

Обучение без учителя

Лекция 7

Обучение без учителя (unsupervised learning)

Просто есть данные, без целевой переменной

Есть множество постановок задачи:

- Обучение ассоциативным правилам
- Кластеризация
- Визуализация
- Понижение размерности

Обучение ассоциативным правилам

	0	1	2	3	4	5	6	7	8	9	10	11	12	13
0	shrimp	almonds	avocado	vegetables mix	green grapes	whole weat flour	yams	cottage cheese	energy drink	tomato juice	low fat yogurt	green tea	honey	salad
1	burgers	meatballs	eggs	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
2	chutney	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
3	turkey	avocado	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
4	mineral water	milk	energy bar	whole wheat rice	green tea	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

Параметры

$$1. \textit{Support}(A) = \frac{\text{Транзакции с } A}{\text{Все транзакции}}; \textit{Support}(A \rightarrow B) = \frac{\text{Транзакции с } A \text{ и } B}{\text{Все транзакции}}$$

$$2. \textit{Confidence}(A \rightarrow B) = \frac{\text{Транзакции с } A \text{ и } B}{\text{Транзакции с } A} = \frac{\textit{Support}(A \rightarrow B)}{\textit{Support}(A)}$$

Обучение ассоциативным правилам

$$3. Lift(A \rightarrow B) = \frac{Confidence(A \rightarrow B)}{Support(B)}$$

$Lift(A \rightarrow B) = 1$ – товары независимы

$Lift(A \rightarrow B) > 1$ – товары зависимы положительно

$Lift(A \rightarrow B) < 1$ – товары зависимы негативно

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift
6	(ground beef)	(mineral water)	0.098254	0.238368	0.040928	0.416554	1.747522
74	(ground beef)	(spaghetti)	0.098254	0.174110	0.039195	0.398915	2.291162
128	(ground beef)	(chocolate)	0.098254	0.163845	0.023064	0.234735	1.432669
175	(ground beef)	(milk)	0.098254	0.129583	0.021997	0.223881	1.727704

Кластеризация

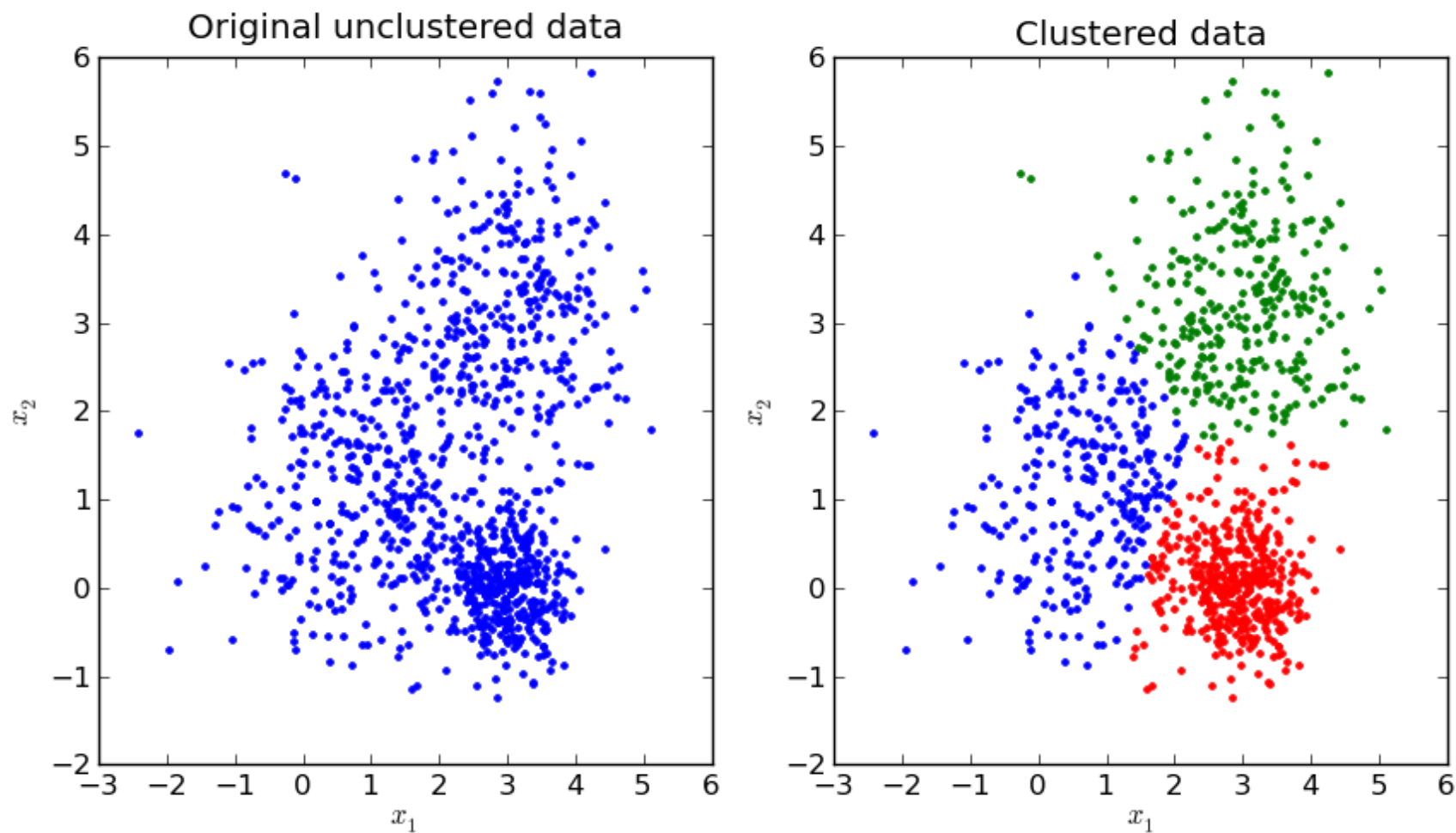
$$X = (x_i)_{i=1}^l$$

Необходимо построить такой алгоритм, который каждую точку соотнесет к одному из кластеров M : $a: X \rightarrow \{1, \dots, M\}$

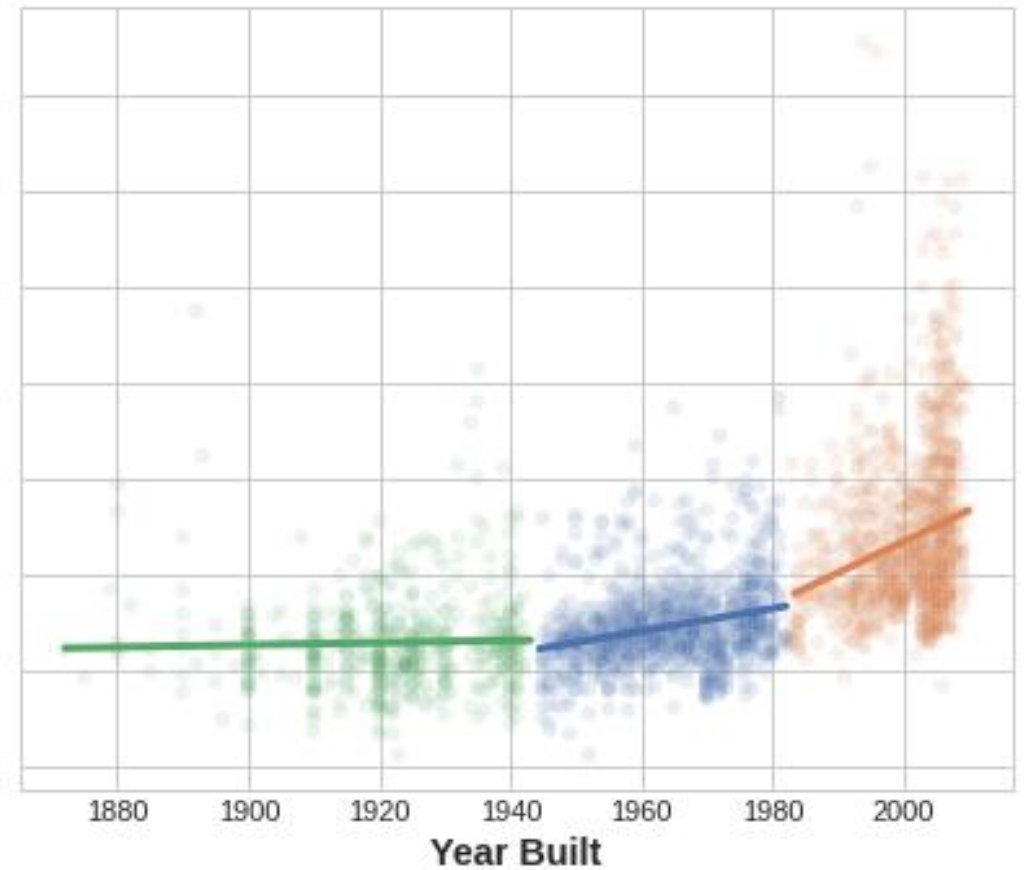
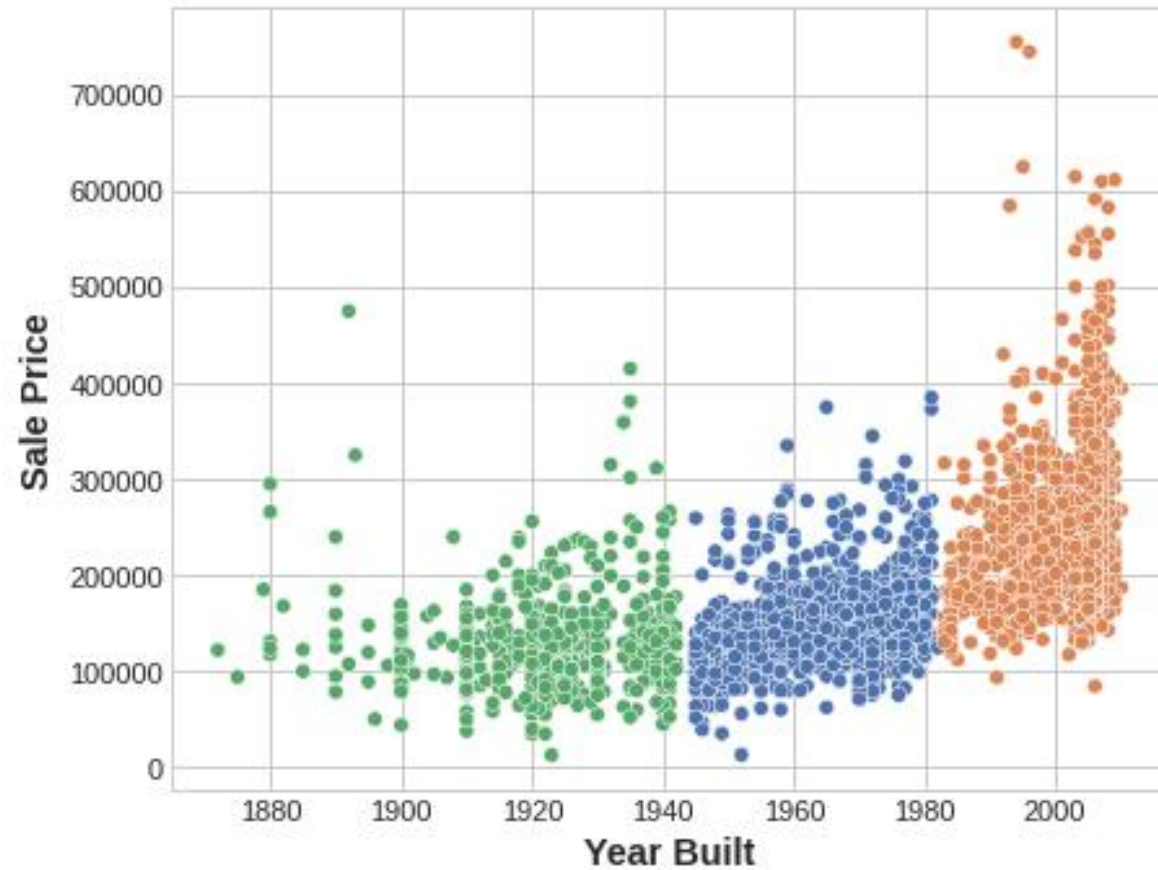
Примеры:

- Сегментация рынка
- Анализ социальных сетей
- Группировка результатов поиска
- Обнаружение аномалий

Пример кластеризации



Пример кластеризации



Кластеризация в общем

Выражение для учета расстояний между всеми парами точек:

$$T = \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N d_{ii'} = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \left(\sum_{C(i')=k} d_{ii'} + \sum_{C(i') \neq k} d_{ii'} \right)$$

$$T = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i')=k} d_{ii'} + \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i') \neq k} d_{ii'}$$

$$T = W(C) + B(C)$$

Кластеризация в общем

$$T = W(C) + B(C)$$

Внутрикластерное расстояние:

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i')=k} d_{ii'}$$

Межкластерное расстояние:

$$B(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i') \neq k} d_{ii'}$$

Кластеризация. Решение

Можно решить оптимизационную задачу, пытаясь найти наилучшее разбиение N точек по K кластерам. Но для этого нужно перебрать все варианты, которых будет:

$$S(N, K) = \frac{1}{K!} \sum_{k=1}^K (-1)^{(K-k)} \binom{K}{k} k^N$$

Так например $S(10, 4) = 34105$.

Но $S(19, 4) \cong 10^{10}$

K-means

Будем считать расстояние как евклидово:

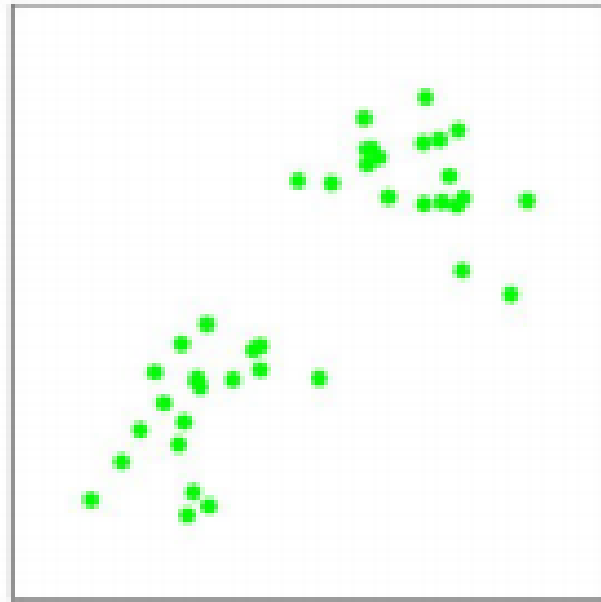
$$d(x_i, x_{i'}) = \sum_{j=1}^p (x_{ij} - x_{i'j})^2 = \|x_i - x_{i'}\|^2$$

Распишем внутрикластерное расстояние как:

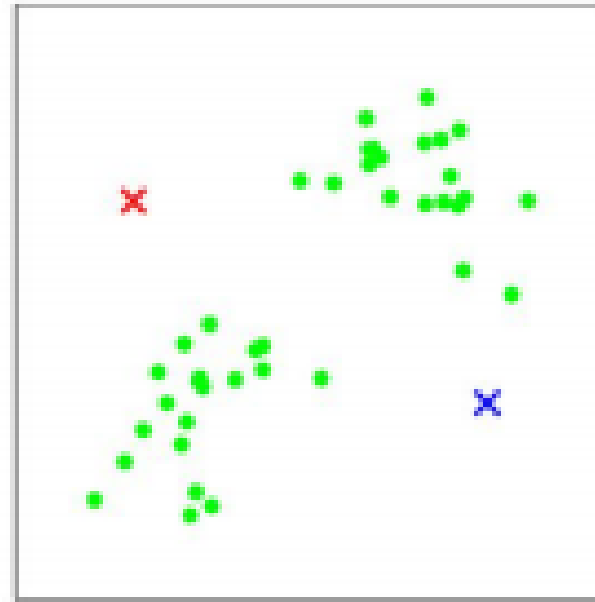
$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i')=k} \|x_i - x_{i'}\|^2 = \sum_{k=1}^K N_k \sum_{C(i)=k} \|x_i - \bar{x}_k\|^2$$

K-means. Шаг 0

Случайно инициализируем центры кластеров C_1, \dots, C_k .



