

# Sketching as a Tool for Numerical Linear Algebra All Lectures

David Woodruff  
IBM Almaden



# Massive data sets

---

## *Examples*

- Internet traffic logs
- Financial data
- etc.

## *Algorithms*

- Want nearly linear time or less
- Usually at the cost of a randomized approximation

# Regression analysis

---

## *Regression*

- Statistical method to study dependencies between variables in the presence of noise.

# Regression analysis

---

## *Linear Regression*

- Statistical method to study **linear** dependencies between variables in the presence of noise.

# Regression analysis

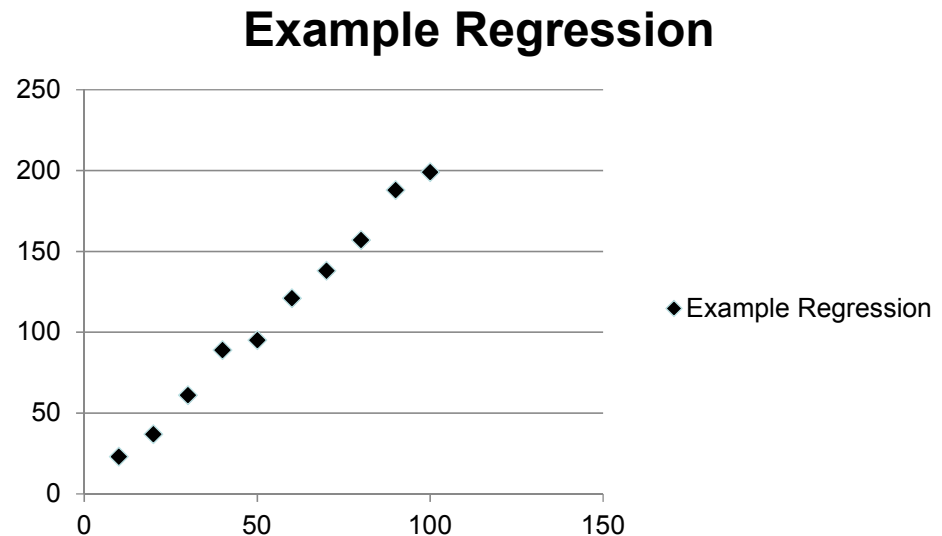
---

## Linear Regression

- Statistical method to study **linear** dependencies between variables in the presence of noise.

## Example

- Ohm's law  $V = R \cdot I$



# Regression analysis

---

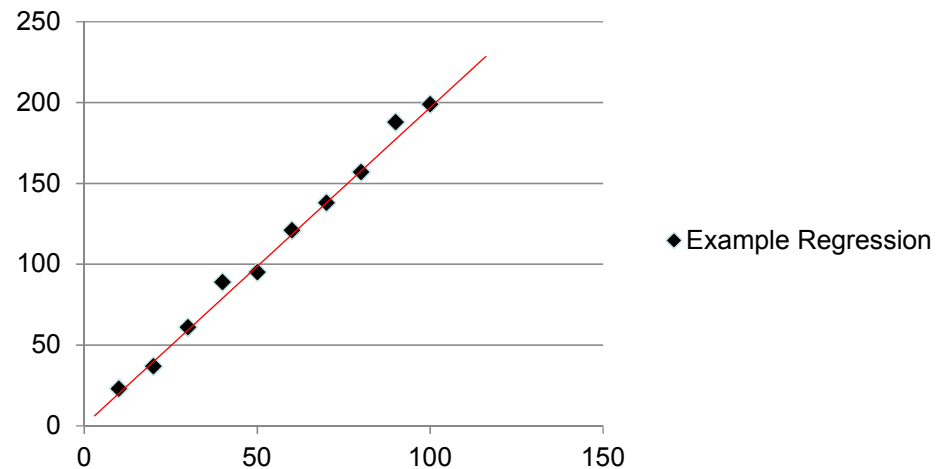
## *Linear Regression*

- Statistical method to study **linear** dependencies between variables in the presence of noise.

## *Example*

- Ohm's law  $V = R \cdot I$
- Find linear function that best fits the data

**Example Regression**



# Regression analysis

---

## *Linear Regression*

- Statistical method to study **linear** dependencies between variables in the presence of noise.

## *Standard Setting*

- One measured variable  $b$
- A set of predictor variables  $a_1, \dots, a_d$
- Assumption:

$$b = x_0 + a_1 x_1 + \dots + a_d x_d + \varepsilon$$

- $\varepsilon$  is assumed to be noise and the  $x_i$  are model parameters we want to learn
- Can assume  $x_0 = 0$
- Now consider  $n$  observations of  $b$

# Regression analysis

---

## *Matrix form*

**Input:**  $n \times d$ -matrix  $A$  and a vector  $b = (b_1, \dots, b_n)$   
 $n$  is the number of observations;  $d$  is the number of predictor variables

**Output:**  $x^*$  so that  $Ax^*$  and  $b$  are close

- Consider the over-constrained case, when  $n \gg d$
- Can assume that  $A$  has full column rank



# Regression analysis

---

## *Least Squares Method*

- Find  $x^*$  that minimizes  $\|Ax-b\|_2^2 = \sum (b_i - \langle A_{i*}, x \rangle)^2$
- $A_{i*}$  is  $i$ -th row of  $A$
- Certain desirable statistical properties

# Regression analysis

---

## Geometry of regression

- We want to find an  $x$  that minimizes  $\|Ax-b\|_2$
- The product  $Ax$  can be written as

$$A_{*1}x_1 + A_{*2}x_2 + \dots + A_{*d}x_d$$

where  $A_{*i}$  is the  $i$ -th column of  $A$

- This is a linear  $d$ -dimensional subspace
- The problem is equivalent to computing the point of the column space of  $A$  nearest to  $b$  in  $l_2$ -norm

# Regression analysis

---

## *Solving least squares regression via the normal equations*

- How to find the solution  $x$  to  $\min_x \|Ax-b\|_2$  ?
- Equivalent problem:  $\min_x \|Ax-b\|_2^2$ 
  - Write  $b = Ax' + b'$ , where  $b'$  orthogonal to columns of  $A$
  - Cost is  $\|A(x-x')\|_2^2 + \|b'\|_2^2$  by Pythagorean theorem
  - Optimal solution  $x$  if and only if  $A^T(Ax-b) = A^T(Ax-Ax') = 0$
  - Normal Equation:  $A^T Ax = A^T b$  for any optimal  $x$
  - $x = (A^T A)^{-1} A^T b$
- If the columns of  $A$  are not linearly independent, the Moore-Penrose pseudoinverse gives a minimum norm solution  $x$

# Moore-Penrose Pseudoinverse

---

## Singular Value Decomposition (SVD)

Any matrix  $A = U \cdot \Sigma \cdot V^T$

- $U$  has orthonormal columns
- $\Sigma$  is diagonal with non-increasing non-negative entries down the diagonal
- $V^T$  has orthonormal rows
  
- Pseudoinverse  $A^- = V \Sigma^{-1} U^T$ 
  - Where  $\Sigma^{-1}$  is a diagonal matrix with  $i$ -th diagonal entry equal to  $1/\Sigma_{ii}$  if  $\Sigma_{ii} > 0$  and is 0 otherwise
  
- $\min_x \|Ax-b\|_2^2$  not unique when columns of  $A$  are linearly independent, but  $x = A^-b$  has minimum norm

# Moore-Penrose Pseudoinverse

---

- Any optimal solution  $x$  has the form  $A^{-}b + (I - V'V'^T)z$ , where  $V'$  corresponds to the rows  $i$  of  $V^T$  for which  $\Sigma_{i,i} > 0$
- **Why?**
- Because  $A(I - V'V'^T)z = 0$ , so  $A^{-}b + (I - V'V'^T)z$  is a solution. This is a  $d - \text{rank}(A)$  dimensional affine space so it spans all optimal solutions
- Since  $A^{-}b$  is in column span of  $V'$ , by Pythagorean theorem,  $|A^{-}b + (I - V'V'^T)z|_2^2 = |A^{-}b|_2^2 + |(I - V'V'^T)z|_2^2 \geq |A^{-}b|_2^2$

# Time Complexity

---

## *Solving least squares regression via the normal equations*

- Need to compute  $x = A^{-1}b$
- Naively this takes  $nd^2$  time
- Can do  $nd^{1.376}$  using fast matrix multiplication
- But we want much better running time!

# Sketching to solve least squares regression

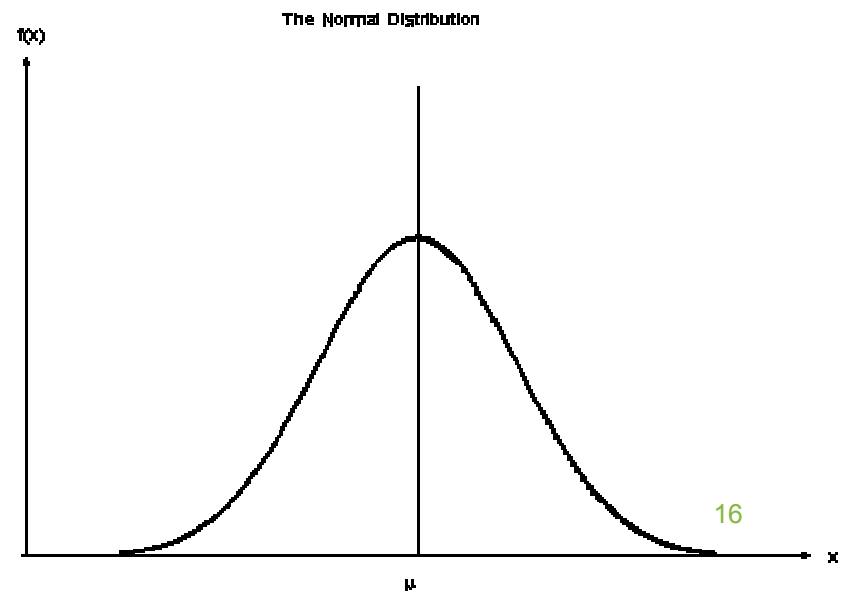
---

- How to find an approximate solution  $x$  to  $\min_x \|Ax-b\|_2$  ?
- **Goal:** output  $x'$  for which  $\|Ax'-b\|_2 \leq (1+\epsilon) \min_x \|Ax-b\|_2$  with high probability
- Draw  $S$  from a  $k \times n$  random family of matrices, for a value  $k \ll n$
- Compute  $S^*A$  and  $S^*b$
- Output the solution  $x'$  to  $\min_{x'} \|(SA)x-(Sb)\|_2$ 
  - $x' = (SA)^{-1}Sb$

# How to choose the right sketching matrix $S$ ?

---

- Recall: output the solution  $x'$  to  $\min_{x'} |(SA)x - (Sb)|_2$
- Lots of matrices work
- $S$  is  $d/\epsilon^2 \times n$  matrix of i.i.d. Normal random variables
- To see why this works, we introduce the notion of a subspace embedding





# Subspace Embeddings

- Let  $k = O(d/\epsilon^2)$
- Let  $S$  be a  $k \times n$  matrix of i.i.d. normal  $N(0, 1/k)$  random variables
- For any fixed  $d$ -dimensional subspace, i.e., the column space of an  $n \times d$  matrix  $A$ 
  - W.h.p., for all  $x$  in  $\mathbb{R}^d$ ,  $|SAx|_2 = (1 \pm \epsilon)|Ax|_2$
- Entire column space of  $A$  is preserved

*Why is this true?*

# Subspace Embeddings – A Proof

- Want to show  $\|SAx\|_2 = (1 \pm \epsilon)\|Ax\|_2$  for all  $x$
- Can assume columns of  $A$  are orthonormal (since we prove this for all  $x$ )
- **Claim:**  $SA$  is a  $k \times d$  matrix of i.i.d.  $N(0, 1/k)$  random variables
  - First property: for two independent random variables  $X$  and  $Y$ , with  $X$  drawn from  $N(0, a^2)$  and  $Y$  drawn from  $N(0, b^2)$ , we have  $X+Y$  is drawn from  $N(0, a^2 + b^2)$

# $X+Y$ is drawn from $N(0, a^2 + b^2)$

- Probability density function  $f_Z$  of  $Z = X+Y$  is convolution of probability density functions  $f_X$  and  $f_Y$

