

RECALL

- Sub-gaussian random vectors X : one-dim marginals are sub-gaussians.

$$\|X\|_{\psi_2} := \sup_{z \in S^{n-1}} \|\langle X, z \rangle\|_{\psi_2}$$

- Special case: vectors X with independent sub-gaussian components

$$K := \max_i \|X_i\|_{\psi_2}$$

(i) $\|X\|_{\psi_2} \leq CK$

(ii) Norm Concentration (Annulus Thm.)

$$\| \|X\| - \sqrt{n} \| \|_{\psi_2} \leq CK^2$$

Johnson-Lindenstrauss (JL) Lemma

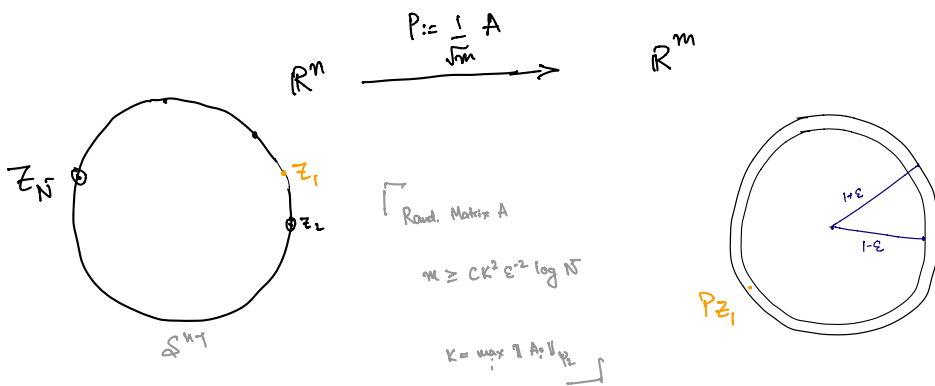
$$X \subseteq \mathbb{R}^n \quad |X| = N$$

with high prob, for all $z \in X$

$$\|Pz\| = (1 \pm \epsilon) \|z\|$$

$$E \left[\max_z \left| \|Pz\| - \|z\| \right| \right] \leq \epsilon$$

$$\text{for } m \gtrsim K^2 \epsilon^{-2} \log N$$



Sparse JL

c.f. [Achlioptas '93]

JL holds as long as the rows of rand. matrix A are independent, mean zero, isotropic, and sub-gaussian.

Hence, the following sparse A suffices

$$(ind) \quad A_{ij} \sim \begin{cases} \sqrt{3} & \text{w.p. } 1/6 \\ -\sqrt{3} & \text{w.p. } 1/6 \\ 0 & \text{w.p. } 2/3 \end{cases}$$

[in expectation, $2/3$ of entries of A are zero.]

$$P = \frac{1}{\sqrt{m}} A$$

FAST JL

[Ailon-Chazelle '06]

Standard JL transform

requires $O\left(\frac{n \log N}{\epsilon^2}\right)$

operations per vector $i \times$

$$\Gamma \begin{cases} A \in \mathbb{R}^{m \times n} \\ \epsilon \in \mathbb{X} \subseteq \mathbb{R}^n \end{cases}$$

$$u = O\left(\frac{\log N}{\epsilon^2}\right)$$

+ Fix unit vector $u \in \mathbb{S}^{n-1}$ (wish to preserve norm of u)

+ Consider sparse matrix $S \in \mathbb{R}^{m \times n}$ with ind. rows S_1, \dots, S_m with $S_i \in \mathbb{R}^n \setminus \{0\} \subseteq \pm e_1, \pm e_2, \dots, \pm e_n$

Ex: Standard basis vector \downarrow

+ Each row is one sparse. Hence, computing Su requires $O(n)$ time (and space)

$$\mathbb{E} S_i = 0 \quad \mathbb{E} S_i S_i^T = I_n$$

(Mean Zero) (Isotropic)

R.V. $S_i^T u$ is bounded

$$\|S_i^T u\|_{\infty} \leq C \|S_i^T u\|_0 = C \cdot \mathbb{1}_{\{n\}} \cdot \|u\|_0$$

+ The vector $Su = \begin{pmatrix} S_1^T u \\ \vdots \\ S_m^T u \end{pmatrix}$ has ind. sub-g components with $\mathbb{E} (S_i^T u)^2 = u^T (\mathbb{E} S_i S_i^T) u = u^T I_n u = 1$

+ Norm Concentration

$$\| \|Su\| - \sqrt{m} \|u\|_2 \leq C m \|u\|_0^2 \quad \text{--- (1)}$$

Lemma Let $u \in \mathbb{S}^{n-1}$ be a fixed unit vector, with $\|u\|_0 \leq \sqrt{\frac{\log n}{n}}$. Also, let $S \in \mathbb{R}^{m \times n}$ be a rand. row sampling matrix with $m \geq \frac{\log^2(n/s)}{\epsilon^2}$.

$$\text{Then, } \Pr \left\{ \left\| \frac{1}{\sqrt{m}} Su \right\| \notin (1 \pm \epsilon) \cdot \|u\| \right\} \leq \delta/2.$$

A JL transform with running time

$$O(n \log n + m), \text{ per vector}$$

[Use Care $N \gg n$]

Key Idea: Randomly "rotate" the basis such that $x \mapsto u$ $\|u\|_0 \leq \sqrt{\frac{\log n}{n}}$ holds with high prob.

