

# Subspace embeddings

Daniel Hsu

COMS 4772

1

# Supremum of simple stochastic processes

2

## Recap: JL lemma

**JL lemma.** For any  $\varepsilon \in (0, 1/2)$ , point set  $S \subset \mathbb{R}^d$  of cardinality  $|S| = n$ , and  $k \in \mathbb{N}$  such that  $k \geq \frac{16 \ln n}{\varepsilon^2}$ , there exists a linear map  $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$  such that

$$(1-\varepsilon)\|\mathbf{x}-\mathbf{y}\|_2^2 \leq \|f(\mathbf{x})-f(\mathbf{y})\|_2^2 \leq (1+\varepsilon)\|\mathbf{x}-\mathbf{y}\|_2^2 \quad \text{for all } \mathbf{x}, \mathbf{y} \in S.$$

## Main probabilistic lemma

$\exists$  random linear map  $\mathbf{M}: \mathbb{R}^d \rightarrow \mathbb{R}^k$  such that, for any  $\mathbf{u} \in S^{d-1}$ ,

$$\mathbb{P}\left(\left|\|\mathbf{M}\mathbf{u}\|_2^2 - 1\right| > \varepsilon\right) \leq 2 \exp(-\Omega(k\varepsilon^2)).$$

JL lemma is consequence of main probabilistic lemma as applied to collection  $T \subset S^{d-1}$  of  $|T| = \binom{n}{2}$  unit vectors (+ union bound):

$$\mathbb{P}\left(\max_{\mathbf{u} \in T} \left|\|\mathbf{M}\mathbf{u}\|_2^2 - 1\right| > \varepsilon\right) \leq |T| \cdot 2 \exp(-\Omega(k\varepsilon^2)).$$

3

## Related question

For  $T \subseteq S^{d-1}$ , *expected* maximum deviation

$$\mathbb{E} \max_{\mathbf{u} \in T} \left|\|\mathbf{M}\mathbf{u}\|_2^2 - 1\right| \leq ?$$

## General questions

For arbitrary collection of zero-mean random variables  $\{X_t : t \in T\}$ :

$$\mathbb{E} \max_{t \in T} X_t \leq ?$$

$$\mathbb{E} \max_{t \in T} |X_t| \leq ?$$

4

## Finite collections

Let  $\{X_t : t \in T\}$  be a *finite* collection of  $v$ -subgaussian and mean-zero random variables. Then

$$\mathbb{E} \max_{t \in T} X_t \leq \sqrt{2v \ln |T|}.$$

- ▶ Doesn't assume independence of  $\{X_t : t \in T\}$ .
  - ▶ (Independent case is the worst.)
- ▶ Get bound on  $\mathbb{E} \max_{t \in T} |X_t|$  as corollary.
  - ▶ Apply result to collection

$$\{X_t : t \in T\} \cup \{-X_t : t \in T\}.$$

5

## Proof

Starting point is identity from two invertible operations ( $\lambda > 0$ ):

$$\mathbb{E} \max_{t \in T} X_t = \frac{1}{\lambda} \ln \exp\left(\mathbb{E} \max_{t \in T} \lambda X_t\right)$$

- ▶ Apply Jensen's inequality:

$$\leq \frac{1}{\lambda} \ln \mathbb{E} \exp\left(\max_{t \in T} \lambda X_t\right) = \frac{1}{\lambda} \ln \mathbb{E}\left(\max_{t \in T} \exp(\lambda X_t)\right)$$

- ▶ Bound max with sum, and use linearity of expectation:

$$\leq \frac{1}{\lambda} \ln \sum_{t \in T} \mathbb{E} \exp(\lambda X_t)$$

- ▶ Exploit  $v$ -subgaussian property:

$$\leq \frac{1}{\lambda} \ln \sum_{t \in T} \exp(v\lambda^2/2) = \frac{\ln |T|}{\lambda} + \frac{v\lambda}{2}$$

- ▶ Choose appropriate  $\lambda$  to conclude. □

6

## Alternative proof

**Integrate tail bound:** for any non-negative random variable  $Y$ ,

$$\mathbb{E}(Y) = \int_0^\infty \mathbb{P}(Y \geq y) dy.$$

For  $Y := \max_{t \in T} |X_t|$ , gives same result up to constants.

