

# CSE 392: Matrix and Tensor Algorithms for Data

Instructor: Shashanka Ubaru

University of Texas, Austin  
Spring 2024

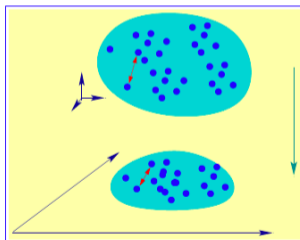
## Lecture 7: JL Lemma and subspace embedding

# Outline

- 1 Near orthogonal vectors and  $\epsilon$ -Net
- 2 Gaussian matrix properties
- 3 Johnson-Lindenstrauss Lemma
- 4 Subspace embedding

# High-dimensional vectors

- Often we deal with data vectors that are high-dimensional.
- **Dimensionality reduction:** One popular approach is to embed these vectors on a low-dimensional space.
- What criteria should we use to compute this low-dimensional embedding? What properties of the data do we wish to preserve?



## Near-orthogonal vectors

Given a  $d$ -dimensional space, what is the largest set of *mutually orthogonal* unit vectors  $\mathbf{x}_1, \dots, \mathbf{x}_t$  we can have? I.e. with the inner products

$$|\mathbf{x}_i^\top \mathbf{x}_j| = 0 \quad \forall i, j$$

## Near-orthogonal vectors

Given a  $d$ -dimensional space, what is the largest set of *mutually orthogonal* unit vectors  $\mathbf{x}_1, \dots, \mathbf{x}_t$  we can have? I.e. with the inner products

$$|\mathbf{x}_i^\top \mathbf{x}_j| = 0 \quad \forall i, j$$

**Answer:**  $d$

Given a  $d$ -dimensional space, what is the largest set of nearly orthogonal unit vectors  $\mathbf{x}_1, \dots, \mathbf{x}_t$ ? I.e. with the inner products

$$|\mathbf{x}_i^\top \mathbf{x}_j| \leq \epsilon \quad \forall i, j$$

Suppose  $\epsilon$  is a constant. E.g.  $\epsilon = 1/10$ .

## Near-orthogonal vectors

Given a  $d$ -dimensional space, what is the largest set of *mutually orthogonal* unit vectors  $\mathbf{x}_1, \dots, \mathbf{x}_t$  we can have? I.e. with the inner products

$$|\mathbf{x}_i^\top \mathbf{x}_j| = 0 \quad \forall i, j$$

**Answer:**  $d$

Given a  $d$ -dimensional space, what is the largest set of nearly orthogonal unit vectors  $\mathbf{x}_1, \dots, \mathbf{x}_t$ ? I.e. with the inner products

$$|\mathbf{x}_i^\top \mathbf{x}_j| \leq \epsilon \quad \forall i, j$$

Suppose  $\epsilon$  is a constant. E.g.  $\epsilon = 1/10$ .

**Answer:**  $2^{\Theta(d)}$

## Near-orthogonal vectors

**Claim:** There is an exponential number of nearly orthogonal unit vectors in  $d$ -dimensional space ( $\sim 2^d$ ).

**Proof approach:** One approach is to use *Probabilistic Argument*. For  $t = 2^{\Theta(d)}$ , define a random process which generates random vectors  $\mathbf{x}_1, \dots, \mathbf{x}_t$  that are unlikely to have large inner product

- Show that, with high probability,  $|\mathbf{x}_i^\top \mathbf{x}_j| \leq \epsilon \quad \forall i, j$ .
- Hence, there must exist some set of unit vectors with all pairwise inner-products bounded by  $\epsilon$ .



**Proof:** Let  $\mathbf{x}_1, \dots, \mathbf{x}_t$  be normalized Radmacher vectors, i.e., have independent random entries, each set to  $\pm 1/\sqrt{d}$  with equal probability.

$$\mathbb{E}[\mathbf{x}_i^\top \mathbf{x}_j] = ?$$

Let  $S = \mathbf{x}_i^\top \mathbf{x}_j = \sum_{i=1}^d c_i$ , where  $c_i$  is random  $\pm 1/d$ .

