

$$J S(A, B) = \Pr [V_i(A) = V_i(B)]$$

$$J S(A, B) = \mathbb{E} \left[\frac{1}{k} \sum_{i=1}^k \mathbb{1}(V_i(A) = V_i(B)) \right]$$

L5: Locality Sensitive Hashing

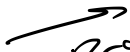
Jeff M. Phillips

January 23, 2019

Raw Data



text



Programs

Abstract Set

min hashing

vector $\begin{bmatrix} m \\ v_1, v_2, \dots, v_m \end{bmatrix}^k$

$n = 1$ million documents $D_1, D_2 \dots D_n$
k-grams: sets $A_1, A_2 \dots A_n$
minhash v_1, v_2, \dots, v_n

Q1: which pairs of objects are similar?

$JS(A_i, A_j) > T$ (eg: = 0.85)
 n^2 distance calculations

Q2: Given a query $D_q \rightarrow A_q \rightarrow v_q$,

which objects are similar?

n distance calculations.

Have n numbers $S = \{s_1, s_2, \dots, s_n\} \subset \mathbb{R}$

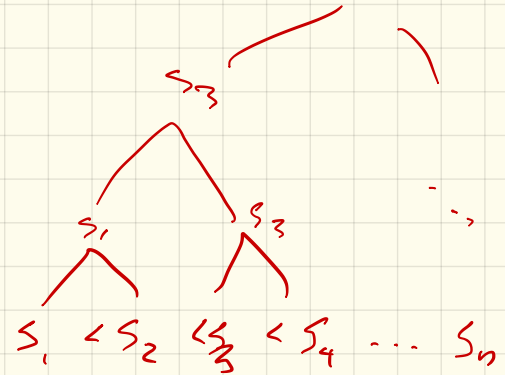
Store S in a data structure

so given query $g \in \mathbb{R}$

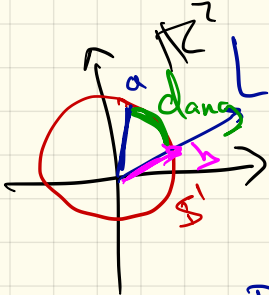
return $\arg \min_{s_i \in S} |g - s_i|$

• Sort in list.

• Build binary tree



Q2: $O(\log n)$ time



Locality Sensitive Hashing

$$h \in \mathcal{H}_{\text{cid}}$$

$$\text{dang}(a, b) = \arccos(\langle a, b \rangle) \quad \Pr[h(a) = h(b)] \approx S(a, b)$$

Jaccard $\Pr[h(A) = h(B)] = JS(A, B)$

Angular Sim $\Pr[h(A) = h(B)] = S_{\text{ang}}(A, B)$

$$a, b \in \mathbb{S}^{d-1} = \{a, b \in \mathbb{R}^d \mid \|a\|=1, \|b\|=1\}$$

$$\bar{a} \leftarrow \frac{v \in \mathbb{R}^d}{\|v\|}$$

$$S_{\text{ang}}(a, b) = 1 - \frac{1}{\pi} \arccos(\langle a, b \rangle) \in [0, 1]$$

