# EE270
# Large scale matrix computation, optimization and learning

Instructor : Mert Pilanci
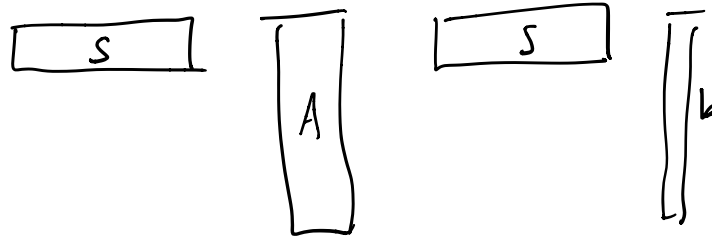
Stanford University

Tuesday, Feb 9 2021

# Randomized Linear Algebra
## Lecture 9: High-dimensional Problems, Least-norm Solutions and Randomized Methods

# 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 $md^2$

# Gaussian Sketch

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

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

- Unbiased $\mathbb{E}[\tilde{x}] = x_{LS}$
  since $\tilde{x} = x_{LS} + \underbrace{(A^T S^T SA)^{-1} A^T S^T S b^\perp}_{\text{zero mean}}$
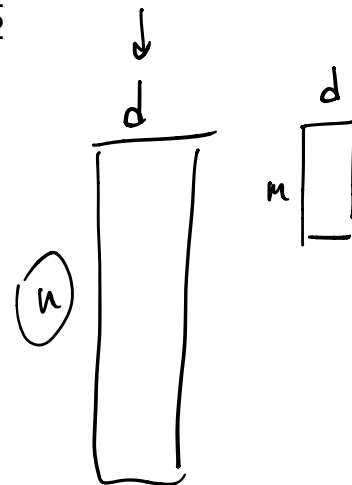
# Gaussian Sketch

$$f(x) = \| Ax - b \|_2^2$$
$$f(x) =$$

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

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

▶ Unbiased $\mathbb{E}[\tilde{x}] = x_{LS}$
since $\tilde{x} = x_{LS} + \underbrace{(A^T S^T SA)^{-1} A^T S^T Sb^\perp}_{\text{zero mean}}$

▶ Variance
$$\mathbb{E}\| A(\tilde{x} - x_{LS}) \|_2^2 = f(x_{LS}) \frac{d}{m-d-1}$$
valid for $m > d + 1$ where $f(x) = \| Ax - b \|_2^2$

# Gaussian Sketch

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

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

- Unbiased $\mathbb{E}[\tilde{x}] = x_{LS}$
  since $\tilde{x} = x_{LS} + \underbrace{(A^T S^T SA)^{-1} A^T S^T Sb^\perp}_{\text{zero mean}}$

- Variance
  $$\mathbb{E}\|A(\tilde{x} - x_{LS})\|_2^2 = f(x_{LS}) \frac{d}{m-d-1}$$
  valid for $m > d + 1$ where $f(x) = \|Ax - b\|_2^2$
- Function value
  $f(\tilde{x}) = \|A\tilde{x} - b\|_2^2 = \|A(\tilde{x} - x_{LS})\|_2^2 + \|Ax_{LS} - b\|_2^2$
- $\mathbb{E}f(\tilde{x}) - f(x_{LS}) = f(x_{LS}) \frac{d}{m-d-1}$

# Variance Reduction by Averaging

- Let $S_1, ..., S_r$ be $\frac{1}{\sqrt{m}} \times$ i.i.d. Gaussian. $\mathbb{E}[S_i^T S_i] = I$

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

- let $\tilde{x} = \frac{1}{r} \sum_{i=1}^{r} x_i$
- Unbiased $\mathbb{E}[\tilde{x}] = x_{LS}$
- Variance is reduced by $\frac{1}{r}$
- $\mathbb{E}\|A(\tilde{x} - x_{LS})\|_2^2 = f(x_{LS}) \frac{1}{r} \frac{d}{m-d-1}$

$$\mathbb{E}\,\tilde{x} = \frac{1}{r} \sum \underbrace{\mathbb{E} x_i}_{x_{LS}} = x_{LS}$$

$$\mathrm{Var}\left(\boxed{\left(\frac{1}{r}\right)} \sum_{i=1}^{r} x_i\right) = \frac{1}{r^2} \cdot \underline{\mathrm{Var}(\sum x_i)}$$