- $f_Z(z) = \int f_Y(z - x)f_X(x) dx$

- $f_X(x) = \frac{1}{a(2\pi)^{.5}} e^{-x^2/2a^2}$  ,  $f_Y(y) = \frac{1}{b(2\pi)^{.5}} e^{-x^2/2b^2}$

- $f_Z(z) = \int \frac{1}{a(2\pi)^{.5}} e^{-(z-x)^2/2a^2} \frac{1}{b(2\pi)^{.5}} e^{-x^2/2b^2} dx$

$$= \frac{1}{(2\pi)^{.5}(a^2+b^2)^{.5}} e^{-z^2/2(a^2+b^2)} \int \frac{(a^2+b^2)^{.5}}{(2\pi)^{.5}ab} e^{-\frac{\left(x - \frac{b^2z}{a^2+b^2}\right)^2}{2\left(\frac{(ab)^2}{a^2+b^2}\right)}} dx$$

$X+Y$  is drawn from  $N(0, a^2 + b^2)$

$$\text{Calculation: } e^{-\frac{(z-x)^2}{2a^2} - \frac{x^2}{2b^2}} = e^{-\frac{z^2}{2(a^2+b^2)} - \frac{\left(x - \frac{b^2z}{a^2+b^2}\right)^2}{2\left(\frac{(ab)^2}{a^2+b^2}\right)}}$$

$$\text{Density of Gaussian distribution: } \int \frac{(a^2+b^2)^{-5}}{(2\pi)^{-5}ab} e^{-\frac{\left(x - \frac{b^2z}{a^2+b^2}\right)^2}{2\left(\frac{(ab)^2}{a^2+b^2}\right)}} dx = 1$$

# Rotational Invariance

- Second property: if  $u, v$  are vectors with  $\langle u, v \rangle = 0$ , then  $\langle g, u \rangle$  and  $\langle g, v \rangle$  are independent, where  $g$  is a vector of i.i.d.  $N(0, 1/k)$  random variables
- **Why?**
- If  $g$  is an  $n$ -dimensional vector of i.i.d.  $N(0, 1)$  random variables, and  $R$  is a fixed matrix, then the probability density function of  $Rg$  is

$$f(x) = \frac{1}{\det(RR^T)(2\pi)^{d/2}} e^{-\frac{x^T(RR^T)^{-1}x}{2}}$$

- $RR^T$  is the covariance matrix
- For a rotation matrix  $R$ , the distribution of  $Rg$  and of  $g$  are the same

# Orthogonal Implies Independent

- Want to show: if  $u, v$  are vectors with  $\langle u, v \rangle = 0$ , then  $\langle g, u \rangle$  and  $\langle g, v \rangle$  are independent, where  $g$  is a vector of i.i.d.  $N(0, 1/k)$  random variables
- Choose a rotation  $R$  which sends  $u$  to  $\alpha e_1$ , and sends  $v$  to  $\beta e_2$
- $\langle g, u \rangle = \langle gR, R^T u \rangle = \langle h, \alpha e_1 \rangle = \alpha h_1$
- $\langle g, v \rangle = \langle gR, R^T v \rangle = \langle h, \beta e_2 \rangle = \beta h_2$   
where  $h$  is a vector of i.i.d.  $N(0, 1/k)$  random variables
- Then  $h_1$  and  $h_2$  are independent by definition

# Where were we?

- **Claim:** SA is a  $k \times d$  matrix of i.i.d.  $N(0, 1/k)$  random variables
- **Proof:** The rows of SA are independent
  - Each row is:  $\langle g, A_1 \rangle, \langle g, A_2 \rangle, \dots, \langle g, A_d \rangle$
  - First property implies the entries in each row are  $N(0, 1/k)$  since the columns  $A_i$  have unit norm
  - Since the columns  $A_i$  are orthonormal, the entries in a row are independent by our second property

# Back to Subspace Embeddings

- Want to show  $\|SAx\|_2 = (1 \pm \epsilon)\|Ax\|_2$  for all  $x$
  - Can assume columns of  $A$  are orthonormal
  - Can also assume  $x$  is a unit vector
  - $SA$  is a  $k \times d$  matrix of i.i.d.  $N(0, 1/k)$  random variables
  
  - Consider any fixed unit vector  $x \in R^d$
  - $\|SAx\|_2^2 = \sum_{i \in [k]} \langle g_i, x \rangle^2$ , where  $g_i$  is  $i$ -th row of  $SA$
  - Each  $\langle g_i, x \rangle^2$  is distributed as  $N\left(0, \frac{1}{k}\right)^2$
  - $E[\langle g_i, x \rangle^2] = 1/k$ , and so  $E[\|SAx\|_2^2] = 1$
- How concentrated is  $\|SAx\|_2^2$  about its expectation?*



# Johnson-Lindenstrauss Theorem

- Suppose  $h_1, \dots, h_k$  are i.i.d.  $N(0,1)$  random variables
- Then  $G = \sum_i h_i^2$  is a  $\chi^2$ -random variable
- Apply known tail bounds to  $G$ :
  - (Upper)  $\Pr[G \geq k + 2(kx)^{.5} + 2x] \leq e^{-x}$
  - (Lower)  $\Pr[G \leq k - 2(kx)^{.5}] \leq e^{-x}$
- If  $x = \frac{\epsilon^2 k}{16}$ , then  $\Pr[G \in k(1 \pm \epsilon)] \geq 1 - 2e^{-\epsilon^2 k/16}$
- If  $k = \Theta(\epsilon^{-2} \log(\frac{1}{\delta}))$ , this probability is  $1 - \delta$
- $\Pr[|SAx|_2^2 \in (1 \pm \epsilon)] \geq 1 - 2^{-\Theta(d)}$

*This only holds for a fixed  $x$ , how to argue for all  $x$ ?*

# Net for Sphere

- Consider the sphere  $S^{d-1}$
- Subset  $N$  is a  $\gamma$ -net if for all  $x \in S^{d-1}$ , there is a  $y \in N$ , such that  $\|x - y\|_2 \leq \gamma$
- Greedy construction of  $N$ 
  - While there is a point  $x \in S^{d-1}$  of distance larger than  $\gamma$  from every point in  $N$ , include  $x$  in  $N$
- The sphere of radius  $\gamma/2$  around every point in  $N$  is contained in the sphere of radius  $1 + \gamma$  around  $0^d$
- Further, all such spheres are disjoint
- Ratio of volume of  $d$ -dimensional sphere of radius  $1 + \gamma$  to  $d$ -dimensional sphere of radius  $\gamma/2$  is  $(1 + \gamma)^d / (\gamma/2)^d$ , so  $|N| \leq (1 + \gamma)^d / (\gamma/2)^d$

# Net for Subspace

- Let  $M = \{Ax \mid x \text{ in } N\}$ , so  $|M| \leq (1 + \gamma)^d / (\gamma/2)^d$
- Claim: For every  $x$  in  $S^{d-1}$ , there is a  $y$  in  $M$  for which  $|Ax - y|_2 \leq \gamma$
- Proof: Let  $x'$  in  $S^{d-1}$  be such that  $|x - x'|_2 \leq \gamma$   
Then  $|Ax - Ax'|_2 = |x - x'|_2 \leq \gamma$ , using that the columns of  $A$  are orthonormal. Set  $y = Ax'$

# Net Argument

- For a fixed unit  $x$ ,  $\Pr[|SAx|_2^2 \in (1 \pm \epsilon)] \geq 1 - 2^{-\Theta(d)}$
- For a fixed pair of unit  $x, x'$ ,  $|SAx|_2^2, |SAx'|_2^2, |SA(x - x')|_2^2$  are all  $1 \pm \epsilon$  with probability  $1 - 2^{-\Theta(d)}$
- $|SA(x - x')|_2^2 = |SAx|_2^2 + |SAx'|_2^2 - 2 \langle SAx, SAx' \rangle$
- $|A(x - x')|_2^2 = |Ax|_2^2 + |Ax'|_2^2 - 2 \langle Ax, Ax' \rangle$ 
  - So  $\Pr[\langle Ax, Ax' \rangle = \langle SAx, SAx' \rangle \pm 0(\epsilon)] = 1 - 2^{-\Theta(d)}$
- Choose a  $1/2$ -net  $M = \{Ax \mid x \text{ in } N\}$  of size  $5^d$
- By a union bound, for all pairs  $y, y'$  in  $M$ ,
$$\langle y, y' \rangle = \langle Sy, Sy' \rangle \pm 0(\epsilon)$$
- Condition on this event
- By linearity, if this holds for  $y, y'$  in  $M$ , for  $\alpha y, \beta y'$  we have
$$\langle \alpha y, \beta y' \rangle = \alpha\beta \langle Sy, Sy' \rangle \pm 0(\epsilon \alpha\beta)$$

# Finishing the Net Argument

- Let  $y = Ax$  for an arbitrary  $x \in S^{d-1}$
- Let  $y_1 \in M$  be such that  $|y - y_1|_2 \leq \gamma$
- Let  $\alpha$  be such that  $|\alpha(y - y_1)|_2 = 1$ 
  - $\alpha \geq 1/\gamma$  (could be infinite)
- Let  $y_2' \in M$  be such that  $|\alpha(y - y_1) - y_2'|_2 \leq \gamma$
- Then  $\left|y - y_1 - \frac{y_2'}{\alpha}\right|_2 \leq \frac{\gamma}{\alpha} \leq \gamma^2$
- Set  $y_2 = \frac{y_2'}{\alpha}$ . Repeat, obtaining  $y_1, y_2, y_3, \dots$  such that for all integers  $i$ ,
$$|y - y_1 - y_2 - \dots - y_i|_2 \leq \gamma^i$$
- Implies  $|y_i|_2 \leq \gamma^{i-1} + \gamma^i \leq 2\gamma^{i-1}$

# Finishing the Net Argument

- Have  $y_1, y_2, y_3, \dots$  such that  $y = \sum_i y_i$  and  $|y_i|_2 \leq 2\gamma^{i-1}$
- $|Sy|_2^2 = |S \sum_i y_i|_2^2$ 
  - $= \sum_i |Sy_i|_2^2 + 2 \sum_{i,j} \langle Sy_i, Sy_j \rangle$
  - $= \sum_i |y_i|_2^2 + 2 \sum_{i,j} \langle y_i, y_j \rangle \pm O(\epsilon) \sum_{i,j} |y_i|_2 |y_j|_2$
  - $= |\sum_i y_i|_2^2 \pm O(\epsilon)$
  - $= |y|_2^2 \pm O(\epsilon)$
  - $= 1 \pm O(\epsilon)$
- Since this held for an arbitrary  $y = Ax$  for unit  $x$ , by linearity it follows that for all  $x$ ,  $|SAx|_2 = (1 \pm \epsilon)|Ax|_2$

# Back to Regression

- We showed that  $S$  is a subspace embedding, that is, simultaneously for all  $x$ ,

$$|SAx|_2 = (1 \pm \epsilon)|Ax|_2$$

*What does this have to do with regression?*

# Subspace Embeddings for Regression

- Want  $x$  so that  $\|Ax-b\|_2 \leq (1+\varepsilon) \min_y \|Ay-b\|_2$
- Consider subspace  $L$  spanned by columns of  $A$  together with  $b$
- Then for all  $y$  in  $L$ ,  $\|Sy\|_2 = (1 \pm \varepsilon) \|y\|_2$
- Hence,  $\|S(Ax-b)\|_2 = (1 \pm \varepsilon) \|Ax-b\|_2$  for all  $x$
- Solve  $\operatorname{argmin}_y \|(SA)y - (Sb)\|_2$
- Given  $SA$ ,  $Sb$ , can solve in  $\operatorname{poly}(d/\varepsilon)$  time

*Only problem is computing  $SA$  takes  $O(nd^2)$  time*



# How to choose the right sketching matrix $S$ ? [S]

---

- $S$  is a Subsampled Randomized Hadamard Transform
  - $S = P^*H^*D$
  - $D$  is a diagonal matrix with  $+1, -1$  on diagonals
  - $H$  is the Hadamard transform
  - $P$  just chooses a random (small) subset of rows of  $H^*D$
  - $S^*A$  can be computed in  $O(nd \log n)$  time

*Why does it work?*

# Why does this work?

---

- We can again assume columns of  $A$  are orthonormal
- It suffices to show  $|SAx|_2^2 = |PHDAx|_2^2 = 1 \pm \epsilon$  for all  $x$
- $HD$  is a rotation matrix, so  $|HDAx|_2^2 = |Ax|_2^2 = 1$  for any  $x$ 
  - Notation: let  $y = Ax$
- Flattening Lemma: For any fixed  $y$ ,

$$\Pr [ |HDy|_\infty \geq C \frac{\log^5 nd/\delta}{n \cdot 5} ] \leq \frac{\delta}{2d}$$

