# CS-GY 6763: Lecture 13
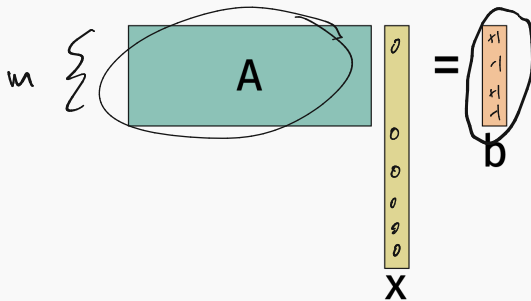# Finish Sparse Recovery and Compressed Sensing, Introduction to Spectral Sparsification

NYU Tandon School of Engineering, Prof. Christopher Musco

## SPARSE RECOVERY/COMPRESSED SENSING PROBLEM SETUP

- Design a matrix $A \in \mathbb{R}^{m \times n}$ with $m < n$, $b \in \mathbb{R}^m$.
- "Measure" $b = Ax$ for some <u>$k$-sparse</u> $x \in \mathbb{R}^n$.



$$m = O(k \log n)$$

- Recover $x$ from $b$.

Sample complexity: Can achieve $m = O(k \log n)$ or similar.

- Usually corresponds to some application-dependent cost (eg. length of time to acquire MRI, space complexity for heavy hitters problem)

Computational complexity: Naive methods take $O(n k)$ time to recover $k$-sparse x from b.

Typically design **A** with as few rows as possible that fulfills some desired property.

- **A** has Kruskal rank $r$. All sets of $r$ columns in **A** are linearly independent.
  - Recover vectors **x** with sparsity $k = r/2$.
- **A** is $\mu$-incoherent. $|A_i^T A_j| \leq \mu \|A_i\|_2 \|A_j\|_2$ for all columns $A_i, A_j, i \neq j$.
  - Recover vectors **x** with sparsity $k = 1/\mu$.

**A** obeys the $(q, \epsilon)$-Restricted Isometry Property.
  - Recover vectors **x** with sparsity $k = O(q)$.

Definition ($(q, \epsilon)$-Restricted Isometry Property)

A matrix **A** satisfies $(q, \epsilon)$-RIP if, for all **x** with $\|x\|_0 \leq q$,

$$(1 - \epsilon)\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + \epsilon)\|x\|_2^2.$$

Argued this holds for random matrices (JL matrices) and subsampled Fourier matrices with roughly $m = O\left(\frac{k \log n}{\epsilon^2}\right)$ rows.

### Theorem ($\ell_0$-minimization)

*Suppose we are given $A \in \mathbb{R}^{m \times n}$ and $b = Ax$ for an unknown k-sparse $x \in \mathbb{R}^n$. If A is $(2k, \epsilon)$-RIP for any $\epsilon < 1$ then x is the unique minimizer of:*

$$\min \|z\|_0 \qquad \text{subject to} \qquad Az = b.$$

- Establishes that information theoretically we can recover x in $O(n^k)$ time from $O(k \log n)$ measurements.

Proof by contradiction:

$$\underline{Ay = Ax} = b \qquad \text{but} \qquad \|y\|_0 \le \|x\|_0 = k$$

$$y - x = \Delta \quad 0 = \|Ax - Ay\| = \|A\Delta\| \ge (1-\epsilon)\|\Delta\| \ne 0$$

Convex relaxation of the $\ell_0$ minimization problem:

**Problem (Basis Pursuit, i.e. $\ell_1$ minimization.)**

$$\min_z \|z\|_1 \qquad \textit{subject to} \qquad Az = b.$$

- Objective is convex.

$$O(n^4)$$

- Optimizing over convex set.

$$\min \|Az - b\|_2^2$$
$$+ \lambda \|z\|_1$$

7

### Theorem

$\mathcal{U}$

$\frac{3}{2} \cdot \frac{1}{.17}$

*If* $A$ *is* $(3k, \epsilon)$*-RIP for* $\epsilon \leq .17$ *and* $\|x\|_0 = k$*, then* x *is the unique optimal solution of the Basis Pursuit LP).*
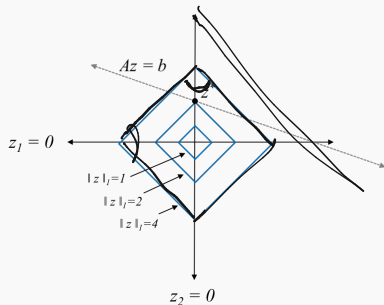
### Two surprising things about this result:

- Exponentially improve computational complexity with only a constant factor overhead in measurement complexity.
- Typical "relax-and-round" algorithm, but rounding is not even necessary! Just return the solution of the relaxed problem.
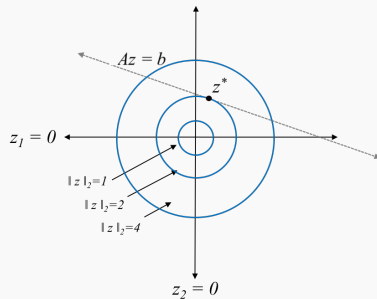
8

Suppose $A$ is $2 \times 1$, so $b$ is just a scalar and $x$ is a 2-dimensional vector.



Vertices of level sets of $\ell_1$ norm correspond to sparse solutions.

This is not the case e.g. for the $\ell_2$ norm.

## Theorem

*If A is (3k, $\epsilon$)-RIP for $\epsilon < .17$ and $\|x\|_0 = k$, then x is the unique optimal solution of the Basis Pursuit LP).*

Similar proof to $\ell_0$ minimization:

- By way of contradiction, assume x is not the optimal solution. Then there exists some non-zero $\Delta$ such that:
  - $\|x + \Delta\|_1 \leq \|x\|_1$
  - $A(x + \Delta) = Ax$. I.e $A\Delta = 0$.