$S$  is sum of i.i.d random variables. Lets use Hoeffding's inequality:

### Hoeffding Inequality

Let  $c_1, \dots, c_d$  be independent random variables with each  $c_i \in [a_i, b_i]$ . Let  $\mathbb{E}[c_i] = \mu_i$  and  $\text{Var}[c_i] = \sigma_i^2$ . Let  $\mu = \sum_i \mu_i$  and  $\sigma^2 = \sum_i \sigma_i^2$ . Then, for and  $\alpha > 0$ ,  $S = \sum_i c_i$  satisfies

$$\Pr[|S - \mu| \geq \alpha] \leq 2e^{-\frac{2\alpha^2}{\sum_i (a_i - b_i)^2}}.$$

## Hoeffding Inequality

Let  $c_1, \dots, c_d$  be independent random variables with each  $c_i \in [a_i, b_i]$ . Let  $\mathbb{E}[c_i] = \mu_i$  and  $\text{Var}[c_i] = \sigma_i^2$ . Let  $\mu = \sum_i \mu_i$  and  $\sigma^2 = \sum_i \sigma_i^2$ . Then, for any  $\alpha > 0$ ,  $S = \sum_i c_i$  satisfies

$$\Pr[|S - \mu| \geq \alpha] \leq 2e^{-\frac{2\alpha^2}{\sum_i (a_i - b_i)^2}}.$$

Here,  $a_i = -1/d, b_i = 1/d$ .  $\mu_i = ?$

## Hoeffding Inequality

Let  $c_1, \dots, c_d$  be independent random variables with each  $c_i \in [a_i, b_i]$ . Let  $\mathbb{E}[c_i] = \mu_i$  and  $\text{Var}[c_i] = \sigma_i^2$ . Let  $\mu = \sum_i \mu_i$  and  $\sigma^2 = \sum_i \sigma_i^2$ . Then, for any  $\alpha > 0$ ,  $S = \sum_i c_i$  satisfies

$$\Pr[|S - \mu| \geq \alpha] \leq 2e^{-\frac{2\alpha^2}{\sum_i (a_i - b_i)^2}}.$$

Here,  $a_i = -1/d, b_i = 1/d$ .  $\mu_i = ?$

We have

$$\Pr[|\mathbf{x}_i^\top \mathbf{x}_j| \geq \epsilon] \leq 2e^{-\epsilon^2 d/2}$$

For any pair  $i, j$ , we have  $\Pr[|\mathbf{x}_i^\top \mathbf{x}_j| < \epsilon] > 1 - 2e^{-\epsilon^2 d/2}$ . Taking union bound over all possible pairs, we get

$$\Pr[|\mathbf{x}_i^\top \mathbf{x}_j| < \epsilon] > 1 - \binom{t}{2} 2e^{-\epsilon^2 d/2}$$

## Near-orthogonal vectors

- **Result:** In  $d$ -dimensional space, there are  $t = 2^{\Theta(\epsilon^2 d)}$  unit vectors with all pairwise inner products  $\leq \epsilon$ .
- **Alternate point of view :** Random vectors tend to be far apart (and roughly equidistant) in high-dimensions.
- **Curse of dimensionality:** If our data distribution is truly random, suppose we want to use say  $k$ -nearest neighbors to learn a function or classify points in  $\mathbb{R}^d$ , we typically need an exponential amount of data.
- Hope is that there exists low dimensional structure in our data.

## Alternate approach: $\epsilon$ -Nets

Some definitions:

- **Unit sphere:** Let  $\mathcal{S}_p^{d-1} \equiv \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x}\|_p = 1\}$ .

We will omit  $p$ , when  $p = 2$ , and  $d$  when in context.

- **Semi-norms from sets:** For symmetric matrix  $\mathbf{W} \in \mathbb{R}^{d \times d}$  and non-empty  $\mathcal{N} \subset \mathbb{R}^d$ , let

$$\|\mathbf{W}\|_{\mathcal{N}} \equiv \sup\{|\mathbf{x}^\top \mathbf{W} \mathbf{x}| / \|\mathbf{x}\|^2 \mid \mathbf{x} \in \mathcal{N}, \mathbf{x} \neq 0\}$$

so when  $\mathcal{N} \subset \mathcal{S}$ ,  $\|\mathbf{W}\|_{\mathcal{N}} \equiv \sup_{\mathbf{x} \in \mathcal{N}} |\mathbf{x}^\top \mathbf{W} \mathbf{x}|$ .

