

L6: Locality Sensitive Hashing

Jeff M. Phillips

January 27, 2025

Documents

D_1

s_1



D_2

\vdots

D_n

s_n

min hash
vectors

m_1

m_2

\vdots

m_n

$$m_j \in \mathbb{R}^t$$

1 : if True

0 : if False

$$JS(s_1, s_2) = E \left[\frac{1}{t} \sum_{i=1}^t \mathbb{1}(m_1(i) = m_2(i)) \right]$$

$t \approx 1000$

$$P_c \left[|JS(s_1, s_2) - \tilde{JS}| > \epsilon \right]$$

$< \delta$

Q₁: Have n documents and query q

Find similar docs to q

time \ll linear in n .

Ex: search engine

RAG

Q₂: Have n documents,
find all similar pairs

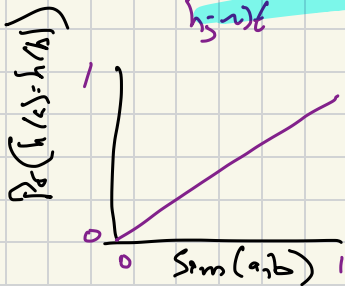
time \ll quadratic in n

Ex: check plagiarism.

Locality Sensitive Hashing

$h_1, h_2, \dots, h_t \sim \mathcal{H}$ $h_j: \mathcal{X} \rightarrow [0]$
 ~~$h_i \in \mathcal{H}$, random hashes~~ ~~$\Pr[h_j(a) = h_j(b)] = \frac{1}{v}$~~

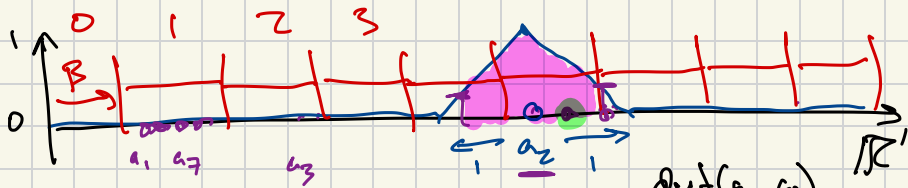
LSH: $\Pr[h_j(a) = h_j(b)] = \text{Sim}(a, b)$
 $h_j \sim \mathcal{H}$ $\in [0, 1]$



min hash

$$\Pr[m_a(i) = m_b(i)] = \text{JS}(a, b)$$

Triangle Similarity objects $a_1, a_2, \dots, a_n \in \mathbb{R}^1$



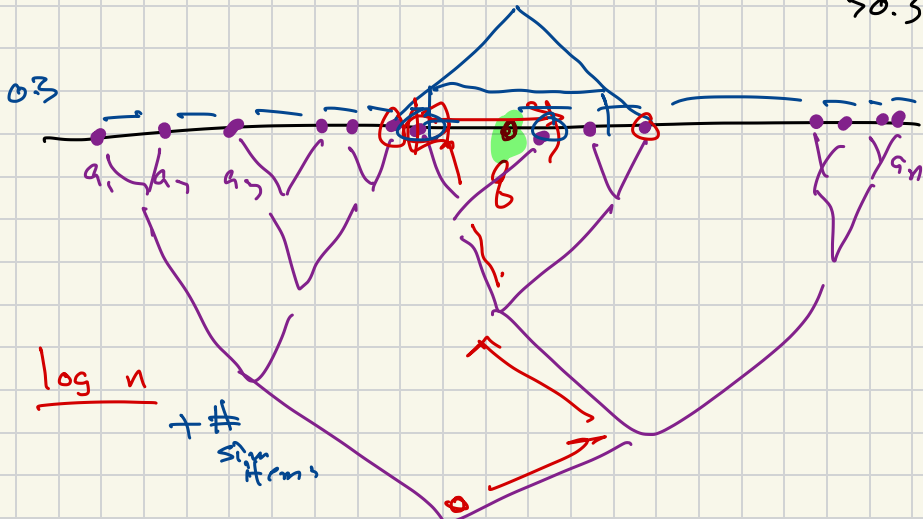
$$\text{Sim}_{\Delta}(a_1, a_2) = 1 - \min\left\{1, \underbrace{\text{dist}(a_1, a_2)}\right\}$$

