

Last time

- short reviews from probability
 - variance & tail bounds
 - Gaussians & CLT
- Distinct elts
- Similarity search & dimension reduction

Today

Locality sensitive hashing (LSH)
 => efficient approx. similarity search

Similarity search, NNS

large collection of data items

- $x^{(1)}, x^{(2)}, \dots, x^{(n)} \in \mathbb{R}^k$
- notion of distance (or equivalently similarity)
 $d(x, y) = \text{dist between } x \text{ \& } y$

2 problems:

① Find all pairs i, j s.t.
 $d(x^{(i)}, x^{(j)}) \leq r$

Naive:
 $O(n^2 k)$

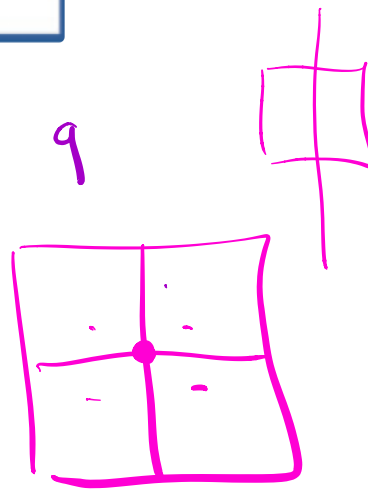
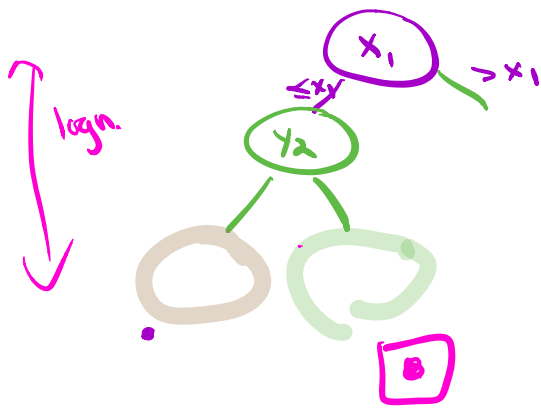
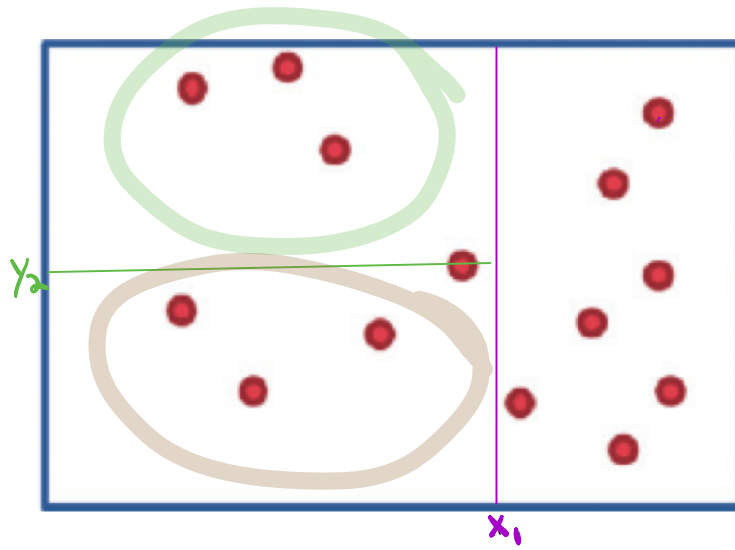
② Preprocess database & then efficiently respond to queries

★ query $q \in \mathbb{R}^k$
 return all pts $x^{(i)}$ in DB s.t.
 $d(q, x^{(i)}) \leq r$

$O(nk)$ / query

If k small, say < 15 , there are space partitioning data structures that are reasonably efficient

k-d trees

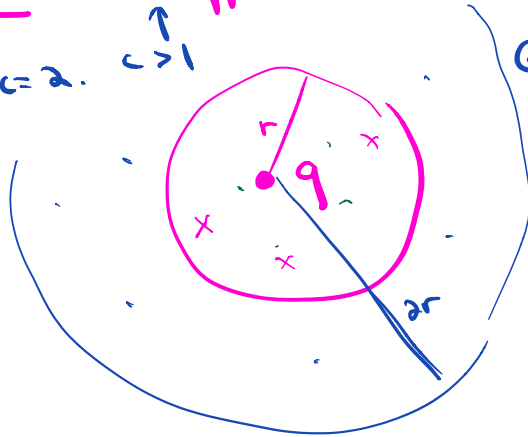


curse of dimensionality.

LSH

c-approx r-nearest neighbors problem

think $c=2$.