7

## Infinite collections

For *infinite* collection of zero-mean random variables  $\{X_t : t \in T\}$ :

$$\mathbb{E} \sup_{t \in T} X_t \leq ?$$

- ▶ In general, can go  $\rightarrow \infty$ .
- ▶ To bound, must exploit *correlations* among the  $X_t$ .
  - ▶ E.g., in  $\left\{ \left| \|\mathbf{M}\mathbf{u}\|_2^2 - 1 \right| : \mathbf{u} \in T \right\}$  for  $T \subseteq S^{d-1}$ , the random variables for  $\mathbf{u}$  and  $\mathbf{u} + \boldsymbol{\delta}$ , for small  $\boldsymbol{\delta}$ , are highly correlated.

8

## Convex hulls of linear functionals

Let  $T \subset \mathbb{R}^d$  be a finite set of vectors, and let  $\mathbf{X}$  be a random vector in  $\mathbb{R}^d$  such that  $\langle \mathbf{w}, \mathbf{X} \rangle$  is  $\nu$ -subgaussian for every  $\mathbf{w} \in T$ . Then

$$\mathbb{E} \max_{\tilde{\mathbf{w}} \in \text{conv}(T)} \langle \tilde{\mathbf{w}}, \mathbf{X} \rangle \leq \sqrt{2\nu \ln |T|}.$$

**Proof:**

- ▶ Write  $\tilde{\mathbf{w}} \in \text{conv}(T)$  as  $\tilde{\mathbf{w}} = \sum_{\mathbf{w} \in T} p_{\mathbf{w}} \mathbf{w}$  for some  $p_{\mathbf{w}} \geq 0$  that sum to one.
- ▶ Observe that

$$\langle \tilde{\mathbf{w}}, \mathbf{x} \rangle = \sum_{\mathbf{w} \in T} p_{\mathbf{w}} \langle \mathbf{w}, \mathbf{x} \rangle \leq \max_{\mathbf{w} \in T} \langle \mathbf{w}, \mathbf{x} \rangle.$$

- ▶ So max over  $\tilde{\mathbf{w}} \in \text{conv}(T)$  is at most max over  $\mathbf{w} \in T$ .
- ▶ Conclude by applying previous result for finite collections.  $\square$

9

## Euclidean norm

Let  $\mathbf{X}$  be a random vector such that  $\langle \mathbf{u}, \mathbf{X} \rangle$  is  $\nu$ -subgaussian for every  $\mathbf{u} \in S^{d-1}$ . Then

$$\mathbb{E} \|\mathbf{X}\|_2 = \mathbb{E} \max_{\mathbf{u} \in S^{d-1}} \langle \mathbf{u}, \mathbf{X} \rangle \leq 2\sqrt{2\nu \ln 5^d} = O(\sqrt{\nu d}).$$

**Key step of proof:**

- ▶ For any  $\varepsilon > 0$ , there is a finite subset  $\mathcal{N} \subset S^{d-1}$  of cardinality  $|\mathcal{N}| \leq (1 + 2/\varepsilon)^d$  such that, for every  $\mathbf{u} \in S^{d-1}$ , there exists  $\mathbf{u}_0 \in \mathcal{N}$  with

$$\|\mathbf{u} - \mathbf{u}_0\|_2 \leq \varepsilon.$$

- ▶ Such a set  $\mathcal{N}$  is called an  $\varepsilon$ -net for  $S^{d-1}$ .
- ▶ We need a  $1/2$ -net, of cardinality at most  $5^d$ .

