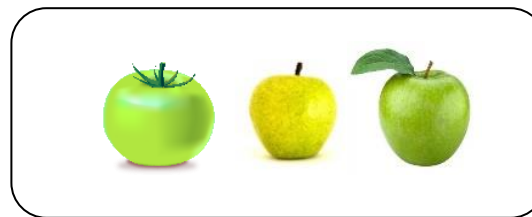
Особенность: Субъективность кластеризации.

Задача: Разложить объекты на две группы

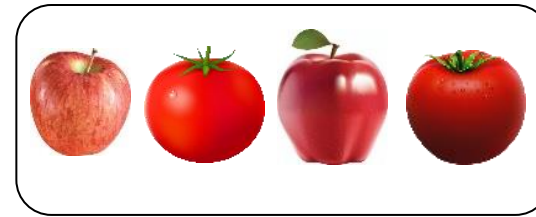


Решение А

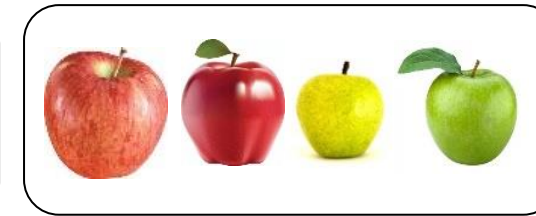
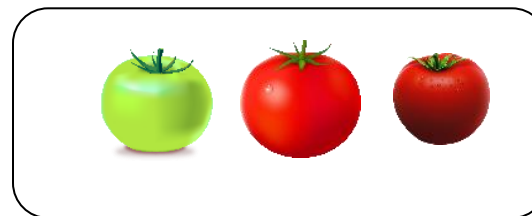
Группа 1



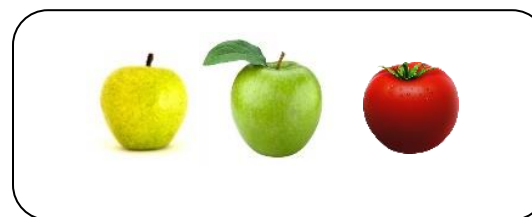
Группа 2



Решение В



Решение С

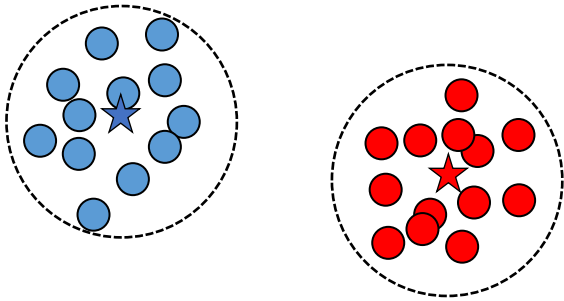


Разные решения!!!

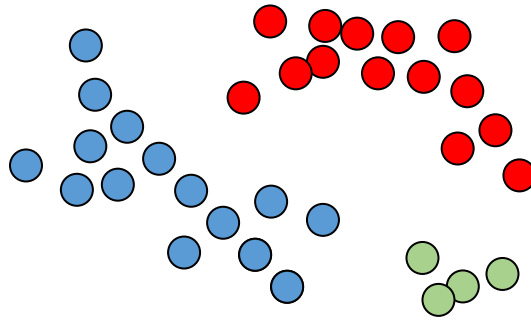
Цели кластеризации

1. Подготовка групп данных для последующей обработки
2. Нахождение аномалий
3. Снижение объема обрабатываемых данных
4. Построение иерархии объектов

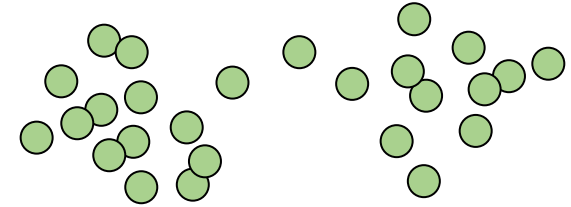
Примеры кластеров



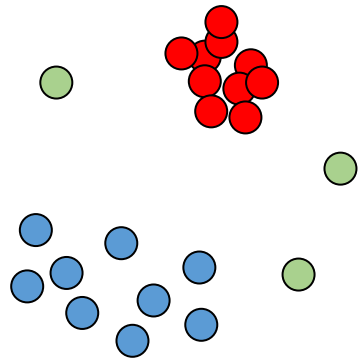
Кластеры с центроидами



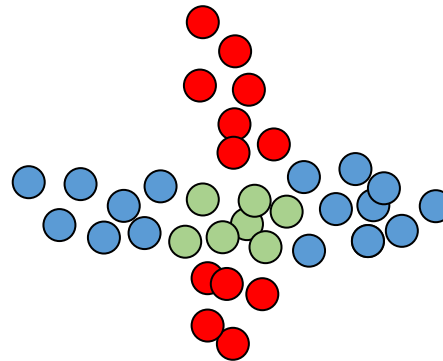
Кластеры произвольной формы



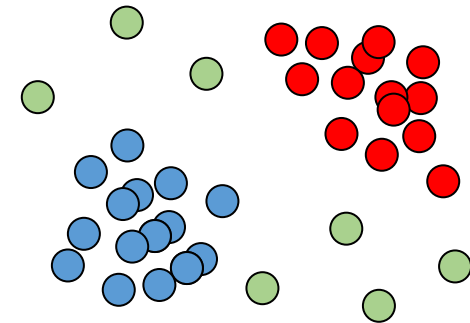
Кластеры с перемычками



Кластеры с разной
плотностью объектов



Кластеры перекрываются



Кластеры с шумами

Проблемы кластеризации

1. Не существует однозначного критерия качества кластеризации.

- Для хорошего решения нужно проводить оценку качества с помощью нескольких критериев

2. Число кластеров, как правило, заранее неизвестно и выбирается по субъективным критериям.

- Число кластеров подбирается путём изменения параметров алгоритма(ов)

3. Результат кластеризации существенно зависит от метрики.