Given query pt q , return

- all pts $x^{(i)}$ s.t. $d(x^{(i)}, q) \leq r$ (w.h.p.)
- may return some pts $x^{(i)}$ s.t. $d(x^{(i)}, q) \leq c \cdot r$

Primitive

\mathcal{H} : family of hash fns that map pts $\in \mathbb{R}^k \rightarrow$

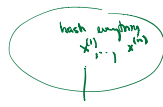
\mathcal{H} is (r, cr, p_1, p_2) -sensitive if $\forall x, y \in \mathbb{R}^k$

If $d(x, y) \leq r \Rightarrow \Pr_{h \in \mathcal{H}} (h(x) = h(y)) \geq p_1$

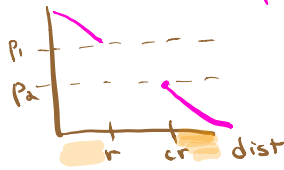
If $d(x, y) \geq cr \Rightarrow \Pr_{h \in \mathcal{H}} (h(x) = h(y)) \leq p_2$

$p_2 < p_1$

$d(x, y) = \begin{cases} 0 & \text{identical } x=y \\ \infty & \text{o.w.} \end{cases}$



} defn.



Example: Suppose pts $\in \{0, 1\}^k$

$\{0, 1\}^k \rightarrow \{0, 1\}$

$\mathcal{H} = \{h(x) = x_i \mid 1 \leq i \leq k\}$

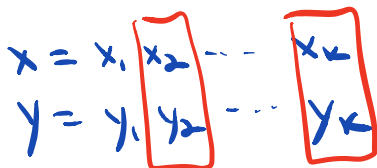
$d(x, y) = \text{Hamming dist} = \sum_{i=1}^k \mathbb{1}_{x_i \neq y_i}$

(# bits where different)

$d(x, y) \leq r$
 $\underline{p_1} = \Pr(h(x) = h(y)) \geq 1 - \frac{r}{k} \approx e^{-\frac{r}{k}}$

$d(x, y) \geq cr$

$\underline{p_2} = \Pr(h(x) = h(y)) \leq 1 - \frac{cr}{k} \approx e^{-\frac{cr}{k}}$



$(r, cr, 1 - \frac{r}{k}, 1 - \frac{cr}{k})$ family of hash fns.

combine these hash fns to amplify difference between p_1 & p_2

$$h(x) \in \{0,1\}^d$$

$$g(x) \in \{0,1\}^d$$

① $g(x) = [h_1(x), h_2(x), \dots, h_\ell(x)]$

Purpose: reduce chance that pts that are far away map to same value.

ANO step $q, y \quad d(q, y) \geq cr$

$$\Pr(g(q) = g(y)) \leq p_2^d$$

② $g_1(x) = [h_{11}(x), h_{12}(x), \dots, h_{1\ell}(x)]$

$$g_2(x) = [h_{21}(x), h_{22}(x), \dots, h_{2\ell}(x)]$$

$$\vdots$$

$$g_\ell(x) = [h_{\ell 1}(x), h_{\ell 2}(x), \dots, h_{\ell \ell}(x)]$$

③ On query q , compute $g_1(q), \dots, g_\ell(q)$

OR $S_q = \{x^{(i)} \mid g_j(x^{(i)}) = g_j(q) \text{ for some } 1 \leq j \leq \ell\}$

④ compute $d(q, x) \quad \forall x \in S_q$

output all of those that have $d(q, x) \leq r$

$$d(x, q) < r$$

$$\Pr(x \notin S_q) =$$

$$\Pr(g_1(x) \neq g_1(q), g_2(x) \neq g_2(q), \dots, g_\ell(x) \neq g_\ell(q))$$

$$= (\Pr(g_1(x) \neq g_1(q)))^\ell$$

$$= 1 - \Pr(g_1(x) = g_1(q))^\ell$$

$$= 1 - \Pr(h_{11}(x) = h_{11}(q), \dots, h_{1\ell}(x) = h_{1\ell}(q))^\ell$$

$$\leq (1 - p_1^d)^\ell \approx e^{-p_1^d \cdot \ell}$$

$$d = \frac{\log n}{\log(\frac{1}{p_2})}$$

$$\ell = n^2$$

where we are:

- preprocessing time: $n \cdot \ell \cdot d$ hash fn computations

- space: $n \cdot \ell + \text{actual pts}$

- exp time to process query:

ℓd (hash fn computations)

+ $E(\# \text{ far } > cr \text{ pts that are in } S_q)$

$$\leq \ell \cdot n \cdot p_2^d$$

- Prob miss a close pt $\Rightarrow (1 - p_1^d)^\ell$

\mathcal{X} is (r, cr, p_1, p_2) -sensitive $\forall x, y \in \mathbb{R}^k$

If $d(x, y) \leq cr \Rightarrow \Pr_{\text{next}}(h(x) = h(y)) \geq p_1$

If $d(x, y) \geq cr \Rightarrow \Pr_{\text{next}}(h(x) = h(y)) \leq p_2$

dist fn.