10

## Proof

- ▶ Write  $\mathbf{u} \in S^{d-1}$  as

$$\mathbf{u} = \mathbf{u}_0 + \delta \mathbf{q},$$

where  $\mathbf{u}_0 \in \mathcal{N}$ ,  $\mathbf{q} \in S^{d-1}$ ,  $\delta \in [0, 1/2]$ , so

$$\langle \mathbf{u}, \mathbf{X} \rangle = \langle \mathbf{u}_0, \mathbf{X} \rangle + \delta \langle \mathbf{q}, \mathbf{X} \rangle.$$

- ▶ Observe that

$$\begin{aligned} \max_{\mathbf{u} \in S^{d-1}} \langle \mathbf{u}, \mathbf{X} \rangle &\leq \max_{\mathbf{u}_0 \in \mathcal{N}} \langle \mathbf{u}_0, \mathbf{X} \rangle + \max_{\delta \in [0, 1/2]} \max_{\mathbf{q} \in S^{d-1}} \delta \langle \mathbf{q}, \mathbf{X} \rangle \\ &\leq \max_{\mathbf{u}_0 \in \mathcal{N}} \langle \mathbf{u}_0, \mathbf{X} \rangle + \frac{1}{2} \max_{\mathbf{q} \in S^{d-1}} \langle \mathbf{q}, \mathbf{X} \rangle. \end{aligned}$$

- ▶ So max over  $S^{d-1}$  is at most twice max over  $\mathcal{N}$ .
- ▶ Conclude by applying previous result for finite collections.  $\square$

11

## $\varepsilon$ -nets for unit sphere

There is an  $\varepsilon$ -net for  $S^{d-1}$  of cardinality at most  $(1 + 2/\varepsilon)^d$ .

### Proof:

- ▶ Repeatedly select points from  $S^{d-1}$  so that each selected point has distance more than  $\varepsilon$  from all previously selected points.
- ▶ Equivalent: repeatedly select points from  $S^{d-1}$  as long as balls of radius  $\varepsilon/2$ , centered at selected points, are disjoint.
  - ▶ (Process must eventually stop.)
- ▶ When process stops, every  $\mathbf{u} \in S^{d-1}$  is at distance at most  $\varepsilon$  from selected points.
  - ▶ I.e., selected points form an  $\varepsilon$ -net for  $S^{d-1}$ .
- ▶ If select  $N$  points, then the  $N$  balls of radius  $\varepsilon/2$  are disjoint, and they are contained in a ball of radius  $1 + \varepsilon/2$ . So

$$N \text{vol}((\varepsilon/2)B^d) \leq \text{vol}((1 + \varepsilon/2)B^d).$$

- ▶ This implies  $N \leq (1 + 2/\varepsilon)^d$ .  $\square$

12

## Remarks

- ▶ All previous results also hold with random variables are  $(v, c)$ -subexponential (possibly with  $c > 0$ ), with a slightly different bound: e.g.,

$$\mathbb{E} \max_{t \in T} X_t \leq \max \left\{ \sqrt{2v \ln |T|}, 2c \ln |T| \right\}.$$

- ▶ Also easy to get probability tail bounds (rather than expectation bounds).

13

## Subspace embeddings

14

## Subspace JL lemma

Consider  $k \times d$  random matrix  $\mathbf{M}$  whose entries are iid  $N(0, 1/k)$ .  
For a  $W \subseteq \mathbb{R}^d$  be a subspace of dimension  $r$ ,

$$\mathbb{E} \max_{\mathbf{u} \in S^{d-1} \cap W} \left| \|\mathbf{M}\mathbf{u}\|_2^2 - 1 \right| \leq O\left(\sqrt{\frac{r}{k}} + \frac{r}{k}\right).$$

Bound is at most  $\varepsilon$  when  $k \geq O\left(\frac{r}{\varepsilon^2}\right)$ .

Implies existence of mapping  $\mathbf{M}: \mathbb{R}^d \rightarrow \mathbb{R}^k$  that approximately preserves all distances between points in  $W$ .