Difference is that we can no longer assume that $\Delta$ is sparse.
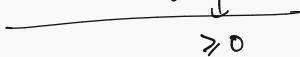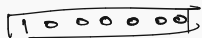
We will argue that $\Delta$ is approximately sparse.

10

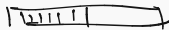First tool: $\|\omega\|_1 = s^T \omega$ where $s : \text{sign}(\omega)$

$\qquad \qquad \qquad \hookrightarrow \underline{s^T \omega} \leq \underline{\|s\|_2 \|\omega\|_2}$

For any $q$-sparse vector $\underline{\mathbf{w}}$, $\quad \|\mathbf{w}\|_2 \leq \underline{\underline{\|\mathbf{w}\|_1}} \leq \sqrt{q}\|\mathbf{w}\|_2$

$$\left(\sum_{i=1}^n |\omega_i|\right)^2 = \sum_{i=1}^n |\omega_i|^2 + \underbrace{\sum_{i \neq j} |\omega_i| |\omega_j|}_{\downarrow}$$

$\qquad \qquad \qquad \qquad \qquad \qquad \geq 0$

$\boxed{1 \ 0 \ 0 \ 0 \ 0 \ 0}$

Second tool: $\boxed{1 \ 2 \ 1 \ 1 \ 1 \ }$

For any norm and vectors $\mathbf{a}, \mathbf{b}$, $\quad \underline{\underline{\|\mathbf{a} + \mathbf{b}\| \geq \|\mathbf{a}\| - \|\mathbf{b}\|}}$

$$\|a\| = \|a + b - b\| \leq \|a + b\| + \|-b\|$$
$$= \|a + b\| + \|b\|$$

11

Some definitions:



$j \geq 2$

$k \left\{ \phantom{x} \right.$ S

$n-k$ $\bar{\bar{S}}$

$T_1$

$T_2$

$2k$

$\sqrt{n-k}$

$\frac{n-k}{2k} \cdot \sqrt{2k}$

$T_{(n-k)/2k}$

$\Delta$

X

Claim 1: $\|\underline{\underline{\Delta_S}}\|_1 \geq \|\underline{\underline{\Delta_{\bar{S}}}}\|_1$

$$\geq \|x_S\|_1 - \|\Delta_S\|_1$$

$$\boxed{\|x + \Delta\|_1} \leq \|x\|_1$$

$$\|x_S + \Delta_S\|_1 + \|\Delta_{\bar{S}}\|_1 \leq \|x\|_1$$

$$\sum_{i \in S} x_i + \Delta_i + \sum_{i \notin S} \Delta_i + \cancel{x_i} = 0$$

$$= \|x_S + \Delta_S\|_1 + \|\Delta_{\bar{S}}\|_1$$

$$\|\cancel{x_S}\|_1 - \|\Delta_S\|_1 + \|\Delta_{\bar{S}}\|_1 \leq \|\cancel{x_S}\|_1$$

Claim 2: $\|\Delta_S\|_2 \geq \sqrt{2} \sum_{j>2} \|\Delta_{T_j}\|_2$: $\quad \geq \frac{1}{\sqrt{2}} \| \Delta_S \|_2$

$$\|\Delta_S\|_2 \geq \frac{1}{\sqrt{k}} \|\Delta_S\|_1 \geq \frac{1}{\sqrt{k}} \|\Delta_{\bar{S}}\|_1 = \frac{1}{\sqrt{k}} \sum_{j\geq 1}^{m} \|\Delta_{T_j}\|_1.$$

Claim: $\|\Delta_{T_j}\|_1 \geq \sqrt{2k} \boxed{\|\Delta_{T_{j+1}}\|_2}$

$\quad \geq \frac{1}{\sqrt{k}} \sum_{j\geq 1}^{m-1} \|\Delta_{T_j}\|_2$

$\ell := \min\left(|\Delta_{T_j}|\right)$

$u = \max\left(|\Delta_{T_{j+1}}|\right)$

$\boxed{\|\Delta_{T_j}\|_1} \geq 2k \cdot \ell$

$\sqrt{2k} \|T_{j+1}\|_2 \leq \sqrt{2k \cdot u^2} = \sqrt{2k} \cdot u$

$\leq 2k \cdot \ell$

14

**Finish up proof by contradiction:** Recall that $A$ is assumed to have the $(3h, \epsilon)$ RIP property.

$$(1-\epsilon) - (1+\epsilon)\frac{1}{\sqrt{2}} = 0$$

$$0 = \|A\Delta\|_2 \geq \|A\Delta_{S\cup T_1}\|_2 - \sum_{j\geq 2}\|A\Delta_{T_j}\|_2$$

$$1 - \frac{1}{\sqrt{2}} = \epsilon - \frac{\epsilon}{\sqrt{2}}$$

$$\geq \|A\Delta_{S\cup T_1}\|_2 - \frac{1}{\sqrt{2}}\|A\Delta_S\|_2$$

$$\geq (1-\epsilon)\|\Delta_{S\cup T_1}\|_2 - (1+\epsilon)\frac{1}{\sqrt{2}}\|\Delta_S\|_2$$