(a)



(b)

K-means. Шаг 1

Обновляем принадлежности точек к кластерам

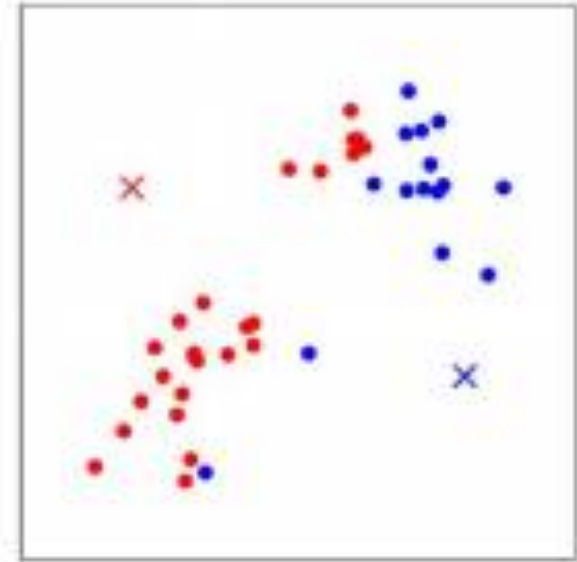
$$C(i) = \operatorname{argmin}_{1 \leq k \leq K} \|x_i - m_k\|^2$$



(a)



(b)

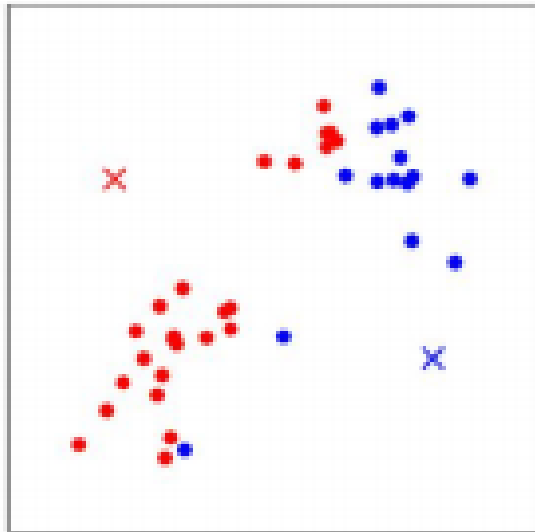


(c)

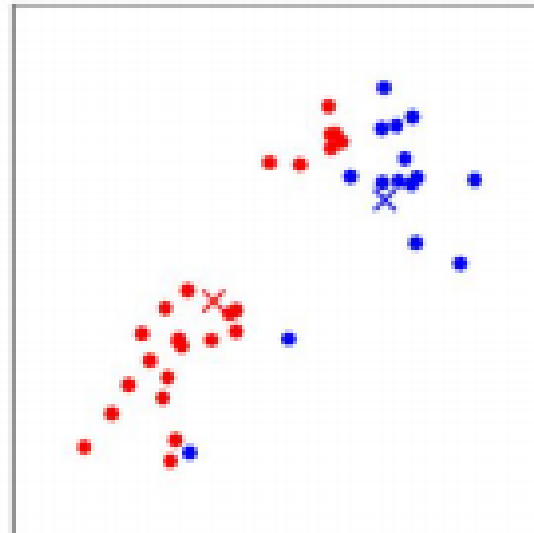
K-means. Шаг 2

Обновляем центры кластеров

$$\bar{x}_S = \operatorname{argmin}_m \sum_{i \in S} \|x_i - m_k\|^2$$



(c)



(d)

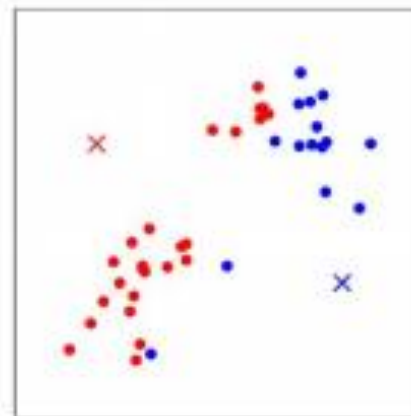
K-means. Повторяем шаги



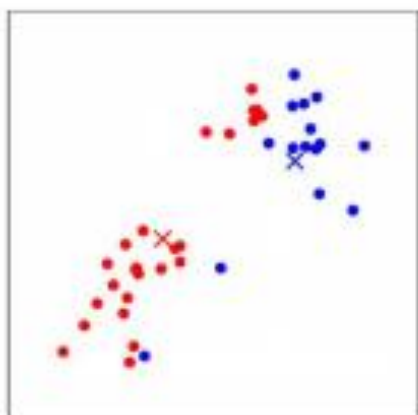
(a)



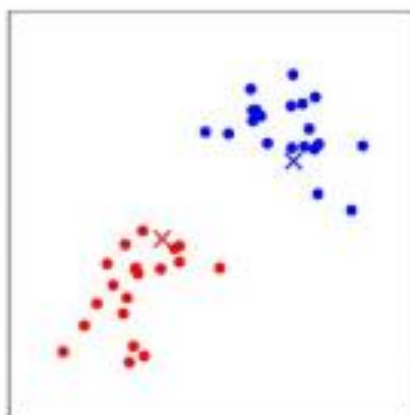
(b)



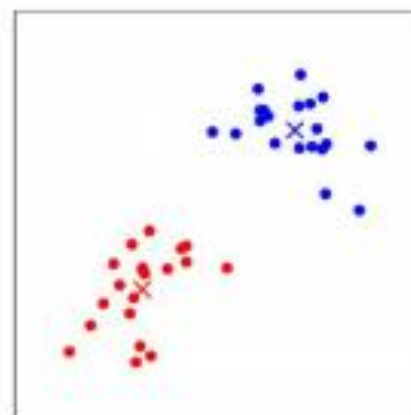
(c)



(d)



(e)

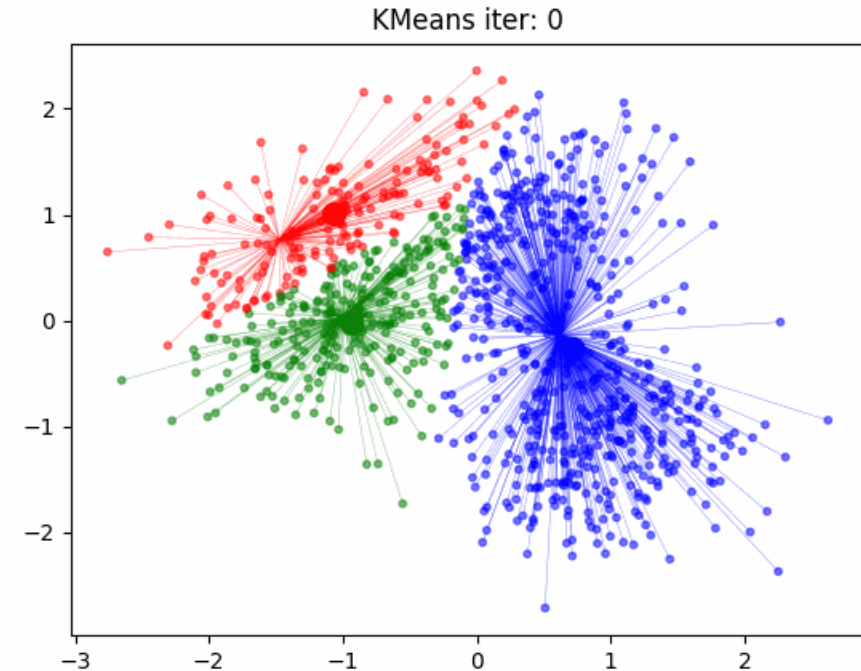
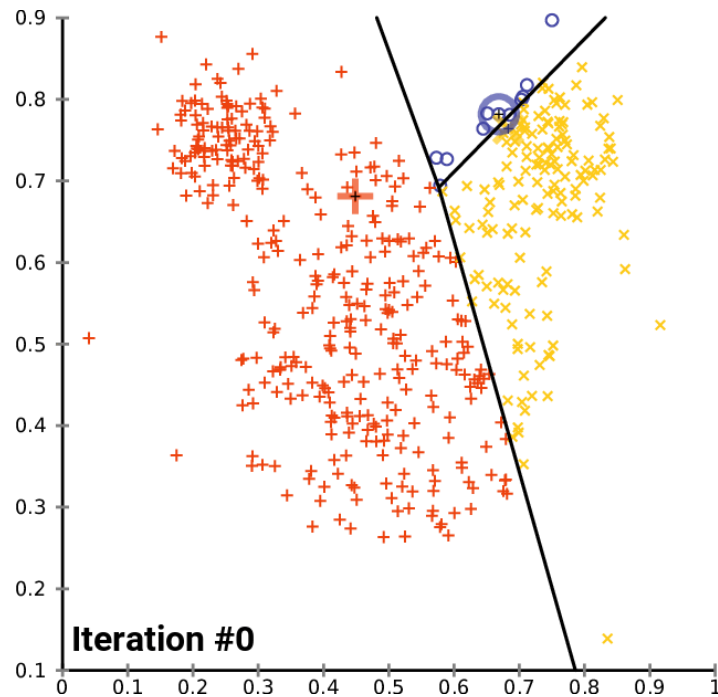


(f)

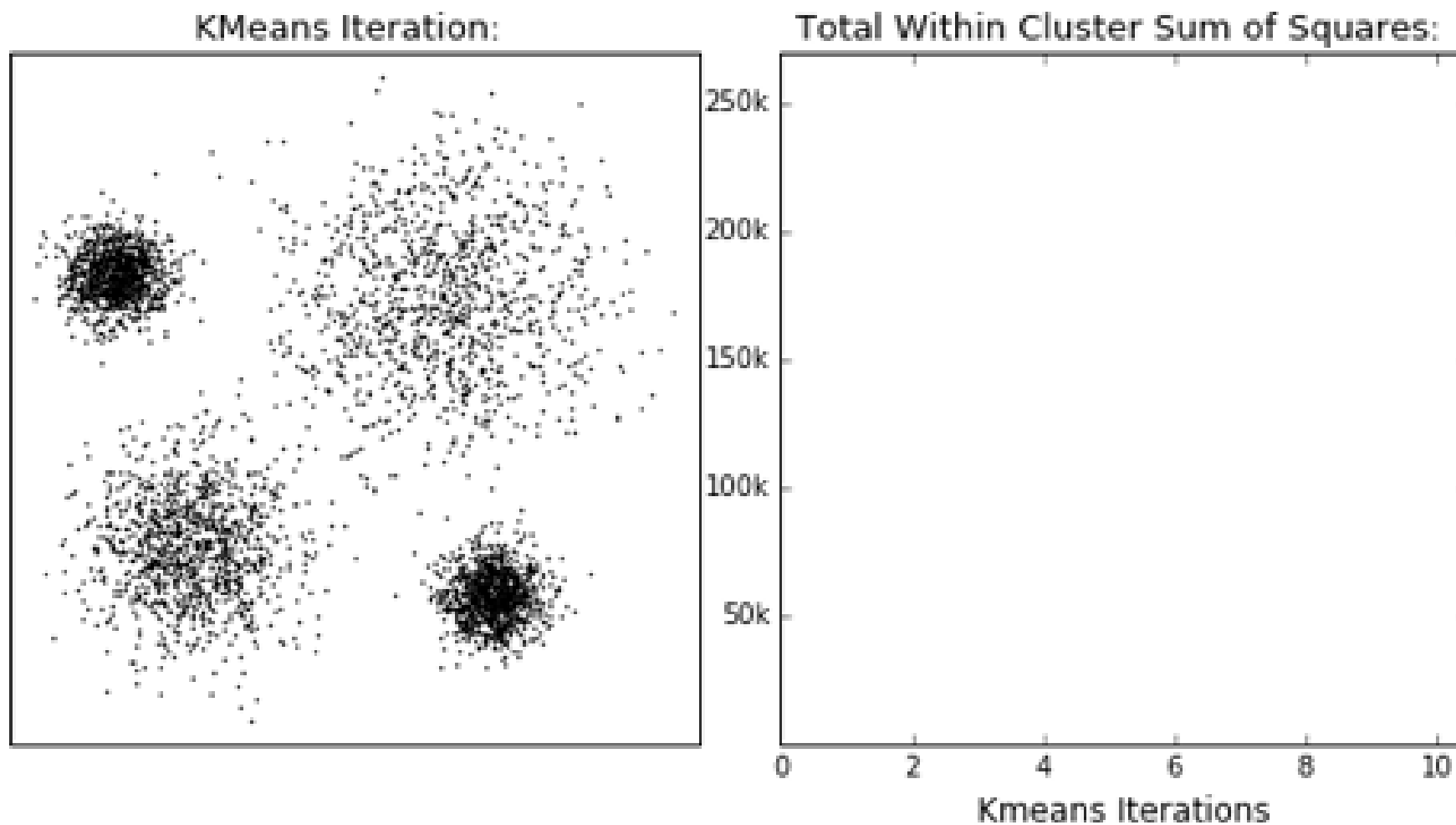
K-means. Повторяем шаги

В общем решаем задачу:

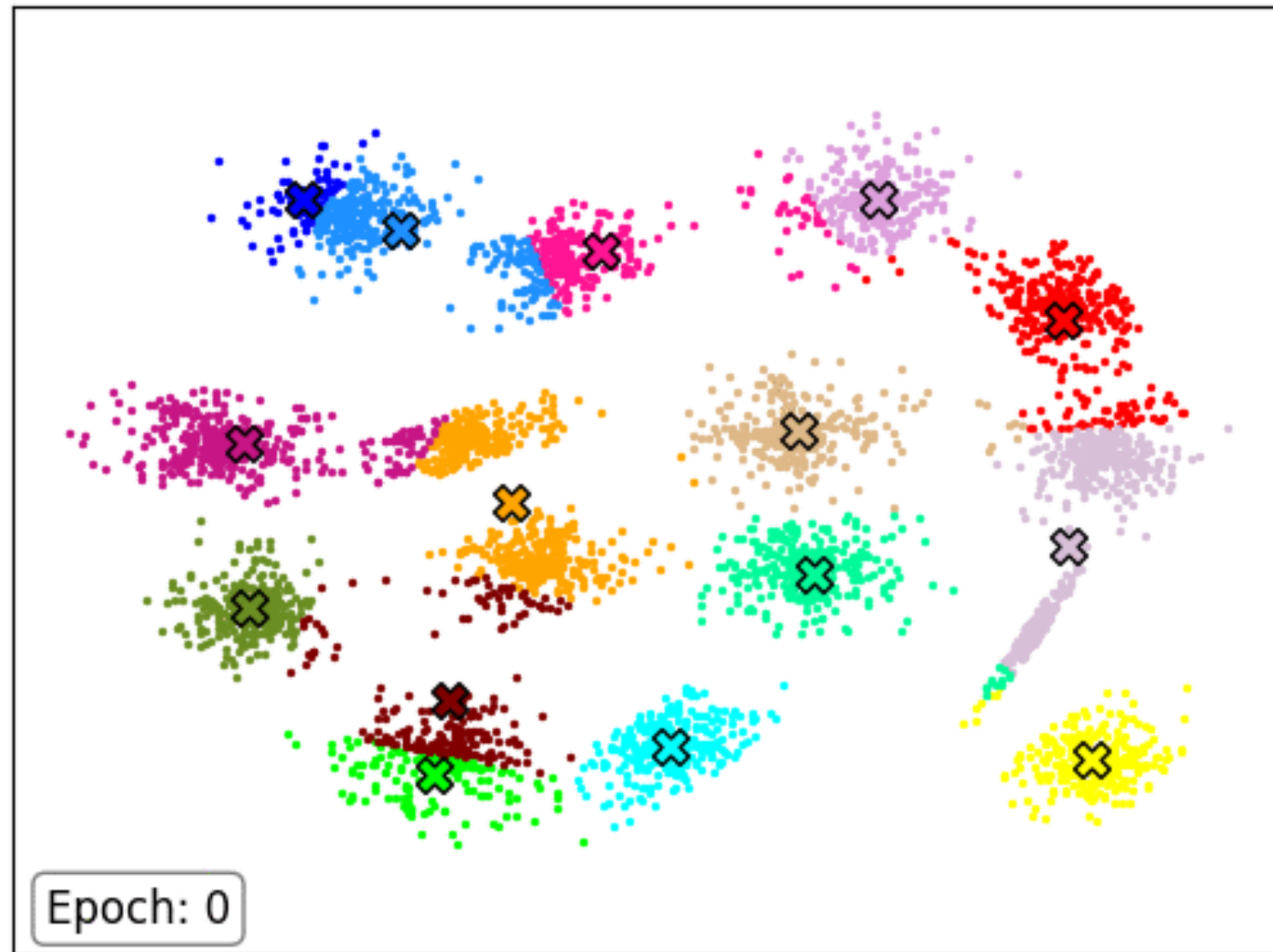
$$C^* = \min_{C, \{m_k\}_1^K} \sum_{k=1}^K N_k \sum_{C(i)=k} \|x_i - m_k\|^2$$



K-means. В картинках

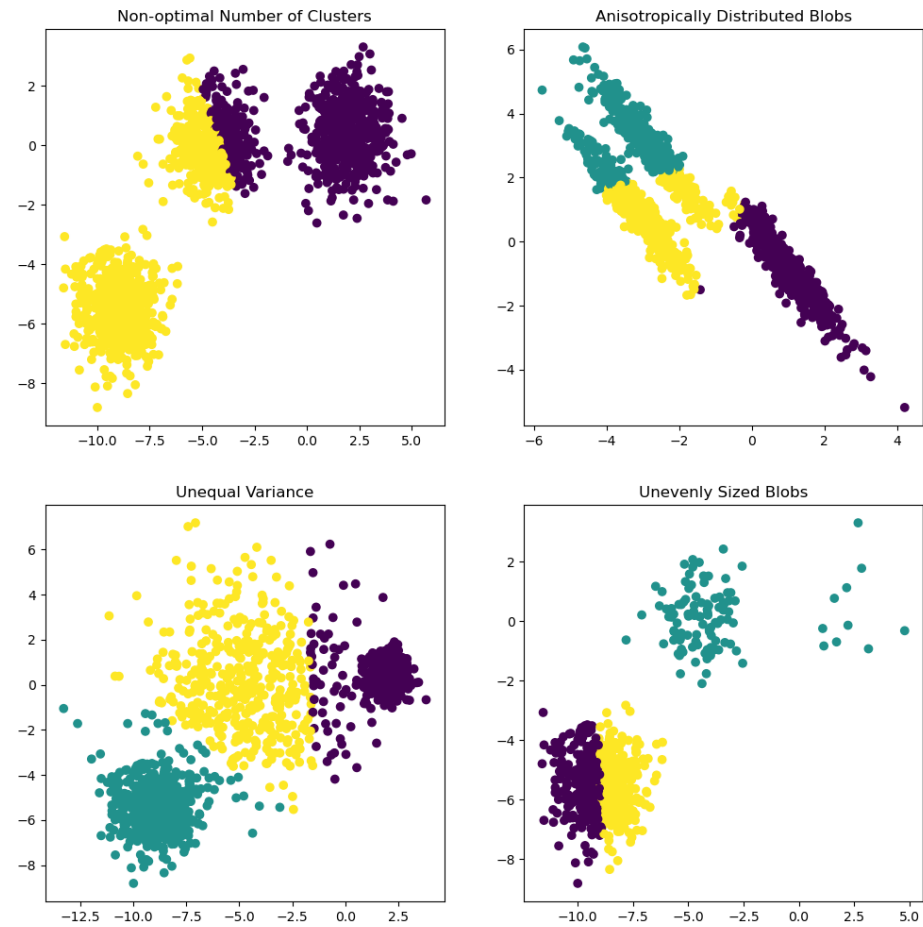


K-means++



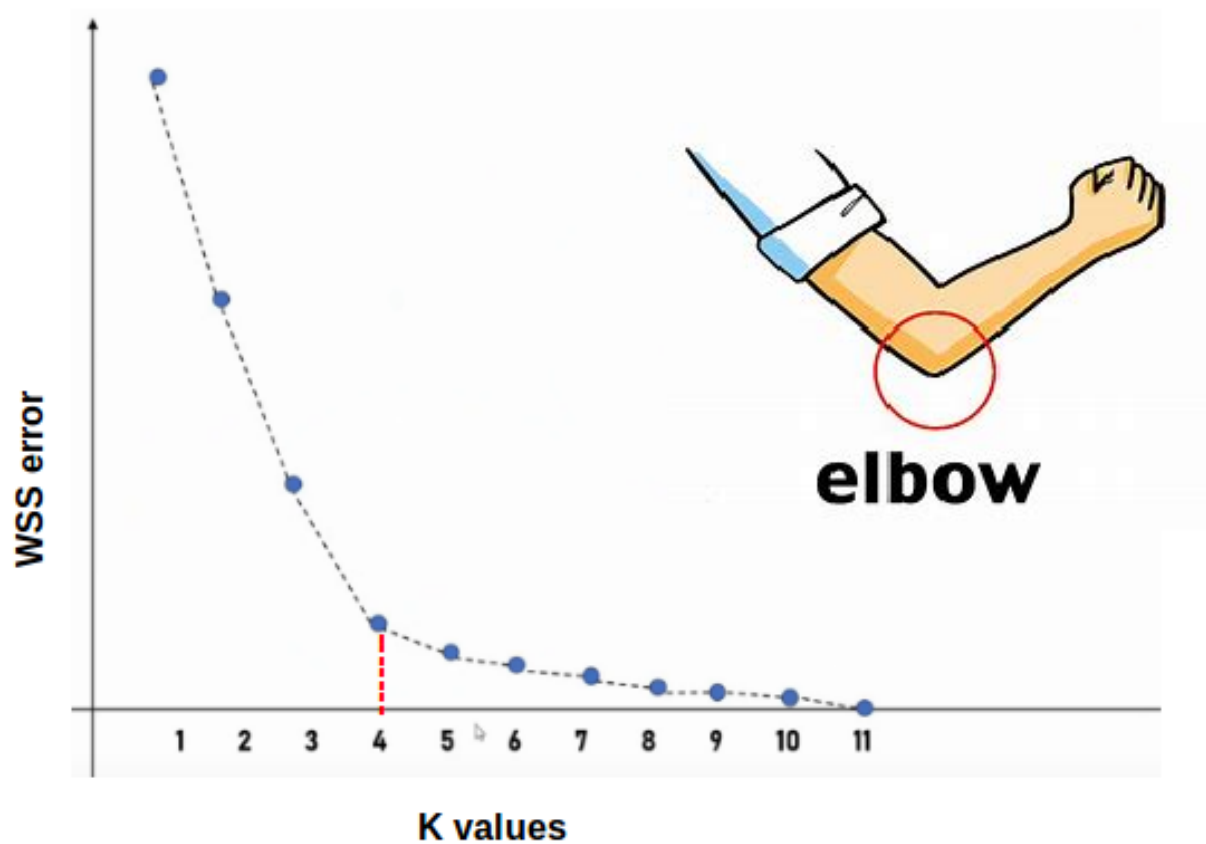
K-means. Проблемы

Unexpected KMeans clusters



K-means. Как найти K

Elbow method



Метрики в кластеризации

Внутрикластерное расстояние:

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i')=k} d_{ii'}$$

Межкластерное расстояние:

$$B(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i') \neq k} d_{ii'}$$

Метрики в кластеризации

Коэффициент силуэта:

$$S(x_i) = \frac{B(x_i) - A(x_i)}{\max(B(x_i), A(x_i))},$$

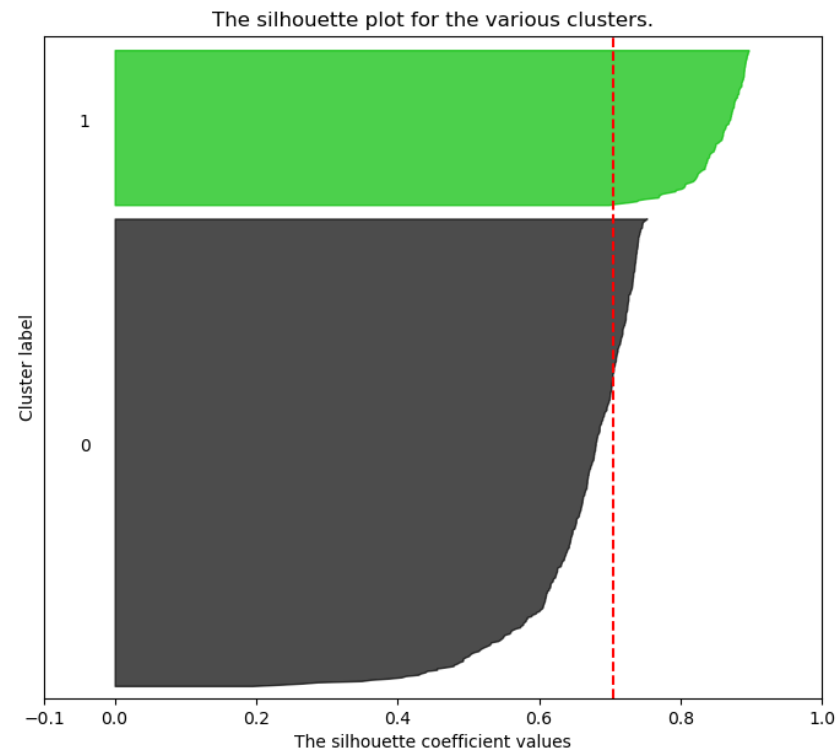
где $A(x_i)$ – среднее расстояние между x_i и точками его кластера

$B(x_i)$ – среднее расстояние между x_i и точками ближайшего кластера

Само значение $S(x_i) \in [-1; 1]$ и чем больше, тем лучше

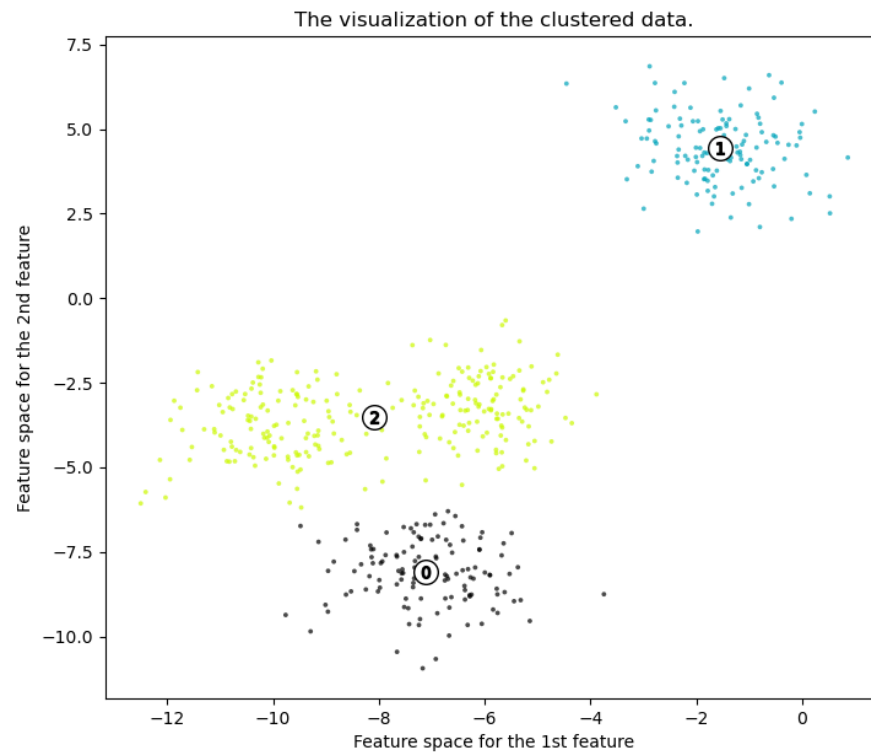
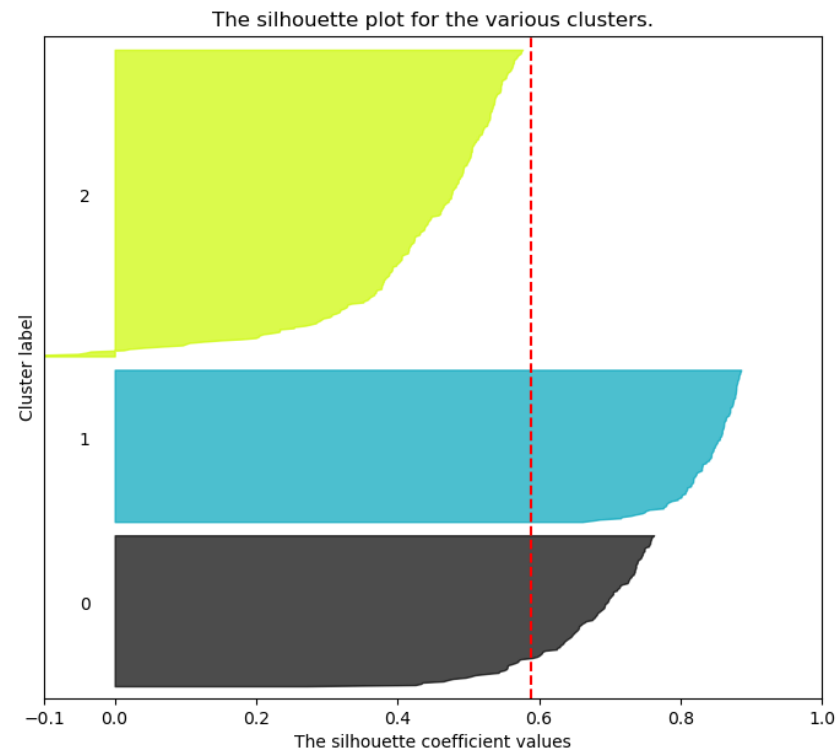
Коэффициент силуэта, 2 кластера

