

# Математические основы информационной безопасности

Груздев Дмитрий Николаевич

Обучение без учителя

# Кластеризация

$x_1, x_2, \dots, x_m$  – обучающая выборка

$\rho(x_i, x_j)$  – расстояние между объектами

Построить:

- $Y$  – множество кластеров
- $A: X \rightarrow Y$  – алгоритм кластеризации объектов

Требования к кластерам:

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

# Кластеризация

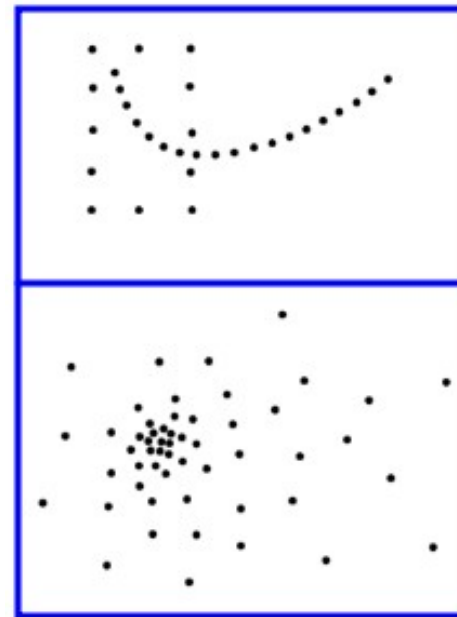
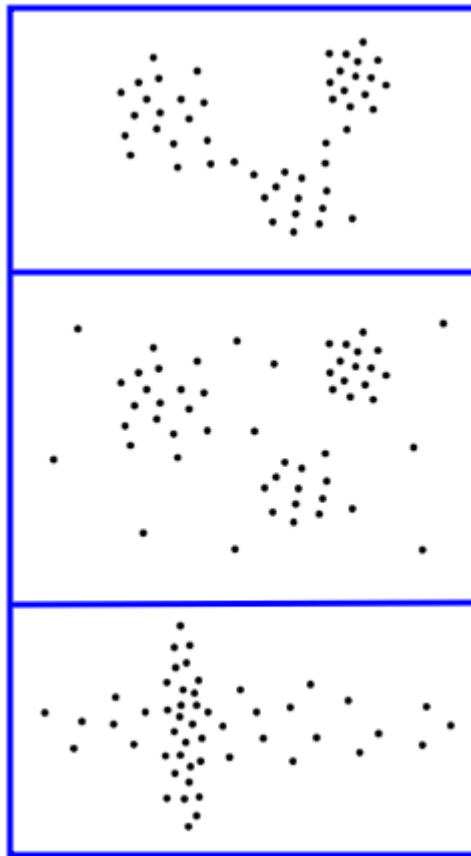
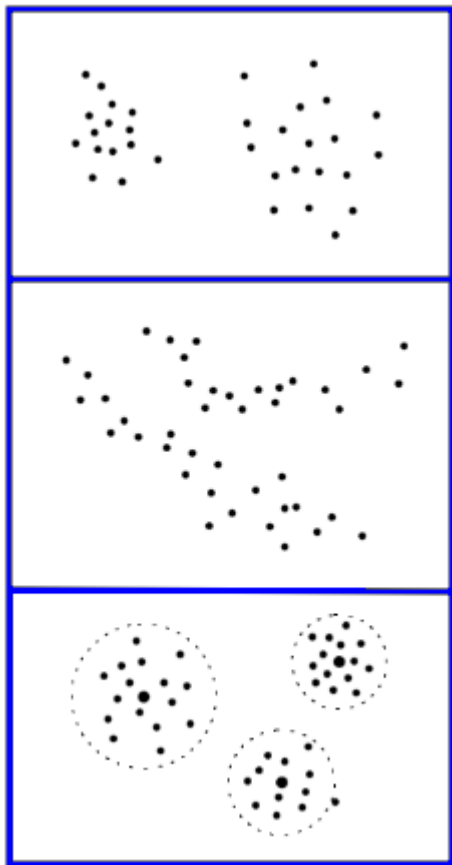
Трудности при формализации задачи

- существует много критериев качества кластеризации
- число кластеров обычно неизвестно заранее

