

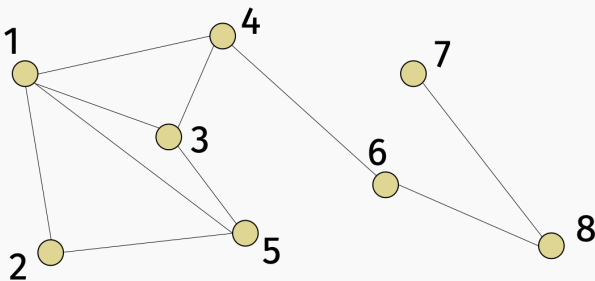
CS-GY 6763: Lecture 12

Stochastic Block Model, subspace embeddings + ϵ -net arguments

NYU Tandon School of Engineering, Prof. Christopher Musco

SPECTRAL GRAPH THEORY

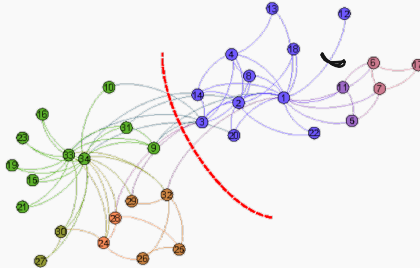
Main idea: Understand graph data by constructing natural matrix representations, and studying that matrix's spectrum (eigenvalues/eigenvectors).



$G = (V, E)$ is an undirected, unweighted graph with n nodes.

Goal: Given a graph $G = (V, E)$, partition nodes along a cut that:

- Has few crossing edges: $|\{(u, v) \in E : u \in S, v \in T\}|$ is small.
- Separates large partitions: $|S|, |T|$ are not too small.



(a) Zachary Karate Club Graph

Applications: Understanding community structure in social networks, partitioning finite element meshes, non-linear clustering in machine learning, data visualization, etc. etc.

β -Balanced Cut:

$$\min_S \text{cut}(S, V \setminus S) \text{ such that } \min(|S|, |V \setminus S|) \geq \beta \cdot n \text{ for } \beta \leq .5$$

Handwritten annotations: A red arrow points from the subscript S to the cut function. Above the min operator, there are two handwritten arrows: one pointing to the first min and another pointing to the second min.

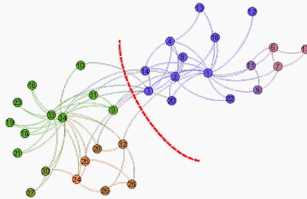
Last class we focused on the (extreme case where $\beta = 1/2$.)

Basic spectral clustering method:

- (Compute second smallest eigenvector of graph, \mathbf{v}_{n-1} .
- (\mathbf{v}_{n-1}) has an entry for every node i in the graph.
- If the i^{th} entry is positive, put node i in T .
- Otherwise if the i^{th} entry is negative, put i in S .

→ Laplacian

THE LAPLACIAN VIEW



(a) Zachary Karate Club Graph

For a cut indicator vector $\mathbf{c} \in \{-1, 1\}^n$ with $\mathbf{c}(i) = -1$ for $i \in S$ and $\mathbf{c}(i) = 1$ for $i \in T$:

$$\mathbf{c}^T \mathbf{L} \mathbf{c} = 4 \cdot \text{cut}(S, T).$$

$$\mathbf{c}^T \mathbf{1} = |T| - |S| = 0$$

$$\mathbf{c} \in \{-1/\sqrt{n}, 1/\sqrt{n}\}$$

Want to minimize both $\mathbf{c}^T \mathbf{L} \mathbf{c}$ (cut size) and $|\mathbf{c}^T \mathbf{1}|$ (imbalance).

Perfectly balanced balanced cut problem:

$$\min_{c \in \left\{ -\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}} \right\}^n} c^T L c \text{ such that } c^T \mathbf{1} = 0.$$

$\frac{1}{\sqrt{n}}$ is cut value

$$\downarrow$$

$$|\mathbf{S}| = |\mathbf{T}|$$

Relaxed perfectly balanced balanced cut problem:

$$\left(\min_{\|c\|_2=1} c^T L c \text{ such that } c^T \mathbf{1} = 0. \right)$$

Main result: The relaxed problem is exactly minimized by the second smallest eigenvector \mathbf{v}_{n-1} of L .

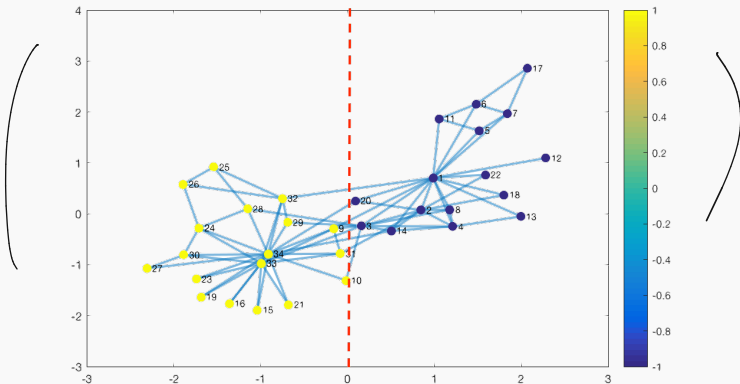
CUTTING WITH THE SECOND LAPLACIAN EIGENVECTOR

Final relax and round algorithm: Compute

$$\mathbf{v}_{n-1} = \arg \min_{\mathbf{v} \in \mathbb{R}^n \text{ with } \|\mathbf{v}\|=1, \mathbf{v}^T \mathbf{1}=0} \mathbf{v}^T \mathbf{L} \mathbf{v}$$



Set S to be all nodes with $\mathbf{v}_{n-1}(i) < 0$, and T to be all with $\mathbf{v}_{n-1}(i) \geq 0$. I.e. set $\mathbf{c} = \text{sign}(\mathbf{v}_{n-1})$



So far: Showed that spectral clustering partitions a graph along a small cut between large pieces.

$$O(n^3)$$

- No formal guarantee on the ‘quality’ of the partitioning.
- Can fail for worst case input graphs.)

Common approach: Design a natural (generative model) that produces random but realistic inputs and analyze how the algorithm performs on inputs drawn from this model.

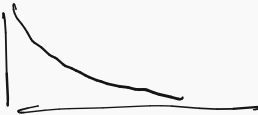
- Very common in algorithm design and analysis. Great way to start approaching a problem. Often our best way to understand why some algorithms “just work” in practice.
- (Similar approach to Bayesian modeling in machine learning.

STOCHASTIC BLOCK MODEL

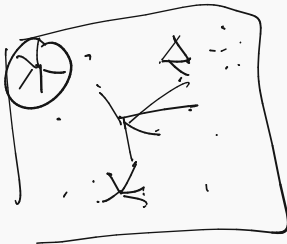
n nodes

Ideas for a generative model for **social network graphs** that would allow us to understand partitioning?