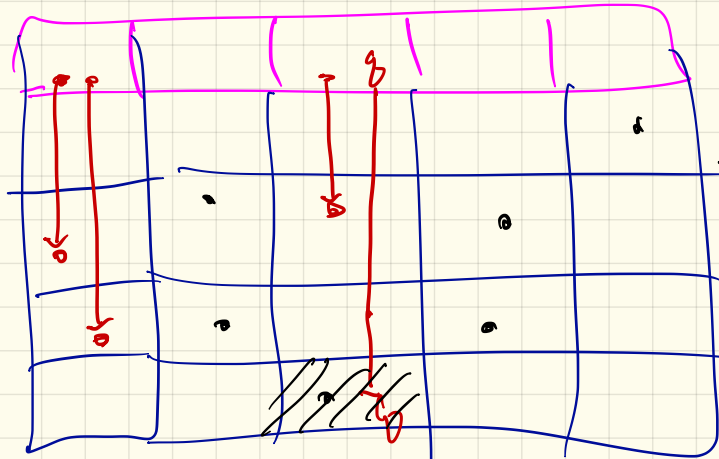
Euclidean \hookrightarrow dot sim

$$\Pr_S[h(A) = h(B)] \approx \langle A, B \rangle$$

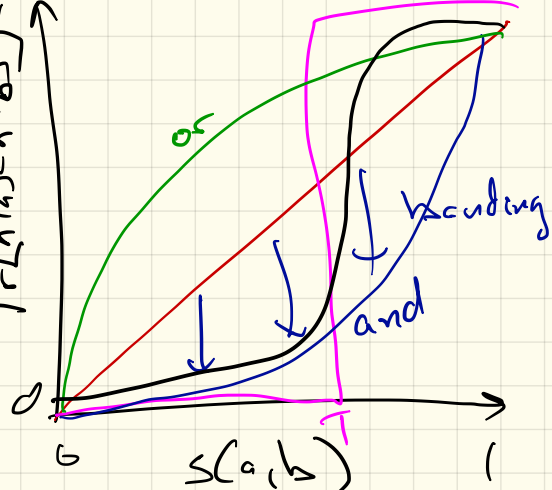
h_2

$$Pr[h(a) = h(b)] = S(a, b)$$

h_1



$Pr[h(a) = h(b)]$

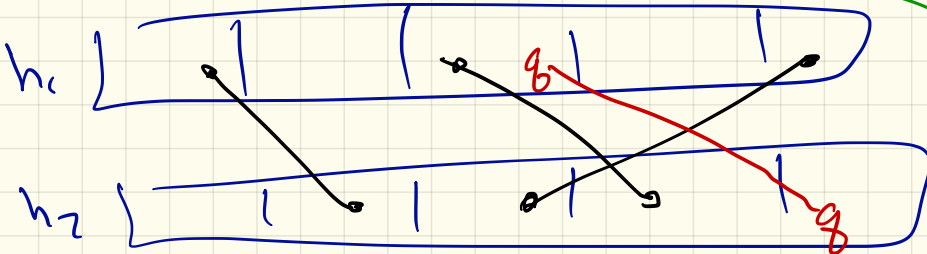


band \equiv one has hash table

papa bear

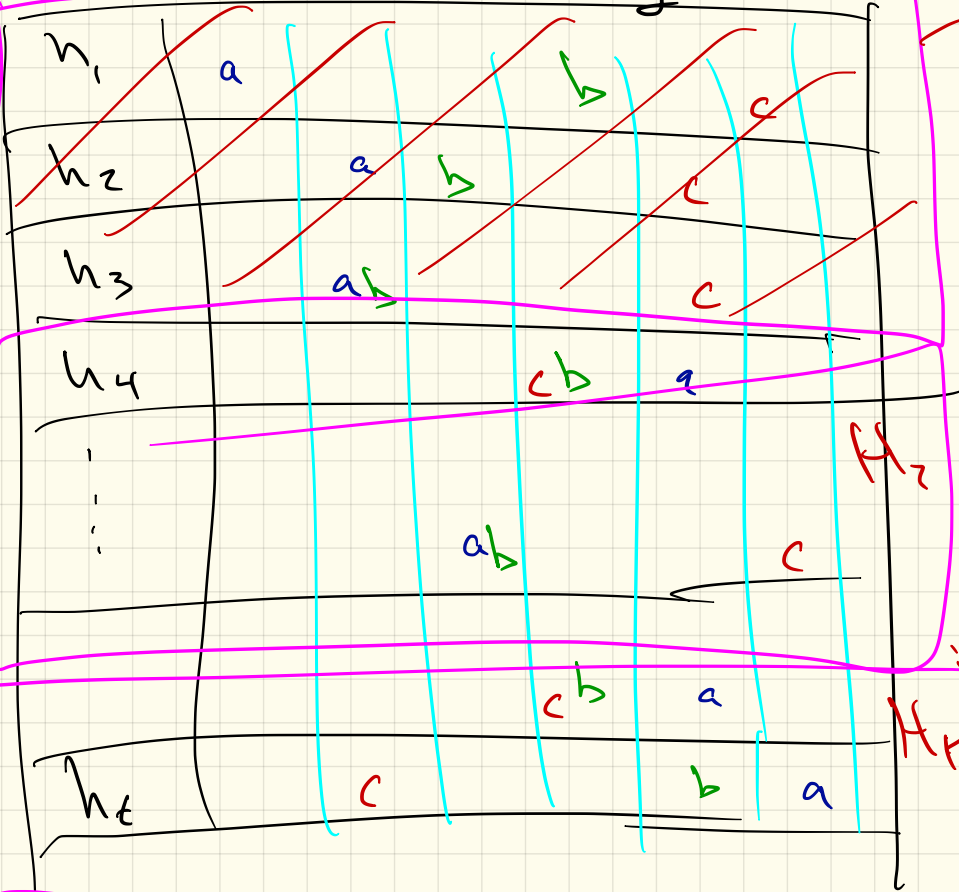
$S(a, b)$
hobby bear

mama bear



banding

one hash table
 H_i



band of
 hashes
 AND
 Bands

Bands

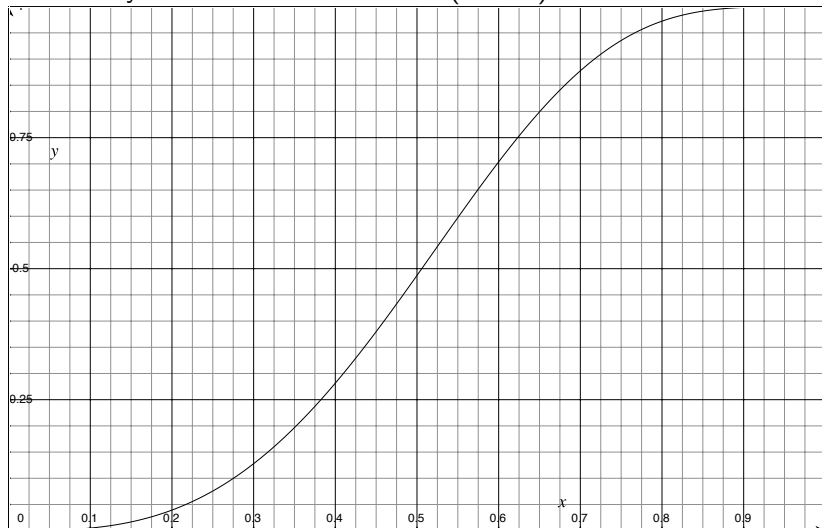
⋮

Bands

Bands
