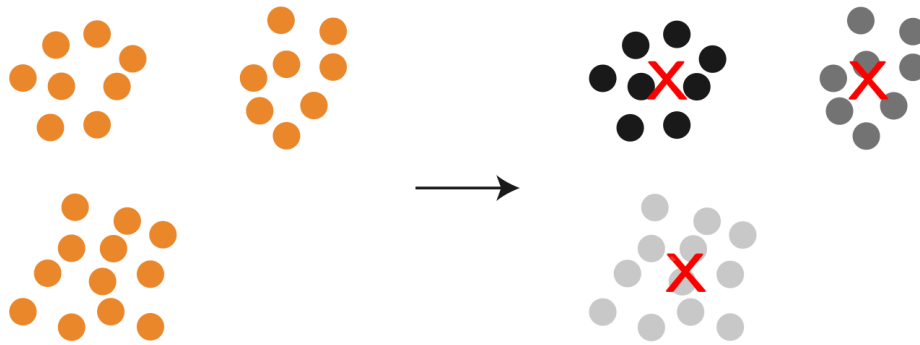


Задача размещения



Дано:

$S \subset \mathbb{R}^d$ ← *набор точек*

← *размер n*

$k \in \mathbb{N}$ ← *число объектов*

Вспомог: $T \subset \mathbb{R}^d$ $|T| = k$ ← *множество объектов*

$$\min_T \text{cost}(S, T) = \sum_{x \in S} \min_{z \in T} \|x - z\|_2^2$$

← *для каждой точки*

← *найти*

← *квадрат расстояния*

Обозначения:

$$\text{cost}(S, T) = \sum_{x \in S} \min_{z \in T} \|x - z\|_2^2$$

$$= \sum_{z \in T} \sum_{x \in C_z} \|x - z\|_2^2$$

группа точек S

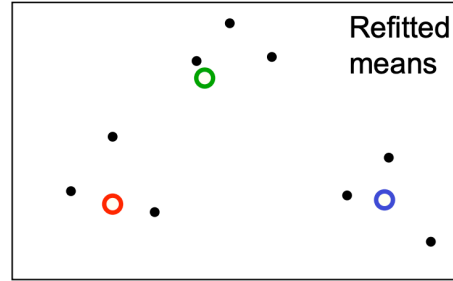
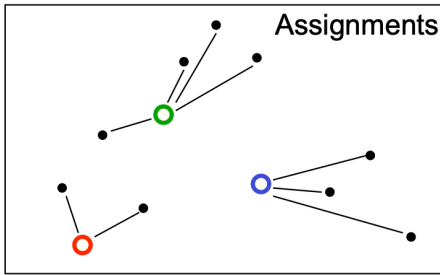
↓

← *объект c которым, сумм. к z*

$$\text{cost}(C, T) = \sum_{x \in C} \min_{z \in T} \|z - x\|_2^2$$

$$\text{cost}(C, \underset{\substack{\uparrow \\ \text{given} \\ \text{user}}}{z}) = \sum_{x \in C} \|z - x\|_2^2$$

K-means
user:



Lemma $\forall C \subset \mathbb{R}^d \quad \forall z \in \mathbb{R}^d$

$$\text{cost}(C, z) = \text{cost}(C, \text{mean}(C)) + |C| \cdot \|z - \text{mean}(C)\|_2^2$$

$z = \text{mean}(C)$ gives $\min \text{cost}(C, z)$

Der. by $\text{cost}(C, z) = \left(\frac{1}{|C|} \sum_{x \in C} \|x - z\|_2^2 \right) \cdot |C|$

$\underbrace{\qquad\qquad\qquad}_{\mathbb{E}_{x \sim U_C} \|x - z\|_2^2}$

$$\cancel{|C|} \cdot \mathbb{E}_x \|x - z\|_2^2 = \mathbb{E} \|x - \mathbb{E}x\|_2^2 \cdot \cancel{|C|} + \cancel{|C|} \cdot \|z - \mathbb{E}x\|_2^2$$

↙

$$\mathbb{E} \|x\|_2^2 + \|\mathbb{E}x\|_2^2 - \underbrace{2 \mathbb{E}(x \cdot \mathbb{E}x)}_{2\|\mathbb{E}x\|_2^2}$$

$$+ \|z\|_2^2 + \|EX\|_2^2 - 2z \cdot EX$$

$$= E \|X\|_2^2 + \|z\|_2^2 - 2z \cdot EX = E \|X - z\|_2^2$$

Algorithm K-means:

