

# *Кластеризация данных*

Курс «Интеллектуальные информационные системы»

Кафедра управления и информатики НИУ «МЭИ»

Осень 2018 г.

# Что такое кластеризация?

**Кластеризация** - задача разбиения заданной выборки *объектов* на непересекающиеся подмножества, называемые кластерами, так, чтобы каждый кластер состоял из схожих объектов, а объекты разных кластеров существенно отличались.

Под схожестью обычно понимается близость друг к другу относительно выбранной метрики.

Задача кластеризации относится к разделу задач обучения без учителя.

**Обучение без учителя** (Unsupervised learning) — один из разделов машинного обучения. Изучает широкий класс задач обработки данных, в которых известны только описания множества объектов (обучающей выборки), и требуется обнаружить внутренние взаимосвязи, зависимости, закономерности, существующие между объектами.

Обучение без учителя часто противопоставляется обучению с учителем, когда для каждого обучающего объекта задаётся «правильный ответ», и требуется найти зависимость между объектами и ответами.

# Постановка задачи кластеризации

## Дано:

$X$  – пространство объектов

$\vec{X}_l$  – обучающая выборка;  $l = 1 \dots L$

$\rho$  - функция расстояния между объектами

## Найти:

$Y$  – множество кластеров и

а:  $X \rightarrow Y$  – алгоритм кластеризации, такие, что:

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

# Особенности задачи кластеризации

Решение задачи кластеризации принципиально неоднозначно:

- Точной постановки задачи кластеризации нет
- Существует множество критериев качества кластеризации
- Существует множество эвристических методов кластеризации
- Число кластеров  $|Y|$  заранее, как правило, не известно
- Результат кластеризации существенно зависит от метрики  $\rho$ , которую эксперт задает субъективно

# Цели кластеризации

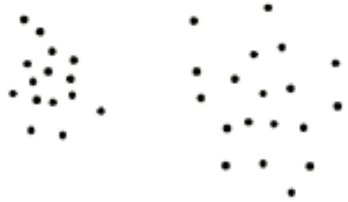
**Понимание данных путём выявления кластерной структуры.** Разбиение выборки на группы схожих объектов позволяет упростить дальнейшую обработку данных и принятия решений, применяя к каждому кластеру свой метод анализа (стратегия «разделяй и властвуй»).

**Сжатие данных.** Если исходная выборка избыточно большая, то можно сократить её, оставив по одному наиболее типичному представителю от каждого кластера.

**Обнаружение новизны (novelty detection).** Выделяются нетипичные объекты, которые не удаётся присоединить ни к одному из кластеров.

**Построение иерархии множества объектов (задача таксономии)**

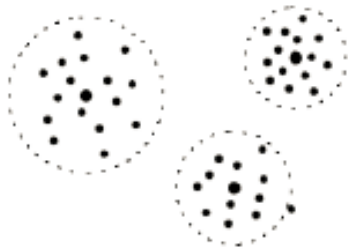
# Примеры кластерных структур



внутрикластерные расстояния, как правило,  
меньше межкластерных



ленточные кластеры



кластеры с центром

# Примеры кластерных структур



кластеры могут соединяться перемычками

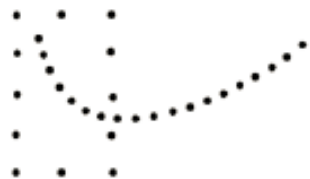


кластеры могут накладываться на разреженный фон из редко расположенных объектов



кластеры могут перекрываться

# Примеры кластерных структур



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

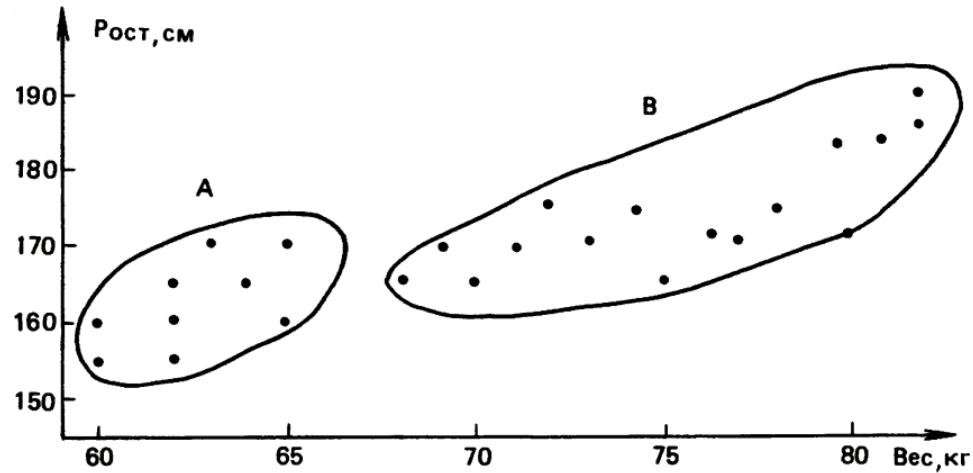


кластеры могут вообще отсутствовать

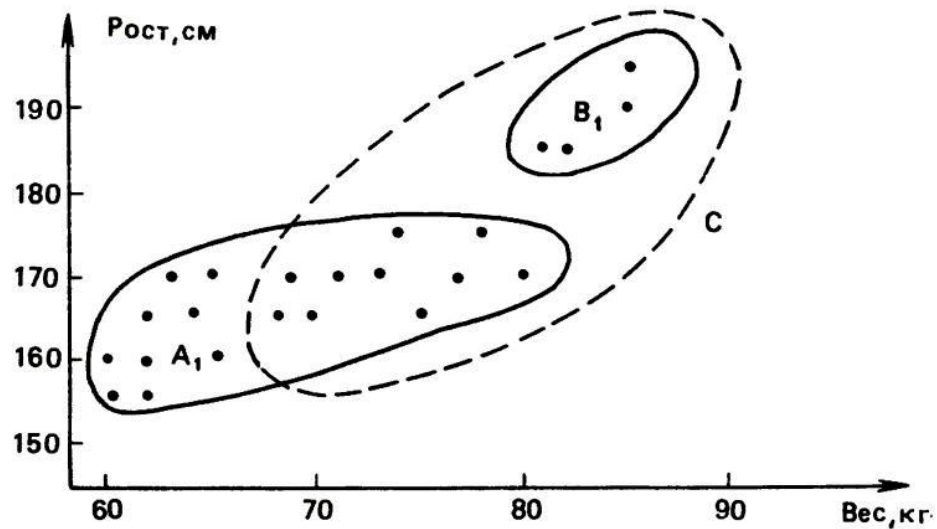
- Каждый метод кластеризации имеет свои ограничения и выделяет кластеры лишь некоторых типов
- Понятие «тип кластерной структуры» зависит от метода и не имеет формального определения