15

## Proof of subspace JL lemma

Let columns of  $\mathbf{Q}$  be ONB for  $W$ . Then

$$\begin{aligned} \max_{\mathbf{u} \in S^{d-1} \cap W} \left| \|\mathbf{M}\mathbf{u}\|_2^2 - 1 \right| &= \max_{\mathbf{u} \in S^{r-1}} \left| \mathbf{u}^\top \mathbf{Q}^\top (\mathbf{M}^\top \mathbf{M} - \mathbf{I}) \mathbf{Q} \mathbf{u} \right| \\ &= \max_{\mathbf{u}, \mathbf{v} \in S^{r-1}} \mathbf{u}^\top \mathbf{Q}^\top (\mathbf{M}^\top \mathbf{M} - \mathbf{I}) \mathbf{Q} \mathbf{v}. \end{aligned}$$

**Lemma.** For any  $\mathbf{u}, \mathbf{v} \in S^{r-1}$ ,

$$X_{\mathbf{u}, \mathbf{v}} := \mathbf{u}^\top \mathbf{Q}^\top (\mathbf{M}^\top \mathbf{M} - \mathbf{I}) \mathbf{Q} \mathbf{v}$$

is  $(O(1/k), O(1/k))$ -subexponential.

16



## Proof of subspace JL lemma (continued)

For  $\mathbf{u}, \mathbf{v} \in S^{r-1}$ ,  $X_{\mathbf{u}, \mathbf{v}} := \mathbf{u}^\top \mathbf{Q}^\top (\mathbf{M}^\top \mathbf{M} - \mathbf{I}) \mathbf{Q} \mathbf{v}$ .

Let  $\mathcal{N}$  be  $1/4$ -net for  $S^{r-1}$ .

- ▶ Write  $\mathbf{u}, \mathbf{v} \in S^{r-1}$  as

$$\mathbf{u} = \mathbf{u}_0 + \varepsilon \mathbf{p}, \quad \mathbf{v} = \mathbf{v}_0 + \delta \mathbf{q},$$

where  $\mathbf{u}_0, \mathbf{v}_0 \in \mathcal{N}$ ,  $\mathbf{p}, \mathbf{q} \in S^{r-1}$  and  $\varepsilon, \delta \in [0, 1/4]$ , so

$$X_{\mathbf{u}, \mathbf{v}} = X_{\mathbf{u}_0, \mathbf{v}_0} + \varepsilon X_{\mathbf{p}, \mathbf{v}} + \delta X_{\mathbf{u}_0, \mathbf{q}}.$$

- ▶ Therefore

$$\max_{\mathbf{u}, \mathbf{v} \in S^{r-1}} X_{\mathbf{u}, \mathbf{v}} \leq \max_{\mathbf{u}_0, \mathbf{v}_0 \in \mathcal{N}} X_{\mathbf{u}_0, \mathbf{v}_0} + \frac{1}{2} \max_{\mathbf{p}, \mathbf{q} \in S^{r-1}} X_{\mathbf{p}, \mathbf{q}},$$

which implies

$$\max_{\mathbf{u}, \mathbf{v} \in S^{r-1}} X_{\mathbf{u}, \mathbf{v}} \leq 2 \max_{\mathbf{u}_0, \mathbf{v}_0 \in \mathcal{N}} X_{\mathbf{u}_0, \mathbf{v}_0}.$$

- ▶ Conclude by applying previous result for finite collections.  $\square$

17

## Application to least squares

18

## Big data least squares

- ▶ **Input:** matrix  $\mathbf{A} \in \mathbb{R}^{n \times d}$ , vector  $\mathbf{b} \in \mathbb{R}^n$  ( $n \gg d$ ).
- ▶ **Goal:** find  $\mathbf{x} \in \mathbb{R}^d$  so as to (approx.) minimize  $\|\mathbf{Ax} - \mathbf{b}\|_2^2$ .
  
- ▶ Computation time:  $O(nd^2)$ .
- ▶ Can we speed this up?

19

## Simple approach

- ▶ Pick  $m \ll n$ .
- ▶ Let  $\mathbf{M}$  be random  $m \times n$  matrix (e.g., entries iid  $N(0, 1/m)$ , Fast JL Transform).
- ▶ Let  $\tilde{\mathbf{A}} := \mathbf{MA}$  and  $\tilde{\mathbf{b}} := \mathbf{Mb}$ .
- ▶ Obtain solution  $\hat{\mathbf{x}}$  to least squares problem on  $(\tilde{\mathbf{A}}, \tilde{\mathbf{b}})$ .