# Proving the Flattening Lemma

---

- **Flattening Lemma:**  $\Pr [|\text{HDy}|_\infty \geq C \frac{\log^5 nd/\delta}{n^5}] \leq \frac{\delta}{2d}$
- Let  $C > 0$  be a constant. We will show for a fixed  $i$  in  $[n]$ ,
 
$$\Pr [|(HDy)_i| \geq C \frac{\log^5 nd/\delta}{n^5}] \leq \frac{\delta}{2nd}$$
- If we show this, we can apply a union bound over all  $i$
- $|(HDy)_i| = \sum_j H_{i,j} D_{j,j} y_j$
- (Azuma-Hoeffding)  $\Pr [|\sum_j Z_j| > t] \leq 2e^{-\frac{t^2}{2 \sum_j \beta_j^2}}$ , where  $|Z_j| \leq \beta_j$  with probability 1
  - $Z_j = H_{i,j} D_{j,j} y_j$  has 0 mean
  - $|Z_j| \leq \frac{|y_j|}{n^5} = \beta_j$  with probability 1
  - $\sum_j \beta_j^2 = \frac{1}{n}$
- $\Pr \left[ |\sum_j Z_j| > \frac{C \log^5 \left(\frac{nd}{\delta}\right)}{n^5} \right] \leq 2e^{-\frac{C^2 \log^2 \left(\frac{\delta}{nd}\right)}{2}} \leq \frac{\delta}{2nd}$

# Consequence of the Flattening Lemma

---

- Recall columns of  $A$  are orthonormal
- HDA has orthonormal columns
- Flattening Lemma implies  $|HDAe_i|_\infty \leq C \frac{\log^5 nd/\delta}{n^{.5}}$  with probability  $1 - \frac{\delta}{2d}$  for a fixed  $i \in [d]$
- With probability  $1 - \frac{\delta}{2}$ ,  $|e_j HDAe_i| \leq C \frac{\log^5 nd/\delta}{n^{.5}}$  for all  $i, j$
- Given this,  $|e_j HDA|_2 \leq C \frac{d^{.5} \log^5 nd/\delta}{n^{.5}}$  for all  $j$

(Can be optimized further)

# Matrix Chernoff Bound

---

- Let  $X_1, \dots, X_s$  be independent copies of a symmetric random matrix  $X \in \mathbb{R}^{d \times d}$  with  $E[X] = 0$ ,  $|X|_2 \leq \gamma$ , and  $|E[X^T X]|_2 \leq \sigma^2$ . Let  $W = \frac{1}{s} \sum_{i \in [s]} X_i$ . For any  $\epsilon > 0$ ,

$$\Pr[|W|_2 > \epsilon] \leq 2d \cdot e^{-s\epsilon^2 / (\sigma^2 + \frac{\gamma\epsilon}{3})}$$

(here  $|W|_2 = \sup |Wx|_2 / |x|_2$ )

- Let  $V = HDA$ , and recall  $V$  has orthonormal columns
- Suppose  $P$  in the  $S = \text{PHD}$  definition samples uniformly with replacement. If row  $i$  is sampled in the  $j$ -th sample, then  $P_{j,i} = n$ , and is 0 otherwise
- Let  $Y_i$  be the  $i$ -th sampled row of  $V = HDA$
- Let  $X_i = I_d - n \cdot Y_i^T Y_i$ 
  - $E[X_i] = I_d - n \cdot \sum_j \left(\frac{1}{n}\right) V_i^T V_i = I_d - V^T V = 0^d$
  - $|X_i|_2 \leq |I_d|_2 + n \cdot \max |e_j HDA|_2^2 = 1 + n \cdot C^2 \log\left(\frac{nd}{\delta}\right) \cdot \frac{d}{n} = \Theta(d \log\left(\frac{nd}{\delta}\right))$

# Matrix Chernoff Bound

---

- Recall: let  $Y_i$  be the  $i$ -th sampled row of  $V = HDA$
- Let  $X_i = I_d - n \cdot Y_i^T Y_i$
- $$E[X^T X + I_d] = I_d + I_d - 2n E[Y_i^T Y_i] + n^2 E[Y_i^T Y_i Y_i^T Y_i]$$

$$= 2I_d - 2I_d + n^2 \sum_i \left(\frac{1}{n}\right) \cdot v_i^T v_i v_i^T v_i = n \sum_i v_i^T v_i \cdot |v_i|_2^2$$
- Define  $Z = n \sum_i v_i^T v_i C^2 \log\left(\frac{nd}{\delta}\right) \cdot \frac{d}{n} = C^2 d \log\left(\frac{nd}{\delta}\right) I_d$
- Note that  $X^T X + I_d$  and  $Z$  are real symmetric, with non-negative eigenvalues
- Claim: for all vectors  $y$ , we have:  $y^T X^T X y + y^T y \leq y^T Z y$
- Proof:  $y^T X^T X y + y^T y = n \sum_i y^T v_i^T v_i y |v_i|_2^2 = n \sum_i \langle v_i, y \rangle^2 |v_i|_2^2$  and
 
$$y^T Z y = n \sum_i y^T v_i^T v_i y C^2 \log\left(\frac{nd}{\delta}\right) \cdot \frac{d}{n} = n \sum_i \langle v_i, y \rangle^2 C^2 \log\left(\frac{nd}{\delta}\right)$$
- Hence,  $|E[X^T X]|_2 \leq |E[X^T X] + I_d|_2 + |I_d|_2 = |E[X^T X + I_d]|_2 + 1$ 

$$\leq |Z|_2 + 1 \leq C^2 d \log\left(\frac{nd}{\delta}\right) + 1$$
- Hence,  $|E[X^T X]|_2 = O\left(d \log\left(\frac{nd}{\delta}\right)\right)$

# Matrix Chernoff Bound

---

- Hence,  $|E[X^T X]|_2 = O\left(d \log\left(\frac{nd}{\delta}\right)\right)$
- Recall: (Matrix Chernoff) Let  $X_1, \dots, X_s$  be independent copies of a symmetric random matrix  $X \in \mathbb{R}^{d \times d}$  with  $E[X] = 0$ ,  $|X|_2 \leq \gamma$ , and  $|E[X^T X]|_2 \leq \sigma^2$ . Let  $W = \frac{1}{s} \sum_{i \in [s]} X_i$ . For any  $\epsilon > 0$ ,  $\Pr[|W|_2 > \epsilon] \leq 2d \cdot e^{-s\epsilon^2/(\sigma^2 + \frac{\gamma\epsilon}{3})}$

$$\Pr\left[|I_d - (\text{PHDA})^T(\text{PHDA})|_2 > \epsilon\right] \leq 2d \cdot e^{-s\epsilon^2/(\Theta(d \log(\frac{nd}{\delta})))}$$

- Set  $s = d \log\left(\frac{nd}{\delta}\right) \frac{\log(\frac{d}{\delta})}{\epsilon^2}$ , to make this probability less than  $\frac{\delta}{2}$

# SRHT Wrapup

---

- Have shown  $\|I_d - (\text{PHDA})^T(\text{PHDA})\|_2 < \epsilon$  using Matrix Chernoff Bound and with  $s = d \log\left(\frac{nd}{\delta}\right) \frac{\log\left(\frac{d}{\delta}\right)}{\epsilon^2}$  samples
- Implies for every unit vector  $x$ ,  
$$|1 - |\text{PHDA}x|_2^2| = |x^T x - x(\text{PHDA})^T(\text{PHDA})x| < \epsilon,$$
so  $|\text{PHDA}x|_2^2 \in 1 \pm \epsilon$  for all unit vectors  $x$
- Considering the column span of  $A$  adjoined with  $b$ , we can again solve the regression problem
- The time for regression is now only  $O(nd \log n) + \text{poly}\left(\frac{d \log(n)}{\epsilon}\right)$ . Nearly optimal in matrix dimensions ( $n \gg d$ )



# Faster Subspace Embeddings S [CW,MM,NN]

---

- CountSketch matrix
- Define  $k \times n$  matrix  $S$ , for  $k = O(d^2/\epsilon^2)$
- $S$  is really sparse: single randomly chosen non-zero entry per column

$$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 & 0 & -1 & 0 \\ 0 & -1 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Can compute  $S \cdot A$  in  $\text{nnz}(A)$  time!

- $\text{nnz}(A)$  is number of non-zero entries of  $A$

# Simple Proof [Nguyen]

---

- Can assume columns of  $A$  are orthonormal
- Suffices to show  $|SAx|_2 = 1 \pm \varepsilon$  for all unit  $x$ 
  - For regression, apply  $S$  to  $[A, b]$
- $SA$  is a  $2d^2/\varepsilon^2 \times d$  matrix
- Suffices to show  $\|A^T S^T SA - I\|_2 \leq \|A^T S^T SA - I\|_F \leq \varepsilon$
- Matrix product result shown below:  
$$\Pr[|CS^TSD - CD|_F^2 \leq [3/(\delta(\# \text{ rows of } S))] * |C|_F^2 |D|_F^2] \geq 1 - \delta$$
- Set  $C = A^T$  and  $D = A$ .
- Then  $|A|_F^2 = d$  and  $(\# \text{ rows of } S) = 3 d^2/(\delta\varepsilon^2)$

# Matrix Product Result [Kane, Nelson]

---

- Show:  $\Pr[|CS^TSD - CD|_F^2 \leq [3/(\delta(\# \text{ rows of } S))] * |C|_F^2 |D|_F^2] \geq 1 - \delta$
- (JL Property) A distribution on matrices  $S \in \mathbb{R}^{k \times n}$  has the  $(\epsilon, \delta, \ell)$ -JL moment property if for all  $x \in \mathbb{R}^n$  with  $|x|_2 = 1$ ,
 
$$E_S \left| |Sx|_2^2 - 1 \right|^\ell \leq \epsilon^\ell \cdot \delta$$
- (From vectors to matrices) For  $\epsilon, \delta \in \left(0, \frac{1}{2}\right)$ , let  $D$  be a distribution on matrices  $S$  with  $k$  rows and  $n$  columns that satisfies the  $(\epsilon, \delta, \ell)$ -JL moment property for some  $\ell \geq 2$ . Then for  $A, B$  matrices with  $n$  rows,

$$\Pr_S \left[ |A^T S^T S B - A^T B|_F \geq 3 \epsilon |A|_F |B|_F \right] \leq \delta$$

# From Vectors to Matrices

---

- (From vectors to matrices) For  $\epsilon, \delta \in \left(0, \frac{1}{2}\right)$ , let  $D$  be a distribution on matrices  $S$  with  $k$  rows and  $n$  columns that satisfies the  $(\epsilon, \delta, \ell)$ -JL moment property for some  $\ell \geq 2$ . Then for  $A, B$  matrices with  $n$  rows,

$$\Pr_S \left[ \left| A^T S^T S B - A^T B \right|_F \geq 3 \epsilon |A|_F |B|_F \right] \leq \delta$$

- Proof: For a random scalar  $X$ , let  $|X|_p = (E|X|^p)^{1/p}$ 
  - Sometimes consider  $X = |T|_F$  for a random matrix  $T$
  - $\| |T|_F \|_p = \left( E \left[ |T|_F^p \right] \right)^{1/p}$
- Can show  $|\cdot|_p$  is a norm
  - Minkowski's Inequality:  $|X + Y|_p \leq |X|_p + |Y|_p$
- For unit vectors  $x, y$ , need to bound  $|\langle Sx, Sy \rangle - \langle x, y \rangle|_\ell$

