

L5: Locality Sensitive Hashing

Jeff M. Phillips

January 24, 2018

Set $X = \{P_1, P_2, \dots, P_n\}$

Q₁: Given query doc q
Find all in X close $\ll O(n)$
time

Q₂: Find all pairs in X close.
 $\ll O(n^2)$
time.



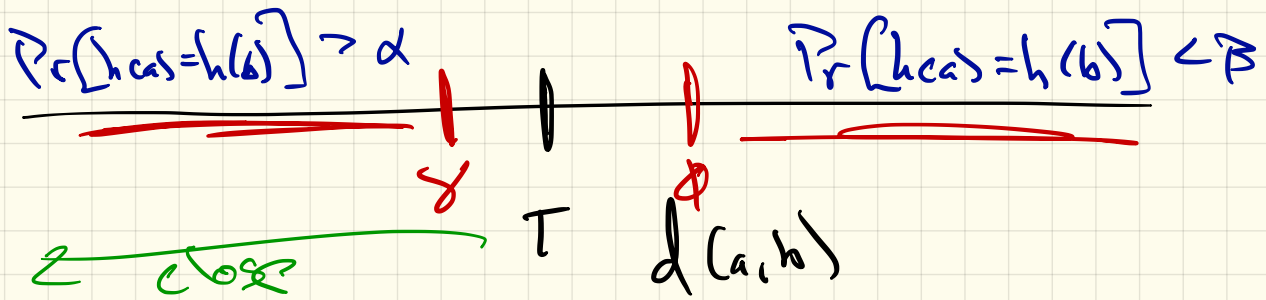
LSH

hash fns $\{h_1, h_2, \dots, h_k\} \in \mathcal{H}$

$h \in \mathcal{H}$ is $(\delta, \phi, \alpha, \beta)$ -sensitive

$$\bullet \Pr[h(a) = h(b)] > \alpha \quad \text{if} \quad d(a, b) < \delta$$

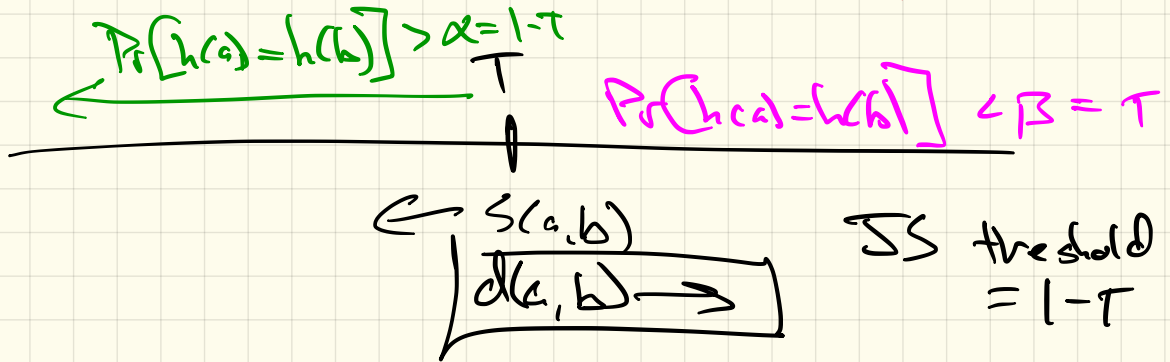
$$\bullet \Pr[h(a) = h(b)] < \beta \quad \text{if} \quad d(a, b) > \phi$$



Min-Hashing \Rightarrow $(T, T, 1-T, T)$ -sensitivity

Documents	D_1	D_2	...	D_n
$h_1 \rightarrow m_1$	1	7		1
$h_2 \rightarrow m_2$	2	2		8