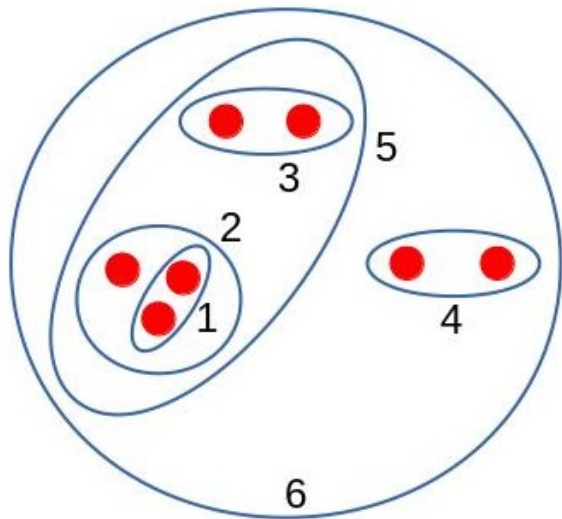
Цели кластеризации:

- упростить обработку данных
- сократить объем хранимых данных
- выделить нетипичные признаки
- построить иерархию множества объектов

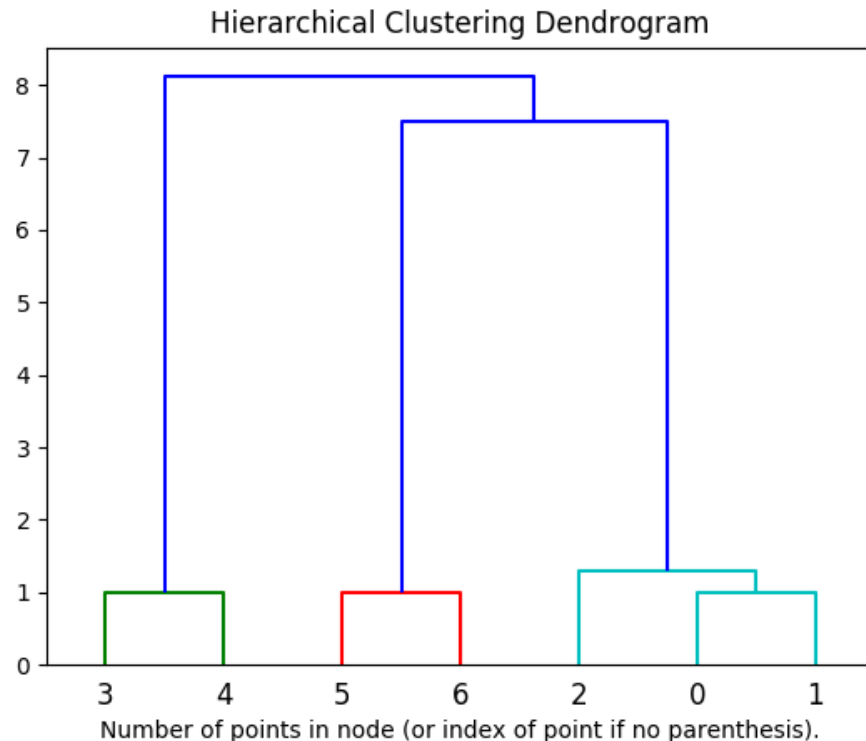
# Виды кластеров



# Иерархическая кластеризация



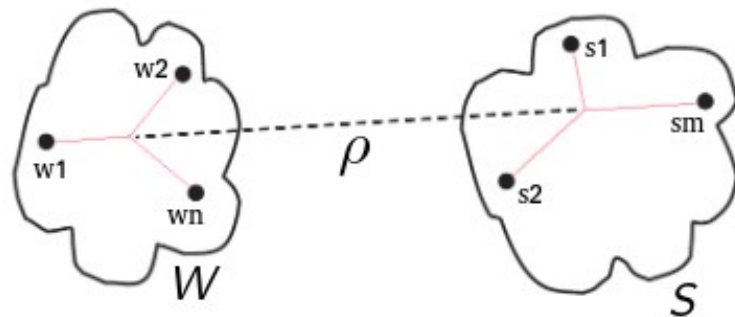
Последовательно объединяем  
два ближайших друг к другу  
кластера в новый.



