

# L5: Locality Sensitive Hashing

Jeff M. Phillips

January 22, 2020

# Family hash functions $\mathcal{H}$

$$\Pr_{h \in \mathcal{H}} [h(p) = h(q)] \approx \text{sim}(p, q)$$

1. 1 hash function

$$\hat{J}_S(p, q) = \begin{cases} 1 & h(p) = h(q) \\ 0 & \text{otherwise} \end{cases}$$

hash table

2.  $k$  hash functions

$$\hat{J}_S^k(p, q) = \frac{1}{k} \sum_{j=1}^k \mathbb{1}(h_j(p) = h_j(q))$$

Algo.?

Jaccard

Triangle

$\approx$

Euclidean  
(dot product)

Indicator Function

$$\mathbb{1}(b) = \begin{cases} 1 & \text{if } b = \text{True} \\ 0 & \text{if } b = \text{False} \end{cases}$$

Large Number of objects  $X$

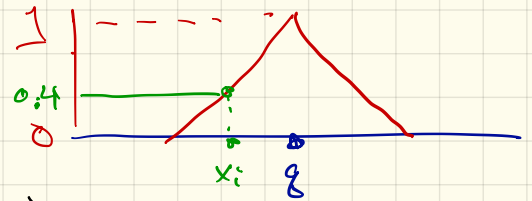
$$X = \{x_1, x_2, \dots, x_n\}$$

(documents (e-mails), IP addresses, customers)

Q1: Which pairs are similar?  
 $n^2$  time

Q2: Given query  $q$ , which  
 $x_i \in X$  are similar to  $q$ ?  
 $n$  time

$$x_1, x_2, \dots, x_n \in \mathbb{R}$$



similarity  $S_{\Delta}(g, x_i) = \max\{0, 1 - |g - x_i|\}$

1. Sort  $x_1, x_2, \dots$

2. Build binary tree  $T$



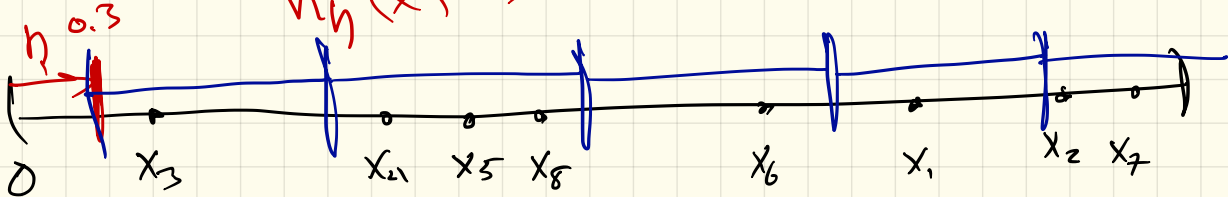
3. Find  $g$  in  $T$

$h \in \text{unif}(0, 1)$

$h_h(x) \rightarrow$  a bin  $= S_{\Delta}(x, x')$

$P[h(x) = h(x')] \log n + 12$

# similar items



Banding  $\rightarrow$  How to combine hash functions

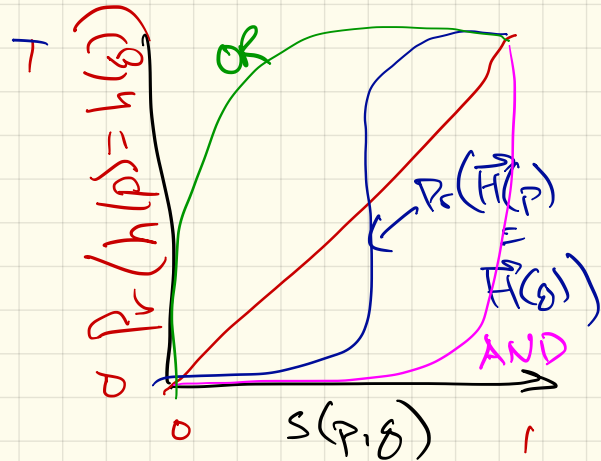
$$H = \{h_1, h_2, \dots, h_r\} \in \mathcal{H}$$