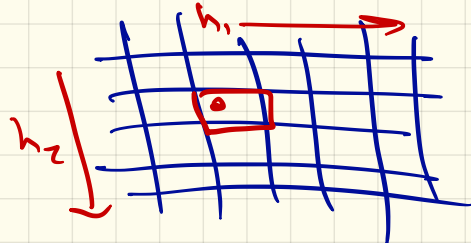
Choose threshold $T = \phi = \delta$



1 hash fxn

for query doc g

Return all $D_i \in X$ sth. $h(g) = h(D_i)$



Papa

• Apply b hash fxns, and return $D_i \in X$ sth. all collide w/ $h(g)$

↳ Super hash table $\vec{h} = h_1 \times h_2 \times \dots \times h_b$

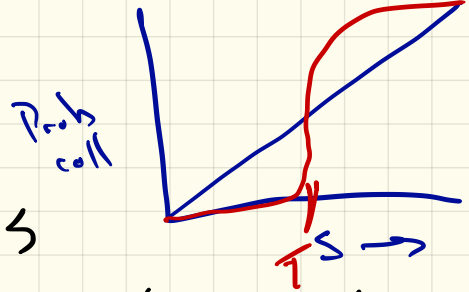
Mama

• Apply r hash fxns \rightarrow r hash tables
return union of collisions

Banding

Use $b \cdot r$ hash fns

r superhash tables, each w/ b hash fns



$$s = SS(a, b)$$

Prob return \rightarrow
on query a

s^b = Prob collide 1 super hash

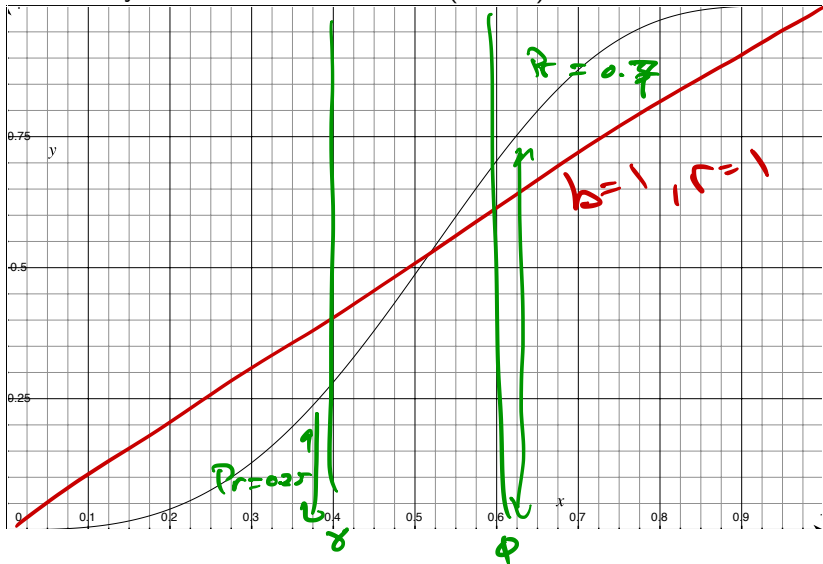
$(1-s^b)$ = Prob don't collide, 1 s.h.

$(1-s^b)^r$ = Prob no s.h.t. collisions

$f(s) = 1 - (1-s^b)^r$ = Prob at least 1 s.h.t. collisions

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

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



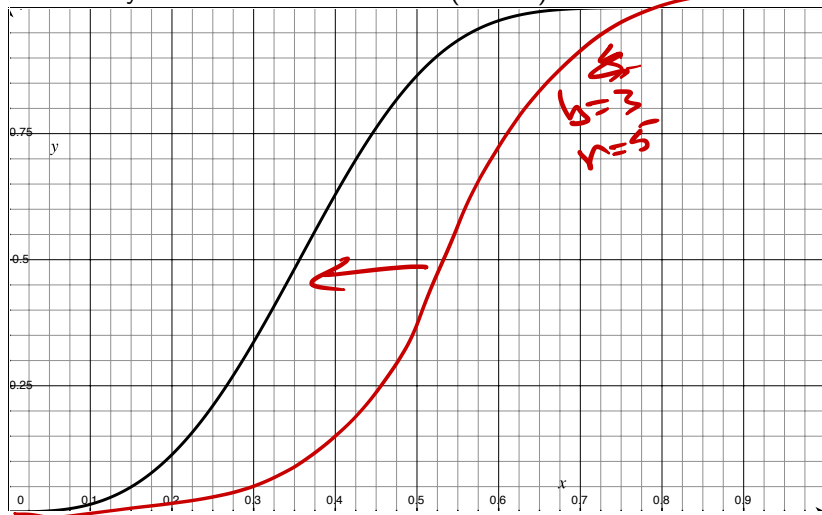
LSH $b = 3$ and $r = 15$

Increase # s.h.t.

