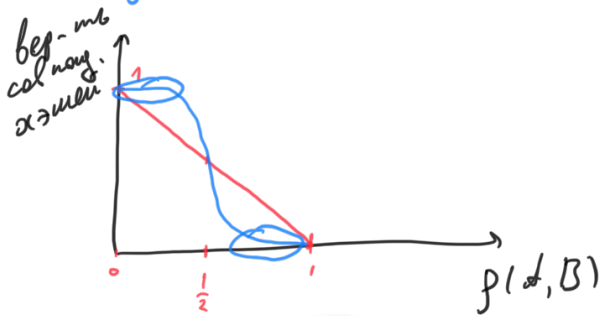


Locality-sensitive hashing (LSH)

① Расстояние Махара

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

$$P(A, B) = d \Rightarrow P[f_{\sigma}(A) = f_{\sigma}(B)] = \underline{1-d}$$



$$x_1, \dots, x_k$$

$$f(x_1), \dots, f(x_k)$$

x	x	x	x	x
x	x	x	x	x
1	2	3	4	5

$$u, f(u) = 3$$

t - порог, м.з.

$$P(x, z) \leq t \Rightarrow \text{с большой вер-тью } f(x) = f(z)$$

$$P(x, z) > t \Rightarrow \text{с большой вер-тью } f(x) \neq f(z)$$

② Композиция хэш-функций

1) $f_1(x), \dots, f_m(x)$ - хэш-функции из \mathcal{F}
(выбираем случайным)

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

$$P_{f \in \mathcal{F}} [f(x) = f(z)] = p$$

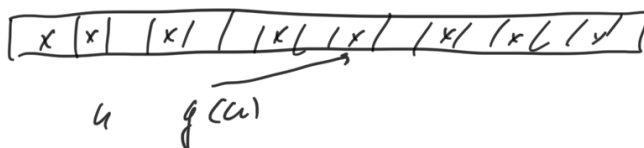
$$(f_1(x), \dots, f_m(x)) = g(x)$$

$$(f_1(z), \dots, f_m(z)) = g(z)$$

$$P[g(x) = g(z)] = \underline{\underline{p^m}}$$

$$f \in \{1, \dots, k\}^m$$

$$g \in \{1, \dots, k\}^{m \cdot m}$$



Может быть много копий, может не найти близких соседей

2) $g_1(x), \dots, g_L(x)$ - L слуг. векторных
 $x \in \mathcal{X}$ - функций

$u \rightarrow g_1(u), \dots, g_L(u)$
 \uparrow
 новый объект
 $C(u) = \{x \in \mathcal{X} \mid g_1(x) = g_1(u) \text{ или } g_2(x) = g_2(u) \text{ или } \dots g_L(x) = g_L(u)\}$
 \uparrow
 кандидаты
 ищем u среди $C(u)$

$$P[g_1(x) = g_1(z)] = p^m$$

$$P[g_1(x) = g_1(z) \text{ или } \dots \text{ или } g_L(x) = g_L(z)] = 1 - (1 - p^m)^L$$

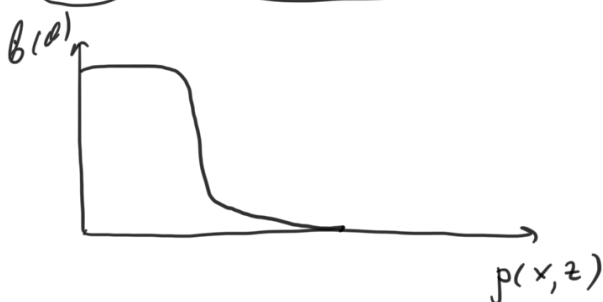
p^m - вер-ть совпадения по g_1

$(1 - p^m)$ - вер-ть несовпадения по g_1

$(1 - p^m)^L$ - вер-ть, что по всем g_1, \dots, g_L нет совпадений

$1 - (1 - p^m)^L$ - вер-ть совпадения хотя бы по одной g_1, \dots, g_L
 -интересно

$$(m, L) \rightarrow 1 - (1 - p^m)^L = b(d)$$



Некоторая теория про LSH

Алгоритм решает задачу поиска

c -ближайшего соседа, если он с

вер-тью $1 - \varepsilon$ для объекта u

возвращает x_u : $p(u, x_u) \leq c \cdot p(u, x_*)$

\uparrow
 истинный

Теорема: можно подобрать m и L так, ^{данный сосед}
 что LSH будет решать задачу поиска
 ϵ -ближ. соседа за $O(d \cdot l^r \log l)$,
 где r зависит от метрики и
 обычно $r \sim \frac{1}{\epsilon}$.

