

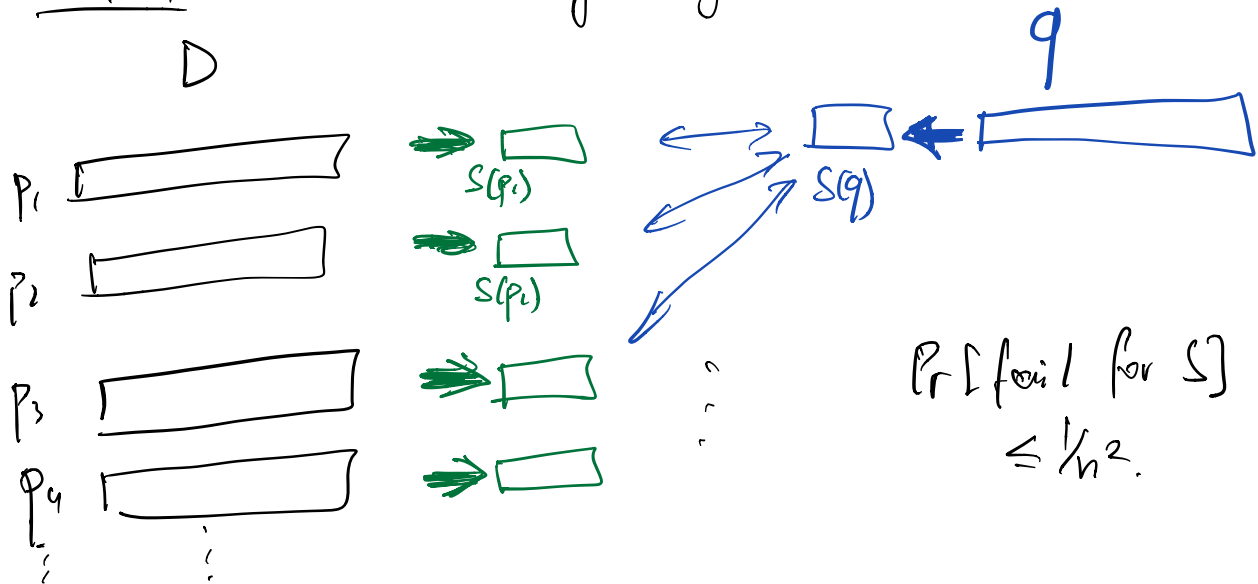
AA Lecture 9 (2/9/21)

Nearest Neighbor Search

c-ANN: for any $r > 0$, given $D \subset \mathbb{R}^d$ (or $\{0,1\}^d$), preprocess D , s.t., for given query q

- if $\exists p^* \in D, \|p^* - q\| \leq r$, output p with $\|p - q\| \leq cr$.

Sol 1: via sketching (generalization of DR)



Thm [Kushilevitz - Ostrovsky - Rabani '98]

sketch for $\{0,1\}^d$: $\exists \theta > 0$ s.t. $\forall r \geq 1$,

\exists distribution over $\varphi: \{0,1\}^d \rightarrow \{0,1\}^k, k = O(\frac{dn}{\epsilon^2})$

s.t.: - if $\|p - q\| \leq r \Rightarrow \Pr[\|\varphi(p) - \varphi(q)\|_1 \leq \theta k] \geq 1 - \frac{1}{n^3}$

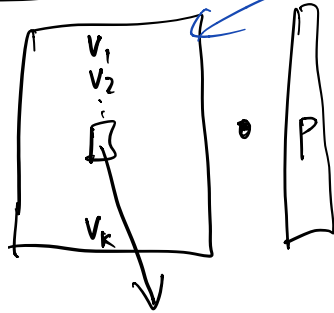
- if $\|p - q\| > cr \Rightarrow \Pr[\|\varphi(p) - \varphi(q)\|_1 > \theta k] \geq 1 - \frac{1}{n^3}$

where $c = 1 + O(\epsilon)$.

like DR for scale r .

proof [sketch]:

$$\varphi(p) =$$



size = $\underline{d \cdot k}$

$$+ \Rightarrow \oplus \mathbb{F}_2$$

= 1 with prob. $d = \frac{1}{2r}$.

