

# **EE270**

## **Large scale matrix computation, optimization and learning**

Instructor : Mert Pilanci

Stanford University

Tuesday, Feb 2 2020

# Randomized Linear Algebra

## Lecture 8: Randomized Least Squares Bias and Variance, Streaming Data

# Least Squares Problems and Random Projection

- Given  $A \in \mathbb{R}^{n \times d}$  and  $b \in \mathbb{R}^d$   
find the best linear fit  $Ax \approx b$  according to

$$\min_{x \in \mathbb{R}^d} \|Ax - b\|_2^2 + \lambda x^T x$$

$$\begin{matrix} d \\ \left[ \begin{matrix} A \end{matrix} \right] \end{matrix} \quad \begin{matrix} \text{rank}(A) = d \\ \text{if } A \text{ is} \\ \text{full column rank} \end{matrix}$$

- no regularization, i.e.,  $\lambda = 0$
- If  $A$  is full column rank then

$$\begin{aligned} x_{LS} &= (A^T A)^{-1} A^T b = A^\dagger b \\ &\approx (A^T S^T S A)^{-1} A^T S^T S b = (SA)^\dagger (Sb) \end{aligned}$$

$S$ : sampling matrix (need a distribution)

$S$ : iid distribution or

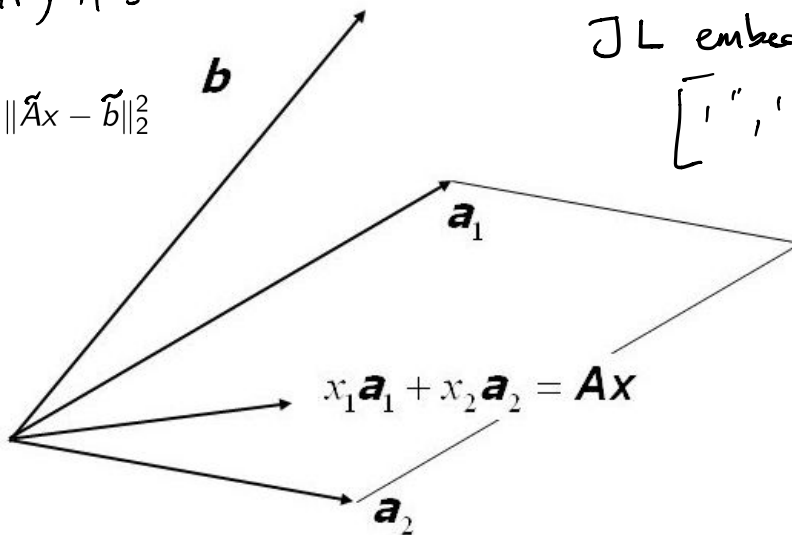
JL embedding  
 $\begin{bmatrix} | & | & | & | \\ \cdot & \cdot & \cdot & \cdot \\ | & | & | & | \end{bmatrix}$

$$= \arg \min_{x \in \mathbb{R}^d} \|\tilde{A}x - \tilde{b}\|_2^2$$

$$\tilde{A} = SA$$

$$\tilde{b} = Sb$$

Left sketching



# Faster Least Squares Optimization: Random Projection



## ► Left-sketching

Form  $SA$  and  $Sb$  where  $S \in \mathbb{R}^{m \times n}$  is a random projection matrix

## ► Solve the smaller problem

$$\min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2^2$$

## ► using any classical method.

Direct method complexity  $md^2$

# Approximation Result

- ▶ Suppose that  $n \gg d$
- ▶ Let  $S \in \mathbb{R}^{m \times d}$  be a Johnson-Lindenstrauss Embedding

$$x_{LS} = \arg \min_{x \in \mathbb{R}^d} \underbrace{\|Ax - b\|_2^2}_{f(x)}$$



$$\tilde{x} = \arg \min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2^2$$

$$\text{rank}(A) \leq d$$

- ▶ **Lemma** If  $m \geq \text{constant} \times \frac{\text{rank}(A)}{\epsilon^2}$  then,
- ▶  $f(x_{LS}) \leq f(\tilde{x}) \leq (1 + \epsilon^2)f(x_{LS})$
- ▶  $\|A(x_{LS} - \tilde{x})\|_2^2 \leq \epsilon^2$  with high probability

error prob. is exponentially small.

