

Lecture #6

Last Time: The Johnson-Lindenstrauss lemma,
and sparse / fast variants

unbiased estimator of l_2 -distance via Gaussians

Today: Nearest neighbor search and LSH

Setup: Given a set P of n points in d -dimensions

(1) construct data structure

(2) ^{be able to} answer query point q , with closest point in P

e.g. spam classification

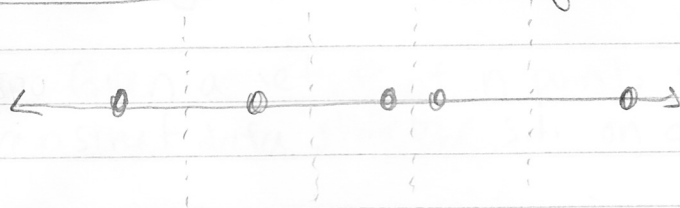
Some initial ideas:

Approach #1: No preprocessing, on query, search through P ^{linearly}

space: $O(dn)$ query time: $O(dn)$

Can we improve the query time? And at what cost?

Approach #2 ($d=1$): Binary search tree



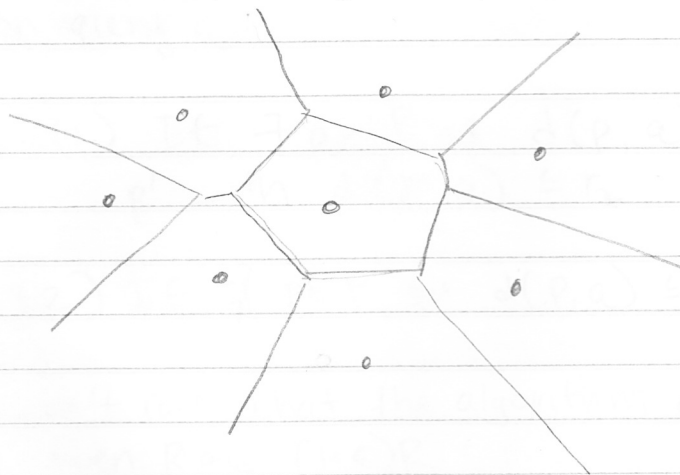
more generally,
think of P as
partitioning space

construct binary search tree on P , on query search for q

space: $O(n)$ query time: $O(\log n)$

Approach #2 ($d=2$) Voronoi Diagram

partition of plane into regions based on which queries map to given point



As dimension increases, becomes much more complex

Other kd-trees [Bentley, 1975]

All approaches in high-dimension have exponential space or query time or both!

"curse of dimensionality"

computational geometry
coined by Bellman

we will study relaxation

c -APX

Setup: Given a set P of n points in d -dimensions
construct data structure s.t. on query q :

Return a point $p \in P$ with

$$d(p, q) \leq c \min_{p' \in P} d(p', q) \leftarrow c\text{-approx nearest neighbor}$$

Even easier:

(r_1, r_2) -
PLEB
point
location
in equal
balls

Setup: Given a set P of n points in d -dimensions and radii r_1, r_2 , construct data structure s.t. on query q

(1) If $\exists p \in P$ st. $d(p, q) \leq r_1$, return any p' with $d(p', q) \leq r_2$ YES and

(2) If $\nexists p \in P$ st. $d(p, q) \leq r_2$ - return No

We don't care what the algorithm does if closest point is between R and $(1+\epsilon)R$

Let D_{max}/D_{min} be max/min interpoint distances in P

Lemma: If for every r , there is a datastructure with space S and query time T that solves $(r, (1+\epsilon)r)$ -PLEB then there is an algorithm for $(1+\epsilon)^2$ -ANN

with space $O(S \log_{1+\epsilon} \frac{D_{max}}{D_{min}})$ and query time $O(T \log_{1+\epsilon} \frac{D_{max}}{D_{min}})$

Pf: Construct a data structure for $(r, (1+\epsilon)r)$ -PLEB for radii

$\frac{D_{min}}{2}, (1+\epsilon) \frac{D_{min}}{2}, (1+\epsilon)^2 \frac{D_{min}}{2}, \dots, \approx D_{max}$
 $(r, (1+\epsilon)r)$

Use binary search to find minimum r st. $(r, (1+\epsilon)r)$ -PLEB returns YES, let p be returned point and let r^* be radius. Then

(1) $d(p, q) \leq (1+\epsilon)r^*$ since $(r^*, (1+\epsilon)r^*)$ -PLEB said YES

(2) $\forall p', d(p', q) \geq \frac{r^*}{(1+\epsilon)}$ since $(\frac{r^*}{(1+\epsilon)}, r^*)$ -PLEB said NO

Thus p is a $(1+\epsilon)^2$ -ANN \square

ring-cover
trees

More efficient reductions in [Indyk, Motwani 1998]
and [Har-Peled, 2001] $(r, (1+\epsilon)r)$ -PLEB \Rightarrow $(1+\epsilon)$ -ANN

