

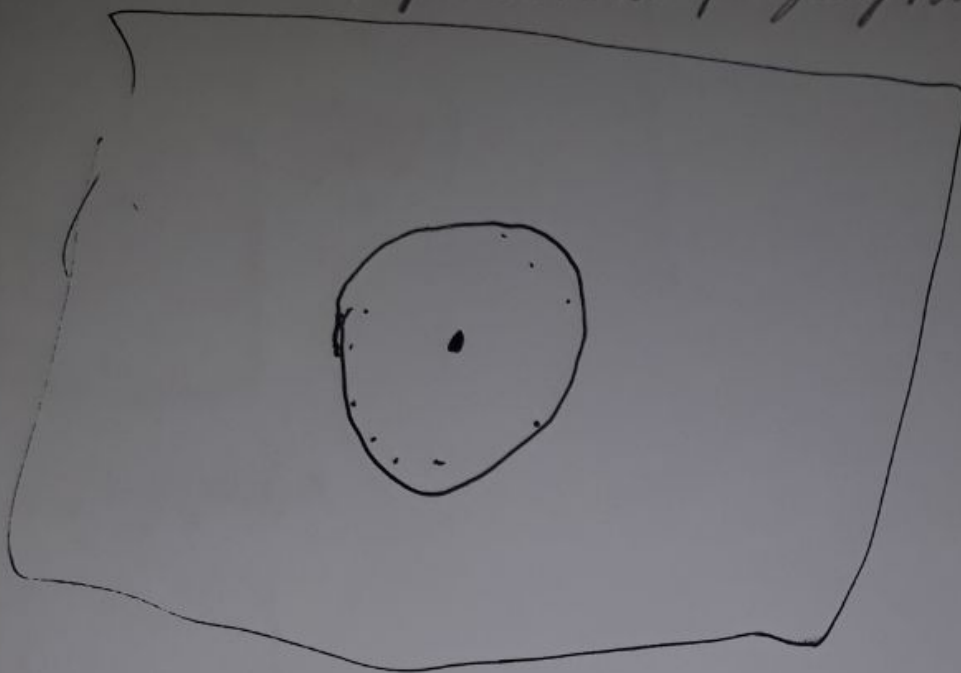


- +
- 1) Много классов, мало представлений
 - 2) Не обязательно выделять признаки
 - 3) Надстройка над более сложными моделями

-
- 1) Нужна вся выборка
 - 2) Проклятые размерности
 - 3) Далеко работает

± Чувствителен к шумам

Прямые размерности.