# From Vectors to Matrices

---

- For unit vectors  $x, y$ ,  $|\langle Sx, Sy \rangle - \langle x, y \rangle|_\ell$

$$= \frac{1}{2} |(|Sx|_2^2 - 1) + (|Sy|_2^2 - 1) - (|S(x-y)|_2^2 - |x-y|_2^2)|_\ell$$

$$\leq \frac{1}{2} (||Sx|_2^2 - 1|_\ell + ||Sy|_2^2 - 1|_\ell + ||S(x-y)|_2^2 - |x-y|_2^2|_\ell)$$

$$\leq \frac{1}{2} (\epsilon \cdot \delta^\ell + \epsilon \cdot \delta^\ell + |x-y|_2^2 \epsilon \cdot \delta^\ell)$$

$$\leq 3 \epsilon \cdot \delta^\ell$$
- By linearity, for arbitrary  $x, y$ ,  $\frac{|\langle Sx, Sy \rangle - \langle x, y \rangle|_\ell}{|x|_2 |y|_2} \leq 3 \epsilon \cdot \delta^\ell$
- Suppose  $A$  has  $d$  columns and  $B$  has  $e$  columns. Let the columns of  $A$  be  $A_1, \dots, A_d$  and the columns of  $B$  be  $B_1, \dots, B_e$
- Define  $X_{i,j} = \frac{1}{|A_i|_2 |B_j|_2} \cdot (\langle SA_i, SB_j \rangle - \langle A_i, B_j \rangle)$
- $$|A^T S^T S B - A^T B|_F^2 = \sum_i \sum_j |A_i|_2^2 \cdot |B_j|_2^2 X_{i,j}^2$$

# From Vectors to Matrices

---

- Have shown: for arbitrary  $x, y$ ,  $\frac{|\langle Sx, Sy \rangle - \langle x, y \rangle|}{\|x\|_2 \|y\|_2} \leq 3\epsilon \cdot \delta^{\frac{1}{\ell}}$
- For  $X_{i,j} = \frac{1}{\|A_i\|_2 \|B_j\|_2} \cdot (\langle SA_i, SB_j \rangle - \langle A_i, B_j \rangle)$ :  $\|A^T S^T S B - A^T B\|_F^2 = \sum_i \sum_j \|A_i\|_2^2 \cdot \|B_j\|_2^2 X_{i,j}^2$
- $$\begin{aligned} \|A^T S^T S B - A^T B\|_F^2 &= \left| \sum_i \sum_j \|A_i\|_2^2 \cdot \|B_j\|_2^2 X_{i,j}^2 \right|_{\ell/2} \\ &\leq \sum_i \sum_j \|A_i\|_2^2 \cdot \|B_j\|_2^2 |X_{i,j}|_{\ell/2}^2 \\ &= \sum_i \sum_j \|A_i\|_2^2 \cdot \|B_j\|_2^2 |X_{i,j}|_{\ell}^2 \\ &\leq \left(3\epsilon\delta^{\frac{1}{\ell}}\right)^2 \sum_i \sum_j \|A_i\|_2^2 \|B_j\|_2^2 \\ &= \left(3\epsilon\delta^{\frac{1}{\ell}}\right)^2 \|A\|_F^2 \|B\|_F^2 \end{aligned}$$
- Since  $E \left[ \|A^T S^T S B - A^T B\|_F^{\ell} \right] = \left\| \|A^T S^T S B - A^T B\|_F^2 \right\|_{\ell/2}^{\ell/2}$ , by Markov's inequality,
- $\Pr \left[ \|A^T S^T S B - A^T B\|_F > 3\epsilon \|A\|_F \|B\|_F \right] \leq \left( \frac{1}{3\epsilon \|A\|_F \|B\|_F} \right)^{\ell} E \left[ \|A^T S^T S B - A^T B\|_F^{\ell} \right] \leq \delta$

# Result for Vectors

---

- Show:  $\Pr[|CS^TSD - CD|_F^2 \leq [3/(\delta(\# \text{ rows of } S))] * |C|_F^2 |D|_F^2] \geq 1 - \delta$

- (JL Property) A distribution on matrices  $S \in \mathbb{R}^{k \times n}$  has the  $(\epsilon, \delta, \ell)$ -JL moment property if for all  $x \in \mathbb{R}^n$  with  $|x|_2 = 1$ ,

$$E_S \left| |Sx|_2^2 - 1 \right|^\ell \leq \epsilon^\ell \cdot \delta$$

- (From vectors to matrices) For  $\epsilon, \delta \in \left(0, \frac{1}{2}\right)$ , let  $D$  be a distribution on matrices  $S$  with  $k$  rows and  $n$  columns that satisfies the  $(\epsilon, \delta, \ell)$ -JL moment property for some  $\ell \geq 2$ . Then for  $A, B$  matrices with  $n$  rows,

$$\Pr_S \left[ \left| A^T S^T S B - A^T B \right|_F \geq 3 \epsilon |A|_F |B|_F \right] \leq \delta$$

- Just need to show that the CountSketch matrix  $S$  satisfies JL property and bound the number  $k$  of rows

# CountSketch Satisfies the JL Property

---

- (JL Property) A distribution on matrices  $S \in \mathbb{R}^{k \times n}$  has the  $(\epsilon, \delta, \ell)$ -JL moment property if for all  $x \in \mathbb{R}^n$  with  $\|x\|_2 = 1$ ,

$$\mathbb{E}_S \left| \|Sx\|_2^2 - 1 \right|^\ell \leq \epsilon^\ell \cdot \delta$$

- We show this property holds with  $\ell = 2$ . First, let us consider  $\ell = 1$

- For CountSketch matrix  $S$ , let

- $h: [n] \rightarrow [k]$  be a 2-wise independent hash function
- $\sigma: [n] \rightarrow \{-1, 1\}$  be a 4-wise independent hash function

- Let  $\delta(E) = 1$  if event  $E$  holds, and  $\delta(E) = 0$  otherwise

- $$\begin{aligned} \mathbb{E}[\|Sx\|_2^2] &= \sum_{j \in [k]} \mathbb{E} \left[ \left( \sum_{i \in [n]} \delta(h(i) = j) \sigma_i x_i \right)^2 \right] \\ &= \sum_{j \in [k]} \sum_{i_1, i_2 \in [n]} \mathbb{E} [\delta(h(i_1) = j) \delta(h(i_2) = j) \sigma_{i_1} \sigma_{i_2}] x_{i_1} x_{i_2} \\ &= \sum_{j \in [k]} \sum_{i \in [n]} \mathbb{E} [\delta(h(i) = j)^2] x_i^2 \\ &= \left( \frac{1}{k} \right) \sum_{j \in [k]} \sum_{i \in [n]} x_i^2 = \|x\|_2^2 \end{aligned}$$



# CountSketch Satisfies the JL Property

---

- $E[|Sx|_2^4] = E[\sum_{j \in [k]} \sum_{j' \in [k]} (\sum_{i \in [n]} \delta(h(i) = j) \sigma_i x_i)^2 (\sum_{i' \in [n]} \delta(h(i') = j') \sigma_{i'} x_{i'})^2] =$
- $\sum_{j_1, j_2, i_1, i_2, i_3, i_4} E[\sigma_{i_1} \sigma_{i_2} \sigma_{i_3} \sigma_{i_4} \delta(h(i_1) = j_1) \delta(h(i_2) = j_1) \delta(h(i_3) = j_2) \delta(h(i_4) = j_2)] x_{i_1} x_{i_2} x_{i_3} x_{i_4}$
- We must be able to partition  $\{i_1, i_2, i_3, i_4\}$  into equal pairs
- Suppose  $i_1 = i_2 = i_3 = i_4$ . Then necessarily  $j_1 = j_2$ . Obtain  $\sum_j \frac{1}{k} \sum_i x_i^4 = |x|_4^4$
- Suppose  $i_1 = i_2$  and  $i_3 = i_4$  but  $i_1 \neq i_3$ . Then get  $\sum_{j_1, j_2, i_1, i_3} \frac{1}{k^2} x_{i_1}^2 x_{i_3}^2 = |x|_2^4 - |x|_4^4$
- Suppose  $i_1 = i_3$  and  $i_2 = i_4$  but  $i_1 \neq i_2$ . Then necessarily  $j_1 = j_2$ . Obtain  $\sum_j \frac{1}{k^2} \sum_{i_1, i_2} x_{i_1}^2 x_{i_2}^2 \leq \frac{1}{k} |x|_4^4$ . Obtain same bound if  $i_1 = i_4$  and  $i_2 = i_3$ .
- Hence,  $E[|Sx|_2^4] \in [ |x|_2^4, |x|_2^4 (1 + \frac{2}{k}) ] = [1, 1 + \frac{2}{k}]$
- So,  $E_S ||Sx|_2^2 - 1|^2 \leq (1 + \frac{2}{k}) - 2 + 1 = \frac{2}{k}$ . Setting  $k = \frac{1}{2\epsilon^2 \delta}$  finishes the proof

# Where are we?

---

- (JL Property) A distribution on matrices  $S \in \mathbb{R}^{k \times n}$  has the  $(\epsilon, \delta, \ell)$ -JL moment property if for all  $x \in \mathbb{R}^n$  with  $\|x\|_2 = 1$ ,

$$\mathbb{E}_S \left| \|Sx\|_2^2 - 1 \right|^\ell \leq \epsilon^\ell \cdot \delta$$

- (From vectors to matrices) For  $\epsilon, \delta \in \left(0, \frac{1}{2}\right)$ , let  $D$  be a distribution on matrices  $S$  with  $k$  rows and  $n$  columns that satisfies the  $(\epsilon, \delta, \ell)$ -JL moment property for some  $\ell \geq 2$ . Then for  $A, B$  matrices with  $n$  rows,

$$\Pr_S \left[ \left| \|A^T S^T S B\|_F^2 - \|A^T B\|_F^2 \right| \geq 3 \epsilon^2 \|A\|_F^2 \|B\|_F^2 \right] \leq \delta$$

- We showed CountSketch has the JL property with  $\ell = 2$ , and  $k = \frac{2}{\epsilon^2 \delta}$

- Matrix product result we wanted was:

$$\Pr \left[ \|CS^TSD - CD\|_F^2 \leq (3/(\delta k)) * \|C\|_F^2 \|D\|_F^2 \right] \geq 1 - \delta$$

- We are now down with the proof CountSketch is a subspace embedding

# Course Outline

---

- Subspace embeddings and least squares regression
  - Gaussian matrices
  - Subsampled Randomized Hadamard Transform
  - CountSketch
- Affine embeddings
  - Application to low rank approximation
- High precision regression
- Leverage score sampling

# Affine Embeddings

---

- Want to solve  $\min_X |AX - B|_F^2$ , A is tall and thin with d columns, but B has a large number of columns
- Can't directly apply subspace embeddings
- Let's try to show  $|SAX - SB|_F = (1 \pm \epsilon)|AX - B|_F$  for all X and see what properties we need of S
- Can assume A has orthonormal columns
- Let  $B^* = AX^* - B$ , where  $X^*$  is the optimum
- $$|S(AX - B)|_F^2 - |SB^*|_F^2 = |SA(X - X^*) + S(AX^* - B)|_F^2 - |SB^*|_F^2$$

$$= |SA(X - X^*)|_F^2 - 2\text{tr}[(X - X^*)^T A^T S^T SB^*]$$

$$\in |SA(X - X^*)|_F^2 \pm 2|X - X^*|_F |A^T S^T SB^*|_F \text{ (use } \text{tr}(CD) \leq |C|_F |D|_F \text{)}$$

$$\in |SA(X - X^*)|_F^2 \pm 2\epsilon |X - X^*|_F |B^*|_F \text{ (if we have approx. matrix product)}$$

$$\in |A(X - X^*)|_F^2 \pm \epsilon (|A(X - X^*)|_F^2 + 2|X - X^*|_F |B^*|) \text{ (subspace embedding for A)}$$

# Affine Embeddings

---

- We have  $|S(AX - B)|_F^2 - |SB^*|_F^2 \in |A(X - X^*)|_F^2 \pm \epsilon(|A(X - X^*)|_F^2 + 2|X - X^*|_F|B^*|)$
- Normal equations imply that  $|AX - B|_F^2 = |A(X - X^*)|_F^2 + |B^*|_F^2$
- $|S(AX - B)|_F^2 - |SB^*|_F^2 - (|AX - B|_F^2 - |B^*|_F^2) \in \epsilon(|A(X - X^*)|_F^2 + 2|X - X^*|_F|B^*|_F) \in \pm\epsilon(|A(X - X^*)|_F + |B^*|_F)^2 \in \pm 2\epsilon(|A(X - X^*)|_F^2 + |B^*|_F^2) = \pm 2\epsilon|AX - B|_F^2$
- $|SB^*|_F^2 = (1 \pm \epsilon)|B^*|_F^2$  (this holds with constant probability<sup>53</sup>)

# Affine Embeddings

---

- Know:  $|S(AX - B)|_F^2 - |SB^*|_F^2 - (|AX - B|_F^2 - |B^*|_F^2) \in \pm 2\epsilon|AX - B|_F^2$
- Know:  $|SB^*|_F^2 = (1 \pm \epsilon)|B^*|_F^2$
- $|S(AX - B)|_F^2 = (1 \pm 2\epsilon)|AX - B|_F^2 + \epsilon|B^*|_F^2$   
 $= (1 \pm 3\epsilon)|AX - B|_F^2$
- Completes proof of affine embedding!