# Расстояние между кластерами

расстояние между центрами

$$\rho(W, S) = \frac{|S||W|}{|S|+|W|} \rho^2 \left( \sum_{w \in W} \frac{w}{|W|}, \sum_{s \in S} \frac{s}{|S|} \right) - \text{расстояние Уорда}$$

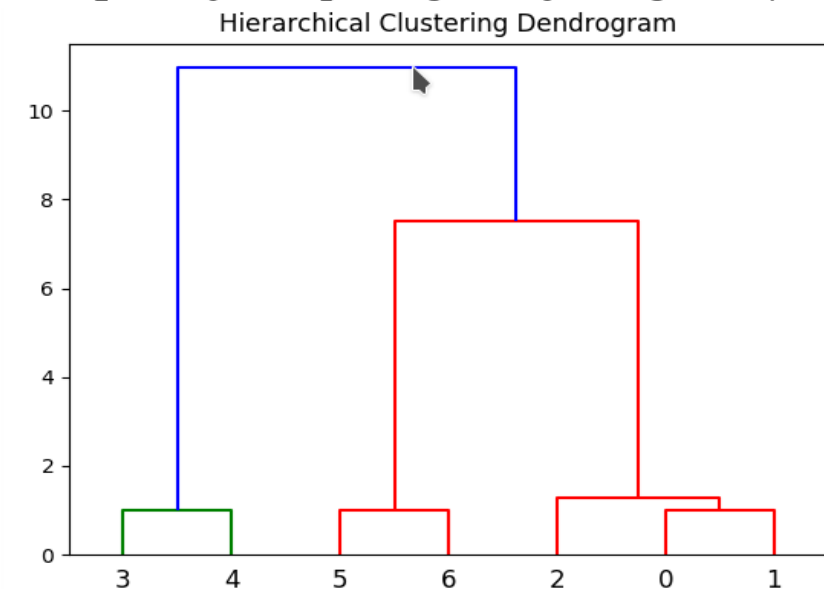
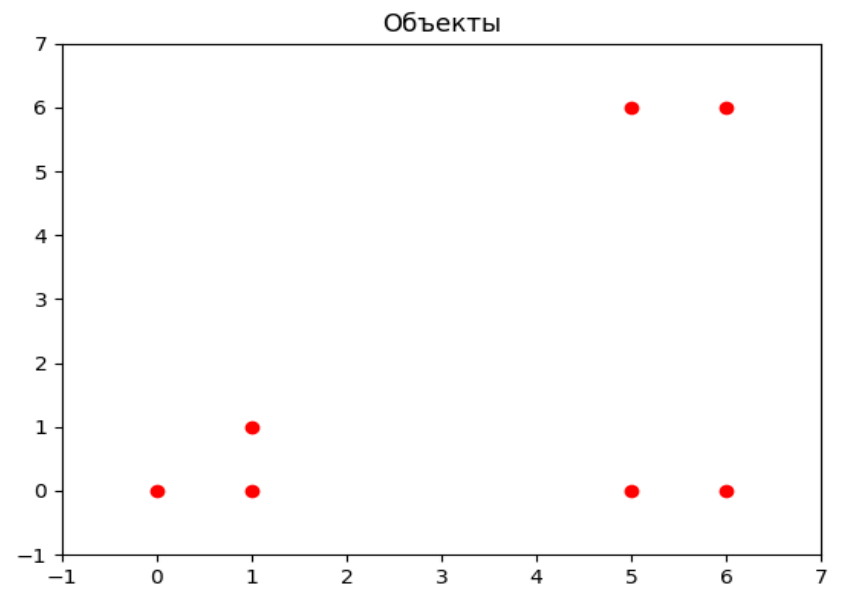
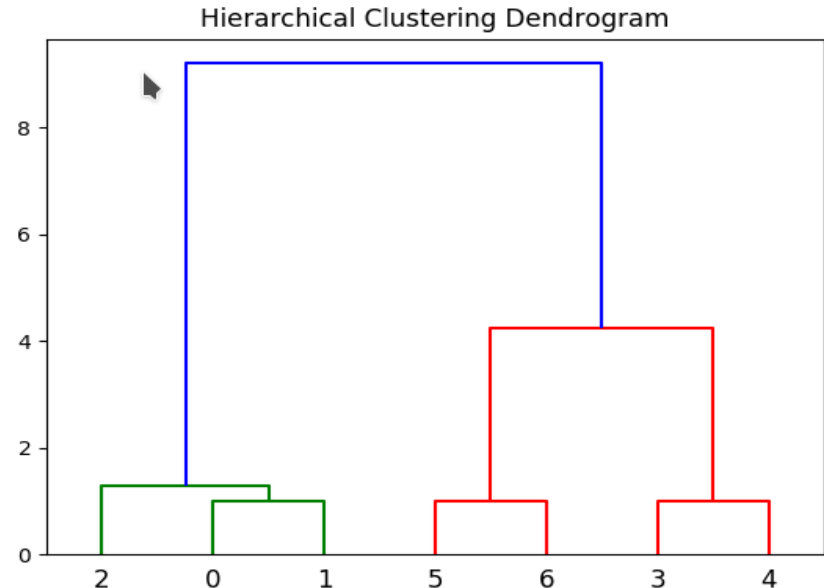
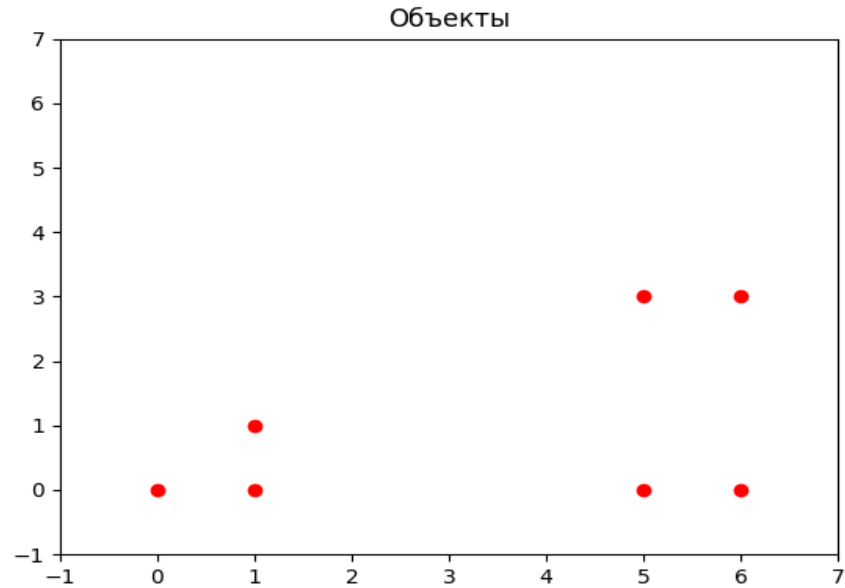