Обозначения

X - Кластеризуемое множество

N - Количество элементов в X

c

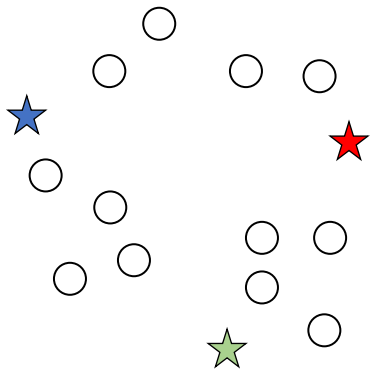
n_{c_j} - Число элементов в кластере c_j

v_j - Центр кластера c_j :
$$v_j = \sum_{x_i \in c_j} x_i / n_{c_j}$$

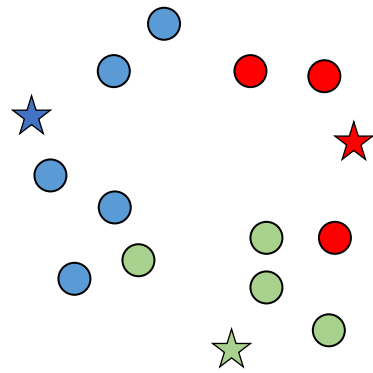
dim - Размерность множества X

u_{ij} - Степень принадлежности x_i кластеру c_j

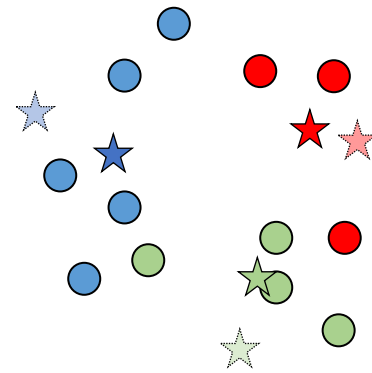
K-Средних: Пример



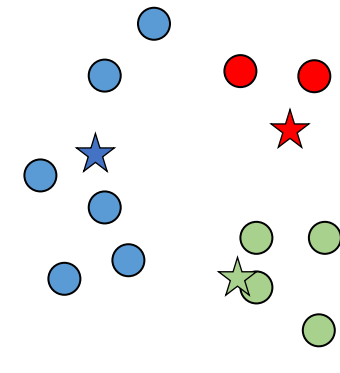
Инициализация



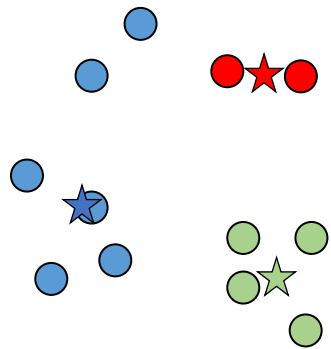
Связывание наблюдений



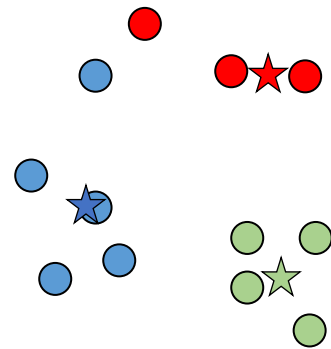
Сдвиг центроидов



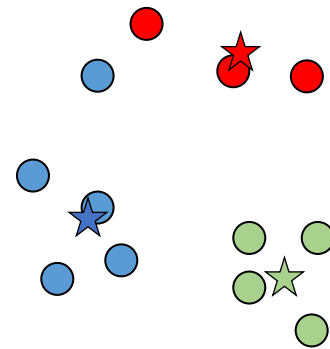