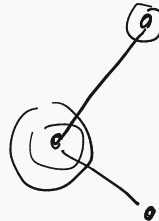
Graphs
(Geometric)



$$X^{-2}$$



0

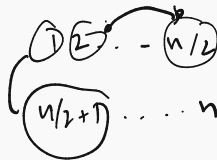
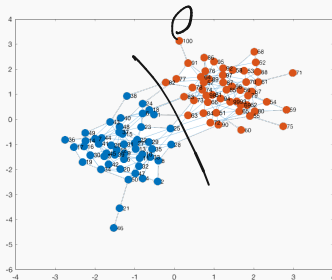
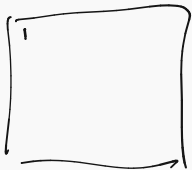


STOCHASTIC BLOCK MODEL

Stochastic Block Model (Planted Partition Model): $p, q \in [0, 1]$

Let $G_n(\underline{p}, \underline{q})$ be a distribution over graphs on n nodes, split equally into two groups B and C , each with $n/2$ nodes.

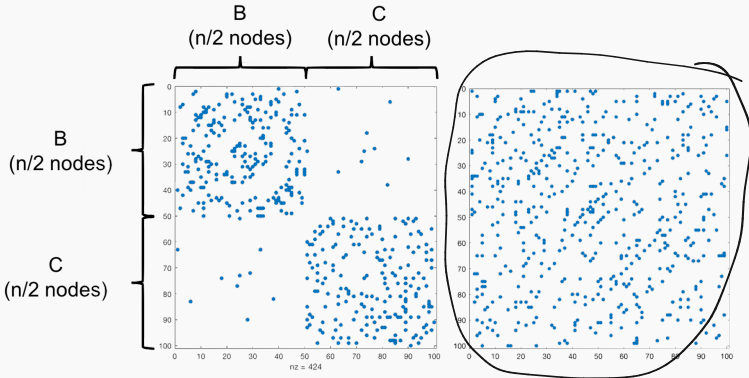
- Any two nodes in the **same group** are connected with probability \underline{p} (including self-loops).
- Any two nodes in **different groups** are connected with prob. $q < p$.



LINEAR ALGEBRAIC VIEW

Let \underline{G} be a stochastic block model graph drawn from $G_n(p, q)$.

- Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ denote the adjacency matrix of G .

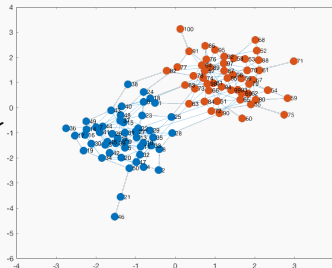


Note that we are arbitrarily ordering the nodes in \mathbf{A} by group. In reality \mathbf{A} would look “scrambled” as on the right.

STOCHASTIC BLOCK MODEL

Goal is to find the “ground truth” balanced partition $\underline{B}, \underline{C}$ using our standard spectral method.

\textcircled{L} → second
smallest
eigenvector



$\tilde{B} \quad \tilde{C}$

To do so, we need to understand the second smallest eigenvector of $L = D - A$. We will start by considering the expected value of these matrices:

$$E[L]_{ij} = E[L_{ij}]$$

is matrix

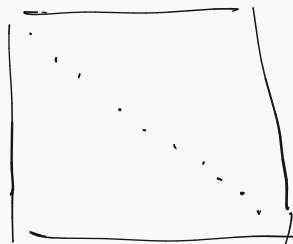
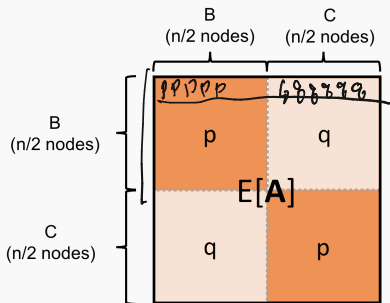
$$\textcircled{E[L]} = \underline{E[D]} - \underline{E[A]}$$

EXPECTED ADJACENCY SPECTRUM

$$\mathbb{E}[L] = \mathbb{E}[D] - \mathbb{E}[A]$$

Letting G be a stochastic block model graph drawn from $G_n(p, q)$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$ be its adjacency matrix. $(\mathbb{E}[\mathbf{A}])_{i,j} = p$ for i, j in same group, $(\mathbb{E}[\mathbf{A}])_{i,j} = q$ otherwise.

$$\mathbb{E}[D] = \frac{(p+q)n}{2} \cdot \mathbf{I}$$



$$\mathbb{E} \left[\sum_{i=1}^n \mathbb{I}(\text{node } i \text{ connects to node } j) \right]$$

EXPECTED LAPLACIAN

What is the expected Laplacian of $G_n(p, q)$?

$$\mathbb{E}[L] = \mathbb{E}[D] - \mathbb{E}[A] = \underline{cI} - \mathbb{E}[A]$$

Suppose v is eigenvector of $\mathbb{E}[L]$.

$$\mathbb{E}[L]v = \lambda v \quad \text{for some } \lambda \in \mathbb{R}.$$

$$\begin{aligned} (\mathbb{E}[D] - \mathbb{E}[A])v &= \lambda v &= -\mathbb{E}[A]v &= \lambda v - \mathbb{E}[D]v \\ & &= \lambda v - cv \end{aligned}$$

$$\lambda_1 \rightarrow c - \lambda_1$$

$$\lambda_2 \rightarrow c - \lambda_2$$

$$\vdots$$

$$\lambda_n \rightarrow c - \lambda_n$$

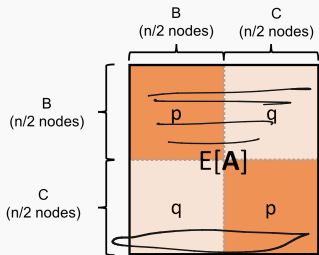
$$\mathbb{E}[A]v = \underline{\underline{(c - \lambda)v}}.$$

$\mathbb{E}[A]$ and $\mathbb{E}[L]$ have the same eigenvectors and eigenvalues are equal up to a shift/inversion. So (second largest eigenvector) of $\mathbb{E}[A]$ is the same as the second smallest of $\mathbb{E}[L]$

EXPECTED ADJACENCY SPECTRUM

Letting G be a stochastic block model graph drawn from $G_n(p, q)$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$ be its adjacency matrix, what are the eigenvectors and eigenvalues of $\mathbb{E}[\mathbf{A}]$?

$$\mathbf{M}\mathbf{v} = \lambda\mathbf{v}.$$



$$\begin{bmatrix} 1 \\ \vdots \\ -1 \\ \vdots \\ -1 \end{bmatrix}$$

$$= \begin{bmatrix} (p-q) n/2 \\ (p-q) n/2 \\ \vdots \\ (q-p) n/2 \end{bmatrix}$$

$$= \frac{(p-q)n}{2} \begin{bmatrix} 1 \\ \vdots \\ -1 \\ \vdots \\ -1 \end{bmatrix}$$

$n-2$ eigenvalues

$= 0$

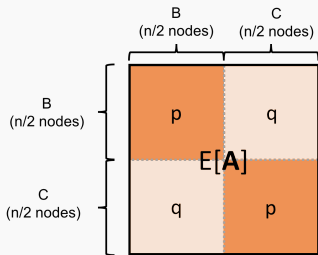
$$\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$= \begin{bmatrix} (p+q) n/2 \\ \vdots \\ (p+q) n/2 \end{bmatrix}$$

$$= \frac{(p+q)n}{2} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

EXPECTED ADJACENCY SPECTRUM

Letting G be a stochastic block model graph drawn from $G_n(p, q)$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$ be its adjacency matrix, what are the eigenvectors and eigenvalues of $\mathbb{E}[\mathbf{A}]$?