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 2$



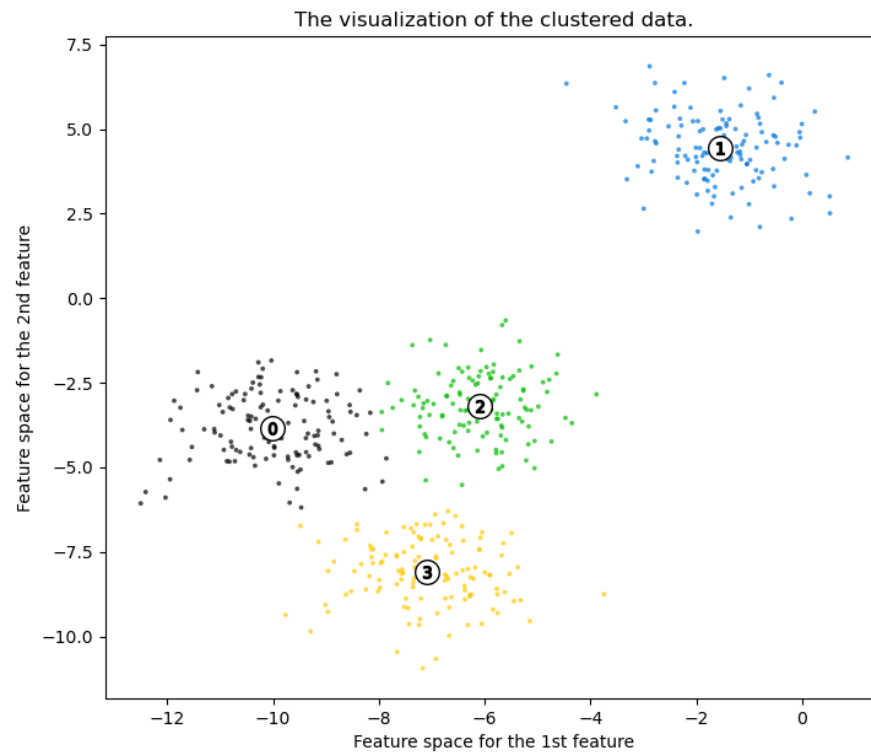
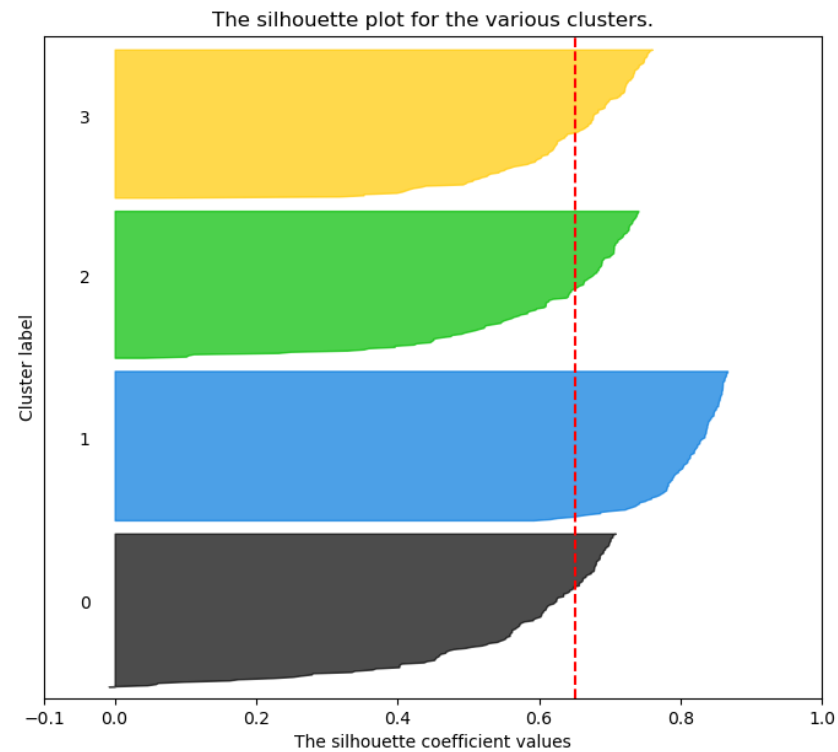
Коэффициент силуэта, 3 кластера

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 3$



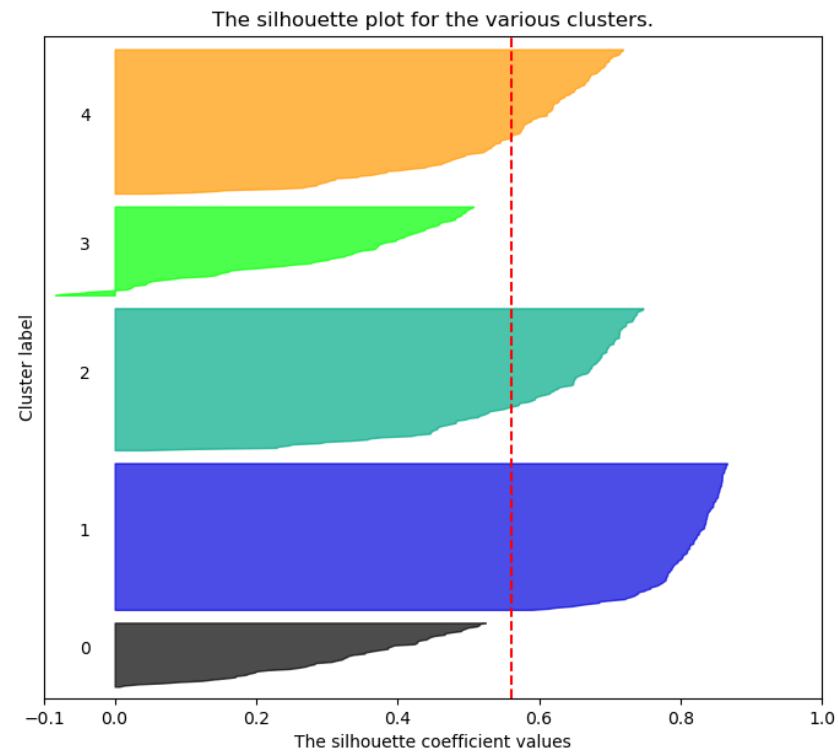
Коэффициент силуэта, 4 кластера

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 4$



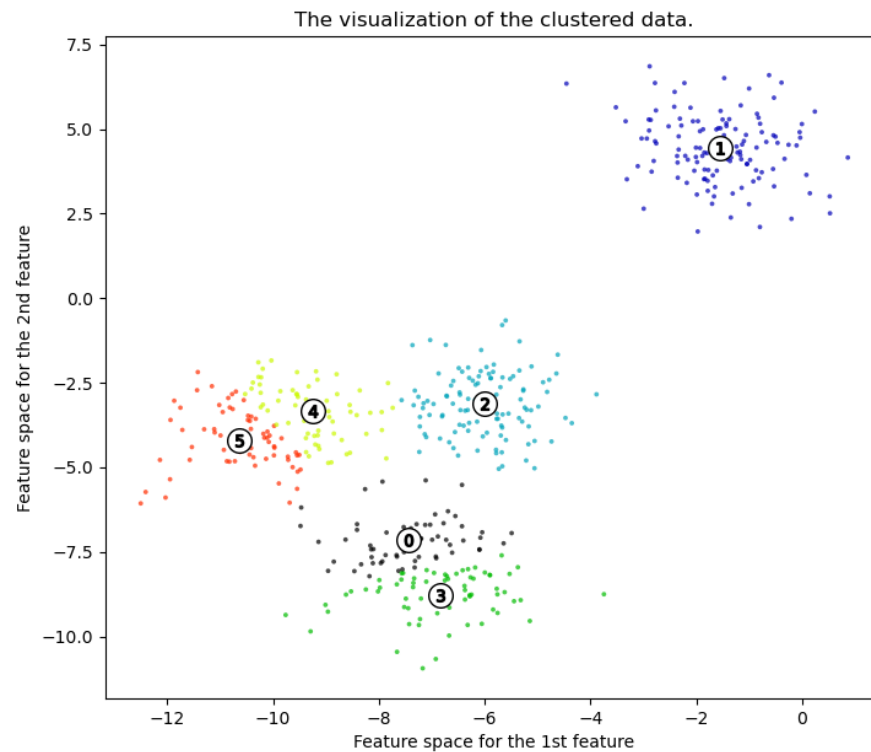
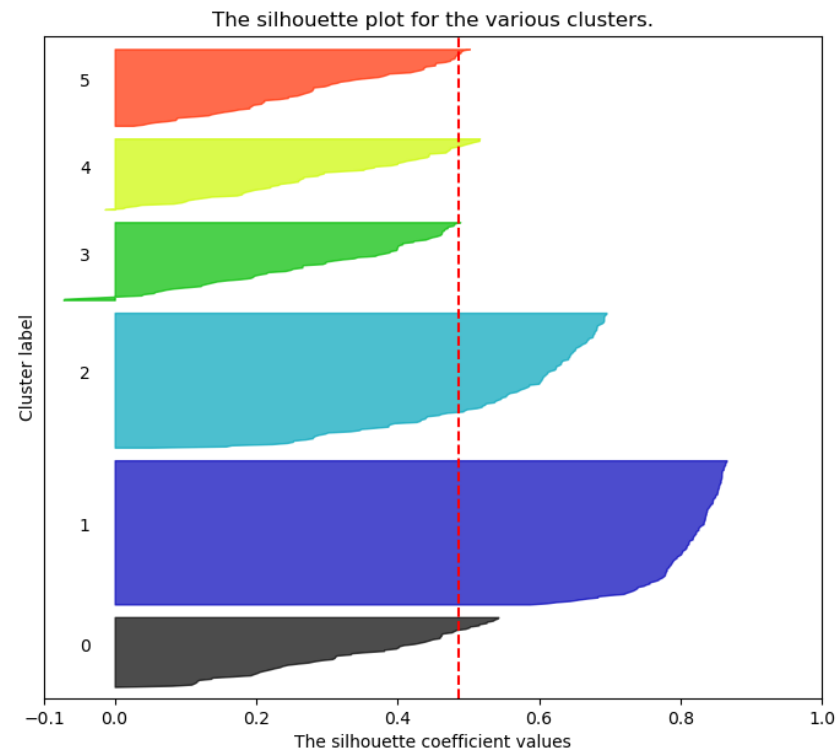
Коэффициент силуэта, 5 кластера

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 5$



Коэффициент силуэта, 6 кластера

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 6$



Метрики в кластеризации

Индекс Данна:

$$Dunn Index = \frac{\min_{1 \leq k < k' \leq K} d(k, k')}{\max_{1 \leq k \leq K} d(k)},$$

где $d(k, k')$ – межкластерное расстояние

$d(k)$ – внутрикластерное расстояние

Чем больше, тем лучше

Метрики в кластеризации. Есть разметка

Введем формулы энтропии:

$$H_{class} = - \sum_{c=1}^C \frac{m_c}{n} \log \frac{m_c}{n}$$

$$H_{clust} = - \sum_{k=1}^K \frac{n_k}{n} \log \frac{n_k}{n}$$

$$H_{class|clust} = - \sum_{c=1}^C \sum_{k=1}^K \frac{n_{ck}}{n_k} \log \frac{n_{ck}}{n_k}$$

Метрики в кластеризации. Есть разметка

Гомогенность:

$$Homogeneity = 1 - \frac{H_{class|clust}}{H_{class}}$$

Лучший вариант – каждый кластер содержит элементы только одного класса

Полнота:

$$Completeness = 1 - \frac{H_{clust|class}}{H_{clust}}$$

Лучший вариант – все объекты класса содержатся в одном кластере

Метрики в кластеризации. Есть разметка

V-мера:

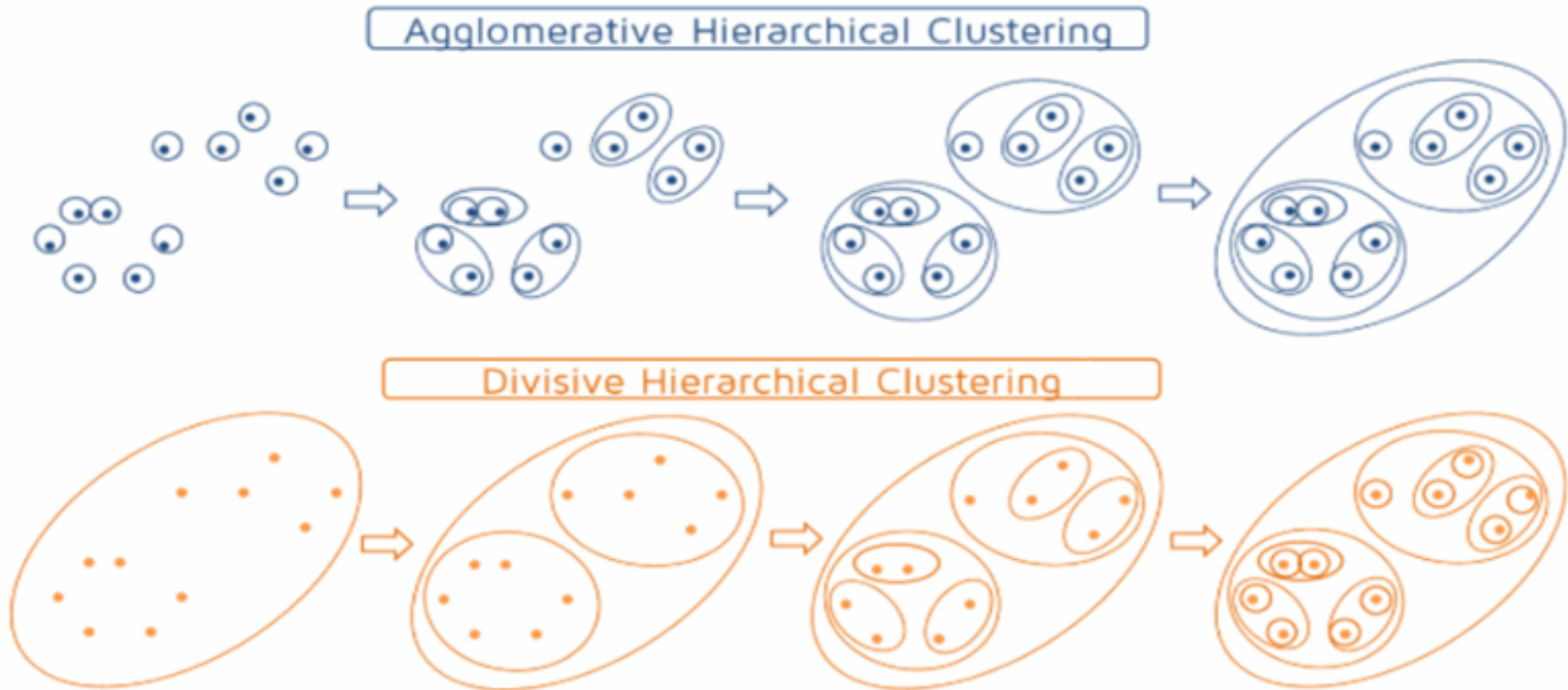
$$V_{\beta} = \frac{(1 + \beta) \cdot \text{Homogeneity} \cdot \text{Completeness}}{\beta \cdot \text{Homogeneity} + \text{Completeness}}$$

$$V_1 = 2 \frac{\text{Homogeneity} \cdot \text{Completeness}}{\text{Homogeneity} + \text{Completeness}}$$

