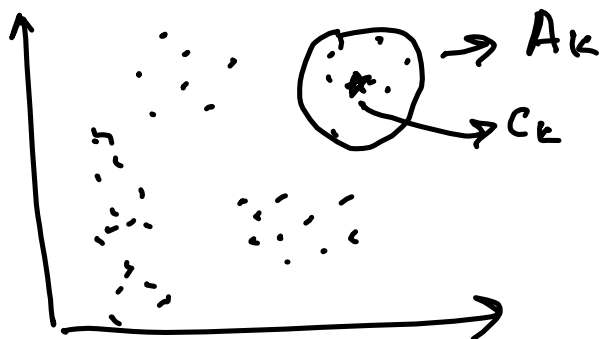# Кластеризация



# k-means

$$Q = \sum_{k=1}^{k} \sum_{i: x_i \in A_k} \| x_i - c_k \|^2 \to \min_{c_1, \ldots, c_k}$$

$$A_k = \{ x_i \mid \rho(x_i, c_k) \leq \rho(x_i, c_j) \; \forall j \neq k \}$$
$$\forall j \neq k$$

$$i = 1, \ldots, n$$
$$j = 1, \ldots, k$$

$$Q \to \min_{c_k}$$

$$\sum_{i: x_i \in A_k} \| x_i - c_k \|^2 \to \min_{c_k}$$

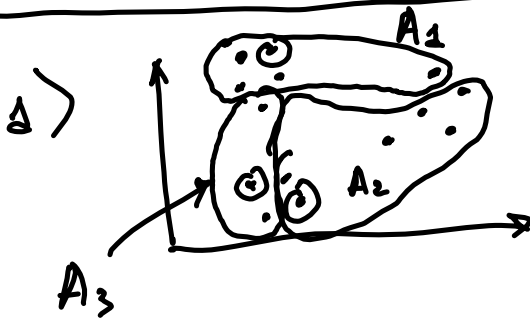$$\|x_i - c_k\|^2 = (x_i^T - c_k^T)(x_i - c_k) =$$

$$= (x_i^T x_i - 2 x_i^T c_k + c_k^T c_k)$$

$$\nabla_{c_k} = -2x_i + 2c_k$$

$$\nabla_{c_k} \sum_{i: x_i \in A_k} \|x_i - c_k\|^2 = \sum_{i: x_i \in A_k} (-2x_i + 2c_k) = 0$$

$$-\sum_{i: x_i \in A_k} x_i + |A_k| \cdot c_k = 0$$

$$\boxed{c_k = \frac{1}{|A_k|} \cdot \sum_{i: x_i \in A_k} x_i}$$



а)

$A_1$

$A_2$

$A_3$

2) Пересчитываем центры.

Критерии останова:

1) Число набл. (доля), сменивших кластер.

2) $\| c_k^{old} - c_k^{new} \| < \varepsilon$

3) Изменение $Q$

$$\left| \frac{Q_{new} - Q_{old}}{Q_{old}} \right|$$

---

0) $\overline{\rho}_{внутр.} > \overline{\rho}_{межкл.}$

1) Индекс. Данна = $\dfrac{\overline{\rho}_{внутр.}}{\overline{\rho}_{межкл.}} \to min$
   Duna Index

2) Silouette$_i = \dfrac{-(d(x_i, c_k) - d(x_i, c_j))_{j \neq k}}{d(x_i, c_k)_{j \neq k}}$

$x_i: \quad d(x_i, c_k) < d(x_i, c_{k+1})$

TSNE

$$x_i \in \mathbb{R}^D \longrightarrow y_i \in \mathbb{R}^d$$

$$\sum_{i=1}^{n}\sum_{j=1}^{n}\left(\rho(x_i,x_j)-\rho(y_i,y_j)\right)^2 \to \min_{y_1,\dots,y_n}$$

MDS (crossed out)

$$\frac{\rho(x_i,x_j)}{\rho(x_i,x_k)}=\alpha \approx \frac{\rho(y_i,y_j)}{\rho(y_i,y_k)}$$

$$P_{j\mid i}=\frac{\exp\left(-\dfrac{\|x_i-x_j\|^2}{2\delta^2}\right)}{\sum_{k\neq j}\exp\left(-\dfrac{\|x_i-x_k\|^2}{2\delta^2}\right)}$$



$$P_{ij}=\frac{P_{i\mid j}+P_{j\mid i}}{2n}$$

$$q_{ij}=\frac{q_{i\mid j}+q_{j\mid i}}{2n}$$

$$q_{j\mid i}=\frac{\left(1+\|y_i-y_j\|^2\right)^{-1}}{\sum_{k\neq j}\left(1+\|y_i-y_k\|^2\right)^{-1}}$$

