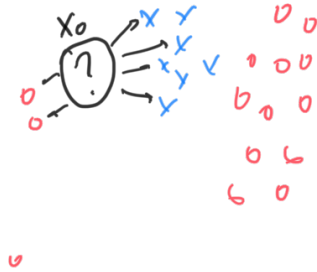


## Быстрый поиск соседей

к NN



Поиск и дист. соседей:

```
① линейный поиск  
idx_best, dist_best = 0,  $\infty$   
for i in range(l):  
    dist = 0  
    for j in range(d):  
        dist += (x0j - xij)2  
    if dist < dist_best:  
        idx_best = i  
        dist_best = dist
```

$\approx$  l.d операциям  
(goes!!!)

Применения:

1) отбор кандидатов в поиск / рекоменд.  
системах

в рекоменд.: много айтемов ( $|I| > 10^6$ )  
нужно уменьшить до  $\approx 10^3$

$p_u$  - вектор пользователя

кандидаты:  $10^3$  айтемов

с минимальными  $p(p_u, p_i)$

в поиске:  $v_q$  - вектор запроса

$w_i$  - векторы документов ( $10^9$ )

отбираем  $i \in I$ :  $p(v_q, w_i) \rightarrow \min$

2) image retrieval

ищем картинки, похожие на заданную

$v_i$  - векторы для картинок

ищем картинки с  $p(v_0, v_i) \rightarrow \min$

Способы ускорения:

- ① Случайный поиск  
(сохраняем  $X$ )
- ② Кластеризовать  $X$   
 $C_1, \dots, C_k$  - центры кластеров  
тогда: ищем ближайший центр  
↓  
ищем соседей только в этом кластере
- ③ Найти важные признаки,  
сначала считать расст.  
только по ним  
или  
сначала PCA
- ④ kd-trees, ball trees  
 $d \approx 10-20 \Rightarrow \log k$  операций  
при  $d \rightarrow \infty$  скорость поиска  $\rightarrow kd$
- ⑤ Отбор эталонов (STOLP)
- ⑥ Приближенный поиск соседей  
( $\delta NN$ )



ищем соседей, которые дальше, чем  
ближайшие, не более чем в  $\delta$  раз

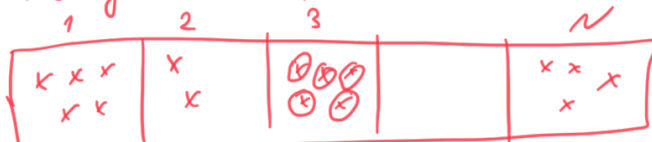
### Locality-sensitive hashing (LSH)

$f(x)$  - хэш-функция

$f(x) \in \{1, \dots, N\}$  - корзины (buckets)

$f(x_1) = 2$   
 $f(x_2) = 10$

$X$  - набор объектов



Опр.  $\mathcal{F}$  - семейство хэм-функций.

$f - (d_1, d_2, p_1, p_2)$  - уравнительное, если

$$2) \quad p(x, z) \geq d_2 \Rightarrow |P_{f \in \mathcal{F}} [f(x) - f(z)]| \leq p_2$$

1) Косинусное расстояние

$$f_w(x) = \text{sign} \langle w, x \rangle$$

$(d_1, d_2, 1-d_1, 1-d_2)$  - избыточные узлы



$$f_{w,b}(x) = \left\lfloor \frac{\langle w, x \rangle + b}{r} \right\rfloor$$

The diagram shows a set of vectors (marked with 'x') in a plane. These vectors are being projected onto a line (marked with 'o'). The orthogonal components of these vectors relative to the line are marked with 'v'. This illustrates the geometric interpretation of the Gram-Schmidt process, where the orthogonal components are the vectors being orthogonalized.



$f_1(x), \dots, f_m(x)$  —  $x \in M$  — функции

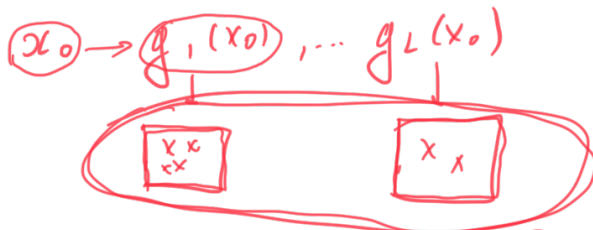
[illegible]

$g(x) = (f_1(x), \dots, f_m(x))$  - координаты  $x$

$f_i \in \{0, 1\}$        $m=3$  :  $\begin{matrix} 000 \\ 001 \\ 010 \\ 011 \\ \dots \\ 111 \end{matrix}$  } 8 ячеек

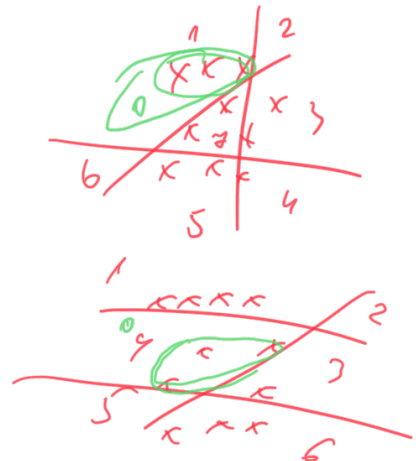
$|g| = 2^m$

$g_1(x), \dots, g_L(x)$  - случайные функции



ищем  $\delta$ -мат. соседей

$m, L$  - гиперпараметры



Можно подобрать  $m$  и  $L$  так,  
что алгоритм  $\delta$ -дгдет находить  
 $\delta$ -мат. соседей за  $O(d \cdot l^r \log l)$ ,  
 $r \approx \frac{1}{\delta}$

$$d \cdot l^r \log l \ll d \cdot l$$

NSW (navigable small world)

Строим граф, вершины соответствуют узлам  
базиса

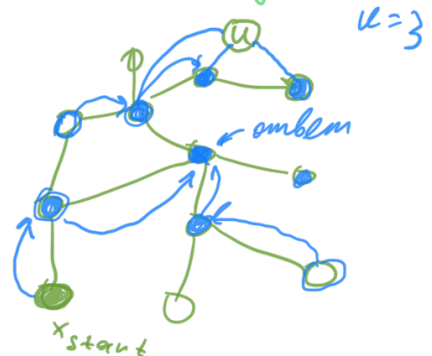
Поиск соседей

$u$  - запрос

Берем случайную вершину  $x$  в графе  
в цикле:

если среди соседей  $x$

есть  $x_i$  :  $p(u, x_i) < p(u, x)$ ,



то  $x = x_i$   
(переходим в лучшего соседа)

Будем запускать несколько раз

### Построение

Добавление в граф нового объекта и

Запускаем поиск соседей  $m$  раз

$S = \{x_1, \dots, x_m\}$  - кандидаты

дополняем  $S$  соседями  $x_1, \dots, x_m$

соединяем и редуцируем  $S$  к

лучш. соседям  $y \in S$

HNSW

$O(\log l)$

'hierarchical