Связывание наблюдений



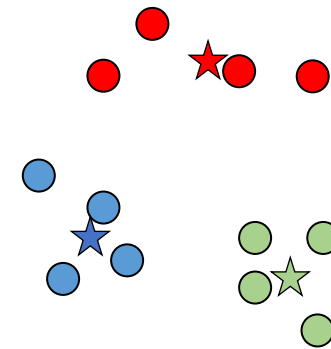
Сдвиг центроидов



Связывание наблюдений

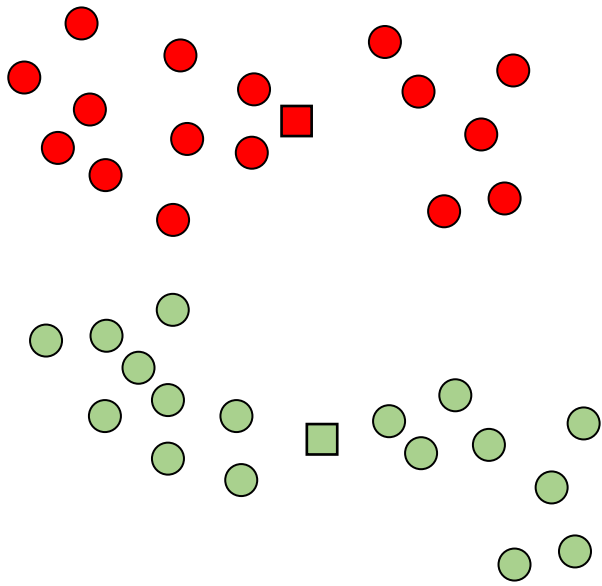


Сдвиг центроидов

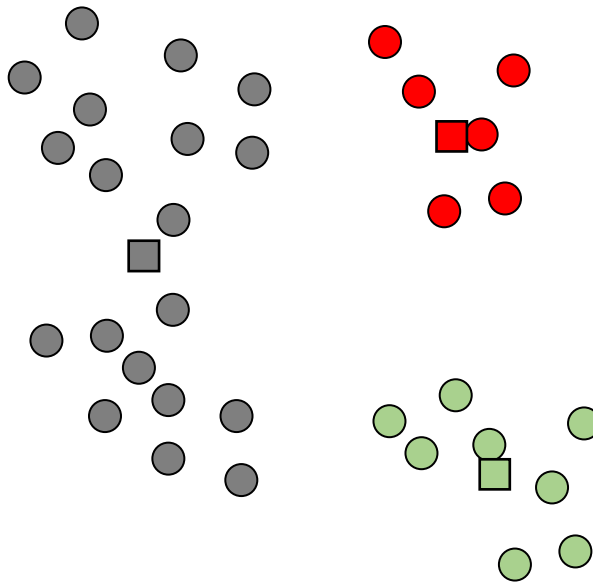


Сдвиг центроидов

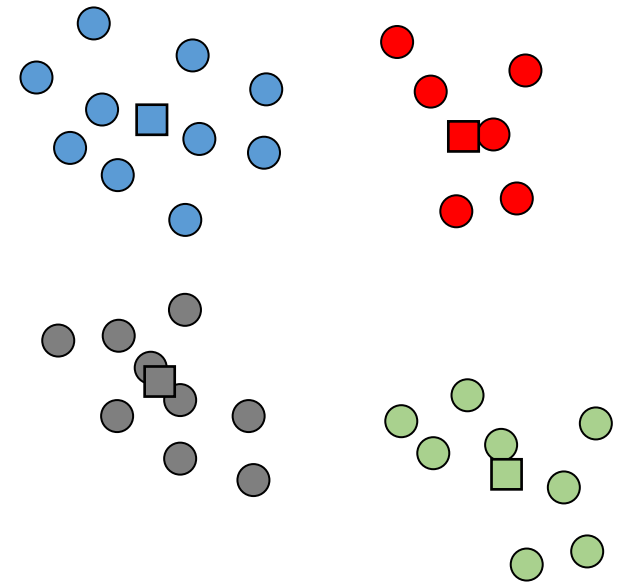
Число кластеров



2 Кластера



3 Кластера



4 Кластера

Кластеризация, основанная на плотности DBSCAN

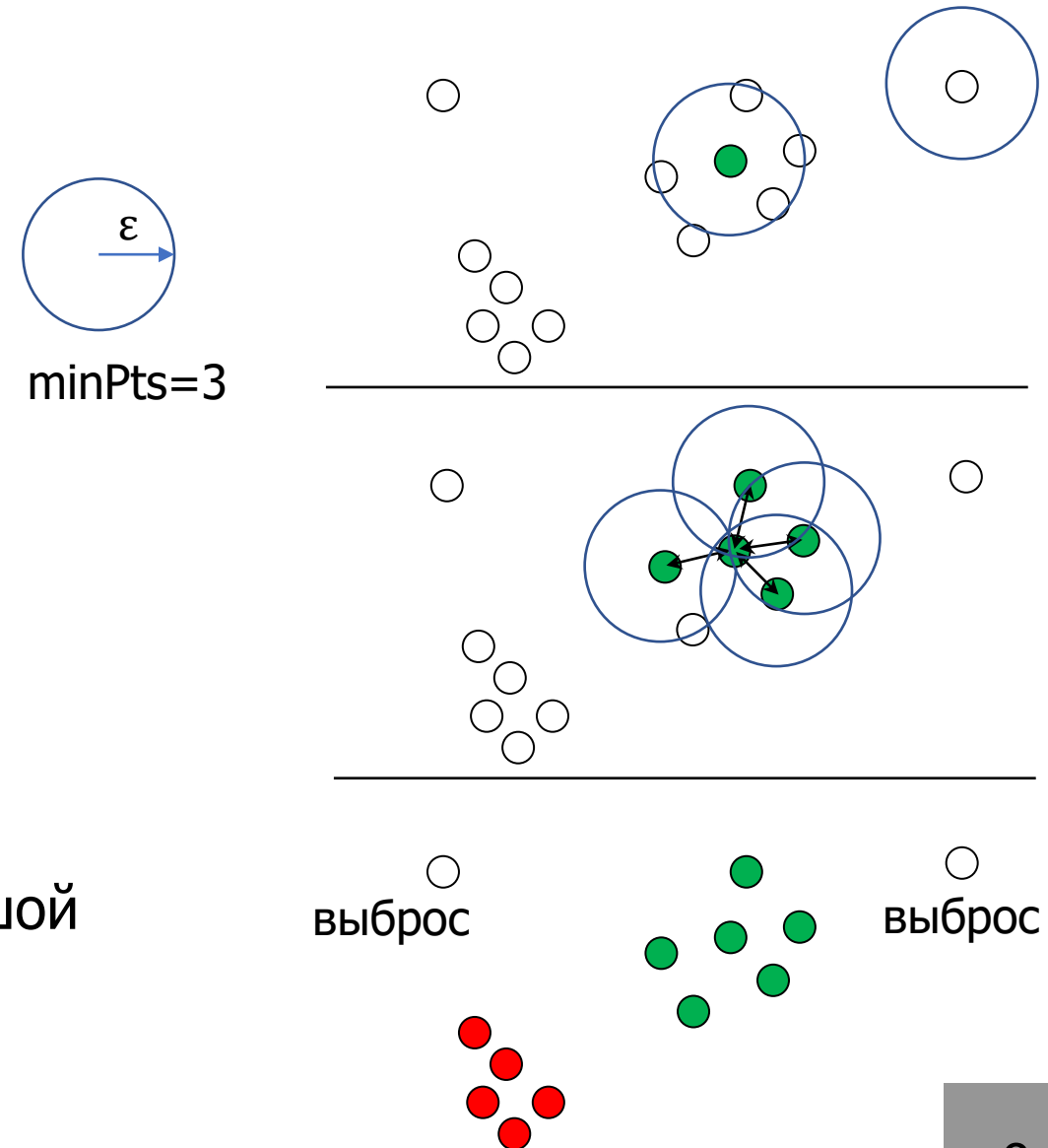
Параметры: minPts – число объектов, для формирования точки кластера, ϵ – радиус соседства

Преимущества:

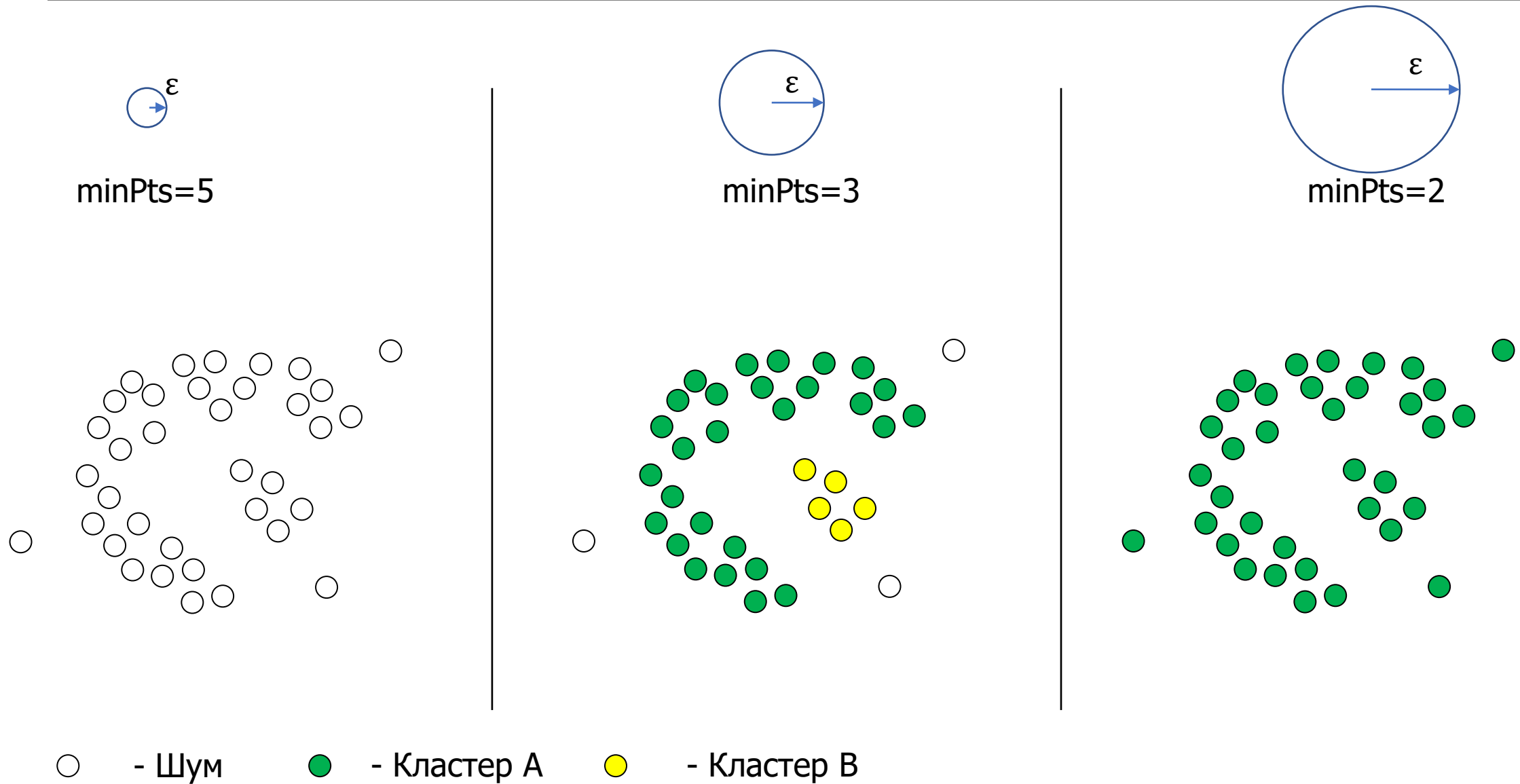
- Кластеры произвольной формы
- Не требуется задавать число кластеров
- Имеется понятие шума

