

KNN

$$\underline{f}: X \times X \rightarrow \underline{[0, +\infty)}$$

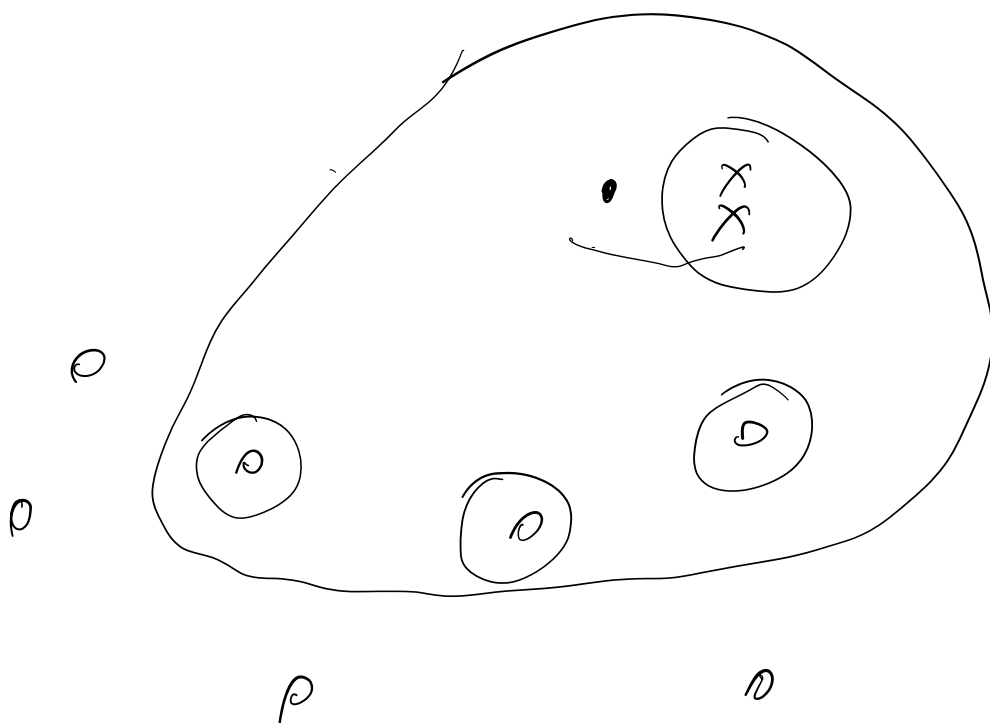
Обучение: Запоминаем X

Прерсказание:

u - новый объект

$$\underline{f(u, x_1) \leq f(u, x_2) \dots \leq f(u, x_c)}$$

$$a(u) = \underset{y \in Y}{\operatorname{argmax}} \sum_{k=1}^K \underbrace{w(k, u, x_{(k)})}_{[y_{(k)} = y]}$$



$$w(k, y, x_k) = \frac{K \left(\frac{P(y, x_k)}{K} \right)}{\frac{1}{P(y, x_k)}}$$

пересчет:

$$a(y) = \frac{\sum_{k=1}^K \frac{K \left(\frac{P(y, x_k)}{K} \right)}{1}}{\sum_{k=1}^K \frac{K \left(\frac{P(y, x_k)}{K} \right)}{1}}$$

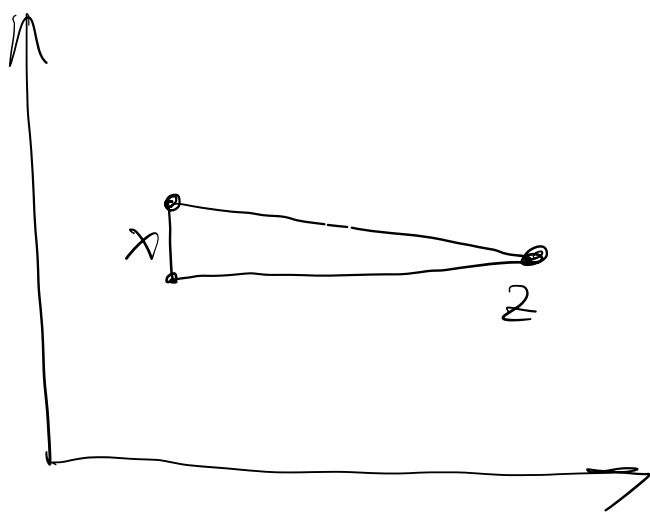
Зачем?