$$= \frac{1}{r^2} \cdot \sum_{i=1}^{r} \mathrm{Var}(x_i)$$

$$= \frac{r}{r^2} \cdot \mathrm{Var}(x_1)$$

# Variance Reduction by Averaging

$$b = Ax_{LS} + b^\perp$$

$$f(x) = \|Ax - b\|_2^2 = \|\widehat{A(x - x_{LS})} - b^\perp\|_2^2$$

$$= \|A(x - x_{LS})\|_2^2 + \underline{\|b^\perp\|_2^2}$$

$$\underbrace{f(x_{LS})}$$

$$\|Ax_{LS} - b\|_2^2$$

▶ Let $S_1, ..., S_r$ be $\frac{1}{m} \times$ i.i.d. Gaussian. $\mathbb{E}[S^T S] = I$

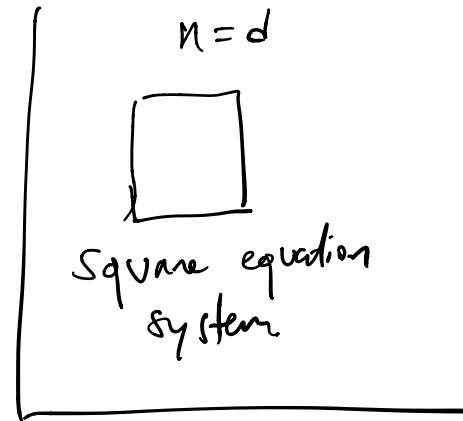$$\tilde{x}_i = \arg\min_{x \in \mathbb{R}^d} \|S_i A x - S_i b\|_2^2$$

▶ let $\tilde{x} = \frac{1}{r} \sum_{i=1}^{r} x_i$
▶ Unbiased $\mathbb{E}[\tilde{x}] = x_{LS}$
▶ Variance is reduced by $\frac{1}{r}$
▶ $\mathbb{E}\|A(\tilde{x} - x_{LS})\|_2^2 = f(x_{LS}) \frac{1}{r} \frac{d}{m - d - 1}$
▶ $\mathbb{E}f(\tilde{x}) - \underbrace{f(x_{LS})} = f(x_{LS}) \frac{1}{r} \frac{d}{m - d - 1}$

$$f(x_{LS}) \leq \mathbb{E} f(\tilde{x}) \leq f(x_{LS}) \cdot \left(1 + \frac{1}{r} \frac{d}{m - d - 1}\right)$$

# High-dimensional Least Squares Problems

Right sketching.

$A \cdot S$

$n = d$

Square equation system

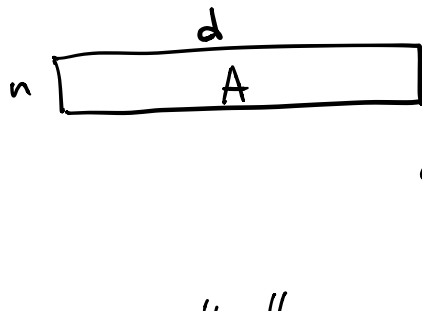- $A \in \mathbb{R}^{n \times d}$ where $d > n$
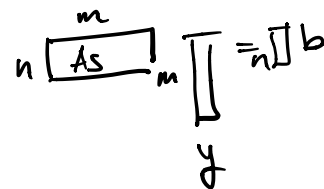- no unique solution

tall LS

wide LS (high-dim)

