## MO с учителем

$$X = \{(x_i, y_i)\}_{i=1}^{\ell}$$ — обучающая выборка

$$a(x_i) \simeq y_i \qquad\qquad Q(a, X) \Rightarrow \min_{a \in A}$$

## MO без учителя    Unsupervised learning.

$$X = \{x_i\}_{i=1}^{\ell}$$    нужно что-то сделать с $X$

## Кластеризация

$$X = \{x_i\}_{i=1}^{\ell} \qquad\text{хотим } a: X \Rightarrow \{1, \dots, k\}$$
модель кластеризации
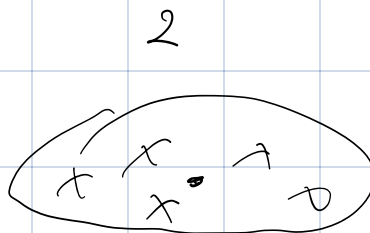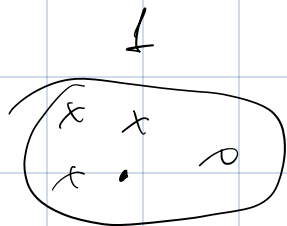
$$x_i \text{ и } x_j \text{ похожи} \iff a(x_i) = a(x_j)$$

## Зачем

1) Исследование данных

- социологический опрос
- телефонный оператор
- кластеризация текстов

2) Генерация новых признаков

- добавить категориальный признак
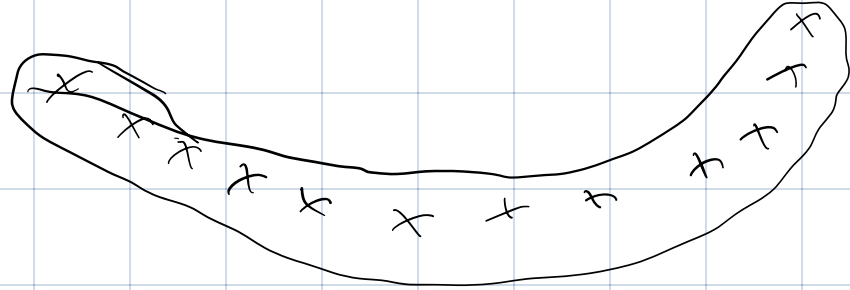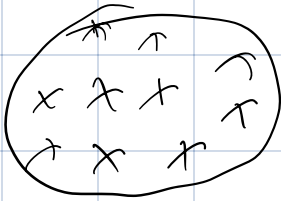- добавить расстояние 1..... к.
  до центра кластеров



## Метрики качества кластеризации

а) Посмотреть глазами

б) Внутрикластерное расстояние
   требует "компактность кластеров"

$$\sum_{k=1}^{K}\sum_{i=1}^{\ell}\left[a(x_i)=k\right]\rho(x_i,\,c_k)\Rightarrow min$$

$c_k$ — центр $k$-го кластера



2) Межкластерное расстояние

Требуем чтобы кластеры были "удалены

друг от друга"

$$\sum_{i\neq j}\left[a(x_i)\neq a(x_j)\right]\rho(x_i,\,x_j)\Rightarrow max$$
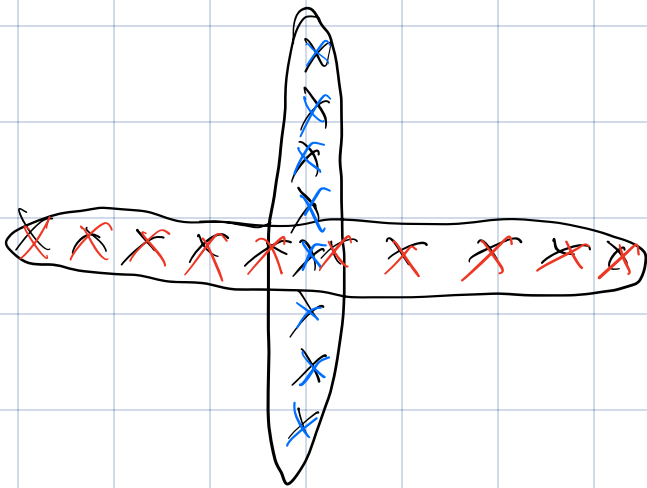