# Application: Streaming data

- Suppose that  $n \gg d$
- Let  $S \in \mathbb{R}^{m \times d}$  be a Johnson-Lindenstrauss Embedding

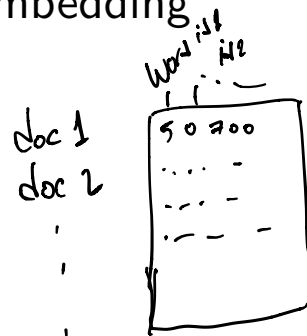
$$A \begin{matrix} \xrightarrow{t=1} \\ \boxed{\phantom{0}} \end{matrix} \quad A + \Delta_t \begin{matrix} \xrightarrow{t=2} \\ \boxed{\phantom{0}} \end{matrix} \quad b + S_t$$

Randomized Linear Algebra

$$x_{LS} = \arg \min_{x \in \mathbb{R}^d} \underbrace{\|Ax - b\|_2^2}_{f(x)}$$

$O(nd)$  space +  $O(nd^2)$  compute per time instant

$$\tilde{x} = \arg \min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2^2$$



Idea: Update  $SA \leftarrow SA + S\Delta$  and  $Sb \leftarrow Sb + S\delta$

- A and b are dynamically updated and we need to find  $x_{LS}$  at any time

$$A_{t+1} = A_t + \Delta_t \text{ and } y_{t+1} = y_t + \Delta_t$$

Can we form and update  $A_t^T A_t \in \mathbb{R}^{d \times d}$  ?

$$x_{LS} = (\bar{A}^T \bar{A})^{-1} \bar{A}^T \bar{b}$$

Classical Linear Algebra

$$(A + \Delta)^T (A + \Delta) = \underbrace{\bar{A}^T \bar{A}} + \underbrace{\Delta^T \bar{A}} + \dots$$

$O(nd)$  space and  $O(nd^2)$  compute per time instant

## Application: Streaming data

- ▶ Suppose that  $n \gg d$
- ▶ Let  $S \in \mathbb{R}^{m \times d}$  be a Johnson-Lindenstrauss Embedding

$$x_{LS} = \arg \min_{x \in \mathbb{R}^d} \underbrace{\|Ax - b\|_2^2}_{f(x)}$$

$$\tilde{x} = \arg \min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2^2$$

- ▶  $A$  and  $b$  are dynamically updated and we need to find  $x_{LS}$  at any time

$$A_{t+1} = A_t + \Delta_t \text{ and } y_{t+1} = y_t + \Delta_t$$

Can we form and update  $A_t^T A_t \in \mathbb{R}^{d \times d}$  ?

- ▶ Linear sketch can be updated on the fly

$$SA_{t+1} = SA_t + S\Delta_t \text{ and } Sy_{t+1} = Sy_t + S\Delta_t$$

## Gaussian Sketch

- ▶ Let  $S$  be  $\frac{1}{m} \times$  i.i.d. Gaussian.  $\mathbb{E}[S^T S] = I$

$$\tilde{x} = \arg \min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2^2$$

- ▶ Is  $\mathbb{E}[\tilde{x}]$  equal to  $x_{LS}$ ?



## Gaussian Sketch

- ▶ Let  $S$  be  $\frac{1}{m} \times$  i.i.d. Gaussian.  $\mathbb{E}[S^T S] = I$

$$\tilde{x} = \arg \min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2^2$$

- ▶ Is  $\mathbb{E}[\tilde{x}]$  equal to  $x_{LS}$ ? Yes! Only Gaussian.
- ▶ Assuming  $A^T S^T SA$  is invertible, we have