 or
 Bands
 or
 ⋮
 or
 Bands

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

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



Banding b hashes in band
 r bands (meta hash fn)

$$s = S(A, B) \leftarrow \Pr(h(a) = h(b)) \quad h \in \mathcal{H}$$

s^b = prob all hash fn in one band collide

$(1 - s^b)$ = prob not all hash " "

$(1 - s^b)^r$ = probability that in no band, do all hash collide

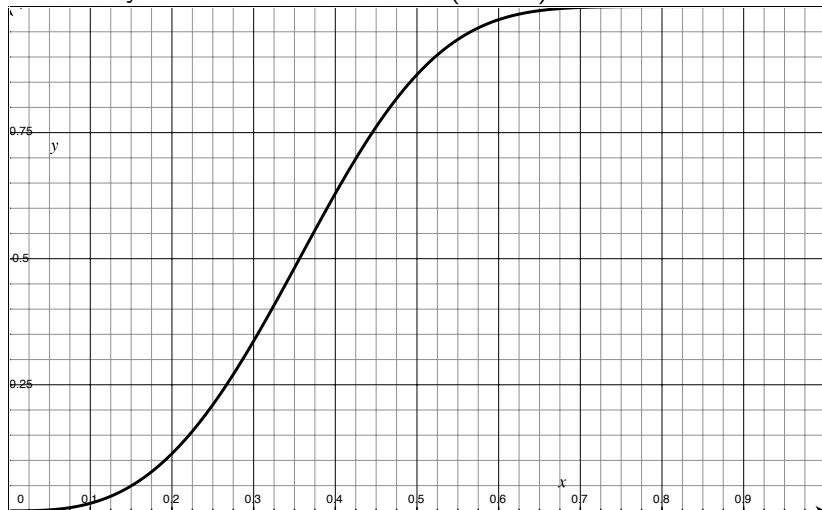
$f(s) = 1 - (1 - s^b)^r$ = prob of at least one band finds collision.