# Проблема чувствительности к метрике



А – девушки  
В – молодые люди



После перенормировки  
(сжали ось «Вес» вдвое)

## Качество кластеризации

**Сумма средних внутрикластерных расстояний:**

$$F_0 = \sum_{k \in Y} \frac{1}{N_k} \sum_{l=1}^{N_k} \rho(\vec{X}_l, \mu_k) \rightarrow \min$$

**Сумма межкластерных расстояний:**

$$F_1 = \sum_{j,k \in Y} \rho(\mu_j, \mu_k) \rightarrow \max$$

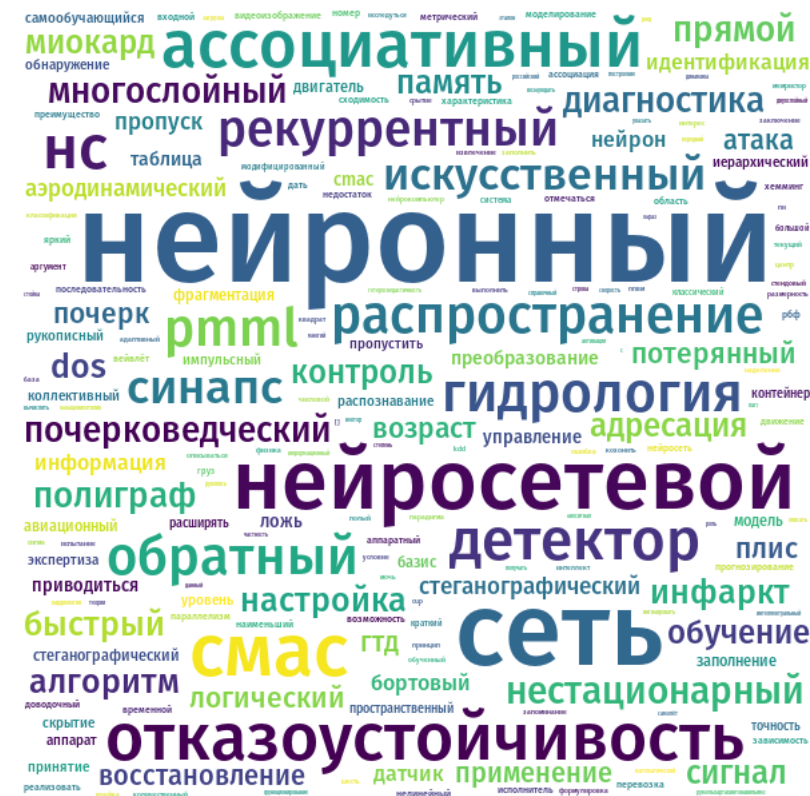
### Обобщенный функционал:

$$F_2 = \frac{F_0}{F_1} \rightarrow \min$$

$\mu_k$  - Центр масс кластера k

$N_k$  - Размер кластера  $k$

В задачах кластеризации текстов качество кластеризации можем косвенно оценить по наиболее частотным терминам, встречающимся в классе («Облако тэгов»). Т.е. мы могли бы дать название каждому кластеру исходя из наиболее частотных терминов :



# Алгоритмы кластеризации

## Иерархические

- Агломеративная кластеризация
- Дивизимная кластеризация

## Статистические

- К-средних (k-means)
- ЕМ-алгоритмы
- Алгоритм FOREL

## Сети Кохонена

# Иерархические алгоритмы кластеризации

Среди алгоритмов иерархической кластеризации различаются два основных типа. Дивизимные или нисходящие алгоритмы разбивают выборку на всё более и более мелкие кластеры. Более распространены агломеративные или восходящие алгоритмы, в которых объекты объединяются во всё более и более крупные кластеры

Сначала каждый объект считается отдельным кластером. Для одноэлементных кластеров естественным образом определяется функция расстояния  $\rho(x_j, x_k)$

Затем запускается процесс слияний. На каждой итерации вместо пары самых близких кластеров  $U$  и  $V$  образуется новый кластер  $W = U \cup V$

Расстояние от нового кластера  $W$  до любого другого кластера  $S$  вычисляется по расстояниям  $R(U, V)$ ,  $R(U, S)$  и  $R(V, S)$ :

$$R(U \cup V, S) = \alpha_u R(U, S) + \alpha_v R(V, S) + \beta R(U, V) + \gamma |R(U, S) - R(V, S)|$$

где  $\alpha_u, \alpha_v, \beta, \gamma$ - числовые параметры

Эта универсальная формула обобщает практически все разумные способы определить расстояние между кластерами. Она была предложена Лансом и Уильямсом в 1967 году.