# Affine Embeddings: Missing Proofs

---

- Claim:  $|A + B|_F^2 = |A|_F^2 + |B|_F^2 + 2\text{Tr}(A^T B)$

- Proof:  $|A + B|_F^2 = \sum_i |A_i + B_i|_2^2$

$$= \sum_i |A_i|_2^2 + \sum_i |B_i|_2^2 + 2\langle A_i, B_i \rangle$$

$$= |A|_F^2 + |B|_F^2 + 2\text{Tr}(A^T B)$$

# Affine Embeddings: Missing Proofs

---

- Claim:  $\text{Tr}(AB) \leq |A|_F |B|_F$
- Proof:  $\text{Tr}(AB) = \sum_i \langle A^i, B_i \rangle$  for rows  $A^i$  and columns  $B_i$   
 $\leq \sum_i |A^i|_2 |B_i|_2$  by Cauchy-Schwarz for each  $i$   
 $\leq \left( \sum_i |A^i|_2^2 \right)^{\frac{1}{2}} \left( \sum_i |B_i|_2^2 \right)^{\frac{1}{2}}$  another Cauchy-Schwarz  
 $= |A|_F |B|_F$



# Affine Embeddings: Homework Proof

---

- Claim:  $\|SB^*\|_F^2 = (1 \pm \epsilon)\|B^*\|_F^2$  with constant probability if CountSketch matrix  $S$  has  $k = O(\frac{1}{\epsilon^2})$  rows
- Proof:
- $\|SB^*\|_F^2 = \sum_i \|SB_i^*\|_2^2$
- By our analysis for CountSketch and linearity of expectation,  $E[\|SB^*\|_F^2] = \sum_i E[\|SB_i^*\|_2^2] = \|B^*\|_F^2$
- $E[\|SB^*\|_F^4] = \sum_{i,j} E[\|SB_i^*\|_2^2 \|SB_j^*\|_2^2]$
- By our CountSketch analysis,  $E[\|SB_i^*\|_2^4] \leq \|B_i^*\|_2^4 (1 + \frac{2}{k})$
- For cross terms see Lemma 40 in [CW13]

# Low rank approximation

---

- A is an  $n \times d$  matrix
  - Think of  $n$  points in  $\mathbb{R}^d$
- E.g., A is a customer-product matrix
  - $A_{i,j}$  = how many times customer  $i$  purchased item  $j$
- A is typically well-approximated by low rank matrix
  - E.g., high rank because of noise
- **Goal:** find a low rank matrix approximating A
  - Easy to store, data more interpretable

# What is a good low rank approximation?

---

## Singular Value Decomposition (SVD)

Any matrix  $A = U \cdot \Sigma \cdot V$

- $U$  has orthonormal columns
  - $\Sigma$  is diagonal with non-increasing positive entries down the diagonal
  - $V$  has orthonormal rows
- 
- Rank-k approximation:  $A_k = U_k \cdot \Sigma_k \cdot V_k$ 
    - rows of  $V_k$  are the top k principal components

$$\begin{pmatrix} \mathbf{A} \end{pmatrix} = \begin{pmatrix} \mathbf{U}_k \end{pmatrix} \begin{pmatrix} \Sigma_k \end{pmatrix} \begin{pmatrix} \mathbf{V}_k \end{pmatrix} + \begin{pmatrix} \mathbf{E} \end{pmatrix}$$

# What is a good low rank approximation?

---

$$A_k = \operatorname{argmin}_{\text{rank } k \text{ matrices } B} \|A-B\|_F$$

$$(\|C\|_F = (\sum_{i,j} C_{i,j}^2)^{1/2})$$

Computing  $A_k$  exactly is expensive

$$\begin{pmatrix} \mathbf{A} \end{pmatrix} = \begin{pmatrix} \mathbf{U}_k \end{pmatrix} \begin{pmatrix} \Sigma_k \end{pmatrix} \begin{pmatrix} \mathbf{V}_k \end{pmatrix} + \begin{pmatrix} \mathbf{E} \end{pmatrix}$$

# Low rank approximation

---

- **Goal:** output a rank  $k$  matrix  $A'$ , so that

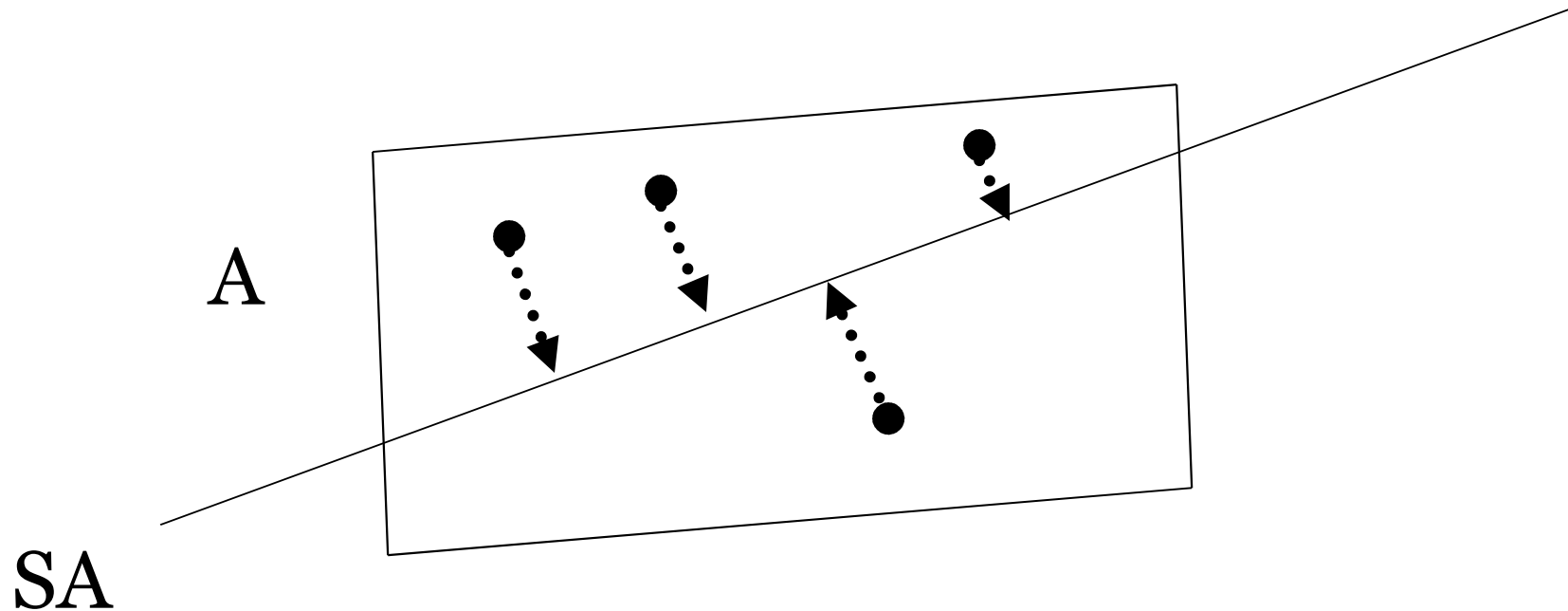
$$\|A-A'\|_F \leq (1+\varepsilon) \|A-A_k\|_F$$

- Can do this in  $\text{nnz}(A) + (n+d) \cdot \text{poly}(k/\varepsilon)$  time [S,CW]
  - $\text{nnz}(A)$  is number of non-zero entries of  $A$

# Solution to low-rank approximation [S]

---

- Given  $n \times d$  input matrix  $A$
- Compute  $S^*A$  using a random matrix  $S$  with  $k/\epsilon \ll n$  rows.  $S^*A$  takes random linear combinations of rows of  $A$



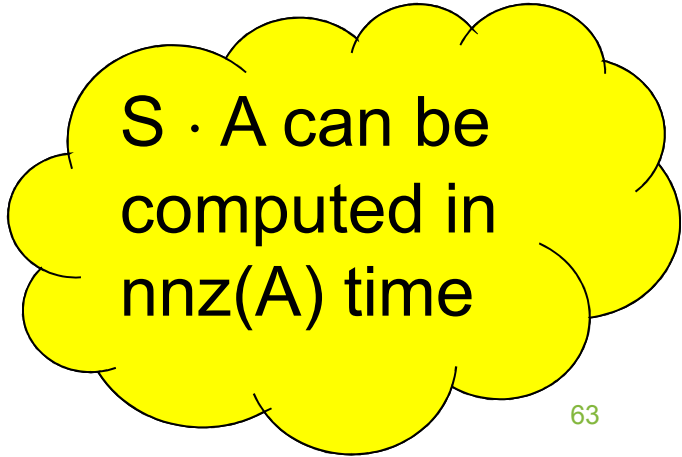
- Project rows of  $A$  onto  $SA$ , then find best rank- $k$  approximation to points inside of  $SA$ .

# What is the matrix $S$ ?

---

- $S$  can be a  $k/\epsilon \times n$  matrix of i.i.d. normal random variables
- [S]  $S$  can be a  $k/\epsilon \times n$  Fast Johnson Lindenstrauss Matrix
  - Uses Fast Fourier Transform
- [CW]  $S$  can be a  $\text{poly}(k/\epsilon) \times n$  CountSketch matrix

$$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 & 0 & -1 & 0 \\ 0 & -1 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$



$S \cdot A$  can be computed in  $\text{nnz}(A)$  time

# Why do these Matrices Work?

---

- Consider the regression problem  $\min_X |A_k X - A|_F$
- Write  $A_k = WY$ , where  $W$  is  $n \times k$  and  $Y$  is  $k \times d$
- Let  $S$  be an affine embedding
- Then  $|SA_k X - SA|_F = (1 \pm \epsilon) |A_k X - A|_F$  for all  $X$
- By normal equations,  $\operatorname{argmin}_X |SA_k X - SA|_F = (SA_k)^{-1} SA$
- So,  $|A_k (SA_k)^{-1} SA - A|_F \leq (1 + \epsilon) |A_k - A|_F$
- But  $A_k (SA_k)^{-1} SA$  is a rank- $k$  matrix in the row span of  $SA$ !
- **Let's formalize why the algorithm works now...**



# Why do these Matrices Work?

---

- $$\min_{\text{rank-}k X} \|XSA - A\|_F^2 \leq \|A_k(SA_k)^{-1}SA - A\|_F^2 \leq (1 + \epsilon) \|A - A_k\|_F^2$$
- By the normal equations,

$$\|XSA - A\|_F^2 = \|XSA - A(SA)^{-1}SA\|_F^2 + \|A(SA)^{-1}SA - A\|_F^2$$
- Hence,

$$\min_{\text{rank-}k X} \|XSA - A\|_F^2 = \|A(SA)^{-1}SA - A\|_F^2 + \min_{\text{rank-}k X} \|XSA - A(SA)^{-1}SA\|_F^2$$
- Can write  $SA = U\Sigma V^T$  in its SVD, where  $U, \Sigma$  are  $k \times k$  and  $V^T$  is  $k \times d$
- Then,

$$\begin{aligned} \min_{\text{rank-}k X} \|XSA - A(SA)^{-1}SA\|_F^2 &= \min_{\text{rank-}k X} \|XU\Sigma - A(SA)^{-1}U\Sigma\|_F^2 \\ &= \min_{\text{rank-}k Y} \|Y - A(SA)^{-1}U\Sigma\|_F^2 \end{aligned}$$
- Hence, we can just compute the SVD of  $A(SA)^{-1}U\Sigma$
- But how do we compute  $A(SA)^{-1}U\Sigma$  quickly?

