



Лекция 07

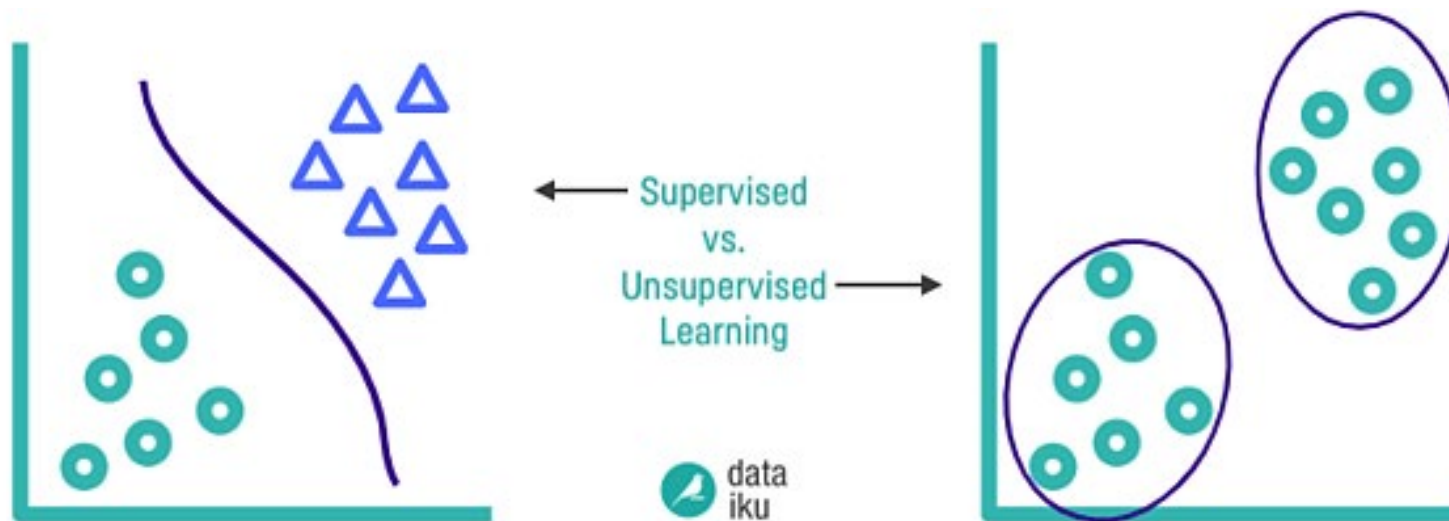
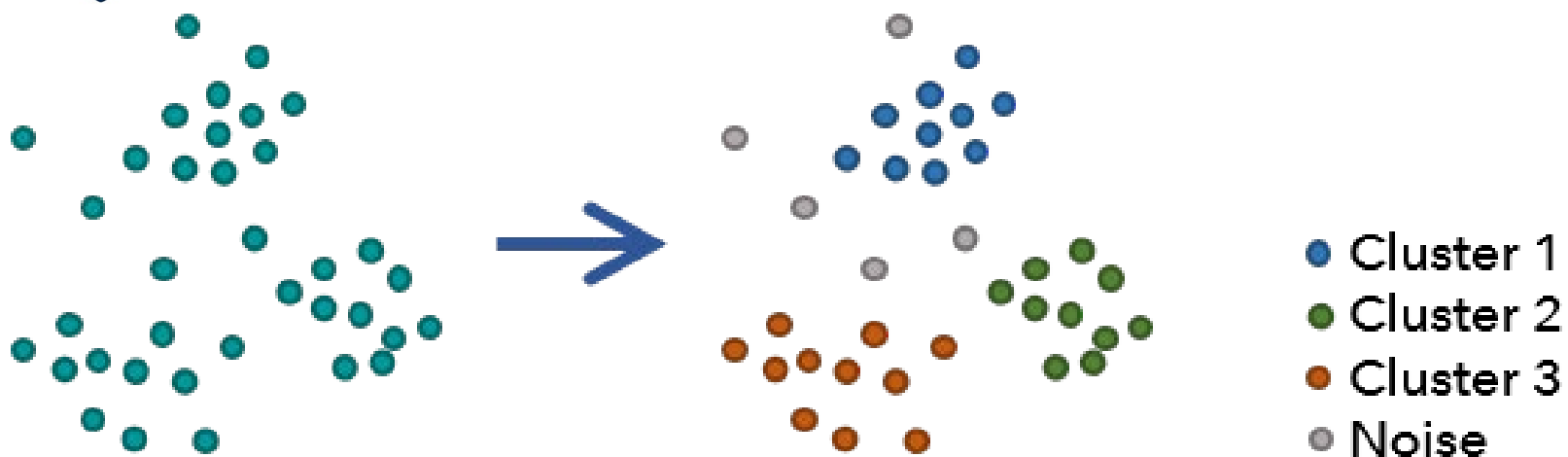
Обучение без учителя (неподконтрольное)

- А. Кластеризация: К-средних, DBSCAN, Иерархическая
- Б. Карта Кохонена и ее варианты
- В. Уменьшение размерности и визуализация: PCA, t-SNE