# Иерархические алгоритмы кластеризации

На практике используются следующие способы вычисления расстояний  $R(W, S)$  между кластерами  $W$  и  $S$ . Для каждого из них доказано соответствие формуле Ланса-Вильямса при определённых сочетаниях параметров:

Расстояние ближнего соседа (single linkage):

$$R^b(W, S) = \min_{w \in W, s \in S} \rho(w, s); \quad \alpha_u = \alpha_v = \frac{1}{2}, \quad \beta = 0, \quad \gamma = -\frac{1}{2}$$

Расстояние дальнего соседа (complete linkage):

$$R^d(W, S) = \max_{w \in W, s \in S} \rho(w, s); \quad \alpha_u = \alpha_v = \frac{1}{2}, \quad \beta = 0, \quad \gamma = \frac{1}{2}$$

Расстояние до центра:

$$R^c(W, S) = \rho^2 \left( \sum_{w \in W} \frac{w}{|W|}, \sum_{s \in S} \frac{s}{|S|} \right); \quad \alpha_u = \frac{|U|}{|W|}, \quad \alpha_v = \frac{|V|}{|W|}, \quad \beta = -\alpha_u \alpha_v, \quad \gamma = 0$$

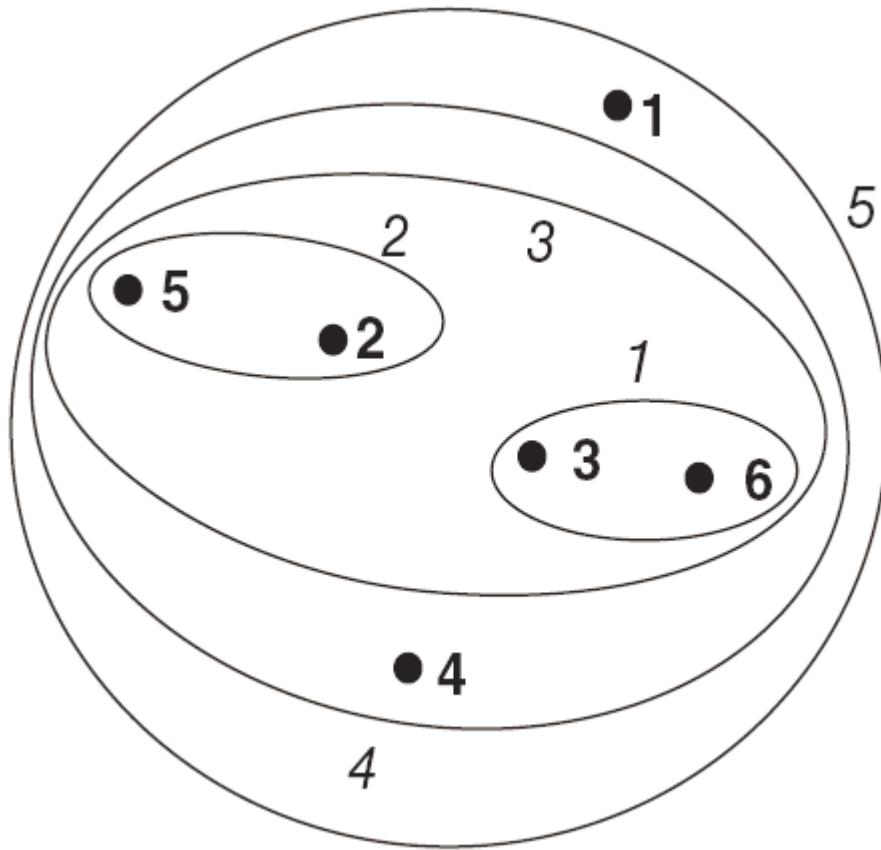
Расстояние Уорда (Варда, Ward):

$$R^y(W, S) = \frac{|S||W|}{|S| + |W|} \rho^2 \left( \sum_{w \in W} \frac{w}{|W|}, \sum_{s \in S} \frac{s}{|S|} \right); \quad \alpha_u = \frac{|S| + |U|}{|S| + |W|}, \quad \alpha_v = \frac{|S| + |V|}{|S| + |W|}, \quad \beta = -\frac{|S|}{|S| + |W|}, \quad \gamma = 0$$