- **Embedding of  $\mathcal{N}$ :** For  $\mathcal{N} \subset \mathbb{R}^d$ ,  $\mathbf{B} \in \mathbb{R}^{m \times d}$ , and  $\beta \in (0, 1]$ ,  $\|\mathbf{B}^\top \mathbf{B} - \mathbf{I}\|_{\mathcal{N}} \leq \beta \implies \mathbf{B}$  is a  $\beta$ -embedding of  $\mathcal{N}$ .
- $\mathbf{B}^\top \mathbf{B} - \mathbf{I}$  is called the centered Grammian of  $\mathbf{B}$ .
- If  $\|\mathbf{B}^\top \mathbf{B} - \mathbf{I}\|_{\mathcal{S}} \leq \beta$ , then  $\mathbf{B}$  is a  $\beta$ -embedding of  $\mathbb{R}^d$ .

- $\mathcal{N} = \mathcal{N}(\epsilon)$  is an  $\epsilon$ -net of set  $\mathcal{P}$  if it is both:

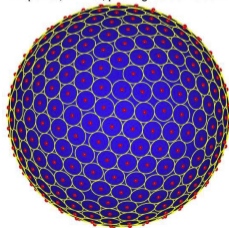
- ▶  $\epsilon$ -packing: all  $p \in \mathcal{N}$  at least  $\epsilon$  from  $\mathcal{N}$

$$d(p, \mathcal{N} \setminus \{p\}) \geq \epsilon \text{ for } p \in \mathcal{N}$$

- ▶  $\epsilon$ -covering: all  $p \in \mathcal{P}$  at most  $\epsilon$  from  $\mathcal{N}$

$$d(p, \mathcal{N}) \leq \epsilon \text{ for } p \in \mathcal{P}$$

PK points, N = 400, packing radius = 0.0924



## Sphere covering number

The unit sphere  $\mathcal{S}$  in  $\mathbb{R}^d$  has an  $\epsilon$ -net of size at most  $(1 + 2/\epsilon)^d$ .

Proof is through a volume argument. Since the points in  $\mathcal{N}(\epsilon)$  are  $\epsilon$ -separated, the balls of radii  $\epsilon/2$  centered at the points in  $\mathcal{N}(\epsilon)$  are disjoint. Also, all such balls lie in  $(1 + \epsilon/2)B_2^d$  where  $B_2^d$  denotes the unit Euclidean ball centered at the origin. So, we have

$$\text{vol}\left(\frac{\epsilon}{2}B_2^d\right) \cdot |\mathcal{N}(\epsilon)| \leq \text{vol}\left(\left(1 + \frac{\epsilon}{2}\right)B_2^d\right)$$

Since,  $\text{vol}(rB_2^d) = r^d \text{vol}(B_2^d)$ , we get

$$|\mathcal{N}(\epsilon)| \leq \left(1 + \frac{\epsilon}{2}\right)^d / \left(\frac{\epsilon}{2}\right)^d = \left(1 + \frac{2}{\epsilon}\right)^d.$$

## $\epsilon$ -Net bound

For  $\mathcal{N}_\epsilon$  an  $\epsilon$ -net of unit sphere  $\mathcal{S}$  in  $\mathbb{R}^d$  and  $\epsilon < 1$ , if matrix  $\mathbf{W}$  is symmetric, then

$$(1 - 2\epsilon)\|\mathbf{W}\|_2 \leq \|\mathbf{W}\|_{\mathcal{N}_\epsilon} \leq \|\mathbf{W}\|_{\mathcal{S}} = \|\mathbf{W}\|_2$$

and so if  $\mathbf{B}$  is a  $\beta$ -embedding of  $\mathcal{N}_\epsilon$ , then it is a  $\beta/(1 - 2\epsilon)$ - embedding of  $\mathcal{S}$ , and so of  $\mathbb{R}^d$ .