Input: C_1^0, \dots, C_k^0 $t=0$

while the algorithm terminates

$$1) \text{ for } i: z_i^{t+1} = \text{mean}(C_i^t)$$

$$2) \text{ for } i: C_i^{t+1} = \{x \in S \mid i = \arg \min_j \|z_j^{t+1} - x\|_2^2\}$$

Theorem K-means minimizes the cost

$$\text{Cost} = \sum_{z \in T} \sum_{x \in C_z} \|x - z\|_2^2 = \sum_{i=1}^k \sum_{x \in C_i} \|x - z_i\|_2^2$$

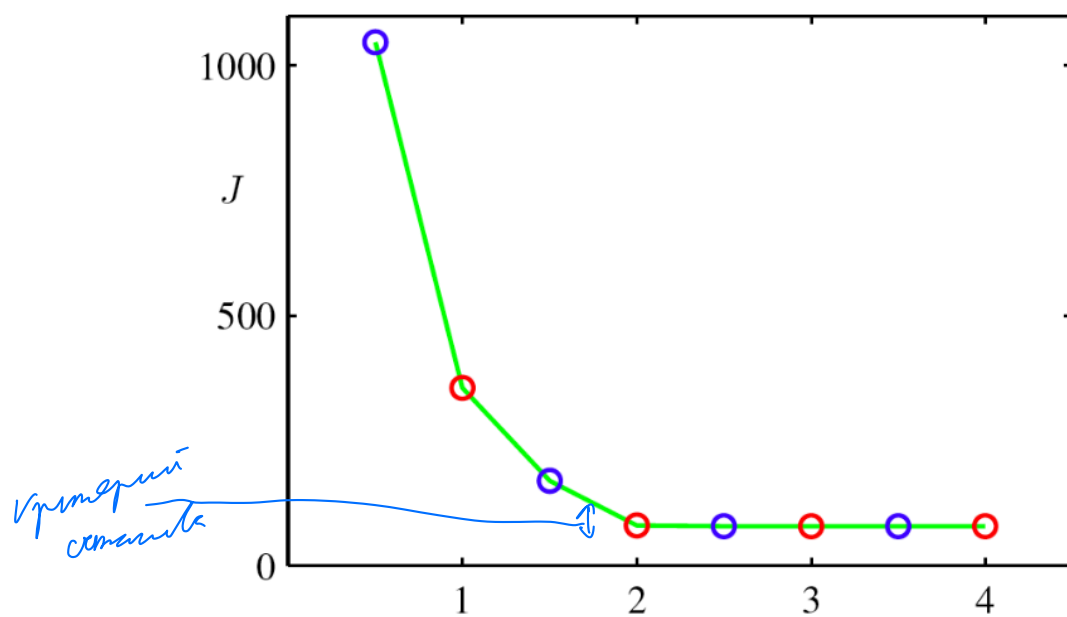
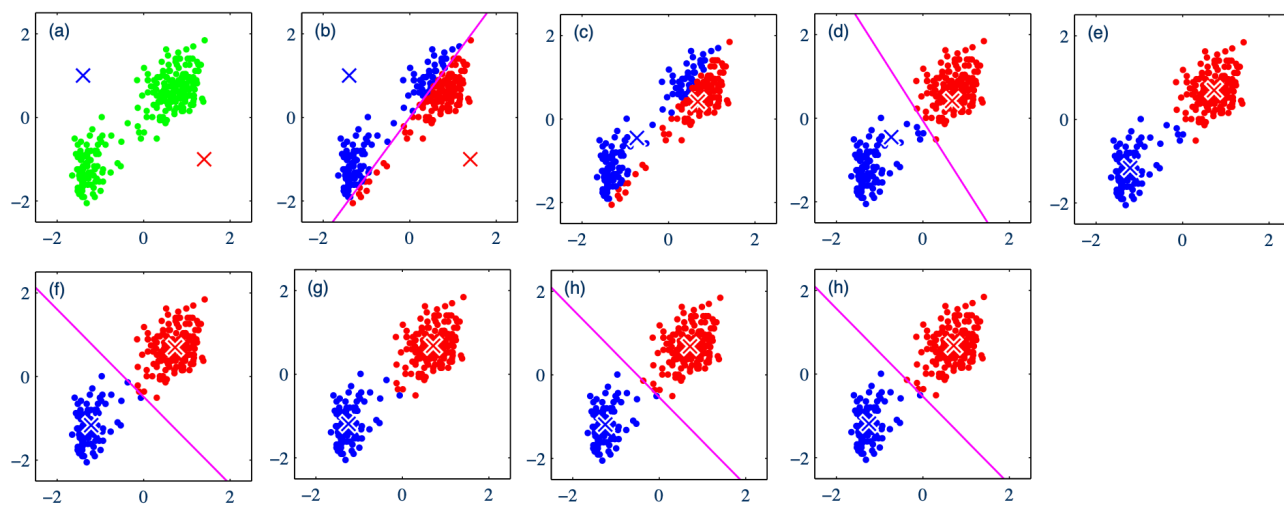
$$1) \text{ cost}(C_i^t, z_i^{t+1}) \leq \text{cost}(C_i^t, z_i^t)$$

