# L7: Locality Sensitive Hashing & Distribution Distances

Jeff M. Phillips

September 10, 2025

# Min Hashing    $h \to \mathcal{H}$    $S, S'$ sets

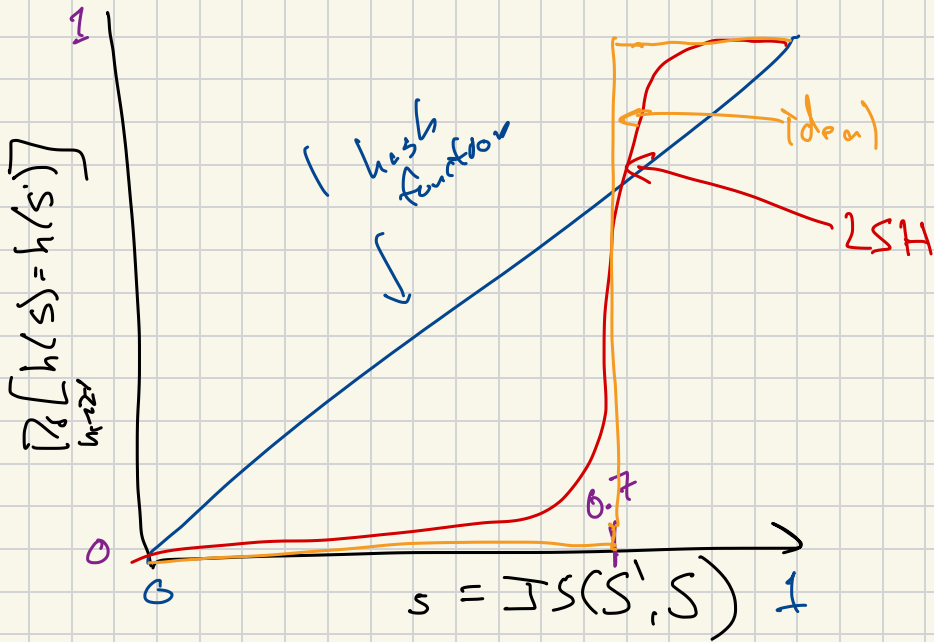$$\Pr_{h \to \mathcal{H}} \left[ h(S) = h(S') \right] = JS(S, S')$$

$$E_{h \to \mathcal{H}} \left[ \mathbb{1} \left( h(S) = h(S') \right) \right] =$$

$$h_1, h_2, \ldots h_t \stackrel{\sim}{\text{iid}} \mathcal{H}$$

$$\begin{cases} 1 & \text{if True} \\ 0 & \text{if False} \end{cases}$$

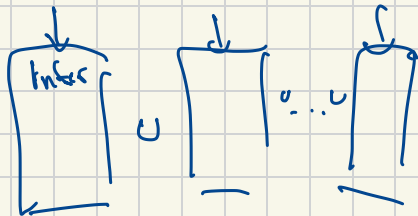$$\hat{JS}(S, S') = \frac{1}{t} \sum_{j=1}^{t} \mathbb{1} \left( h_j(S) = h_j(S') \right)$$

$$E_{h_1 \ldots h_t \to \mathcal{H}} \left[ \hat{JS}(S, S') \right] = JS(S, S')$$

# Aggresive (few false negatives)

Put guess $g \rightarrow h_1(g), h_2(g) .. h_t(g)$

Union of all
intersections!



Inters $\cup$ $\cup ... \cup$

# Conservative (few false positives)

Concatenate $\left[ h_1(g) \; h_2(g) .. h_t(g) \right] \rightarrow$ bits
hash
table.