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

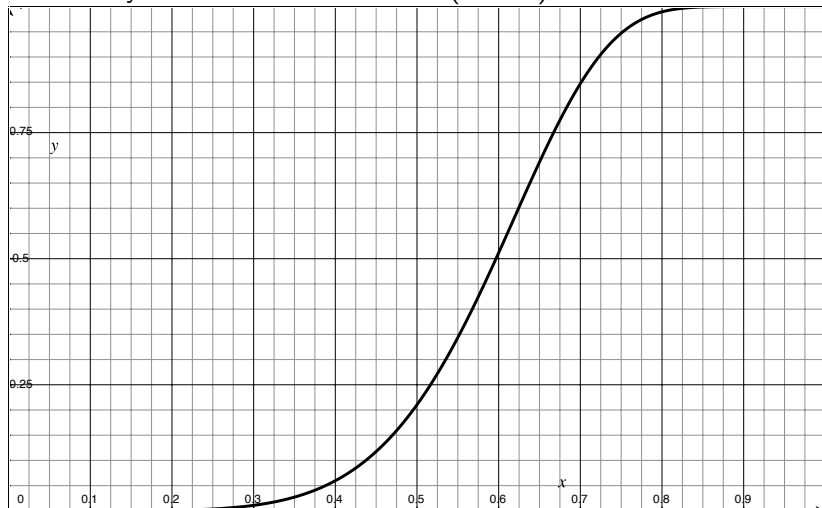


LSH $b = 6$ and $r = 15$

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

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

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

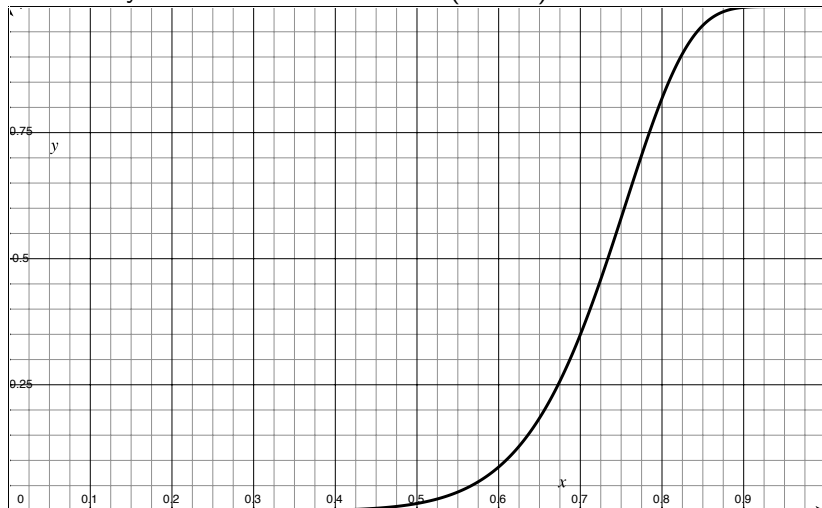


LSH $b = 10$ and $r = 15$

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

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

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



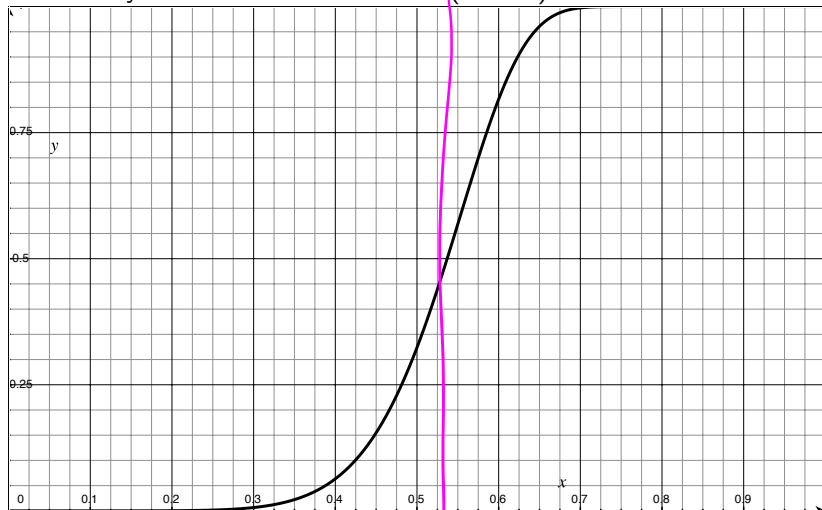
LSH $b = 8$ and $r = 100$

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

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

$$t = r \cdot b = 100 \cdot 8 = 800$$

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

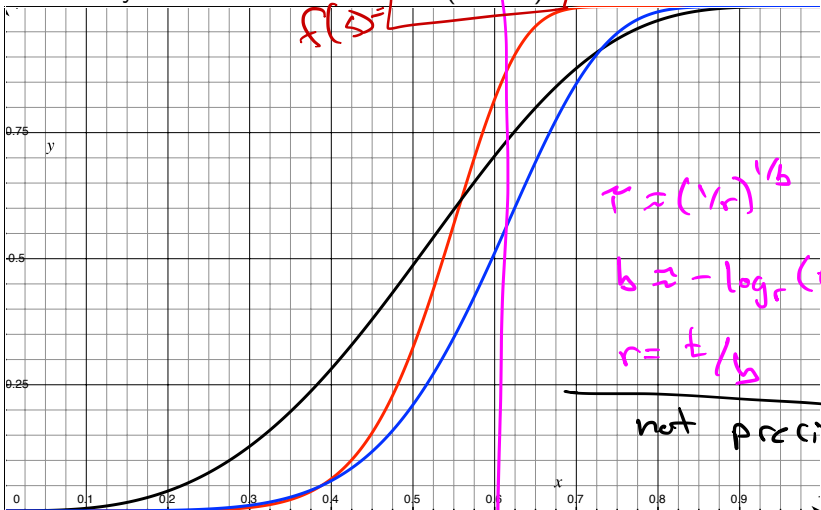
LSH ($b = 3, r = 5$) & ($b = 6, r = 15$) & ($b = 8, r = 100$)

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

$t = 1000$

$f(s) =$

$f(s)$



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

$b \approx -\log_r(t)$

$r = t/b$

not precise

s

$T = 0.6$

$$S_{\text{ang}} : \mathbb{S}^{d-1} \times \mathbb{S}^{d-1} \rightarrow [0, 1]$$

