

Johnson-Lindstrauss Lemma

THM: Let X be a set of N points in \mathbb{R}^n and $0 < \epsilon < 1$. Consider an $m \times n$ matrix A whose rows are independent, mean-zero, isotropic, and sub-gaussian rand. vectors in \mathbb{R}^n . Rescale A by defining (random projection) $P := \frac{1}{\sqrt{m}} A$.

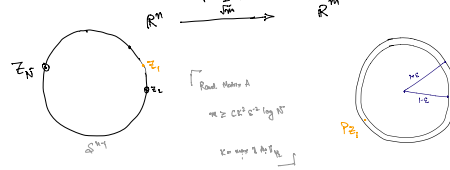
Assume that $m \geq C\epsilon^{-2} \log N$, where C is a large enough constant (which depends on the sub-g norm of A 's rows).

Then, w.h.p., the map P preserves all pairwise distances in X , up to ϵ relative error.

$$(1-\epsilon) \|x-y\| \leq \|P_x - P_y\|_2 \leq (1+\epsilon) \|x-y\| \quad \text{for all } x, y \in X.$$

Johnson-Lindstrauss (JL) Lemma

$$X \subseteq \mathbb{R}^n \quad |X| = N$$



with high prob.,
for all $x \in X$
 $\|P_x\| = (1 \pm \epsilon) \|x\|$

$$\mathbb{E} \left[\max_{x \in X} \|P_x\|^2 - \|x\|^2 \right] \leq \epsilon$$

$$m \geq \frac{1}{\epsilon^2} \sum_{x \in X} \|x\|^2 \log N$$

$$\mathbb{E} \left[\max_{x \in X} \|P_x - P_y\| \right] \leq \epsilon \sqrt{\log N}$$

Applications

- ① Pairwise Euclidean Distances
- ② Streaming Alg. for F_2 -estimation
- ③ Approximate Nearest Neighbor

} This Lecture

Streaming Algorithm for F_2 -estimation

• Given a stream of N indices i_1, i_2, \dots, i_N $i_t \in \{1, 2, \dots, m\}$

• Maintain fq. vector $f_t := \{z_t : i_t = j\}$

• Goal: Approximate the L_2 norm $\|f_t\|_2$

[Space Constrained Environment]

Pioneering Work of
Alan, Matias, & Szegedy
* Streaming Model of
Computation

$$E_m \quad m=3 \quad N=6$$

$$[1, 1, 3, 3, 1, 2]$$

$$f = \begin{pmatrix} 3 \\ 1 \\ 2 \end{pmatrix}$$

[Trivial in $O(m \log m)$ space
Here, will use
 $O(\log \log m)$ space]

Idea: Instead of storing f , maintain a projection of f

$$y = P \cdot f$$

At time t , alg receives $i_t = 3 \rightarrow f \leftarrow f + e_3$

Update $y \leftarrow y + P \cdot e_3$

(add the 3th column of P to y)

Analysis

Via JL we have that, w.h.p.,

$$(1-\epsilon)\|f\| \leq \|Pf\| \leq (1+\epsilon)\|f\|$$

Setting $m \geq \frac{1}{\epsilon^2}$, we have this bound.

(with $m \geq \frac{1}{\epsilon^2} \log N$ this guarantee is achieved for all N steps).

Approximate Nearest Neighbor

Nearest Neighbor (Exact)

Input: Points $P = \{p_1, p_2, \dots, p_N\} \subset \mathbb{R}^n$
Objective: Preprocess P (find a data structure) such that given any query point $q \in \mathbb{R}^n$,

we can quickly find arg min $\|p - q\|$
 $p \in P$

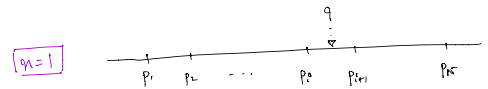
Considerations

- Space (required to store data structure)
- Query Time

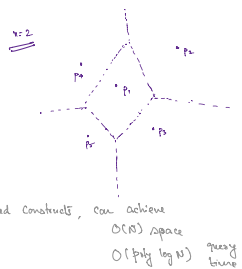
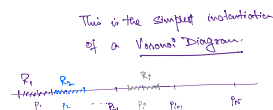
Exhaustive Search

- $O(nN)$ time

For low dimensions ($n=1,2$) efficient solutions are known.



- Store given points $P = \{p_1, p_2, \dots, p_n\}$ in increasing order. Space Required $O(N)$
- Binary Search to find $p_i \leq q < p_{i+1}$. Query Time $O(\log N)$. Output the closer of p_i, p_{i+1} .



In high dimensions no significant improvements over the naive method. (Relevant Data Structure: Kd-Trees)

Curse of Dimensionality: The failure of low-dimensional methods when applied in high dimensions

Solution by allowing randomization & approximation

ϵ -Approximate Nearest Neighbor (ϵ -ANN)

Objective: Given a query point $q \in \mathbb{R}^n$, find $p \in P$ such that

$$\|p - q\| \leq (1 + \epsilon) \min_{p' \in P} \|p' - q\|$$

We solve ϵ -ANN via a reduction to a "gap" problem.