- 1) Легко задать расстояния, но сложно придумать признаки
- 2) если объектов одного класса мало

Метрика Минковского

$$\rho(x, z) = \left(\sum_{j=1}^d |x_j - z_j|^p \right)^{1/p}$$

1) $p=2$ Евклидова метрика



2) $p=1$

Манхэттенское
расстояние

$$3) p=\infty \quad \rho_{\infty}(x, z) = \max_{j=1, 2, \dots, d} |x_j - z_j|$$

метрика Чебышева

x x x x x x x

5) Косинусная метрика

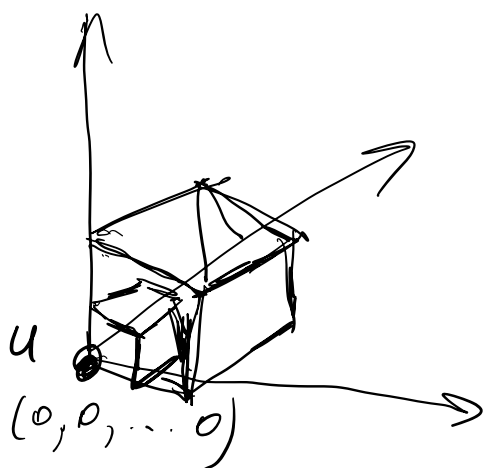
$$\underline{\langle x, z \rangle} = \|x\| \cdot \|z\| \cdot \cos \theta$$

$$\rho_{\cos}(x, y) = \arccos \left(\frac{\langle x, y \rangle}{\|x\| \|y\|} \right)$$

6) Расстояние Жаккарда

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

Проклятые размерности



$$[0, 1]^d$$

5000 объектов

5 объектов

$$[0, \varepsilon]^d \subset [0, 1]^d \quad \varepsilon \in (0, 1]$$

$$\delta = \varepsilon^d \quad \text{найти } \delta$$

при котором в подкубе $[0, \varepsilon]^d$ окажется
хотим δ 5 объектов с вер-тью больше
0.95.

$$\min \delta : \sum_{k=5}^{5000} \binom{5000}{k} \delta^k (1-\delta)^{5000-k} \geq 0.95$$

$$\delta = 0.0018$$

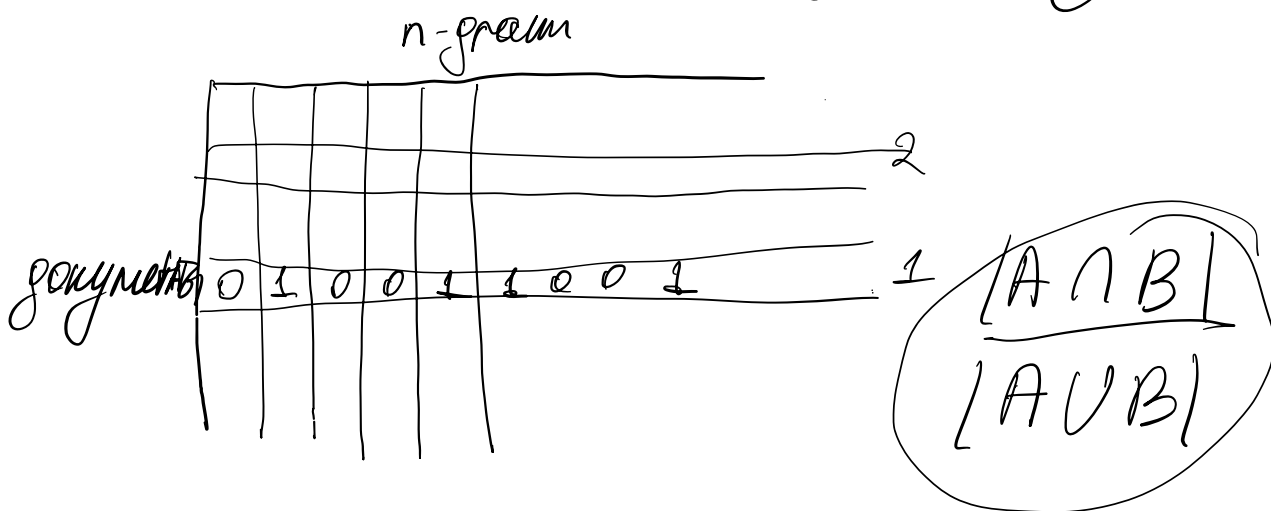
$$d = 10, \quad \varepsilon = 0.53$$

$$d = 100 \quad \varepsilon = 0.94.$$

LSH

Min Hash LSH.