Недостатки:

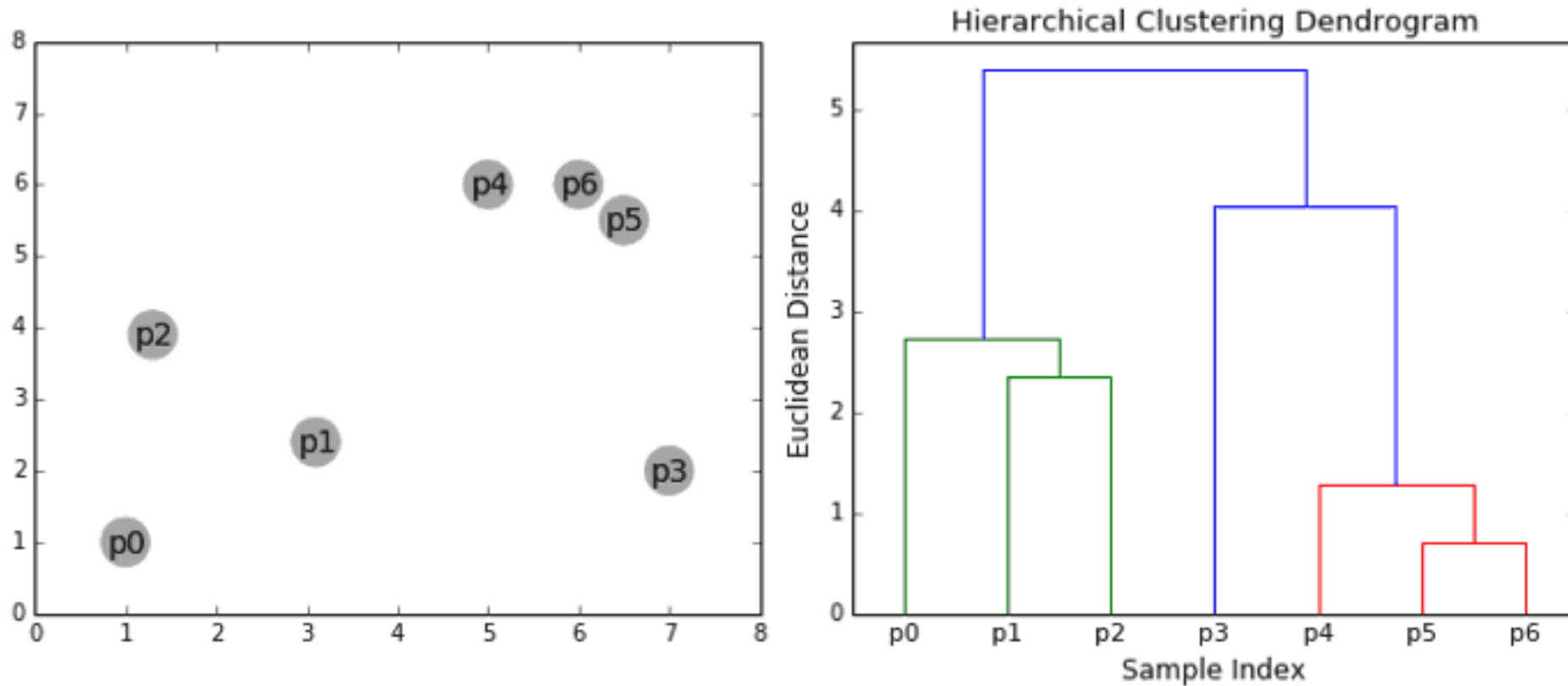
- Проблема краевых точек
- Трудность подбора параметров
- Плохо кластеризует наборы данных с большой разницей в плотности



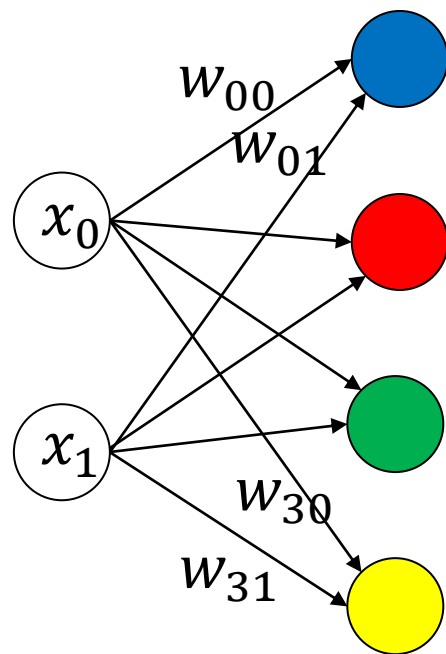
Пример кластеризации DBSCAN



Иерархическая кластеризация



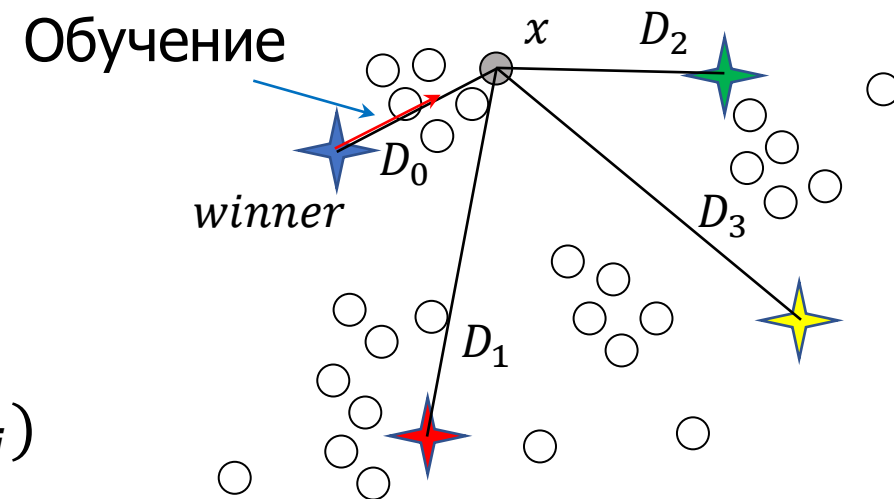
Сеть Кохонена (Kohonen network)



