

Locality Sensitive Hashing (LSH)

Note Title

1/27/2016

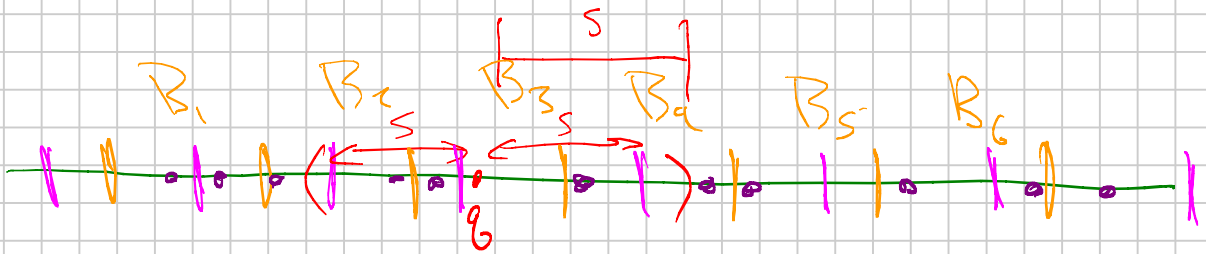
Doc \rightarrow Sets \rightarrow Vectors

$$E[S(v_1, v_2)] = SS(s_1, s_2)$$

Many Documents $\# = n = 1$ million

(Q1) What items are similar? n^2

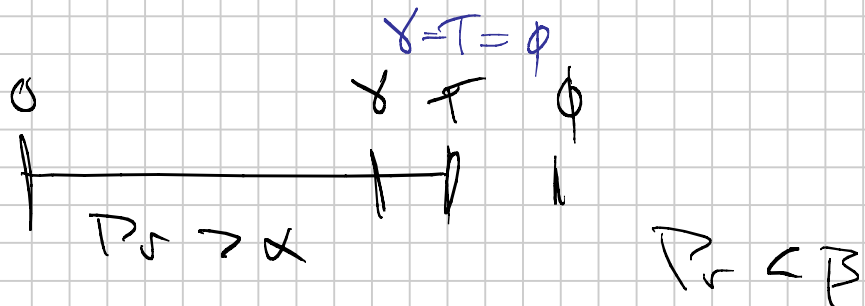
(Q2) Given query, which in set are similar? n



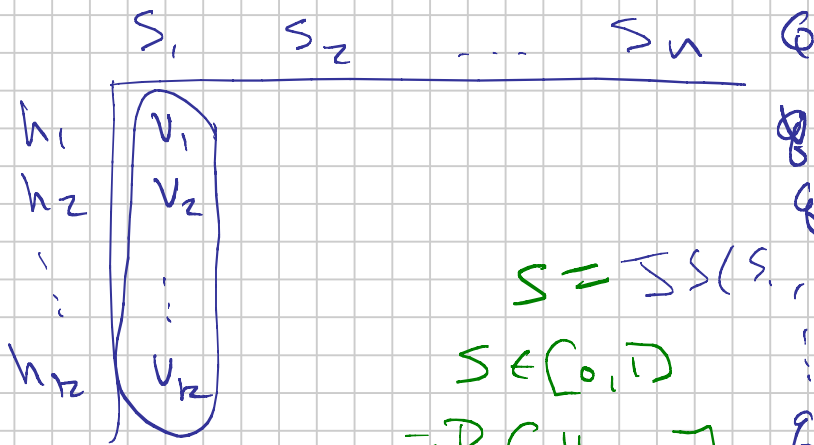
$h \in \mathcal{H}$ $(\gamma, \phi, \alpha, \beta)$ - similar

• $\Pr[h(a) = h(b)] \geq \alpha$ if $d(a, b) < \gamma$

• $\Pr[h(a) = h(b)] < \beta$ if $d(a, b) \geq \phi$



$$h_i: [n] \rightarrow [n \log n]$$



$$S = \Pr(s_1, s_2) = \mathbb{E} \left[\frac{1}{r} \sum_{i=1}^r \mathbb{1}(h_i(s_1) = h_i(s_2)) \right]$$

$$S \in [0, 1]$$

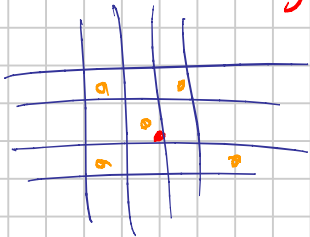
$$= \Pr[\text{collision}]$$

- similar if all hash fcn collide
too harsh \rightarrow Papa bear
- similar if any hash fcn collide
too nice \rightarrow Momma bear