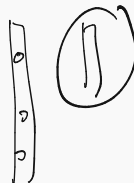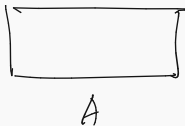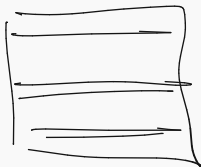$$\geq (1-\epsilon)\|\Delta_S\|_2 - (1+\epsilon)\frac{1}{\sqrt{2}}\|\Delta_S\|_2$$

$$= \left((1-\epsilon) - (1+\epsilon)\frac{1}{\sqrt{2}}\right)\|\Delta_S\|_2$$

15

A lot of interest in developing even faster algorithms that avoid using the "heavy hammer" of linear programming and run in even faster than $O(n^{3.5})$ time.
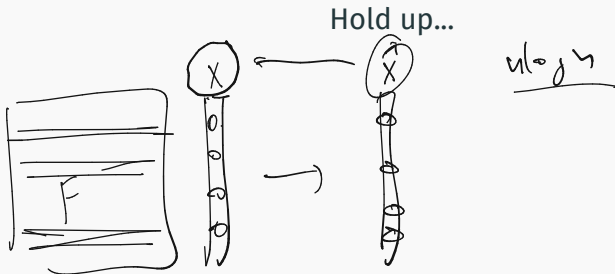
- **Iterative Hard Thresholding**: Looks a lot like projected gradient descent. Solve $\min_z \|Az - b\|$ with gradient descent while continually projecting $z$ back to the set of $k$-sparse vectors. Runs in time $\sim O(nk \log n)$ for Gaussian measurement matrices and $O(n \log n)$ for subsampled Fourer matrices.

- Other "first order" type methods: Orthogonal Matching Pursuit, CoSaMP, Subspace Pursuit, etc.

When **A** is a subsampled Fourier matrix, there are now methods that run in $O(k \log^c n)$ time [Hassanieh, Indyk, Kapralov, Katabi, Price, Shi, etc. 2012+].
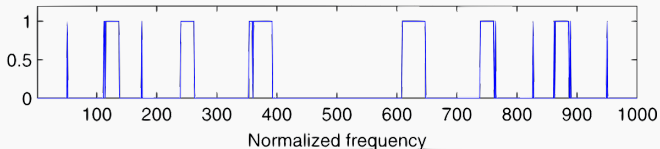
Hold up...

**Corollary:** When **x** is $k$-sparse, we can compute the inverse Fourier transform $\mathbf{F^*Fx}$ of $\mathbf{Fx}$ in $O(k \log^c n)$ time!

- Randomly subsample **Fx**.
- Feed that input into our sparse recovery algorithm to extract **x**.

Fourier and inverse Fourier transforms in <u>sublinear time</u> when the output is sparse.



Normalized frequency

**Applications in:** Wireless communications, GPS, protein imaging, radio astronomy, etc. etc.

A LITTLE ABOUT MY RESEARCH

### Theorem (Subspace Embedding)

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

$$(1 - \epsilon)\|Ax\|_2^2 \le \|\mathbf{\Pi}Ax\|_2^2 \le (1 + \epsilon)\|Ax\|_2^2$$

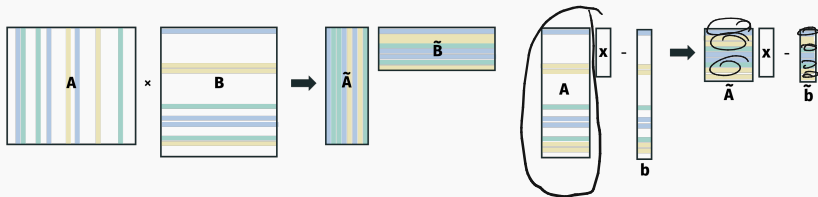*for all $x \in \mathbb{R}^d$, as long as $m = O\left(\frac{d + \log(1/\delta)}{\epsilon^2}\right)$.*

Implies regression result, and more.

**Example:** The any singular value $\tilde{\sigma}_i$ of $\mathbf{\Pi}A$ is a $(1 \pm \epsilon)$ approximation to the true singular value $\sigma_i$ of B.

**Recurring research interest:** Replace random projection methods with random sampling methods. Prove that for essentially all problems of interest, can obtain same asymptotic runtimes.
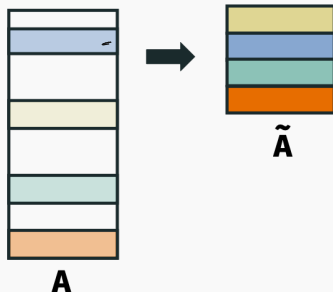


Sampling has the added benefit of preserving matrix sparsity or structure, and can be applied in a wider variety of settings where random projections are too expensive.

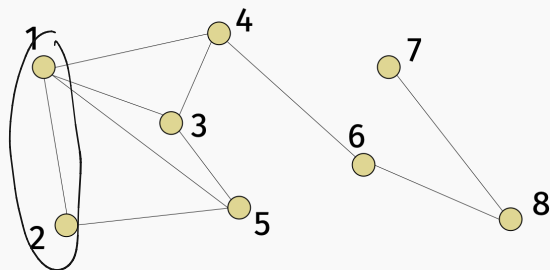Goal: Can we use sampling to obtain subspace embeddings? I.e. for a given $A$ find $\tilde{A}$ whose rows are a (weighted) subset of rows in $A$ and:

$O(d/\epsilon^2)$

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



$\tilde{A}$

$A$

## EXAMPLE WHERE STRUCTURE MATTERS

Let $B$ be the edge-vertex incidence matrix of a graph $G$ with vertex set $V$, $|V| = d$. Recall that $B^T B = L$.



Recall that if $x \in \{-1, 1\}^n$ is the cut indicator vector for a cut $S$ in the graph, then $\frac{1}{4}\|Bx\|_2^2 = \text{cut}(S, V \setminus S)$.

$$x = [1, 1, 1, -1, 1, -1, -1, -1]$$



$x \in \{-1, 1\}^d$ is the <u>cut indicator vector</u> for a cut $S$ in the graph, then $\frac{1}{4}\|Bx\|_2^2 = \text{cut}(S, V \setminus S)$

Extends to weighted graphs, as long as square root of weights is included in $\mathbf{B}$. Still have the $\mathbf{B}^T\mathbf{B} = \underline{\mathbf{L}}$.