$$v_i \in \{0,1\}^d, \quad i=1 \dots k.$$

$$\varphi(p) = (v_1 \cdot p, v_2 \cdot p, \dots, v_k \cdot p).$$

$$\|\varphi(p) - \varphi(q)\|_1 = \sum_{i=1}^k \mathbb{1}[v_i \cdot p \neq v_i \cdot q].$$

$$\mathbb{E}[\|\varphi(p) - \varphi(q)\|_1] = \mathbb{E}[\sum_{i=1}^k \mathbb{1}[v_i \cdot p \neq v_i \cdot q]] = k \cdot \Pr[v_i \cdot p \neq v_i \cdot q].$$

$$\Pr[v_i \cdot p \neq v_i \cdot q] = \Pr\left[\bigoplus_{i=1}^d u_i \cdot p_i \neq \bigoplus_{i=1}^d u_i \cdot q_i\right]$$

u_i
each coord
= 1 w/Pr = d

$$= \Pr\left[\bigoplus_{i=1}^d u_i \cdot (p_i \oplus q_i) \neq 0\right]$$

$u_i=1, Pr=d$ $\sum_{i=1}^d u_i \neq 0$ is $= \|p - q\|_1$

= Prob that $\|p-q\|$, random q_i s
 where $\Pr=1$ is d , xor to $\neq 0$.

$$= 1 - \Pr_{u_i} \left[\bigoplus_{i=1}^d u_i \cdot (p_i \oplus q_i) = 0 \right]$$

u_i : with Prob = $2d$, $u_i \in \{0,1\}$
 $= 0$ otherwise.

← "randomizing coordinates"

$$= 1 - (1-2d)^{\|p-q\|} - (1 - (1-2d)^{\|p-q\|}) \cdot \frac{1}{2}$$

$$= \frac{1}{2} (1 - (1-2d)^{\|p-q\|}).$$

$$\boxed{d = \frac{1}{2r}}$$

1) if $\|p-q\| < r \Rightarrow \Pr[u_p \neq u_q] \leq \frac{1}{2} (1 - (1 - \frac{2}{2r})^r)$

$$\approx \frac{1}{2} (1 - e^{-1}) = \frac{1}{2} - \frac{1}{2e}.$$

$(1-x) \approx e^{-x}$
 for x small.

2) if $\|p-q\| \geq (1+\epsilon)r \Rightarrow$

$$\Pr[u_p \neq u_q] \geq \frac{1}{2} (1 - (1 - \frac{2}{2r})^{r \cdot (1+\epsilon)})$$

$$\approx \frac{1}{2} (1 - e^{-1-\epsilon})$$

$$\approx \frac{1}{2} (1 - e^{-1} \cdot (1 - \epsilon))$$

$$\approx \frac{1}{2} - \frac{1}{2e} + \epsilon \cdot \frac{1}{2e}$$

Using Chernoff bounds (like HW2/P3 hint),
can prove that: 1) $\frac{\| \varphi(p) - \varphi(q) \|}{k} \leq \frac{1}{2} - \frac{1}{2e} + \frac{\epsilon}{4e}$

2) $\| - \| \geq \frac{1}{2} - \frac{2}{2e} + \frac{\epsilon}{2e} - \frac{\epsilon}{5e}$

(with prob. $\geq 1 - 1/n^3$). □

Corollary: can solve $(1 \leq \epsilon) - AMN$ in $\{0,1\}^d$
using spaces $O(n \cdot k) = O(n \cdot \frac{\lg n}{\epsilon^2})$.