Кластеризация vs классификация

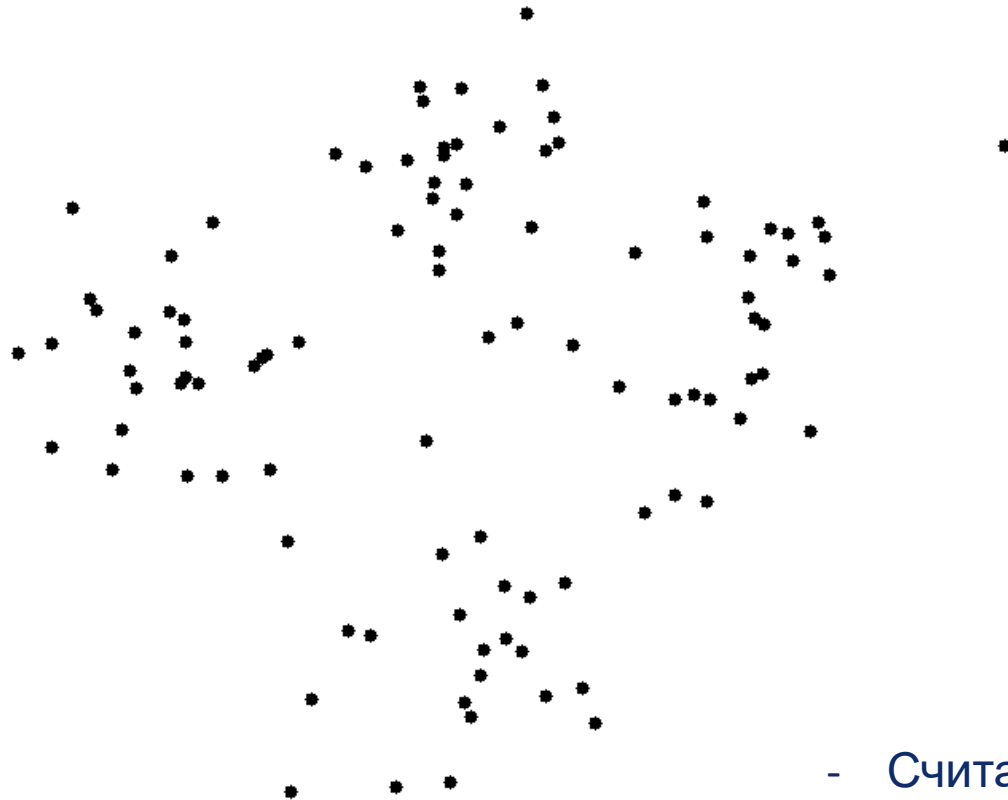
2





Метод К-средних

3



0. Назначить К «центроидов»
1. Посчитать расстояния от «центроидов» до всех точек
2. Распределить точки к ближайшим «центроидам»
3. Посчитать геометрический центр точек, принадлежащих каждому «центроиду»
4. Передвинуть «центроиды» в геометрический центр его точек
5. Повторить с п.1.

- Считать медианы
- Другая метрика расстояния, ускорение расчетов расстояния
- Ввести взаимодействие между центроидами
- Инициализация центроидов



Метод DBSCAN

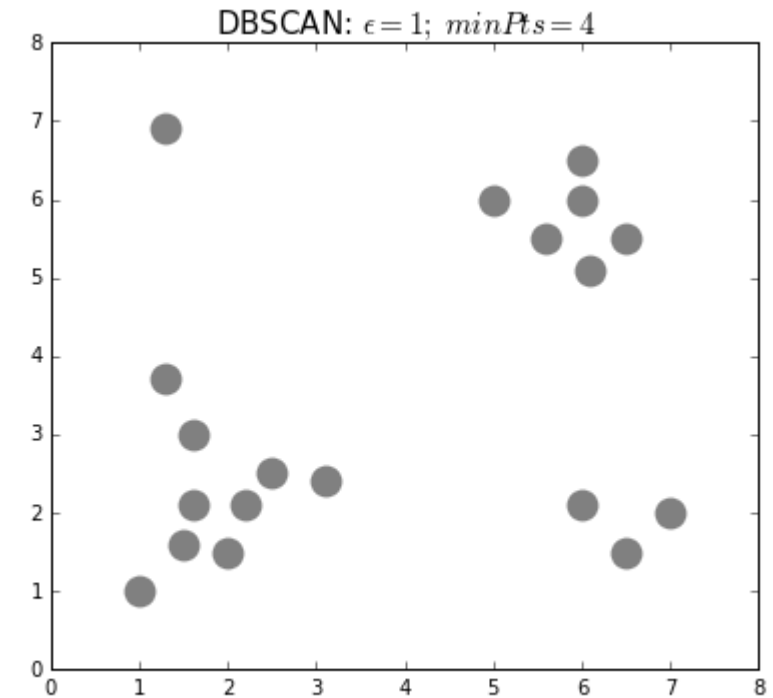
4

```
DBSCAN(DB, distFunc, eps, minPts) {  
    C=0  
    for each point P in database DB {  
        if label(P) ≠ undefined then continue  
        Neighbors N=RangeQuery(DB, distFunc, P, eps)  
        if |N| < minPts then {  
            label(P)=Noise  
            continue  
        }  
        C=C + 1  
        label(P)=C  
        Seed set S=N \ {P}  
        for each point Q in S {  
            if label(Q)=Noise then label(Q)=C  
            if label(Q) ≠ undefined then continue  
            label(Q)=C  
            Neighbors N=RangeQuery(DB, distFunc, Q, eps)  
            if |N| ≥ minPts then {  
                S=S ∪ N  
            }  
        }  
    }  
}
```