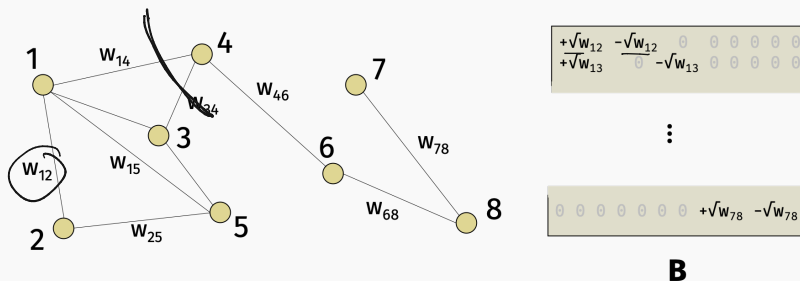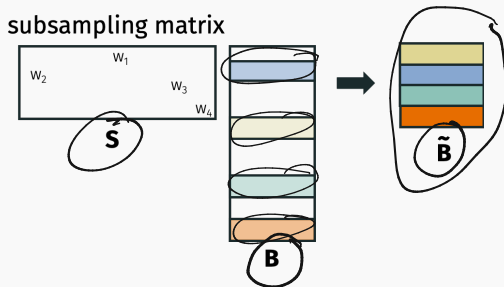And still have that if $\mathbf{x} \in \{-1, 1\}^d$ is the <u>cut indicator vector</u> for a cut $S$ in the graph, then $\frac{1}{4}\|\underline{\mathbf{Bx}}\|_2^2 = \text{cut}(S, V \setminus S)$.

**Goal:** Approximate $\mathbf{B}$ by a weighted subsample. I.e. by $\tilde{\mathbf{B}}$ with $m \ll |E|$ rows, each of which is a scaled copy of a row from $\mathbf{B}$.
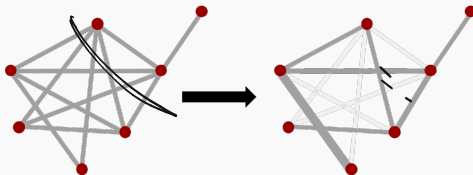


subsampling matrix

**Natural goal:** $\tilde{\mathbf{B}}$ is a subspace embedding for $\mathbf{B}$. In other words, $\tilde{\mathbf{B}}$ has $\approx O(d)$ rows and for all $\mathbf{x}$,

$$(1 - \epsilon)\|\mathbf{Bx}\|_2^2 \leq \|\tilde{\mathbf{Bx}}\|_2^2 \leq (1 + \epsilon)\|\mathbf{Bx}\|_2^2.$$

25

$\tilde{B}$ is itself an edge-vertex incidence matrix for some <u>sparser</u> graph $\tilde{G}$, which preserves many properties about $G$! $\tilde{G}$ is called a <u>spectral sparsifier</u> for $G$.



For example, we have that for any set $S$,

$$(1 - \epsilon)\,\mathsf{cut}_G(S, V \setminus S) \le \mathsf{cut}_{\tilde{G}}(S, V \setminus S) \le (1 + \epsilon)\,\mathsf{cut}_G(S, V \setminus S).$$

So $\tilde{G}$ can be used in place of $G$ in solving e.g. max/min cut problems, balanced cut problems, etc.

In contrast $\mathbf{\Pi B}$ would look nothing like an edge-vertex incidence matrix if $\mathbf{\Pi}$ is a JL matrix.

Spectral sparsifiers were introduced in 2004 by Spielman and Teng in an influential paper on faster algorithms for solving Laplacian linear systems.

$$O\left(\frac{d \log^L \delta}{u^2}\right)$$

- Generalize the cut sparsifiers of Benczur, Karger '96.
- Further developed in work by Spielman, Srivastava + Batson, '08.
- Have had huge influence in algorithms, and other areas of mathematics – this line of work lead to the 2013 resolution of the Kadison-Singer problem in functional analysis by Marcus, Spielman, Srivastava.

Rest of class: Learn about an important random sampling algorithm for constructing spectral sparsifiers, and subspace embeddings for matrices more generally.

**Goal:** Find $\tilde{\mathbf{A}}$ such that $\|\tilde{\mathbf{A}}x\|_2^2 = (1 \pm \epsilon)\|\mathbf{A}x\|_2^2$ for all x.

**Possible Approach:** Construct $\tilde{\mathbf{A}}$ by uniformly sampling rows from $\mathbf{A}$.



Can check that this approach fails even for the special case of a graph vertex-edge incidence matrix.

**Key idea:** <u>Importance sampling</u>. Select some rows with higher probability.

Suppose $A$ has $n$ rows $a_1 \ldots, a_n$. Let $p_1, \ldots, p_n \in [0, 1]$ be sampling probabilities. Construct $\tilde{A}$ as follows:

- For $i = 1, \ldots, n$
    - Select $a_i$ with probability $p_i$.
    - If $a_i$ is selected, add the scaled row $\frac{1}{\sqrt{p_i}} a_i$ to $\tilde{A}$.

Remember, ultimately want that $\|\tilde{A}x\|_2^2 = (1 \pm \epsilon)\|Ax\|_2^2$ for all x.

**Claim 1:** $\mathbb{E}[\|\tilde{A}x\|_2^2] = \|Ax\|_2^2$.

$$\sum_{i=1}^{n} \left( \frac{1}{\sqrt{p_i}} \, a_i^{\top} x \right)^2 \cdot \mathbb{1} \, [i \text{ was selected}]$$

**Claim 2:** Expected number of rows in $\tilde{A}$ is $\sum_{i=1}^{n} p_i$.

29

How should we choose the probabilities $p_1, \ldots, p_n$?