$$\sum_i \leq \text{no lemma 1}$$

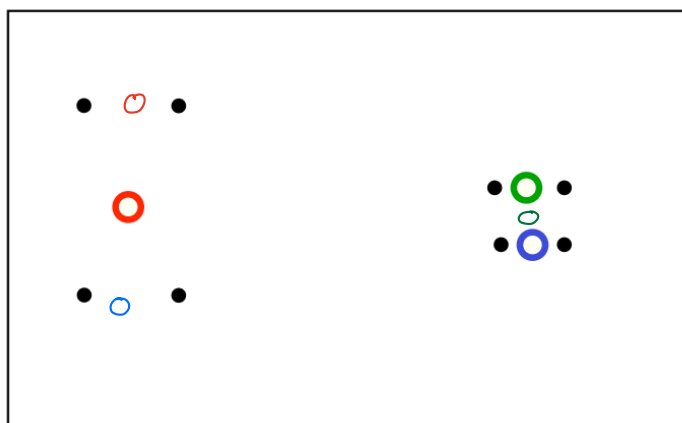
2) if you know

$$\text{cost}(C_i^{t+1}, z_i^{t+1}) \leq \text{cost}(C_i^t, z_i^{t+1})$$

$$\text{cost}(C_i^{t+1}, z_i^{t+1}) \leq \text{cost}(C_i^t, z_i^t)$$



Триггер поиска результатов:



использованы

A var minimierungslösung K-means?

K-means ++

Gegeben $x \in S$ - geo repräs von z_1

for $i = 1 \dots k-1$

Gegeben $x \in S$ berechne

$$P\{x\} \sim \text{cost}(x, T_i) = \min_{z \in T_i} \|x - z\|_2^2$$

$\{z_1, \dots, z_i\}$

$$P\{x\} = \frac{\min_{z \in T_i} \|x - z\|_2^2}{\sum_{x' \in S} \min_{z \in T_i} \|x' - z\|_2^2}$$

Bezug 1 moving - b modern approach von b. V. K-means

Lemma 2

$\forall C \in \mathbb{R}^d, \quad z \sim U(C), \text{ mean}$

$$E[\text{cost}(C, z)] = 2 \text{cost}(C, \text{mean}(C))$$

Beweis:

$$E[\text{cost}(C, z)] = \sum_{z \in C} \frac{1}{|C|} \underbrace{\text{cost}(C, z)}_{\text{lemma 1}}$$

$$= \frac{1}{|C|} \sum_{z \in C} (\text{cost}(C, \text{mean}(C)) + |C| \cdot \|z' - \text{mean}(C)\|_2^2)$$

$$\begin{aligned}
&= \text{cost}(C, \text{mean}(C)) \\
&\quad + \sum_{z' \in C} \|z' - \text{mean}(C)\|_2^2 \\
&= 2 \text{cost}(C, \text{mean}(C))
\end{aligned}$$

Lemma 3

T_i - множество точек, а z_{i+1}^* - значение центра

$z_{i+1} \in C_{i+1}^*$ - значение центра

$$\mathbb{E} [\text{cost}(C_{i+1}^*, T_i \cup \{z_{i+1}^*\})] \leq 8 \text{cost}(C_{i+1}^*, \text{mean}(C_{i+1}^*))$$