/ Счётчик кластеров */*

/ Точка была просмотрена во внутреннем цикле */*
/ Находим соседей */*
/ Проверка плотности */*
/ Помечаем как шум */*

/ следующая метка кластера */*
/ Помечаем начальную точку */*
/ Соседи для расширения */*
/ Обрабатываем каждую зачаточную точку */*
/ Заменяем метку Шум на Край */*
/ Была просмотрена */*
/ Помечаем соседа */*
/ Находим соседей */*
/ Проверяем плотность */*
/ Добавляем соседей в набор зачаточных точек */*

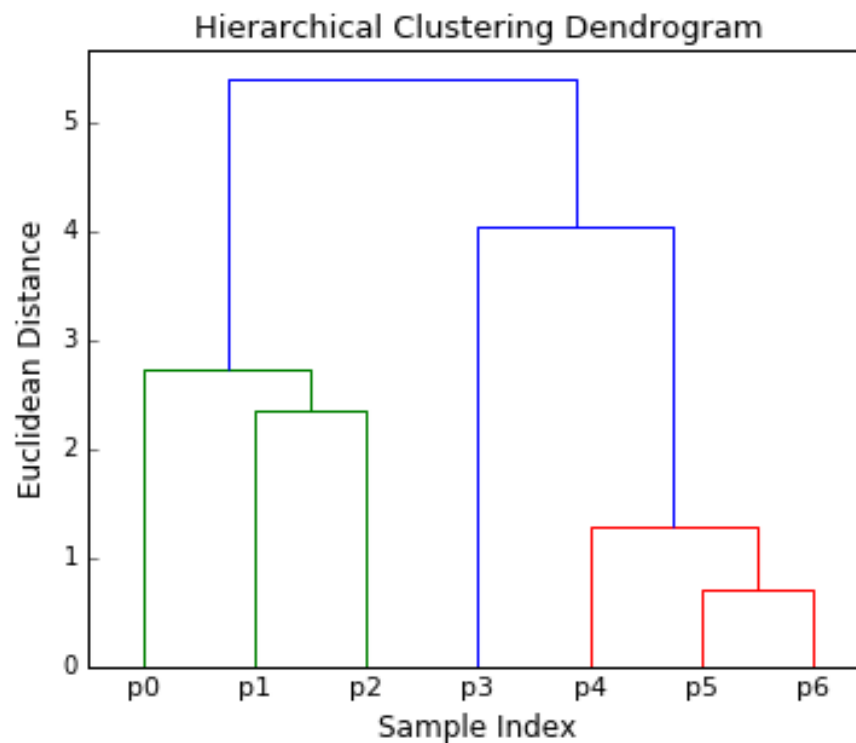
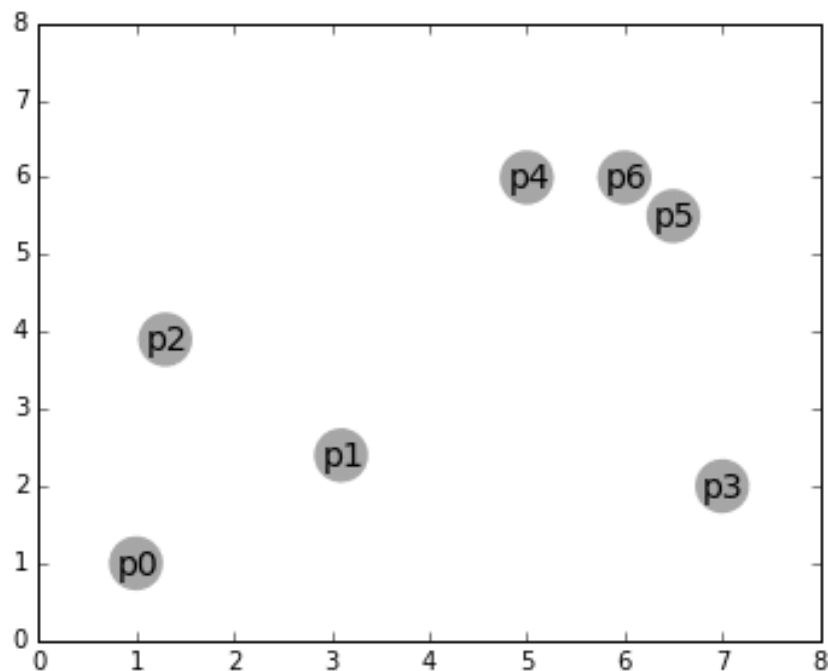


- Другая метрика расстояния, ускорение расчетов расстояния
- Иерархия кластеров по плотности
- Оценка гиперпараметров



Иерархическая кластеризация

5



0. Каждый пример – свой кластер. Считаем расстояние между всеми кластерами.
1. Выбрать ближайшие кластеры.
2. Слить их в один кластер.
3. Удалить слитые из списка кластеров.
4. Добавить новый кластер в список, его координаты.
- Отметить расстояние, на котором произошло слияние.
5. Пересчитать расстояние между этим кластером и другими.
6. Повторить с п.1.