ϵ -PLEB (ϵ -Point Location in Equal Balls)

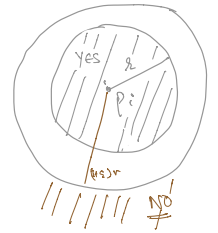
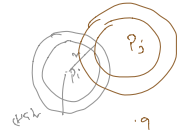
ϵ -PLEB (ϵ -Point Location in Equal Balls)

Input: N Euclidean balls of radius r , centered at points $P = \{p_1, p_2, \dots, p_N\} \subset \mathbb{R}^n$

Output: Given a query $q \in \mathbb{R}^n$ we must answer as follows

- If there exists $p_i \in P$ such that $q \in B(p_i, r)$, then we must answer YES and output any point p_i with $q \in B(p_i, (1 + \epsilon)r)$
- If there does not exist p_i with $q \in B(p_i, (1 + \epsilon)r)$ we must say NO
- Otherwise (i.e., $2r < \|p_i - q\| \leq (1 + \epsilon)r$) we can say either YES or NO (if we say YES, we must output a point p_i with $q \in B(p_i, (1 + \epsilon)r)$)

$B(p, r)$: Euclidean ball of radius r centered at p



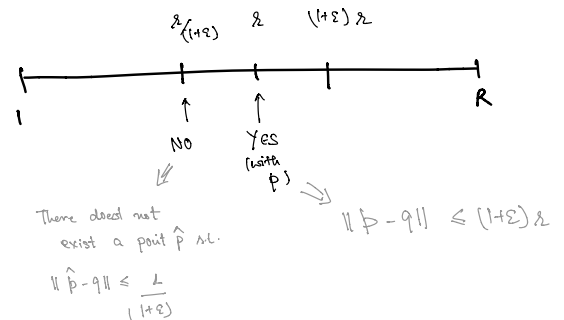
YES always returned with a certificate p_i , which can be verified \downarrow

Reducing ϵ -ANN to ϵ -PLEB

- For every $r = (1 + \epsilon)^0, (1 + \epsilon)^1, (1 + \epsilon)^2, \dots, R$
Construct data structure for PLEB(r)

- Given a query point $q \in \mathbb{R}^n$, via binary search, find min r s.t. PLEB(r) returns a YES for q .
(Let p be the point returned by PLEB(r))

$R :=$ max inter-point distance of P
 $1 \leq \|p - p'\| \leq R$
Scaling.



Therefore, p satisfies

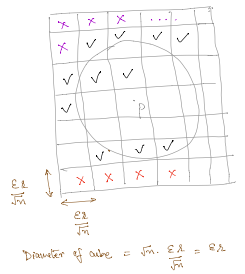
$$\|p - q\| \leq (1 + \epsilon)^2 \min_{p' \in P} \|p' - q\|$$

✓ p is a solution for ϵ -ANN

Solving ϵ -PLEB

We discretize the space & use a hash table.

- Partition the space into n -dimensional cubes of side length $\frac{\epsilon R}{\sqrt{n}}$ [Each cube is identified by the min point i it]



We create a hash table (initially empty)

For each cuboid C & $p_i \in P$ if $B(p_i, R) \cap C \neq \emptyset$ then insert (C, p_i) into the hash table

Preprocessing

Query Time:

1. Given query $q \in \mathbb{R}^n$, find C that contains it

2. Look up i in C in the hash table.

No match implies that $\nexists B(p_i, r)$ which intersects with C

Therefore, $q \notin B(p_i, r)$ for any $p_i \in P$.

Hence, we can safely report NO (following PLEB requirement)

3. Say C is in the hash table with $B(p_j, r)$, (i.e. $B(p_j, r) \cap C \neq \emptyset$)

By triangle ineq.

$$\|p_j - q\| \leq r + \text{diam}(C) = r + \epsilon r = (1 + \epsilon)r$$

Therefore, we can return p_j with

$$\|p_j - q\| \leq (1 + \epsilon)r$$

Time and Space Analysis

Query Time $O(n)$ - Find the cube
- Hash.

Space Requirement Exponential in n \times

(& Preprocessing time)

Curse of dimensionality

Volume of a ball of radius $R = \frac{2^n \cdot R^n}{n^{n/2}}$

Volume of a cuboid with side length $\frac{\epsilon R}{\sqrt{n}} = \left(\frac{\epsilon R}{\sqrt{n}}\right)^n$

Number of cubes that intersect with the ball $\approx \left(\frac{1}{\epsilon}\right)^n$

JL allows us to circumvent the curse of dimensionality

Apply JL to map all given points P & query q
into $m = O(\epsilon^{-2} \log N)$ space

• Now, the space (and preprocessing time) requirement is
$$\left(\frac{1}{\epsilon}\right)^m \sim N^{-\epsilon^2 \log(1/\epsilon)}$$

┌ PLEB \longrightarrow ANN
Time complexity goes up
by a factor of
 $O\left(\frac{\log R}{\epsilon}\right)$. ┘

• Query time $O(mn) + O(m)$
JL Find cube & hash.
 $= O(\epsilon^{-2} n \log N)$