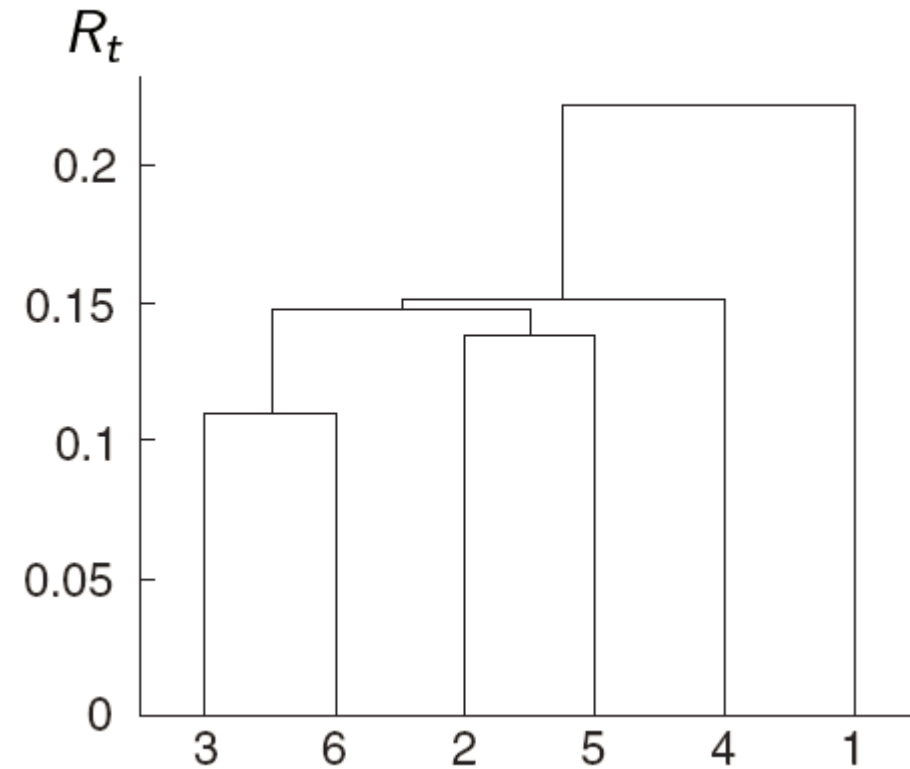
# Иерархические алгоритмы кластеризации

## Расстояние ближнего соседа

Диаграмма вложения



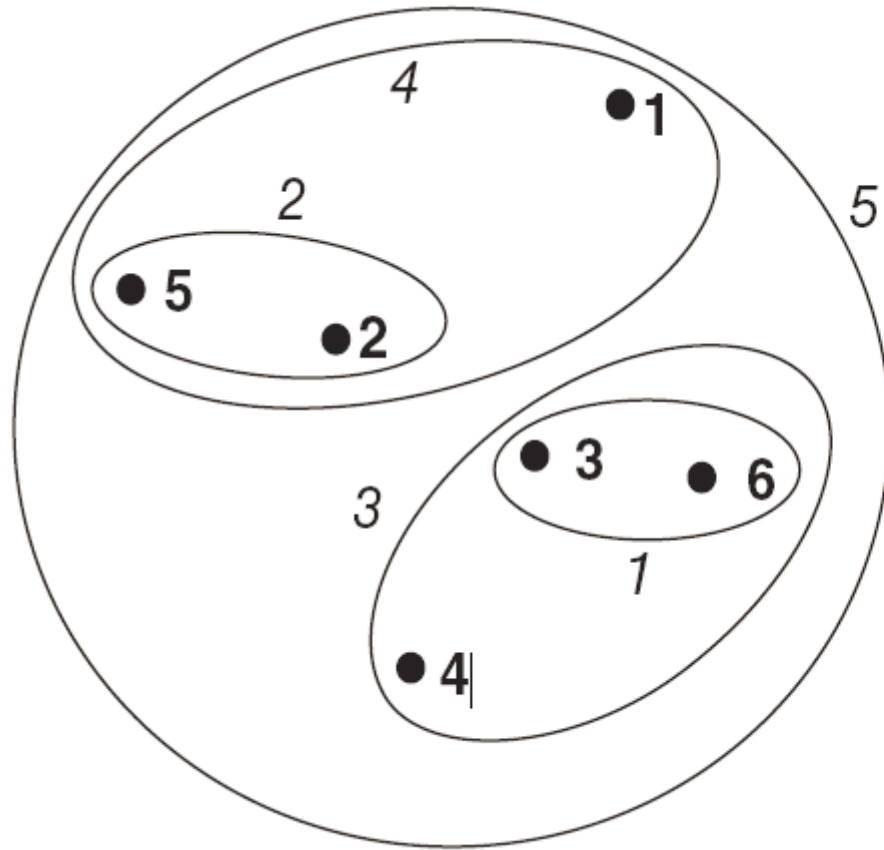
Дендрограмма



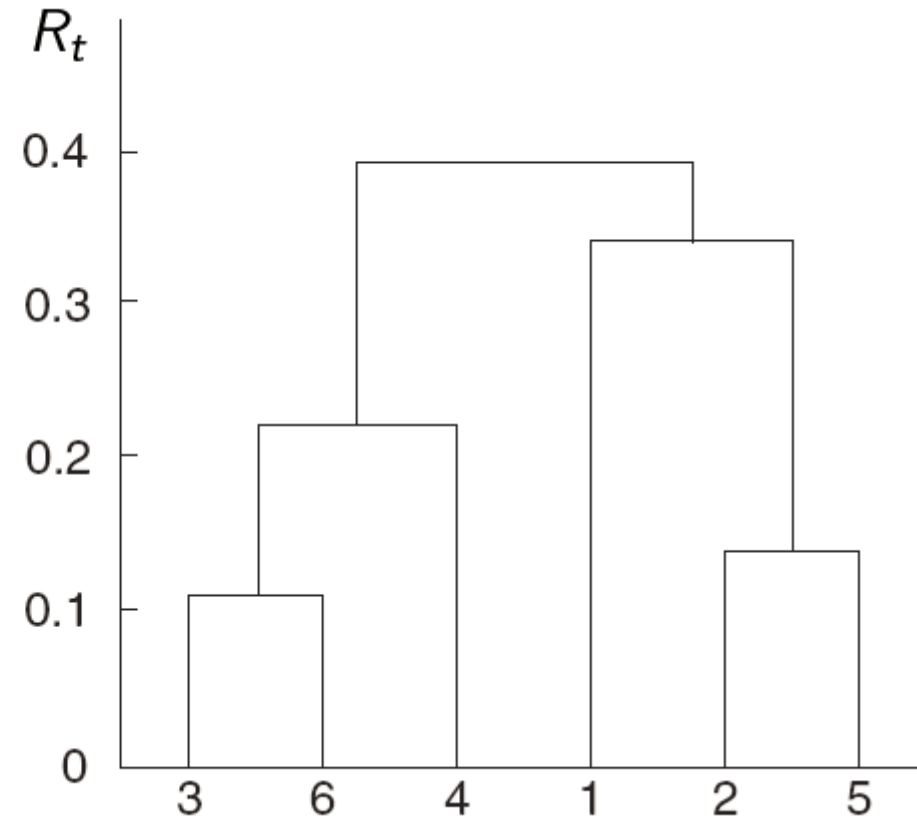
# Иерархические алгоритмы кластеризации

## Расстояние дальнего соседа

Диаграмма вложения



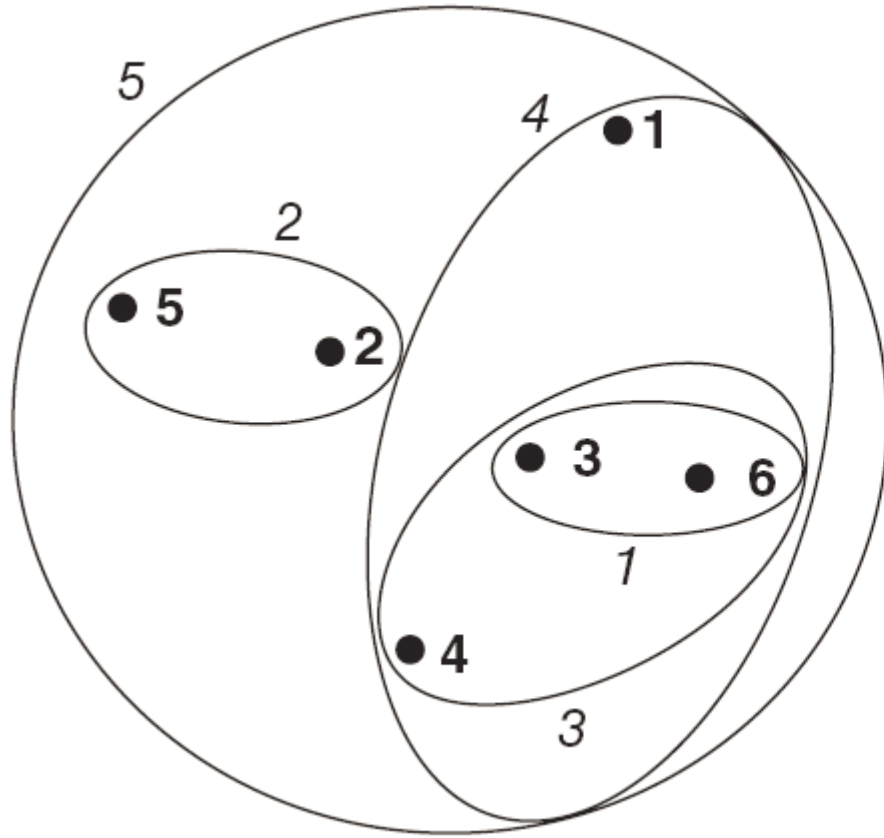
Дендрограмма



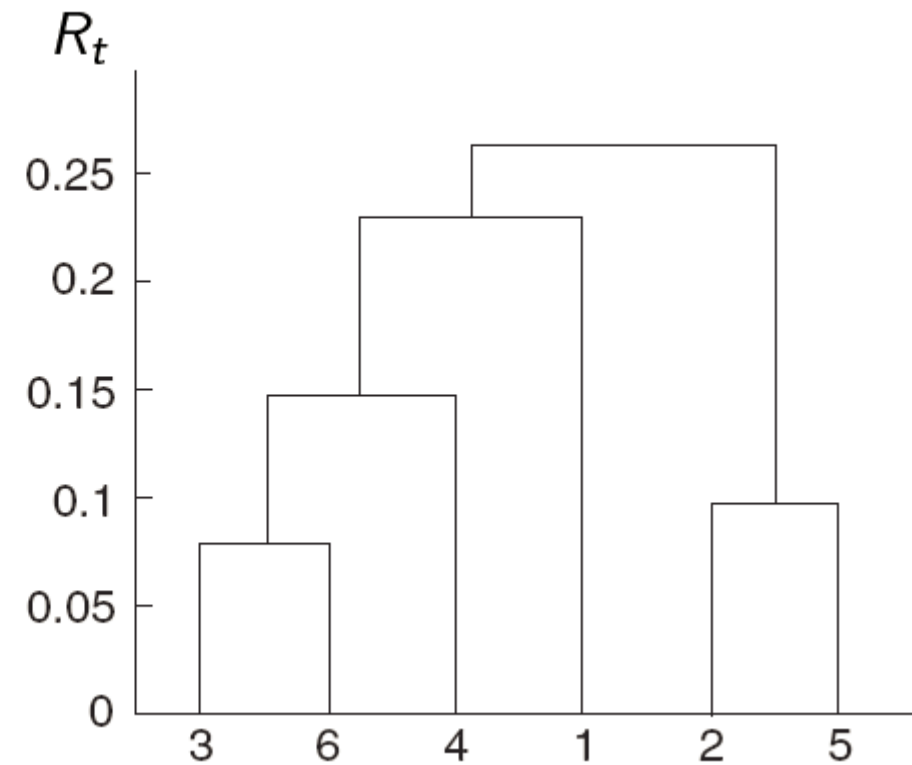
# Иерархические алгоритмы кластеризации

## Расстояние Уорда

Диаграмма вложения



Дендрограмма





