

CSE 392: Matrix and Tensor Algorithms for Data

Instructor: Shashanka Ubaru

University of Texas, Austin
Spring 2024

Lecture 10: Sampling and preconditioning for least squares

- 1 Sketch and solve - Proof
- 2 Sampling for least squares
- 3 Preconditioning for least squares

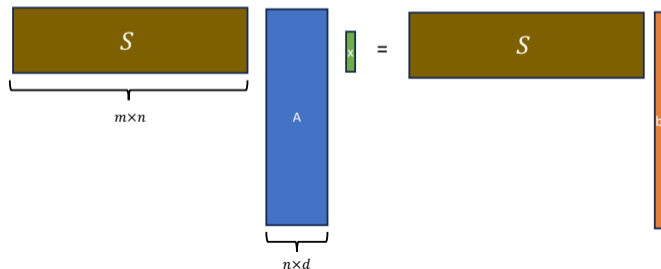
Sketch and solve

Recall:

- Generate a sketching matrix $\mathbf{S} \in \mathbb{R}^{m \times n}$.
- Compute sketches \mathbf{SA} and \mathbf{Sb} .
- Solve:

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

- Typically, $m = \text{poly}(d/\epsilon)$.



Subspace embedding for sketch and solve

Sketch and solve

Suppose $\mathbf{S} \in \mathbb{R}^{m \times n}$ is a subspace ϵ -embedding for $\text{span}([\mathbf{A} \ \mathbf{b}])$.

Let,

$$\mathbf{x}^* = \min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$$

$$\tilde{\mathbf{x}} = \min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{S}(\mathbf{A}\mathbf{x} - \mathbf{b})\|_2,$$

for $\epsilon \leq 1/3$, we have

$$\|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{b}\|_2 \leq (1 + 3\epsilon)\|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|_2$$

Implies, we have $O(1/\epsilon^2)$ dependency on the error tolerance.

Alternate proof

Sketch and solve

If $\mathbf{S} \in \mathbb{R}^{m \times n}$ is a Countsketch matrix with $m = O(d^2/\epsilon)$ or SRHT with $m = O(d \log d/\epsilon)$, or Gaussian sketch with $m = O(d/\epsilon)$, then

$$\|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{b}\|_2 \leq (1 + \epsilon)\|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|_2$$

Alternate proof

Sketch and solve

If $\mathbf{S} \in \mathbb{R}^{m \times n}$ is a Countsketch matrix with $m = O(d^2/\epsilon)$ or SRHT with $m = O(d \log d/\epsilon)$, or Gaussian sketch with $m = O(d/\epsilon)$, then

$$\|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{b}\|_2 \leq (1 + \epsilon)\|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|_2$$

Proof: Let us consider an orthonormal basis \mathbf{U} for \mathbf{A} .

Let, $\mathbf{U}\tilde{\mathbf{y}} = \mathbf{A}\tilde{\mathbf{x}}$ and $\mathbf{U}\mathbf{y}^* = \mathbf{A}\mathbf{x}^*$. Then,

$$\|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{b}\|_2^2 = \|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|_2^2 + \|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{A}\mathbf{x}^*\|_2^2$$

and

$$\|\mathbf{U}\tilde{\mathbf{y}} - \mathbf{b}\|_2^2 = \|\mathbf{U}\mathbf{y}^* - \mathbf{b}\|_2^2 + \|\mathbf{U}\tilde{\mathbf{y}} - \mathbf{U}\mathbf{y}^*\|_2^2$$

Need to show that $\|\mathbf{U}(\tilde{\mathbf{y}} - \mathbf{y}^*)\|_2^2 = \|\tilde{\mathbf{y}} - \mathbf{y}^*\|_2^2 = O(\epsilon)\|\mathbf{U}\mathbf{y}^* - \mathbf{b}\|_2^2$.