Индукционное вычисление расстояния Уорда:

$$W = A \cup B$$

$$\rho(A \cup B, S) = \alpha * \rho(A, S) + \beta * \rho(B, S) + \gamma * \rho(A, B), \text{ где}$$

$$\alpha = (|A|+|S|)/(|S|+|W|), \beta = (|B|+|S|)/(|S|+|W|), \gamma = -|S|/(|S|+|W|)$$





# Метод k-средних

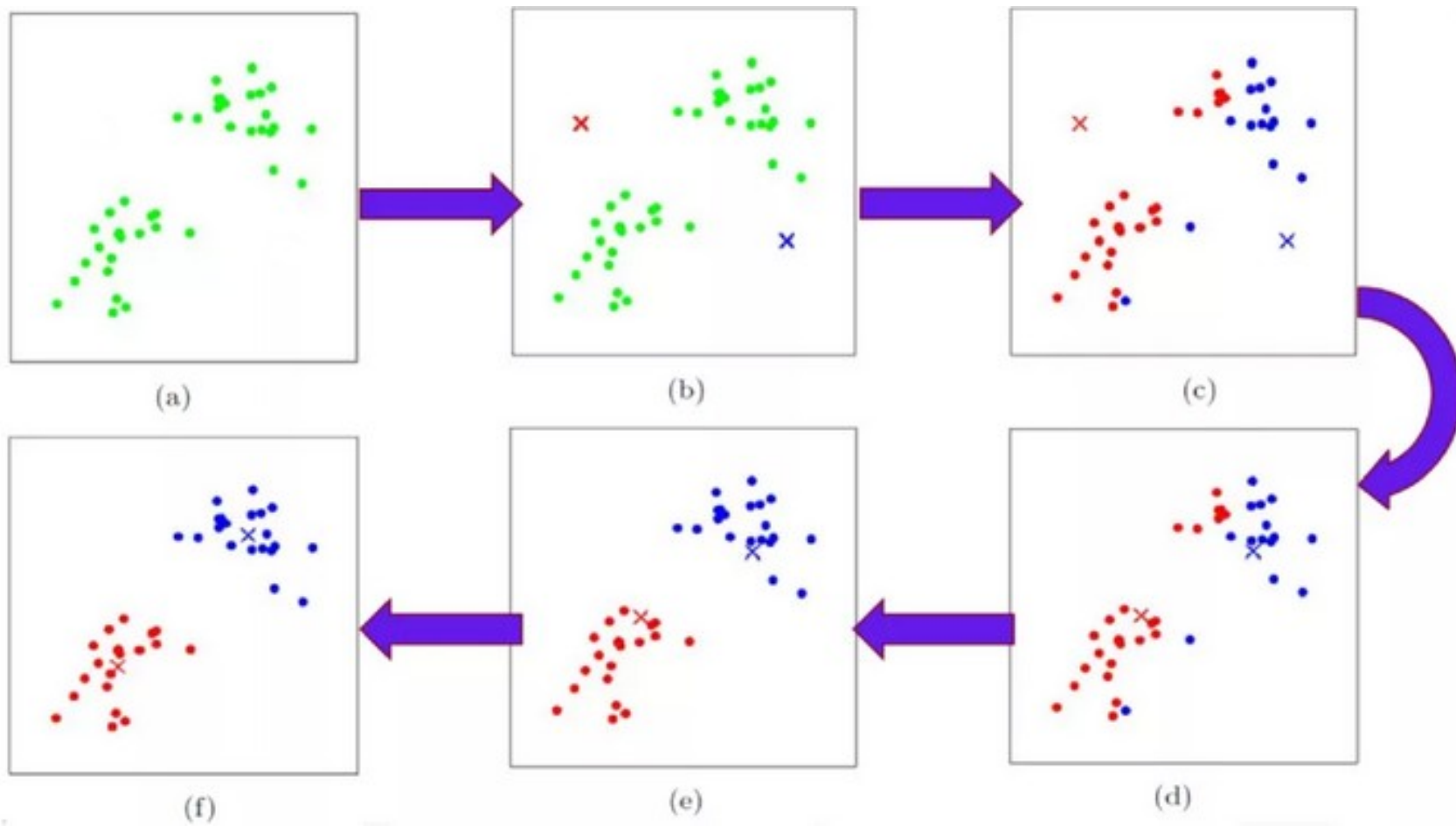
$x_1, x_2, \dots, x_m$  – обучающая выборка