EXPECTED ADJACENCY SPECTRUM

Diagram illustrating the spectral decomposition of the expected adjacency matrix $E[A]$.

The matrix $E[A]$ is a 2×2 block matrix with blocks p and q , representing communities B and C (each with $n/2$ nodes).

The decomposition is shown as:

$$E[A] = V \Lambda V^T$$

Where:

- V is a matrix of eigenvectors (columns).
- Λ is a diagonal matrix of eigenvalues.
- V^T is the transpose of V .

The eigenvalues are:

- $\lambda_1 = \frac{(p+q)n}{2}$
- $\lambda_2 = \frac{(p-q)n}{2}$

The eigenvectors are:

- $\bar{v}_1 \sim \mathbf{1}$ (all ones)
- $\bar{v}_2 \sim \chi_{B,C}$ (difference between communities)

- $\bar{v}_1 \sim \mathbf{1}$ with eigenvalue $\lambda_1 = \frac{(p+q)n}{2}$.
- $\bar{v}_2 \sim \chi_{B,C}$ with eigenvalue $\lambda_2 = \frac{(p-q)n}{2}$.

If we compute \bar{v}_2 then we exactly recover the communities B and C !

Upshot: The second smallest eigenvector of $\mathbb{E}[L]$, equivalently the second largest of $\mathbb{E}[A]$, is exactly $\chi_{B,C}$ – the indicator vector for the cut between the communities.

- If the random graph G (equivalently A and L) were exactly equal to its expectation, partitioning using this eigenvector would exactly recover communities B and C .

How do we show that a matrix (e.g. A) is close to its expectation? (Matrix concentration inequalities.)

- Analogous to scalar concentration inequalities like Markovs, Chebyshevs, Bernsteins.

$$\|A - \mathbb{E}[A]\|_2$$

Alon, Krivelevich, Vu, 2002:



Matrix Concentration Inequality: If $p \geq O\left(\frac{\log^4 n}{n}\right)$, then with high probability $1 - \frac{1}{p^2 n^2}$

$$\|A - \mathbb{E}[A]\|_2 \leq O(\sqrt{pn}). \quad = \frac{\|A\|_2}{\sqrt{np}}$$

where $\|\cdot\|_2$ is the matrix **spectral** norm (operator norm).

Recall that $\|X\|_2 = \max_{Z \in \mathbb{R}^n: \|Z\|_2=1} \|XZ\|_2 = \underline{\underline{\sigma_1(X)}}$.

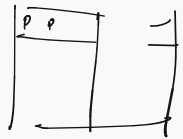
$\|A\|_2$ is on the order of $O(pn)$ so another way of thinking about the right hand side is $\frac{\|A\|_2}{\sqrt{np}}$. I.e. get's better with p and n .

$$\|A\|_2 \approx O(pn)$$

MATRIX CONCENTRATION

$\|A\|_2$ is on the order of $O(pn)$ so another way of thinking about the right hand side is $\frac{\|A\|_2}{\sqrt{np}}$. I.e. get's better with p and n .

$$\frac{\|A_1\|_2}{\|I\|_2} \leq \|A\|_2 = \max_z \frac{\|A_2\|_2}{\|z\|_2}$$




$$\sqrt{\frac{(p+\delta)^2 \cdot n^2 \cdot n}{n}}$$

$$\frac{(p+\delta) n^{1.5}}{\sqrt{n}} = \underline{\underline{(p+\delta)n}}$$

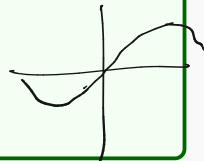
EIGENVECTOR PERTURBATION

For the stochastic block model application, we want to show that the second eigenvectors of \mathbf{A} and $\mathbb{E}[\mathbf{A}]$ are close. How does this relate to their difference in spectral norm?

Davis-Kahan Eigenvector Perturbation Theorem: Suppose $\underline{\mathbf{A}}, \underline{\bar{\mathbf{A}}} \in \mathbb{R}^{d \times d}$ are symmetric with $\|\underline{\mathbf{A}} - \underline{\bar{\mathbf{A}}}\|_2 \leq \underline{\epsilon}$ and eigenvectors $\underline{\mathbf{v}}_1, \underline{\mathbf{v}}_2, \dots, \underline{\mathbf{v}}_n$ and $\underline{\bar{\mathbf{v}}}_1, \underline{\bar{\mathbf{v}}}_2, \dots, \underline{\bar{\mathbf{v}}}_n$. Letting $\theta(\underline{\mathbf{v}}_i, \underline{\bar{\mathbf{v}}}_i)$ denote the angle between $\underline{\mathbf{v}}_i$ and $\underline{\bar{\mathbf{v}}}_i$, for all i :


$$\sin[\theta(\underline{\mathbf{v}}_i, \underline{\bar{\mathbf{v}}}_i)] \leq \frac{\epsilon}{\min_{j \neq i} |\lambda_i - \lambda_j|}$$

where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of $\underline{\bar{\mathbf{A}}}$.



We will apply with $\underline{\bar{\mathbf{A}}} = \mathbb{E}[\mathbf{A}]$.