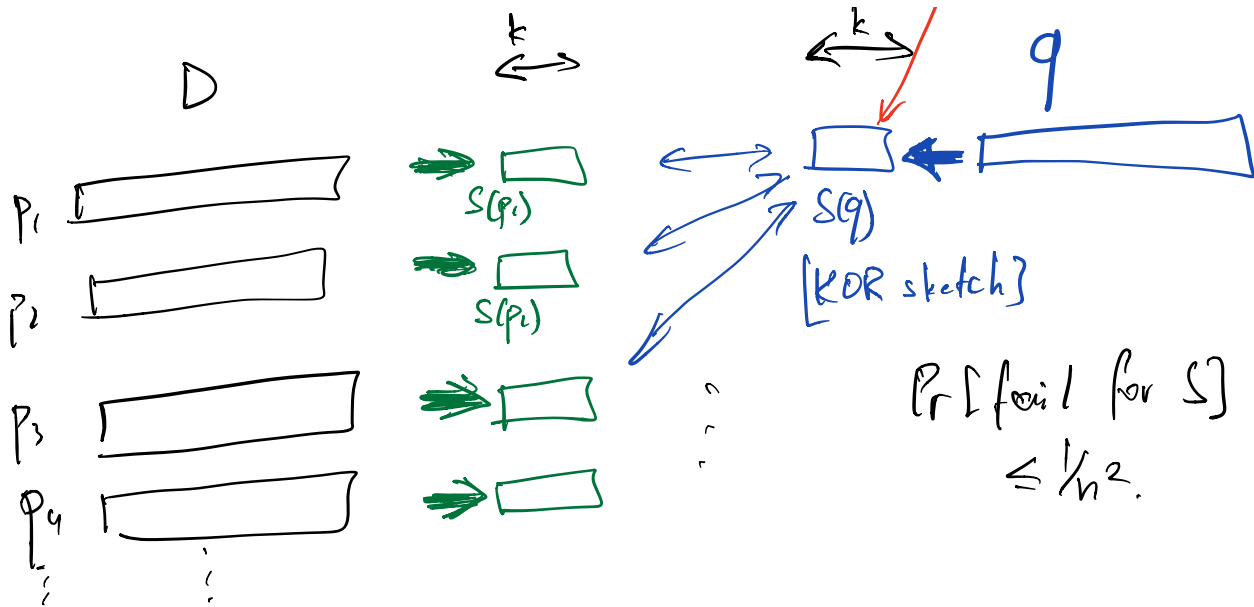
q.t.: $O(d \cdot k) + O(n \cdot \frac{\lg n}{\epsilon^2})$
[compute]
 $\varphi(q)$

Sol 2: space: $n^{O(1/\epsilon^2)}$.

q.t.: $O(d \cdot \frac{\lg n}{\epsilon^2})$.

in $\{0,1\}^d$.

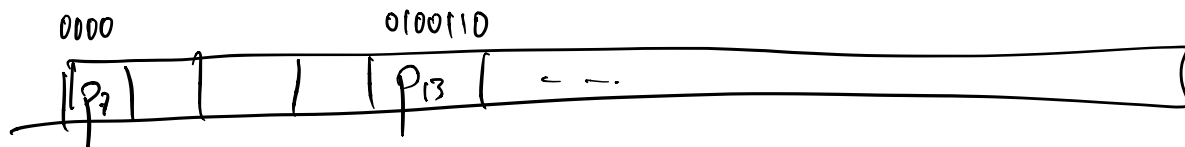
binary string
on length
k.



$$k = O\left(\frac{\log n}{\epsilon^2}\right).$$

Idea: $S(q) = \varphi(q) \in \{0,1\}^k$, $k =$ is all the info we need from q to determine sol to ANN.

Algorithm: for each possible string $\sigma \in \{0,1\}^k$, prepare what would be the ANN answer if $S(q) = \sigma$.



space: $2^k = 2^{O(\frac{\lg n}{\epsilon^2})} = n^{O(1/\epsilon^2)}$ (cells).

q.f.: $O(dk) + O(1) = O\left(\frac{\lg n}{\epsilon^2}\right)$

compute $S(q)$

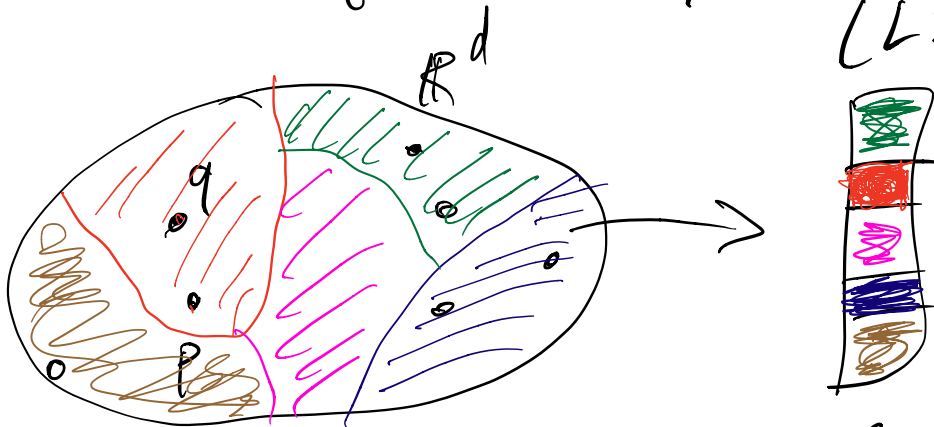
for $\epsilon < 1/2$.