For  $\mathcal{N}_\epsilon$  an  $\epsilon$ -net of unit sphere  $\mathcal{S}$  in  $\mathbb{R}^d$  and  $\epsilon < 1$ , if matrix  $\mathbf{W}$  is symmetric, then

$$(1 - 2\epsilon)\|\mathbf{W}\|_2 \leq \|\mathbf{W}\|_{\mathcal{N}_\epsilon} \leq \|\mathbf{W}\|_{\mathcal{S}} = \|\mathbf{W}\|_2$$

and so if  $\mathbf{B}$  is a  $\beta$ -embedding of  $\mathcal{N}_\epsilon$ , then it is a  $\beta/(1 - 2\epsilon)$ -embedding of  $\mathcal{S}$ , and so of  $\mathbb{R}^d$ .

**Proof:** Let unit  $\mathbf{y}$  be such that  $|\mathbf{y}^\top \mathbf{W} \mathbf{y}| = \|\mathbf{W}\|_2 = \|\mathbf{W}\|_{\mathcal{S}}$ .

Since  $\mathcal{N}_\epsilon$  is an  $\epsilon$ -net, there is  $\mathbf{z}$  with  $\|\mathbf{z}\| \leq \epsilon$  and  $(\mathbf{y} - \mathbf{z}) \in \mathcal{N}_\epsilon$ .

Next,

$$\begin{aligned} \|\mathbf{W}\|_2 &= |\mathbf{y}^\top \mathbf{W} \mathbf{y}| = |(\mathbf{y} - \mathbf{z})^\top \mathbf{W} (\mathbf{y} - \mathbf{z}) + \mathbf{z}^\top \mathbf{W} \mathbf{y} + \mathbf{z}^\top \mathbf{W} (\mathbf{y} - \mathbf{z})| \\ &\leq |(\mathbf{y} - \mathbf{z})^\top \mathbf{W} (\mathbf{y} - \mathbf{z})| + |\mathbf{z}^\top \mathbf{W} \mathbf{y}| + |\mathbf{z}^\top \mathbf{W} (\mathbf{y} - \mathbf{z})| \\ &\leq \|\mathbf{W}\|_{\mathcal{N}_\epsilon} + \|\mathbf{z}\| \cdot \|\mathbf{W} \mathbf{y}\| + \|\mathbf{z}\| \cdot \|\mathbf{W} (\mathbf{y} - \mathbf{z})\| \\ &\leq \|\mathbf{W}\|_{\mathcal{N}_\epsilon} + 2\epsilon \|\mathbf{W}\|_2. \end{aligned}$$

# Independent Gaussians

Recall the norm estimation random vectors.

- **Gaussians are stable:** Given  $\mathbf{y} \in \mathbb{R}^d$ , if  $\mathbf{g} \in \mathbb{R}^d$  has entries i.i.d  $\mathcal{N}(0, 1)$ , then

$$\mathbf{g}^\top \mathbf{y} \sim \mathcal{N}(0, \|\mathbf{y}\|^2)$$

- A sum of independent Gaussians is Gaussian, and a scalar multiple of a Gaussian is Gaussian.

# Independent Gaussians

Recall the norm estimation random vectors.

- **Gaussians are stable:** Given  $\mathbf{y} \in \mathbb{R}^d$ , if  $\mathbf{g} \in \mathbb{R}^d$  has entries i.i.d  $\mathcal{N}(0, 1)$ , then

$$\mathbf{g}^\top \mathbf{y} \sim \mathcal{N}(0, \|\mathbf{y}\|^2)$$

- A sum of independent Gaussians is Gaussian, and a scalar multiple of a Gaussian is Gaussian.
- **Vector embedding:** Given a unit vector  $\mathbf{y} \in \mathbb{R}^d$ ,  $\epsilon \in (0, 1]$ . If  $\mathbf{G} \in \mathbb{R}^{m \times d}$  has independent entries  $g_{ij} \sim \mathcal{N}(0, 1/m)$ , then

$$\Pr\{|\|\mathbf{G}\mathbf{y}\|_2^2 - 1| \geq \epsilon\} \leq 2 \exp(-\epsilon^2 m/16).$$

We know  $\sqrt{m}\mathbf{G}\mathbf{y} \sim \mathcal{N}(0, 1)$  and squared norm is a  $\chi_m^2$  distribution. Using the standard bounds for concentration of a  $\chi_m^2$ , we get the above.