# Caveat: projecting the points onto SA is slow

---

- Current algorithm:
  1. Compute  $S^*A$
  2. Project each of the rows onto  $S^*A$
  3. Find best rank-k approximation of projected points inside of rowspace of  $S^*A$

- Bottleneck is step 2

$$\min_{\text{rank-}k \times} |X(SA)R-AR|_F^2$$

Can solve with affine embeddings

- [CW] Approximate the projection
  - Fast algorithm for approximate regression

$$\min_{\text{rank-}k \times} |X(SA)-A|_F^2$$

- Want  $\text{nnz}(A) + (n+d) \cdot \text{poly}(k/\epsilon)$  time

# Using Affine Embeddings

---

- We know we can just output  $\arg \min_{\text{rank-}k X} \|XSA - A\|_F^2$

- Choose an affine embedding R:

$$\|XSAR - AR\|_F^2 = (1 \pm \epsilon) \|XSA - A\|_F^2 \text{ for all } X$$

- Note: we can compute AR and SAR in  $\text{nnz}(A)$  time

- Can just solve  $\min_{\text{rank-}k X} \|XSAR - AR\|_F^2$

- $\min_{\text{rank-}k X} \|XSAR - AR\|_F^2 = \|AR(SAR)^-(SAR) - AR\|_F^2 + \min_{\text{rank-}k X} \|XSAR - AR(SAR)^-(SAR)\|_F^2$

- Compute  $\min_{\text{rank-}k Y} \|Y - AR(SAR)^-(SAR)\|_F^2$  using SVD which is  $(n + d)\text{poly}\left(\frac{k}{\epsilon}\right)$  time

- Necessarily,  $Y = XSAR$  for some  $X$ . Output  $Y(SAR)^-SA$  in factored form. We're done!

# Low Rank Approximation Summary

---

1. Compute SA
2. Compute SAR and AR
3. Compute  $\min_{\text{rank-}k Y} \|Y - AR(SAR)^{-1}(SAR)\|_F^2$  using SVD
4. Output  $Y(SAR)^{-1}SA$  in factored form

Overall time:  $\text{nnz}(A) + (n+d)\text{poly}(k/\epsilon)$

# Course Outline

---

- Subspace embeddings and least squares regression
  - Gaussian matrices
  - Subsampled Randomized Hadamard Transform
  - CountSketch
- Affine embeddings
  - Application to low rank approximation
- High precision regression
- Leverage score sampling

# High Precision Regression

---

- **Goal:** output  $x'$  for which  $|Ax'-b|_2 \leq (1+\varepsilon) \min_x |Ax-b|_2$  with high probability
- Our algorithms all have running time  $\text{poly}(d/\varepsilon)$
- **Goal:** Sometimes we want running time  $\text{poly}(d) \cdot \log(1/\varepsilon)$
- Want to make  $A$  well-conditioned
  - $\kappa(A) = \sup_{|x|_2=1} |Ax|_2 / \inf_{|x|_2=1} |Ax|_2$
- Lots of algorithms' time complexity depends on  $\kappa(A)$
- Use sketching to reduce  $\kappa(A)$  to  $O(1)$ !

# Small QR Decomposition

---

- Let  $S$  be a  $1 + \epsilon_0$  - subspace embedding for  $A$
- Compute  $SA$
- Compute QR-factorization,  $SA = QR^{-1}$
- Claim:  $\kappa(AR) = \frac{(1+\epsilon_0)}{1-\epsilon_0}$
- For all unit  $x$ ,  $(1 - \epsilon_0)|ARx|_2 \leq |SARx|_2 = 1$
- For all unit  $x$ ,  $(1 + \epsilon_0)|ARx|_2 \geq |SARx|_2 = 1$
- So  $\kappa(AR) = \sup_{|x|_2=1} |ARx|_2 / \inf_{|x|_2=1} |ARx|_2 \leq \frac{1+\epsilon_0}{1-\epsilon_0}$

# Finding a Constant Factor Solution

---

- Let  $S$  be a  $1 + \epsilon_0$  - subspace embedding for  $AR$
- Solve  $x_0 = \underset{x}{\operatorname{argmin}} |SARx - Sb|_2$
- Time to compute  $AR$  and  $x_0$  is  $\operatorname{nnz}(A) + \operatorname{poly}(d)$  for constant  $\epsilon_0$
- $x_{m+1} \leftarrow x_m + R^T A^T (b - ARx_m)$
- $$\begin{aligned} AR(x_{m+1} - x^*) &= AR(x_m + R^T A^T (b - ARx_m) - x^*) \\ &= (AR - ARR^T A^T AR)(x_m - x^*) \\ &= U(\Sigma - \Sigma^3)V^T(x_m - x^*), \end{aligned}$$
where  $AR = U \Sigma V^T$  is the SVD of  $AR$
- $|AR(x_{m+1} - x^*)|_2 = |(\Sigma - \Sigma^3)V^T(x_m - x^*)|_2 = O(\epsilon_0)|AR(x_m - x^*)|_2$
- $|ARx_m - b|_2^2 = |AR(x_m - x^*)|_2^2 + |ARx^* - b|_2^2$



# Course Outline

---

- Subspace embeddings and least squares regression
  - Gaussian matrices
  - Subsampled Randomized Hadamard Transform
  - CountSketch
- Affine embeddings
  - Application to low rank approximation
- High precision regression
- Leverage score sampling

# Leverage Score Sampling

---

- This is another subspace embedding, but it is based on sampling!
  - If  $A$  has sparse rows, then  $SA$  has sparse rows!
- Let  $A = U \Sigma V^T$  be an  $n \times d$  matrix with rank  $d$ , written in its SVD
- Define the  $i$ -th leverage score  $\ell(i)$  of  $A$  to be  $\|U_{i,*}\|_2^2$
- What is  $\sum_i \ell(i)$ ?
  - Let  $(q_1, \dots, q_n)$  be a distribution with  $q_i \geq \frac{\beta \ell(i)}{d}$ , where  $\beta$  is a parameter
- Define sampling matrix  $S = D \cdot \Omega^T$ , where  $D$  is  $k \times k$  and  $\Omega$  is  $n \times k$ 
  - $\Omega$  is a sampling matrix, and  $D$  is a rescaling matrix
  - For each column  $j$  of  $\Omega$ ,  $D$ , independently, and with replacement, pick a row index  $i$  in  $[n]$  with probability  $q_i$ , and set  $\Omega_{i,j} = 1$  and  $D_{j,j} = (q_i k)^{.5}$

# Leverage Score Sampling

---

- Note: leverage scores do not depend on choice of orthonormal basis  $U$  for columns of  $A$
- Indeed, let  $U$  and  $U'$  be two such orthonormal bases
- Claim:  $\|e_i U\|_2^2 = \|e_i U'\|_2^2$  for all  $i$
- Proof: Since both  $U$  and  $U'$  have column space equal to that of  $A$ , we have  $U = U'Z$  for change of basis matrix  $Z$
- Since  $U$  and  $U'$  each have orthonormal columns,  $Z$  is a rotation matrix (orthonormal rows and columns)
- Then  $\|e_i U\|_2^2 = \|e_i U'Z\|_2^2 = \|e_i U'\|_2^2$

# Leverage Score Sampling gives a Subspace Embedding

---

- Want to show for  $S = D \cdot \Omega^T$ , that  $|SAx|_2^2 = (1 \pm \epsilon)|Ax|_2^2$  for all  $x$
- Writing  $A = U \Sigma V^T$  in its SVD, this is equivalent to showing  $|SUy|_2^2 = (1 \pm \epsilon)|Uy|_2^2 = (1 \pm \epsilon)|y|_2^2$  for all  $y$
- As usual, we can just show with high probability,  $|U^T S^T S U - I|_2 \leq \epsilon$
- How can we analyze  $U^T S^T S U$ ?
- (Matrix Chernoff) Let  $X_1, \dots, X_k$  be independent copies of a symmetric random matrix  $X \in \mathbb{R}^{d \times d}$  with  $E[X] = 0$ ,  $|X|_2 \leq \gamma$ , and  $|E[X^T X]|_2 \leq \sigma^2$ . Let  $W = \frac{1}{k} \sum_{j \in [k]} X_j$ . For any  $\epsilon > 0$ ,

$$\Pr[|W|_2 > \epsilon] \leq 2d \cdot e^{-k\epsilon^2 / (\sigma^2 + \frac{\gamma\epsilon}{3})}$$

(here  $|W|_2 = \sup \frac{|Wx|_2}{|x|_2}$ . Since  $W$  is symmetric,  $|W|_2 = \sup_{|x|_2=1} x^T W x$ .)

# Leverage Score Sampling gives a Subspace Embedding

---

- Let  $i(j)$  denote the index of the row of  $U$  sampled in the  $j$ -th trial
- Let  $X_j = I_d - \frac{U_{i(j)}^T U_{i(j)}}{q_{i(j)}}$ , where  $U_{i(j)}$  is the  $j$ -th sampled row of  $U$
- The  $X_j$  are independent copies of a symmetric matrix random variable
- $E[X_j] = I_d - \sum_i q_i \left( \frac{U_i^T U_i}{q_i} \right) = I_d - I_d = 0^d$
- $|X_j|_2 \leq |I_d|_2 + \frac{|U_{i(j)}^T U_{i(j)}|_2}{q_{i(j)}} \leq 1 + \max_i \frac{|U_i|_2^2}{q_i} \leq 1 + \frac{d}{\beta}$
- $$E[X^T X] = I_d - 2E \left[ \frac{U_{i(j)}^T U_{i(j)}}{q_{i(j)}} \right] + E \left[ \frac{U_{i(j)}^T U_{i(j)} U_{i(j)}^T U_{i(j)}}{q_{i(j)}^2} \right]$$

$$= \sum_i \frac{U_i^T U_i U_i^T U_i}{q(i)} - I_d \leq \left( \frac{d}{\beta} \right) \sum_i U_i^T U_i - I_d \leq \left( \frac{d}{\beta} - 1 \right) I_d,$$

where  $A \leq B$  means  $x^T A x \leq x^T B x$  for all  $x$
- Hence,  $|E[X^T X]|_2 \leq \frac{d}{\beta} - 1$

# Applying the Matrix Chernoff Bound

---

- (Matrix Chernoff) Let  $X_1, \dots, X_k$  be independent copies of a symmetric random matrix  $X \in \mathbb{R}^{d \times d}$  with  $E[X] = 0$ ,  $|X|_2 \leq \gamma$ , and  $|E[X^T X]|_2 \leq \sigma^2$ . Let  $W = \frac{1}{k} \sum_{j \in [k]} X_j$ . For any  $\epsilon > 0$ ,

$$\Pr[|W|_2 > \epsilon] \leq 2d \cdot e^{-k\epsilon^2 / (\sigma^2 + \frac{\gamma\epsilon}{3})}$$

(here  $|W|_2 = \sup_{|x|_2=1} |Wx|_2$ . Since  $W$  is symmetric,  $|W|_2 = \sup_{|x|_2=1} x^T W x$ .)