$$x_i = (x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(n)})$$

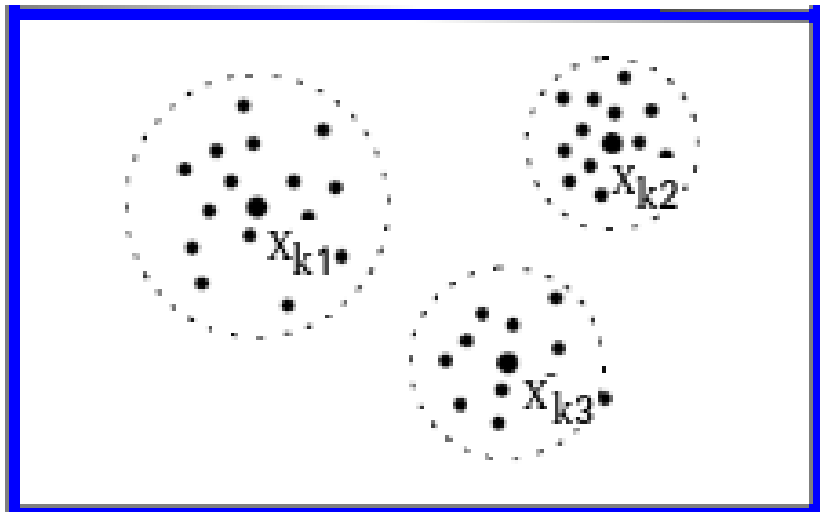
Алгоритм:

1. Задать количество кластеров и их центры  $x_{k1}, \dots, x_{km}$ .
2. Отнести каждый объект к ближайшему центру.
3. Перенести координаты центров кластеров в центры тяжести соответствующих групп объектов.
4. Повторять 2 и 3, пока происходят переходы объектов между кластерами.