$\vec{H} \leftarrow$  single super hash function

Band  $b=2$

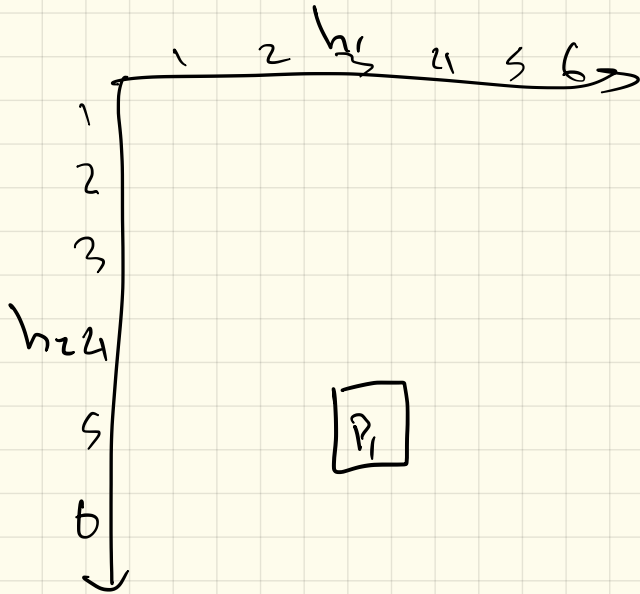
	$h_1$	$h_2$	$h_3$	$h_4$	$h_5$	$h_6$
$P_1$	3	5	0	2	4	3
$P_2$	3	4	1	0	2	2
$P_3$	0	2	4	3	5	1
$P_4$	$H_1$		$H_2$		$H_3$	
$P_5$						
$P_6$						
$g$	3	5	1	2	3	2

# bands  $= r$



$$\vec{H}(p, g) = \text{OR}(H_1, H_2, H_3)$$

$$P_1 \rightarrow [3, 5]$$



$$P_\sigma [h_1, h_2(p) = h_1, h_2(s)]$$

$$= S^2$$

Much more selective

$r$  bands, each with  $b$  hash functions

$t = \#$  hash functions  $t \geq r \cdot b$

$$S(P, g) = s$$

$s^b = P_r$   $P, g$  collide in one band.

$(1 - s^b) = P_r$   $P, g$  don't collide

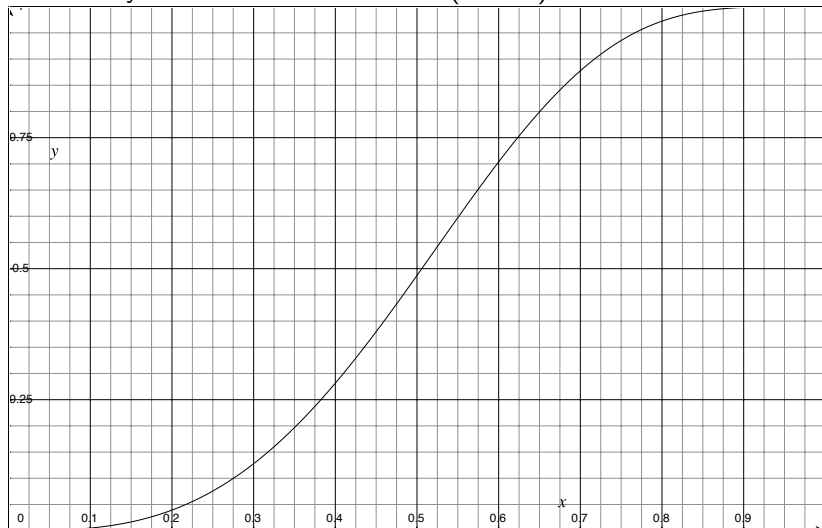
$(1 - s^b)^r = P_r$   $P, g$  don't collide in  $r$  bands

$f(s) = 1 - (1 - s^b)^r = P_r$   $P, g$  collide in at least one band.