- $\gamma = 1 + \frac{d}{\beta}$ , and  $\sigma^2 = \frac{d}{\beta} - 1$

- $X_j = I_d - \frac{U_{i(j)}^T U_{i(j)}}{q_{i(j)}}$ , and recall how we generated  $S = D \cdot \Omega^T$ : For each column  $j$  of  $\Omega$ ,  $D$ , independently, and with replacement, pick a row index  $i$  in  $[n]$  with probability  $q_i$ , and set  $\Omega_{i,j} = 1$  and  $D_{j,j} = (q_i k)^{.5}$

- Implies  $W = I_d - U^T S^T S U$

- $\Pr[|I_d - U^T S^T S U|_2 > \epsilon] \leq 2d \cdot e^{-k\epsilon^2 \Theta(\frac{\beta}{d})}$ . Set  $k = \Theta(\frac{d \log d}{\beta \epsilon^2})$  and we're done. 78

# Fast Computation of Leverage Scores

---

- Naively, need to do an SVD to compute leverage scores
- Suppose we compute SA for a subspace embedding S
- Let  $SA = QR^{-1}$  be such that Q has orthonormal columns
- Set  $\ell'_i = |e_i AR|_2^2$
- Since AR has the same column span of A,  $AR = UT^{-1}$ 
  - $(1 - \epsilon)|ARx|_2 \leq |SARx|_2 = |x|_2$
  - $(1 + \epsilon)|ARx|_2 \geq |SARx|_2 = |x|_2$
  - $(1 \pm O(\epsilon))|x|_2 = |ARx|_2 = |UT^{-1}x|_2 = |T^{-1}x|_2,$
- $\ell_i = |e_i ART|_2^2 = (1 \pm O(\epsilon))|e_i AR|_2^2 = (1 \pm O(\epsilon))\ell'_i$
- But how do we compute AR? We want  $\text{nnz}(A)$  time

# Fast Computation of Leverage Scores

---

- $\ell_i = (1 \pm O(\epsilon))\ell_i'$
- Suffices to set this  $\epsilon$  to be a constant
- Set  $\ell_i' = |e_i A R|_2^2$ 
  - This takes too long
- Let  $G$  be a  $d \times O(\log n)$  matrix of i.i.d. normal random variables
  - For any vector  $z$ ,  $\Pr[|zG|_2^2 = \left(1 \pm \frac{1}{2}\right) |z|^2] \geq 1 - \frac{1}{n^2}$
- Instead set  $\ell_i' = |e_i A R G|_2^2$ .
  - Can compute in  $(\text{nnz}(A) + d^2) \log n$  time
- Can solve regression in  $\text{nnz}(A) \log n + \text{poly}(d(\log n)/\epsilon)$  time



# Course Outline

---

- Subspace embeddings and least squares regression
  - Gaussian matrices
  - Subsampled Randomized Hadamard Transform
  - CountSketch
- Affine embeddings
  - Application to low rank approximation
- High precision regression
- Leverage score sampling
- Distributed Low Rank Approximation

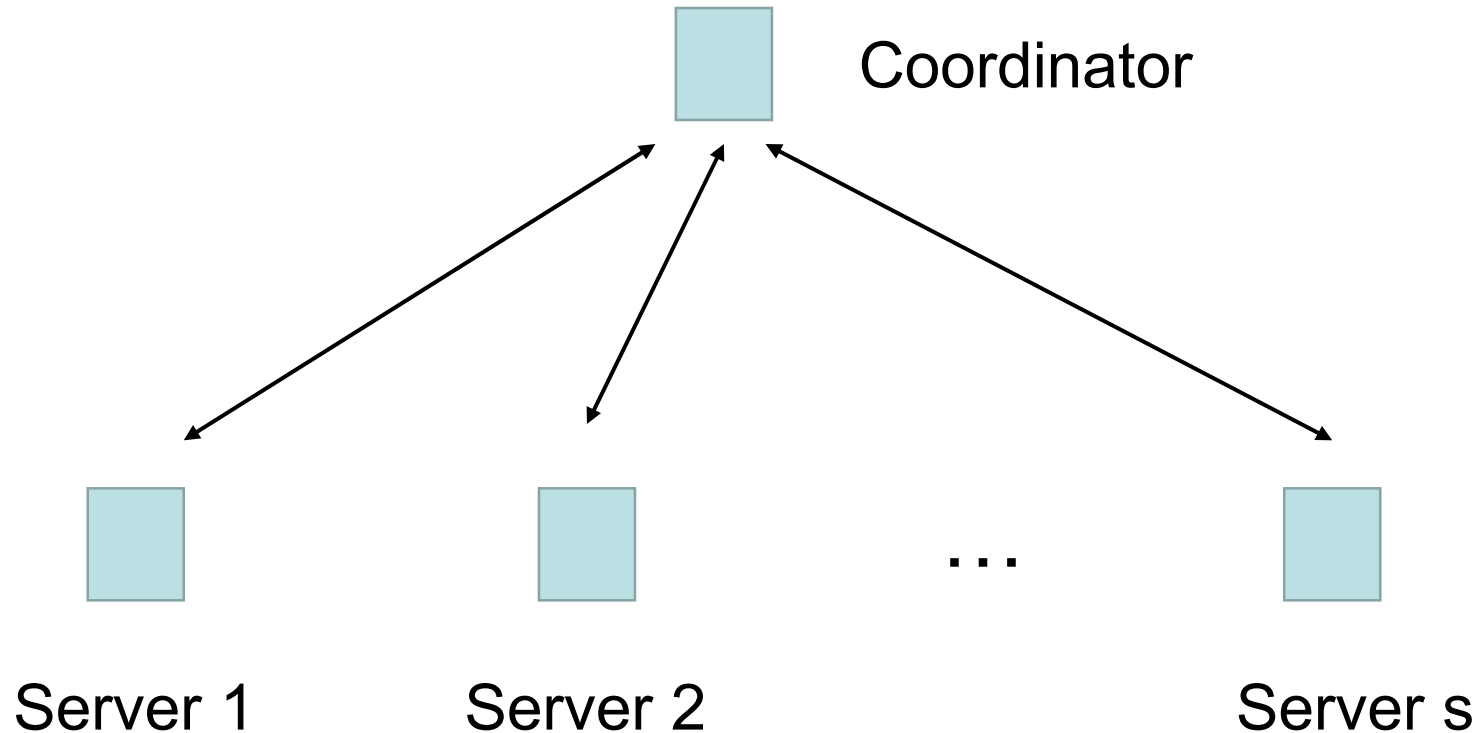
# Distributed low rank approximation

---

- *We have fast algorithms for low rank approximation, but can they be made to work in a distributed setting?*
- Matrix A distributed among s servers
- For  $t = 1, \dots, s$ , we get a customer-product matrix from the t-th shop stored in server t. Server t's matrix =  $A^t$
- Customer-product matrix  $A = A^1 + A^2 + \dots + A^s$ 
  - Model is called the **arbitrary partition model**
- More general than the **row-partition model** in which each customer shops in only one shop

# The Communication Model

---



- Each player talks only to a Coordinator via 2-way communication
- Can simulate arbitrary point-to-point communication up to factor of 2 (and an additive  $O(\log s)$  factor per message)

# Communication cost of low rank approximation

---

- **Input:**  $n \times d$  matrix  $A$  stored on  $s$  servers
  - Server  $t$  has  $n \times d$  matrix  $A^t$
  - $A = A^1 + A^2 + \dots + A^s$
  - Assume entries of  $A^t$  are  $O(\log(nd))$ -bit integers
- **Output:** Each server outputs the same  $k$ -dimensional space  $W$ 
  - $C = A^1 P_W + A^2 P_W + \dots + A^s P_W$ , where  $P_W$  is the projector onto  $W$
  - $|A-C|_F \leq (1+\epsilon)|A-A_k|_F$
  - Application:  $k$ -means clustering
- **Resources:** Minimize total communication and computation.  
Also want  $O(1)$  rounds and input sparsity time

# Work on Distributed Low Rank Approximation

---

- [FSS]: First protocol for the row-partition model.
  - $O(sdk/\epsilon)$  real numbers of communication
  - Don't analyze bit complexity (can be large)
  - SVD Running time, see also [BKLW]
- [KVW]:  $O(sdk/\epsilon)$  communication in arbitrary partition model
- [BWZ]:  $O(skd) + \text{poly}(sk/\epsilon)$  words of communication in arbitrary partition model. Input sparsity time
  - Matching  $\Omega(skd)$  words of communication lower bound
- Variants: kernel low rank approximation [BLSWX], low rank approximation of an implicit matrix [WZ], sparsity [BWZ]

# Outline of Distributed Protocols

---

- [FSS] protocol
- [KVW] protocol
- [BWZ] protocol

# Constructing a Coreset [FSS]

---

- Let  $A = U \Sigma V^T$  be its SVD
- Let  $m = k + k/\epsilon$
- Let  $\Sigma_m$  agree with  $\Sigma$  on the first  $m$  diagonal entries, and be 0 otherwise
- Claim: For all projection matrices  $Y=I-X$  onto  $(n-k)$ -dimensional subspaces,

$$|\Sigma_m V^T Y|_F^2 = (1 \pm \epsilon) |AY|_F^2 + c,$$

where  $c = |A - A_m|_F^2$  does not depend on  $Y$

- We can think of  $S$  as  $U_m^T$  so that  $SA = U_m^T U \Sigma V^T = \Sigma_m V^T$  is a sketch

# Constructing a Coreset

---

- Claim: For all projection matrices  $Y=I-X$  onto  $(n-k)$ -dimensional subspaces,

$$|\Sigma_m V^T Y|_F^2 + c = (1 \pm \epsilon) |AY|_F^2,$$

where  $c = |A - A_m|_F^2$  does not depend on  $Y$

- Proof:  $|AY|_F^2 = |U\Sigma_m V^T Y|_F^2 + |U(\Sigma - \Sigma_m)V^T Y|_F^2$   
 $\leq |\Sigma_m V^T Y|_F^2 + |A - A_m|_F^2 = |\Sigma_m V^T Y|_F^2 + c$

Also,  $|\Sigma_m V^T Y|_F^2 + |A - A_m|_F^2 - |AY|_F^2$

$$= |\Sigma_m V^T|_F^2 - |\Sigma_m V^T X|_F^2 + |A - A_m|_F^2 - |A|_F^2 + |AX|_F^2$$

$$= |AX|_F^2 - |\Sigma_m V^T X|_F^2$$

$$= |(\Sigma - \Sigma_m)V^T X|_F^2$$

$$\leq |(\Sigma - \Sigma_m)V^T|_2^2 \cdot |X|_F^2$$

$$\leq \sigma_{m+1}^2 k \leq \epsilon \sigma_{m+1}^2 (m - k + 1) \leq \epsilon \sum_{i \in \{k+1, \dots, m+1\}} \sigma_i^2 \leq \epsilon |A - A_k|_F^2$$



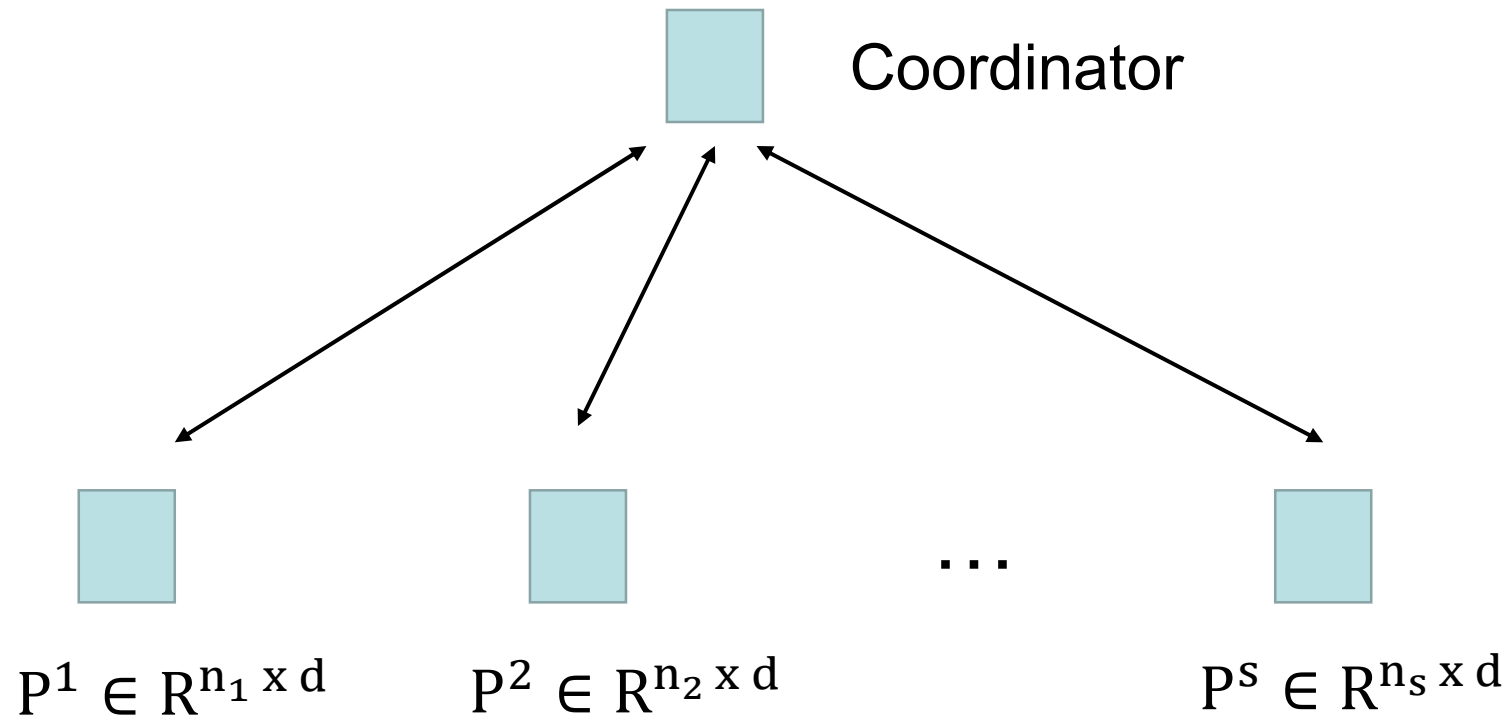
# Unions of Coresets

---

- Suppose we have matrices  $A^1, \dots, A^s$  and construct  $\Sigma_m^1 V^{T,1}, \Sigma_m^2 V^{T,2}, \dots, \Sigma_m^s V^{T,s}$  as in the previous slide, together with  $c_1, \dots, c_s$
- Then  $\sum_i |\Sigma_m^i V^{T,i} Y|_F^2 + c_i = (1 \pm \epsilon) |AY|_F^2$ , where  $A$  is the matrix formed by concatenating the rows of  $A^1, \dots, A^s$
- Let  $B$  be the matrix obtained by concatenating the rows of  $\Sigma_m^1 V^{T,1}, \Sigma_m^2 V^{T,2}, \dots, \Sigma_m^s V^{T,s}$
- Suppose we compute  $B = U \Sigma V^T$  and compute  $\Sigma_m V^T$  and  $|B - B_m|_F^2$
- Then  $|\Sigma_m V^T Y|_F^2 + c + \sum_i c_i = (1 \pm \epsilon) |BY|_F^2 + \sum_i c_i = (1 \pm O(\epsilon)) |AY|_F^2$
- So  $\Sigma_m V^T$  and the constant  $c + \sum_i c_i$  are a coreset for  $A$