start w/ r, c

get $\mathcal{X}(r, cr, p_1, p_2)$

need to select

- d
- ℓ

$$\text{set } n p_2^d = 1$$

$\equiv \leq 1$ bad pt per table

$$d = \frac{\log n}{\log(\frac{1}{p_2})}$$

$$\Pr(\text{missing close pt}) = \frac{1}{2}$$

$$p_1^d \cdot \ell = 1$$

$$\ell = n \cdot \frac{\log(\frac{1}{p_1})}{\log(\frac{1}{p_2})}$$

$$\frac{\log(\frac{1}{p_1})}{\log(\frac{1}{p_2})} \triangleq \rho$$

$$p_1 = \frac{1}{2} \quad p_2 = \frac{1}{4}$$

$$l = n^{\frac{1}{2}}$$

Where we are:

- preprocessing time: $n \cdot l \cdot d$ hash fn computations $\leftarrow n \cdot n^{\frac{1}{2}} \cdot \log n / \log(\frac{1}{p_2})$
 $n^{1.5} + \text{pts}$

- space: $n \cdot l + \text{actual pts}$

- exp time to process query:
 ld (hash fn computations)
 $+ E(\# \text{ far pts that are in } S_q)$

$$= n^{\frac{1}{2}} \frac{\log n}{\log(\frac{1}{p_2})} + n^{\frac{1}{2}}$$

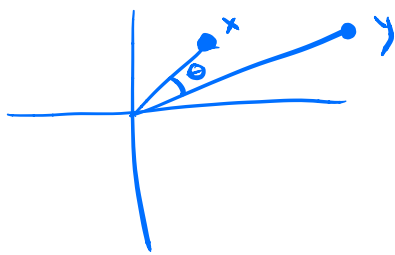
$$\leq l \cdot n \cdot p_2^d$$

- Prob miss a close pt $\rightarrow 1$
 $\Rightarrow (1 - p_1^d)^k$

$$\frac{1}{e}$$

Cosine similarity

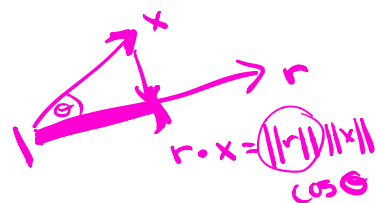
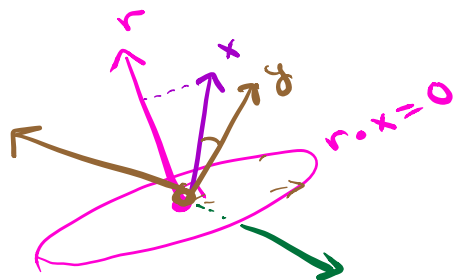
$$x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^k$$



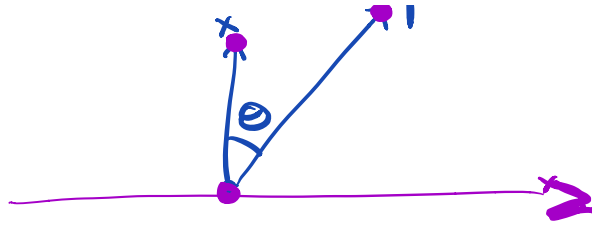
similar (close) if $\theta \leq \theta_1$ $\geq p_1$
 for $\theta \geq \theta_2$ $\leq p_2$

$$h: \mathbb{R}^k \rightarrow \{0, 1\}$$

$$h(x) = \text{sign}(r \cdot x) \mid \begin{matrix} r = (r_1, \dots, r_k) \\ r_i \sim N(0, 1) \\ \text{indep} \end{matrix}$$



$$\text{sign}(r \cdot x) = \begin{cases} 1 & \text{if } r \cdot x \geq 0 \\ 0 & \text{if } r \cdot x < 0 \end{cases}$$



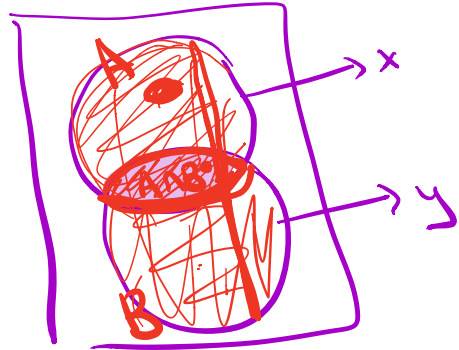
Jaccard Similarity

$$x = (x_1, \dots, x_n)$$

$$x_i = \begin{cases} \# \text{ occurrences} & \text{word } i \text{ is in doc} \\ 0 & \text{o.w.} \end{cases}$$

$$J(x, y) = \frac{\sum_i \min(x_i, y_i)}{\sum_i \max(x_i, y_i)}$$

$$= \frac{|\text{intersection}|}{|\text{union}|}$$



$$x = (0, 0, 1, 0, 0, 1, 0)$$

$$\pi = 5, 3, 2, 4, 7, 1, 6$$

$$h_\pi(x) = 0, 1$$

$$\Pr(h_\pi(x) = h_\pi(y)) = \frac{|A \cap B|}{|A \cup B|}$$

$$\pi^l = 7, 2, 4, \dots$$

$$h_{\pi^l}(x) = 4$$