EIGENVECTOR PERTURBATION

$$\begin{array}{ccc}
 \mathbf{A} & \mathbf{\bar{A}} & \|\mathbf{A} - \mathbf{\bar{A}}\|_2 = \epsilon \\
 \begin{bmatrix} \underline{1+\epsilon} & 0 \\ 0 & \underline{1} \end{bmatrix} & \begin{bmatrix} \underline{1} & 0 \\ 0 & \underline{1+\epsilon} \end{bmatrix} & = \begin{bmatrix} \epsilon & 0 \\ 0 & \epsilon \end{bmatrix} \\
 \begin{array}{c} \left[\begin{array}{c|c} 1 & 0 \\ 0 & 1 \end{array} \right] \end{array} & \begin{array}{c} \left[\begin{array}{c|c} 1 & 0 \\ 0 & 1 \end{array} \right] \end{array} & \begin{array}{cc} \left[\begin{array}{c} 1 \\ 0 \end{array} \right] & \left[\begin{array}{c} 0 \\ 1 \end{array} \right] \end{array} \\
 \begin{array}{cc} \left(\begin{bmatrix} 1 \\ 0 \end{bmatrix}, (1+\epsilon) \right) & \left(\begin{bmatrix} 0 \\ 1 \end{bmatrix}, 1 \right) \\ \downarrow & \downarrow \\ v_1 & v_2 \end{array} & \begin{array}{cc} \left(\begin{bmatrix} 0 \\ 1 \end{bmatrix}, 1 \right), \left(\begin{bmatrix} 1 \\ 0 \end{bmatrix}, 1+\epsilon \right) \\ \downarrow & \downarrow \\ \bar{v}_2 & \bar{v}_1 \end{array}
 \end{array}$$

APPLICATION TO STOCHASTIC BLOCK MODEL

Claim 1 (Matrix Concentration): For $p \geq O\left(\frac{\log^4 n}{n}\right)$, $\mathbb{E}[A] = \bar{A}$

$$\|A - \mathbb{E}[A]\|_2 \leq O(\sqrt{pn}).$$

Recall: $\mathbb{E}[A]$, has eigenvalues $\lambda_1 = \frac{(p+q)n}{2}$, $\lambda_2 = \frac{(p-q)n}{2}$, $\lambda_i = 0$ for $i \geq 3$.

$$\min_{j \neq i} |\lambda_i - \lambda_j| = \min \left(qn, \frac{(p-q)n}{2} \right).$$

$$\left((p+q) - (p-q) \right) \frac{n}{2}$$

Assume $\frac{(p-q)n}{2}$ will be the minimum of these two gaps.

Claim 2 (Davis-Kahan): For $p \geq O\left(\frac{\log^4 n}{n}\right)$,

$$O(1/\sqrt{n})$$

$$\sin \theta(\mathbf{v}_2, \bar{\mathbf{v}}_2) \leq \frac{O(\sqrt{pn})}{\min_{j \neq i} |\lambda_i - \lambda_j|} \leq \frac{O(\sqrt{pn})}{(p-q)n/2} = O\left(\frac{\sqrt{p}}{(p-q)\sqrt{n}}\right)$$

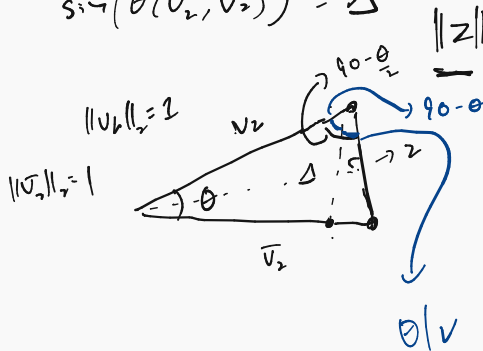
(A slightly trickier analysis can remove the qn term entirely.)

APPLICATION TO STOCHASTIC BLOCK MODEL

So far: $\sin \theta(\underline{v}_2, \bar{v}_2) \leq O\left(\frac{\sqrt{p}}{(p-q)\sqrt{n}}\right)$. $\rightarrow = \Delta$ What does this give us?

- Can show that this implies $\|\underline{v}_2 - \bar{v}_2\|_2^2 \leq O\left(\frac{p}{(p-q)^2 n}\right)$.

$$\sin(\theta(\underline{v}_2, \bar{v}_2)) = \Delta$$



$$\|z\|_2 \|\underline{v}_2 - \bar{v}_2\|_2$$

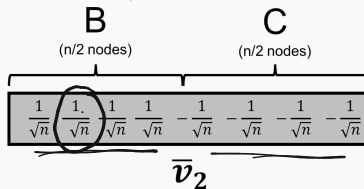
$$\|z\|_2 \leq \Delta + \frac{\Delta}{2}$$

$$\|z\|_2^2 = O(\Delta^2)$$

APPLICATION TO STOCHASTIC BLOCK MODEL

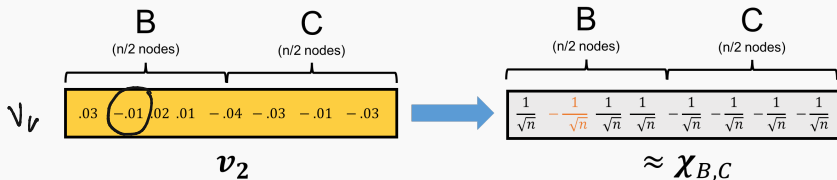
So far: $\|\underline{v_2} - \bar{v}_2\|_2^2 \leq O\left(\frac{p}{(p-q)^2 n}\right)$ $\frac{1}{\sqrt{n}}$ sign(v_2)

- \bar{v}_2 is $\frac{1}{\sqrt{n}} \chi_{B,C}$: the community indicator vector.



$\approx 1/\sqrt{n}$

- We want to show that $\text{sign}(v_2)$ and \bar{v}_2 are close. They only differ at locations where v_2 and \bar{v}_2 differ in sign.



$1/\sqrt{n}$ $-1/\sqrt{n}$ $1/\sqrt{n} \dots$

Main argument:

- Every i where $\underline{v}_2(i), \underline{\bar{v}}_2(i)$ differ in sign contributes $\geq \frac{1}{n}$ to $\|\mathbf{v}_2 - \bar{\mathbf{v}}_2\|_2^2$.
- We know that $\|\mathbf{v}_2 - \bar{\mathbf{v}}_2\|_2^2 \leq O\left(\frac{p}{(p-q)^2 n}\right)$
- So \mathbf{v}_2 and $\bar{\mathbf{v}}_2$ differ in sign in at most $O\left(\frac{p}{(p-q)^2}\right)$ positions.

$$p = \frac{250}{n}$$

Upshot: If G is a stochastic block model graph with adjacency matrix A , if we compute its second largest eigenvector \mathbf{v}_2 and assign nodes to communities according to the sign pattern of this vector, we will correctly assign all but $O\left(\frac{p}{(p-q)^2}\right)$ nodes.

- **Hard case:** Suppose $q = .8p$ so $\frac{p}{(p-q)^2} = 25/p$. Even if p is really small, i.e. $p = 250/n$, then we assign roughly 90% of nodes to the right partition.

$$\frac{p}{(p-.8p)^2} = \frac{p}{(p-.8p)^2} = \frac{p}{(0.2p)^2} = \frac{25}{p} = \frac{25}{250/n} = .1 \cdot n.$$

Forget about the previous problem, but still consider the matrix $\mathbf{M} = \mathbb{E}[\mathbf{A}]$.