20

## Simple (somewhat loose) analysis

- ▶ Let  $W$  be subspace spanned by columns of  $\mathbf{A}$  and  $\mathbf{b}$ .
  - ▶ Dimension is at most  $d + 1$ .
- ▶ If  $m \geq O(d/\varepsilon^2)$ , then  $\mathbf{M}$  is subspace embedding for  $W$ :

$$(1 - \varepsilon)\|\mathbf{x}\|_2^2 \leq \|\mathbf{M}\mathbf{x}\|_2^2 \leq (1 + \varepsilon)\|\mathbf{x}\|_2^2 \quad \text{for all } \mathbf{x} \in W.$$

- ▶ Let  $\mathbf{x}_* := \arg \min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$ .

▶

$$\begin{aligned} \|\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}\|_2^2 &\leq \frac{1}{1 - \varepsilon} \|\mathbf{M}(\mathbf{A}\hat{\mathbf{x}} - \mathbf{b})\|_2^2 \\ &\leq \frac{1}{1 - \varepsilon} \|\mathbf{M}(\mathbf{A}\mathbf{x}_* - \mathbf{b})\|_2^2 \\ &\leq \frac{1 + \varepsilon}{1 - \varepsilon} \|\mathbf{A}\mathbf{x}_* - \mathbf{b}\|_2^2. \end{aligned}$$

- ▶ Running time (using FJLT):  $O((m + n)d \log n + md^2)$ . □

21

## Another perspective: random sampling

- ▶ Pick random sample of  $m \ll n$  of rows of  $(\mathbf{A}, \mathbf{b})$ ; obtain solution  $\hat{\mathbf{x}}$  for least squares problem on the sample.
- ▶ Hope  $\hat{\mathbf{x}}$  is also good for the original problem.
- ▶ In statistics, this is the *random design* setting for regression.
  - ▶ Random sample of covariates  $\tilde{\mathbf{A}} \in \mathbb{R}^{m \times d}$  and responses  $\tilde{\mathbf{b}} \in \mathbb{R}^m$  from full population  $(\mathbf{A}, \mathbf{b})$ .
  - ▶ Least squares solution  $\hat{\mathbf{x}}$  on  $(\tilde{\mathbf{A}}, \tilde{\mathbf{b}})$  is *MLE* for linear regression coefficients under linear model with Gaussian noise.
  - ▶ Can also regard  $\hat{\mathbf{x}}$  as *empirical risk minimizer* among all linear predictors under squared loss.

22

## Simple random design analysis

- ▶ Let  $\mathbf{x}_* := \arg \min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$ .
- ▶ With high probability over choice of random sample,

$$\|\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}\|_2^2 \leq \left(1 + O\left(\frac{\kappa}{m}\right)\right) \cdot \|\mathbf{A}\mathbf{x}_* - \mathbf{b}\|_2^2$$

(up to lower-order terms), where

$$\kappa := n \cdot \max_{i \in [n]} \|(\mathbf{A}^\top \mathbf{A})^{-1/2} \mathbf{A}^\top \mathbf{e}_i\|_2^2$$

and  $\mathbf{e}_i$  is  $i$ -th coordinate basis vector.

- ▶ Write thin SVD of  $\mathbf{A}$  as  $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ , where  $\mathbf{U} \in \mathbb{R}^{n \times d}$ . Then

$$(\mathbf{A}^\top \mathbf{A})^{-1/2} \mathbf{A}^\top = (\mathbf{V}\mathbf{S}^2\mathbf{V}^\top)^{-1/2} \mathbf{V}\mathbf{S}\mathbf{U}^\top = \mathbf{V}\mathbf{U}^\top.$$

- ▶ So  $\kappa = n \cdot \max_{i \in [n]} \|\mathbf{U}^\top \mathbf{e}_i\|_2^2$ .
  - ▶  $\|\mathbf{U}^\top \mathbf{e}_i\|_2^2$  is *statistical leverage score* for  $i$ -th row of  $\mathbf{A}$ : measures how much “influence”  $i$ -th row has on least squares solution.