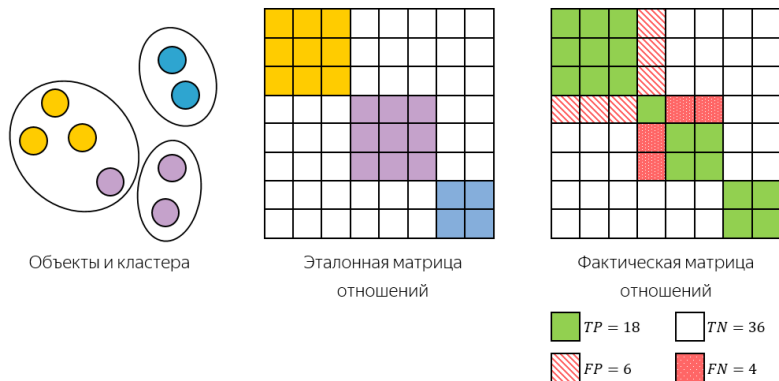
- Другая метрика расстояния, ускорение расчетов расстояния



Метрики кластеризации

ВНЕШНИЕ

Есть эталонные кластеры (классы)
(редкость!)



$$Rand = \frac{TP + TN}{TP + TN + FP + FN}$$

$$Jaccard = \frac{TP}{TP + FN + FP}$$

$$\widehat{ARI} = \frac{\overbrace{\sum_{ij} \binom{n_{ij}}{2}}^{\text{Index}} - \overbrace{[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}]/\binom{n}{2}}^{\text{Expected Index}}}{\underbrace{\frac{1}{2} [\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2}]}_{\text{Max Index}} - \underbrace{[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}]/\binom{n}{2}}_{\text{Expected Index}}}$$

$$E = - \sum_i p_i \left(\sum_j \frac{p_{ij}}{p_i} \log \left(\frac{p_{ij}}{p_i} \right) \right)$$

однородность

$$Q\left(\begin{array}{c} \bullet \bullet \bullet \bullet \\ \bullet \bullet \bullet \bullet \end{array}\right) < Q\left(\begin{array}{c} \bullet \bullet \bullet \bullet \\ \bullet \bullet \bullet \bullet \end{array}\right)$$

$X \backslash Y$	Y_1	Y_2	\dots	Y_s	Sums
X_1	n_{11}	n_{12}	\dots	n_{1s}	a_1
X_2	n_{21}	n_{22}	\dots	n_{2s}	a_2
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
X_r	n_{r1}	n_{r2}	\dots	n_{rs}	a_r
Sums	b_1	b_2	\dots	b_s	n

$$\text{Пусть } p_{ij} = \frac{n_{ij}}{n}, p_i = \frac{a_i}{n}, p_j = \frac{b_j}{n}.$$

- Элементы принадлежат одному кластеру и одному классу — TP
- Элементы принадлежат одному кластеру, но разным классам — FP
- Элементы принадлежат разным кластерам, но одному классу — FN
- Элементы принадлежат разным кластерам и разным классам — TN

Пары примеров. Число перестановок: $\binom{n}{k} = \frac{n!}{k!(n-k)!}$



Метрики кластеризации

ВНУТРЕННИЕ

Нет эталонных кластеров.

Сравниваем кластеры между собой и их центроиды

Компактность

$$WSS = \sum_{j=1}^M \sum_{i=1}^{|C_j|} (x_{ij} - \bar{x}_j)^2$$

$$a(x_i, c_k) = \frac{1}{|c_k|} \sum_{x_j \in c_k} \|x_i - x_j\|$$

Отделимость

$$BSS = n \cdot \sum_{j=1}^M (\bar{x}_j - \bar{x})^2$$

$$b(x_i, c_k) = \min_{c_l \in C \setminus c_k} \left\{ \frac{1}{|c_l|} \sum_{x_j \in c_l} \|x_i - x_j\| \right\}$$

Силуэт

$$Sil(c) = \frac{1}{N} \sum_{c_k \in C} \sum_{x_i \in c_k} \frac{b(x_i, c_k) - a(x_i, c_k)}{\max\{a(x_i, c_k), b(x_i, c_k)\}},$$

Индекс Дэвиса-Болдуина

$$DB(C) = \frac{1}{K} \sum_{c_k \in C} \max_{c_l \in C \setminus c_k} \left\{ \frac{S(c_k) + S(c_l)}{\|\bar{c}_k - \bar{c}_l\|} \right\},$$

где:

$$S(c_k) = \frac{1}{|c_k|} \sum_{x_i \in c_k} \|x_i - \bar{c}_k\|$$



Карта Кохонена

8



0. Назначить K «центроидов»
1. Посчитать расстояния от «центроидов» до всех точек
2. Распределить точки к ближайшим «центроидам»
3. Посчитать геометрический центр точек, принадлежащих каждому «центроиду»
4. Передвинуть «центроиды» в геометрический центр его точек
5. Повторить с п.1.