- With high probability,  $\mathbf{G}$   $\epsilon$ -embeds unit vectors  $\mathbf{y} \in \mathbb{R}^d$ . Also, for any fixed  $\mathbf{y} \in \mathbb{R}^d$ .

- **Gaussian width:** Given  $\mathcal{R} \subset \mathbb{R}^d$ , the Gaussian width of  $\mathcal{R}$  is

$$w(\mathcal{R}) \equiv \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(0, \mathbf{I})} \left[ \sup_{\mathbf{y}, \mathbf{x} \in \mathcal{R}} \mathbf{g}^\top (\mathbf{y} - \mathbf{x}) \right].$$

- Alternatively, the Gaussian width of  $\mathcal{R}$  is

$$w(\mathcal{R}) \equiv \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(0, \mathbf{I})} \left[ \sup_{\mathbf{y} \in \mathcal{R}} \mathbf{g}^\top \mathbf{y} / \|\mathbf{y}\| \right].$$

- Gaussian widths:

- ▶  $w(\mathbb{R}^d) \leq \sqrt{d}$
- ▶  $w(\mathcal{L}) \leq \sqrt{k}$  for  $\mathcal{L}$  a  $k$ -dimensional subspace.
- ▶  $w(\mathcal{R}) \leq \sqrt{2 \log |\mathcal{R}|}$  for finite  $\mathcal{R}$ .

# Gordon's theorem

## Gordon's theorem [G88]

For given  $\mathcal{R} \subset \mathbb{R}^d$ , if  $\mathbf{G} \in \mathbb{R}^{m \times d}$  has independent entries  $g_{ij} \sim \mathcal{N}(0, 1/m)$ , then

$$\Pr\{\|\mathbf{G}^\top \mathbf{G} - \mathbf{I}\|_{\mathcal{R}} \geq 2\beta + \beta^2\} \leq 2 \exp(-t^2/2),$$

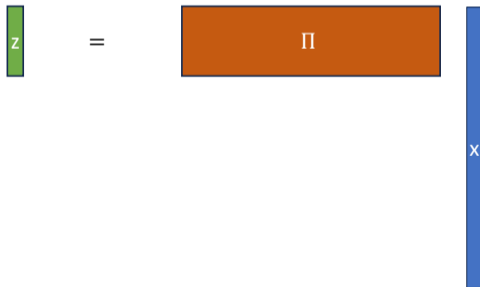
where  $\beta \equiv \frac{w(\mathcal{R})+1+t}{\sqrt{m}}$ .

# Euclidean dimensionality reduction

Johnson-Lindenstrauss, 1984

For any set of  $n$  data points  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$  there exists a *linear map*  $\Pi : \mathbb{R}^d \rightarrow \mathbb{R}^m$  where  $m = O(\frac{\log n}{\epsilon^2})$  such that for all  $i, j$ ,

$$(1 - \epsilon)\|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq \|\Pi\mathbf{x}_i - \Pi\mathbf{x}_j\|_2 \leq (1 + \epsilon)\|\mathbf{x}_i - \mathbf{x}_j\|_2$$



## Johnson-Lindenstrauss, 1984

For any set of  $n$  data points  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$  there exists a *linear map*  $\Pi : \mathbb{R}^d \rightarrow \mathbb{R}^m$  where  $m = O(\frac{\log n}{\epsilon^2})$  such that for all  $i, j$ ,

$$(1 - \epsilon) \|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq \|\Pi \mathbf{x}_i - \Pi \mathbf{x}_j\|_2 \leq (1 + \epsilon) \|\mathbf{x}_i - \mathbf{x}_j\|_2$$

### Proof:

- We show that for a Gaussian matrix  $\mathbf{G} \in \mathbb{R}^{m \times d}$  has independent entries  $g_{ij} \sim \mathcal{N}(0, 1/m)$ , the result holds.
- Use the vector embedding result from before (squared norm  $\|\mathbf{G}(\mathbf{x}_i - \mathbf{x}_j)\|^2$  is  $\chi_m^2$  distribution with mean  $\|\mathbf{x}_i - \mathbf{x}_j\|^2$ ).
- Set the probability to  $1/n^2$ . Since we have  $< n^2$  possible pairs  $i, j$ , using union bound, we get the result.
- For vectors in finite  $\mathcal{R} \subset \mathbb{R}^d$ , we can use Gordon's theorem to prove similar result.