Use a "structured" matrix M such that (i) $\|Mx\|_0 \leq \sqrt{\frac{\log n}{n}}$ (ii) M admits fast matrix-vector mult.

Hadamard Matrices

Recursive Defn. $H_d := \frac{1}{\sqrt{2}} \begin{pmatrix} H_{d/2} & H_{d/2} \\ H_{d/2} & -H_{d/2} \end{pmatrix}$

$$H_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$$

(a) Rows of H are orthonormal

(b) Recursive defn. ensures H_{2^k} can be computed in $O(n \log n)$ time.

(c) Cannot expect $\|Hx\|_0 \leq \sqrt{\frac{\log n}{n}}$ for all x (e.g., x is a row of H).

write $H := H_n$

$$H = \frac{1}{\sqrt{m}} \begin{pmatrix} \pm 1 & \dots \\ \vdots & \end{pmatrix}$$

Idea: Randomize over column signs

Set diagonal matrix D such that $D_{ii} \in \mathbb{R} \setminus \{0\}$

For fixed unit vector x , consider HDx

Lemma: let $x \in \mathbb{S}^{n-1}$ and $u := HDx$, where HD is a random Hadamard matrix.

Then, $\Pr_D \left\{ \|u\|_0 \geq \sqrt{\frac{\log(m/s)}{n}} \right\} \leq \delta/2.$

$$\Pr_D \left\{ H_{ij} \in \mathbb{R} \setminus \left\{ \frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}} \right\} \right\} \text{ i.i.d. Bernoulli (obviously ind.)}$$

$$\|H_{ij}\|_{\infty} \leq \frac{C}{\sqrt{m}}$$

$$\|H_i^T x\|_{\infty}^2 \leq C \cdot \|x\|_2^2 \cdot \frac{C}{n} = \frac{C}{n} \quad (\text{Hoeffding's Ineq.})$$

$$\|H_i^T x\|_{\infty} \leq \frac{C}{\sqrt{n}}$$

$$\mathbb{E} \max_i |H_i^T x| \leq \frac{C}{\sqrt{n}} \cdot \sqrt{\log n}$$

Using tail prob. we get the desired dependence on δ .

Thm (Fast JL): There exists a random matrix $M \in \mathbb{R}^{m \times n}$

with $m \geq \frac{1}{\epsilon^2} \log\left(\frac{n}{\delta}\right)$ such that for any $x \in \mathbb{R}^n$

$$\|Mx\| = (1 \pm \epsilon) \cdot \|x\| \text{ holds w.p. at least } (1 - \delta).$$

In addition, matrix-vector multiplication with M takes $O(n \log n + m)$ time.

Pf: Fix $x \in S^{n-1}$. Two "bad" events
 $M = S \cdot HD$
 $B_1 := \{ \|HDx\|_\infty > \sqrt{\frac{\log(n/\delta)}{n}} \}$
 $B_2 := \{ \|Mx\| \notin (1 \pm \epsilon) \|x\| \}$

To prove the thm, suffices to bound $\Pr\{B_2\} \leq \delta$.

$$\Pr\{B_1\} \leq \frac{\delta}{2} \quad (\text{lemma 2})$$

$$\Pr\{B_2 | \neg B_1\} \leq \frac{\delta}{2} \quad (\text{lemma 1})$$

$$\begin{aligned} \Pr\{B_2\} &= \Pr\{B_2 | B_1\} \cdot \Pr\{B_1\} + \Pr\{B_2 | \neg B_1\} \cdot \Pr\{\neg B_1\} \\ &\leq 1 \cdot \frac{\delta}{2} + \frac{\delta}{2} \cdot 1 = \delta \end{aligned}$$

□

Applications of JL

① Pairwise Distances (Euclidean)

$$X \subseteq \mathbb{R}^n$$

- Naive, Exact Computation $O(N^2 n)$
- Approx Computation JL

JL Matrix $P: \mathbb{R}^n \rightarrow \mathbb{R}^m$

Computation time
 $O(\text{JL time} + N^2 m)$

Using Fast JL we get $m = O(\epsilon^{-2} \log^2 n)$
 with JL time $(Nm \log n)$

Overall runnig time
 $O(Nn \log n + N^2 \epsilon^{-2} \log^2 n)$

② Streaming Algo. for ℓ_2 estimation

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

• Maintain fq. vector $f_j := |\{t : i_t = j\}|$

• Goal: Approximate the ℓ_2 norm $\|f\|_2$

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

$$y = P \cdot f \quad \left[P \leftarrow \text{Fast JL transform} \right]$$

At time t , alg receives $i_t = j$

$$\text{Update } y \leftarrow y + P \cdot e_j$$

(add the j th column of P to y)

Space Constrained Environment

Pioneering Work of:
Alon, Matias, & Szegedy
* Streaming Model of Computation

$$\underline{\underline{E}}_m \quad n=3 \quad N=6$$

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

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

Trivial in $O(n \log n)$ space
Here, will use $O(\text{poly} \log n)$ space

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).