

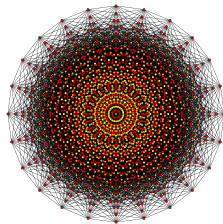
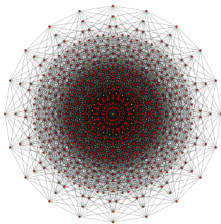
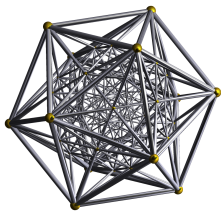
Computational Principles for High-dim Data Analysis

(Lecture Six)

Yi Ma

EECS Department, UC Berkeley

September 14, 2021



Convex Methods for Sparse Signal Recovery

(Matrices with Restricted Isometry Property)

- 1 The Johnson-Lindenstrauss Lemma
- 2 RIP of Gaussian Matrices
- 3 RIP of Non-Gaussian Matrices

“Algebra is but written geometry; geometry is but drawn algebra.”
– Sophie Germain

Restricted Isometry Property (Recap)

Definition (Restricted Isometry Property)

The matrix \mathbf{A} satisfies the *restricted isometry property (RIP)* of order k , with constant $\delta \in [0, 1)$, if

$$\forall \mathbf{x} \text{ } k\text{-sparse}, \quad (1 - \delta) \|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta) \|\mathbf{x}\|_2^2. \quad (1)$$

The *order- k restricted isometry constant* $\delta_k(\mathbf{A})$ is the smallest number δ such that the above inequality holds.

Example of Gaussian Matrices: If \mathbf{A}_1 is a large $m \times k$ ($k < m$) matrix with entries independent $\mathcal{N}(0, 1/m)$,

$$\sigma_{\min}(\mathbf{A}_1^* \mathbf{A}_1) \approx (\sqrt{1} - \sqrt{k/m})^2 \geq 1 - 2\sqrt{k/m},$$

$$\sigma_{\max}(\mathbf{A}_1^* \mathbf{A}_1) \approx (\sqrt{1} + \sqrt{k/m})^2 \leq 1 + 3\sqrt{k/m}.$$

Length Concentration of Gaussian Vectors

Lemma

Let $\mathbf{g} = [g_1, \dots, g_m]^* \in \mathbb{R}^m$ be an m -dimensional random vector whose entries are iid $\mathcal{N}(0, 1/m)$. Then for any $t \in [0, 1]$,

$$\mathbb{P} \left[\left| \|\mathbf{g}\|_2^2 - 1 \right| > t \right] \leq 2 \exp \left(-\frac{t^2 m}{8} \right). \quad (2)$$

This result can be obtained via the Cramer-Chernoff exponential moment method or Bernstein inequality (see Appendix E).¹

¹High-Dimensional Probability, Roman Vershynin, Cambridge University Press, 2018.

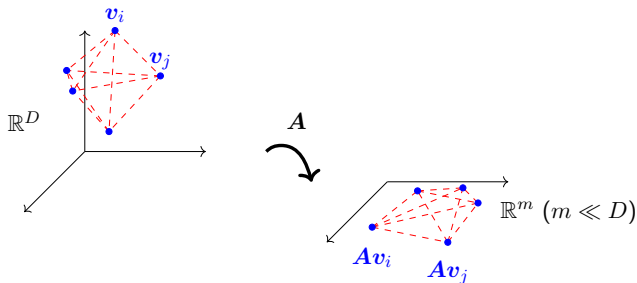
The JL Lemma: Distance Preserving Random Projections

Theorem (Johnson-Lindenstrauss Lemma)

Let $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^D$ and $\mathbf{A} \in \mathbb{R}^{m \times D}$ be a random matrix whose entries are i.i.d. $\mathcal{N}(0, 1/m)$. Then for any $\epsilon \in (0, 1)$, with probability at least $1 - 1/n^2$, the following holds:

$$\forall i \neq j, \quad (1 - \epsilon) \|\mathbf{v}_i - \mathbf{v}_j\|_2^2 \leq \|\mathbf{A}\mathbf{v}_i - \mathbf{A}\mathbf{v}_j\|_2^2 \leq (1 + \epsilon) \|\mathbf{v}_i - \mathbf{v}_j\|_2^2, \quad (3)$$

provided $m > 32 \frac{\log n}{\epsilon^2}$.



The JL Lemma

Proof.

Finite cases: let $\mathbf{g}_{ij} = A \frac{\mathbf{v}_i - \mathbf{v}_j}{\|\mathbf{v}_i - \mathbf{v}_j\|_2}$ for any $i \neq j \in \{1, \dots, n\}$.