$$D_j = \sum_{i=1}^{dim} (x_i - w_{ij})^2$$

$$winner = k, \text{ где } D_k = \min_j (D_j)$$

$0 < \eta < 1$ - Скорость настройки



Алгоритм обучения

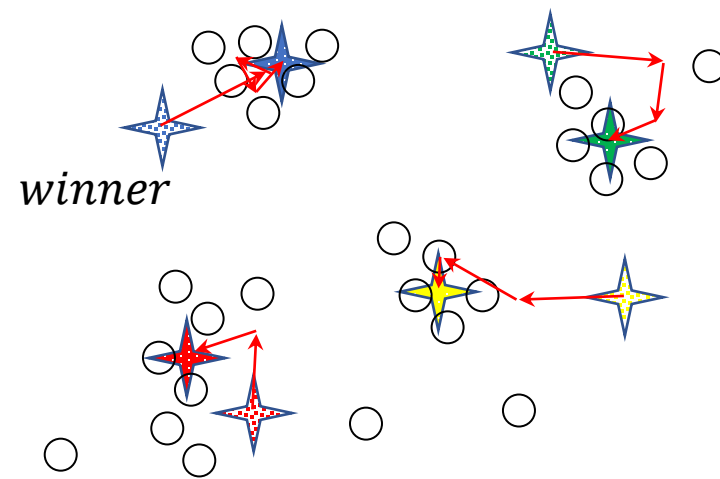
Шаг 1. Установка случайных весов w_{ij}

Шаг 2. Активация сети примером из выборки (расчет D_j)

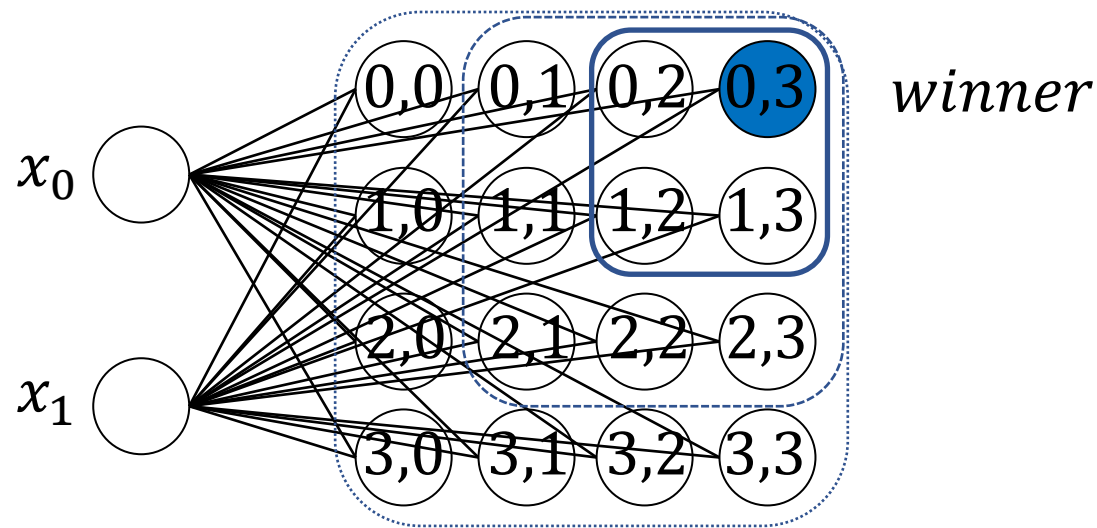
Шаг 3. Корректировка весов победителя по правилу:

$$w_{ki}(t+1) = w_{ki}(t) + \eta \cdot [x_i - w_{ki}(t)]$$

Шаг 4. Оценка критерия останова. В противном случае на шаг 2.



Самоорганизующаяся карта признаков (SOM)



Корректировка всех весов сети на этапе настройки.

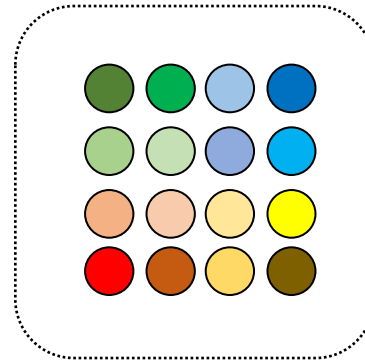
$$w_{ji}(t + 1) = w_{ji}(t) + \theta(j, winner, t) \cdot \eta \cdot [x_i - w_{ji}(t)]$$

$\theta(j, winner, t)$ - Функция расстояния, при $j = winner$: $\theta(j, j, t) = 1$

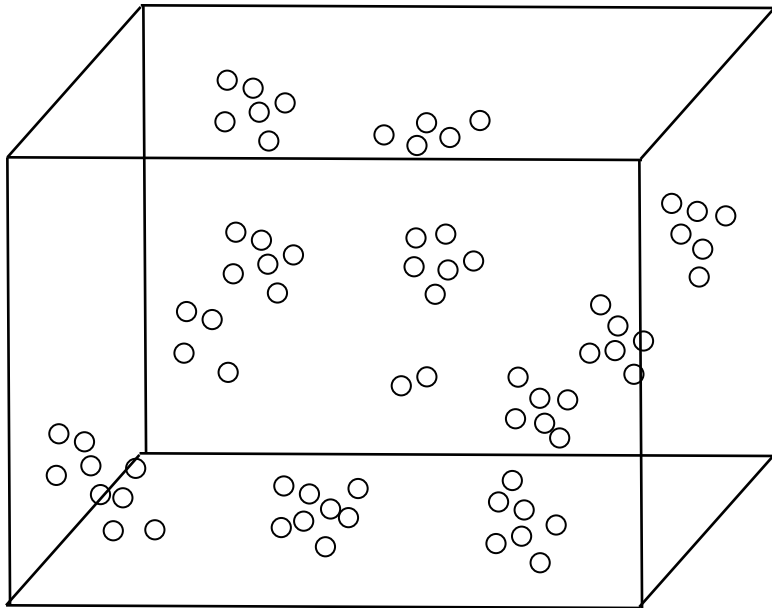
при $j \neq winner$: $0 < \theta(j, j, t) < 1$

Чем дальше настраиваемый узел от победителя, тем ближе $\theta(j, winner, t)$ к 0.