- Dense $n \times n$ matrix.
- Computing top eigenvectors takes $\approx \overline{O(n^2/\sqrt{\epsilon})}$ time.

$$\begin{pmatrix} p & q \\ q & p \end{pmatrix}$$

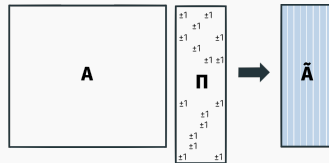
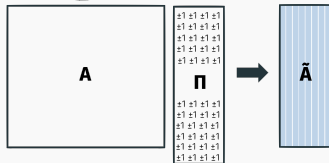
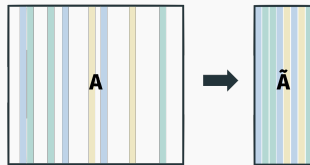
If someone asked you to speed this up and return approximate top eigenvectors, what could you do?

$$\text{num}(\mathbf{A}) \cdot (\# \text{ iterations})$$

RANDOMIZED NUMERICAL LINEAR ALGEBRA

Main idea: If you want to compute singular vectors, multiply two matrices, solve a regression problem, etc.: *5:30pm*

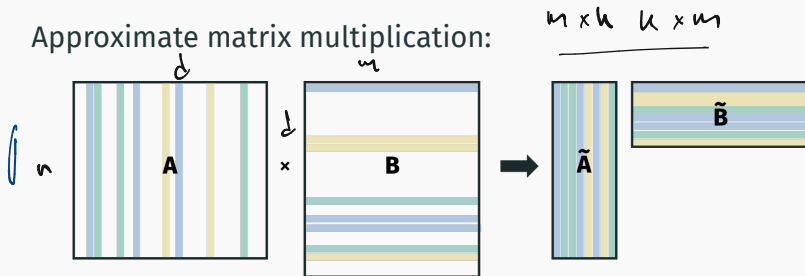
1. Compress your matrices using a randomized method (e.g. subsampling).
2. Solve the problem on the smaller or sparser matrix.
 - \tilde{A} called a “sketch” or “coreset” for A .



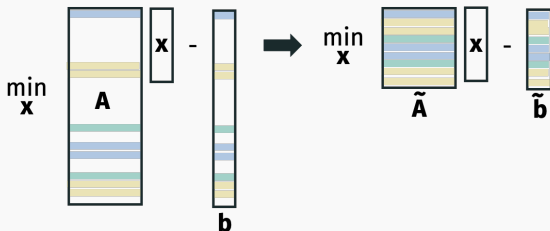
BREAK

RANDOMIZED NUMERICAL LINEAR ALGEBRA

Approximate matrix multiplication:



(Approximate regression)



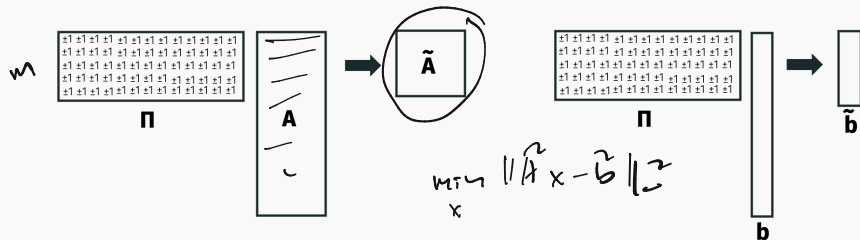
SKETCHED REGRESSION

$$O(nd^2)$$

$$O(m d^2)$$

$$m \approx O(d)$$

Today's example: Randomized approximate regression using a Johnson-Lindenstrauss matrix for compression.



Input: $A \in \mathbb{R}^{n \times d}$, $b \in \mathbb{R}^n$.

Goal: Let $x^* = \arg \min_x \|Ax - b\|_2^2$. Let $\tilde{x} = \arg \min_x \|\Pi Ax - \Pi b\|_2^2$

$$\text{Want: } \underbrace{\|A\tilde{x} - b\|_2^2}_{\leq (1+\epsilon) \underbrace{\|Ax^* - b\|_2^2}}$$

Theorem (Randomized Linear Regression)

Let Π be a JL matrix (random Gaussian, sign, sparse random, etc.) with $m = O\left(\frac{d}{\epsilon^2}\right)$ rows¹. Then with probability 9/10, for any $\mathbf{A} \in \mathbb{R}^{n \times d}$ and $\mathbf{b} \in \mathbb{R}^n$,

$$\|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{b}\|_2^2 \leq (1 + \epsilon) \|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|_2^2$$

where $\tilde{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\Pi \mathbf{A} \mathbf{x} - \Pi \mathbf{b}\|_2^2$.

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} \|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2^2$$

¹This can be improved to $O(d/\epsilon)$ with a tighter analysis

- Prove this theorem using an (ϵ -net argument), which is a popular technique for applying our standard concentration inequality + union bound argument to an (infinite number of events.)
- These sort of arguments appear all the time in theoretical algorithms and ML research, so this part of lecture is as much about the technique as the final result.

SKETCHED REGRESSION

Claim: Suffices to prove that for all $x \in \mathbb{R}^d$

\tilde{x} : argmin $\|\pi Ax - \pi b\|_2^2$
x

$$(1 - \epsilon) \|Ax - b\|_2^2 \leq \|\pi Ax - \pi b\|_2^2 \leq (1 + \epsilon) \|Ax - b\|_2^2$$

$$\|A\tilde{x} - b\|_2^2 \leq (1 + \epsilon) \|Ax^* - b\|_2^2$$

$$\|A\tilde{x} - b\|_2^2 \leq \frac{1}{1 - \epsilon} \|\pi A\tilde{x} - \pi b\|_2^2 \leq \frac{1}{1 - \epsilon} \|\pi Ax^* - \pi b\|_2^2$$

↑
optimality of \tilde{x}

$$\leq \frac{(1 + \epsilon)}{(1 - \epsilon)} \|Ax^* - b\|_2^2 \approx \underline{(1 + 2\epsilon)} \|Ax^* - b\|_2^2$$

Lemma (Distributional JL)

If Π is chosen to a random Gaussian matrix, sign matrix, sparse random matrix, etc. (scaled by $1/\sqrt{m}$) with $m = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ rows then for any fixed y ,

$$(1 - \epsilon)\|y\|_2^2 \leq \|\Pi y\|_2^2 \leq (1 + \epsilon)\|y\|_2^2$$

with probability $(1 - \delta)$.

Corollary: For any fixed x , with probability $(1 - \delta)$,

$$(1 - \epsilon)\|\underline{Ax - b}\|_2^2 \leq \|\Pi Ax - \Pi b\|_2^2 \leq (1 + \epsilon)\|Ax - b\|_2^2.$$

$$= \|\Pi(Ax - b)\|_2^2$$

FOR ANY TO FOR ALL

$$\|Ax - b\|$$

How do we go from “for any fixed x ” to “for all $x \in \mathbb{R}^d$ ”.