# [FSS] Row-Partition Protocol

---



- Server  $t$  sends the top  $k/\varepsilon + k$  principal components of  $P^t$ , scaled by the top  $k/\varepsilon + k$  singular values  $\Sigma^t$ , together with  $c^t$
- Coordinator returns top  $k$  principal components of  $[\Sigma^1 V^1; \Sigma^2 V^2; \dots; \Sigma^s V^s]$

# [FSS] Row-Partition Protocol

---

[KVV] protocol  
will handle 2, 3,  
and 4

## Problems:

1.  $\text{sdk}/\epsilon$  real numbers of communication
2. bit complexity can be large
3. running time for SVDs [BLKW]
4. doesn't work in arbitrary partition model

*This is an SVD-based protocol. Maybe  
our random matrix techniques can  
improve communication just like they  
improved computation?*

# [KVW] Arbitrary Partition Model Protocol

---

- Inspired by the sketching algorithm presented earlier
- Let  $S$  be one of the  $k/\epsilon \times n$  random matrices discussed
  - $S$  can be generated pseudorandomly from small seed
  - Coordinator sends small seed for  $S$  to all servers
- Server  $t$  computes  $SA^t$  and sends it to Coordinator
- Coordinator sends  $\sum_{t=1}^s SA^t = SA$  to all servers
- There is a good  $k$ -dimensional subspace inside of  $SA$ . If we knew it,  $t$ -th server could output projection of  $A^t$  onto it

# [KVW] Arbitrary Partition Model Protocol

---

## Problems:

- Can't output projection of  $A^t$  onto SA since the rank is too large
- Could communicate this projection to the coordinator who could find a  $k$ -dimensional space, but communication depends on  $n$

# [KVW] Arbitrary Partition Model Protocol

---

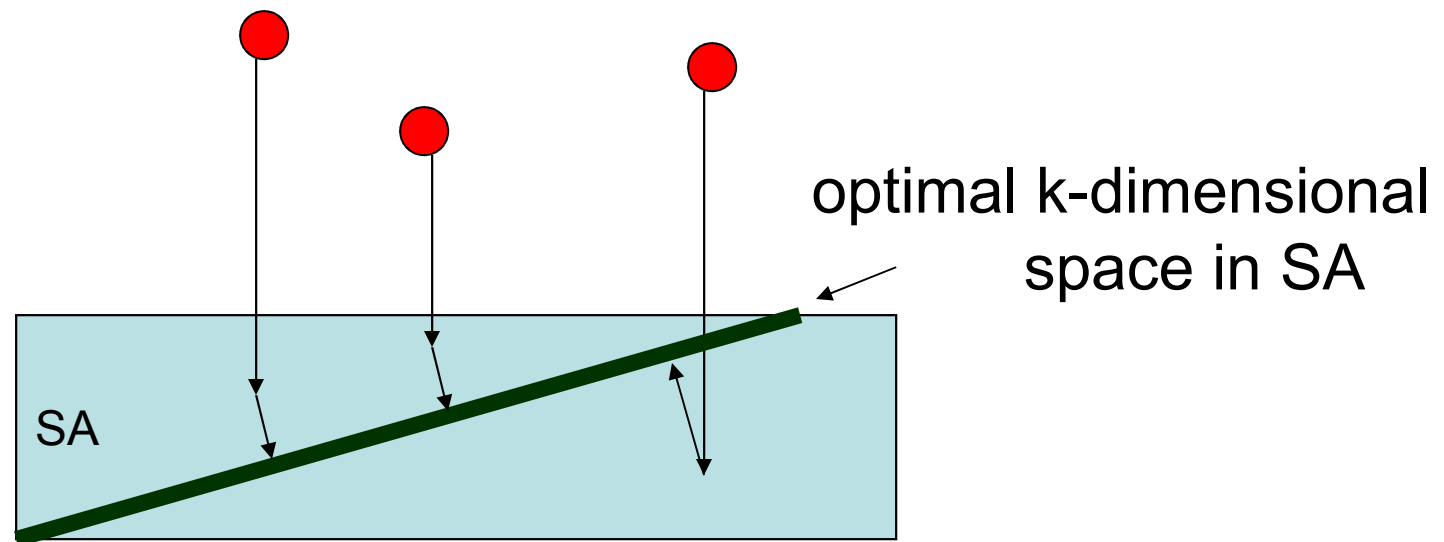
## Fix:

- Instead of projecting  $A$  onto  $SA$ , recall we can solve  $\min_{\text{rank-}k X} \|A(SA)^T XSA - A\|_F^2$
- Let  $T_1, T_2$  be affine embeddings, solve  $\min_{\text{rank-}k X} \|T_1 A(SA)^T XSA T_2 - T_1 A T_2\|_F^2$   
(optimization problem is small and has a closed form solution)
- Everyone can then compute  $XSA$  and then output  $k$  directions

# [KVW] protocol

---

- Phase 1:
- Learn the row space of SA

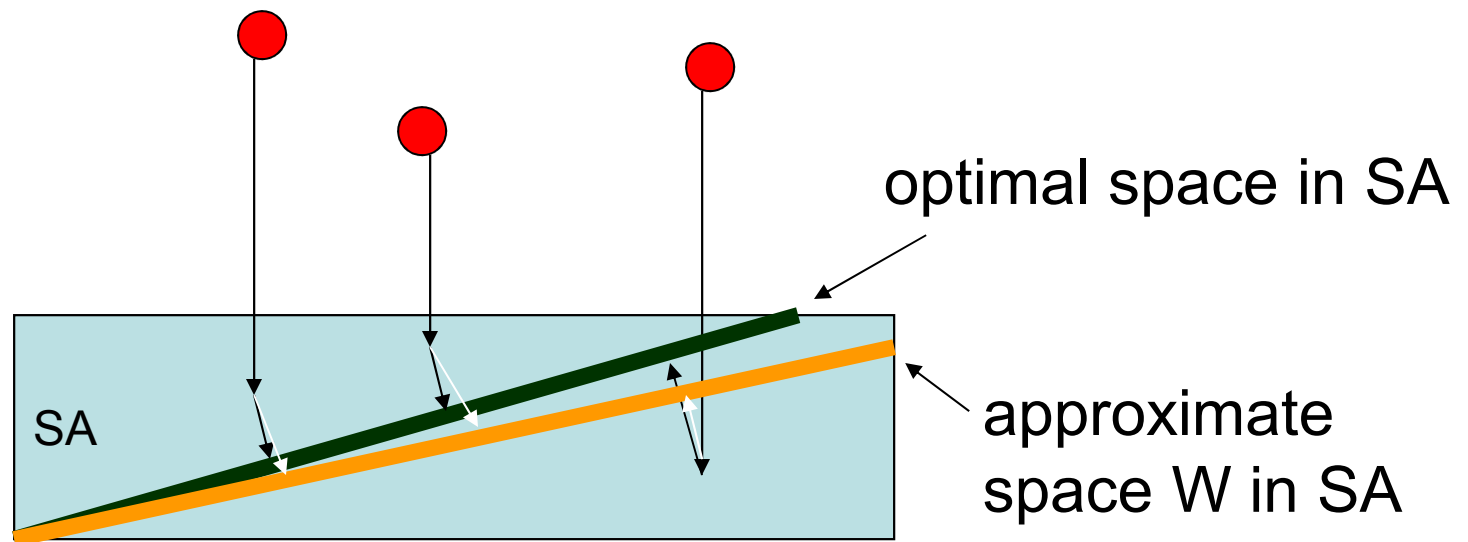


$$\text{cost} \leq (1+\varepsilon)|A-A_k|_F$$

# [KVW] protocol

---

- Phase 2:
- Find an approximately optimal space  $W$  inside of  $SA$



$$\text{cost} \leq (1+\varepsilon)^2 |A - A_k|_F$$



# [BWZ] Protocol

---

- Main Problem: communication is  $O(\text{skd}/\epsilon) + \text{poly}(\text{sk}/\epsilon)$
- We want  $O(\text{skd}) + \text{poly}(\text{sk}/\epsilon)$  communication!
- Idea: use **projection-cost preserving sketches** [CEMMP]
- Let  $A$  be an  $n \times d$  matrix
- If  $S$  is a random  $k/\epsilon^2 \times n$  matrix, then there is a constant  $c \geq 0$  so that for all  $k$ -dimensional projection matrices  $P$ :  
$$|SA(I - P)|_F + c = (1 \pm \epsilon)|A(I - P)|_F$$

## [BWZ] Protocol

Intuitively,  $U$  looks like top  $k$  left singular vectors of  $SA$

- Let  $S$  be a  $k/\varepsilon^2 \times n$  projection-cost preserving sketch
- Let  $T$  be a  $d \times k/\varepsilon^2$  projection-cost preserving sketch
- Server  $t$  sends  $SA^tT$  to Coordinator
- Coordinator sends back  $SAT = \sum_t SA^tT$  to servers
- Each server computes  $k/\varepsilon^2 \times k$  matrix  $U$  of top  $k$  left singular vectors of  $SAT$

Thus,  $U^TSA$  looks like top  $k$  right singular vectors of  $SA$

- Server  $t$  sends  $U^TSA^t$  to Coordinator
- Coordinator returns the space  $U^TSA = \sum_t U^TSA^t$  to output

Top  $k$  right singular vectors of  $SA$  work because  $S$  is a projection-cost preserving sketch!

## [BWZ] Analysis

---

- Let  $W$  be the row span of  $U^T SA$ , and  $P$  be the projection onto  $W$
- Want to show  $|A - AP|_F \leq (1 + \epsilon)|A - A_k|_F$
- Since  $T$  is a projection-cost preserving sketch,

$$(*) \quad |SA - SAP|_F \leq |SA - UU^T SA|_F + c_1 \leq (1 + \epsilon)|SA - [SA]_k|_F$$

- Since  $S$  is a projection-cost preserving sketch, there is a scalar  $c > 0$ , so that for all  $k$ -dimensional projection matrices  $P$ ,

$$|SA - SAP|_F + c = (1 \pm \epsilon)|A - AP|_F$$

- Add  $c$  to both sides of  $(*)$  to conclude  $|A - AP|_F \leq (1 + \epsilon)|A - A_k|_F$  99

# Conclusions for Distributed Low Rank Approximation

---

- [BWZ] Optimal  $O(sdk) + \text{poly}(sk/\epsilon)$  communication protocol for low rank approximation in arbitrary partition model
  - Handle bit complexity by adding Tao/Vu noise
  - Input sparsity time
  - 2 rounds, which is optimal [W]
  - Optimal data stream algorithms improves [CW, L, GP]
- Communication of other optimization problems?
  - Computing the rank of an  $n \times n$  matrix over the reals
  - Linear Programming
  - Graph problems: Matching
  - etc.

# Additional Time-Permitting Topics

---

- Will cover some recent topics at a research-level (many details omitted)
- Weighted Low Rank Approximation
- Regression and Low Rank Approximation with M-Estimator Loss Functions
- Finding Heavy Hitters in a Data Stream optimally