LSH  $b = 3$  and  $r = 5$

$$t = 15$$

Probability of found collision =  $1 - (1 - s^b)^r$



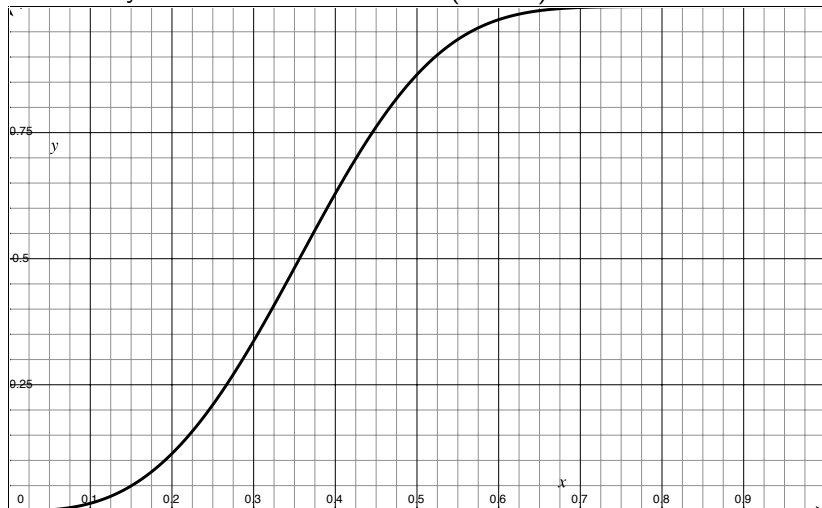


LSH  $b = 3$  and  $r = 15$

Probability of found collision =  $1 - (1 - s^b)^r$

LSH  $b = 3$  and  $r = 15$

Probability of found collision =  $1 - (1 - s^b)^r$



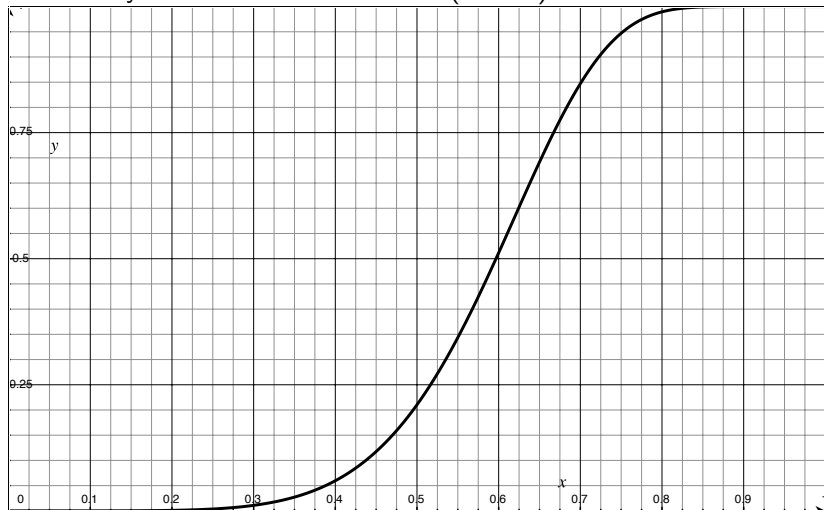
LSH  $b = 6$  and  $r = 15$

Probability of found collision =  $1 - (1 - s^b)^r$

LSH  $b = 6$  and  $r = 15$

$$t = 90$$

Probability of found collision =  $1 - (1 - s^b)^r$

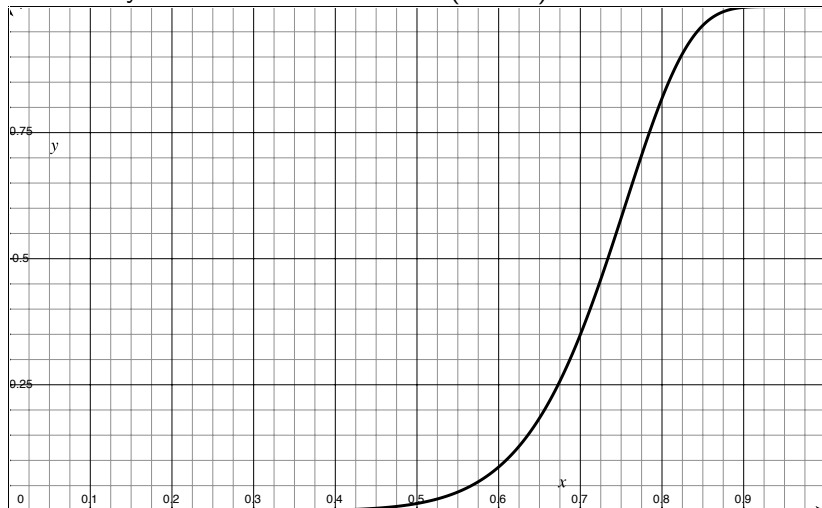


LSH  $b = 10$  and  $r = 15$

Probability of found collision =  $1 - (1 - s^b)^r$

LSH  $b = 10$  and  $r = 15$

Probability of found collision =  $1 - (1 - s^b)^r$

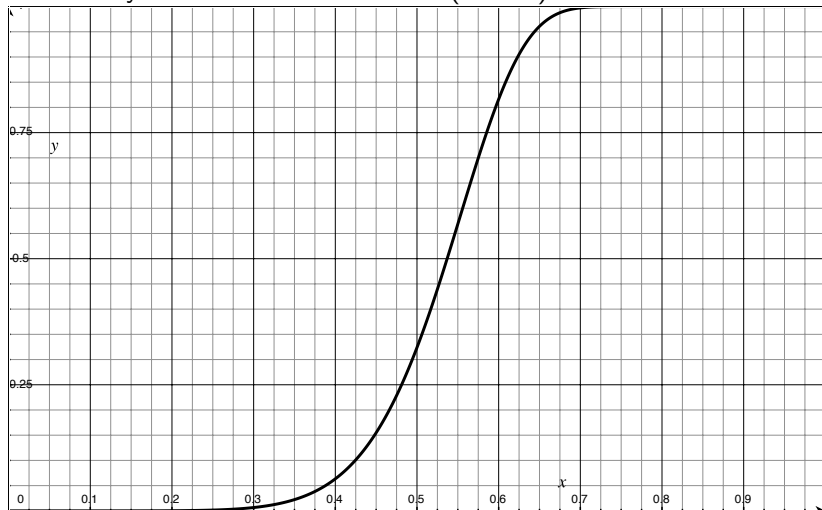


LSH  $b = 8$  and  $r = 100$