23

## Statistical leverage

- ▶  $i$ -th *statistical leverage score*:  $\ell_i := \|\mathbf{U}^\top \mathbf{e}_i\|_2^2$ , where  $\mathbf{U} \in \mathbb{R}^{n \times d}$  is matrix of left singular vectors of  $\mathbf{A}$ .
- ▶ Two extreme cases:

$$\mathbf{U} = \begin{bmatrix} \mathbf{I}_{d \times d} \\ \mathbf{0}_{(n-d) \times d} \end{bmatrix} \Rightarrow n \cdot \max_{i \in [n]} \ell_i = n.$$

$$\mathbf{U} = \frac{1}{\sqrt{n}} \begin{bmatrix} \mathbf{H}_n \mathbf{e}_1 & \mathbf{H}_n \mathbf{e}_2 & \cdots & \mathbf{H}_n \mathbf{e}_d \end{bmatrix} \Rightarrow n \cdot \max_{i \in [n]} \ell_i = d,$$

where  $\mathbf{H}_n$  is  $n \times n$  Hadamard matrix.

- ▶ First case: first  $d$  rows are the only rows that matter.
- ▶ Second case: all  $n$  rows equally important.

24

## Ensuring small statistical leverage

- ▶ To ensure situation is more like second case, apply random rotation (e.g., randomized Hadamard transform) to  $\mathbf{A}$  and  $\mathbf{b}$ .
  - ▶ Randomly mixes up rows of  $(\mathbf{A}, \mathbf{b})$  so no single row is (much) more important than another.
  - ▶ Get  $n \cdot \max_{i \in [n]} \ell_i = O(d + \log n)$  with high probability.
- ▶ To get  $1 + \varepsilon$  approximation ratio, i.e.,

$$\|\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}\|_2^2 \leq (1 + \varepsilon) \cdot \|\mathbf{A}\mathbf{x}_* - \mathbf{b}\|_2^2,$$

suffices to have

$$m \geq O\left(\frac{d + \log n}{\varepsilon}\right).$$

25

Application to compressed sensing

26

## Under-determined least squares

- ▶ **Input:** matrix  $\mathbf{A} \in \mathbb{R}^{n \times d}$ , vector  $\mathbf{b} \in \mathbb{R}^n$  ( $n \ll d$ ).
- ▶ **Goal:** find *sparsest*  $\mathbf{x} \in \mathbb{R}^d$  so as to minimize  $\|\mathbf{Ax} - \mathbf{b}\|_2^2$ .
  
- ▶ NP-hard in general.
- ▶ Suppose  $\mathbf{b} = \mathbf{A}\bar{\mathbf{x}}$  for some  $\bar{\mathbf{x}} \in \mathbb{R}^d$  with  $\text{nnz}(\bar{\mathbf{x}}) \leq k$ .
  - ▶ I.e.,  $\bar{\mathbf{x}}$  is  $k$ -sparse.
  - ▶ Is  $\bar{\mathbf{x}}$  the (unique) sparsest solution?
  - ▶ If so, how to find it?

27

## Null space property

**Lemma.** Null space of  $\mathbf{A}$  does not contain any non-zero  $2k$ -sparse vectors  $\iff$  every  $k$ -sparse vector  $\bar{\mathbf{x}} \in \mathbb{R}^d$  is the unique solution to  $\mathbf{Ax} = \mathbf{A}\bar{\mathbf{x}}$ .