# Основные свойства иерархической кластеризации

**Монотонность:** дендрограмма не имеет самопересечений, при каждом слиянии расстояние между объединяемыми кластерами увеличивается:  $R_2 \leq R_3 \leq R_4 \dots$

$R^{\text{ц}}$  — не монотонна,  $R^{\text{б}}$   $R^{\text{д}}$   $R^{\text{у}}$  — монотонны

**Сжимаемость и растягиваемость:**

$R_t \leq \rho(\mu_u, \mu_v), \forall t$  — сжимающее расстояние

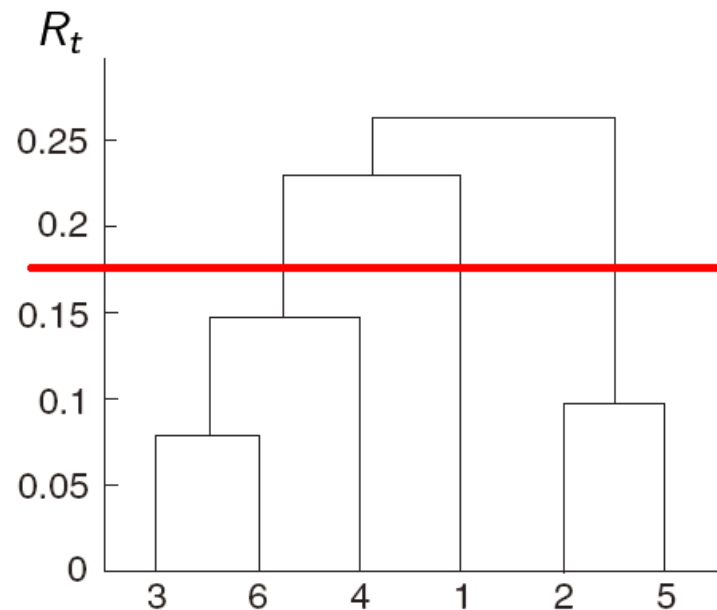
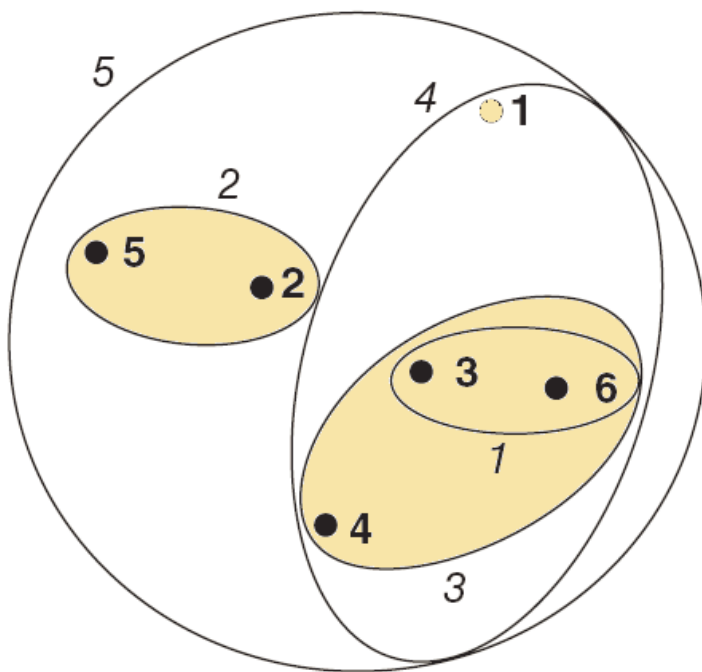
$R_t \geq \rho(\mu_u, \mu_v), \forall t$  — растягивающее расстояние