This statement requires establishing a Johnson-Lindenstrauss type bound for an infinity of possible vectors ($Ax - b$) which can't be tackled directly with a union bound argument.



Note that all vectors of the form $(Ax - b)$ lie in a low dimensional subspace: spanned by $d + 1$ vectors, where d is the width of A . **So even though the set is infinite, it is “simple” in some way. Parameterized by just $d + 1$ numbers.**

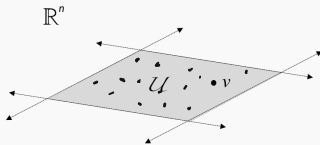
SUBSPACE EMBEDDINGS

Theorem (Subspace Embedding from JL) [Sarlos, 2006]

Let $\mathcal{U} \subset \mathbb{R}^n$ be a d -dimensional linear subspace in \mathbb{R}^n . If $\Pi \in \mathbb{R}^{m \times d}$ is chosen from any distribution \mathcal{D} satisfying the Distributional JL Lemma, then with probability $1 - \delta$,

$$(1 - \epsilon) \|\underline{\mathbf{v}}\|_2^2 \leq \|\underline{\Pi \mathbf{v}}\|_2^2 \leq (1 + \epsilon) \|\underline{\mathbf{v}}\|_2^2$$

for all $\mathbf{v} \in \mathcal{U}$, as long as $m = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)$.² $\approx O\left(\frac{d}{\epsilon^2}\right)$



²It's possible to obtain a slightly tighter bound of $O\left(\frac{d + \log(1/\delta)}{\epsilon^2}\right)$. It's a nice challenge to try proving this.

SUBSPACE EMBEDDING TO APPROXIMATE REGRESSION

Corollary: If we choose Π and properly scale, then with $O(d/\epsilon^2)$ rows,

$$\left((1 - \epsilon) \|Ax - b\|_2^2 \leq \|\Pi Ax - \Pi b\|_2^2 \leq (1 + \epsilon) \|Ax - b\|_2^2 \right)$$

for all x and thus

$$Ax - b \quad d/r$$

$$\left(\|A\tilde{x} - b\|_2^2 \leq (1 + O(\epsilon)) \min_x \|Ax - b\|_2^2 \right)$$

i.e., our main theorem is proven.

Proof: Apply Subspace Embedding Thm. to the $(d + 1)$ dimensional subspace spanned by A 's d columns and b . Every vector $Ax - b$ lies in this subspace.

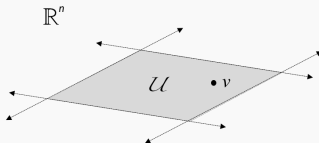
SUBSPACE EMBEDDINGS

Theorem (Subspace Embedding from JL) †

Let $\mathcal{U} \subset \mathbb{R}^n$ be a d -dimensional linear subspace in \mathbb{R}^n . If $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ is chosen from any distribution \mathcal{D} satisfying the Distributional JL Lemma, then with probability $1 - \delta$,

$$(1 - \epsilon)\|\mathbf{v}\|_2^2 \leq \|\mathbf{\Pi}\mathbf{v}\|_2^2 \leq (1 + \epsilon)\|\mathbf{v}\|_2^2 \quad (1)$$

for all $\mathbf{v} \in \mathcal{U}$, as long as $m = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)$



Subspace embeddings have tons of other applications!

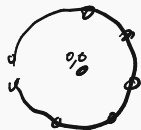
SUBSPACE EMBEDDING PROOF

$$w = c_1 e_1 + c_2 e_2 + \dots + c_d e_d \quad \Pi w = c_1 \Pi e_1 \dots c_d \Pi e_d$$

$$(1 - \epsilon) \|\underline{v}\|_2^2 \leq \|\underline{\Pi v}\|_2^2 \leq (1 + \epsilon) \|\underline{v}\|_2^2 \quad (2)$$

First Observation: The theorem holds as long as (1) holds for all w on the unit sphere in \mathcal{U} . Denote the sphere $S_{\mathcal{U}}$:

$$\underline{S_{\mathcal{U}}} = \{\underline{w} \mid \underline{w} \in \mathcal{U} \text{ and } \|\underline{w}\|_2 = 1\}.$$



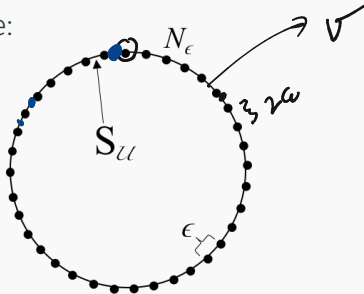
Follows from linearity: Any point $v \in \mathcal{U}$ can be written as cw for some scalar c and some point $w \in S_{\mathcal{U}}$.

$$v = c \cdot w$$

- If $(1 - \epsilon) \|\underline{w}\|_2 \leq \|\underline{\Pi w}\|_2 \leq (1 + \epsilon) \|\underline{w}\|_2$.
 - then $c(1 - \epsilon) \|\underline{w}\|_2 \leq c \|\underline{\Pi w}\|_2 \leq c(1 + \epsilon) \|\underline{w}\|_2$,
 - and thus $(1 - \epsilon) \|\underline{cw}\|_2 \leq \|\underline{\Pi cw}\|_2 \leq (1 + \epsilon) \|\underline{cw}\|_2$.
- $\|\Pi w\| \leq \|c_1 \Pi e_1\| + \dots + \|c_d \Pi e_d\|$
- \downarrow \downarrow \downarrow
 \checkmark \checkmark \checkmark

SUBSPACE EMBEDDING PROOF

Intuition: There are not too many “different” points on a d -dimensional sphere:



Goal: Find a set \underline{N}_ϵ such that, for every $\underline{v} \in \underline{S}_U$, there is some point $\underline{w} \in \underline{N}_\epsilon$ such that $\|\underline{w} - \underline{v}\|_2 \leq \epsilon$. N_ϵ is called an “ ϵ ”-net. If we can prove

$$(1 - \epsilon)\|\underline{w}\|_2 \leq \|\Pi \underline{w}\|_2 \leq (1 + \epsilon)\|\underline{w}\|_2$$

for all points $\underline{w} \in \underline{N}_\epsilon$, we can hopefully extend to all of S_U .