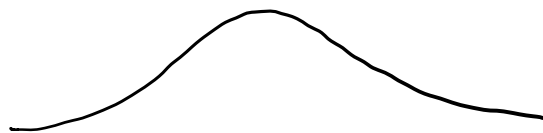
$$\tilde{x} = (A^T S^T SA)^{-1} A^T S^T Sb$$

let  $b = Ax_{LS} + b^\perp$  where  $b^\perp \perp \text{Range}(A)$

$$\begin{aligned}\tilde{x} &= (A^T S^T SA)^{-1} A^T S^T S(Ax_{LS} + b^\perp) \\ &= x_{LS} + (A^T S^T SA)^{-1} A^T S^T Sb^\perp\end{aligned}$$

- ▶  $\mathbb{E}(A^T S^T SA)^{-1} A^T S^T Sb^\perp = 0$  since  $Sb^\perp$  and  $SA$  are uncorrelated zero mean Gaussian.

# Gaussian Sketch: Variance



- Let  $S$  be i.i.d. Gaussian

$$f(x_{LS}) = \min_x \|Ax - b\|_2^2$$

$$\begin{aligned}\tilde{x} &= \arg \min_{x \in \mathbb{R}^d} \underbrace{\|SAx - Sb\|_2^2}_{f(x)} = x_{LS} + (A^T S^T SA)^{-1} A^T S^T Sb^\perp \\ &= x_{LS} + (SA)^\dagger Sb^\perp\end{aligned}$$

- Analyzing the variance  $\mathbb{E}\|A\tilde{x} - x_{LS}\|_2^2$
- **Lemma (a)** Conditioned on the matrix  $SA$

$$\tilde{x} \sim N\left(x_{LS}, \frac{f(x_{LS})}{m} (A^T S^T SA)^{-1}\right)$$

# Gaussian Sketch: Variance

$$\mathbb{E} \tilde{S} \tilde{S}^T = I$$

$$S$$

$$b^\perp$$

$$\mathbb{E}(S b^\perp) = \mathbb{E} S_i^\top b^\perp = 0$$

- Let  $S$  be i.i.d. Gaussian  $\times \frac{1}{\sqrt{m}}$

$$\begin{aligned} \tilde{x} &= \arg \min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2^2 = x_{LS} + (A^T S^T SA)^{-1} A^T S^T S b^\perp \\ &= x_{LS} + (SA)^\dagger S b^\perp \end{aligned}$$

independent

- Analyzing the variance  $\mathbb{E} \|A\tilde{x} - x_{LS}\|_2^2$

- Lemma (a)** Conditioned on the matrix  $(SA)$

$$\tilde{x} \sim N\left(x_{LS}, \frac{f(x_{LS})}{m} (A^T S^T SA)^{-1}\right)$$

$$\begin{aligned} \mathbb{E} S b^\perp (S b^\perp)^\top &= \mathbb{E} S_i^\top b^\perp \cdot S_i^\top b^\perp \\ &= \frac{\mathbb{E} S_i^\top b^\perp \cdot S_i^\top b^\perp}{0} - \frac{\mathbb{E} S_i^\top b^\perp \cdot S_i^\top b^\perp}{0} \\ &= \frac{b^\perp \cdot \mathbb{E} S_i^\top S_i^\top b^\perp}{I \cdot \frac{1}{m}} = \frac{b^\perp b^\perp^\top}{m} \end{aligned}$$

$$S b^\perp \sim N\left(0, \frac{\|b^\perp\|_2^2}{m} I\right)$$

$$\mathbb{E}(\tilde{x} - x_{LS})(\tilde{x} - x_{LS})^T = (SA)^\dagger ((SA)^\dagger)^T = (A^T S^T SA)^{-1} \frac{\|b^\perp\|_2^2}{m}$$

$$= \mathbb{E} \left( \underbrace{(A^T S^T SA)^{-1}}_{\text{matrix}} \underbrace{A^T S^T}_{\text{matrix}} \underbrace{S b^\perp}_{\text{vector}} \left( \underbrace{(A^T S^T SA)^{-1}}_{\text{matrix}} \underbrace{A^T S^T}_{\text{matrix}} \underbrace{S b^\perp}_{\text{vector}} \right)^\top \right)$$