Original result used rows of a random orthogonal matrix. Random sign matrix, where rows are Radamacher vectors, is an example.

# Oblivious subspace embedding

- For real  $\mathbf{x}, \mathbf{y}$  and  $\epsilon$ , by  $\mathbf{x} = (1 \pm \epsilon)\mathbf{y}$  we mean that  $|\mathbf{x} - \mathbf{y}| \leq \epsilon|\mathbf{y}|$ .
- **Embedding:** A matrix  $\mathbf{S} \in \mathbb{R}^{m \times n}$  is an  $\epsilon$ -embedding of set  $\mathcal{P} \subset \mathbb{R}^n$  if, for all  $\mathbf{y} \in \mathcal{P}$ ,

$$\|\mathbf{S}\mathbf{y}\|_2 = (1 \pm \epsilon)\|\mathbf{y}\|_2.$$

We will call  $\mathbf{S}$  a “sketching matrix”.

## Subspace embedding

For  $\mathbf{A} \in \mathbb{R}^{n \times d}$ , a matrix  $\mathbf{S} \in \mathbb{R}^{m \times n}$  is a subspace  $\epsilon$ -embedding for  $\mathbf{A}$  if  $\mathbf{S}$  is an  $\epsilon$ -embedding for  $\text{span}(\mathbf{A}) = \{\mathbf{A}\mathbf{x} \mid \mathbf{x} \in \mathbb{R}^d\}$ . I.e., for all  $\mathbf{x} \in \mathbb{R}^d$ ,

$$\|\mathbf{S}\mathbf{A}\mathbf{x}\|_2 = (1 \pm \epsilon)\|\mathbf{A}\mathbf{x}\|_2.$$

We will call  $\mathbf{S}\mathbf{A}$  a “sketch”.



An *Oblivious* subspace embedding is:

- A probability distribution  $\mathcal{D}$  over matrices  $\mathbf{S} \in \mathbb{R}^{m \times n}$ , so that
- For any unknown but fixed matrix  $\mathbf{A}$ ,  $\mathbf{S}$  is a subspace  $\epsilon$ -embedding for  $\mathbf{A}$  with high probability.

**Advantages:**

- Distribution  $\mathcal{D}$  does not depend on input data. Construct  $\mathbf{S}$  without knowing  $\mathbf{A}$ .
- *Streaming*: when entries of  $\mathbf{A}$  change,  $\mathbf{S}\mathbf{A}$  is easy to update.
- *Distributed*: If each  $p$  processor has matrix  $\mathbf{A}^{(p)}$  and  $\mathbf{A} = \sum_p \mathbf{A}^{(p)}$ , compute sketch at each processor.
- *Analysis*: If  $\mathbf{U}$  has  $\text{span}(\mathbf{U}) = \text{span}(\mathbf{A})$ , then the embedding condition holds for  $\text{span}(\mathbf{A})$  iff it holds for  $\text{span}(\mathbf{U})$ . So, we can assume  $\mathbf{A}$  is orthonormal.

## Subspace embedding

Given  $\epsilon, \delta > 0$ ,  $\mathbf{A} \in \mathbb{R}^{n \times d}$ , and unit vector  $\mathbf{y} \in \mathbb{R}^n$ . There is  $m = O\left(\frac{d \log(1/\delta)}{\epsilon^2}\right)$  so that if  $\mathbf{S} \in \mathbb{R}^{m \times n}$  is randomly chosen from a fixed (oblivious to  $\mathbf{A}$ ) distribution with the property that  $\mathbf{S}$  is an  $\epsilon/6$ -embedding of  $\mathbf{y}$  (JL property) with failure probability  $\delta' = K_1 \exp(-K_2 \epsilon^2 m)$ , for some  $K_1, K_2 > 0$ , then  $\mathbf{S}$  is a *subspace*  $\epsilon$ -embedding for  $\mathbf{A}$  with failure probability  $\delta$ .