$n >> d$

$\min \|x\|_2$

s.t. $Ax = b$

$\min \|y\|_2$

$ASy = b$

# High-dimensional Least Squares Problems

$A = U\Sigma V^\top \quad A^+ = A^\top (AA^\top)^{-1}$

$$\min_x \quad \tfrac{1}{2}\|x\|_2^2 + \lambda^\top \cdot (Ax - b) = \tfrac{1}{2}\|A^\top \lambda\|_2^2 - \tfrac{1}{2}\lambda^\top AA^\top \lambda - \lambda^\top b$$
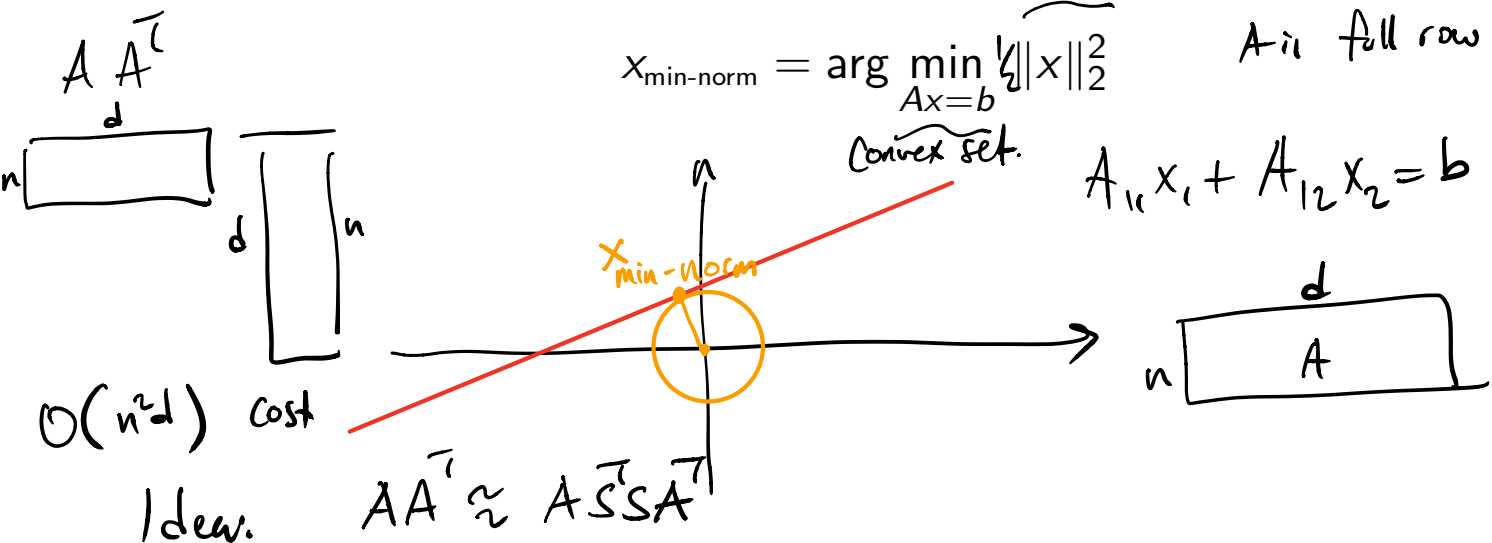
$$x + A^\top \lambda = 0 \implies x = -A^\top \lambda = A^\top (AA^\top)^{-1} b = A^+ b$$

- $A \in \mathbb{R}^{n \times d}$ where $d > n$
- no unique solution
- minimum ($\ell_2$) norm solution is unique

$$-AA^\top \lambda - b = 0$$

$$\lambda = -(AA^\top)^{-1} b$$

invertible iff
$A$ is full row rank

$$x_{\text{min-norm}} = \arg\min_{Ax=b} \tfrac{1}{2}\|x\|_2^2$$

Convex set.

$A A^\top$

$$n \boxed{\phantom{AAAA}} \; d$$

$$d \left[ \begin{matrix} \phantom{A} \\ \phantom{A} \end{matrix} \right] n$$

$x_{\text{min-norm}}$

$O(n^2 d)$ cost

Idea.

$$A A^\top \approx A S^\top S A^\top$$

$$A_{11} x_1 + A_{12} x_2 = b$$
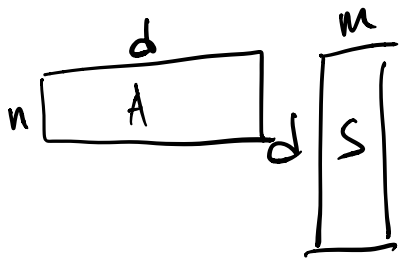
$d$

$n \boxed{\quad A \quad}$

# Minimum norm solution and SVD

$$x_{\text{min-norm}} = \arg \min_{Ax=b} \|x\|_2^2$$

# Random projection to reduce dimension: Right Sketch

$$x_{\text{min-norm}} = \arg \min_{Ax=b} \|x\|_2^2$$

▶ We can right multiply $A$ and form $AS$ where $S \in \mathbb{R}^{d \times m}$ and solve