Свойство растяжения желательно, так как оно способствует более четкому отделению кластеров

$R^{\text{б}}$  — сильно сжимающее,  $R^{\text{д}}$   $R^{\text{у}}$  — растягивающие,  $R^{\text{ц}}$  — сохраняет метрику пространства

# Выводы и рекомендации

- Рекомендуется пользоваться расстоянием Уорда.
- Обычно строят несколько вариантов и выбирают лучший визуально по дендрограмме.
- Определять число кластеров рекомендуется по максимальной высоте участка  $|R_{t+1} - R_t|$  на дендрограмме.



# ЕМ-алгоритм. Предпосылки

## Гипотеза о вероятностной природе данных:

Обучающая выборка  $X$  случайна и независима, состоит из смеси распределений

$$p(x) = \sum_{y \in Y} w_y p_y(x) \quad \sum_{y \in Y} w_y = 1$$

$p_y(x)$  – плотность,  $w_y$  - априорная вероятность кластера  $y$

## Гипотеза о пространстве объектов и форме кластеров:

Кластеры  $n$ -мерные, гауссовские

$$p_y(x) = (2\pi)^{-n/2} (\sigma_{y1} \cdots \sigma_{yn})^{-1} \exp(-1/2 \rho_y^2(x, \mu_y))$$

$\mu_y = (\mu_{y1}, \dots, \mu_{yn})$  – центр кластера  $y$

$\Sigma_y = \text{diag}(\sigma_{y1}^2, \dots, \sigma_{yn}^2)$  – диагональная матрица ковариаций

$$\rho_y^2(x, x') = \sum_{j=1}^n \sigma_{yj}^{-2} |f_j(x) - f_j(x')|^2$$

# ЕМ-алгоритм

1: начальное приближение  $w_y, \mu_y, \Sigma_y$  для всех  $y \in Y$ ;

2: **повторять**

3: Е-шаг (expectation):

$$g_{iy} := P(y|x_i) \equiv \frac{w_y p_y(x_i)}{\sum_{z \in Y} w_z p_z(x_i)}, \quad y \in Y, \quad i = 1, \dots, \ell;$$

4: М-шаг (maximization):

$$w_y := \frac{1}{\ell} \sum_{i=1}^{\ell} g_{iy}, \quad y \in Y;$$

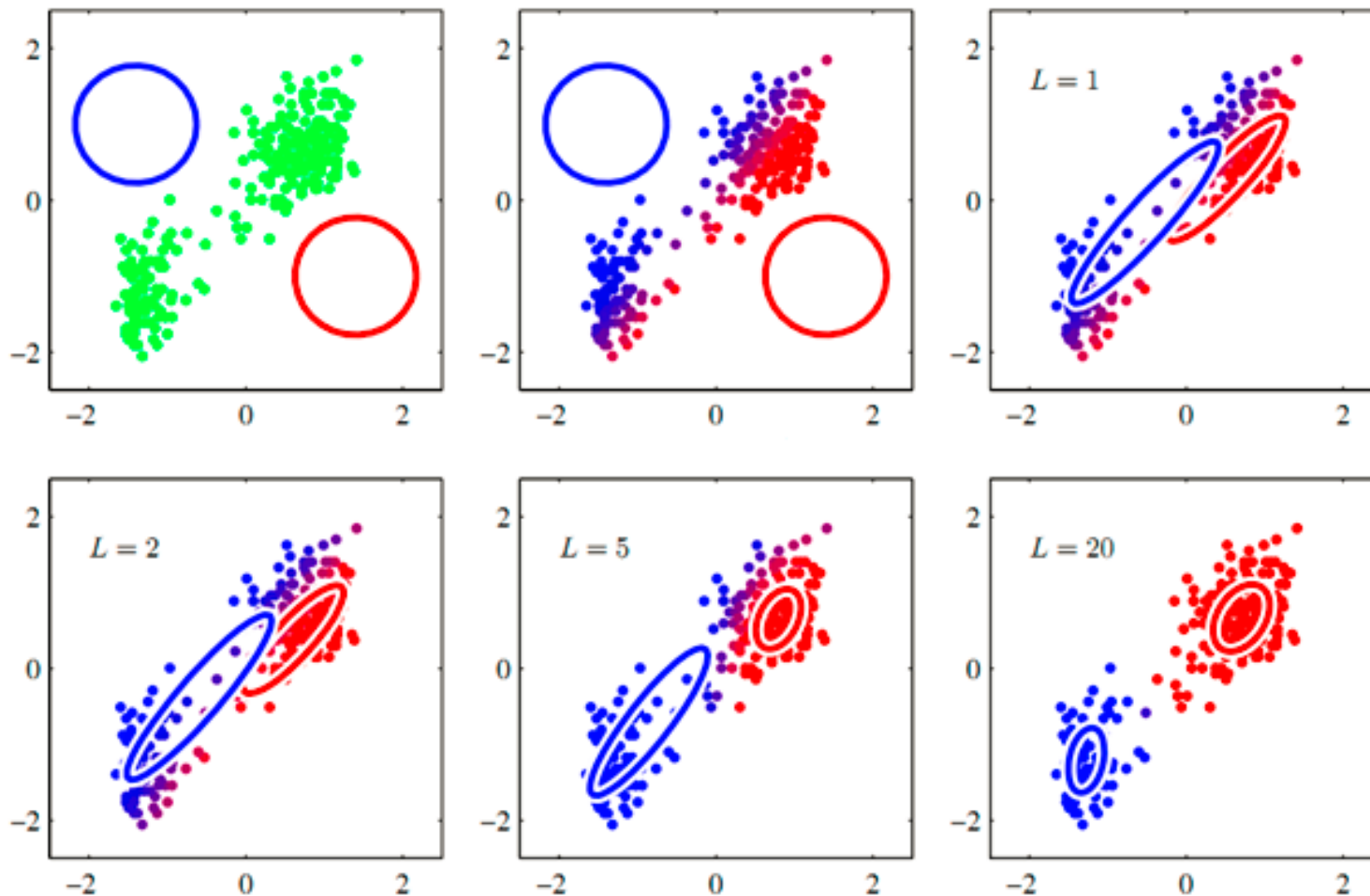
$$\mu_{yj} := \frac{1}{\ell w_y} \sum_{i=1}^{\ell} g_{iy} f_j(x_i), \quad y \in Y, \quad j = 1, \dots, n;$$

