

Nearest Neighbor Search

$$x_1 = [11011001]$$

← d bits

⋮

$$x_n = [01000010]$$

$$\text{query vector } q = [0011101]$$

← d bits

find x_i

$$\text{s.t. } x_i = \underset{?}{\operatorname{argmin}} \operatorname{dist}(x_i, q)$$

Hamming dist = # diff. bits.

$$\operatorname{dist}(x_1, q) = \begin{matrix} x_1 & [1] & [1] & [0] & [1] & [1] & [0] & [0] & [1] \\ q & [0] & [0] & [1] & [1] & [1] & [0] & [1] & [1] \end{matrix}$$

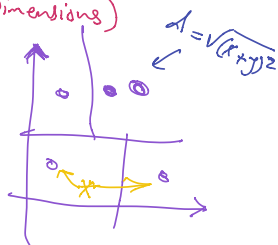
= 4

Find the nearest: $O(nd)$

Can we do it faster:
what data structures / Preprocessing
can help?

Range Trees (low Dimensions)

↳ works for
low dimensions
 $d = O(1)$
↳ is small



How about High Dimensions (HD)

* d is not small

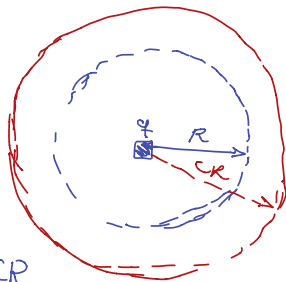
if $\exists x_i$ s.t.

$$\operatorname{dist}(q, x_i) \leq R$$

return

x_j s.t.

$$\operatorname{dist}(q, x_j) \leq CR$$



(X, d)

e.g.
↳ distance = Hamming dist.
↳ Points = Binary vectors
e.g.

A (Binary) hash maps X to
a (Binary) number (Bucket)

$$x \rightarrow \boxed{H} \rightarrow \begin{matrix} 1 \\ 0 \end{matrix} \quad H(x) = \begin{cases} 1 \\ 0 \end{cases}$$

e.g.: for binary vectors

$$H(x) = \begin{cases} 1 & \text{if } x[i] = 1 \\ 0 & \text{if } x[i] = 0 \end{cases}$$

Given a System (X, d) , a hashing
 H belongs to the class of

Local Sensitive Hashing (LSH)

if the following hold for
a pair of points (x, y)

$$\text{if } d(x, y) \leq R$$

$$P(H(x) = H(y)) \geq P_1$$

$$\text{if } d(x, y) \geq CR$$

$$P(H(x) = H(y)) \leq P_2$$

$$P_1 > P_2$$

e.g.

$$H(x) = x_j : \text{The } j\text{-th bit of } x$$

$$P(H(x) = H(y)) = 1 - \frac{d(x, y)}{d}$$

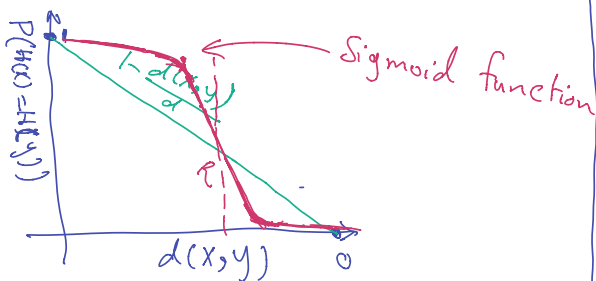
$$\text{if } d(x, y) \leq R$$

$$P(H(x) = H(y)) \geq 1 - \frac{R}{d} = P_1$$

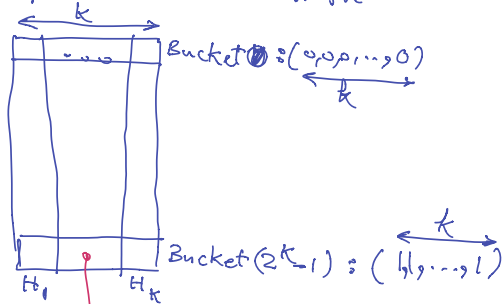
$$\text{if } d(x, y) \geq CR$$

$$P(H(x) = H(y)) \leq 1 - \frac{CR}{d} = P_2$$

$$P_1 > P_2 \quad \checkmark \in \text{LSH}$$



Step 1: define k -Hash



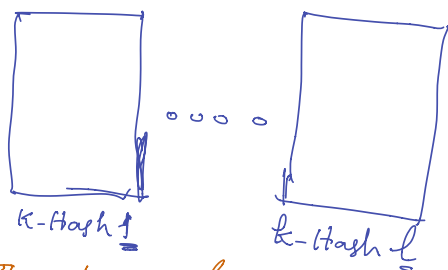
H_i is an indep. Hash function

$$P(\text{Bucket}(X) = \text{Bucket}(Y)) = \left(1 - \frac{d(X, Y)}{d}\right)^k$$