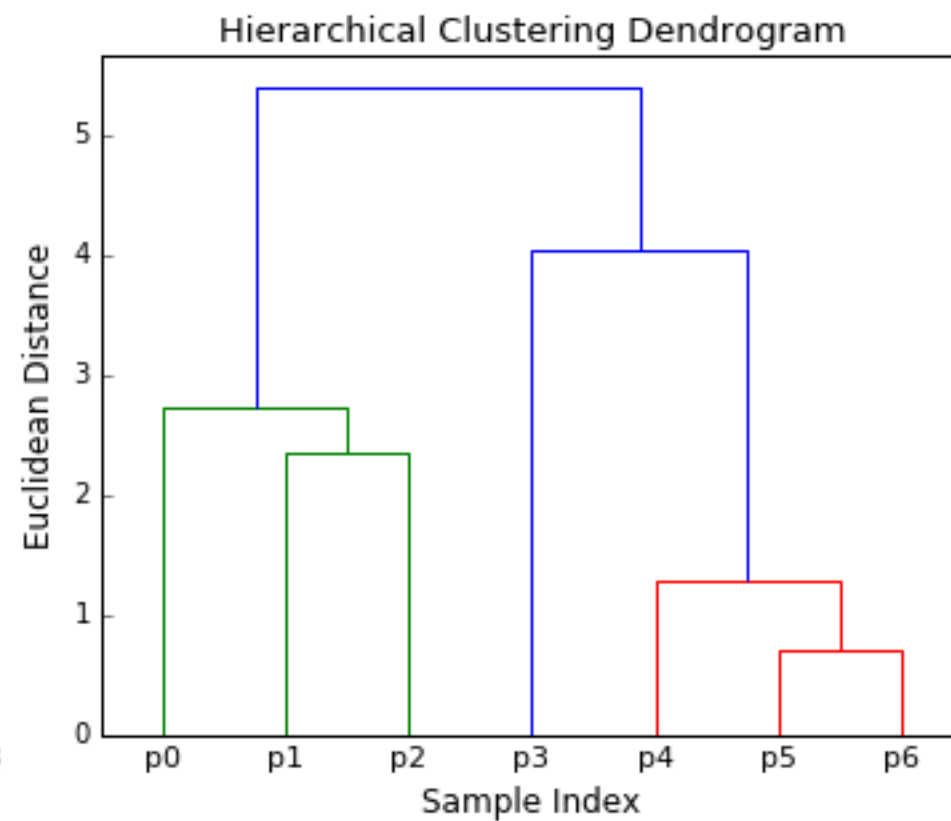
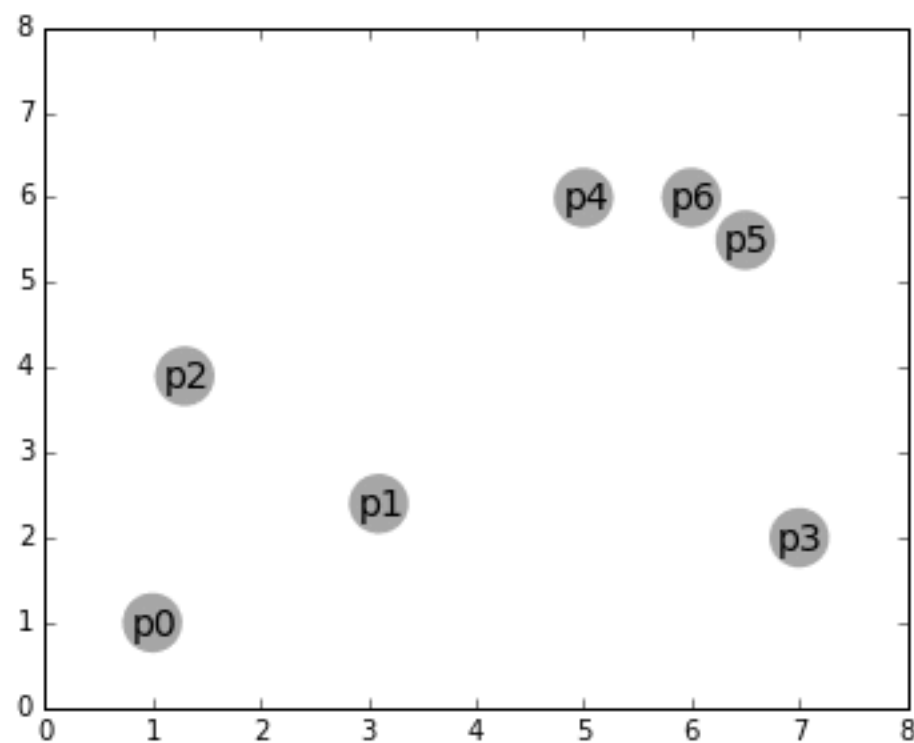
Дополнительно про метрики:

<https://habr.com/ru/companies/yandex/articles/500742/>

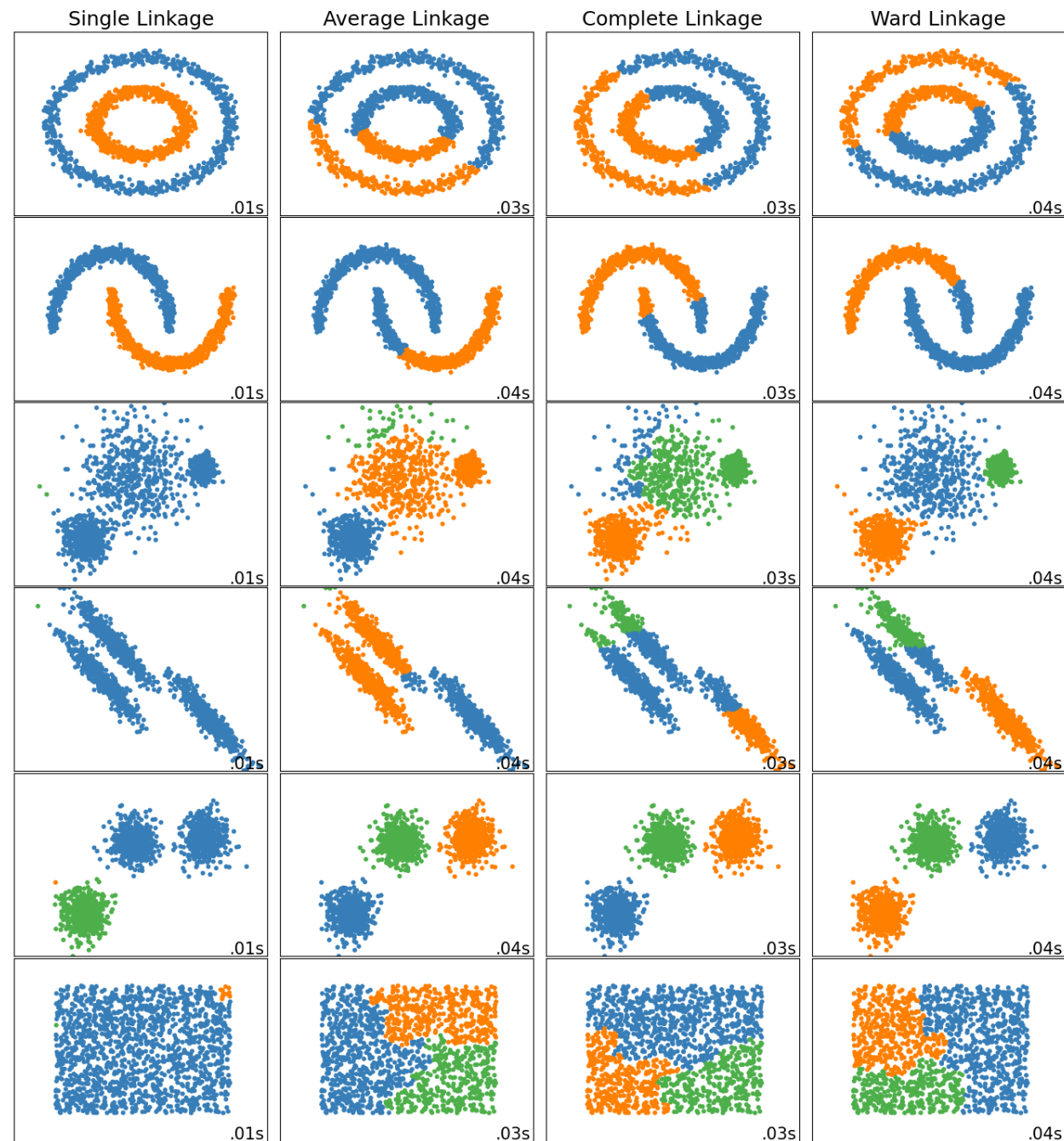
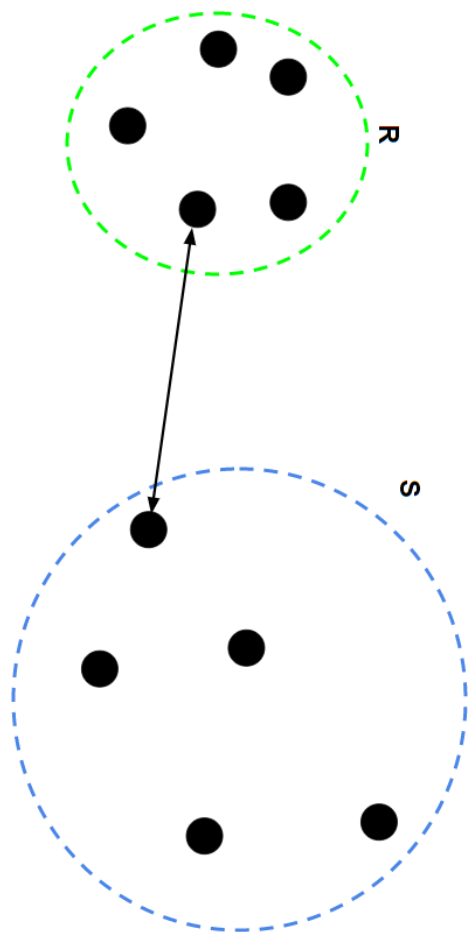
Иерархические алгоритмы кластеризации



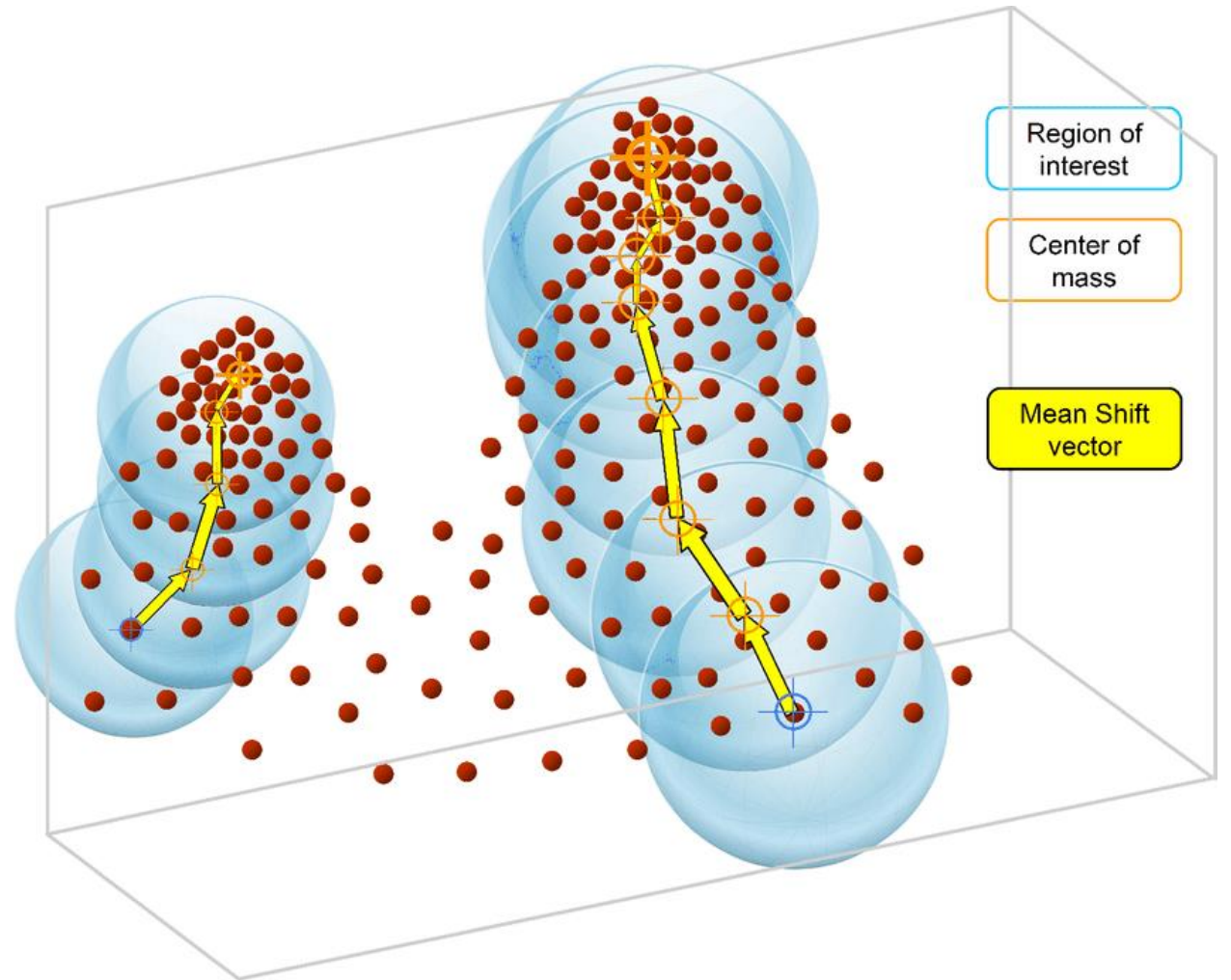
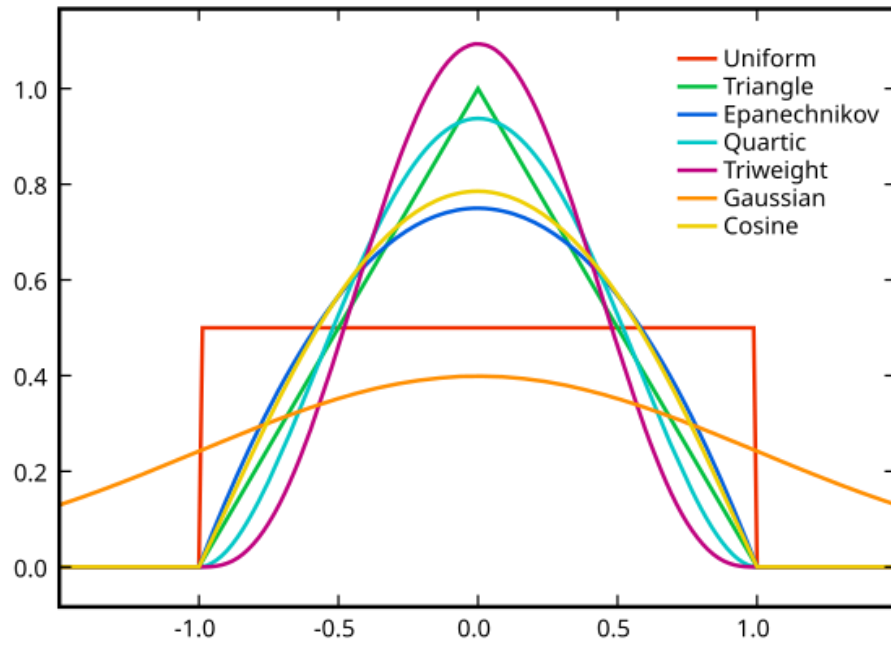
Дендрограммы



