

AA Lecture 8, 2/4/21.

Last time: Dimension Reduction

JL: distribution over $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^k$,
s.t. $\forall x, y \in \mathbb{R}^d$:

$$\Pr_{\varphi} [\|\varphi(x) - \varphi(y)\| = (1 \pm \epsilon) \cdot \|x - y\|] \geq 1 - e^{-\epsilon^2 k/5}$$

$$\Pr [-] \geq 1 - \frac{1}{m^3} \text{ for } k = O\left(\frac{\log n}{\epsilon^2}\right)$$

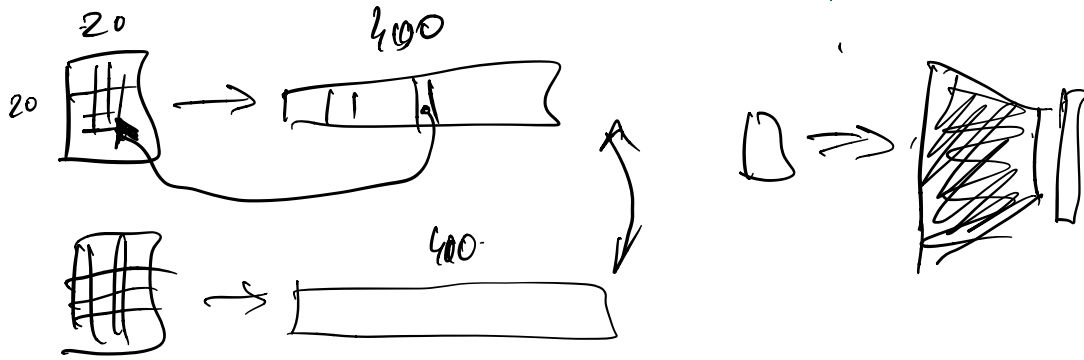
\Rightarrow embed n pts in \mathbb{R}^d into
 \mathbb{R}^k preserving dist $1 \pm \epsilon$,
 $k = O\left(\frac{\log n}{\epsilon^2}\right)$.

(metric embeddings).

Nearest Neighbor Search

Def: (NNS): given $D \subset \mathbb{R}^d$, $|D| = n$,
preprocess D s.t. can answer q's
given $q \in \mathbb{R}^d$, output closest $p^* \in D$,
 $\arg \min_{p \in D} \|q - p\|$.

Notations: look d high



Performance measures of NDS:

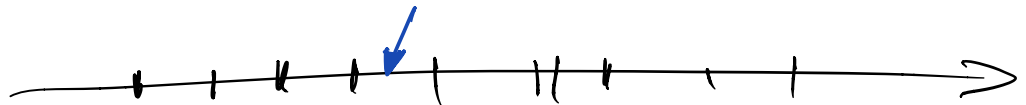
- space
- query time.

Solution #0: no preprocessing "lazy".

- space : $O(n \cdot d)$ ← good.
- q. t. : $O(n \cdot d)$. ~~bad~~

Solution #0.1: $d = 1$

Algo: sort all D. @query q, binary search.



space: $O(n)$.

q.t.s $O(\lg n)$.

← good.

Solution #0.2: $d=2$

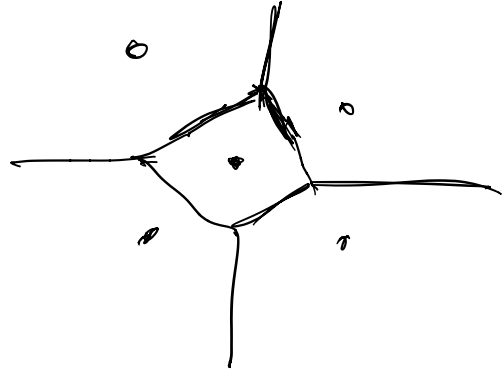
Voronoi diagram:

size = $O(n)$.

+ point location d.s.

- space: $O(n \lg n)$.

- q.t. : $O(\lg n)$.



$d \gg 2$?

Voronoi d. size = $n^{\lfloor \frac{d}{2} \rfloor}$

↗ curse of dimensionality.

Thm: if can solve NNS in dimension $d = (\lg^2 n)$ with $n^{O(1)}$ preproc. time & space,

and query time $n^{0.99}$, then

SETH is false.

↳ Strong Exponential Time Hypothesis.

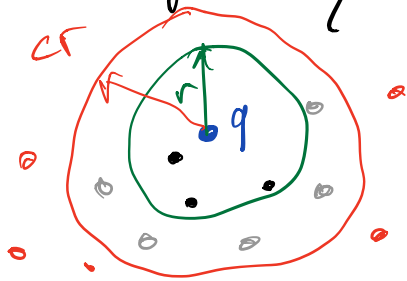
"we can't solve CNF on n vars in time $\ll 2^n$ ".