1. Introduce the idea of row leverage scores.
2. Motivate why these scores make for good sampling probabilities.
3. Prove that sampling with probabilities proportional to these scores yields a subspace embedding (or a spectral sparsifier) with a near optimal number of rows.

Let $a_1, \ldots, a_n$ be $A$'s rows. We define the statistical leverage score $\tau_i$ of row $a_i$ as:

$$\tau_i = a_i^T (A^T A)^{-1} a_i. \qquad \to \in [0, 1]$$

We will show that $\tau_i$ is a natural <u>importance measure</u> for each row in $A$.

We have that $\tau_i \in [0, 1]$ and $\sum_{i=1}^{n} \tau_i = d$ if $A$ has $d$ columns.

For $i = 1, \ldots, n$,

$$\tau_i = \mathbf{a}_i^T (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{a}_i.$$

### Theorem (Subspace Embedding from Subsampling)

*For each i, and fixed constant c, let $p_i = \min\left(1, \frac{c \log d}{\epsilon^2} \tau_i\right)$. Let $\tilde{\mathbf{A}}$ have rows sampled from $\mathbf{A}$ with probabilities $p_1, \ldots, p_n$. With probability 9/10,*

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

*and $\tilde{\mathbf{A}}$ has $O(d \log d / \epsilon^2)$ rows in expectation.*

$$\frac{d}{\epsilon^2}$$

32

How should we choose the probabilities $p_1, \ldots, p_n$?

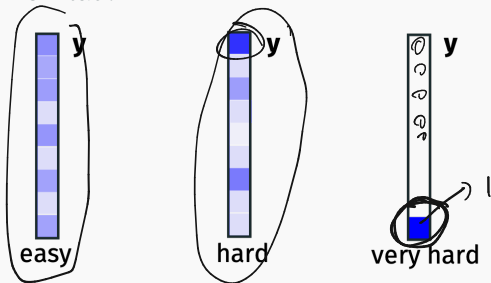As usual, consider a single vector x and understand how to sample to preserve norm of $y = Ax$:

$$\|\tilde{A}x\|_2^2 = \|SAx\|_2^2 = \|Sy\|_2^2 \approx \|y\|_2^2 = \|Ax\|_2^2.$$

Then we can union bound over an $\epsilon$-net to extend to all x.

As discussed a few lectures ago, uniform sampling only works well if $\mathbf{y} = \mathbf{Ax}$ is "flat".



Instead consider sampling with probabilities at least proportional to the magnitude of $\mathbf{y}$'s entries:

$$p_i > c \cdot \frac{y_i^2}{\|y\|_2^2} \text{ for constant } c \text{ to be determined.}$$

$$\sum_{i=1}^{} y_i^2$$

$\tilde{y} = (S)y$

Let $\tilde{y}$ be the subsampled $y$. Recall that, when sampling with probabilities $p_1, \ldots, p_n$, for $i = 1, \ldots, n$ we add $y_i$ to $\tilde{y}$ with probability $p_i$ and reweight by $\frac{1}{\sqrt{p_i}}$.

$$\|\tilde{y}\|_2^2 = \sum_{i=1}^n \frac{y_i^2}{p_i} \cdot Z_i \quad \text{where} \quad Z_i = \begin{cases} 1 \text{ with probability } p_i \\ 0 \text{ otherwise} \end{cases}$$

$$\mathsf{Var}[\|\tilde{y}\|_2^2] = \sum_{i=1}^n \frac{y_i^2}{p_i} \cdot \mathsf{Var}[Z_i] \leq \sum_{i=1}^n \frac{y_i^4}{p_i^2} \cdot p_i = \frac{y_i^4}{p_i}$$

We set $p_i = c \cdot \frac{y_i^2}{\|y\|_2^2}$ so get total variance:

$$\frac{1}{c}\|y\|_2^4$$

35

Using a Bernstein bound (or Chebyshev's inequality if you don't care about the $\delta$ dependence) we have that if $c = \frac{\log(1/\delta)}{\epsilon^2}$ then:

$$\Pr[\left|\|\tilde{\mathbf{y}}\|_2^2 - \|\mathbf{y}\|_2^2\right| \geq \epsilon\|\mathbf{y}\|_2^2] \leq \delta.$$

The number of samples we take in expectation is:

$$\sum_{i=1}^{n} p_i = \sum_{i=1}^{n} c \cdot \frac{y_i^2}{\|y_\ell\|_2^2} = \frac{\log(1/\delta)}{\epsilon^2}.$$

$$a_i^\top (A^\top A)^{-1} a_i$$

We don't know $y_1, \ldots, y_n$! And in fact, these values aren't fixed. We wanted to prove a bound for $\underline{y} = \underline{Ax}$ for any x.

**Idea behind leverage scores:** Sample row $i$ from A using the worst case (largest necessary) sampling probability:

$$\tau_i = \max_x \frac{y_i^2}{\|y\|_2^2} \qquad \text{where} \qquad \underline{y = Ax}.$$

If we sample with probability $p_i = \frac{1}{\epsilon^2} \cdot \tau_i$, then we will be sampling by at least $\frac{1}{\epsilon^2} \cdot \frac{y_i^2}{\|y\|_2^2}$, no matter what **y** is.

**Two concerns:**

1) How to compute $\tau_1, \ldots, \tau_n$?

2) the number of samples we take will be roughly $\sum_{i=1}^{n} \tau_i$. How do we bound this?