→ Reduced the chance of false Positive

→ Increased the false Negative

Step 2: Consider l k -Hash Buckets

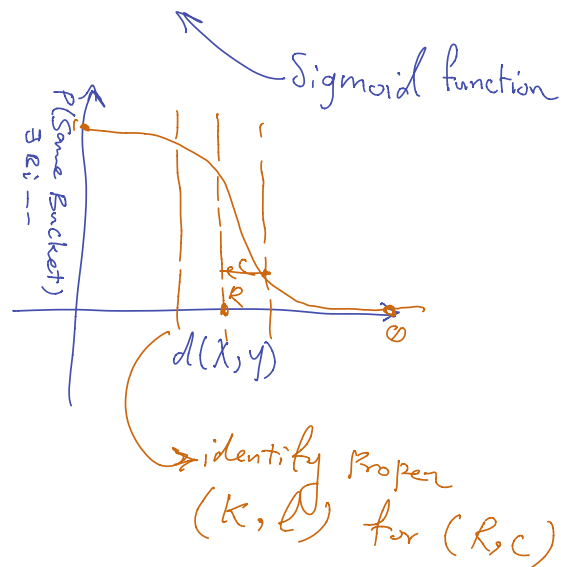


The chance of similar points to fall into the same bucket in at least one of k -Hashes is high

$$P(B_i(X) = B_i(Y)) = \left(1 - \frac{d(X, Y)}{d}\right)^k$$

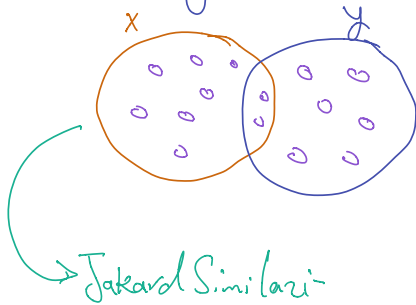
$$P(\exists B_i \text{ s.t. } B_i(X) = B_i(Y)) = \left(1 - \left(1 - \frac{d(X, Y)}{d}\right)^k\right)^l$$

$$\Rightarrow P(\exists B_i \text{ s.t. } B_i(X) = B_i(Y)) = 1 - \left(1 - \left(1 - \frac{d(X, Y)}{d}\right)^k\right)^l$$



Space Complexity:
 $O(ul)$

Set Similarity



Jakard Similarity

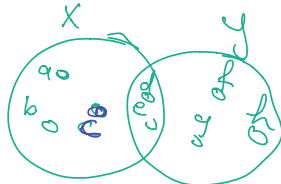
$$S(X, Y) = \frac{|X \cap Y|}{|X \cup Y|} \leftarrow J\text{-Sim}$$

$$d(X, Y) = 1 - \frac{|X \cap Y|}{|X \cup Y|} \leftarrow J\text{-dist.}$$

Estimate the distance (Similarity) of two Sets (X, Y)

Hashing: Randomly shuffle the elements
(Take a random order of elements)
 \Rightarrow The index of the first element in a set is its hash

e.g.



$H: \begin{array}{|c|c|c|c|c|c|c|} \hline c & d & h & f & a & g & b \\ \hline 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline \end{array}$ random order

$$H(X) = H(c) = 0$$

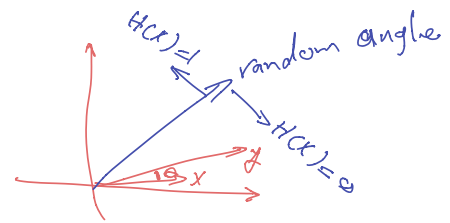
$$H(Y) = H(d) = 1$$

$$P(H(X) = H(Y)) = \frac{|X \cap Y|}{|X \cup Y|} : J\text{-Sim}$$

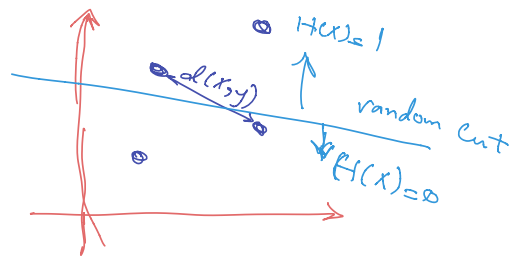
... \rightarrow LSH

Every k -hash for sets is called bottom-k sketches

dist: Cosine Sim.



dist: l_2 (Euclidean dist.)



LSH for dimension Reduction
 $R^d \rightarrow R^k$ s.t. distances are Preserved.

\rightarrow Simply create a k -hash from a LSH family