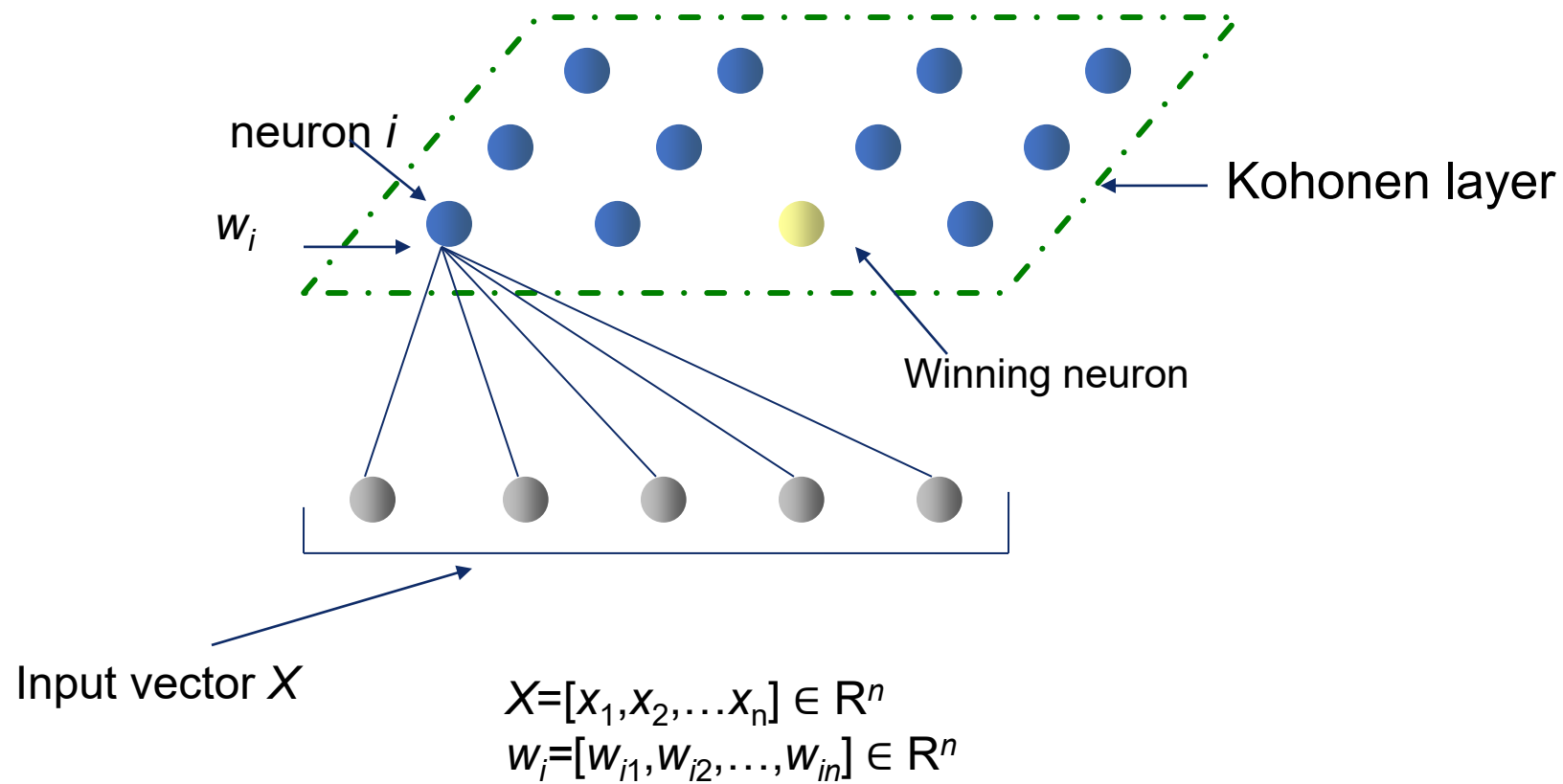
НЕЗАВИСИМЫ

- Считать медианы
- Другая метрика расстояния, ускорение расчетов расстояния
- Ввести взаимодействие между центроидами
- Инициализация центроидов



Карта Кохонена. Структура

9





Карта Кохонена. Обучение

10

- 1) Инициализация центроидов нейронов W_i
- 2) Входной m -мерный вектор X_s подается в карту
- 3) Вычисляем расстояния до всех нейронов $di(W_i, X_s)$:

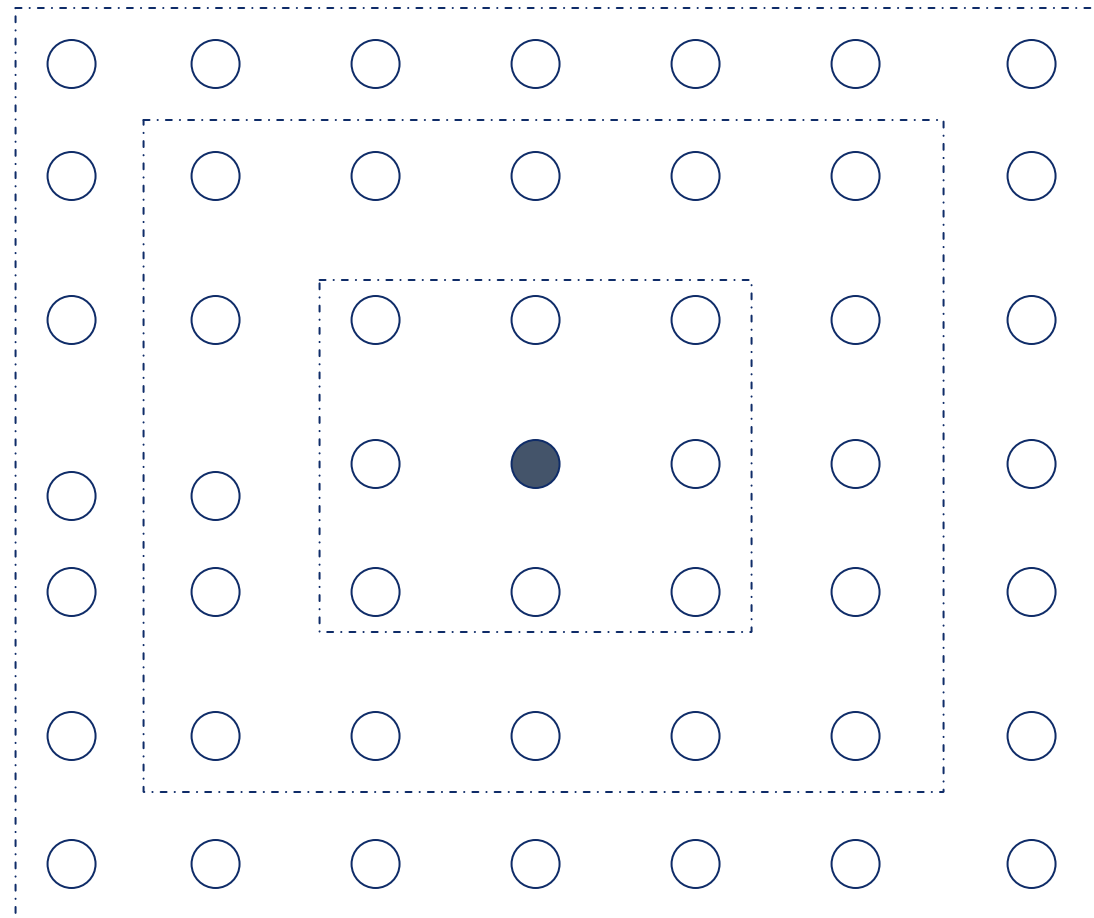
$$d_i(W_i, X_s) = \sum_{j=1}^m (w_j - x_j)^2$$

- 4) Ищем ближайший нейрон “победитель”
- 5) Сдвигаем нейрон-победитель и его соседей в направлении входа:

$$W_{j_{new}} = W_{j_{old}} + \alpha(X_i - W_{j_{old}})$$

α – скорость движения, у победителя больше, у соседей меньше

- 6) Повторяем с новым входом





Карта Кохонена. Применения и варианты

11

- Уменьшение размерности и Визуализация
- Кластеризация (требуется доп. проверка)
- Классификация (аналогично K-средним)

Варианты и модификации:

- Одномерные кольца – задача коммивояжера
- Растущий нейронный газ – добавлены механизмы изменения связей и нейронов – изменение сетки и построение метакластеров