Approximate NNS in high-d.

Def: (r -near neighbor); given $r > 0$, $D \subset \mathbb{R}^d$,
 build data str s.t. given query $q \in \mathbb{R}^d$,
 report any $p \in D$ s.t. $\|q - p\| \leq r$.
"threshold"

Def: (c -approx r -near neighbor, c -ANN):
 given $c > 1$, $r > 0$, dataset $D \subset \mathbb{R}^d$, $n = |D|$,
 build str. on D s.t. can answer:

given $q \in \mathbb{R}^d$: - if $\exists p^* \in D$ s.t. $\|p^* - q\| \leq r$,
 report $p \in D$ s.t.



$\|p - q\| \leq c \cdot r$. } $\frac{c \cdot r}{r} \geq 90\%$
 - if $\nexists p^* \in D$ s.t. $\|p^* - q\| \leq r$,

may or not report k

Remark: 1) formal reduction from
approx nearest neighbor to approx near n.

idea: guess r .

2) usually, algo for ANN can be modified
to report LCD set:

- if $p \in D$ is $\|p - q\| \leq r \Rightarrow p \in L$
with $r \geq 90\%$ per pt.
- if p is $\|p - q\| \geq cr \Rightarrow p \notin L$.

3) randomized algos,

Back to c-ANN:

Improving over:
space: $O(nd)$,
q.t.: $O(nd)$

Solution #1:

Use dimensionality reduction (JL).

Algo - pick φ at random $k = O\left(\frac{\lg n}{\epsilon^2}\right)$.

- store $\varphi(p)$ for all $p \in D$.

- @query q : compute $p(q)$,

compute $\|\varphi(p) - \varphi(q)\|$ for $p \in D$.

Space: $O(nk) = O\left(n \frac{\lg n}{\epsilon^2}\right)$

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

approx: $c = 1 + \epsilon$. (assuming $n \gg d$).

Correctness analysis:

Using Corollary of JL on points

$$D \cup \{q\} \quad (N = n+1)$$

\Rightarrow distances $p \in D, q$ are all preserved up to $(1+\epsilon)$ factor, with prob. $\geq 1 - \frac{1}{N}$.

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

Crucially used φ is oblivious.

\Rightarrow report an approx nearest neighbor with prob. $\geq 1 - \frac{1}{N} \geq 1 - \frac{1}{n}$.

what we used here:

exist φ : original space / metric \mathbb{R}^d
 \downarrow
 "lower complexity" \mathbb{R}^k
 low-dim

s.t. φ is defined on entire metric.

$\forall p, q$, have that
 distance b/w $\varphi(p)$ & $\varphi(q)$
 \approx distance b/w p & q .
 with prob. $\geq 1 - 1/n^3$.

$$\text{JL: } \varphi: (\mathbb{R}^d, \ell_2) \rightarrow (\mathbb{R}^k, \ell_2)$$

Enough to use sketching map
 (instead of dim-reducing map).

Thm (sketching for ℓ_1): \exists distrib over φ

s.t. $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^k$, & estimation algo

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

$$\forall x, y \in \mathbb{R}^d, \text{ -if } \|x - y\|_1 \leq r \Rightarrow E(\varphi(x), \varphi(y)) = 0$$

$$\text{- if } \|x-y\|_1 \geq (1+\epsilon)r \Rightarrow \\ E(\varphi(x), \varphi(y)) = 1.$$

with probability $\geq 1-\delta$,

$$k = O\left(\frac{\lg \frac{1}{\delta}}{\epsilon^2}\right).$$

Thm: $\varphi: \ell_1 \rightarrow (\ell_2)^k$ preserving distances.

Consider $\{0,1\}^d$ with ℓ_1 distance
 \Rightarrow Hamming space.

Thm: Fix $d, r \geq 1$, approx $1 \pm \epsilon$.

\exists distrib. $\varphi: \{0,1\}^d \rightarrow \{0,1\}^k$,

$$k = O\left(\frac{\lg d}{\epsilon^2}\right) \text{ s.t.}$$

- if $\|x-y\|_1 \leq r \Rightarrow \|\varphi(x) - \varphi(y)\|_1 \leq \epsilon k$

- if $\|x-y\|_1 \geq (1+\epsilon)r \Rightarrow \|\varphi(x) - \varphi(y)\|_1 \geq \left(\frac{\epsilon}{2}\right)k$.

ϵ is a universal constant.

$\Rightarrow (L \in \mathcal{E})$ - ANN for Hamming space
with same params as Sol #1
for ℓ_2 .