$$\beta \sim \text{Unif}(0, 1) \quad h_{\beta} = \text{bin}(\text{offset } \beta)$$

$$\Pr[h_{\beta}(a_1) = h_{\beta}(a_2)] = \text{Sim}_{\Delta}(a_1, a_2)$$

Biked

Under Sim_Δ , find all objects
 from a_1, a_2, \dots, a_n with $\text{Sim}_\Delta(a_i, g) > 0.3$



	s_1	s_2	s_3	s_4	...	s_n
h_1	1	1	3	4		1
h_2	2	0	3	1		2
h_3	2	3	2	1		1
h_4	4	4	3	2		3

AND

h_1, h_2

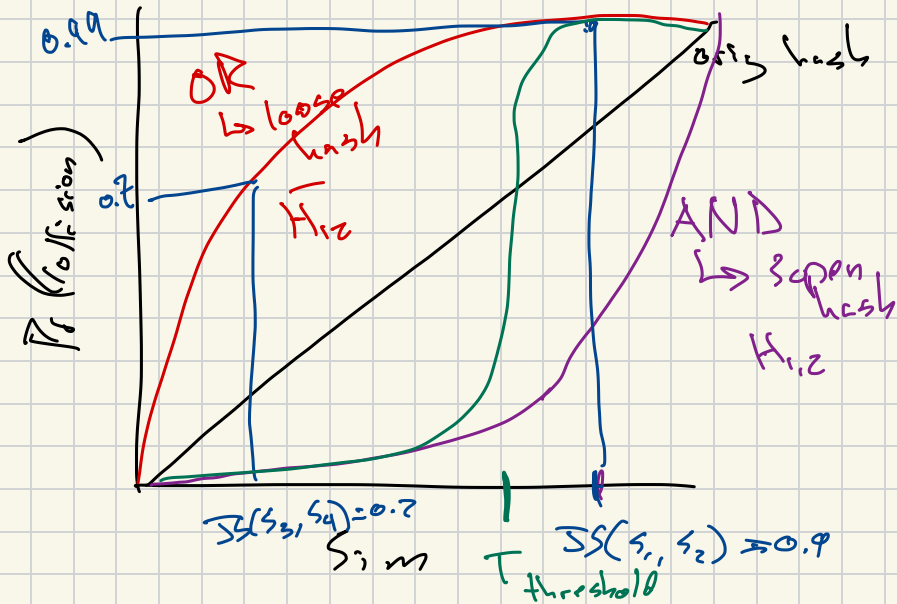
$$H_{1,2}(s_i) = (1, 2)$$

$$H_{h_2}(s_1) \stackrel{?}{=} H_{1,2}(s_4)$$

s_i	0	1	2	3	4	5
1	1	2	3	4	5	
6	7	8	9			

$$P_i [H_{1,2}(s_1) = H_{1,2}(s_2)]$$

$$? = \text{Sim}(s_1, s_2)$$



OR heterogeneous hash cells

$$\overline{H}_{1,2}(s_1) = \overline{H}_{1,2}(s_2) \text{ if}$$

$$\underline{h_1}(s_1) = h_1(s_2) \quad \text{or} \quad \underline{h_2}(s_1) = h_2(s_2)$$

Banding: combines together
AND and OR hashes

t : total # hash functions

b : # within 1 band (~~OR~~ ^{AND} combinations)

s : # of bands (~~AND~~ ^{OR} combo)

$$t = b \cdot s$$

The lecture reversed the order of AND and OR combinations of hashes – on this slide

	s_1	s_2	s_3	s_4	...	s_n
h_1	1	1	3	2		1
h_2	2	0	3	1		2
h_3	2	3	2	1		1
h_4	1	0	1	3		1
\vdots						
h_{n-1}	3	2	4	1		3
h_n	2	2	2	2		2

