# Nearest Neighbor Search

$X_1 = [1 1 0 0 \cdots 0 1]$
$\xleftarrow{\hspace{1cm}} d \xrightarrow{\hspace{1cm}}$

$X_2 = [0 1 0 1 0 \cdots 1 0]$

$\vdots$

$X_n = [1 0 1 1 0 \cdots 1 0]$

**Given** ⟵

a query vector

$q = [1 1 0 \cdots 1]$
$\xleftarrow{\hspace{1cm}} d \xrightarrow{\hspace{1cm}}$

Objective:

find $X_i$

s.t. $X_i = \underset{\forall X_j}{\arg\min} \; dist(q, X_j)$

(assume Hamming distance is Considered)

Brute force: $O(nd)$

---

Question: what Type of Datastructures (indices) to build, So that queries Can be answered faster?
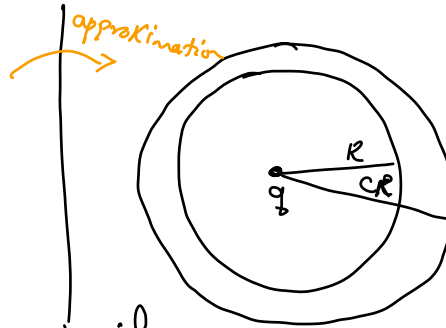
For Hamming Dist.: Range Trees.

Euclidean dist: Voronoi Diag.
↳ $O(d \log n)$

Challenge: These indices work for low-D i.e., $d$ is a Small Constant

---

How to answer NN queries in High D Spaces?
↳ $d$ is not Small

---

*approximation* →



if $\exists \; X_i$ s.t. $dist(q, X_i) \leqslant R$

return any point $X_j$ where

$dist(q, X_j) \leqslant CR$

\* $i \in [1, n]$
$j \in [1, n]$

*Randomized Index* →

---

• A Binary hash $H: X \twoheadrightarrow \{0,1\}$ is a function that given an input $X$ maps it to a Binary value $H(X)$

$X \rightarrow \boxed{H} \nearrow^{1} \searrow_{0}$

• Given a System $(X, d)$
↳ *Universe of Possible inputs* ↳ *distance function*

a binary hash belongs to the class **Local Sensitive hash (LSH)** if

$\forall \; dist(X, y) \leqslant R:$

$\qquad P(H(X) = H(Y)) \gtrsim P_1$

$\forall \; dist(X, y) \gtrsim CR$

$\qquad P(H(X) = H(Y)) \leqslant P_2$

AND $P_1 > P_2$

for binary vectors:

$H(X) = X_\ell$ : The $\ell$-th bit of vector $X$

↳ a random index

e.g.
$$\begin{array}{ccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 \end{array}$$
$$X = [1 \; 0 \; 0 \; 1 \; 0 \; 1 \; 1]$$
$$Y = [0 \; 0 \; 1 \; 1 \; 0 \; 1 \; 1]$$

$H(X) = X_3$

$$H(X) = 0 \quad , \quad H(Y) = 1$$

$$P(H(X) \neq H(Y)) = \frac{dist(X,y)}{d}$$

$$\Rightarrow P(H(X) = H(y)) = 1 - \frac{dist(X,y)}{d}$$

$\forall \, dist(X,y) \leqslant R$

$$\frac{dist(X,y)}{d} \leqslant R/d$$

$$P(H(X) = H(y)) = 1 - \frac{dist(X,y)}{d} \geqslant 1 - R/d \quad \overset{}{\underset{P_1}{\longleftarrow}}$$

$\forall \, dist(X,y) \geqslant CR$

$$P(H(X) = H(y)) \leqslant 1 - \frac{CR}{d} \overset{P_2}{\longleftarrow}$$

$$P_1 - P_2 \overset{?}{>} 0$$

$$P_1 - P_2 = 1 - R/d - 1 + \frac{CR}{d}$$

$$= \frac{(C-1)R}{d} > 0 \quad \checkmark$$

$H \in LSH$



Goal: Using binary hashes from the class of LSH, build a (non-binary) hash that satisfies:

✓ $dist(X,y) \leqslant R$
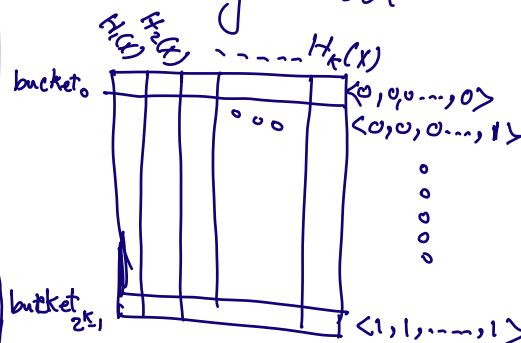$$P(H(X) = H(y)) \approx 1$$