$$\text{Probability of found collision} = 1 - (1 - s^b)^r$$

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

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

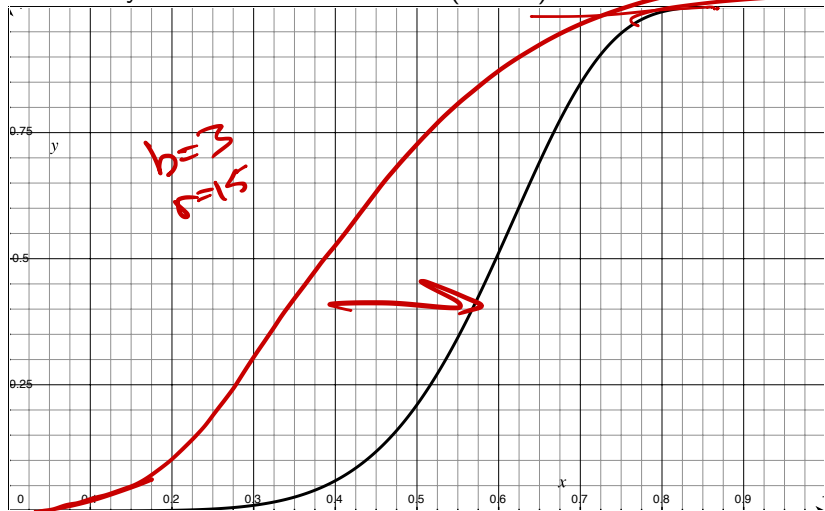


LSH $b = 6$ and $r = 15$

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

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

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

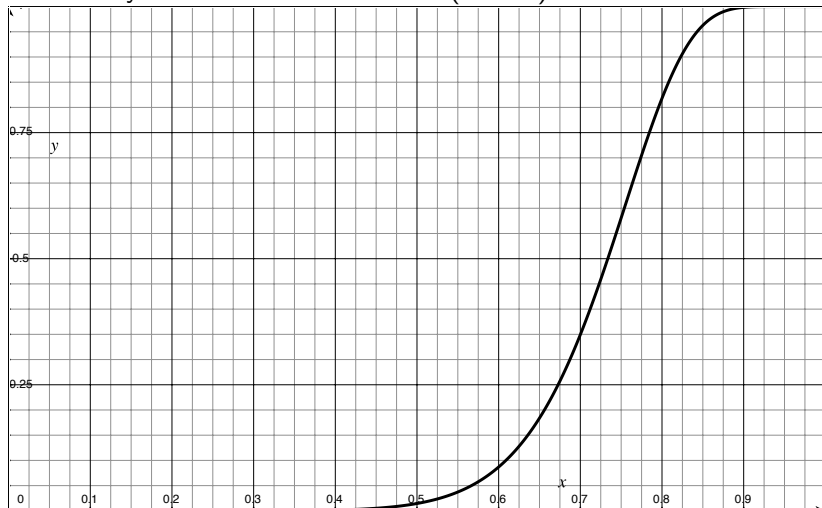


LSH $b = 10$ and $r = 15$

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

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

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



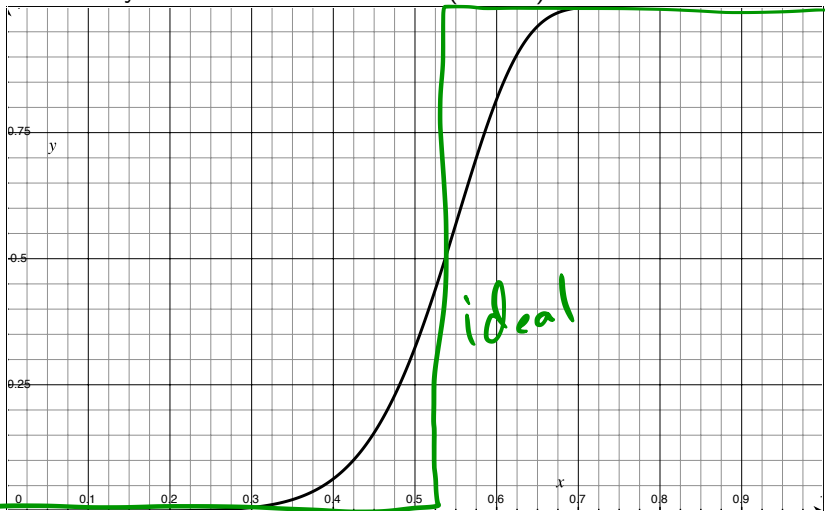
LSH $b = 8$ and $r = 100$

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

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

$h, r \rightarrow$
steeper

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



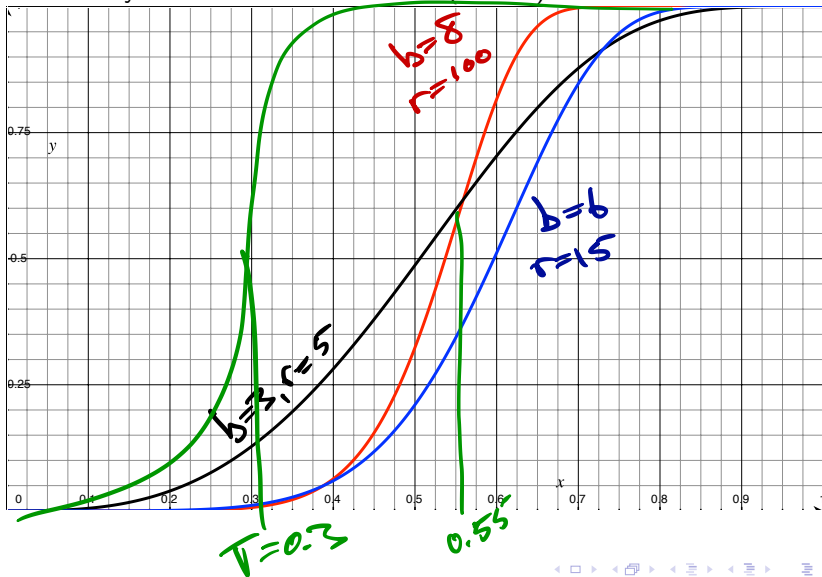
$T = 0.52$

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

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

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

$$\text{Probability of found collision} = 1 - (1 - s^b)^r$$



Choosing r, b so curve
is steepest at T
threshold

$$t = r \cdot b$$

\neq

steepest $T \approx (1/r)^{(1/b)}$

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

$$r = t/b$$

"rule of thumb"