$b = 2$
band

AND
OR
 h_1 and h_2

bands
= 5

$t/b = 5$

$t = 6$

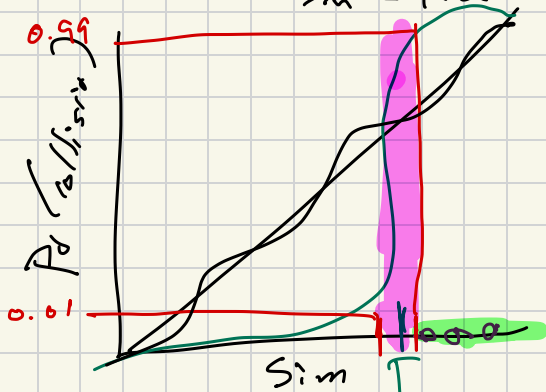
$s = 3$

The lecture reversed the order of AND and OR combinations of hashes – on this slide

Before $P_c[h(a) = h(b)] = \text{Sim}$

After $P_c[h^*(a) = h^*(b)] > 0.99 \quad \text{if} \quad \text{Sim} > T + \alpha$

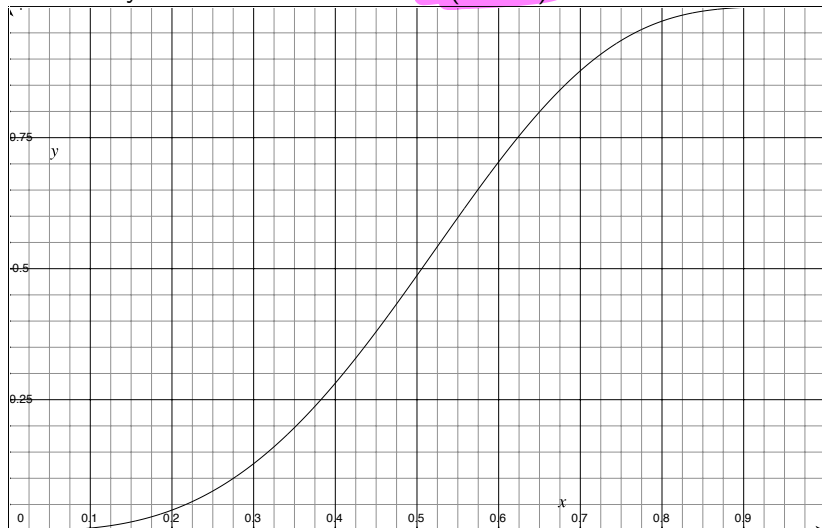
$P_c[\dots] < 0.01 \quad \text{if} \quad \text{Sim} < T - \alpha$