✓ $dist(X,y) \geqslant CR$
$$P(H(X) = H(y)) \approx 0$$

Step 1: Build a bucket of binary hash



$$H(X) = \langle H_1(X), H_2(X), \cdots, H_K(X) \rangle$$

e.g. $X = [\overset{1}{0} \; \overset{2}{0} \; \overset{3}{0} \; 1 \; 1 \; 0 \; 1 \; 1 \; 1]$
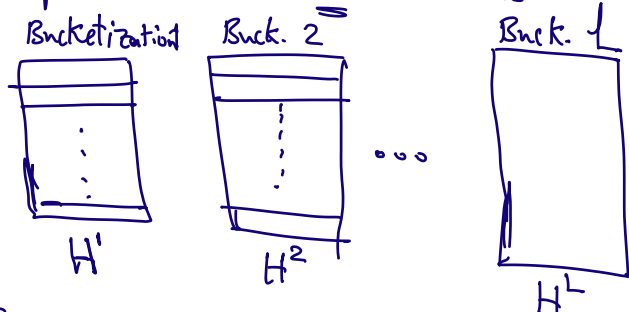$Y = [0 \; 1 \; 1 \; 1 \; 0 \; 1 \; 1 \; 0]$

$H_1 = X_2 \, , \quad H_2 = X_5$

$H(X) = \langle X_2, X_5 \rangle = \langle 0, 1 \rangle$

$H(y) = \langle y_2, y_5 \rangle = \langle 1, 0 \rangle$

P( X and y fall into the Same bucket) = P(H(y) = H(x))

$$= \left(1 - \frac{dist(x,y)}{d}\right)^k$$

**Step 2 : Build L Buckets**

Bucketization   Buck. 2                    Buck. L



H¹              H²              . . .        $H^L$

**Step 3 :**

given the Query point q,
return any $x_j$ that falls into
the bucket of q in at least
one of the bucketizations

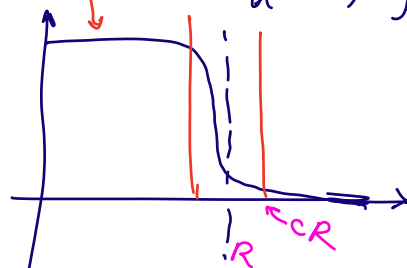$P(\exists B_i \text{ s.t. } H^i(x) = H^i(y)) =$

$P\left((H^1(x) = H^1(y)) \vee (H^2(x) = H^2(y)) \vee \cdots (H^L(x) = H^L(y))\right)$

$$= 1 - P\left(\bigwedge_{i=1}^{L} (H^i(x) \neq H^i(x))\right)$$

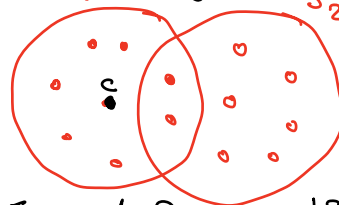$$= 1 - \left(1 - \left(1 - \frac{dist(x,y)}{d}\right)^k\right)^L$$

← A Sigmoid function

---

$$1 - \left(1 - \left(1 - \frac{dist(x,y)}{d}\right)^k\right)^L$$



.R    ←cR

**find Proper k, L**

---

Binary Vectors:  LSH – B – Hash:
                 a random bit of
                 vector

Set Similarity (distance)
e.g.   $S_1$           $S_2$



Jaccard Sim. $= \frac{|S_1 \cap S_1|}{|S_1 \cup S_2|}$

Jaccard dist. $= 1 - \frac{|S_1 \cap S_1|}{|S_1 \cup S_2|}$

LSH – B – Hash: a random element
                from the Union

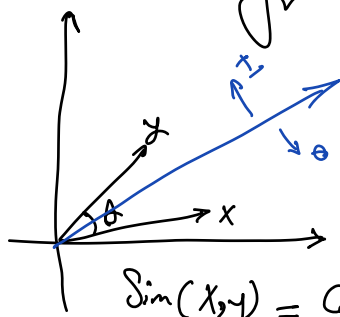$$H_e(S) = \begin{cases} 1 & e \in S \\ 0 & e \notin S \end{cases}$$

e.g.:
$H_c(S_1) = 1, \quad H_c(S_2) = \emptyset$

$P(H_e(S_1) = H_e(S_2))$
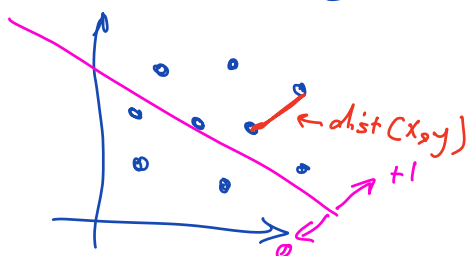
$$= \frac{|S_1 \cap S_2|}{|S_1 \cup S_2|}$$

## Cosine Similarity



$$Sim(X,y) = Cos(\theta)$$

LSH-B-Hash: a random Vector, Partitioning the Space in two "half Spaces"

## Euclidean Distance



$\leftarrow dist(X,y)$

$+1$

$0$

LSH-B-Hash: A random Line (HyperPlane) in the Space

---

## LSH for Dimension Reduction

$$\mathbb{R}^d \xrightarrow{\;t\;} \mathbb{R}^k$$

$$X, y \in \mathbb{R}^d$$

$$dist(X,y) \overset{\Delta}{\simeq} dist(t(X), t(y))$$

$\rightarrow$ Construct a LSH of $k$ binary hash

Intuitively: Since LSH maintains the distances relatively fixed