$$\sigma_{yj}^2 := \frac{1}{\ell w_y} \sum_{i=1}^{\ell} g_{iy} (f_j(x_i) - \mu_{yj})^2, \quad y \in Y, \quad j = 1, \dots, n;$$

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

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

# ЕМ-алгоритм



# Метод $k$ -средних ( $k$ -means)

Упрощенный аналог ЕМ-алгоритма:

Жесткая кластеризация вместо мягкой

1. Начальное приближение центроидов  $\mu_y$ ,  $y \in Y$

**2. Повторять:**

3. Аналог Е-шага:

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. Аналог М-шага:

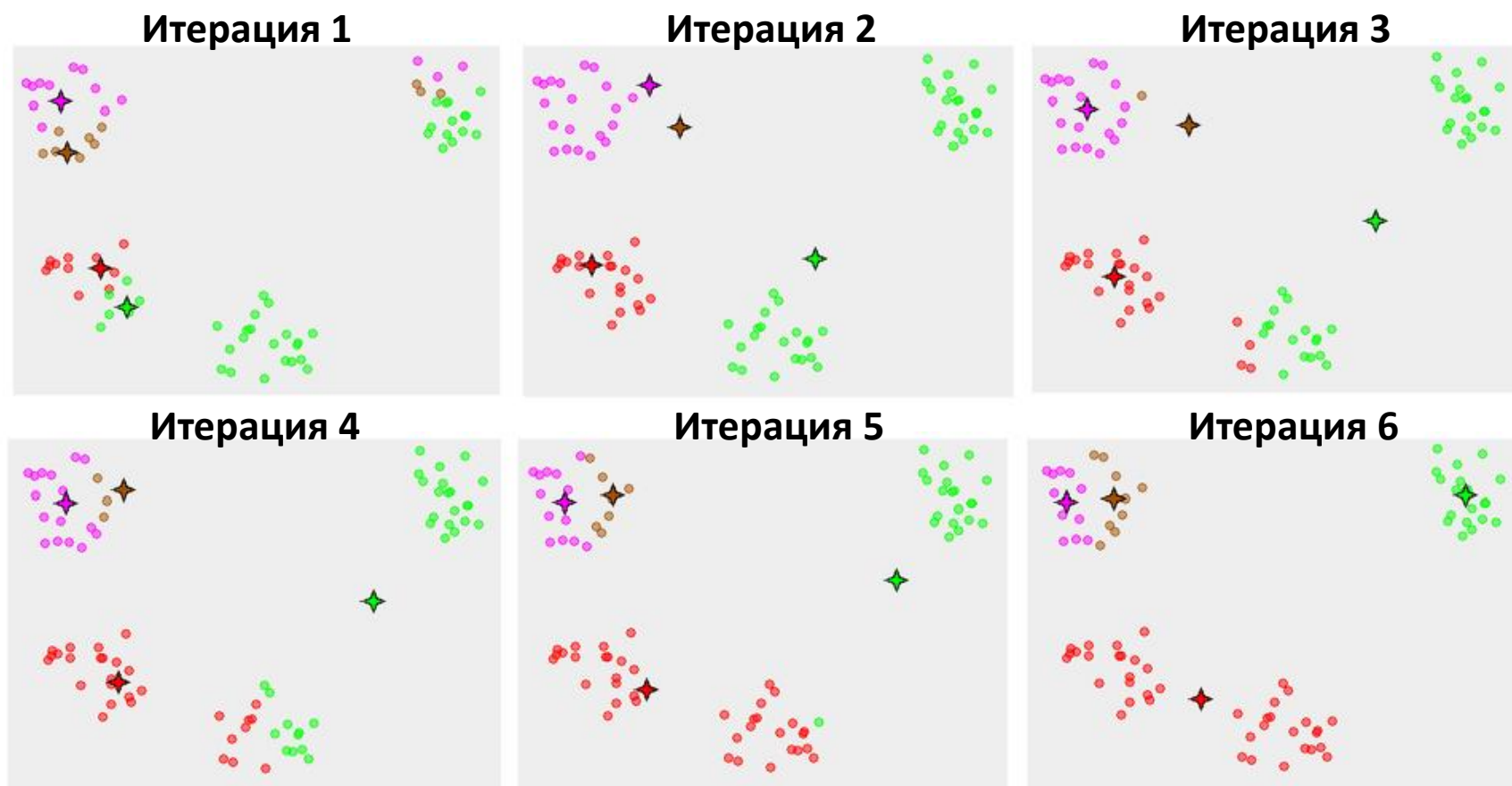
вычислить новые положения центров:

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

5. **Пока**  $y_i$  не перестанут изменяться

# Недостатки метода (k-means)

- Не гарантируется достижение глобального минимума суммарного квадратичного отклонения  $V$ , а только одного из локальных минимумов.
- Результат зависит от выбора исходных центров кластеров, их оптимальный выбор неизвестен.
- Рекомендуется повторная прогонка алгоритма для избежания ситуации «плохой» кластеризации.
- Число кластеров надо знать заранее.