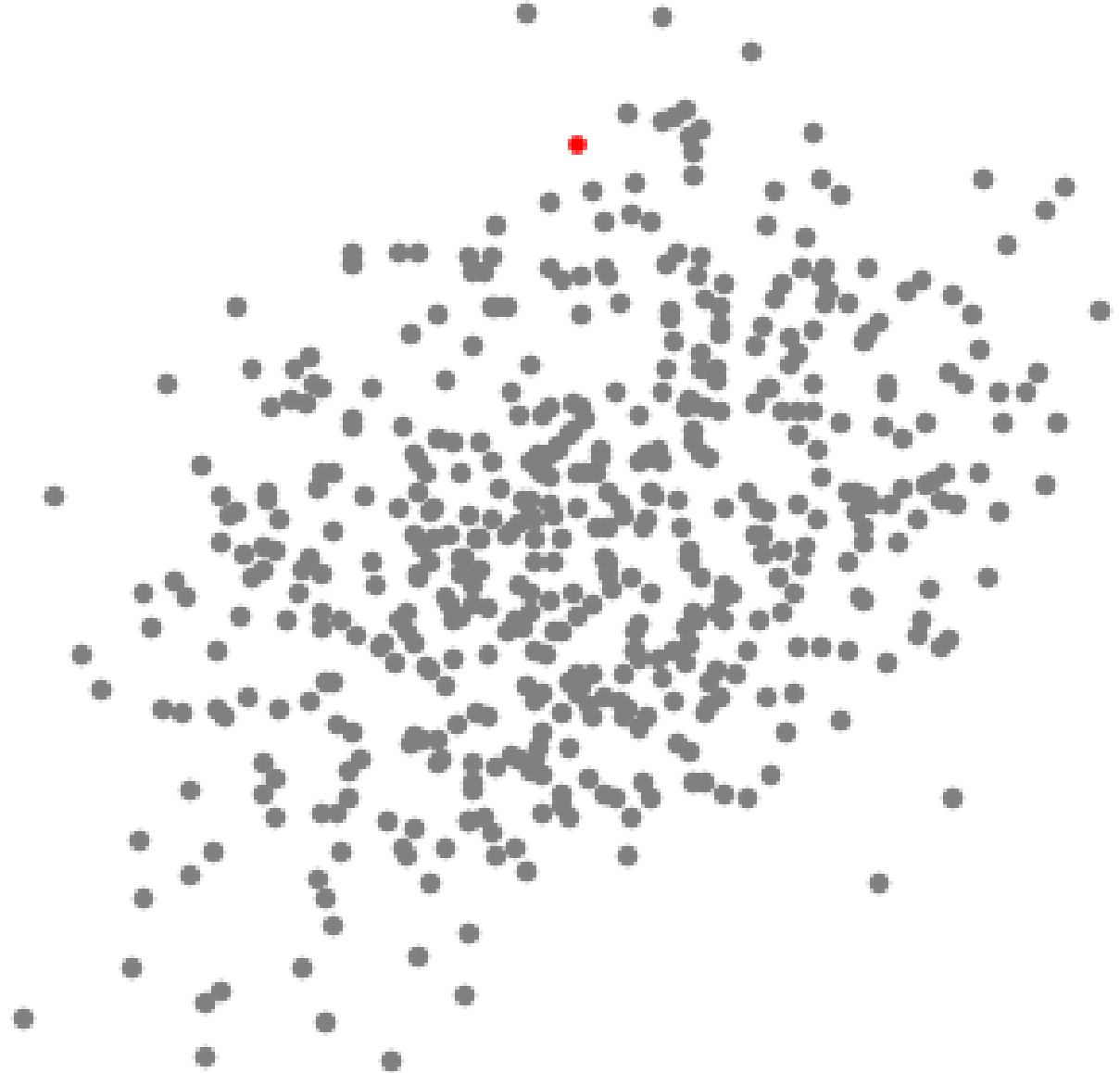
Можно по-разному считать
межкластерные
расстояния



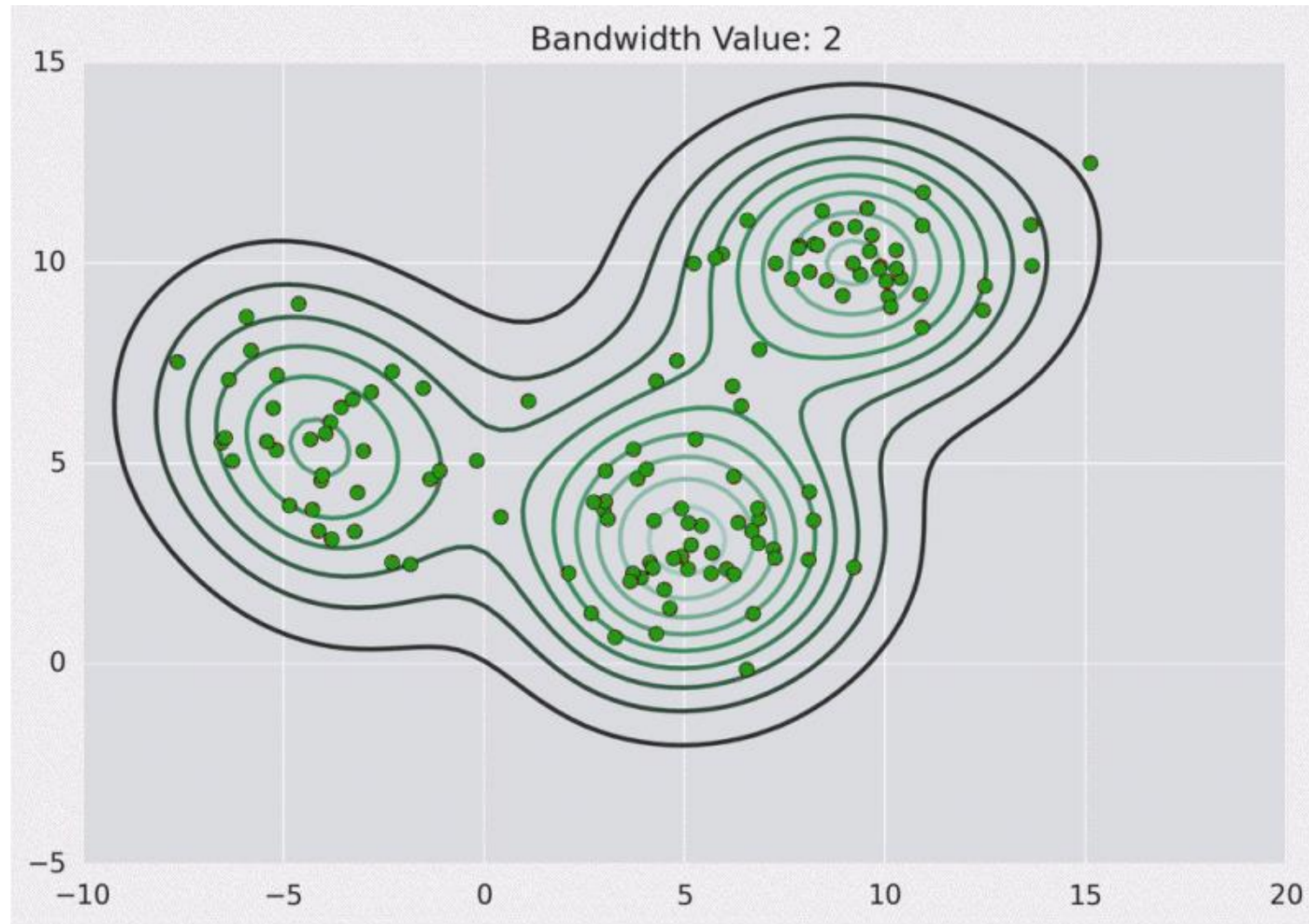
Mean Shift



Mean Shift



Mean Shift

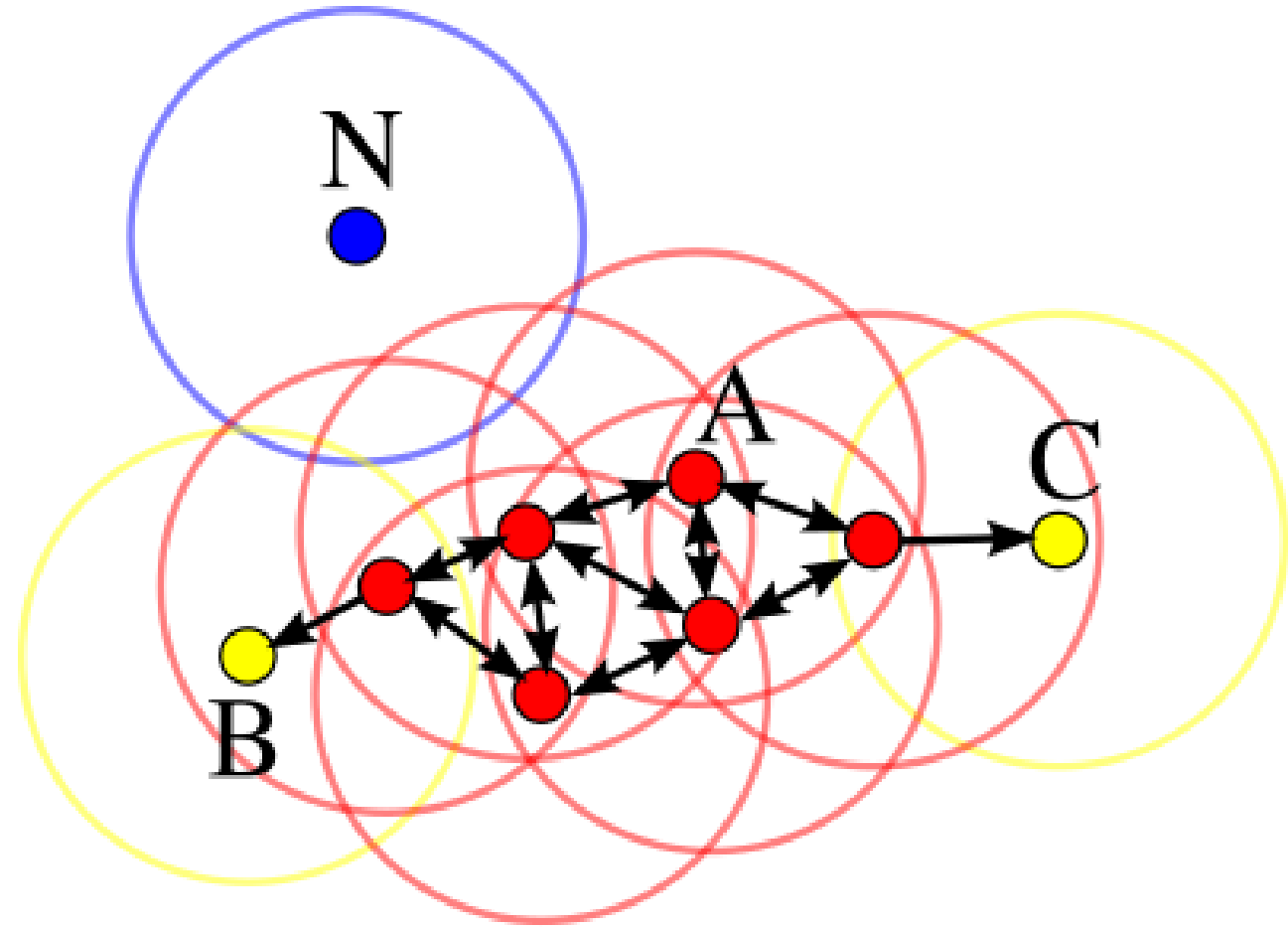


DBSCAN

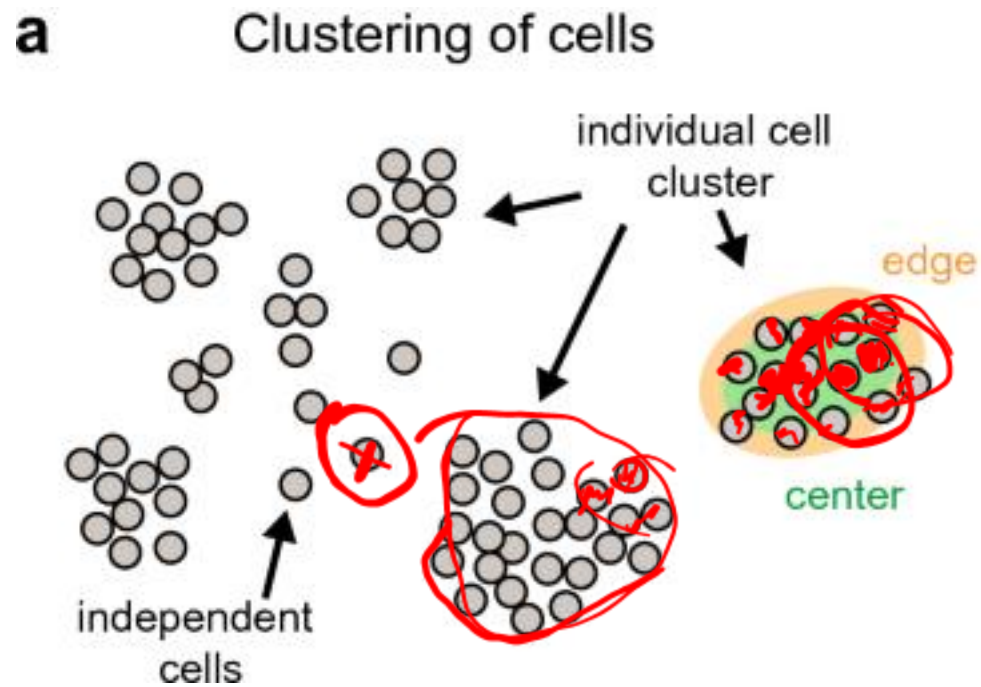
Типы точек:

1. **Основные** – в окрестности ε находится хотя бы N точек
2. **Граничные** – в окрестности ε кол-во **основных** точек $< N$, но хотя бы 1
3. **Шумовые** – в окрестности ε нет **основных** точек и точек $< N$

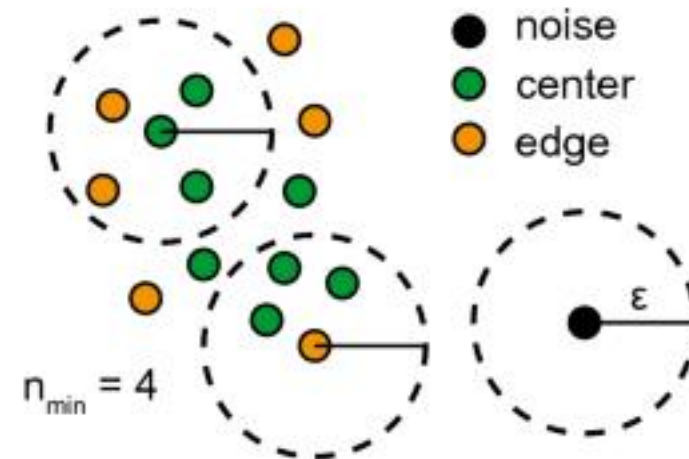
ε и N – гиперпараметры



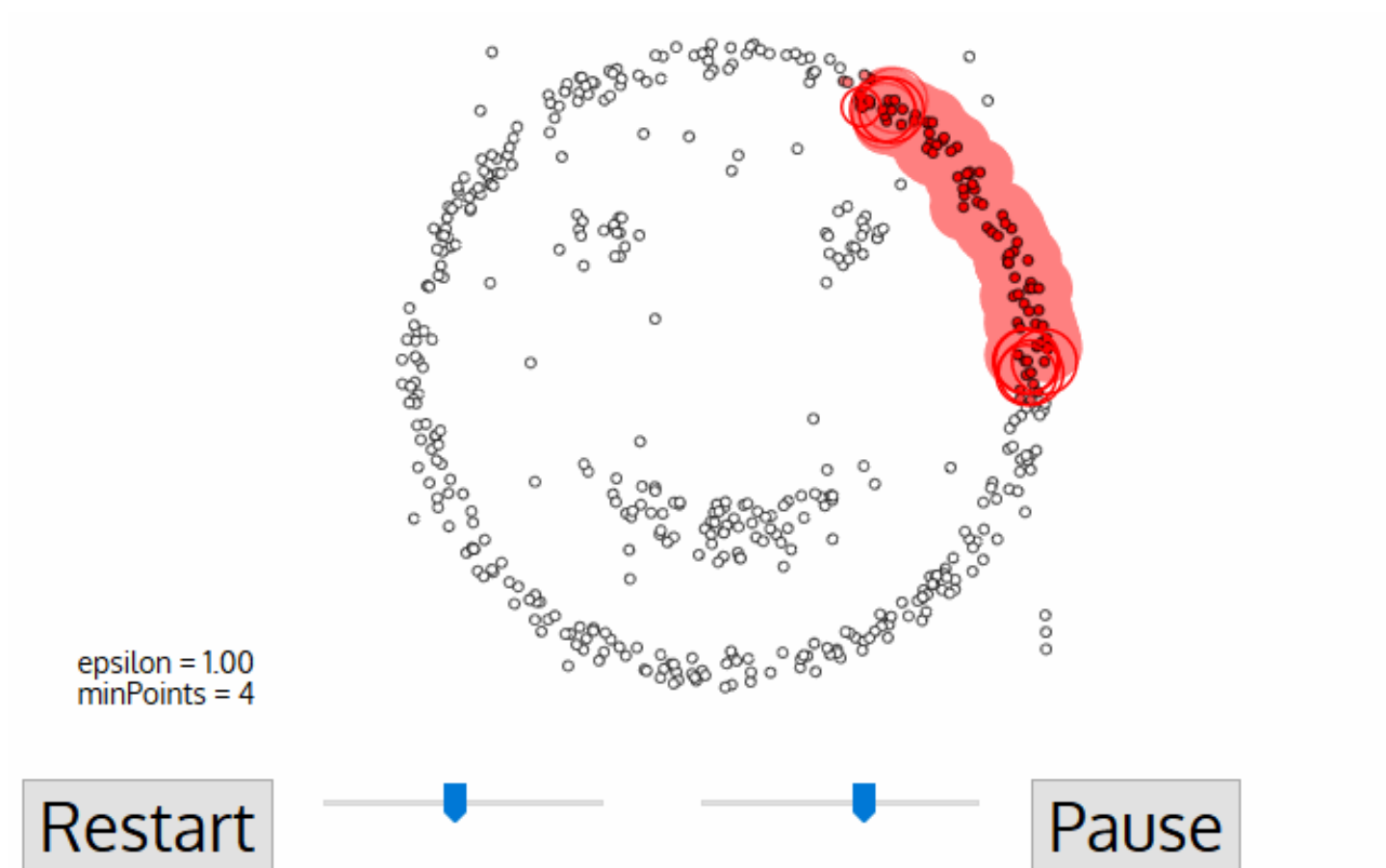
DBSCAN. Пример



b DBSCAN cell classification



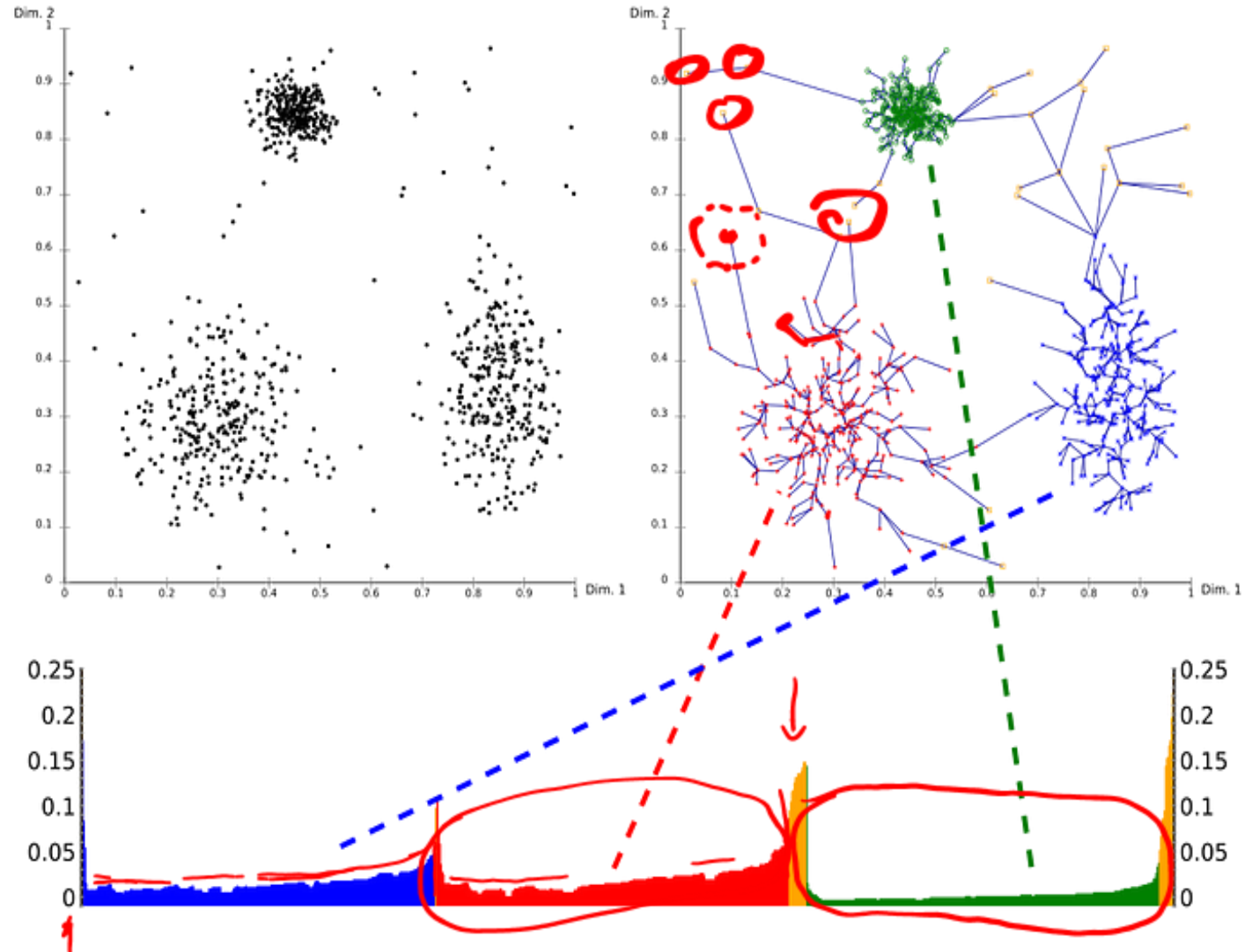
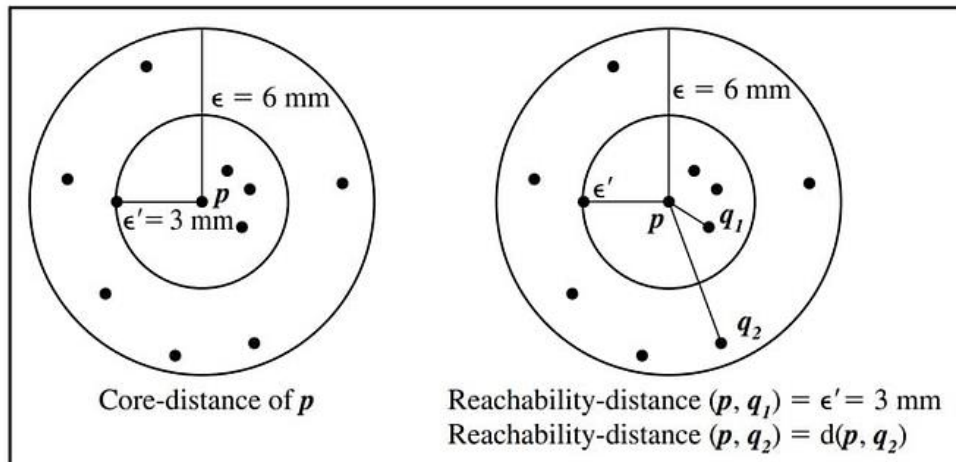
DBSCAN в действии



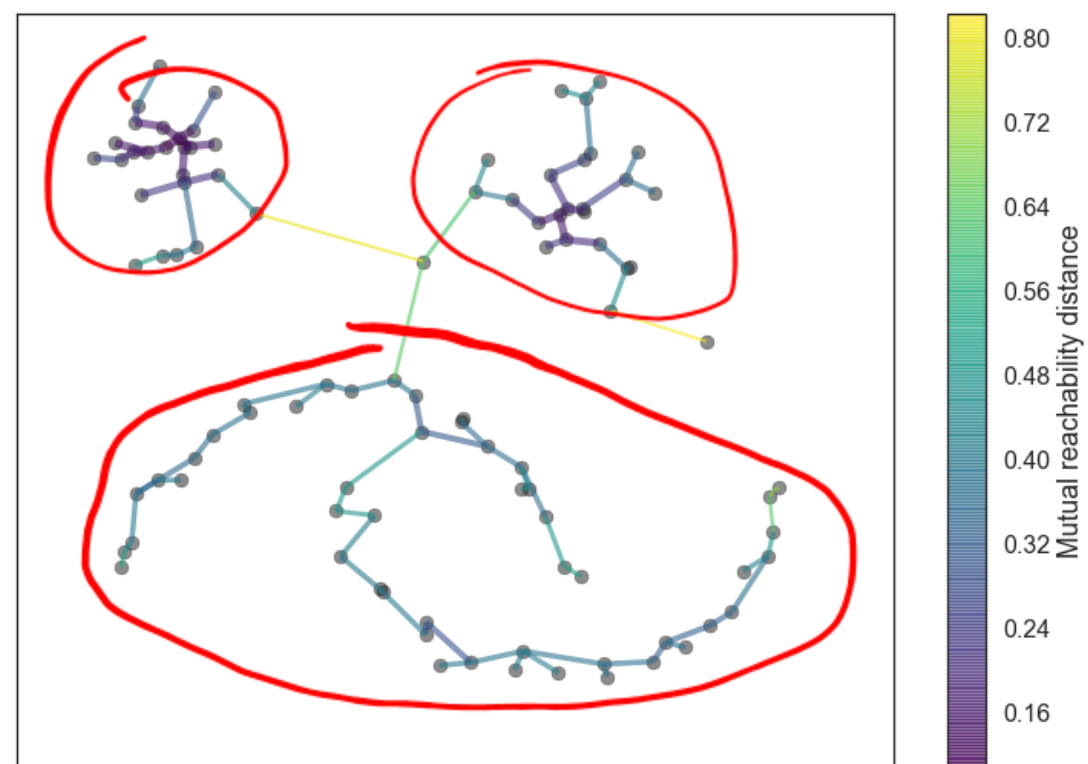
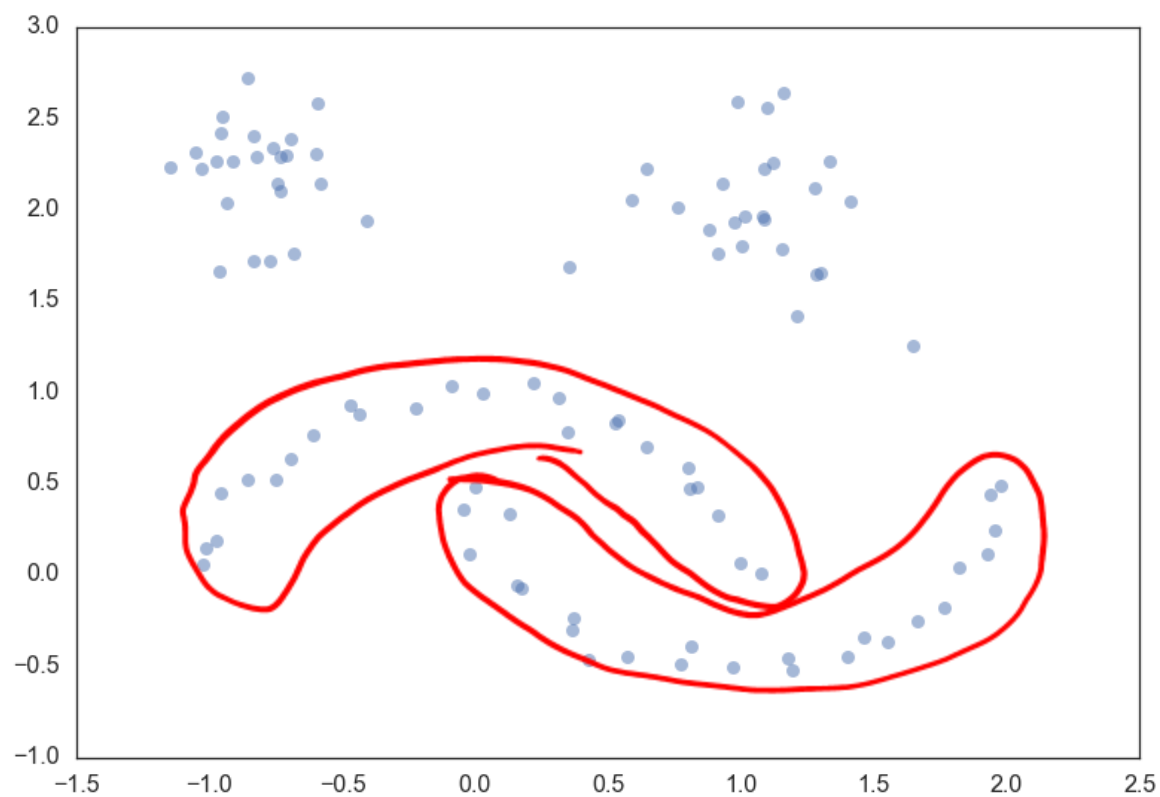
OPTICS