3) Комбинируем подходы

Индекс Данна

$d(k,k')$ — расстояние между

кластерами $k$ и $k'$

$d(k)$ - внутрикластерное расстояние

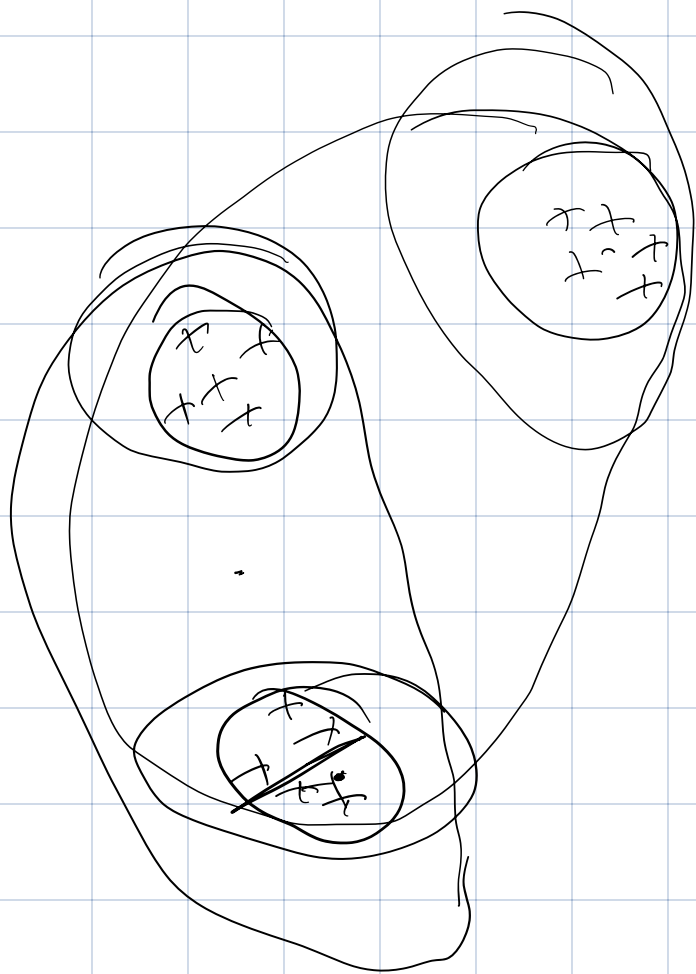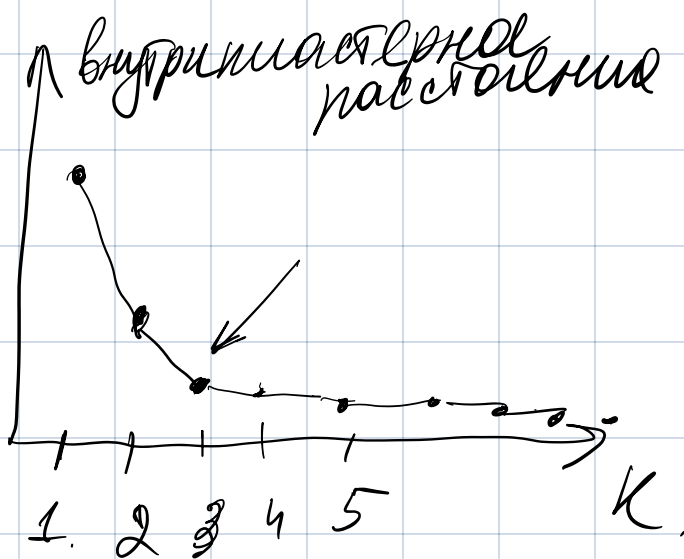$$\frac{\min\limits_{1\leq k \leq k' \leq K} d(k, k')}{\max\limits_{1\leq k \leq K} d(k)} \to max$$



- метод выбора $k$ - кол-во кластеров реализован внутри алгоритма кластеризации

o  $K$ - гиперпараметр.

# Как подобрать $K$?

1) Смотреть глазами

2) elbol method



внутрикластерное расстояние

1.  2  3  4  5          $K$

---

## K-means.

$\rho(x, z)$ - какое-то расстояние

$K$ - гиперпараметр, число кластеров

$$\sum_{k=1}^{K} \sum_{i=1}^{ } \left[ a(x_i) = k \right] \rho(x_i, c_k) \rightarrow \min_{a,\, c_k}$$

$c_k$ — центр кластера $k$.

Шаг 0  Инициализируем центры
$$c_k$$

Шаг (а)  Фиксируем $c_k$
$$a(x_i) = \arg\min \rho(x_i, c_k)$$
$$k = 1, \dots K.$$

Шаг (б)  Фиксируем $a(x_i)$

Ищем центры масс

$$c_k = \arg\min_{c} \sum_{a(x_i) = k} \rho(x_i, c)$$

если $\rho_{(x, c_k)} = \| x - c_k \|^2$

$$c_k = \frac{\sum_{i=1}^{\ell} \left[ a(x_i) = k \right] \cdot x_i}{\sum_{i=1}^{\ell} \left[ a(x_i) = k \right]}$$