$q_{ij} \approx p_{ij}$



$$\# \quad KL(p \| q) = \int_{\mathbb{R}} p(x) \log \frac{p(x)}{q(x)} \, dx =$$

$$= \mathbb{E}_p\left(\log \frac{p}{q}\right)$$
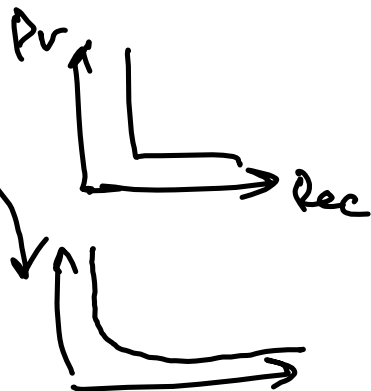


$$KL(N_1 \| N_2)$$

$$N_1 = N(a_1, \delta^2)$$

$$N_2 = N(a_2, \delta^2)$$

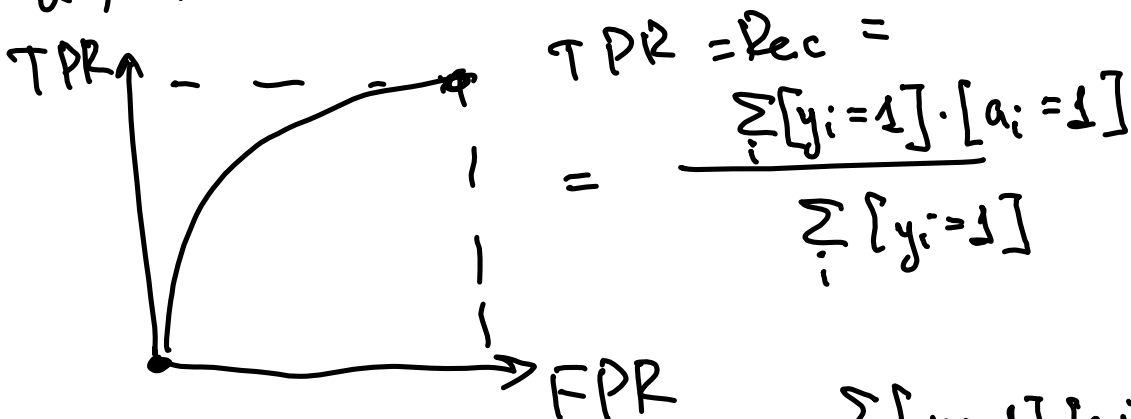$$Q = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}} \to \min_{y_1, \ldots, y_n}$$

---

Семинар:

1) $f_1 = \dfrac{2 Pr \cdot Rec}{Pr + Rec}$

$\max(\min(Pr, Rec))$



2) $AUC\text{-}ROC = TPR(FPR)$



$TPR = Rec = $
$$= \frac{\sum_i [y_i = 1] \cdot [a_i = 1]}{\sum_i [y_i = 1]}$$

$$FPR = \frac{\sum_i [y_i = -1] \cdot [a_i = 1]}{\sum_i [y_i = -1]}$$

$$AUC\text{-}ROC = \frac{\sum_i \sum_j [y_i < y_j] \cdot [a_i < a_j]}{\sum_i \sum_j [y_i < y_j]} = n_+ \, n_-$$

3) $\mathbb{E}(AUC\text{-}ROC) = 0.5$

4)
$$\begin{cases} a_k^T \underset{D \times D}{S} a_k \to \underset{a_k}{max} \\ \|a_k\|^2 = 1 \end{cases}$$

$$S = \frac{1}{n} \sum_{i=1}^{n} x_i \cdot x_i^T \qquad \begin{array}{l} x_i \in \mathbb{R}^D \\ a_k \in \mathbb{R}^D \end{array}$$

5) $x_i \in \mathbb{R}^D$, $A = \underset{D \times d}{\left[ \overset{|}{\underset{|}{a_1}} \dots \overset{|}{\underset{|}{a_d}} \right]}$

$$\underset{\mathbb{R}^d}{\overset{\curvearrowright}{\hat{x}_i}} = \underbrace{\underset{d \times D}{(A^T A)^{-1}} \underset{D \times 1}{A^T x_i}}_{I_d} = A^T x_i$$