

## Дано:

$X$  — пространство объектов;

$X^\ell = \{x_i\}_{i=1}^\ell$  — обучающая выборка;

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

## Найти:

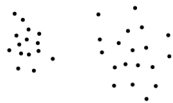
$y_i \in Y$  — метки кластеров объектов:

- каждый кластер состоит из близких объектов;
- объекты разных кластеров существенно различны.

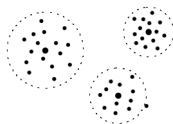
Кластеризация — это *обучение без учителя*.

Решение задачи кластеризации принципиально неоднозначно:

- различные критерии качества кластеризации
- различные эвристические методы кластеризации
- различные варианты функции расстояния  $\rho$



внутрикластерные расстояния, как правило, меньше межкластерных



кластеры с центром



кластеры могут соединяться перемычками



кластеры могут накладываться на разреженный фон из редко расположенных объектов

# Типы кластерных структур



ленточные кластеры



перекрывающиеся кластеры



кластеры могут образовываться не по сходству, а по иным типам регулярностей



кластеры могут вообще отсутствовать

Объекты  $x_i$  задаются векторами признаков  $(f_1(x_i), \dots, f_n(x_i))$ .

**Вход:**  $X^\ell$  — обучающая выборка, параметр  $k$ ;

**Выход:** центры кластеров  $\mu_y, y \in Y$ ;

1: начальное приближение центров  $\mu_y, y \in Y$ ;

2: **повторять**

3: отнести каждый  $x_i$  к ближайшему центру:

$$y_i := \arg \min_{y \in Y} \rho(x_i, \mu_y), \quad i = 1, \dots, \ell;$$

4: вычислить новые положения центров:

$$\mu_{yj} := \frac{\sum_{i=1}^{\ell} [y_i = y] f_j(x_i)}{\sum_{i=1}^{\ell} [y_i = y]}, \quad y \in Y, \quad j = 1, \dots, n;$$

5: **пока**  $y_i$  не перестанут изменяться;

**Вход:**  $X^\ell$  — обучающая выборка, параметр  $k$ ;

**Выход:** центры кластеров  $\mu_y$ ,  $y \in Y$ ;

1: начальное приближение центров  $\mu_y$ ,  $w_y := \frac{1}{|Y|}$ ,  $y \in Y$ ;

2: **повторять**

3: оценить близость каждого  $x_i$  ко всем центрам:

$$g_{iy} := w_y \exp\left(-\frac{1}{2}\rho^2(x_i, \mu_y)\right), \quad i = 1, \dots, \ell, \quad y \in Y;$$

$$g_{iy} := \frac{g_{iy}}{\sum_{z \in Y} g_{iz}} \text{ — нормированные близости};$$

4: отнести каждый  $x_i$  к ближайшему центру:

$$y_i := \arg \max_{y \in Y} g_{iy}, \quad i = 1, \dots, \ell;$$

5: новые положения центров и мощности кластеров:

$$\mu_{yj} := \frac{1}{\ell w_y} \sum_{i=1}^{\ell} g_{iy} f_j(x_i), \quad w_y := \frac{1}{\ell} \sum_{i=1}^{\ell} g_{iy}, \quad y \in Y, \quad j = 1, \dots, n;$$

6: **пока**  $g_{iy}$  не перестанут изменяться;

**Вход:**  $X^\ell$  — обучающая выборка, параметры  $q, \delta, k$

**Выход:**  $U \subset X^\ell$  — начальные приближения центров  $\mu_y, y \in Y$ ;

1: среднее расстояние до  $q$  ближайших соседей:

$$R_i := \frac{1}{q} \sum_{j=1}^q \rho(x_i, x_i^{(j)}), \text{ для всех } i = 1, \dots, \ell,$$

где  $x_i^{(j)}$  —  $j$ -й ближайший сосед объекта  $x_i$ ;

2: отбросить шумовые объекты:

$$X' := \{x_i \in X^\ell \mid R_i \leq \Delta\} \text{ при } \Delta: |X'| = (1 - \delta)\ell;$$

3: выбрать пару самых удалённых объектов:

$$U := \arg \max_{x, x' \in X'} \rho(x, x');$$

далее последовательно присоединять к  $U$  по одному объекту, самому удалённому от уже выбранных:

4: **повторять  $k - 2$  раз**

$$U := U \cup \arg \max_{x \in X'} \min_{u \in U} \rho(x, u);$$

- Чувствительность к выбору начального приближения
- Необходимость задавать  $k$

### Способы устранения этих недостатков:

- Эвристики для выбора начального приближения
- Мягкая кластеризация
- Мультистарт: несколько случайных инициализаций; выбор лучшей кластеризации по функционалу качества.
- Быстрые алгоритмы ( $k$ -means++, сэмплирование)
- Варьирование числа кластеров  $k$  в ходе итераций