Probability of found collision =  $1 - (1 - s^b)^r$

LSH  $b = 8$  and  $r = 100$

Probability of found collision =  $1 - (1 - s^b)^r$



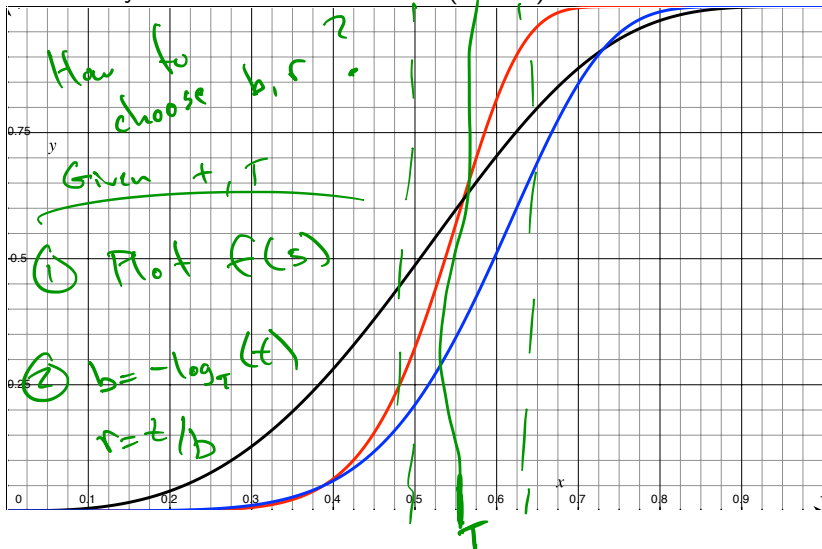


LSH ( $b = 3, r = 5$ ) & ( $b = 6, r = 15$ ) & ( $b = 8, r = 100$ )

Probability of found collision =  $1 - (1 - s^b)^r$

# LSH ( $b = 3, r = 5$ ) & ( $b = 6, r = 15$ ) & ( $b = 8, r = 100$ )

Probability of found collision =  $1 - (1 - s^b)^r$



# LSH for Euclidean Dist.

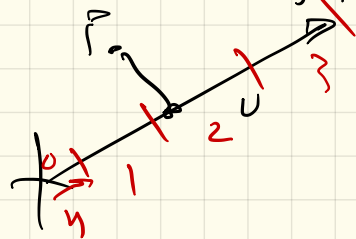
$$d_E(p, q) \iff S_E(p, q)$$

$$h_{h, v}: \mathbb{R}^d \rightarrow [m]$$

$$h \in \text{unif}(0, \tau)$$

$$v \in \mathbb{R}^d, \|v\|=1$$

$$\begin{aligned} \langle p, q \rangle &= \sum_{i=1}^d p_i \cdot q_i \end{aligned}$$



$$h_{h, v}(p) = \left( \underbrace{\lfloor \langle p, v \rangle - h \rfloor}_{\text{floor}} \bmod m \right)$$