37

$$q_i^\top (A^\top A)^{-1} q_i$$

$$y_i = q_i^\top x$$

$$\tau_i = \max_x \frac{y_i^2}{\|y\|_2^2} \qquad \text{where} \qquad y = Ax.$$

$$= \max_x \frac{(q_i^\top x)^2}{\|Ax\|_2^2} = \max_x \frac{\left(q_i^\top (A^\top A)^{-1/2}(A^\top A)^{1/2} x\right)}{\|Ax\|_2^2}$$

$$\leq q_i^\top (A^\top A)^{-1/2}(A^\top A)^{-1/2} q_i \cdot \frac{x^\top (A^\top A)^{1/2}(A^\top A)^{1/2} x}{x^\top A^\top A x}$$

$$x = (A^\top A)^{-1} q_i$$

$$\leq q_i^\top (A^\top A)^{-1} q_i \leq \max_x \frac{y_i^2}{\|y\|_2^2}$$

Recall Cauchy-Schwarz inequality: $(w^\top z)^2 \leq w^\top w \cdot z^\top z$

Leverage score sampling:

- For $i = 1, \ldots, n,$
  - Compute $\tau_i = \mathbf{a}_i^T(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{a}_i$.
  - Set $p_i = \frac{c \log(1/\delta)}{\epsilon^2} \cdot \tau_i$.
  - Add row $\mathbf{a}_i$ to $\tilde{\mathbf{A}}$ with probability $p_i$ and reweight by $\frac{1}{\sqrt{p_i}}$.

For any fixed $\mathbf{x}$, we will have that
$(1 - \epsilon)\|\mathbf{Ax}\|_2^2 \leq \|\tilde{\mathbf{A}}\mathbf{x}\|_2^2 \leq (1 + \epsilon)\|\mathbf{Ax}\|_2^2$ with probability $(1 - \delta)$.

How many rows do we sample in expectation?

$$c \frac{\log(1/\delta)}{\epsilon^2} \cdot \sum_{i=1}^{n} \tau_i \quad = \sum_{i=1}^{n} a_i^T (A^TA)^{-1} a_i$$

**Claim:** No matter how large $n$ is, $\sum_{i=1}^{n} \tau_i = d$ a matrix $\mathbf{A} \in \mathbb{R}^d$.

$$\sum_{i=1}^{n} q_i^\top (A^\top A)^{-1} q_i \;=\; \text{tr}\left(A (A^\top A)^{-1} A^\top\right)$$

$$= \text{tr}\left((A^\top A)^{-1} A^\top A\right)$$

$$= \text{tr}\left(I_{d \times d}\right) \;=\; d$$

"Zero-sum" law for the importance of matrix rows.

Leverage score sampling:

$\|x = y$

- For $i = 1, \ldots, n$,
    - Compute $\tau_i = a_i^T (A^T A)^{-1} a_i$.
    - Set $p_i = \frac{c \log(1/\delta)}{\epsilon^2} \cdot \tau_i$.
    - Add row $a_i$ to $\tilde{A}$ with probability $p_i$ and reweight by $\frac{1}{\sqrt{p_i}}$.

For any fixed x, we will have that
$(1 - \epsilon)\|Ax\|_2^2 \le \|\tilde{A}x\|_2^2 \le (1 + \epsilon)\|Ax\|_2^2$ with high probability.

And since $\sum_{i=1}^n p_i = \frac{c \log(1/\delta)}{\epsilon^2} \cdot \sum_{i=1}^n \tau_i$, $\tilde{A}$ contains $O\left(\frac{d \log(1/\delta)}{\epsilon^2}\right)$ rows in expectation.

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

Last step: need to extend to all x.

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

## MAIN RESULT

Naive $\epsilon$-net argument leads to $d^2$ dependence since we need to set $\delta = c^d$. Getting the right $d \log d$ dependence below requires a "matrix Chernoff bound" (see e.g. Tropp 2015).

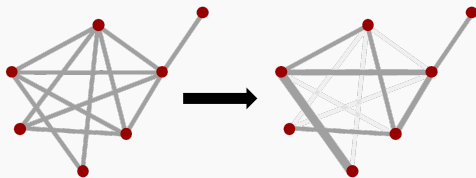### Theorem (Subspace Embedding from Subsampling)

*For each i, and fixed constant c, let* $p_i = \min\left(1, \frac{c \log d}{\epsilon^2} \cdot \tau_i\right)$. *Let* $\tilde{A}$ *have rows sampled from* A *with probabilities* $p_1, \ldots, p_n$. *With probability* 9/10,

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

*and* $\tilde{A}$ *has* $O(d \log d/\epsilon^2)$ *rows in expectation.*

For any graph $G$ with $d$ nodes, there exists a graph $\tilde{G}$ with $O(d\ \cancel{\log d}/\epsilon^2)$ edges such that, for all $\mathbf{x}$, $\|\tilde{\mathbf{B}}\mathbf{x}\|_2^2 = (1 \pm \epsilon)\|\mathbf{B}\mathbf{x}\|_2^2$.
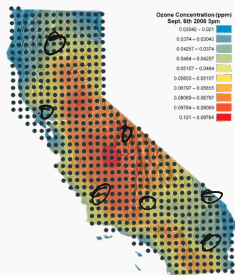


As a result, the value of any cut in $\tilde{G}$ is within a $(1 \pm \epsilon)$ factor of the value in $G$, the Laplacian eigenvalues are with a $(1 \pm \epsilon)$ factors, etc.

In many applications, computational costs are second order to data collection costs. We have a huge range of possible data points $a_1, \ldots, a_n$ that we can collect labels/values $b_1, \ldots, b_n$ for. Goal is to learn $x$ such that:
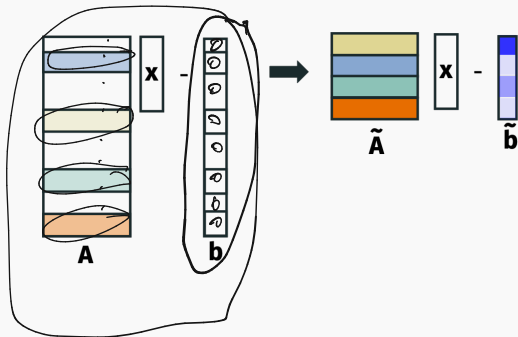
$$a_i^T x \approx b_i.$$

Want to do so after observing as few $b_1, \ldots, b_n$ as possible. Applications include healthcare, environmental science, etc.