Sol 3: space: closer to lin. in n

q.f.: n^β , $\beta < 1$.

still sub-lin. in n .

Idea: hashing: Locality-Sensitive Hashing (LSH).



Def: fix $r > 0$, approx c . Family \mathcal{H} of $h: \mathbb{R}^d \rightarrow \mathcal{U}$ is \checkmark LSH

if: $\|p - q\| \leq r$, $\Pr_{h \in \mathcal{H}} [h(p) = h(q)] \geq P_1$

- $\|p - q\| \geq cr$, $P \cap L \rightarrow \{ \} \subseteq P_2$.

Super-opt: $\exists H$ is $(r, cr, \underline{1}, \frac{1}{n})$ -LSH.

Algo based on LSH:

Preproc: = compute $h(p)$, $p \in D$

- store in a table (dictionary data struct.)

@ query q : = compute $h(q)$.

- look-up all $p \in D$ s.t.

$$h(q) = h(p)$$

- iterate through these p 's until we find one that is

$$\|p - q\| \leq cr.$$

Proof of corr/analysis:

\hookrightarrow : the bucket $h(q)$ must contain

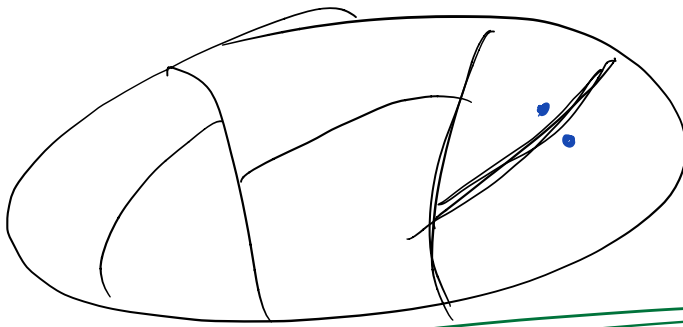
the point p^* (p^* s.t. $\|q - p^*\| \leq r$)

q.t.: $\mathbb{E}[\text{q.t.}] \leq D(f) + \mathbb{E}[\# \text{ pts } p \in D]$

$$\|p - q\| > cr, \\ h(p) = h(q) \} = d$$

$$\leq O(1) + n \cdot P_2 \cdot d \\ = O(dn \cdot P_2) = O(d).$$

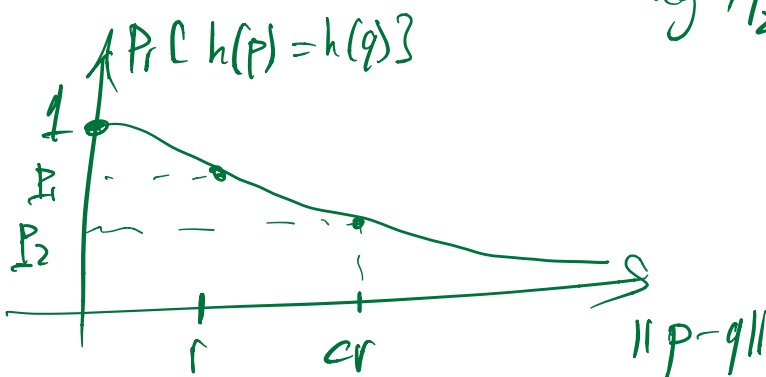
Counter: $(r, cr, \frac{1}{n})$ - LSH is not possible.



P_1 must be < 1 .

$$(r, cr, P_1, P_2)$$

$$\rho = \frac{\lg \frac{1}{P_1}}{\lg \frac{1}{P_2}}$$



$\rightarrow (r, cr, P_1^k, P_2^k) \rightarrow \rho \text{ LSH } \forall k \geq 1.$

for $l_1 / \{0, 1\}^d$: opt. $f = \frac{1}{c}$ UB: [IM '98]

$$\begin{bmatrix} \text{sp.} & n^{1+\epsilon} \\ \text{qt.} & n^\epsilon \end{bmatrix}$$

2B: [MNP'06]
[OWZ'10].

for l_2 : $f = \frac{1}{c^2}$. UB: [AI '06].