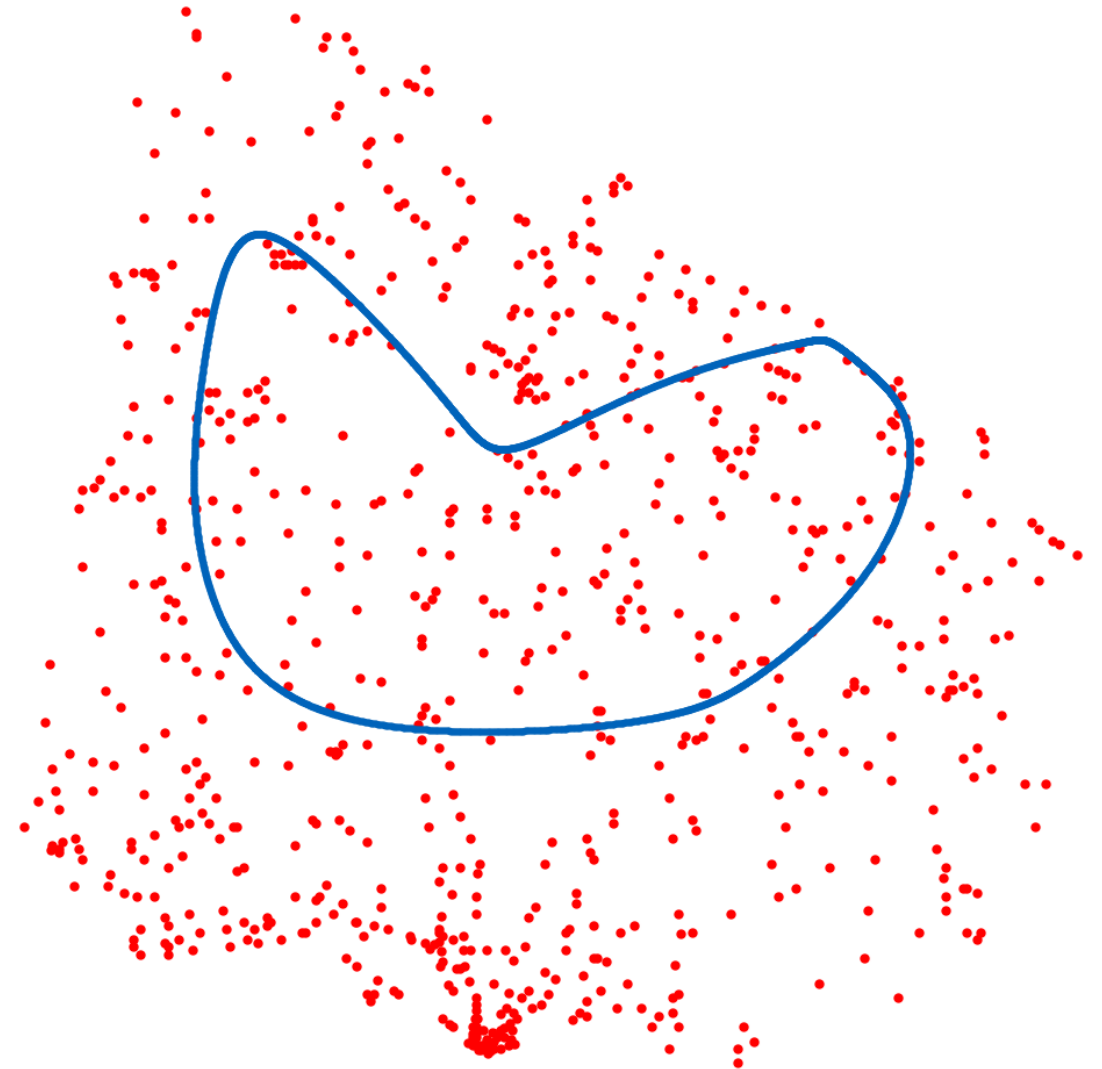
<https://www.demogng.de/>



Маршрутизация и коммивояжер

12

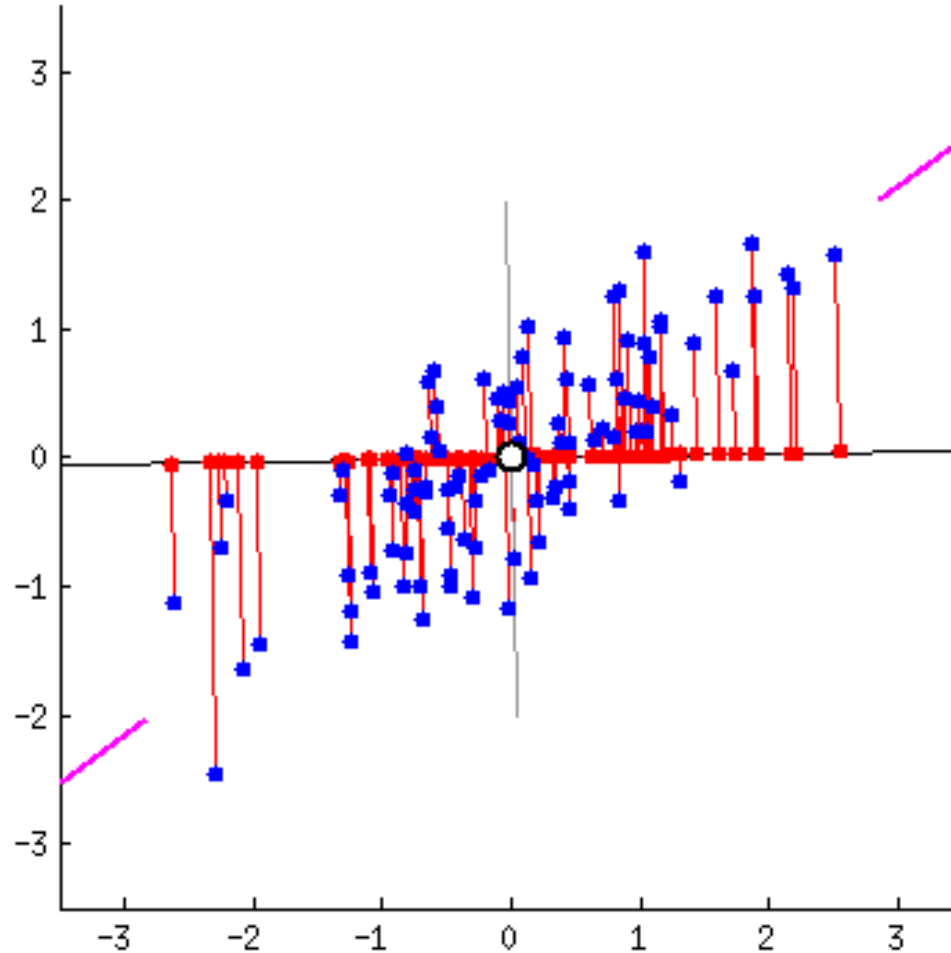
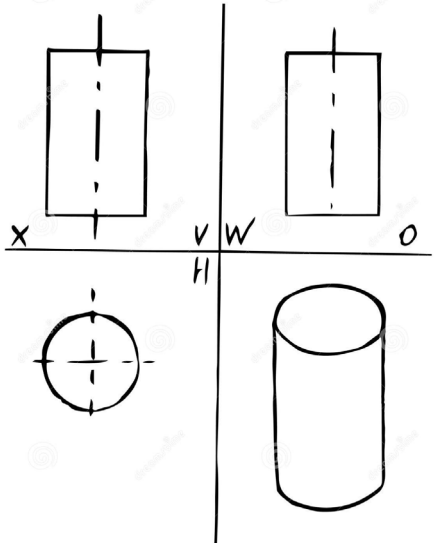
- города это входные данные (точки)
- нейроны замкнуты в **кольцо**
- механизмы удаления нейронс (небыли победителями)
- механизмы дублирования нейронс (победители более 1 раза)