# Banding

$$t = b \cdot r$$

$\#$ new hashes    size of band    $\#$ of bands

$\left.\begin{array}{c} h_1 \\ \vdots \\ h_b \end{array}\right\} \begin{array}{c} h_1 \\ h_2 \end{array}$    $H_1$

$\left.\begin{array}{c} h_{b+1} \\ \vdots \\ h_{2b} \end{array}\right\} \begin{array}{c} h_3 \\ h_4 \end{array}$    $H_2$

$\left.\begin{array}{c} h_{t-b} \\ h_t \end{array}\right\} \begin{array}{c} h_{t-1} \\ h_t \end{array}$    $H_{t/b} = r$

Return union of collisions on $H_1, H_2, \dots H_{t/b}$

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

$t = b \cdot r = 3 \cdot 5 = 15$
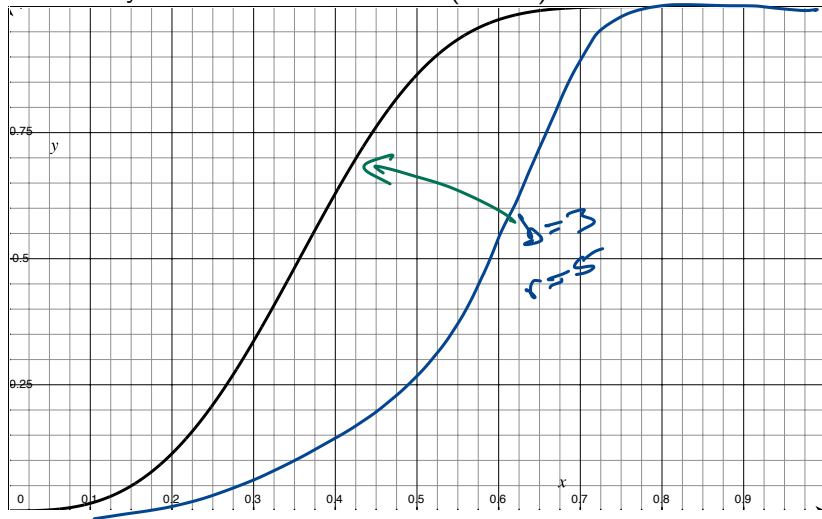
Probability of found collision $= 1 - (1 - s^b)^r$

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

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

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

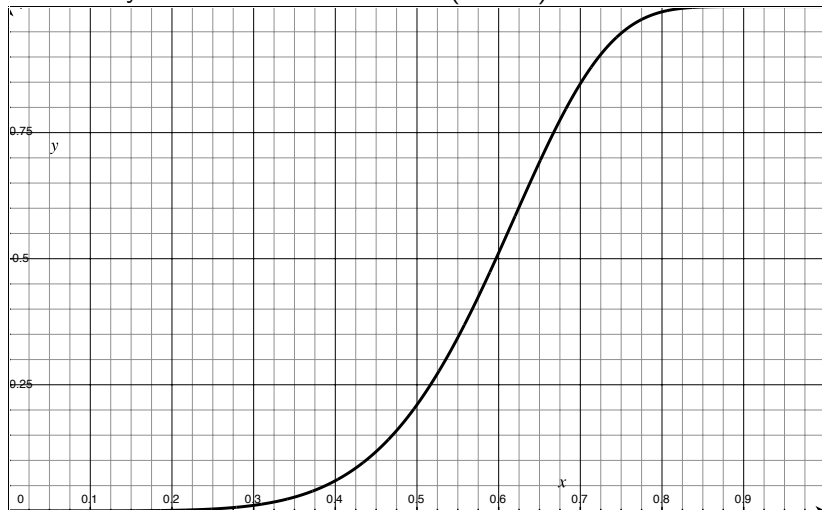Probability of found collision $= 1 - (1 - s^b)^r$