# Метод k-средних



# Оптимальное число кластеров

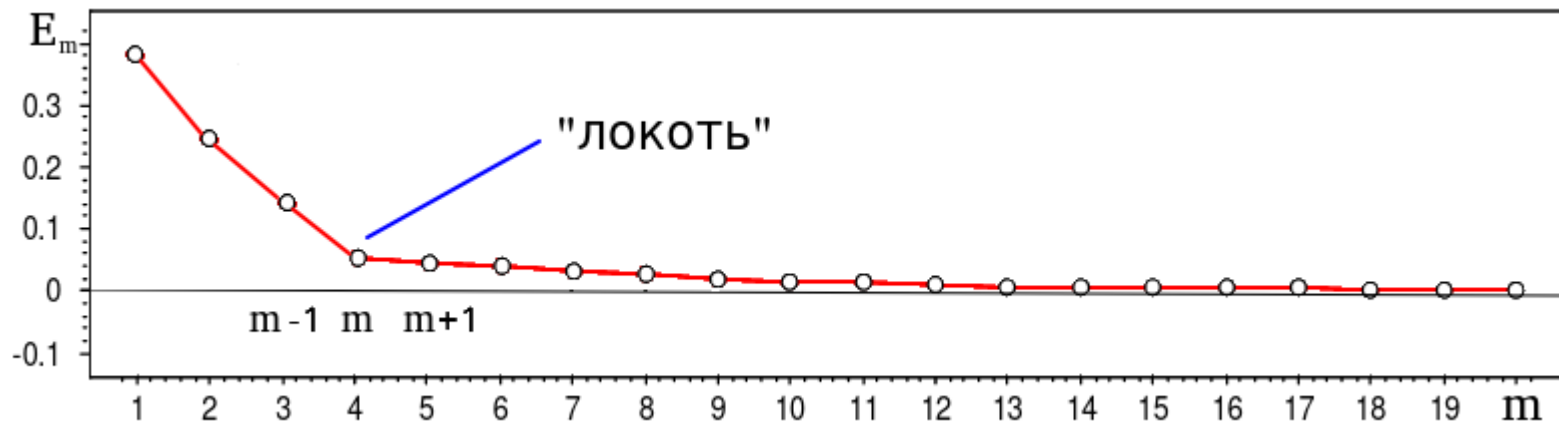


Ошибка кластеризации:

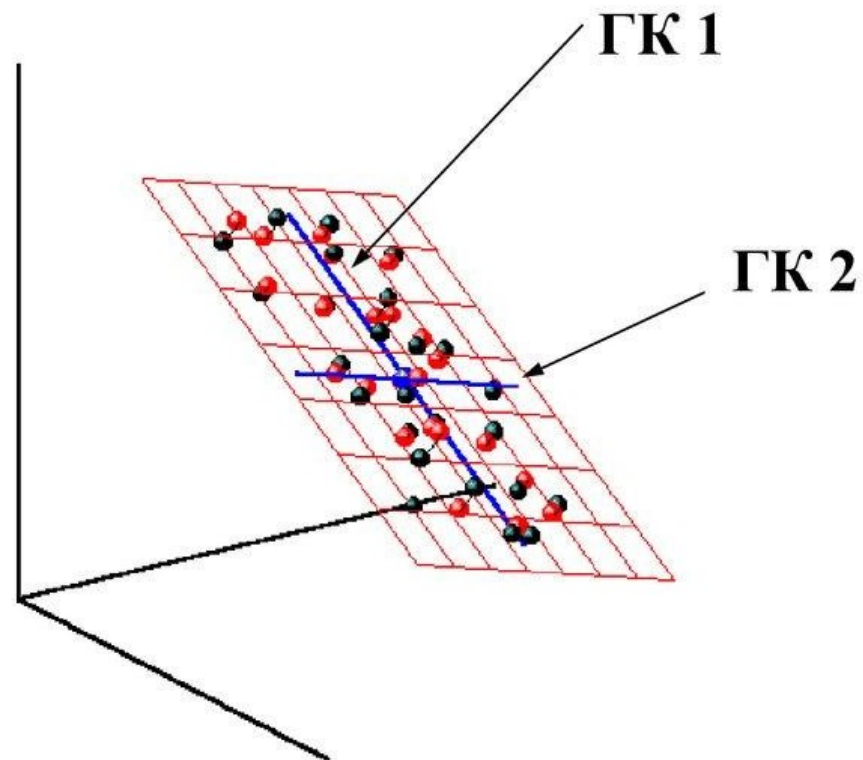
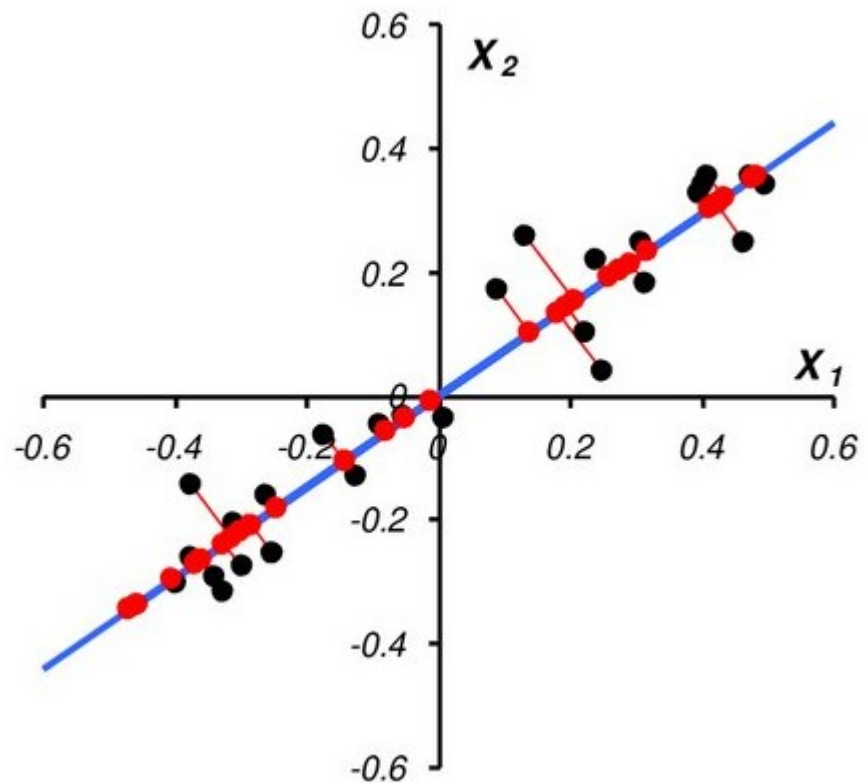
$$E_m = \sum_{1 \leq i \leq m} \sum_{1 \leq j \leq |K_i|} \|x_{K_i}^{(j)} - x_{K_i}\|^2$$