## Subspace embedding

Given  $\epsilon, \delta > 0$ ,  $\mathbf{A} \in \mathbb{R}^{n \times d}$ , and unit vector  $\mathbf{y} \in \mathbb{R}^n$ . There is  $m = O\left(\frac{d \log(1/\delta)}{\epsilon^2}\right)$  so that if  $\mathbf{S} \in \mathbb{R}^{m \times n}$  is randomly chosen from a fixed (oblivious to  $\mathbf{A}$ ) distribution with the property that  $\mathbf{S}$  is an  $\epsilon/6$ -embedding of  $\mathbf{y}$  (JL property) with failure probability  $\delta' = K_1 \exp(-K_2 \epsilon^2 m)$ , for some  $K_1, K_2 > 0$ , then  $\mathbf{S}$  is a *subspace  $\epsilon$ -embedding* for  $\mathbf{A}$  with failure probability  $\delta$ .

**Proof:** We will use the  $\epsilon$ -net argument with the  $\epsilon$ -embedding (JL) property.

- Since  $\mathbf{S}$  is oblivious, assume  $\mathbf{A}$  has orthonormal columns.
- For some  $\epsilon_0 > 0$  (to be determined), we pick an  $\epsilon_0$ -net  $\mathcal{N}_{\epsilon_0} \subset \mathcal{S}$ .
- For  $\mathbf{x} \in \mathcal{N}_{\epsilon_0}$ ,  $\mathbf{y} = \mathbf{A}\mathbf{x} \in \text{span}(\mathbf{A})$  is a unit vector.
- Let  $\mathbf{W} := \mathbf{A}^\top \mathbf{S}^\top \mathbf{S} \mathbf{A} - \mathbf{I}$ .
- Note that, for any  $\beta \in (0, 1]$ ,  $(1 + \beta)^2 \leq (1 + 3\beta)$  and  $(1 - \beta)^2 \geq (1 - 3\beta)$ .

So, we have  $|\|\mathbf{S}\mathbf{y}\|_2^2 - 1| \leq \epsilon/2$ . Also,

$$|\|\mathbf{S}\mathbf{y}\|_2^2 - 1| = |\mathbf{y}^\top \mathbf{S}^\top \mathbf{S} \mathbf{y} - \mathbf{y}^\top \mathbf{y}| = |\mathbf{x}^\top \mathbf{A}^\top \mathbf{S}^\top \mathbf{S} \mathbf{A} \mathbf{x} - \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x}| = |\mathbf{x}^\top \mathbf{W} \mathbf{x}| \leq \epsilon/2$$

with failure probability  $\delta'$ .

Applying this to all vectors in  $\mathcal{N}_{\epsilon_0}$ , and union bound,

$$\|\mathbf{W}\|_{\mathcal{N}_{\epsilon_0}} \leq \epsilon/2 \text{ with failure probability } \leq \delta' |\mathcal{N}_{\epsilon_0}|$$

Using the relation between  $\|\mathbf{W}\|_{\mathcal{S}}$  and  $\|\mathbf{W}\|_{\mathcal{N}_{\epsilon_0}}$  and the bound on net size  $|\mathcal{N}_{\epsilon_0}|$ ,

$$\|\mathbf{W}\|_{\mathcal{S}} \leq \epsilon/2/(1 - \epsilon_0) \text{ with failure probability } \leq \delta' |\mathcal{N}_{\epsilon_0}| \leq (1 + \frac{2}{\epsilon_0})^d K_1 \exp(-K_2 \epsilon^2 m).$$

For fixed  $\epsilon_0$ , there is  $m = O(\frac{d \log(1/\delta)}{\epsilon^2})$ , so that this is at most  $\delta$ .

For  $\epsilon_0 \leq 1/2$ , we have  $\|\mathbf{W}\|_{\mathcal{S}} \leq \epsilon$ .

## Further Reading

- Woodruff, David P. “Sketching as a tool for numerical linear algebra.” *Foundations and Trends® in Theoretical Computer Science* 10.1–2 (2014): 1-157.
- Martinsson, P. G., and Tropp, J. “Randomized numerical linear algebra: foundations and algorithms”. arXiv preprint arXiv:2002.01387 (2020).

Questions?