Изменения:

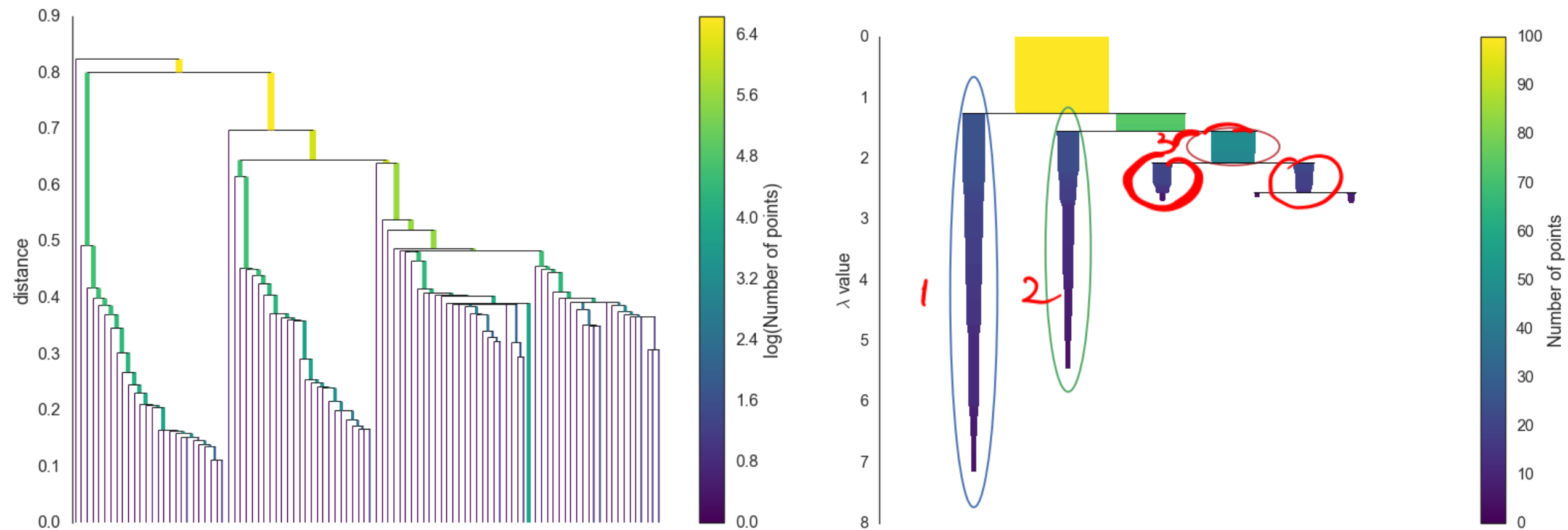
1. График доступности (или дендрограмма) – учитывает иерархию
2. Расстояние Core – ϵ'
3. Расстояние доступности – ϵ

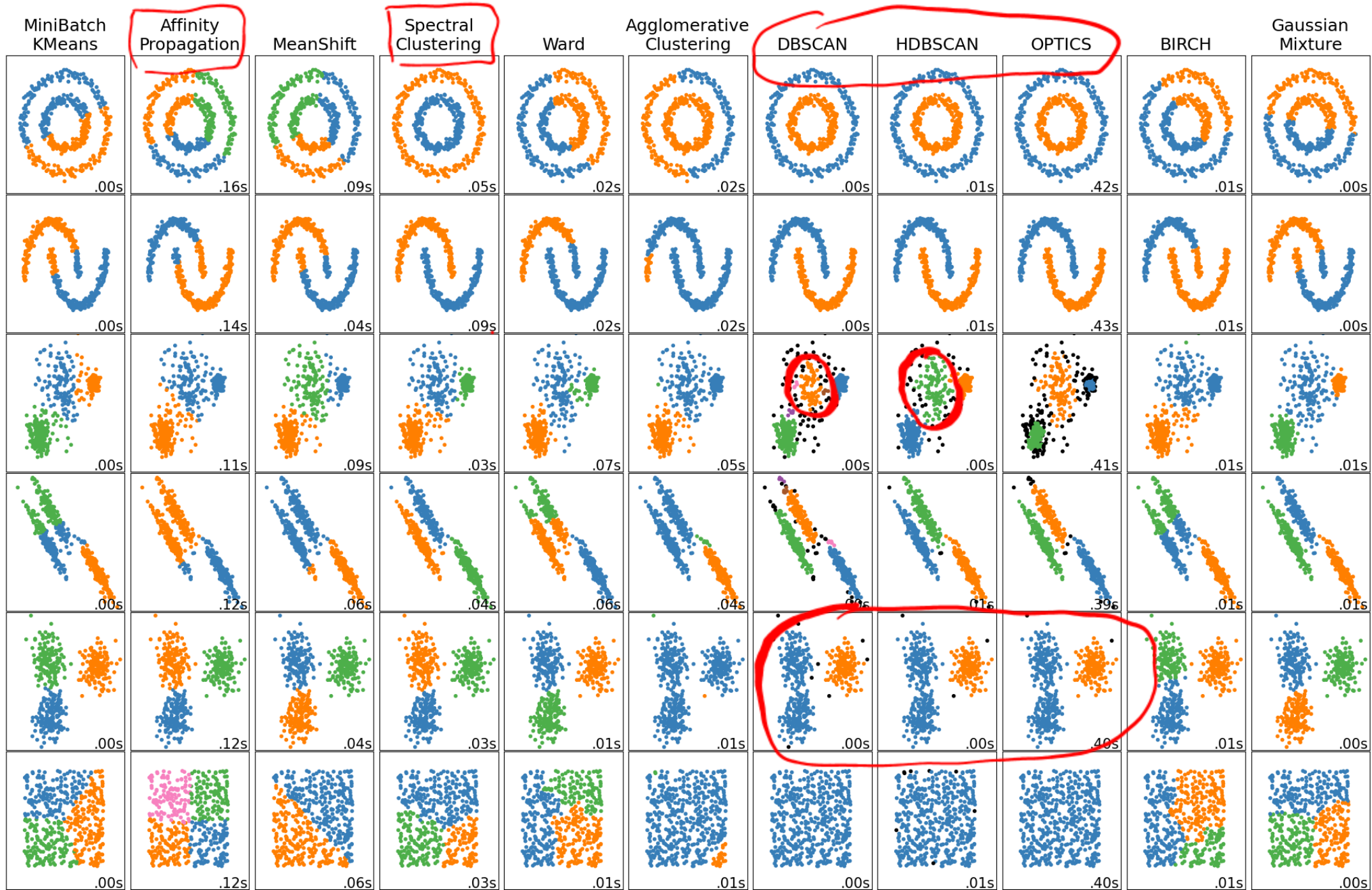


HDBSCAN



HDBSCAN





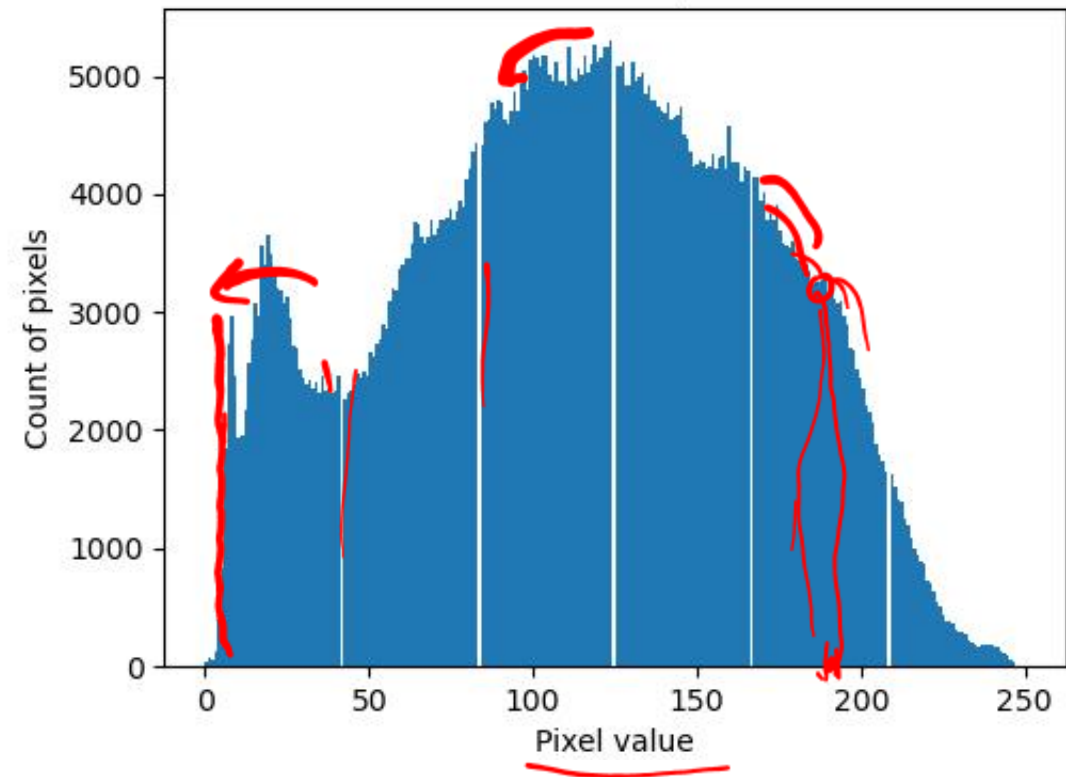
Векторная квантизация

256 8 бит
384 — 9 бит.

Original image of a raccoon face
Rendering of the image



Distribution of the pixel values

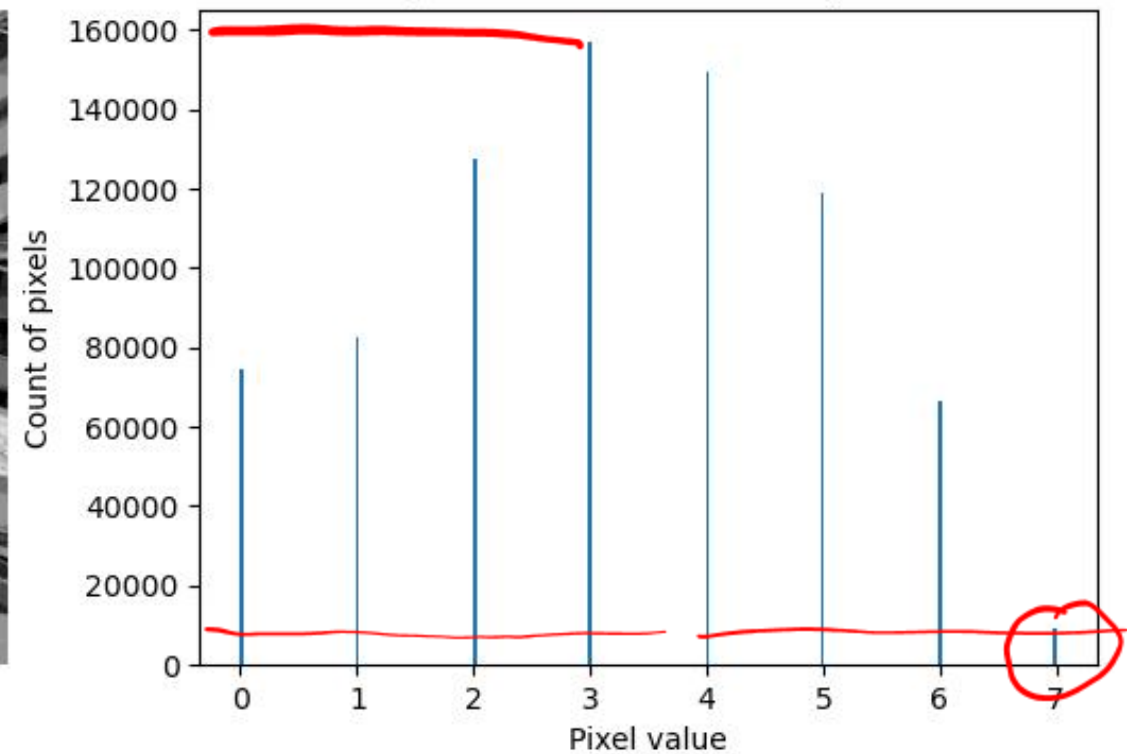


Векторная квантизация

Raccoon face compressed using 3 bits and a uniform strategy
Rendering of the image



Sub-sampled distribution of the pixel values



Векторная квантизация

Raccoon face compressed using 3 bits and a K-means strategy
Rendering of the image



Distribution of the pixel values