Выбор  $m$  – числа кластеров:

$$E_m - E_{m+1} \gg E_{m-1} - E_m$$



# Метод главных компонент



# Метод главных компонент

$x_1, x_2, \dots, x_m$  – обучающая выборка

$$x_i = (x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(n)})$$

$$x_i = (g_i^{(1)}, g_i^{(2)}, \dots, g_i^{(k)}), \quad k \leq n$$

$$x_i^{(j)} \approx \sum g_i^{(l)} u_l = h_i^{(j)}$$

$$\sum_{1 \leq i \leq m} \sum_{1 \leq j \leq n} (h_i^{(j)} - x_i^{(j)})^2 \rightarrow \min_{(u, g)}$$

# Преобразование признаков

$$X = \begin{pmatrix} x_1^{(1)} & \dots & x_1^{(n)} \\ \vdots & & \vdots \\ x_m^{(1)} & \dots & x_m^{(n)} \end{pmatrix} \quad G = \begin{pmatrix} g_1^{(1)} & \dots & g_1^{(k)} \\ \vdots & & \vdots \\ g_m^{(1)} & \dots & g_m^{(k)} \end{pmatrix} \quad U = \begin{pmatrix} u_1^{(1)} & \dots & u_1^{(k)} \\ \vdots & & \vdots \\ u_n^{(1)} & \dots & u_n^{(k)} \end{pmatrix}$$

$$X \approx GU^T; \|GU^T - X\|^2 \rightarrow \min_{G,U}$$

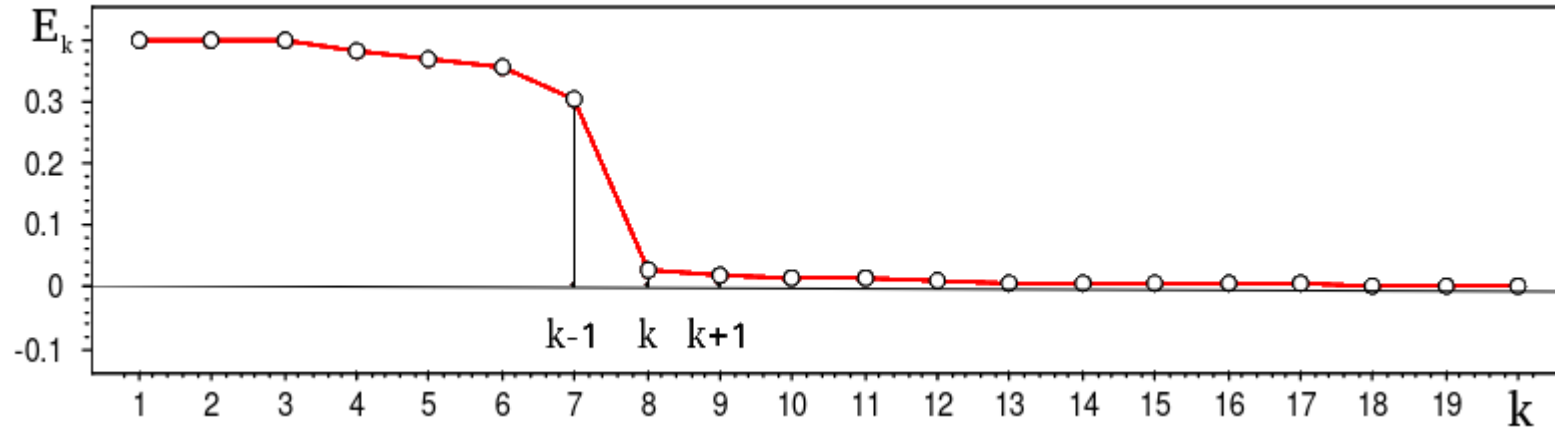
Решение:

$$G = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_k), \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k \geq \dots \geq \lambda_n$$