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

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

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

$t = 6 \cdot 15 = 90$

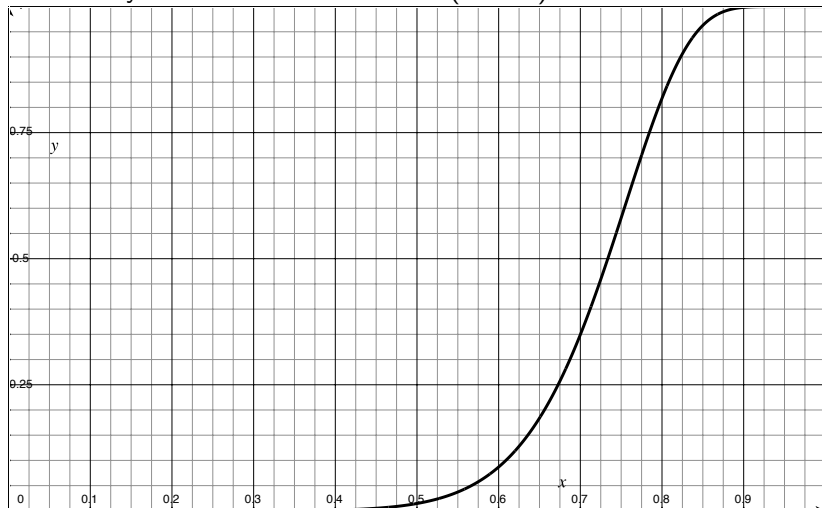Probability of found collision $= 1 - (1 - s^b)^r$