Пример SOM

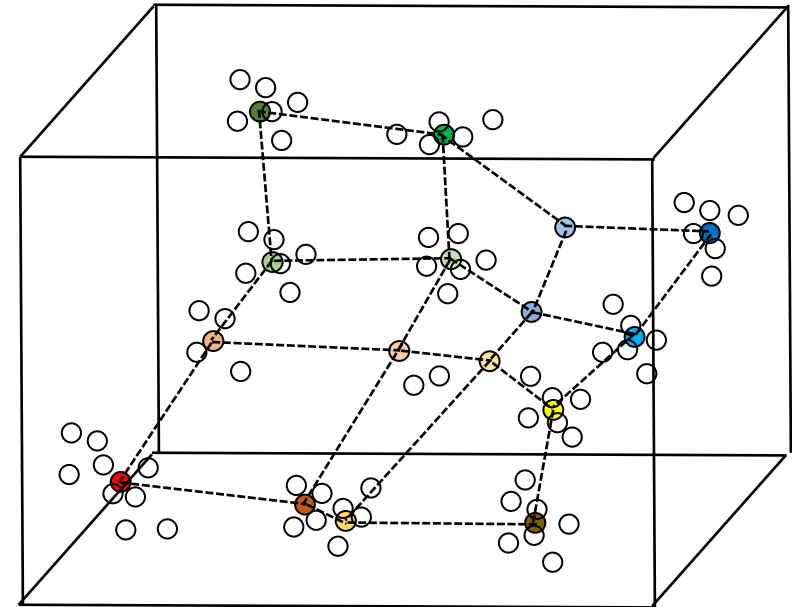


4×4 SOM



Пространство признаков

$3D \rightarrow 2D$



«Мертвые» нейроны: 

Относительные методы оценки качества четкой кластеризации

Оценка решения производится путём сравнения нескольких структур:

- Повторный запуск одного и того же алгоритма
- Запуск одного и того же алгоритма с разными значениями параметров
- Запуск разных алгоритмов

Критерии оценки качества:

Компактность – элементы, находящиеся в одном кластере, должны быть как можно ближе друг к другу

Отделимость – элементы, располагающиеся в разных кластерах, должны быть как можно дальше друг от друга

Индекс Данна (Dunn Index)

$$Dunn_index = \min_{i,j \in \{1, \dots, c\}, i \neq j} \left\{ \frac{d(c_i, c_j)}{\max_{k \in \{1, \dots, c\}} diam(c_k)} \right\}$$

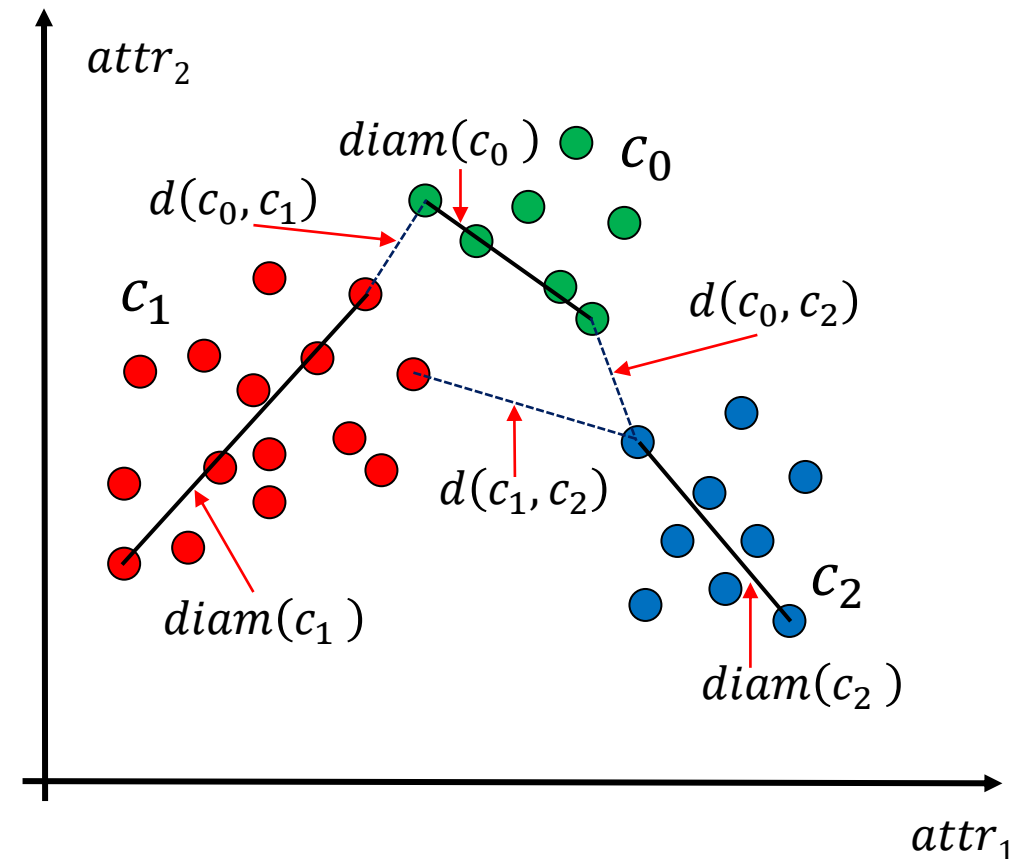
Диаметр кластера:

$$diam(c_i) = \max_{x,y \in c_i} \|x - y\|$$

Расстояние между кластерами:

$$d(c_i, c_j) = \min_{x \in c_i, y \in c_j} \|x - y\|$$

$Dunn_index \rightarrow Max$



Индекс Дэвиса-Болдуина (Davis-Baldwin Index)