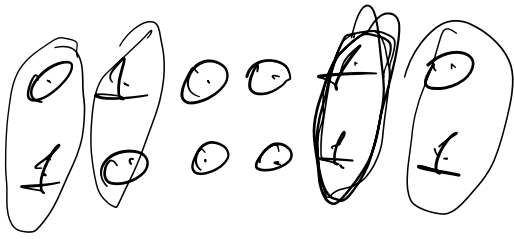
Задача: Поиск дубликатов документов



Bag of Words

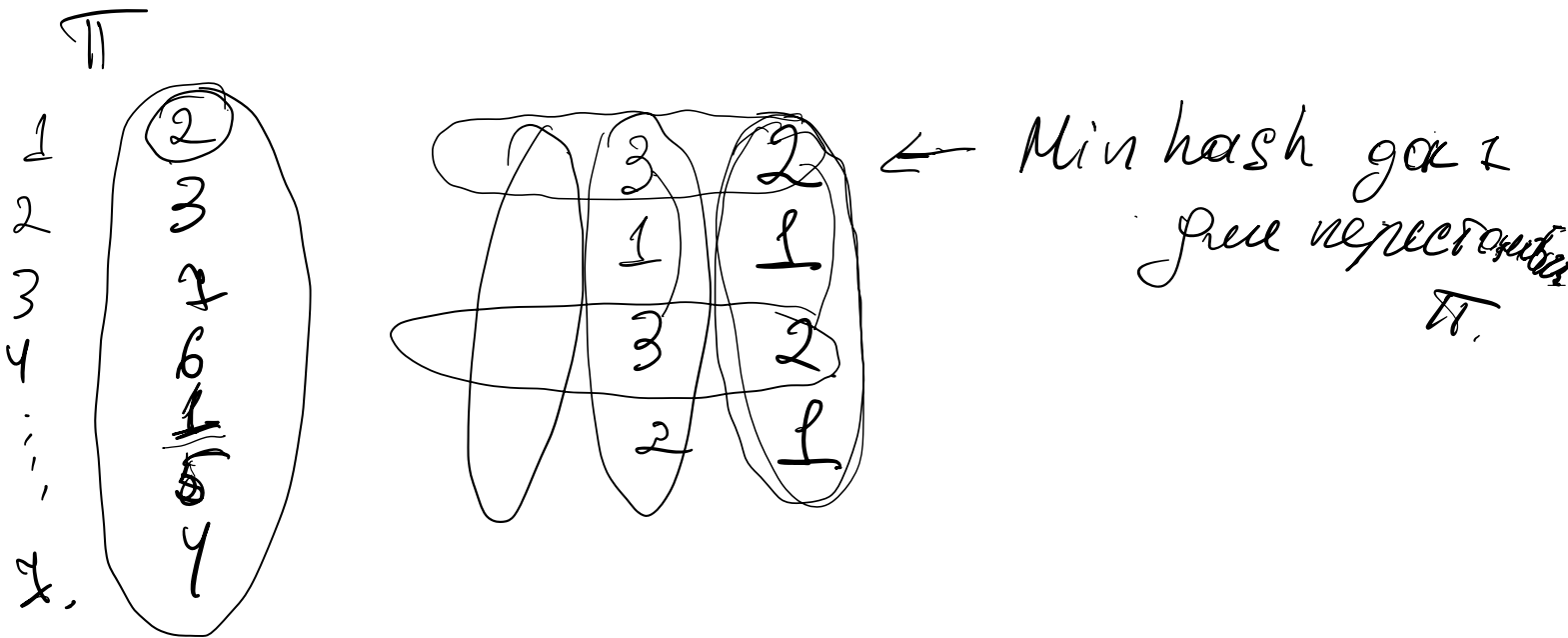
мама ммаа раму.