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

$$t = 15$$

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

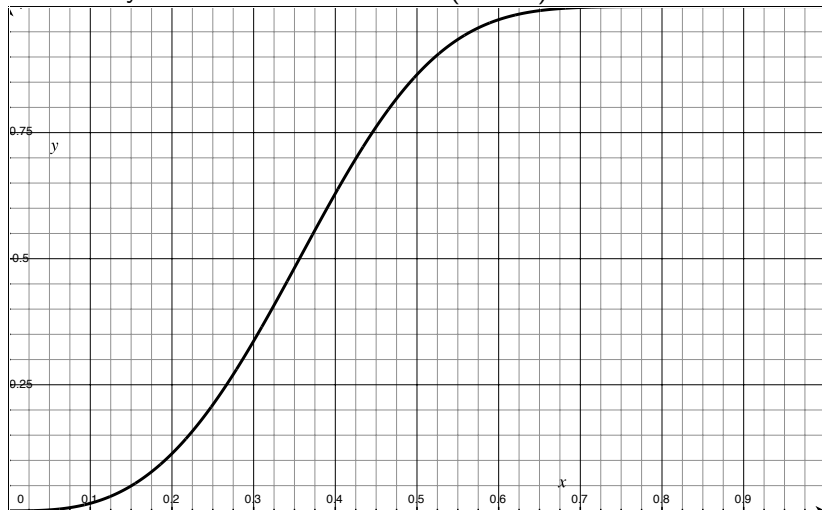


LSH $b = 3$ and $r = 15$

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

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

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

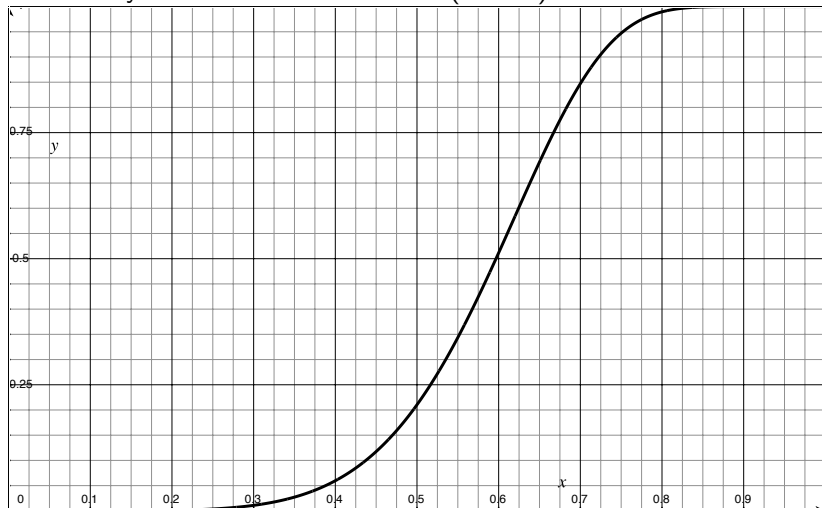


LSH $b = 6$ and $r = 15$

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

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

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

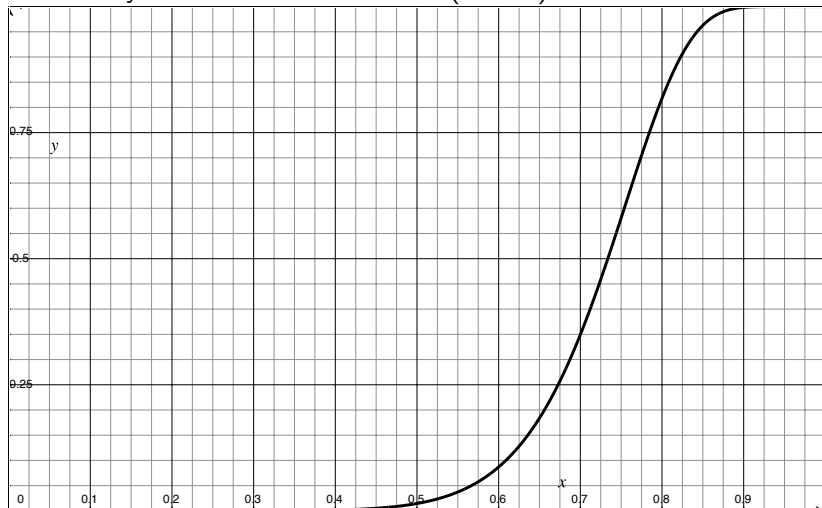


LSH $b = 10$ and $r = 15$

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

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

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

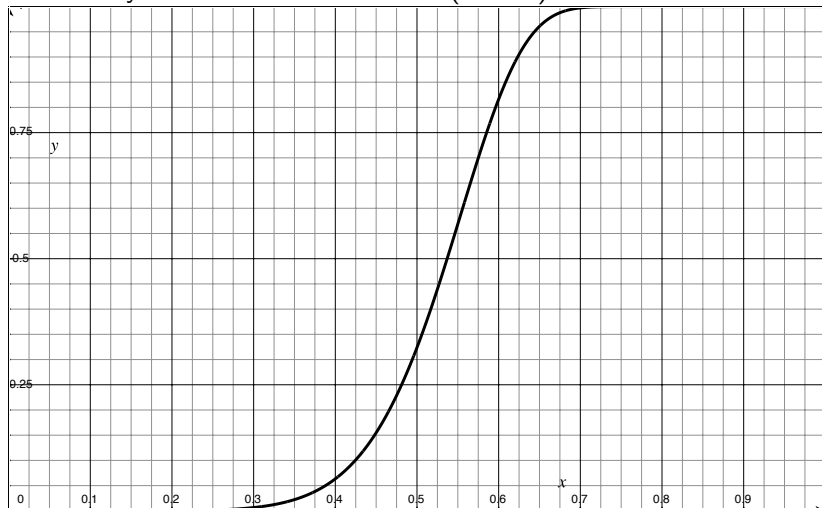


LSH $b = 8$ and $r = 100$

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