Can be solved via **random sampling** for linear models.

## ANOTHER APPLICATION: ACTIVE REGRESSION

Claim: Let $\tilde{A}$ is an $O(1)$-factor subspace embedding for $A$ (obtained via leverage score sampling). Then $\tilde{x} = \arg\min \|\tilde{A}x - \tilde{b}\|_2^2$ satisfies:

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

where $x^* = \arg\min \|Ax - b\|_2^2$. Computing $\tilde{x}$ only requires collecting $O(d \log d)$ labels (independent of $n$).

Lots of applications:

- Robust bandlimited and multiband interpolation [STOC 2019].
- Active learning for Gaussian process regression [NeurIPS 2020].
- Active learning beyond the $\ell_2$ norm [Preprint 2021]
- Active learning for polynomial regression [Preprint 2021]
- DOE Grant on "learning based" algorithms for solving parametric partial differential equations.

**Problem**: Computing leverage scores $\tau_i = a_i^T (A^T A)^{-1} a_j$ is expensive.

$$O(ud^2)$$

**Main algorithmic idea:** Bootstrap leverage score sampling from uniform sampling (ITCS 2015).

**Problem**: Computing leverage scores $\tau_i = \mathbf{a}_i^T (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{a}_i$ is expensive.
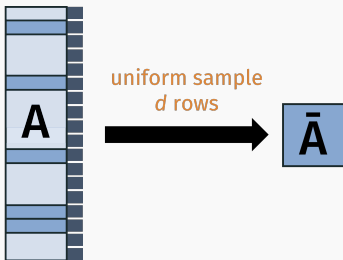


uniform sample
$d$ rows

**Main algorithmic idea:** Bootstrap <u>leverage score</u> sampling from <u>uniform sampling</u> (ITCS 2015).
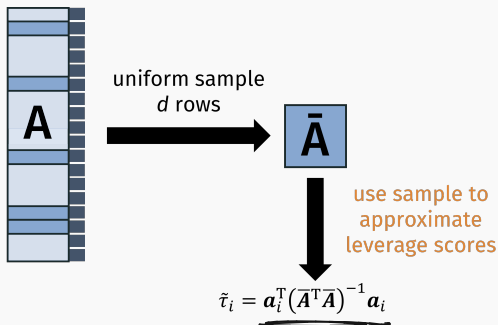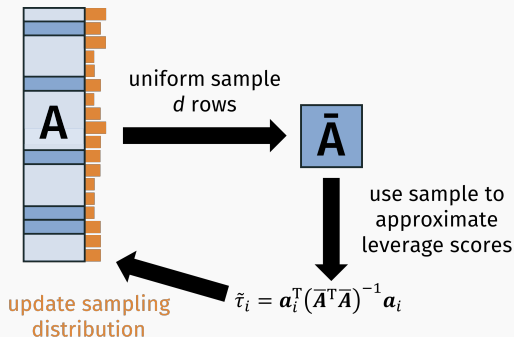
**Problem**: Computing leverage scores $\tau_i = \mathbf{a}_i^T(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{a}_i$ is expensive.



uniform sample
$d$ rows

use sample to
approximate
leverage scores

$$\tilde{\tau}_i = \boldsymbol{a}_i^{\mathrm{T}}\left(\overline{\boldsymbol{A}}^{\mathrm{T}}\overline{\boldsymbol{A}}\right)^{-1}\boldsymbol{a}_i$$

**Main algorithmic idea:** Bootstrap leverage score sampling from uniform sampling (ITCS 2015).

**Problem**: Computing leverage scores $\tau_i = a_i^T (A^T A)^{-1} a_i$ is expensive.



uniform sample
$d$ rows

$\bar{A}$

use sample to
approximate
leverage scores

update sampling
distribution

$\tilde{\tau}_i = a_i^T (\bar{A}^T \bar{A})^{-1} a_i$

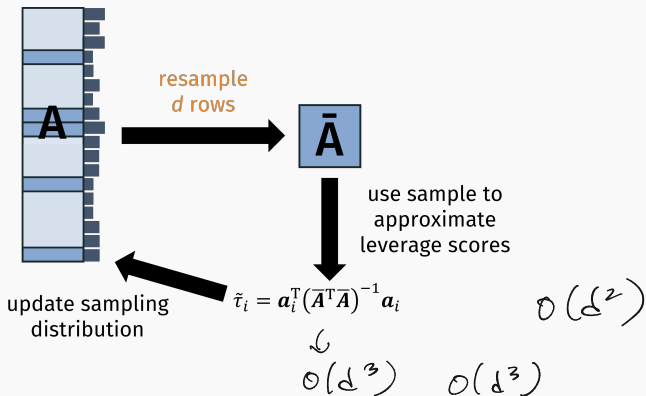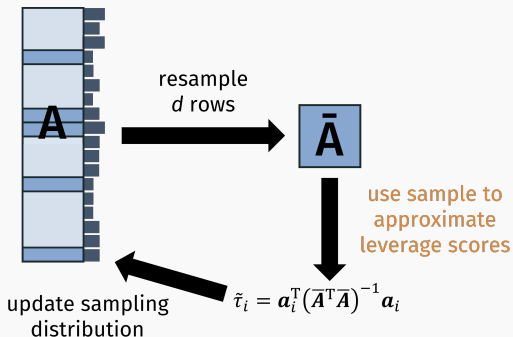**Main algorithmic idea:** Bootstrap leverage score sampling from uniform sampling (ITCS 2015).