Lemma (ϵ -net for the sphere)

For any $\epsilon \leq 1$, there exists a set $N_\epsilon \subset S_{\mathcal{U}}$ with $|N_\epsilon| = \left(\frac{3}{\epsilon}\right)^d$ such that $\forall \mathbf{v} \in S_{\mathcal{U}}$,

$$\min_{\mathbf{w} \in N_\epsilon} \|\mathbf{v} - \mathbf{w}\|_2 \leq \epsilon.$$

Take this claim to be true for now: we will prove later.

$$\log\left(\left(\frac{3}{\epsilon}\right)^d\right) = d \log(3/\epsilon)$$

1. Preserving norms of all points in net N_ϵ .

$1 - \delta'$

$O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$

Set $\delta' = \frac{1}{|N_\epsilon|} \cdot \delta = \left(\frac{\epsilon}{3}\right)^d \cdot \delta$. As long as Π has $O\left(\frac{\log(1/\delta')}{\epsilon^2}\right)$
 $= O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)$ rows, then by a union bound,

$$(1 - \epsilon)\|\mathbf{w}\|_2 \leq \|\Pi \mathbf{w}\|_2 \leq (1 + \epsilon)\|\mathbf{w}\|_2. \quad \text{)}$$

for all $\mathbf{w} \in N_\epsilon$, with probability $1 - \delta$.

$$\delta' = \frac{1}{|N_\epsilon|} \cdot \delta = \left(\frac{\epsilon}{3}\right)^d \cdot \delta$$

SUBSPACE EMBEDDING PROOF

$$v - w_0$$

2. Extending to all points in the sphere.

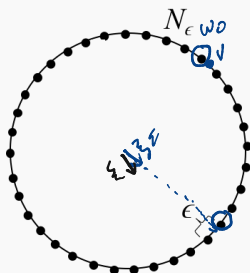
For some $w_0, w_1, w_2 \dots \in N_\epsilon$, any $v \in S_{\mathcal{U}}$ can be written:

$$v = \underline{1} \underline{w_0} + c_1 \underline{w_1} + c_2 \underline{w_2} + \dots$$

for constants c_1, c_2, \dots where $|c_i| \leq \epsilon^i$.

$$r_i = (w_0 + \dots + c_{i-1} w_{i-1}) - v$$

$$\|r_i\|_2 \leq \epsilon^i$$



$$\begin{aligned} |c_1| &\leq \epsilon \\ |c_2| &\leq \epsilon^2 \\ &\vdots \\ |c_i| &\leq \epsilon^i \end{aligned}$$

$$\frac{r_i}{\|r_i\|}$$

\exists

w_i

$$\|w_i - \frac{r_i}{\|r_i\|}\| \leq \epsilon = \|w_i, \|r_i\| - r_i\|_2 \leq \epsilon \|r_i\|$$

2. Extending to all points in the sphere.

For some $\mathbf{w}_0, \mathbf{w}_1, \mathbf{w}_2 \dots \in N_\epsilon$, any $\mathbf{v} \in S_{\mathcal{U}}$ can be written:

$$\mathbf{v} = \mathbf{w}_0 + c_1 \mathbf{w}_1 + c_2 \mathbf{w}_2 + \dots$$

for constants c_1, c_2, \dots where $|c_i| \leq \epsilon^i$.

Greedy construction:

$$\mathbf{w}_0 = \min_{\mathbf{w} \in \mathcal{N}_\epsilon} \|\mathbf{v} - \mathbf{w}\|_2$$

$$\mathbf{r}_0 = \mathbf{v} - \mathbf{w}_0$$

$$\mathbf{w}_1 = \min_{\mathbf{w} \in \mathcal{N}_\epsilon} \left\| \frac{\mathbf{r}_0}{\|\mathbf{r}_0\|} - \mathbf{w} \right\|_2 \quad c_1 = \|\mathbf{r}_0\|_2 \quad \mathbf{r}_1 = \mathbf{v} - \mathbf{w}_0 - c_1 \mathbf{w}_1$$

$$\mathbf{w}_2 = \min_{\mathbf{w} \in \mathcal{N}_\epsilon} \left\| \frac{\mathbf{r}_1}{\|\mathbf{r}_1\|} - \mathbf{w} \right\|_2 \quad c_2 = \|\mathbf{r}_1\|_2 \quad \mathbf{r}_2 = \mathbf{v} - \mathbf{w}_0 - c_1 \mathbf{w}_1 - c_2 \mathbf{w}_2$$

\vdots

SUBSPACE EMBEDDING PROOF

$$(1-\epsilon) \leq \|\Pi v\|_2 \leq (1+\epsilon) \quad \square$$

2. Extending to all points in the sphere.

Applying triangle inequality, we have that: $v = w_0 + c_1 w_1 + c_2 w_2 + \dots$

$$\begin{aligned} \|\Pi v\|_2 &= \|\Pi w_0 + c_1 \Pi w_1 + c_2 \Pi w_2 + \dots\| \\ &\leq \|\Pi w_0\| + c_1 \|\Pi w_1\| + c_2 \|\Pi w_2\| + \dots \\ &\leq \|\Pi w_0\| + \epsilon \|\Pi w_1\| + \epsilon^2 \|\Pi w_2\| + \epsilon^3 \|\Pi w_3\| + \dots \\ &\leq (1+\epsilon) + \underbrace{\epsilon(1+\epsilon)}_{\leq 2} + \epsilon^2(1+\epsilon) + \dots \\ &\leq 1 + 4\epsilon. \end{aligned}$$

$$\begin{aligned} &\leq (1+\epsilon) + (1+\epsilon)(\underbrace{\epsilon + \epsilon^2 + \dots}_{\leq 2\epsilon}) \leq 1 + 4\epsilon \end{aligned}$$

$1 + \epsilon + 2\epsilon + 2\epsilon^2$

3. Preserving norm of v .

Similarly,

$$\begin{aligned}
 \|\Pi v\|_2 &= \|\Pi w_0 + c_1 \Pi w_1 + c_2 \Pi w_2 + \dots\| \\
 &\geq \|\Pi w_0\| - \epsilon \|\Pi w_1\| - \epsilon^2 \|\Pi w_2\| - \dots \\
 &\geq \underbrace{(1 - \epsilon)}_{\geq 1 - 4\epsilon} - \underbrace{\epsilon(1 + \epsilon) - \epsilon^2(1 + \epsilon) - \dots}_{O(\epsilon)}
 \end{aligned}$$

So we have proven

$$(1 - O(\epsilon)) \|\mathbf{v}\|_2 \leq \|\Pi \mathbf{v}\|_2 \leq (1 + O(\epsilon)) \|\mathbf{v}\|_2$$

for all $\mathbf{v} \in S_{\mathcal{U}}$, which in turn implies,

$$(1 - O(\epsilon)) \|\mathbf{v}\|_2^2 \leq \|\Pi \mathbf{v}\|_2^2 \leq (1 + O(\epsilon)) \|\mathbf{v}\|_2^2$$

Adjusting ϵ proves the Subspace Embedding theorem.