So far collisions \equiv bad, but today we will exploit them

Locality Sensitive Hashing: [Indyk, Motwani 1998]
similar items are more likely to map to same bucket

setup: Hash family $\mathcal{H} = \{h: U \rightarrow S\}$ is called
 (r_1, r_2, p_1, p_2) -locality sensitive if for any $p, p' \in U$

(1) If $d(p, p') \leq r_1$, then $\Pr_{h \in \mathcal{H}} [h(p) = h(p')] \geq p_1$

(2) If $d(p, p') \geq r_2$ then $\Pr_{h \in \mathcal{H}} [h(p) = h(p')] \leq p_2$

we will always work with $p_1 > p_2$ and $r_1 < r_2$

Theorem [Indyk, Motwani]: Suppose there is a
 (r_1, r_2, p_1, p_2) -locality sensitive hash family \mathcal{H} .
Then there is an algorithm for (r_1, r_2) -PLEB
which uses

space: $O(dn + n^{1+p})$

query time: $O(n^p)$ evaluations of hash fcn

where $p = \frac{\ln 1/p_1}{\ln 1/p_2}$, and succeeds with constant probability
(for fixed p, α)

The space is polynomial, but query time is sublinear

Let k, ℓ be parameters chosen later. Let

$$G = \left\{ g: U \rightarrow S^k \right\} \quad \text{where}$$

new bucket

amplification

$$g = (h_1(p), h_2(p), \dots, h_\ell(p)) \quad \text{and each } h_i \in \mathcal{H}$$

Preprocessing: (1) Choose g_1, g_2, \dots, g_ℓ independently, uniformly at random

(2) For each $p \in P$, store it in buckets $g_1(p), g_2(p), \dots, g_\ell(p)$

(3) Discard empty buckets

On query: Search $g_1(q), g_2(q), \dots, g_\ell(q)$ and find first 2ℓ points.
 return any point p found with $d(p, q) \leq r_2$ and YES. Else return NO.

Analysis: We want to show, with constant probability, over choice of g_i 's
 the following events hold: for fixed P, q

(1) If $\exists p \in P$ with $d(p, q) \leq r_1$, then $g_j(p) = g_j(q)$ for some j

(2) there are at most 2ℓ points $p \in P$ with $d(p, q) \geq r_2$ and $g_j(p) = g_j(q)$ for some j

Fix j :

Let $k = \log_{1/p_2} n$, then the expected number of points satisfying the conditions in (2) is at most

$$n (p_2)^k = n (p_2)^{\log_{1/p_2} n} \leq 1$$

Thus by Markov, the probability (2) does not hold is at most $\frac{1}{2}$

Fix j : The probability of $g_j(p) = g_j(q)$ in (1) is bounded from below by

$$P_1^k = P_1^{\log_{1/p_2} n} = \left(\underbrace{\left(\frac{1}{p_2} \right)^{\log_{1/p_2} P_1}}_{P_1^{\log_{1/p_2} P_1}} \right)^{\log_{1/p_2} n} = n^{\frac{-\log_{1/p_2} P_1}{\log_{1/p_2} P_2}} = n^{-\beta}$$

$\ln(1+x) = x - \frac{x^2}{2!} \dots$; thus

$$P = \frac{\ln(1 - \frac{r}{d})}{\ln(1 - \frac{cr}{d})} \approx \frac{\ln(1 + \frac{r}{d})}{\ln(1 + \frac{cr}{d})} \approx \frac{\frac{r}{d}}{\frac{cr}{d}} = \frac{1}{c}$$

Thus the probability that (1) holds is at least

$$1 - (1 - n^{-P})^d \geq 1 - 1/e > 1/2$$

and setting $d = n^P$ gives.

Thus with constant probability both (1) and (2) hold, which implies the algorithm succeeds in solving (r_1, r_2) -PLEB \square

Some applications

Let $\mathbb{H}^d = \{0, 1\}^d$ be the d -dimensional Hypercube, then

$$\mathcal{H} = \left\{ h_i \mid h_i(b_1, b_2, \dots, b_d) = b_i \text{ for } i=1, 2, \dots, d \right\} \text{ is}$$

$(r_1, cr_1, 1 - \frac{r_1}{d}, 1 - \frac{cr_1}{d})$ -locality sensitive, for the Hamming distance.

Corollary: There is an algorithm for (r, cr) -PLEB in \mathbb{H}^d that uses space $O(dn + n^{1+\frac{1}{c}})$ and for each query needs $O(n^{\frac{1}{c}})$ evaluations of the hash function, each of which takes $O(d)$ time.

i.e. for the Hamming distance, $P \leq 1/c$

Theorem [Andoni, Indyk, 2006] for Euclidean distance, $P \leq 1/c^2$

These bounds are tight [O'Donnell, Wu, Zhou 2011]