- ▶ **Proof.** ( $\implies$ ) Take any  $k$ -sparse vectors  $\mathbf{x}$  and  $\mathbf{y}$  with  $\mathbf{Ax} = \mathbf{Ay}$ .  
Want to show  $\mathbf{x} = \mathbf{y}$ .
  - ▶ Then  $\mathbf{x} - \mathbf{y}$  is  $2k$ -sparse, and  $\mathbf{A}(\mathbf{x} - \mathbf{y}) = \mathbf{0}$ .
  - ▶ By assumption, null space of  $\mathbf{A}$  does not contain any non-zero  $2k$ -sparse vectors.
  - ▶ So  $\mathbf{x} - \mathbf{y} = \mathbf{0}$ , i.e.,  $\mathbf{x} = \mathbf{y}$ .
- ▶ ( $\impliedby$ ) Take any  $2k$ -sparse vector  $\mathbf{z}$  in the null space of  $\mathbf{A}$ . Want to show  $\mathbf{z} = \mathbf{0}$ .
  - ▶ Write it as  $\mathbf{z} = \mathbf{x} - \mathbf{y}$  for some  $k$ -sparse vectors  $\mathbf{x}$  and  $\mathbf{y}$  with disjoint supports.
  - ▶ Then  $\mathbf{A}(\mathbf{x} - \mathbf{y}) = \mathbf{0}$ , and hence  $\mathbf{x} = \mathbf{y}$  by assumption.
  - ▶ But  $\mathbf{x}$  and  $\mathbf{y}$  have disjoint support, so it must be that  $\mathbf{x} = \mathbf{y} = \mathbf{0}$ , so  $\mathbf{z} = \mathbf{0}$ . □

28

## Null space property from subspace embeddings

If  $\mathbf{A}$  is  $n \times d$  random matrix with iid  $N(0, 1)$  entries, then under what conditions is there no non-zero  $2k$ -sparse vector in its null space?

- ▶ Want: for any  $2k$ -sparse vector  $\mathbf{z}$ ,  $\mathbf{A}\mathbf{z} \neq \mathbf{0}$ , i.e.,  $\|\mathbf{A}\mathbf{z}\|_2^2 > 0$ .
- ▶ Consider a particular choice  $\mathcal{I} \subseteq [d]$  of  $|\mathcal{I}| = 2k$  coordinates, and the corresponding subspace  $W_{\mathcal{I}}$  spanned by  $\{\mathbf{e}_i : i \in \mathcal{I}\}$ .
  - ▶ Every  $2k$ -sparse  $\mathbf{z}$  is in  $W_{\mathcal{I}}$  for some  $\mathcal{I}$ .
- ▶ Sufficient for  $\mathbf{A}$  to be  $1/2$ -subspace embedding for  $W_{\mathcal{I}}$  for all  $\mathcal{I}$ :

$$\frac{1}{2}\|\mathbf{z}\|_2^2 \leq \|\mathbf{A}\mathbf{z}\|_2^2 \leq \frac{3}{2}\|\mathbf{z}\|_2^2 \quad \text{for all } 2k\text{-sparse } \mathbf{z}.$$

29

## Null space property from subspace embeddings (continued)

- ▶ Say  $\mathbf{A}$  fails for  $\mathcal{I}$  if it is not a  $1/2$ -subspace embedding for  $W_{\mathcal{I}}$ .
- ▶ Subspace JL lemma:

$$\mathbb{P}(\mathbf{A} \text{ fails for } \mathcal{I}) \leq 2^{O(k)} \exp(-\Omega(n)).$$

- ▶ Union bound over all choices of  $\mathcal{I}$  with  $|\mathcal{I}| = 2k$ :

$$\mathbb{P}(\mathbf{A} \text{ fails for some } \mathcal{I}) \leq \binom{d}{2k} 2^{O(k)} \exp(-\Omega(n)).$$

- ▶ To ensure this is, say, at most  $1/2$ , just need

$$n \geq O\left(k + \log \binom{d}{2k}\right) = O(k + k \log(d/k)).$$

30

## Restricted isometry property

$(\ell, \delta)$ -restricted isometry property (RIP):

$$(1 - \delta)\|\mathbf{z}\|_2^2 \leq \|\mathbf{Az}\|_2^2 \leq (1 + \delta)\|\mathbf{z}\|_2^2 \quad \text{for all } \ell\text{-sparse } \mathbf{z}.$$

- ▶ Many algorithms can recover unique sparsest solution under RIP (with  $\ell = O(k)$  and  $\delta = \Omega(1)$ ).
  - ▶ E.g., Basis pursuit, Lasso, orthogonal matching pursuit.