## Gaussian Sketch: Variance

- ▶ Let  $S$  be i.i.d. Gaussian

$$\tilde{x} = \arg \min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2^2$$

- ▶ Analyzing the variance  $\mathbb{E}\|A\tilde{x} - x_{LS}\|_2^2$
- ▶ **Lemma (a)** Conditioned on the matrix  $SA$

$$\tilde{x} \sim N\left(x_{LS}, \frac{f(x_{LS})}{m} (A^T S^T SA)^{-1}\right)$$

$$A(\tilde{x} - x_{LS}) \sim N\left(0, \frac{f(x_{LS})}{m} A(A^T S^T SA)^{-1} A\right)$$

$$x \sim N(0, \Sigma)$$

$$Ax \sim N(0, A \Sigma A^T)$$

$$\mathbb{E} \cdot Ax \cdot x^T A^T = A \Sigma A^T$$

# Gaussian Sketch: Variance

- ▶ Let  $S$  be i.i.d. Gaussian

$$\tilde{x} = \arg \min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2^2$$

- ▶ Analyzing the variance  $\mathbb{E}\|A\tilde{x} - x_{LS}\|_2^2$
- ▶ **Lemma (a)** Conditioned on the matrix  $SA$

$$\tilde{x} \sim N\left(x_{LS}, \frac{f(x_{LS})}{m} (A^T S^T SA)^{-1}\right)$$

$$A(\tilde{x} - x_{LS}) \sim N\left(0, \frac{f(x_{LS})}{m} \underbrace{A(A^T S^T SA)^{-1} A^T}_{\approx \frac{1}{m} A(A^T A)^{-1} A^T} \right) \approx \mathbb{I} \cdot \frac{m}{m-d-1}$$

**Lemma (b)** (removing conditioning) for  $m > d + 1$

$$\mathbb{E}[(A^T S^T SA)^{-1}] = (A^T A)^{-1} \frac{m}{m-d-1}$$

$$\mathbb{E}[\underbrace{(A^T \tilde{S}^T \tilde{S} A)^{-1}}_{\approx \frac{1}{m} A^T A}]$$

for  $m > d+1$

Only for  
the Gaussian  
distribution,

→  
Inverse matrices  
are biased but  
they are predictable

# Gaussian Sketch: Variance

- Let  $S$  be i.i.d. Gaussian

$$x \sim N(0, \Sigma)$$

$$\begin{aligned} \mathbb{E} \|x\|_2^2 &= \mathbb{E} \operatorname{tr} x x^T = \mathbb{E} \operatorname{tr} x^T x \\ &= \operatorname{tr} \mathbb{E} x x^T \\ &= \operatorname{tr} \Sigma \end{aligned}$$

$$\tilde{x} = \arg \min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2^2$$

- Analyzing the variance  $\mathbb{E} \|A\tilde{x} - x_{LS}\|_2^2$

- Lemma (a)** Conditioned on the matrix  $SA$

lower prop. of cond. exp.  $\Rightarrow \mathbb{E}[\mathbb{E}[\|C\|_2^2 \mid SA]]$

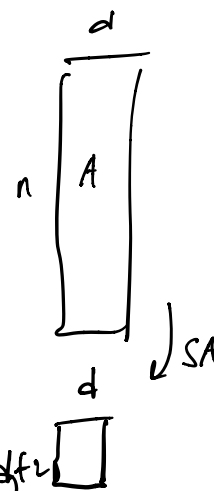
$$\tilde{x} \sim N\left(x_{LS}, \frac{f(x_{LS})}{m} (A^T S^T SA)^{-1}\right)$$

$SA$   
 $= \mathbb{E} \|C\|_2^2$

$$A(\tilde{x} - x_{LS}) \sim N\left(0, \frac{f(x_{LS})}{m} A(A^T S^T SA)^{-1} A\right)$$