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

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

# LSH $b = 10$ and $r = 15$
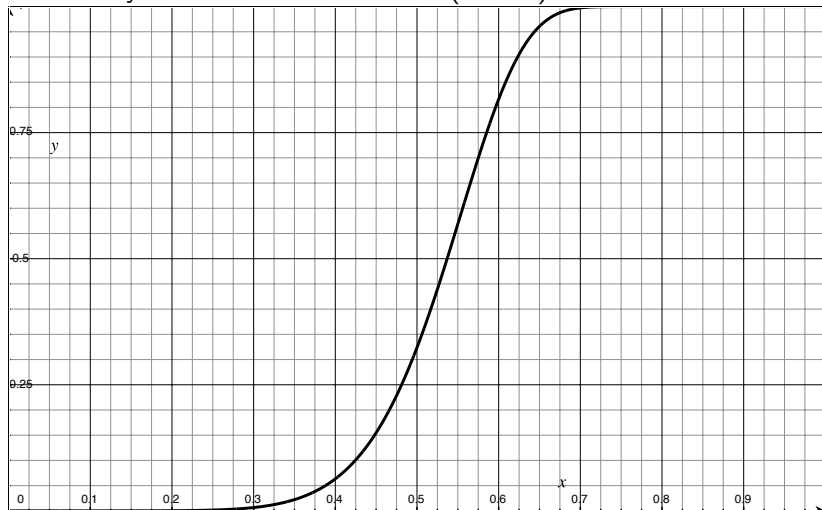
Probability of found collision $= 1 - (1 - s^b)^r$

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

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

# LSH $b = 8$ and $r = 100$
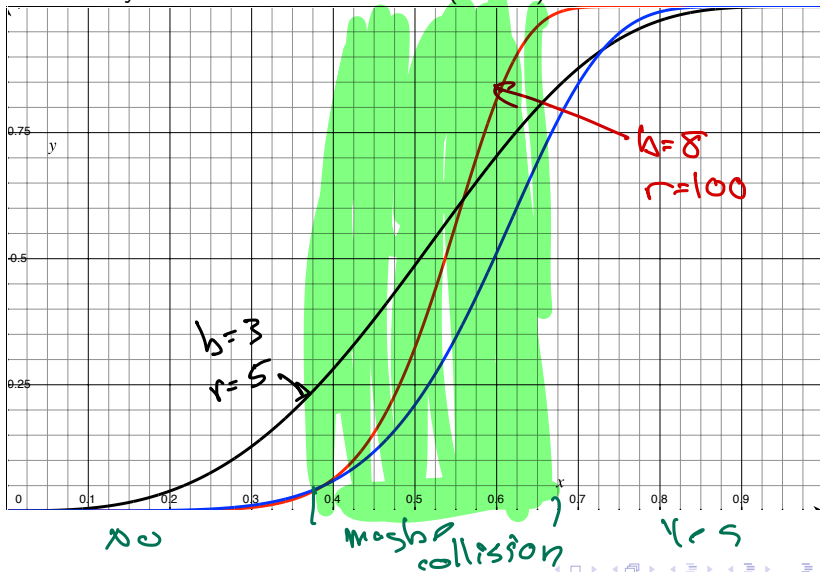
Probability of found collision $= 1 - (1 - s^b)^r$

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

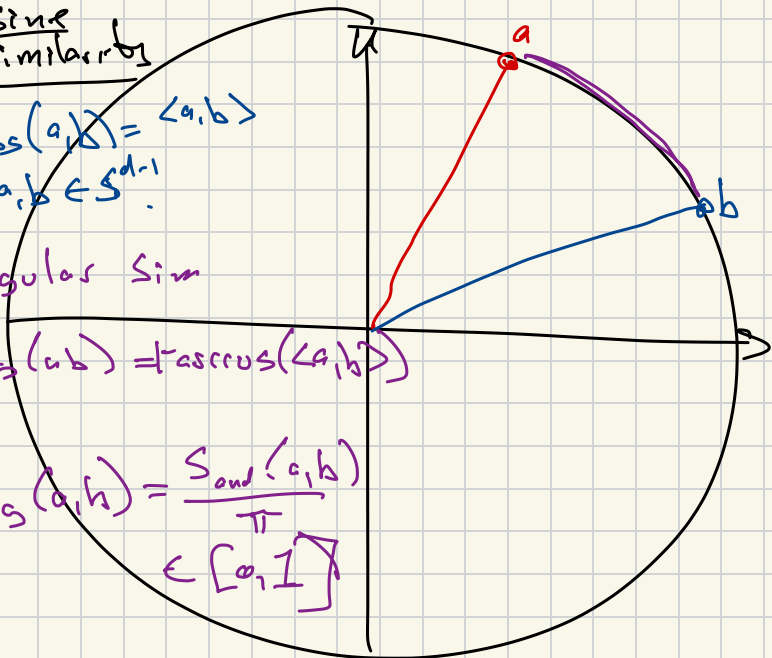# Cosine Similarity

$$S_{cos}(a,b) = \langle a, b \rangle$$

$$a, b \in S^{d-1}.$$

Angular Sim

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

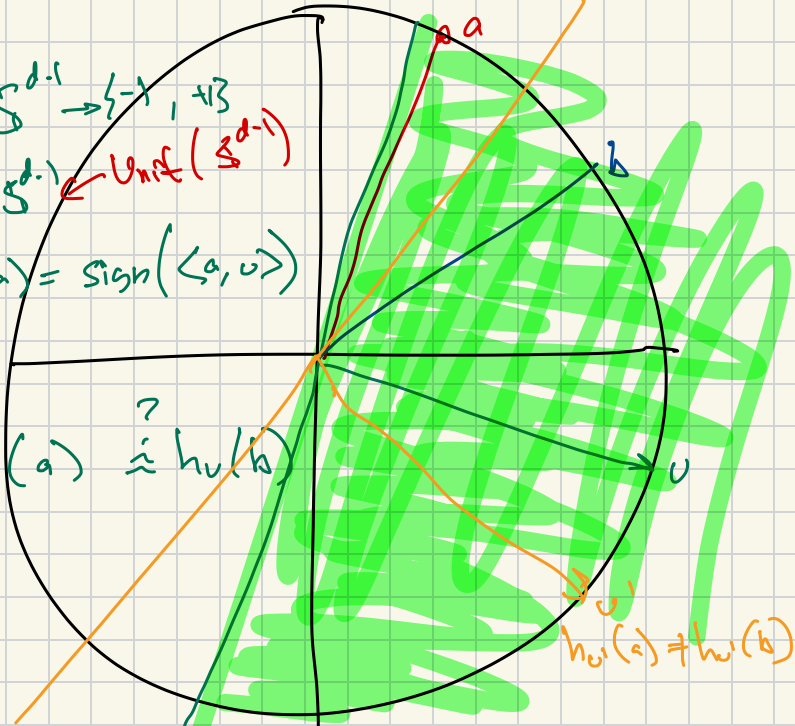$$\overline{S}_{ang}(a,b) = \frac{S_{ang}(a,b)}{\pi}$$