$$\arg \min_{ASz=b} \|z\|_2^2$$



$$z = (AS)^+ b \qquad m \text{ -dimensional}$$

$$\text{Idea 1) let} \qquad \overset{n}{x} = S (AS)^+ b$$

# Random projection to reduce dimension: Right Sketch

$$x_{\text{min-norm}} = \arg \min_{Ax=b} \|x\|_2^2$$

▶ We can right multiply $A$ and form $AS$ where $S \in \mathbb{R}^{d \times m}$ and solve

$$\arg \min_{ASz=b} \|z\|_2^2$$

▶ How do we use $z \in \mathbb{R}^m$?

# Right Sketch

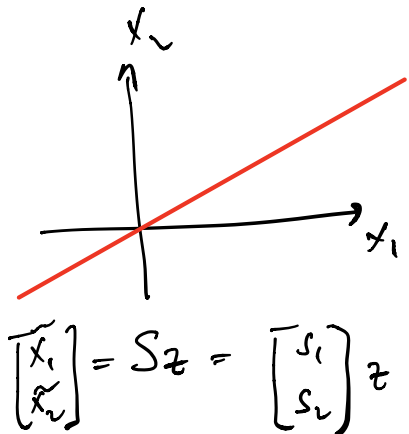$$x_{\text{min-norm}} = \arg \min_{Ax=b} \underbrace{\|x\|_2^2}_{f(x)}$$

approximation $\quad \tilde{x} = S\tilde{z}$

where $\tilde{z} := \arg \min_{ASz=b} \|z\|_2^2$

Feasible estimate for $x$ : $\qquad A\tilde{x} = AS\tilde{z}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad = b$

# Right Sketch



$$x_{\text{min-norm}} = \arg \min_{Ax=b} \underbrace{\|x\|_2^2}_{f(x)}$$

approximation $\quad \tilde{x} = S\tilde{z}$

where $\tilde{z} := \arg \min_{ASz=b} \|z\|_2^2$

- Let $S$ be i.i.d. Gaussian $N(0, \frac{1}{\sqrt{m}})$
- Is $\tilde{x}$ unbiased, i.e., $\mathbb{E}\tilde{x} \stackrel{?}{=} x_{\text{min-norm}}$

# Right Sketch

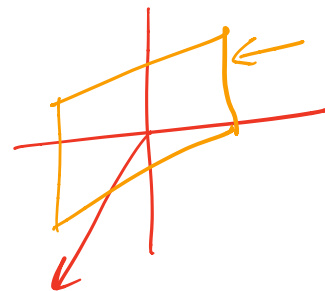$$x_{\text{min-norm}} = \arg \min_{Ax=b} \underbrace{\|x\|_2^2}_{f(x)}$$

approximation $\quad \tilde{x} = S\tilde{z}$

where $\tilde{z} := \arg \min_{ASz=b} \|z\|_2^2$

▶ Let $S$ be i.i.d. Gaussian $N(0, \frac{1}{\sqrt{m}})$
▶ Is $\tilde{x}$ unbiased, i.e., $\mathbb{E}\tilde{x} \stackrel{?}{=} x_{\text{min-norm}}$
▶ Yes, conditioned on $SA$

$$\tilde{x} \sim N(x_{\text{min-norm}}, V_2 V_2^T b^T (AS^T SA^T)^{-1} b)$$

▶ $V_2 V_2^T$ is the projection onto the null space of $A$
▶ error $\tilde{x} - x_{\text{min-norm}} \in \text{Null}(A)$

$Ax = b$

$A\tilde{x} = b$

$A(x - \hat{x}) = 0$

$x - \tilde{x} \in \text{Null}(A)$

$A$

$A = U \Sigma V^T$
$= U [\Sigma \ 0][v_1 \ v_2]^T$

# Right Sketch

$$x_{\text{min-norm}} = \arg \min_{Ax=b} \underbrace{\|x\|_2^2}_{f(x)}$$

$$A^+b = \bar{A}(A\bar{A}^{-1})b$$

$$\text{approximation} \quad \tilde{x} = S\tilde{z}$$

$$\text{where } \tilde{z} := \arg \min_{ASz=b} \|z\|_2^2$$

$$(AS)^+ b = \bar{S}\bar{A}^T(AS\bar{S}^T\bar{A}^T)^{-1}b$$

$n \times n$ matrix

- Let $S$ be i.i.d. Gaussian $N(0, \frac{1}{\sqrt{m}})$
- Is $\tilde{x}$ unbiased, i.e., $\mathbb{E}\tilde{x} =^? x_{\text{min-norm}}$
- Yes, conditioned on $SA$