- Lemma (b)** (removing conditioning) for  $m > d + 1$

$$\mathbb{E}[(A^T S^T SA)^{-1}] = (A^T A)^{-1} \frac{m}{m-d-1}$$



$$\begin{aligned} &\mathbb{E} \|A(\tilde{x} - x_{LS})\|_2^2 \\ &= \frac{f(x_{LS})}{m-d-1} \operatorname{tr} A(A^T A)^{-1} A^T = f(x_{LS}) \frac{d}{m-d-1} \end{aligned}$$

# Expected Inverse of a Random Matrix

$$S = \text{---} \quad A = \begin{vmatrix} \frac{1}{g_1^2 + g_2^2 + \dots} \end{vmatrix}$$

$$\mathbb{E} A^T S^T S A = A^T A$$

$$\mathbb{E} (A^T S^T S A)^{-1} = A^T A \cdot \gamma$$

► Where does the formula

$$\mathbb{E} [(A^T S^T S A)^{-1}] = (A^T A)^{-1} \underbrace{\left( \frac{m}{m-d-1} \right)}_{\gamma}$$

come from?

inverse chi-square dist.

as  $m \rightarrow \infty$  we have  $\gamma \rightarrow 1$

set  $m = \alpha \cdot d$

$$\gamma = \frac{\alpha^d}{(\alpha-1)d-1} \approx \frac{\alpha}{\alpha-1}$$

Randomized algorithm to estimate the inverse!

# Which sketching matrices are good?

deterministic      slow & exact  
random      fast & approx.

- ▶ We need to find conditions to guarantee approximate optimality
- ▶ Let  $A = U\Sigma V^T$  SVD in compact form (drop zero singular vals)

$$\bar{A}^T A : O(nd^2)$$

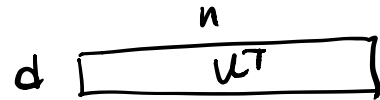
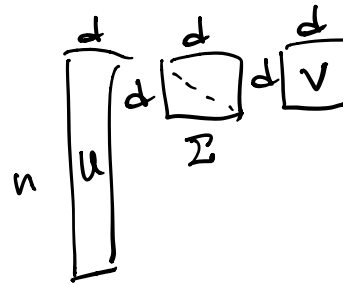
SA is faster

some deterministic options

- ▶  $S = U^T$  is  $d \times n$  (optimal but slow)

$O(nd)$   $\leftarrow S = A^T$  (optimal but slow)

min  $\| \bar{A}^T A x - \bar{A}^T b \|_2$   
 $\times \quad ((\bar{A}^T A)(\bar{A}^T A)^T) \bar{A}^T A \cdot \bar{A}^T b \quad (\bar{A}^T A)^T \bar{A}^T b$   
 since we need SVD



$$A = U\Sigma V^T$$

$$\bar{A}^T A = V\Sigma^T U U^T \Sigma V^T$$

$$= V \Sigma^T V^T$$

- ▶ For random  $S$  matrices  $A^T S^T S A$  needs to be invertible we want it to be close to  $A^T A$

$$x_{LS} = (A^T A)^+ A^T b = A^+ b$$

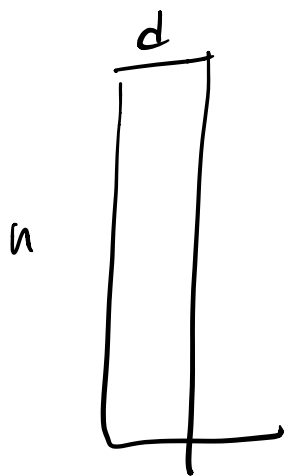
$$\min \| S A x - S b \|_2^2$$

pick  $S = U^T$

$$= \min \| U^T U \Sigma V^T - U^T b \|_2^2 = (\Sigma^T V^T)^+ U^T b = \underbrace{(V \Sigma^T V^T)^+}_{A^T A} \underbrace{V \Sigma U^T b}_{A^T b}$$

=  $x_{LS}$  no error!





$$\bar{A}^T A : O(nd^2)$$

$$\approx \underbrace{\bar{A}^T S^T}_{m \times d} S A : O(nd^2)$$

$$m \approx d$$

Questions?