# L5 Locality Sensitive Hashing

- Family of hash functions $\mathcal{H}$
  ↳ we want similar items to have a larger chance of collision.

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

1.1 hash function
$$\hat{J}s(p, q) = \begin{cases} 1 & h(p) = h(q) \\ 0 & ow \end{cases}$$

hash table ??

= Jaccard Triangle

$\approx$ Euclidean (dot products)

2. K hash functions

$$\hat{J}s_K(p, q) = \frac{1}{K} \sum_{j=1}^{K} I\{h_j(p) = h_j(q)\}$$
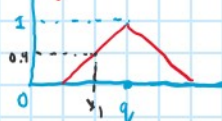
Algo??

## Large Number of Objects

$X = \{x_1, x_2, \ldots, x_n\}$  (maybe documents from k-grams)

Q1: which pairs are similar?  ($n^2$ time)

Q2: given a query 'q', which $x_i \in X$ are similar to q?  (n time)

- Pretend $x_1, \ldots, x_n \in \mathbb{R}$
  ↳ Then similarity $S_\Delta(q, x_i) = max\{0, 1 - |q - x_i|\}$
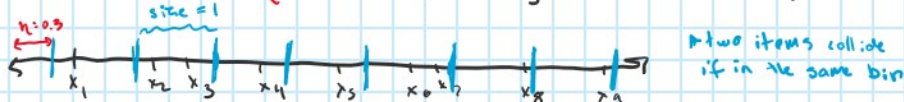


1. Sort $x_1, x_2, \ldots, x_n$
2. Build binary tree
3. Find q in tree

$n \in Unif(0,1)$    $h_n(x) \to a$ bin  Runtime: $\log(n) + K$ ≠ similar items



▸ two items collide if in the same bin

- It can be shown that $P[h(x) = h(x')] = S_\Delta(x, x')$

♢ Though we know there exists data structure solutions in 1D, we will explore the algorithmic solution for problems in higher dimensional spaces!
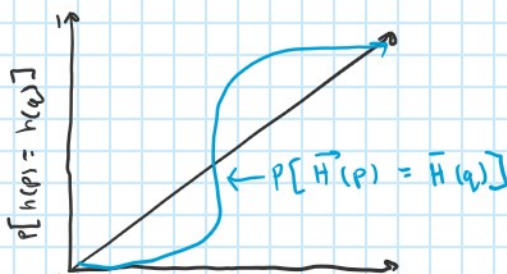
## Banding: How to combine hash functions...

- $H = \{h_1, h_2, \ldots, h_+\} \in \mathcal{H}$
- $\vec{H}$ ← 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 | 4 | 2 |
| $P_3$ | 0 | 2 | 4 | 3 | 5 | 1 |
| $P_4$ | | | | | | |



$\leftarrow P[\vec{H}(p) = \vec{H}(q)]$

$P_4$

$P_5$

$P_6$

q | (3) 5 1 2 3 2



1 2 3 $\overset{h_4}{4}$ 5 6 7

$h_2$ | 1 2 3 4 5 6 7

$P_1$

$P_1 = [3,5]$

$P_2 [h_1, h_2(p) = h_1, h_2(q)]$

$= s^2$

ie, much more selective

$\Rightarrow$



$\leftarrow P[\vec{H}(p) = \vec{H}(q)]$

AND

$s(p,q)$

| | $h_1$ | $h_2$ | | $h_3$ | $h_4$ | | $h_5$ | $h_6$ |
|---|---|---|---|---|---|---|---|---|
| $P_1$ | (3) | 5 | | 0 | 2 | | 4 | 3 |
| $P_2$ | 3 | 4 | | 1 | 0 | | 4 | 2 |
| $P_3$ | 0 | 2 | | 4 | 3 | | 5 | 1 |
| $P_4$ | | | | | | | | |
| $P_5$ | $H_1$ | | $H_2$ | | | $H_3$ | | |
| $P_6$ | | | | | | | | |
| q | (3) | 5 | | 1 | 2 | | 3 | 2 |

# of bands = r

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

$\Rightarrow$



OR

$\leftarrow P[\vec{H}(p) = \vec{H}(q)]$

AND

$s(p,q)$

## Analysis:

- r bands, each with b hash functions.

  ↳ t = # of hash functions.

    ↳ $t \geq r \cdot b$

- $S(p,q) = s$

  ↳ $s^b$ = Probability p,q collide in one band

- $(1 - s^b)$ = Probability p,q do not collide.

- $(1 - s^b)^r$ = Probability p,q don't collide in r bands.

- $f(s) = 1 - (1 - s^b)^r$ = Probability p,q collide in at least 1 band.

  ↳ r↑ increases the OR count (→ ⌐)
  ↳ b↑ increases the AND count (← ⌐)

  ✱ (time is used creating hash functions but lookup is then constant)

Q: How to choose values for $b$ and $r$ given $t, T$

    S1. Plot $f(s)$   (we want the most steep line)

    S2. $\begin{cases} b = -\log_T(t) \\ r = t/b \end{cases}$

## LHS For Euclidean Distance:

- $d_E(p,q) \Longleftrightarrow s_E(p,q)$

$$= \langle p, q \rangle$$
$$= \sum_{j=1}^{a} p_i \cdot q_i$$

- $h : \mathbb{R}^a \to [m]$   where:
    $\begin{cases} h \in \text{Unif}(0,T) & \text{(random shift)} \\ u \in \mathbb{R}^a & \text{s.t. } \|u\| = 1 \quad \text{(random unit vector)} \end{cases}$

- $h_{n,u}(P) = \left( \lfloor \langle p, u \rangle - h \rfloor \bmod m \right)$

    (projecting $\vec{p}$ onto $\vec{u}$)

    = dot product of $p, u$