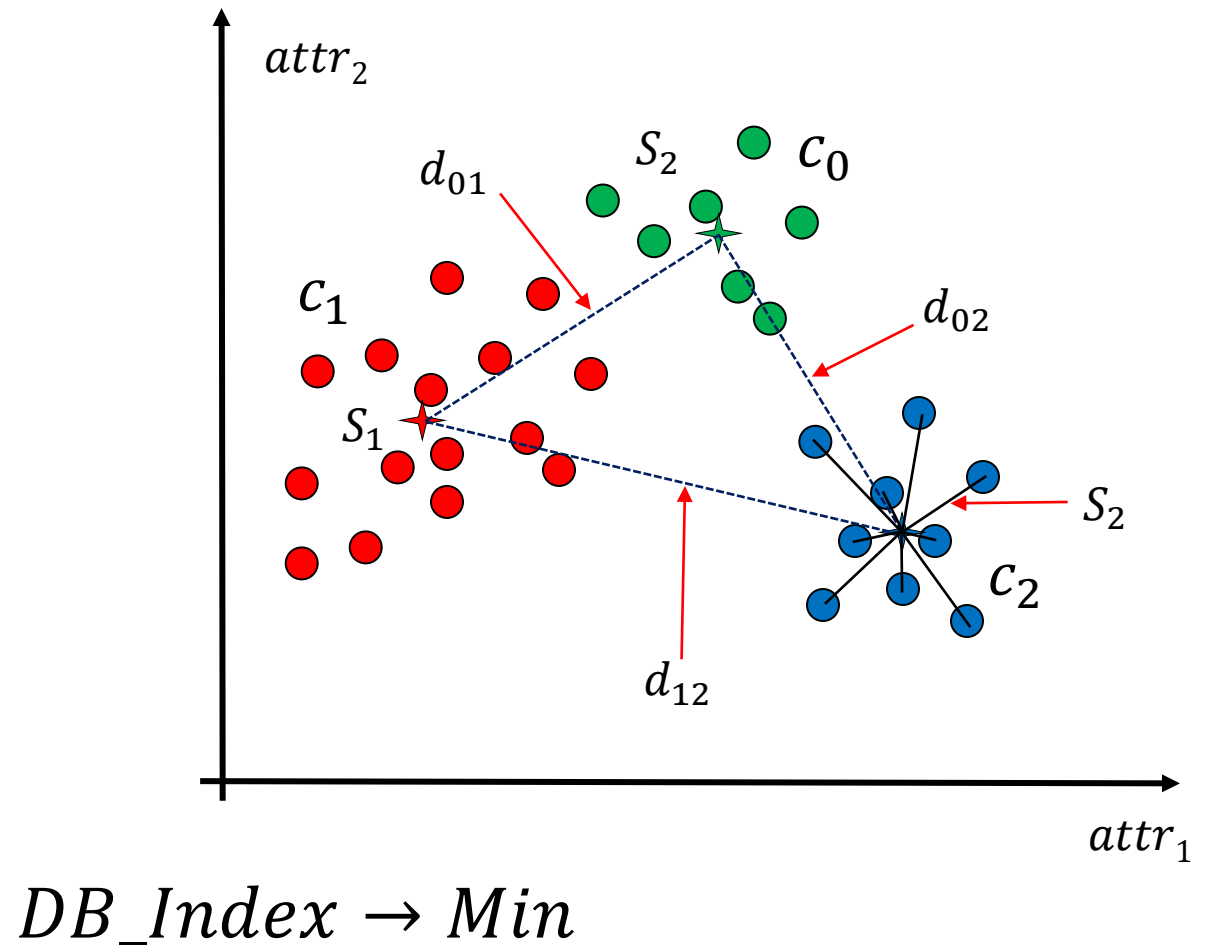
$$S_i = \left\{ \frac{1}{n_{c_i}} \sum_{x \in c_i} \|x - v_i\|^q \right\}^{\frac{1}{q}}$$

$$d_{ij} = \left\{ \sum_{k=1}^{dim} |v_i^k - v_j^k|^p \right\}^{\frac{1}{p}}$$

$$R_{ij} = \frac{S_i + S_j}{d_{ij}}$$

$$R_i = \max_{i,j \in \{1, \dots, c\}} (R_{ij})$$

$$DB_Index = \frac{1}{c} \sum_{i=1}^c R_i$$



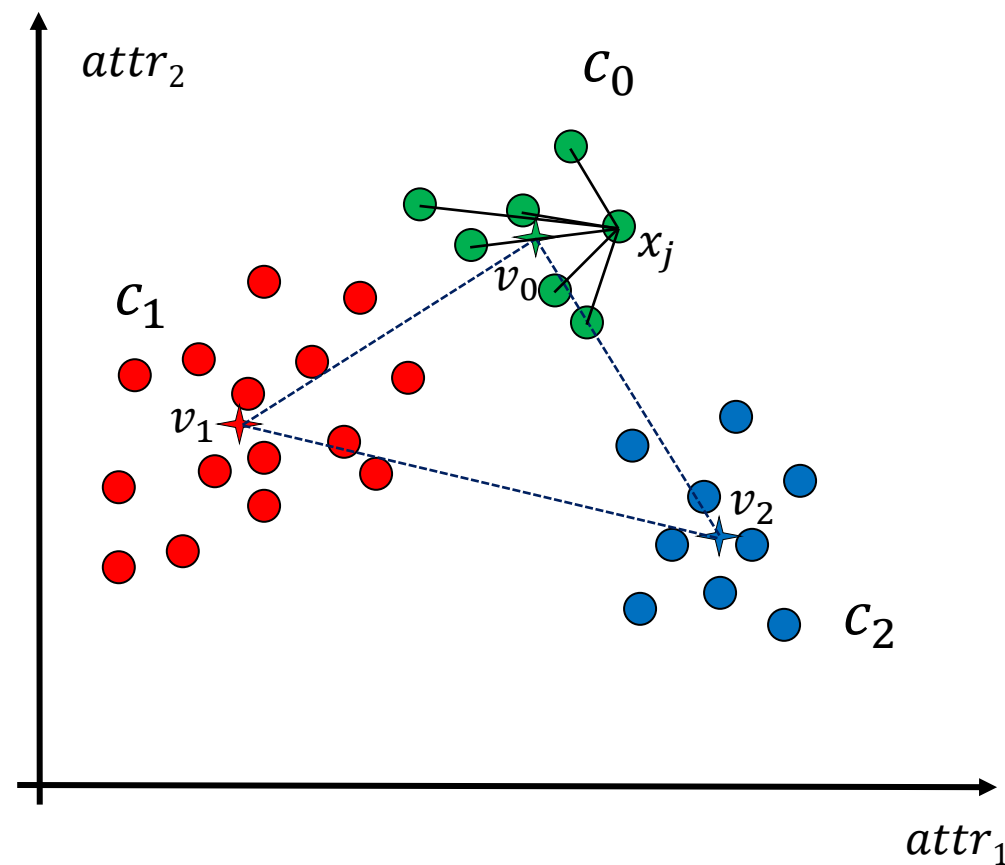
CS Индекс (CS Index)

$$CS_Index = \frac{\sum_{i=1}^c \left\{ \sum_{x_j, x_k \in c_i} \max(\|x_j - x_k\|) \right\}}{\sum_{i=1}^c \min_{i \neq j} (\|v_i - v_j\|)}$$

Компактность – сумма расстояний между всеми объектами в группе.

Для отделимости используется сумма наименьших расстояний между центрами кластеров.