3-граммы мам, ама, маа,
ама, ...



gon.

n-gram							
1	1	0	0	0	1	1	1
0	0	1	1	1	0	0	3
1	0	0	0	0	1	1	1
0	1	1	1	1	0	0	2



Kozop Dnaa.

$$\frac{|A \cap B|}{|A \cup B|}$$

1 C 1 C

$$\frac{A}{A+B+C}$$

	1	2
A	1	1
B	1	0
C	0	1
D	0	0

$$1-1$$

$$1-0$$

$$0-1$$

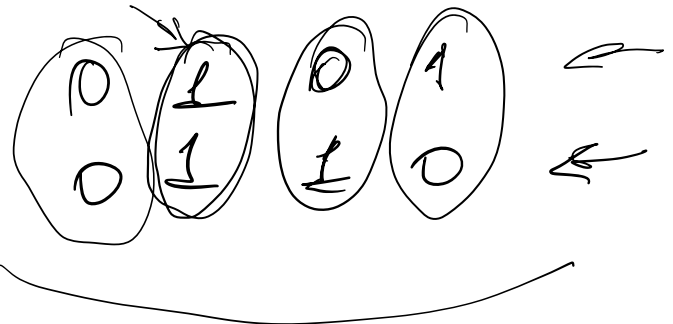
$$0-0$$

$$doc 1 \quad 00110...$$

$$doc 2 \quad 01010...$$

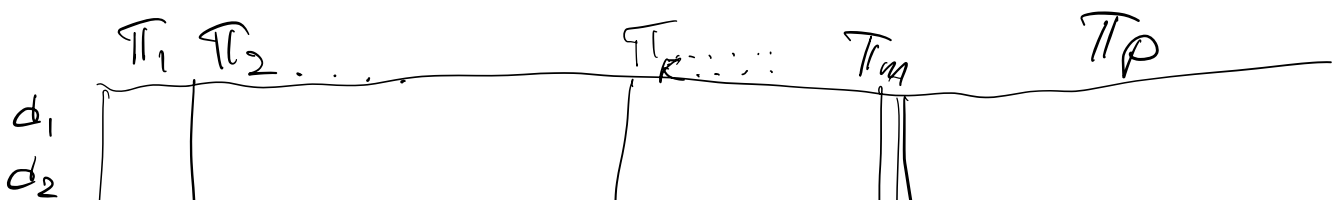
π -перестановка

$$\frac{A}{A+B+C}$$

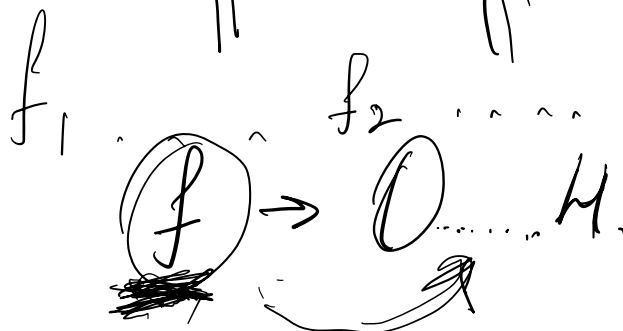
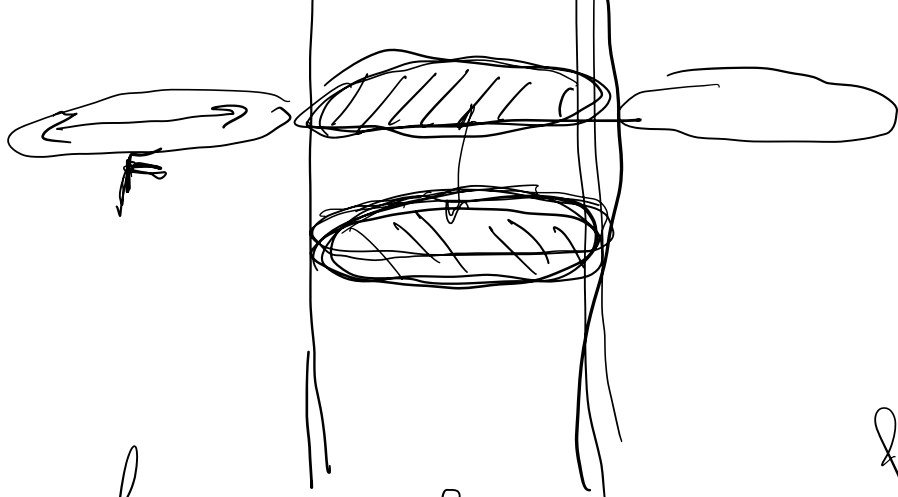


1/3.

LSM



d_0



$$f \left(\begin{matrix} 1, 2, 3, i \\ \sqrt{\quad} \end{matrix} \right)$$