- $UU^T = I_k$  – ортонормирована, и  $G \approx XU$ ,  $X \approx GU^T$

$$\|GU^T - X\|^2 = \lambda_{k+1} + \lambda_{k+2} + \dots + \lambda_n$$

# Эффективная размерность



$$E_k = \|GU^T - F\|^2 / \|F\|^2 = (\lambda_{k+1} + \dots + \lambda_n) / (\lambda_1 + \dots + \lambda_n)$$

Установить допустимую погрешность  $E_k \leq \varepsilon$  или

если  $E_{k-1} \gg E_k$ , то стоит выбрать размерность  $k$ .

# Диагностика аномалий

$x_1, x_2, \dots, x_m$  — обучающая выборка

$$x_i = (x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(n)})$$

Особенности:

- большое количество нормальных результатов,
- малое число (отсутствие) результатов с аномалиями,
- типы аномалий трудно систематизировать,
- возможны неизвестные аномалии.



# Случайные величины

## Теория вероятности

$X, Y$  – случайные величины

$M_X$  – математическое ожидание  $X$

$D_X$  – дисперсия  $X$

$\sigma_X = D_X^{1/2}$  – стандартное отклонение  $X$

$\text{cov}_{XY} = M((X - M_X)(Y - M_Y))$  – ковариация  $X$  и  $Y$

$$r_{XY} = \frac{\text{cov}_{XY}}{\sigma_X \sigma_Y} -$$

коэффициент линейной корреляции  
(коэффициент корреляции Пирсона)

## Математическая статистика

$X = \{x_1, \dots, x_N\}, Y = \{y_1, \dots, y_N\}$  – выборки

$\bar{X} = \frac{1}{n} \sum x_i$  - выборочное среднее

$S_X^2 = \frac{1}{n} \sum (x_i - \bar{X})^2$  - выборочная дисперсия

$\text{cov}_{XY} = \frac{1}{n} \sum (x_i - \bar{X})(y_i - \bar{Y})$  - выборочная ковариация

$$r_{XY} = \frac{\text{cov}_{XY}}{S_X S_Y} = \frac{\sum (x_i - \bar{X})(y_i - \bar{Y})}{\sqrt{\sum (x_i - \bar{X})^2 \sum (y_i - \bar{Y})^2}}$$

коэффициент корреляции Пирсона

# Вероятность появления объекта

$$x(i) \sim N(\mu_i, \sigma_i^2)$$

$$\mu_i = \sum_{1 \leq j \leq m} x_j(i) / m - \text{выборочное среднее}$$

$$\sigma_i^2 = \sum_{1 \leq j \leq m} (x_j(i) - \mu_i)^2 / m - \text{выборочная дисперсия}$$

$$p(x^{*(i)}) = p(x^{*(i)}; \mu_i, \sigma_i^2) = f(x^{*(i)}; \mu_i, \sigma_i^2) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp \left[ -\frac{(x^{*(i)} - \mu_i)^2}{2\sigma_i^2} \right]$$

$$x^* = (x^{*(1)}, \dots, x^{*(n)})$$

Если  $x^{(1)}, x^{(2)}, \dots, x^{(n)}$  — независимы, то

$$p(x^*) = p(x^{*(1)}; \mu_1, \sigma_1^2) * p(x^{*(2)}; \mu_2, \sigma_2^2) * \dots * p(x^{*(n)}; \mu_n, \sigma_n^2)$$

# Диагностика аномалий

Алгоритм диагностики:

- выбрать  $\varepsilon > 0$  пороговое:
- если  $p(x^*) < \varepsilon$ , то объект относится к аномальным.

Настройка  $\varepsilon$ :

- алгоритм должен определять известные аномалии,
- алгоритм не должен относить к аномалиям нормальные объекты
- если не удастся подобрать пороговый параметр, то нужны другие признаки объектов

# scikit-learn

<https://sesc-infosec.github.io/>