Tail bound: \mathbf{g}_{ij} is distributed as an iid Gaussian vector, with entries $\mathcal{N}(0, 1/m)$. Applying Lemma:

$$\mathbb{P} \left[\left| \|\mathbf{g}_{ij}\|_2^2 - 1 \right| > t \right] \leq 2 \exp(-t^2 m / 8). \quad (4)$$

Union bound: Summing the probability of failure over all $i \neq j$, and then plugging in $t = \epsilon$ and $m \geq 32 \log n / \epsilon^2$, we get

$$\mathbb{P} \left[\exists (i, j) : \left| \|\mathbf{g}_{ij}\|_2^2 - 1 \right| > t \right] \leq \frac{n(n-1)}{2} \times 2 \exp(-t^2 m / 8) \leq n^{-2}. \quad (5)$$

Hence $\left| \|\mathbf{g}_{ij}\|_2^2 - 1 \right| \leq \epsilon$ with probability $1 - n^{-2}$.



The JL Lemma: Generalization to ℓ^p Norms

Locality-Sensitive Hashing²: for $p \in (0, 2]$, there exist the so-called *p-stable distributions* such that a random matrix \mathbf{A} drawn from a p -stable distribution will preserve ℓ^p distance between vectors:

$$(1 - \epsilon) \|\mathbf{v}_i - \mathbf{v}_j\|_p^2 \leq \|\mathbf{A}\mathbf{v}_i - \mathbf{A}\mathbf{v}_j\|_p^2 \leq (1 + \epsilon) \|\mathbf{v}_i - \mathbf{v}_j\|_p^2. \quad (6)$$

Example: For ℓ^1 norm, the corresponding distribution is the Cauchy distribution with density:

$$p(x) = \frac{1}{\pi} \cdot \frac{1}{1 + x^2}.$$

²Locality-sensitive hashing scheme based on p -stable distributions, M. Datar, N. Immorlica, P. Indyk, and V. S. Mirrokni. ACM SCG 2004.

The JL Lemma: Fast Nearest Neighbors

Compact Code for Fast Nearest Neighbor³:

- 1: **Goal:** Generate compact binary code for efficient nearest neighbor search of high-dimensional data points.
- 2: **Input:** $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^D$ and $m = O(\log n)$.
- 3: Generate a Gaussian matrix $\mathbf{R} \in \mathbb{R}^{m \times D}$ with entries i.i.d. $\mathcal{N}(0, 1)$.
- 4: **for** $i = 1, \dots, n$ **do**
- 5: Compute $\mathbf{R}\mathbf{x}_i$,
- 6: Set $\mathbf{y}_i = \sigma(\mathbf{R}\mathbf{x}_i)$ where $\sigma(\cdot)$ is the entry-wise binary thresholding.
- 7: **end for**
- 8: **Output:** $\mathbf{y}_1, \dots, \mathbf{y}_n \in \{0, 1\}^m$.

Instead of $O(\log n)$ real numbers, one only needs $O(\log n)$ binary bits!

³Compact projection: Simple and efficient near neighbor search with practical memory requirements, K. Min, J. Wright, and Y. Ma, CVPR 2010.

RIP of Gaussian Matrices

Theorem (RIP of Gaussian Matrices)

There exists a numerical constant $C > 0$ such that if $\mathbf{A} \in \mathbb{R}^{m \times n}$ is a random matrix with entries independent $\mathcal{N}(0, \frac{1}{m})$ random variables, with high probability, $\delta_k(\mathbf{A}) < \delta$, provided

$$m \geq Ck \log(n/k) / \delta^2. \quad (7)$$

Implications: ℓ^1 minimization can successfully recover k -sparse solutions \mathbf{x}_o from about

$$m \geq Ck \log(n/k) \sim \Omega(k)$$

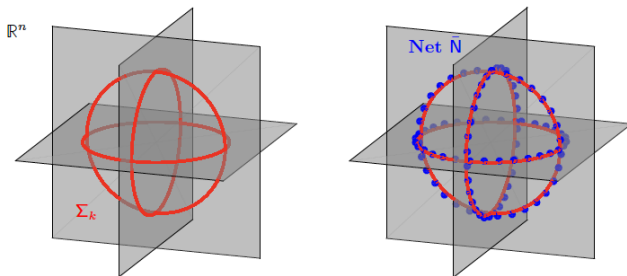
random measurements.

Proof: Step 1. Discretization to Finite Cases

$\delta_k(\mathbf{A}) \leq \delta$ if and only if $\sup_{\mathbf{x} \in \Sigma_k} \left| \|\mathbf{Ax}\|_2^2 - 1 \right| \leq \delta$ where

$$\Sigma_k = \{\mathbf{x} \mid \|\mathbf{x}\|_0 \leq k, \|\mathbf{x}\|_2 = 1\}. \quad (8)$$

Construct a finite (minimal) ϵ -**net** for Σ_k .