47

**Problem**: Computing leverage scores $\tau_i = \mathbf{a}_i^T(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{a}_i$ is expensive.



**Main algorithmic idea:** Bootstrap <u>leverage score</u> sampling from <u>uniform sampling</u> (ITCS 2015).

47

**Problem**: Computing leverage scores $\tau_i = \mathbf{a}_i^T(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{a}_i$ is expensive.



resample
$d$ rows

$\bar{\mathbf{A}}$

use sample to
approximate
leverage scores

update sampling
distribution

$\tilde{\tau}_i = \boldsymbol{a}_i^{\mathrm{T}}(\bar{\boldsymbol{A}}^{\mathrm{T}}\bar{\boldsymbol{A}})^{-1}\boldsymbol{a}_i$

**Main algorithmic idea:** Bootstrap leverage score sampling from uniform sampling (ITCS 2015).
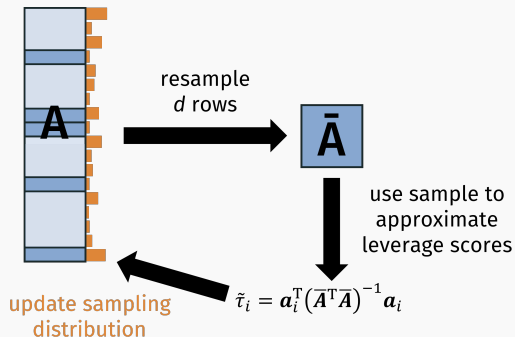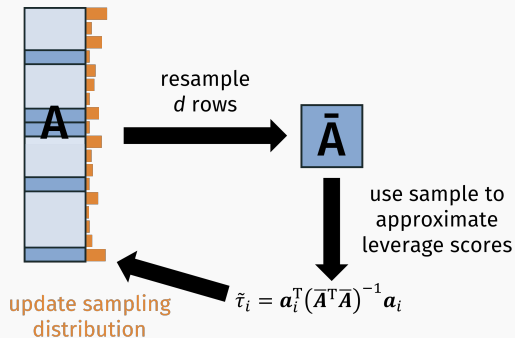
47

**Problem**: Computing leverage scores $\tau_i = \mathbf{a}_i^T (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{a}_i$ is expensive.



**Main algorithmic idea:** Bootstrap leverage score sampling from uniform sampling (ITCS 2015).

47

**Problem**: Computing leverage scores $\tau_i = \mathbf{a}_i^T(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{a}_i$ is expensive.



resample
$d$ rows

use sample to
approximate
leverage scores

$$\tilde{\tau}_i = \boldsymbol{a}_i^{\mathrm{T}}(\overline{\boldsymbol{A}}^{\mathrm{T}}\overline{\boldsymbol{A}})^{-1}\boldsymbol{a}_i$$

update sampling
distribution

After $O(\log n)$ rounds, $\tilde{\tau}_i \approx \tau_i$ for all $i$.
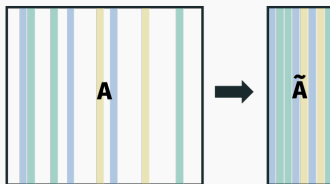
**Main algorithmic idea:** Bootstrap leverage score sampling from uniform sampling (ITCS 2015).

**Problem**: Sometimes we want to compress down to $\ll d$ rows or columns. E.g. we don't need a full subspace embedding, but just want to find a near optimal rank $k$ approximation.

**Approach:** Use "regularized" version of the leverage scores:

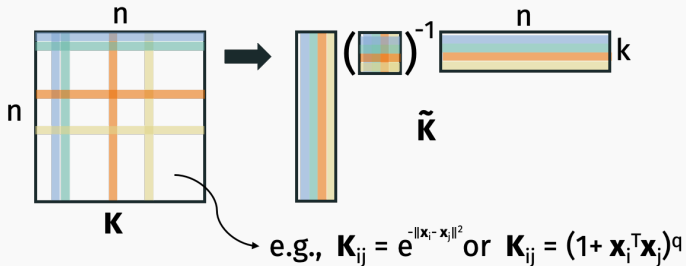$$\bar{\tau}_i = \mathsf{a}_i^T (\mathsf{A}^T \mathsf{A} + \lambda \mathsf{I})^{-1} \mathsf{a}_i$$



**Result:** Sample $O(k \log k / \epsilon)$ columns whose span contains a near-optimal low-approximation to **A** (SODA 2017).

The first $O(nk^2/\epsilon^2)$ time algorithm[1] for near optimal rank-$k$ approximation of any $n \times n$ positive semidefinite kernel matrix:
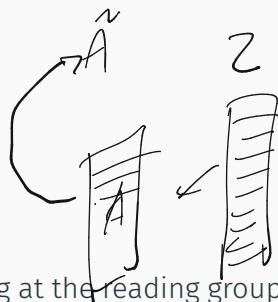


e.g., $K_{ij} = e^{-\|x_i - x_j\|^2}$ or $K_{ij} = (1 + x_i^T x_j)^q$

Based on the classic Nyström method. Importantly, does not even require constructing **K** explicitly, which takes $O(n^2)$ time.

[1]NeurIPS 2017.

Highlights of the semester for me:

- Very active office hours!
- Large number of students presenting at the reading group. Got to learn about a lot of your reseach interests.
- Lots of collaboration between students.

$$\underset{y = Ax}{\max} \frac{(\underline{\quad} y_i)^2}{\|y\|_2^2}$$

$$f : \{1 \dots n\} \to \mathbb{R} \qquad F : \{\underline{Ax} : \lambda \in \mathbb{R}^N\}$$

50