make
integers

→ then experiment

$(\delta/2, 2\delta, 1/2, 1/2)$ - sensitive
 LSH for Euclidean Dist.

$$P = (p_1 \dots p_d)$$

$$g \in (g_1 \dots g_d)$$

$$\|P - g\| = \sqrt{\sum_{i=1}^d (p_i - g_i)^2}$$

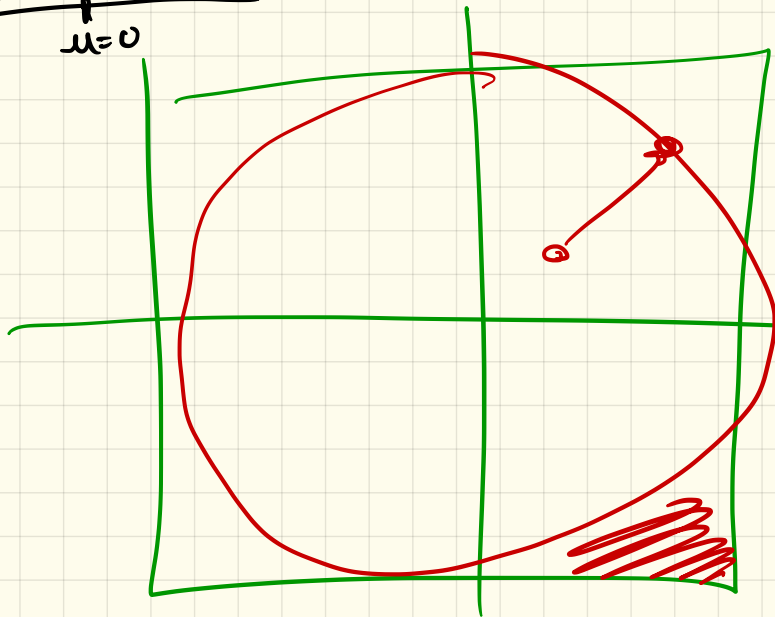
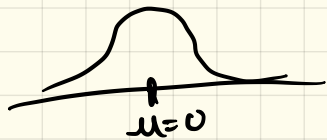
hash fxn

1. project random direction
 how? u unit vector

2. bin distances

$$h(p) = \left\lceil \frac{\langle u, p \rangle \cdot \delta}{\text{bin width}} \pmod{n} \right\rceil$$

↑
random



Generate
d-dim Gaussian
R.V.

$$g \sim G_d$$

$$= g = (g_1, g_2, \dots, g_d)$$

where each

$$g_i \sim G_i = \exp(-x^2)$$

2 uniform $u_1, u_2 \in \text{Unif}(0,1)$

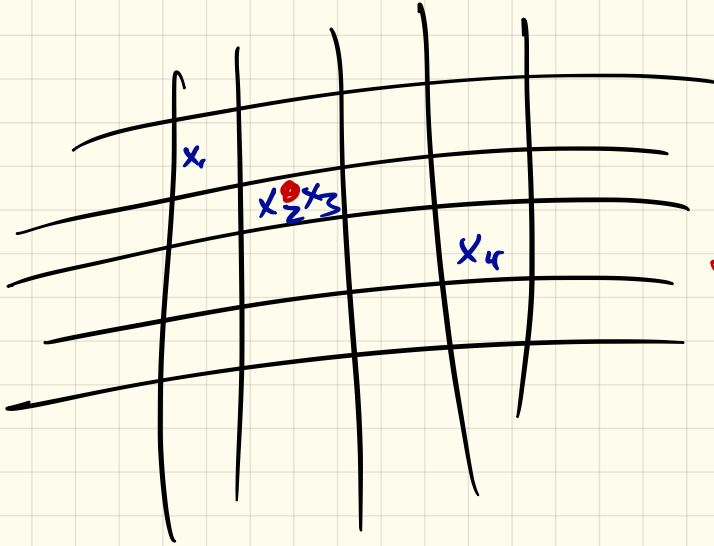
$$\Rightarrow g_1, g_2 \in G_i$$

Box-Mueller
transform

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

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

query
↓



⇒ Union $\{x_2, x_3\}$
 $\{x_3\}$