$$O(d \cdot \sqrt{\epsilon} \cdot \log l)$$

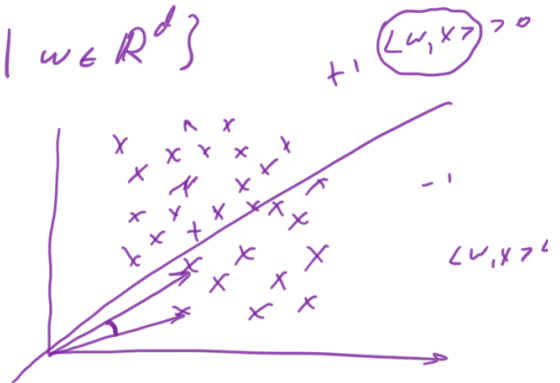
~~$O(d \cdot l)$~~

Пример для косинусного расстояния

$$\mathcal{F} = \{ f_w(x) = \text{sign} \langle w, x \rangle \mid w \in \mathbb{R}^d \}$$

$$f_w(x) \in \{-1, 0, +1\}$$

вер-ть совпадения
 хэшей ^{обратно}
 растет пропорц. углу

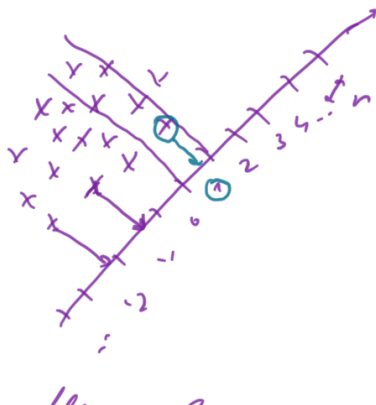


Пример для евклидова расстояния

$$\mathcal{F} = \left\{ f_{w,b}(x) = \left\lfloor \frac{\langle w, x \rangle + b}{r} \right\rfloor \mid w \in \mathbb{R}^d, b \in [0, r) \right\}$$

$$b \sim U[0, r]$$

$$w \sim \mathcal{N}(0, I)$$



работает для раст. лемновского
 $\leq p \in (0, 2]$, если взять правильно
 распределение по w
 $p=1 \Rightarrow w \sim \text{Cauchy}$

Поиск соседей с помощью NSW
 (navigable small world)

$G = (X, E)$

① Поиск соседей

Считаем, что граф у нас есть

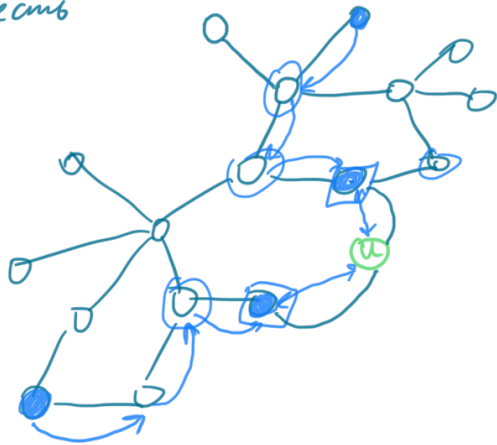
- берем слуг. вершину
 в графе

- в цикле: переходим
 в соседнюю вершину,
 самую близкую
 к u (если есть те,
 которые ближе
 у нас)

v - текущая вершина

$v' : (v, v') \in E$

$p(u, v) > p(u, v')$



делаем мультистарт

$L(u)$ - результаты по
 всем запускам

идем в $L(u)$ или в $L(u) + \text{соседей}$
 этой вершины

② Построение графа

Добавление вершины u :

мультистарт $\rightarrow L(u)$

$\bigcup L(u) \cup \text{соседей}(u) = D(u)$

1. ...
 соединены и с u сумм. соединены
 из $D(u)$

HNSW (Hierarchical NSW)