We need at least $m \geq n$

$$\tilde{x} \sim N(x_{\text{min-norm}}, VV^T b^T (AS^T SA^T)^{-1} b)$$
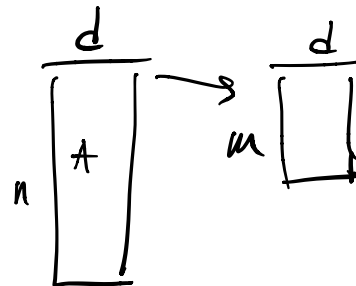
- $VV^T$ is the projection onto the null space of $A$
- error $\tilde{x} - x_{\text{min-norm}} \in \text{Null}(A)$
- Using $\mathbb{E}(AS^T SA^T)^{-1} = (AA^T)^{-1} \frac{m}{m-n-1}$

$$\boxed{m > n+1}$$

$$\mathbb{E}\|\tilde{x} - x_{\text{min-norm}}\|_2^2 = \frac{d-n}{m-n-1} f(x_{\text{min-norm}}) = \frac{d-n}{m-n-1} \|x_{\text{min-norm}}\|_2^2$$

# Left Sketch vs Right Sketch Summary



▶ Both are unbiased using Gaussian projections

▶ $A$ is $n \times d$

▶ Left sketch $n \geq d$

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

Variance: $\mathbb{E}\|A(\tilde{x} - x_{LS})\|_2^2 = f(x_{LS}) \dfrac{d}{m-d-1}$

dimension of the relevant subspace

▶ Right sketch $d > n$

$$\tilde{x} = S\tilde{z} \quad \text{where } \tilde{z} := \arg \min_{ASz=b} \|z\|_2^2$$

Variance: $\mathbb{E}\|\tilde{x} - x_{\text{min-norm}}\|_2^2 = f(x_{\text{min-norm}}) \dfrac{d-n}{m-n-1}$

dimension of the relevant subspace

# Back to Left Sketch: Which sketching matrices are good?

▶ We need to find conditions to guarantee approximate optimality

▶ Let $A = U\Sigma V^T$ SVD in compact form

some deterministic options

▶ $S = U^T$ is $d \times n$

▶ $S = A^T$

high dam

$V_1^T$

$S = A$

▶ For random $S$ matrices $A^T S^T S A$ needs to be invertible
we want it to be close to $A^T A$
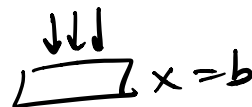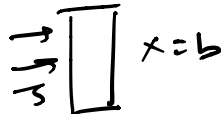
# Approximate Matrix Multiplication

▶ Let the approximate product of $AB$ be $C = AS^T SB$

$$\mathbb{P}\left[\|AB - C\|_F > \epsilon\|A\|_F\|B\|_F\right] \leq \delta$$

▶ Follows from JL Moment property
▶ $S \in \mathbb{R}^{m \times n} \sim \frac{1}{\sqrt{m}} \times$ random i.i.d. sub-Gaussian, e.g., $\pm 1$, or $N(0,1)$ with $m = \frac{c_1}{\epsilon^2} \log \frac{1}{\delta}$
▶ $S \in \mathbb{R}^{m \times n} \sim \frac{1}{\sqrt{m}} \times$ CountSketch matrix (one nonzero per column, which is $\pm 1$ at a uniformly random location) with $m = \frac{c_2}{\epsilon^2 \delta}$
▶ $S \in \mathbb{R}^{m \times n} \sim \frac{1}{\sqrt{m}} \times \underline{\text{Fast JL}}$ Transform with $m = \frac{c_3}{\epsilon} \log \frac{1}{\delta}$