SUBSPACE EMBEDDINGS

Theorem (Subspace Embedding from JL)

Let $\mathcal{U} \subset \mathbb{R}^n$ be a d -dimensional linear subspace in \mathbb{R}^n . If $\Pi \in \mathbb{R}^{m \times d}$ is chosen from any distribution \mathcal{D} satisfying the Distributional JL Lemma, then with probability $1 - \delta$,

$$(1 - \epsilon) \|\mathbf{v}\|_2^2 \leq \|\Pi \mathbf{v}\|_2^2 \leq (1 + \epsilon) \|\mathbf{v}\|_2^2 \quad (3)$$

for all $\mathbf{v} \in \mathcal{U}$, as long as $m = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)$

(Subspace embeddings have many other applications!)

For example, if $m = O(k/\epsilon)$, $(\Pi \mathbf{A})$ can be used to compute an approximate partial SVD, which leads to a $(1 + \epsilon)$ approximate low-rank approximation for \mathbf{A} .

$$\|\mathbf{A} - \mathbf{A} \tilde{\mathbf{V}}_k \tilde{\mathbf{V}}_k^T\|_F \leq (1 + \epsilon) \|\mathbf{A} - \mathbf{A} \mathbf{V}_k \mathbf{V}_k^T\|_F$$

Lemma (ϵ -net for the sphere)

For any $\epsilon \leq 1$, there exists a set $N_\epsilon \subset S_{\mathcal{U}}$ with $|N_\epsilon| = \left(\frac{3}{\epsilon}\right)^d$ such that $\forall \mathbf{v} \in S_{\mathcal{U}}$,

$$\min_{\mathbf{w} \in N_\epsilon} \|\mathbf{v} - \mathbf{w}\| \leq \epsilon.$$

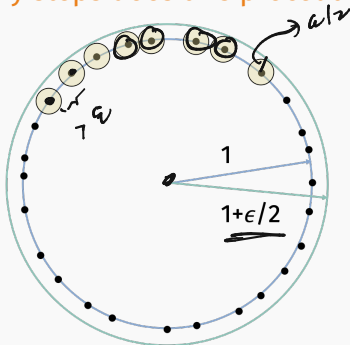
Imaginary algorithm for constructing N_ϵ :

- Set $N_\epsilon = \{\}$
- While such a point exists, choose an arbitrary point $\mathbf{v} \in S_{\mathcal{U}}$ where $\nexists \mathbf{w} \in N_\epsilon$ with $\|\mathbf{v} - \mathbf{w}\| \leq \epsilon$. Set $N_\epsilon = N_\epsilon \cup \{\mathbf{w}\}$.



After running this procedure, we have $N_\epsilon = \{\mathbf{w}_1, \dots, \mathbf{w}_{|N_\epsilon|}\}$ and $\min_{\mathbf{w} \in N_\epsilon} \|\mathbf{v} - \mathbf{w}\| \leq \epsilon$ for all $\mathbf{v} \in S_{\mathcal{U}}$ as desired.

How many steps does this procedure take?



Can place a ball of radius $\epsilon/2$ around each w_i without intersecting any other balls. All of these balls live in a ball of radius $1 + \epsilon/2$.

VOLUME ARGUMENT

Volume of d dimensional ball of radius r is

$$\text{vol}(d, r) = \underline{c} r^d,$$

$$\left(\frac{2}{\epsilon}\right)^d$$

where c is a constant that depends on d , but not r . From previous slide we have:

$$\begin{aligned} \underline{\text{vol}(d, \epsilon/2)} \cdot \underline{|N_\epsilon|} &\leq \underline{\text{vol}(d, 1 + \epsilon/2)} \\ \underline{|N_\epsilon|} &\leq \frac{\underline{\text{vol}(d, 1 + \epsilon/2)}}{\underline{\text{vol}(d, \epsilon/2)}} \\ &\leq \left(\frac{1 + \epsilon/2}{\epsilon/2}\right)^d \leq \left(\frac{3}{\epsilon}\right)^d \\ &\stackrel{\downarrow}{=} \left(\frac{2 + \epsilon}{\epsilon}\right)^d \leq \left(\frac{3}{\epsilon}\right)^d \end{aligned}$$

You can actually show that $m = O\left(\frac{d + \log(1/\delta)}{\epsilon}\right)$ suffices to be a d dimensional subspace embedding, instead of the bound we proved of $m = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon}\right)$.

The trick is to show that a constant factor net is actually all that you need instead of an ϵ factor.

RUNTIME CONSIDERATION

For $\epsilon, \delta = O(1)$, we need Π to have $m = O(d)$ rows.

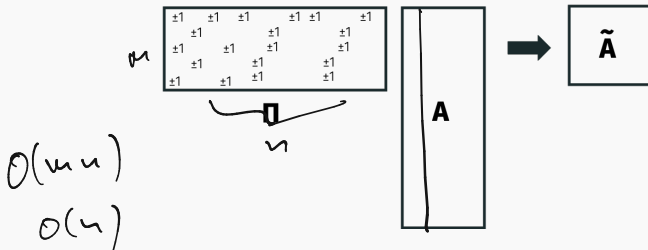
- Cost to solve $\|\mathbf{Ax} - \mathbf{b}\|_2^2$: $O(n d^2)$
 - $O(nd^2)$ time for direct method. Need to compute $(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$.
 - $(O(nd) \cdot (\# \text{ of iterations}))$ time for iterative method (GD, AGD, conjugate gradient method).
- Cost to solve $\|\Pi \mathbf{Ax} - \Pi \mathbf{b}\|_2^2$:
 - $O(d^3)$ time for direct method.
 - $O(d^2) \cdot (\# \text{ of iterations})$ time for iterative method.

$$\underline{\Pi A} \quad (d \times n) (n \times d) = \cancel{O(n d^2)} \\ \underline{O(nd)}$$

RUNTIME CONSIDERATION

But time to compute ΠA is an $(m \times n) \times (n \times d)$ matrix multiply: $O(mnd) = O(nd^2)$ time!

Goal: Develop faster Johnson-Lindenstrauss projections.



Typically using sparse and structured matrices.

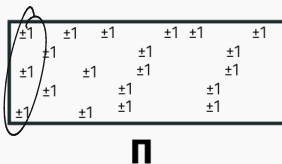
Next class: We will describe a construction where ΠA can be computed in $O(nd \log n)$ time.

RETURN TO SINGLE VECTOR PROBLEM

Goal: Develop methods that reduce a vector $\mathbf{x} \in \mathbb{R}^n$ down to $m \approx \frac{\log(1/\delta)}{\epsilon^2}$ dimensions in $o(mn)$ time and guarantee:

$$(1 - \epsilon) \|\mathbf{x}\|_2^2 \leq \|\mathbf{P}\mathbf{x}\|_2^2 \leq (1 + \epsilon) \|\mathbf{x}\|_2^2$$

$\log(1/\delta)$
 $\frac{1}{\epsilon^2}$
 \approx



There is a truly brilliant method that runs in $O(n \log n)$ time.

Preview: Will involve Fast Fourier Transform in disguise.