$CS_Index \rightarrow Min$



Индекс оценки Силуэта (Silhouette Index)

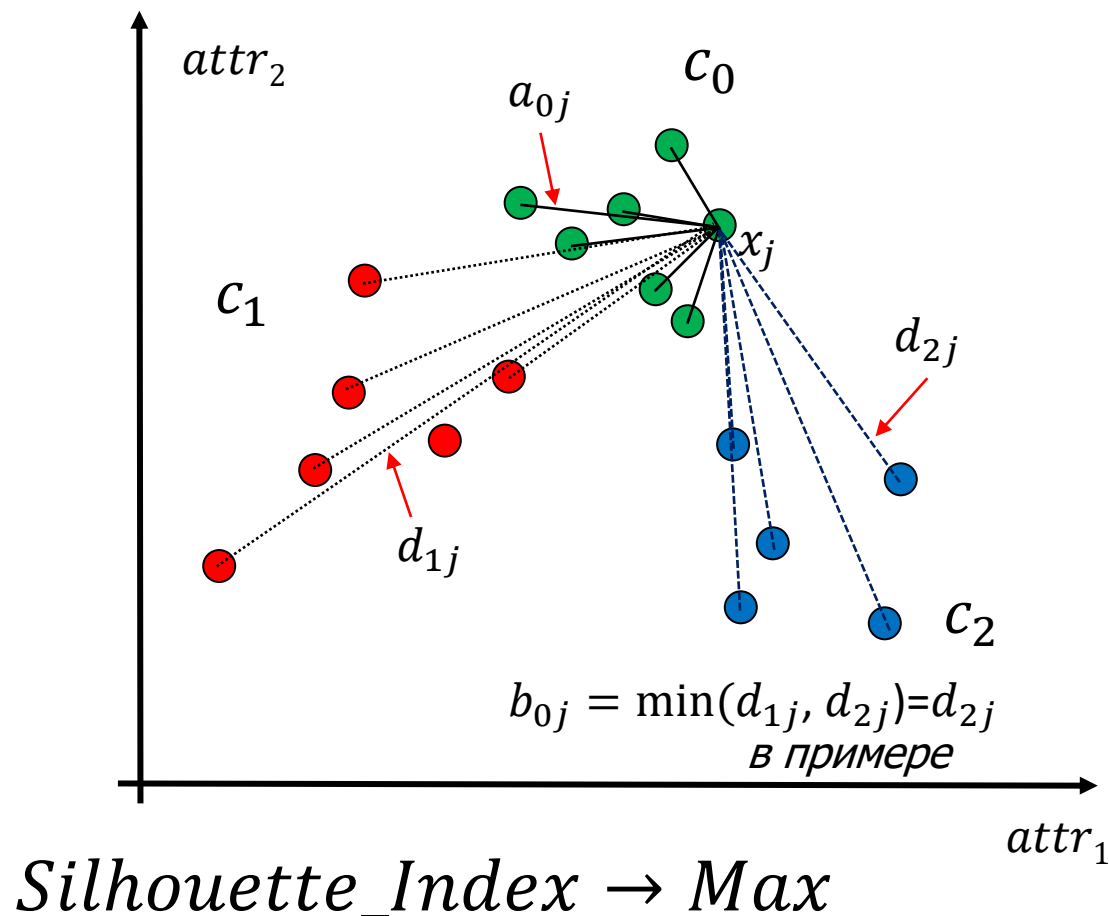
$$a_{pj} = \frac{1}{n_{c_p} - 1} \sum_{x_k \in c_p} \|x_j - x_k\|$$

$$d_{qj} = \frac{1}{n_{c_q}} \sum_{x_k \in c_q} \|x_j - x_k\|$$

$$b_{pj} = \min_{q \neq p} d_{qj}$$

$$S_{x_j} = \frac{b_{pj} - a_{pj}}{\max(a_{pj}, b_{pj})}$$

$$Silhouette_Index = \frac{1}{N} \sum_{j=1}^N S_{x_j}$$



Нечеткая кластеризация: GK и Fuzzy K-means

$$v_j = \frac{\sum_{i=1}^N (u_{ij})^m \cdot x_i}{\sum_{i=1}^N (u_{ij})^m}$$

$$D_{ij} = \sqrt{\|x_i - v_j\|^2}$$

$$A_j = \frac{\sum_{i=1}^N (u_{ij})^m \cdot (x_i - v_j)^T \cdot (x_i - v_j)}{\sum_{i=1}^N (u_{ij})^m}$$

$$D_{A_j} = (x_i - v_j) \cdot \left[\left(\det(A_j) \right)^{\frac{1}{N}} \cdot A_j^{-1} \right] \cdot (x_i - v_j)^T$$

$$u_{ij} = \left((D_{ij})^2 \sum_{k=1}^c \frac{1}{(D_{ik})^2} \right)^{-\frac{1}{m-1}}$$

Индекс Си-Бени (XB - Xie-Beni Index)

$$Xie - Beni_Index = \frac{\sum_{i=1}^c \sum_{j=1}^N (u_{ij})^m \|x_j - v_i\|^2}{N \cdot \min_{l \neq s} \|v_l - v_s\|^2}$$

$$Xie - Beni_Index \rightarrow Min$$

Нечеткий упрощенный Силуэт (Fuzzy Simplified Silhouette)

$$S_j = \frac{b_{pj} - a_{pj}}{\max\{a_{pj}, b_{pj}\}}$$

$$F_S_Silhoutte_Index = \frac{\sum_{j=1}^N (u_{pj} - u_{qj})^\alpha S_j}{\sum_{j=1}^N (u_{pj} - u_{qj})^\alpha}$$

$$F_S_Silhoutte_Index \rightarrow Max$$

Индекс Квона (Kwon Index)

$$Kwon_Index = \frac{\sum_{i=1}^c \sum_{j=1}^N (u_{ij})^m \|x_j - v_i\|^2 + \frac{1}{c} \sum_{i=1}^c \|v_i - \bar{x}\|^2}{\min_{l \neq s} \|v_l - v_s\|^2}$$

$$\bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$Kwon_Index \rightarrow Min$

TSS индекс (Tang-Sun-Sun Index)

$$TSS_Index = \frac{\sum_{i=1}^c \sum_{j=1}^N (u_{ij})^m \|x_j - v_i\|^2 + \frac{1}{c(c-1)} \sum_{i=1}^c \sum_{l=1, l \neq i}^c \|v_i - \bar{x}\|^2}{\min_{l \neq s} \|v_l - v_s\|^2 + 1/k}$$

$TSS_Index \rightarrow Min$

Индекс Фукуямы-Сугено (Fukuyama-Sugeno Index)

$$FS_Index = \sum_{i=1}^c \sum_{j=1}^N (u_{ij})^m \|x_j - v_i\|^2 + \sum_{i=1}^c \sum_{j=1}^N (u_{ij})^m \|v_i - \bar{x}\|^2$$

$FS_Index \rightarrow Min$