$$S_{\text{ang}}(a, b) = 1 - \frac{1}{\pi} \arccos(\langle a, b \rangle)$$

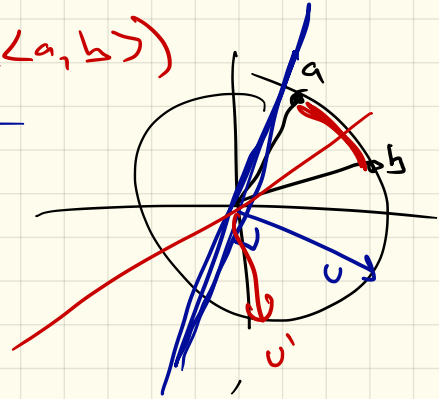
$$h_u \in \{ \pm 1 \} \text{ (unit } \mathbb{S}^{d-1} \text{)}$$

$$h_u : \mathbb{S}^{d-1} \rightarrow \{-1, +1\}$$

$$h_u(a) = \text{sign}(\langle u, a \rangle)$$

$$h_u(a) = +1 \quad h_u(b) = +1$$

$$h_{u'}(a) = -1 \quad h_{u'}(b) = +1$$

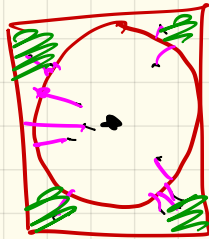


$$v \sim \text{Unif}(\mathbb{S}^{d-1})$$

Guess

$$P \sim (\text{Unif}[-1, 1])^d$$

$$v = \frac{P}{\|P\|}$$



\mathbb{R}^d

$d=2$: rejection sampling

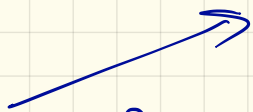
$$P \sim \text{Unif}[-1, 1]^2$$

if $\|P\| > 1 \rightarrow$ start over

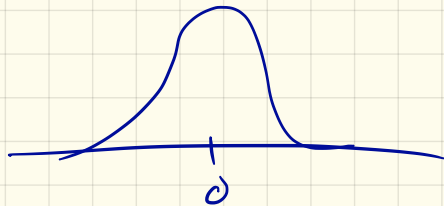
o.w $v = \frac{P}{\|P\|}$

$$v \sim \text{Unif}(\mathbb{S}^{d-1})$$

$$1. \quad g \sim G_d(x) = \frac{1}{(2\pi)^{d/2}} e^{-\|x\|^2/2}$$



$$2. \quad v = \frac{g}{\|g\|}$$



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

$$g_i \stackrel{\text{iid}}{\sim} G_1(x) = \frac{1}{\sqrt{2\pi}} e^{-\|x\|^2/2}$$

Box-Mueller Transform

$$u_1, u_2 \sim \text{Unif}(0, 1) \Rightarrow$$

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

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