Fourier Transform

→ row /colum sampling

$\vec{\rightarrow} \Box \; x = b$

$\Box x = b$

# Approximate Matrix Multiplication

- Let the approximate product of $AB$ be $C = AS^T SB$

$$\mathbb{P}\left[\|AB - C\|_F > \epsilon \|A\|_F \|B\|_F\right] \leq \delta$$
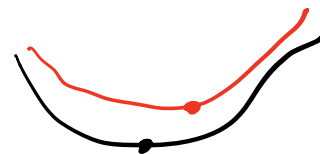
- Follows from JL Moment property
- $S \in \mathbb{R}^{m \times n} \sim \frac{1}{\sqrt{m}} \times$ random i.i.d. sub-Gaussian, e.g., $\pm 1$, or $N(0,1)$ with $m = \frac{c_1}{\epsilon^2} \log \frac{1}{\delta}$
- $S \in \mathbb{R}^{m \times n} \sim \frac{1}{\sqrt{m}} \times$ CountSketch matrix (one nonzero per column, which is $\pm 1$ at a uniformly random location) with $m = \frac{c_2}{\epsilon^2 \delta}$
- $S \in \mathbb{R}^{m \times n} \sim \frac{1}{\sqrt{m}} \times$ Fast JL Transform with $m = \frac{c_3}{\epsilon} \log \frac{1}{\delta}$
- Sparse JL and Fast JL are more efficient
- advantages: doesn't require any knowledge about matrices $A$ and $B$ (**oblivious**)
- optimal sampling probabilities depend on the column/row norms of $A$ and $B$

# Basic Inequality Method

min $f(x)$    min $\tilde{f}(x)$

second

first
ineq.

▶ We minimize $\tilde{x} = \arg\min \|S(Ax - b)\|_2^2$

▶ $x_{LS}$ minimizes $\|Ax - b\|_2^2$

▶ How far is $\tilde{x}$ from $x_{LS}$?

▶ **Step 1**. Establish two optimality (in)equalities for these variables

▶ $\|Ax_{LS} - b\|_2^2 \leq \|Ax' - b)\|_2^2$ for any $x'$, i.e., $A^T(Ax_{LS} - b) = 0$

▶ $\|S(A\tilde{x} - b)\|_2^2 \leq \|S(Ax_{LS} - b)\|_2^2$

13 / 15

# Basic Inequality Method

$$\|SA(\Delta + x_{LS}) - Sb\|_2^2 \le \|SAx_{LS} - Sb\|_2^2$$

$$b = Ax_{LS} + b^\perp$$

$$\|SA\Delta + S\!\!\!\!/Ax_{LS} - S\!\!\!\!/Ax_{LS} - Sb^\perp\|_2^2 \le \|SA\!\!\!\!/x_{LS} - SAx_{LS} - Sb^\perp\|_2^2$$

$$= \|Sb^\perp\|_2^2$$

$$\|SA\Delta\|_2^2 + \|Sb^\perp\|_2^2 - 2b^{\perp T}S^TS A\Delta \le \|\cancel{b^\perp}\|_2^2$$

- We minimize $\tilde{x} = \arg\min \|S(Ax - b)\|_2^2$
- $x_{LS}$ minimizes $\|Ax - b\|_2^2$
- How far is $\tilde{x}$ from $x_{LS}$?
- **Step 1**. Establish two optimality (in)equalities for these variables
- $\|Ax_{LS} - b\|_2^2 \le \|Ax' - b)\|_2^2$ for any $x'$, i.e., $A^T(Ax_{LS} - b) = 0$
- $\|S(A\tilde{x} - b)\|_2^2 \le \|S(Ax_{LS} - b)\|_2^2$
- **Step 2**. Define error $\Delta = \tilde{x} - x_{LS}$ and re-write these inequalities in terms of $\Delta$  $\qquad \tilde{x} = \Delta + x_{LS}$
- $\|SA\Delta\|_2^2 \le 2b^{\perp T}(S^TS - I)A\Delta$
- **Step 3**. Argue $S^TS \approx I$

$$\underbrace{\Delta^T A^T}_{\approx I} S^TS A\Delta = \|SA\Delta\|_2^2 \le 2b^\perp S^TSA\Delta = 2b^\perp\underbrace{(S^TS - I)}_{\approx I}A\Delta \quad \text{since} \quad b^{\perp T}\!\!\cdot A = 0$$

# Leverage Scores

# Questions?