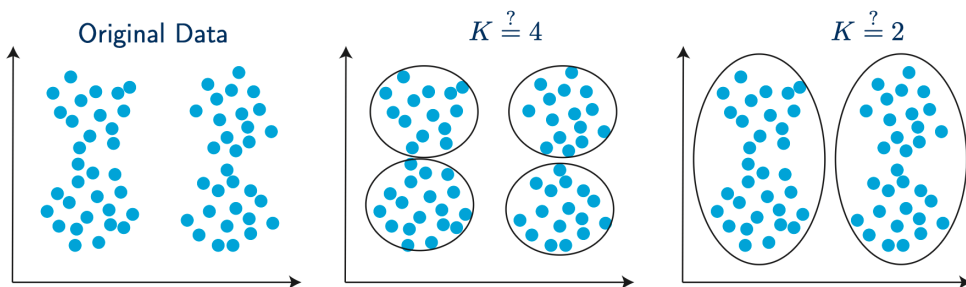
Теорема

T - множество точек k -means++

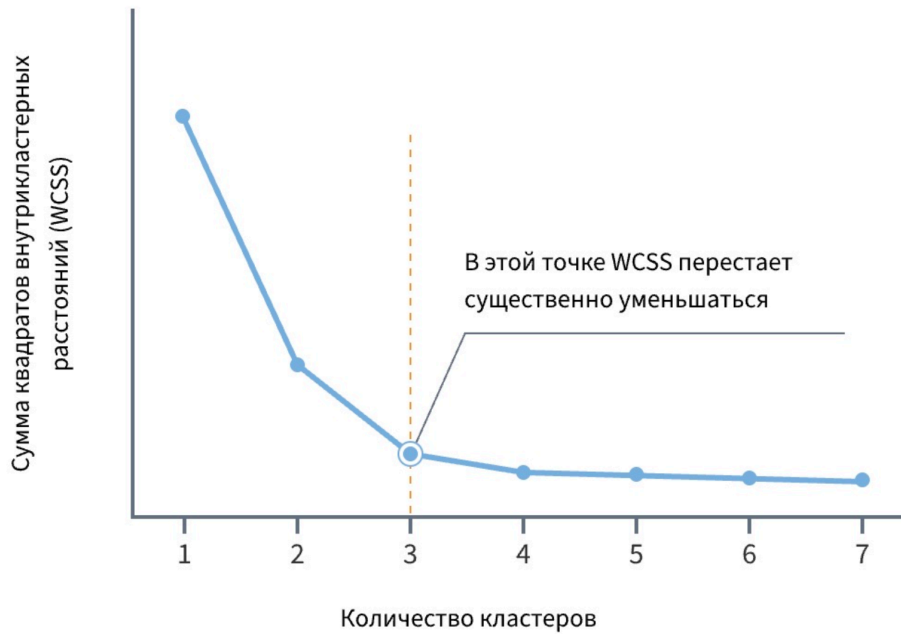
T^* - идеальное множество (cost(S, T*))

$$\mathbb{E} [\text{cost}(S, T)] \leq 8 \cdot (1 + \ln k) \text{cost}(S, T^*)$$

Как выбрать k ?



1) Метод локтя

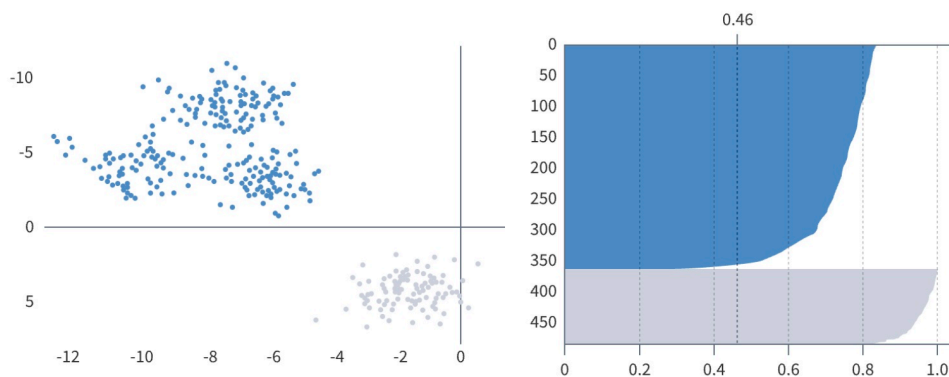


2) Метод силуэтов

$i \rightarrow a(i)$ — среднее расстояние от i до всех остальных точек в кластере