Уменьшение размерности и визуализация. Метод главных компонент

13



$$X_{\text{scaled}} = \frac{X - \mu}{\sigma}$$

$$C = \frac{1}{m - 1} X_{\text{scaled}}^T X_{\text{scaled}}$$

$$Cv = \lambda v$$

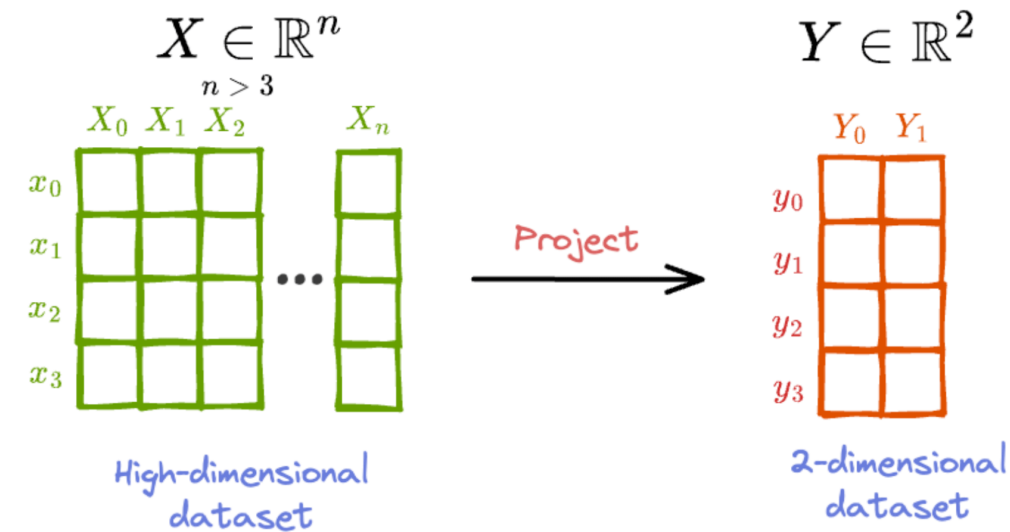
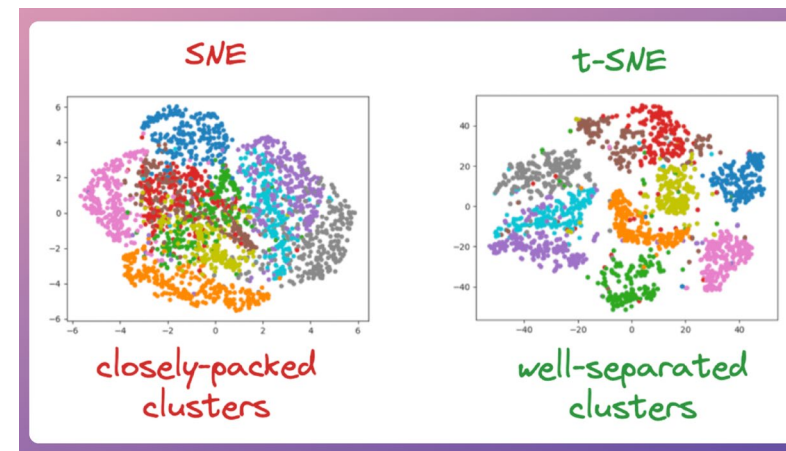
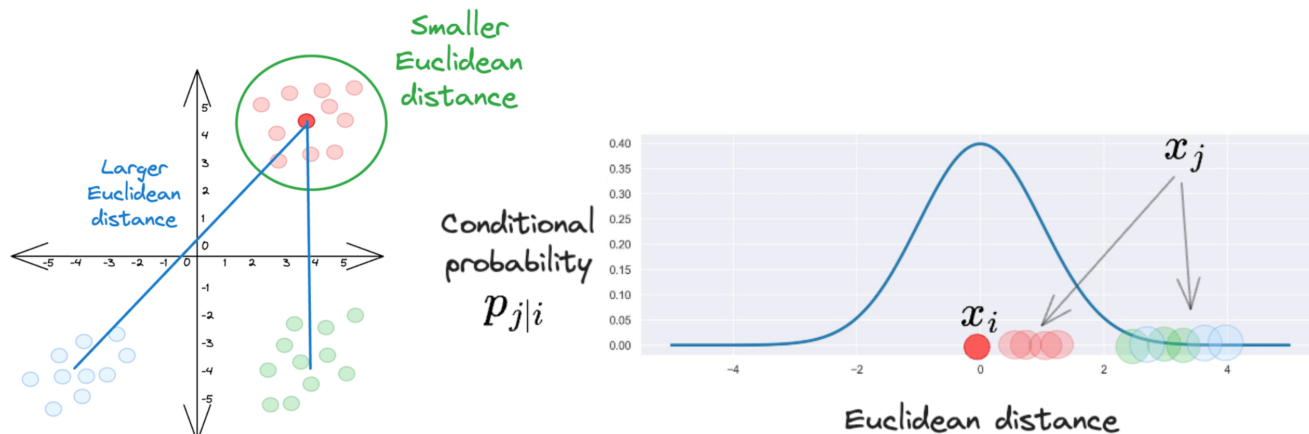
$$W = [v_1, v_2, \dots, v_k]$$

$$X_{\text{PCA}} = X_{\text{scaled}} W$$



Уменьшение размерности и визуализация. SNE, t-SNE

14



$$D_{KL}(P \parallel Q) = \sum_x P(x) \cdot \log \left(\frac{P(x)}{Q(x)} \right)$$



Группа по дисциплине:

<https://t.me/+8dShF1tFSDg0ZmJi>