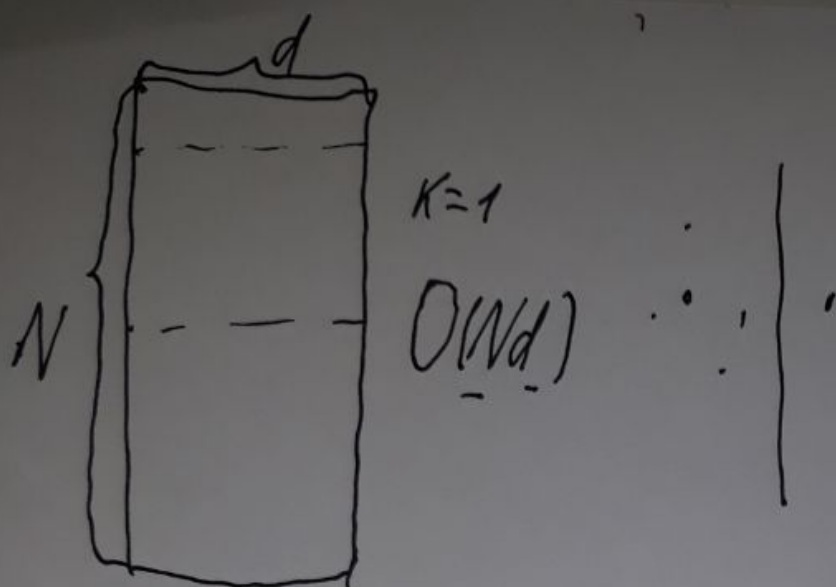


$$V_{\text{шар}}^n = \pi R^2, \frac{4}{3} \pi R^3$$

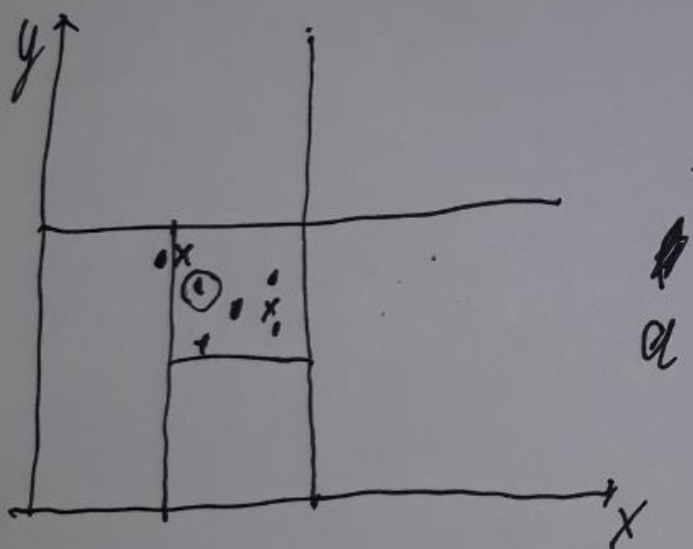
$$C R^n$$

$$R - \varepsilon, \varepsilon > 0.$$

$$\frac{C_1 R - \varepsilon^n}{C_2 R^n} = C_3 \left(\frac{R - \varepsilon}{R} \right)^n \xrightarrow{n \rightarrow \infty} 0$$



KD-геоборд



$d=10-20$ - распространение
время

$d \gg 20$ - гено

Locality sensitive hashing (LSH)

~~$F(x_1, x_2) = 0$~~

$$x_1, x_2 \in X$$

$$F(x_1) = F(x_2)$$

x_1 и x_2 близки

$F(x_1) \neq F(x_2)$ далеко.

Опр: семейство F называется (d_1, d_2, p_1, p_2) -чувствительным если $\forall x_1, x_2 \in X$:

если $p(x_1, x_2) \leq d_1$, то

$$p(f(x_1) = f(x_2)) \geq p_1.$$

если $p(x_1, x_2) \geq d_2$ то

$$p(f(x_1) = f(x_2)) < p_2$$

$$p_1 \geq p_2 \quad d_1 \leq d_2.$$

Метрика Джаммаса
(Jaccard)

~~В~~ A, B - множества

$$P_J(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$



$$U = \{u_1, u_2, \dots, u_n\}$$

$$A, B \subset U$$

$$\pi = \{u_{k_1}, u_{k_2}, u_{k_3}, u_{k_n}\}$$

$$f_\pi(A) = \min \{ \pi(u_i) \mid u_i \in A \}$$

$$F = \{f_\pi \mid \pi \in \text{ПЕРЕСТ}(U)\}$$

Умб: Вероятность что $\forall f_\pi \in F$
будет одинаковой на A, B
равно $\frac{|A \cap B|}{|A \cup B|} = 1 - P_J(A, B)$

1) $u \in A, u \in B$ - р объектов

2) $u \in A, u \notin B$ или $u \notin A, u \in B$. - q объектов

3) $u \notin A, u \notin B$

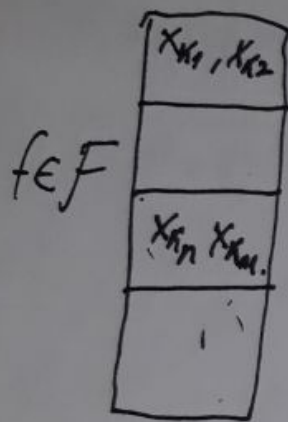
$$\left(\frac{p}{p+q} \right) = 1 - P_J(A, B)$$

$$P_j(A, \theta) \leq d_1.$$

$$1 - P_j(A, \theta) \geq 1 - d_1.$$

$$P_1 > 1 - d_1. \quad P_2 < 1 - d_2.$$

$(d_1, d_2, 1 - d_1, 1 - d_2)$ - zufallsum.



~~100~~ $K_{\text{neozest}} = 5$
 100% genau 4
 umgeklammert?

$$g(x) = (f_1(x), f_2(x), \dots, f_m(x)).$$

$$f_1, \dots, f_m \in F$$

$$g_1(x)$$

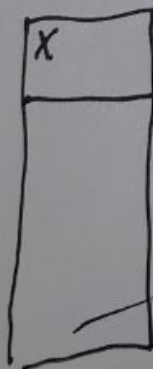
$$g_2(x)$$

\vdots

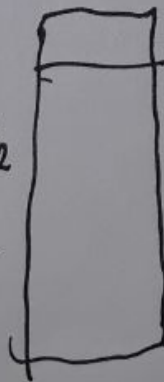
$$g_L(x)$$

$$m, L$$

g_i



g_2



Вспомогательная функция для алгоритма

$$x_1 = 0.1 \quad y_1 = 1$$

$$x_2 = 0.5 \quad y_2 = 2$$

$$u = 0 \quad (0, 0)$$

Вспомогательная функция $\sim U[0, 1]$

$$P(p(u, x_1) > p(u, x_2)) =$$

$$P(\sqrt{0.1^2 + \xi_1^2} > \sqrt{0.5^2 + \xi_2^2}) =$$

$$= P(0.1^2 + \xi_1^2 > 0.5^2 + \xi_2^2) =$$

$$= P(\xi_1^2 \geq 0.24 + \xi_2^2) =$$

$$= \int_0^1 \int_{\sqrt{x_2^2 + 0.24}}^1 dx_1 dx_2 \approx 0.275$$