\rightarrow Baby bear

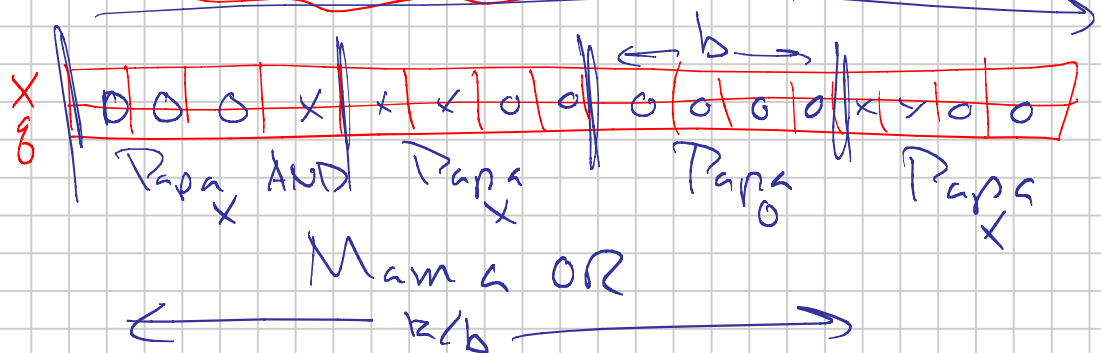


Banding



$$r = t/b$$

Try until 1st collision
 \hookrightarrow then count



Mamma OR

$S = \Pr[\text{collision}]$

$S^b = \Pr[\text{collision in } b \text{ band}]$

$(1 - S^b) = \Pr[\text{no coll in } b \text{ band}]$

$(1 - S^b)^r = \Pr[\text{no coll in } r \text{ bands}]$

$f(s) = 1 - (1 - sb)^r = \text{Prob of coll in any of } r \text{ bands}$



Budgeted to hash fxns

↳ # hash in band b

↳ # bands r $r \cdot b \leq k$

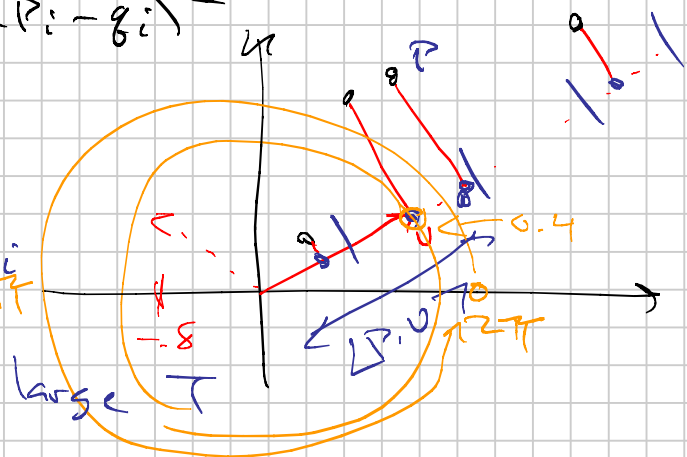
Similarity threshold T

$$b \approx -\log_T(k) \quad r = k/b$$

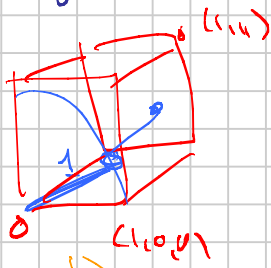
$$d_E(p, q) = \|p - q\| = \sqrt{\sum_{i=1}^d (p_i - q_i)^2}$$

$$p \in \mathbb{R}^d$$

$$h_u(p) = \left\lfloor \frac{\langle p, u \rangle}{\|u\|} \right\rfloor \text{ mod } (T \times \text{large})$$



h_v $v \leftarrow$ rand unit vector in \mathbb{R}^d

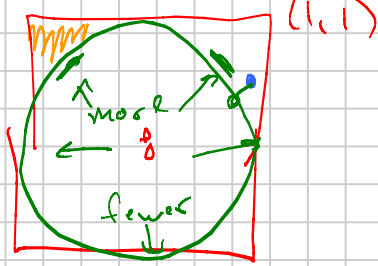
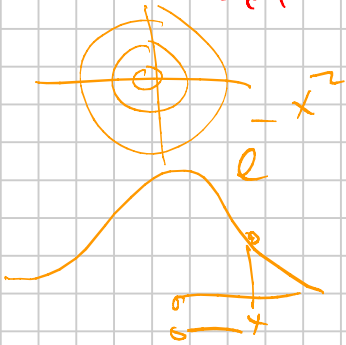


$r \in \mathbb{R}^d$

$r_i \in \text{Unif}(0, 1)$

$a \in \text{Unif}[0, 1]$

$b \in \text{Rand-bit in } \{0, 1\}$



$(-1, -1)$

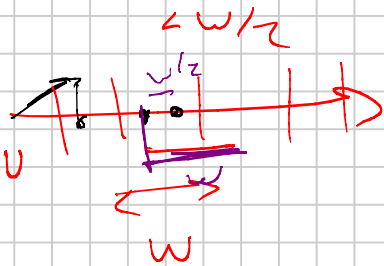
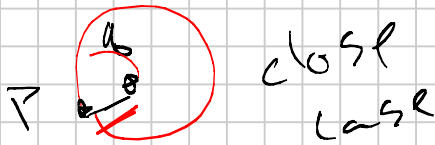
each $g_i \in \mathbb{R}^d$
 $g_i \sim N(0, 1)$

Box-Mueller Transform

$u_1, u_2 \sim \text{Unif}[0, 1]$

$$g_1 = \sqrt{-2 \ln(u_1)} \cos(2\pi u_2)$$

$$g_2 = \sqrt{-2 \ln(u_1)} \sin(2\pi u_2)$$



$$\langle p, v \rangle - \langle g, v \rangle$$

$=$

$$\langle p - g, v \rangle \leq \|p - g\|$$

$$\|p - g\| > 2w$$