Proof: Step 1. Discretization to Finite Cases

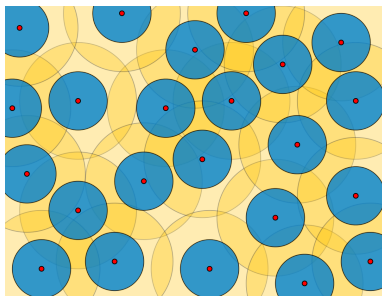
An ϵ -net (or covering) N for a given set S if

$$\forall \mathbf{x} \in S, \quad \exists \bar{\mathbf{x}} \in N \quad \text{such that} \quad \|\mathbf{x} - \bar{\mathbf{x}}\|_2 \leq \epsilon. \quad (9)$$

A set M is **ϵ -separated** if every pair of distinct points \mathbf{x}, \mathbf{x}' in M has distance at least ϵ :

$$\|\mathbf{x} - \mathbf{x}'\|_2 \geq \epsilon. \quad (10)$$

Fact: A maximal ϵ -separated subset $M \subset S$ is a (minimal) ϵ -net of S .



Proof: Step 1. Discretization to Finite Cases

Lemma (ϵ -Nets for the Unit Ball)

There exists an ϵ -net for the unit ball $B(\mathbf{0}, 1) \subset \mathbb{R}^d$ of size at most $(3/\epsilon)^d$.

Proof: Let $N \subset B(\mathbf{0}, 1)$

be a *maximal* ϵ -separated set. The balls $B(\mathbf{x}, \epsilon/2)$ with $\mathbf{x} \in N$ are contained in $B(\mathbf{0}, 1 + \epsilon/2)$. Thus,

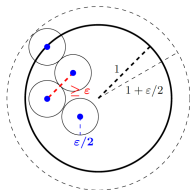
$$|N| \operatorname{vol}(B(\mathbf{0}, \epsilon/2)) \leq \operatorname{vol}(B(\mathbf{0}, 1 + \epsilon/2)). \quad (11)$$

Hence,

$$|N| \leq \frac{\operatorname{vol}(B(\mathbf{0}, 1 + \epsilon/2))}{\operatorname{vol}(B(\mathbf{0}, \epsilon/2))} \quad (12)$$

$$= \left(\frac{1 + \epsilon/2}{\epsilon/2} \right)^d = (1 + 2/\epsilon)^d \quad (13)$$

$$\leq (3/\epsilon)^d \quad (14)$$



Proof: Step 1. Discretization to Finite Cases

Lemma (Discretization)

Suppose we have a set $\bar{N} \subseteq \Sigma_k$ with the following property: for all $\mathbf{x} \in \Sigma_k$, there exists $\bar{\mathbf{x}} \in \bar{N}$ such that

- $|\text{supp}(\bar{\mathbf{x}}) \cup \text{supp}(\mathbf{x})| \leq k$
- $\|\mathbf{x} - \bar{\mathbf{x}}\|_2 \leq \epsilon.$

set

$$\delta_{\bar{N}} = \max_{\bar{\mathbf{x}} \in \bar{N}} \left| \|\mathbf{A}\bar{\mathbf{x}}\|_2^2 - 1 \right|. \quad (15)$$

Then

$$\delta_k(\mathbf{A}) \leq \frac{\delta_{\bar{N}} + 2\epsilon}{1 - 2\epsilon}. \quad (16)$$

Implications: RIP constant δ does not change much if we restrict our calculation to a finite ϵ -covering set \bar{N} .

Proof: Step 1. Discretization to Finite Cases

Lemma (ϵ -Nets for Σ_k)

There exists an ϵ -net \bar{N} for Σ_k satisfying the two properties required in Lemma 6, with

$$|\bar{N}| \leq \exp\left(k \log(3/\epsilon) + k \log(n/k) + k\right). \quad (17)$$

Proof.

Constructing an ϵ -Net for each ball in Σ_k and take the union. Using the Stirling's formula,⁴ we can estimate

$$|\bar{N}| \leq (3/\epsilon)^k \binom{n}{k} \leq (3/\epsilon)^k \left(\frac{ne}{k}\right)^k. \quad (18)$$



⁴Stirling's formula gives the bounds for factorials: $\sqrt{2\pi k} \left(\frac{k}{e}\right)^k \leq k! \leq e\sqrt{k} \left(\frac{k}{e}\right)^k$

Proof: Steps 2 and 3

Step 2: Tail Bound for Probability of Each Failure Case:

For each $\mathbf{x} \in \bar{\mathbf{N}}$, \mathbf{Ax} is a random vector with entries independent $\mathcal{N}(0, 1/m)$. We have

$$\mathbb{P} \left[\left| \|\mathbf{Ax}\|_2^2 - 1 \right| > t \right] \leq 2 \exp(-mt^2/8). \quad (19)$$