# Семейство алгоритмов FOREL (ФОРмальный Элемент)

Алгоритм предложен Загоруйко Н. Г. и Ёлкиной В. Н. в 1967 году.

Задается параметр  $R$  – радиус поиска локальных сгущений.

На каждом шаге мы

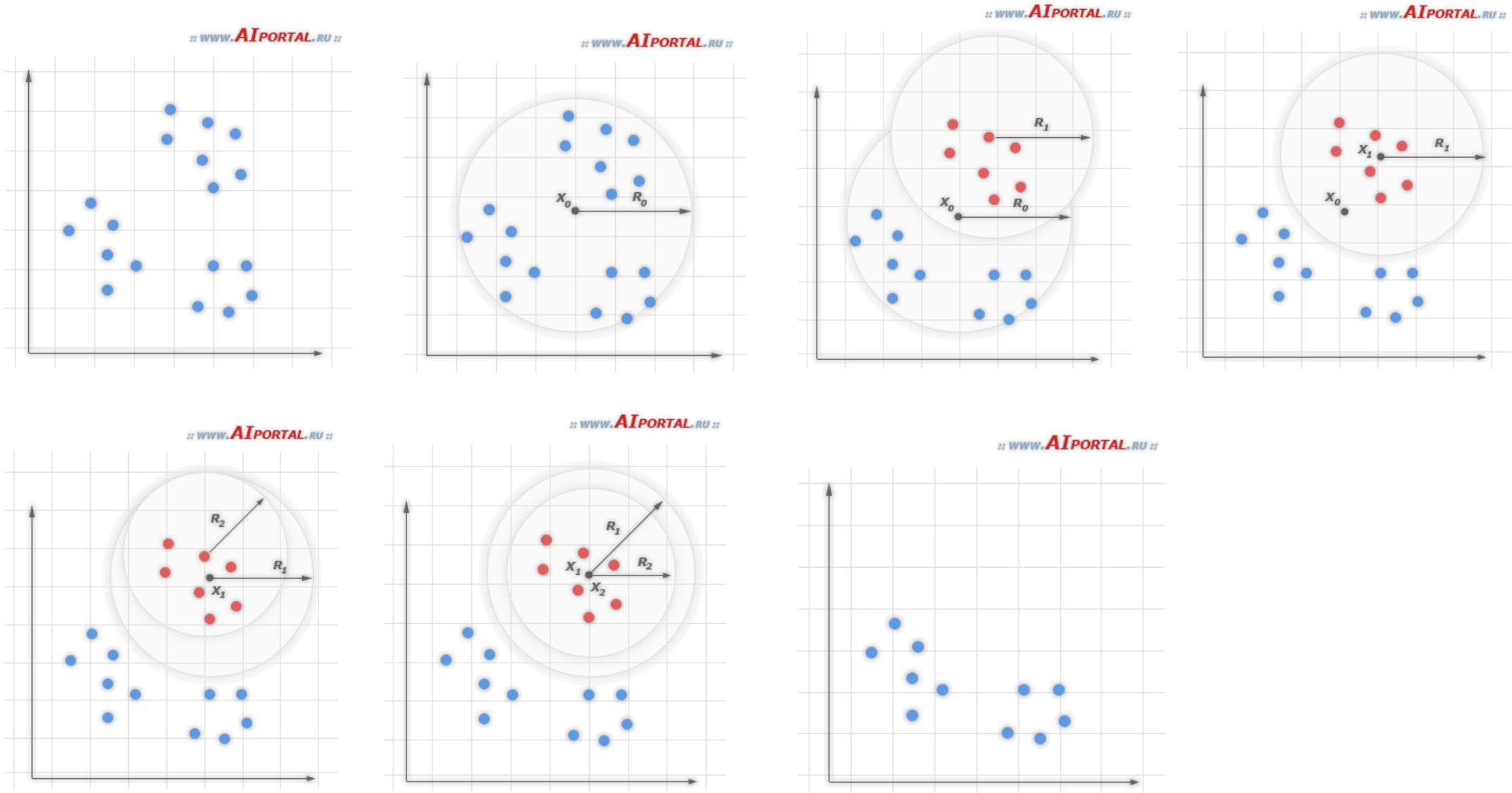
1. случайным образом выбираем объект из выборки,
2. раздуваем вокруг него сферу радиуса  $R$ ,
3. внутри этой сферы выбираем центр тяжести и делаем его центром новой сферы.

Таким образом, мы на каждом шаге двигаем сферу в сторону локального сгущения объектов выборки, т.е. стараемся захватить как можно больше объектов выборки сферой фиксированного радиуса.

4. После того как центр сферы стабилизируется, все объекты внутри сферы с этим центром мы помечаем как кластеризованные и выкидываем их из выборки. Этот процесс мы повторяем до тех пор, пока вся выборка не будет кластеризована.



# Визуализация алгоритма семейства FOREL



# Алгоритм FOREL (ФОРмальный Элемент)

## Преимущества:

- Точность минимизации функционала качества (при удачном подборе параметра  $R$ )
- Наглядность визуализации кластеризации
- Сходимость алгоритма
- Возможность подсчета промежуточных функционалов качества, например, длины цепочки локальных сгущений

## Недостатки:

- Относительно низкая производительность
- Плохая применимость алгоритма при плохой делимости выборки на кластеры
- Неустойчивость алгоритма (зависимость от выбора начального объекта)
- Произвольное по количеству разбиение на кластеры
- Необходимость априорных знаний о ширине (диаметре) кластеров