$$\in [0, 1]$$

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

$$u \in S^{d-1} \leftarrow \text{Unif}\left(S^{d-1}\right)$$

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

$$h_u(a) \overset{?}{=} h_u(b)$$

$a$

$b$

$u$

$u'$

$h_u(a) \neq h_u(b)$

$v \sim \text{Unif}\left(\, \$^{d-1}\right)$

proposal
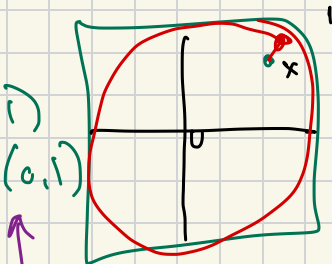
$x \in \text{Unif}[-1, 1] \times \text{Unif}[-1 \times 1] \times \ldots \text{Unif}(-1, 1)$

$u = \dfrac{x}{\|x\|}$

$g \sim N(0,1) \times N(0,1)$
$\times \ldots \times N(0,1)$

$v = \dfrac{g}{\|g\|}$

this
works



Box-Mueller Transform

# LSH for Euclidean in $\mathbb{R}^1$

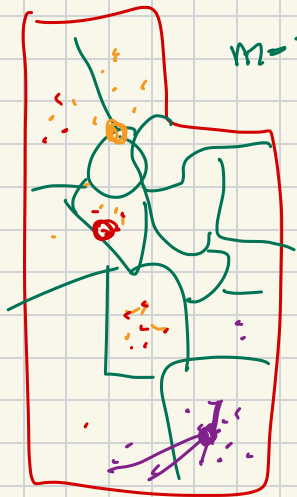threshold for similarity

$\vdash\!\!\!\!\dashv = \sigma$



$$Pr\left[h(a) = h(b)\right] = 1 - \frac{\|a-b\|}{\sigma} \quad \text{if } \|a-b\| < \sigma$$

$$= 0 \quad \text{otherwise}$$

# LSH for $\mathbb{R}^d$ Euclidean

$$h_{v,\beta}(a) = \left\lfloor -\beta + \langle a, v \rangle \right\rfloor$$

# Distances for Distributions
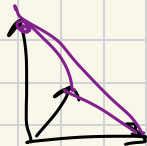


$m = 29$ counties

Map to counts rounds

vector $= x \in \mathbb{R}^m$

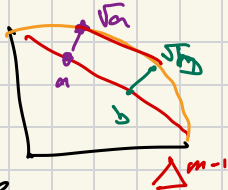$x_j = \dfrac{\# \text{ strikes in round } i}{N = \text{tot } \# \text{ strikes}}$

$x, x', x'' \in \triangle^{m-1} \subset \mathbb{R}^d$

$x_j \in [0, 1]$

$\sum_j x_j = 1$

$x, \acute{x} \in \Delta^{m-1}$

Kullbach-Leibler Divergence

$$D_{KL}(a,b) = D(a \| b)$$

$$= \sum_{j=1}^{m} a_j \ln\left(\frac{a_j}{b_j}\right)$$

Hellinger Dist

$$D_H(a,b) = \frac{1}{\sqrt{2}} \sqrt{\sum_{j=1}^{m}\left(\sqrt{a_j} - \sqrt{b_j}\right)^2}$$

$D$   $D'$



Wasserstein Distance

Earth Mover Dist

Find best matching $\gamma: D \to D'$

$$W_1(D, D') = \frac{1}{N} \sum_{x \in D} \| x - \gamma(x) \|$$

Runtime $N^3$

$D''$