Step 3: Union Bound for Probability of All Failure Cases:

Summing over all elements of $\bar{\mathbf{N}}$, we have

$$\mathbb{P} [\delta_{\bar{\mathbf{N}}} > t] \leq 2 |\bar{\mathbf{N}}| \exp(-mt^2/8) \quad (20)$$

$$\leq 2 \exp \left(-\frac{mt^2}{8} + k \log \left(\frac{n}{k} \right) + k \left(\log \left(\frac{3}{\epsilon} \right) + 1 \right) \right). \quad (21)$$

On the complement of the event $\delta_{\bar{\mathbf{N}}} > t$, we have

$$\delta_k(\mathbf{A}) \leq \frac{2\epsilon + t}{1 - 2\epsilon}. \quad (22)$$

Setting $\epsilon = \delta/8$, $t = \delta/4$, and ensuring that $m \geq Ck \log(n/k)/\delta^2$ for sufficiently large numerical constant C , we obtain the result.



RIP of Order k for Gaussian Matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$

From the above derivation, especially from equation (21), we see that a **slight more tight bound** for m is of the form

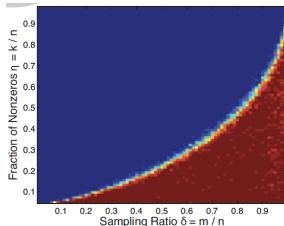
$$m \geq 128k \log(n/k)/\delta^2 + (\log(24/\delta) + 1)k/\delta^2 \doteq C_1 k \log(n/k) + C_2 k.$$

For a small δ , the constants C_1 and C_2 can be rather large.

A much tighter bound (one of the best known) for m is given as⁵:

$$m \geq 8k \log(n/k) + 12k.$$

A precise (phase transition) expression of m as function of k, n exists (section 3.6 or Chapter 6).



⁵On sparse reconstruction from Fourier and Gaussian measurements, M. Rudelson and R. Vershynin. Comm. on Pure and Applied Mathematics, 2008.

RIP of Random Unitary Matrices

Motivating example: recall the MRI sensing model:

$$\mathbf{y} = \mathbf{F}_\Omega \mathbf{\Psi} \mathbf{x}, \quad \text{with } \mathbf{F} \text{ Fourier and } \mathbf{\Psi} \text{ wavelet.}$$

Theorem

Let $\mathbf{U} \in \mathbb{C}^{n \times n}$ be unitary ($\mathbf{U}^* \mathbf{U} = \mathbf{I}$) and Ω is a random subset of m elements from $\{1, \dots, n\}$. Suppose that

$$\|\mathbf{U}\|_\infty \leq \zeta / \sqrt{n}. \quad (23)$$

If

$$m \geq \frac{C\zeta^2}{\delta^2} k \log^4(n), \quad (24)$$

then with high probability, $\mathbf{A} = \sqrt{\frac{n}{m}} \mathbf{U}_{\Omega, \bullet}$ satisfies the RIP of order k , with constant $\delta_k(\mathbf{A}) \leq \delta$.

Circulant Convolution Matrices

A (random) circulant convolution:

$$\mathbf{r} * \mathbf{x} = \begin{bmatrix} r_0 & r_{n-1} & \dots & r_2 & r_1 \\ r_1 & r_0 & r_{n-1} & & r_2 \\ \vdots & r_1 & r_0 & \ddots & \vdots \\ r_{n-2} & & \ddots & \ddots & r_{n-1} \\ r_{n-1} & r_{n-2} & \dots & r_1 & r_0 \end{bmatrix} \mathbf{x} \doteq \mathbf{R}\mathbf{x}. \quad (25)$$

Fact: any circulant matrix can be diagonalized by the discrete Fourier transform:

$$\mathbf{R} = \mathbf{F}\mathbf{D}\mathbf{F}^*.$$

Select a (random) subset of the measurements:

$$\mathbf{y} = \mathcal{P}_\Omega[\mathbf{r} * \mathbf{x}] = \mathbf{A}\mathbf{x}, \quad (26)$$

RIP of Random Circulant Convolution Matrices

Theorem

Let $\Omega \subseteq \{1, \dots, n\}$ be any fixed subset of size $|\Omega| = m$. Then if

$$m \geq \frac{Ck \log^2(k) \log^2(n)}{\delta^2}, \quad (27)$$

then with high probability, \mathbf{A} has RIP of order k with $\delta_k(\mathbf{A}) \leq \delta$.

Approximate isometric property is the key to deep convolution neural networks!⁶

⁶Deep Isometric Learning for Visual Recognition, H. Qi, C. You, X. Wang, Yi Ma, and J. Malik, ICML 2020.

Assignments

- Reading: Section 3.4 of Chapter 3.