$\rightarrow b(i)$ — среднее расстояние от i до всех точек во втором ближайшем кластере

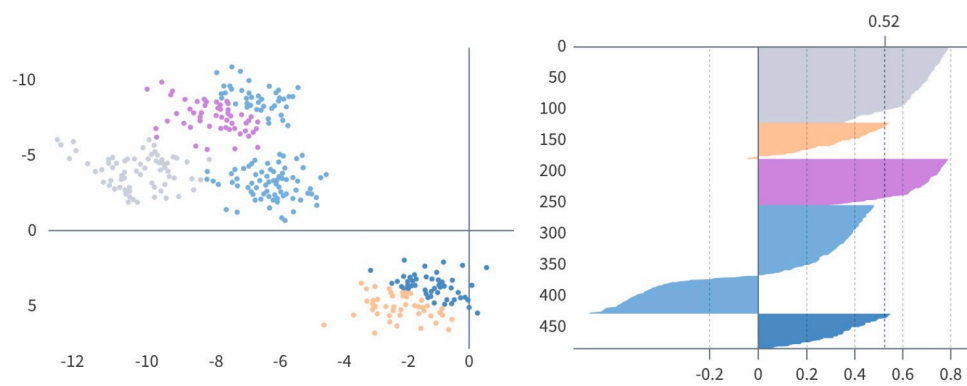
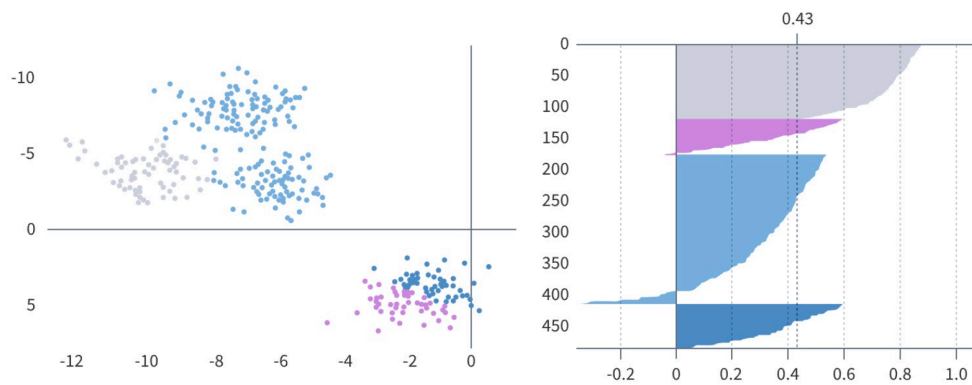
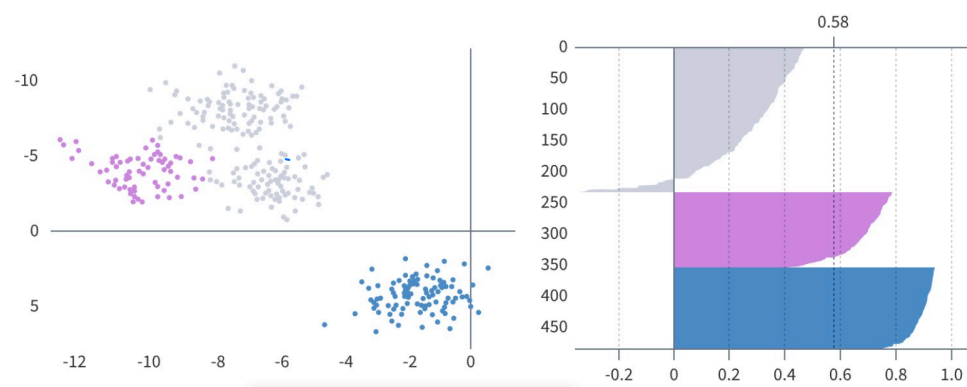
$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \in [-1; 1]$$



$s(i) \approx 1$ — уверен

$s(i) \approx 0$ — не уверен

$s(i) \approx -1$ — скорее всего не в этом кластере



3) Загор стаменик