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

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

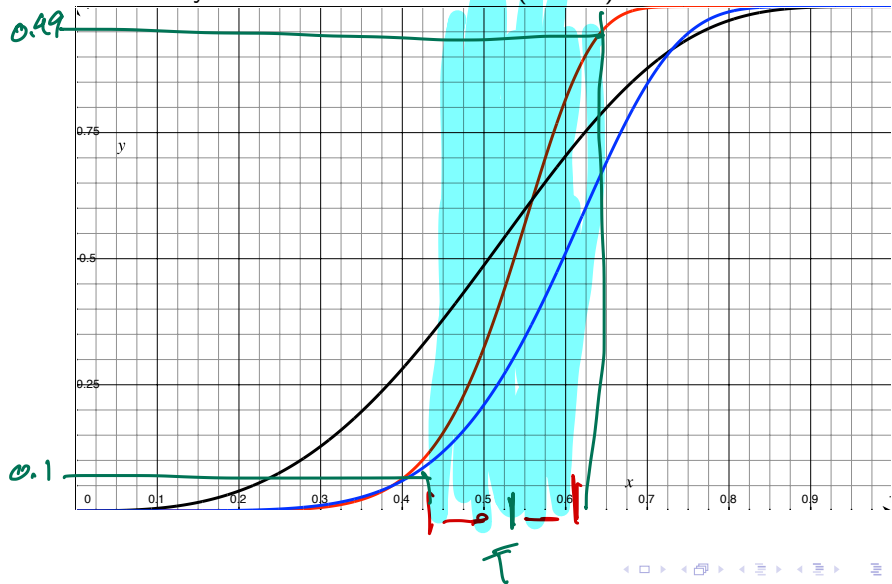


LSH ($b = 3, r = 5$) & ($b = 6, r = 15$) & ($b = 8, r = 100$)

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 = $1 - (1 - s^b)^r$

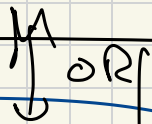




h_1	1	1	1
h_2	2	2	3

AND # inside band

band
diff all
hashes
in one
band



h_3	3	2	3
h_4	2	3	2

of bands
 r

bands of size b | $s = \text{Sim}(a, b)$
bands r

1 hash, Prob coll = s

b hashes in band, Prob ^{not} coll = $(1 - \underline{s^b})$

r bands, each w/ b hashes
Prob not coll = $(1 - s^b)^r$

r bands, each w/ b hashes

Prob collision = $1 - (1 - s^b)^r$

If have t hash functions,
how to split into bands

$t \geq b \cdot r$, Give threshold T

(1) Plot $1 - (1 - sb)^r$ for diff choices
compare to T

(2) Rule of thumb: steepest slope
 $1 - (1 - sb)^r$

$$b = -\log_r(t)$$

$$\text{so } r = \frac{t}{b}$$

good
first
succ.