For a subspace embedding \mathbf{S} , we have

$$\|\mathbf{U}^\top \mathbf{S}^\top \mathbf{S} \mathbf{U} - \mathbf{I}\|_2 \leq \frac{1}{2}.$$

Hence,

$$\|\tilde{\mathbf{y}} - \mathbf{y}^*\|_2 \leq$$

For a subspace embedding \mathbf{S} , we have

$$\|\mathbf{U}^\top \mathbf{S}^\top \mathbf{S} \mathbf{U} - \mathbf{I}\|_2 \leq \frac{1}{2}.$$

Hence,

$$\|\tilde{\mathbf{y}} - \mathbf{y}^*\|_2 \leq$$

By normal equation, we have

$$\mathbf{U}^\top \mathbf{S}^\top \mathbf{S} \mathbf{U} \tilde{\mathbf{y}} = \mathbf{U} \mathbf{S}^\top \mathbf{S} \mathbf{b},$$

so,

$$\|\tilde{\mathbf{y}} - \mathbf{y}^*\|_2 \leq 2\|\mathbf{U}^\top \mathbf{S}^\top \mathbf{S}(\mathbf{U} \mathbf{y}^* - \mathbf{b})\|_2.$$

For a subspace embedding \mathbf{S} , we have

$$\|\mathbf{U}^\top \mathbf{S}^\top \mathbf{S} \mathbf{U} - \mathbf{I}\|_2 \leq \frac{1}{2}.$$

Hence,

$$\|\tilde{\mathbf{y}} - \mathbf{y}^*\|_2 \leq$$

By normal equation, we have

$$\mathbf{U}^\top \mathbf{S}^\top \mathbf{S} \mathbf{U} \tilde{\mathbf{y}} = \mathbf{U} \mathbf{S}^\top \mathbf{S} \mathbf{b},$$

so,

$$\|\tilde{\mathbf{y}} - \mathbf{y}^*\|_2 \leq 2\|\mathbf{U}^\top \mathbf{S}^\top \mathbf{S}(\mathbf{U} \mathbf{y}^* - \mathbf{b})\|_2.$$

For \mathbf{S} with the choice of m , we have

$$\Pr \left[\|\mathbf{U}^\top \mathbf{S}^\top \mathbf{S}(\mathbf{U} \mathbf{y}^* - \mathbf{b})\|_F \geq 3 \frac{\sqrt{\epsilon}}{d} \|\mathbf{U}\|_F \|\mathbf{U} \mathbf{y}^* - \mathbf{b}\|_F \right] \leq \delta.$$

Sampling for least squares

- We can consider sampling rows of $[\mathbf{A} \ \mathbf{b}]$.
- Recall leverage scores.

Leverage scores

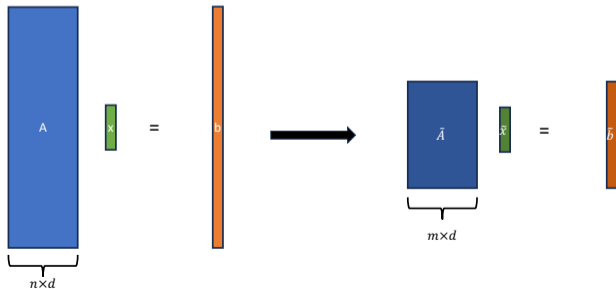
Given $\mathbf{A} \in \mathbb{R}^{n \times d}$, and an orthonormal basis \mathbf{U} for $\text{span}(\mathbf{A})$, for $i \in [n]$, the i th *leverage score*

$$\ell_i(\mathbf{A}) = \sup_{\mathbf{x}} \frac{(\mathbf{A}_{i*} \mathbf{x})^2}{\|\mathbf{A} \mathbf{x}\|^2} = \|\mathbf{U}_{i*}\|^2.$$

Sampling for least squares

Algorithm:

- Compute the row-leverage scores of \mathbf{A} , l_i , $i = 1, \dots, n$.
- Pick m rows of \mathbf{A} and the corresponding elements of \mathbf{b} with respect to the probabilities $p_i = l_i/d$ to $i \in [n]$.
- Rescale sampled rows of \mathbf{A} and sampled elements of \mathbf{b} by $1/\sqrt{mp_i}$.
- Solve the induced problem.



Leverage score sampling is subspace embedding

Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ with $r = \text{rank}(\mathbf{A})$, and $\mathbf{S} \in \mathbb{R}^{m \times n}$ be a sampling matrix with probabilities $p_i = \ell_i/r$, and $\mathbf{S}_{i*} = \mathbf{e}_j/\sqrt{mp_j}$ with $\Pr(j = i) = p_i$. If $m = O(r \log(r/\delta)/\epsilon^2)$, then \mathbf{S} is ϵ -subspace embedding of $\text{span}(\mathbf{A})$ with probability $1 - \delta$.

Leverage score sampling is subspace embedding

Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ with $r = \text{rank}(\mathbf{A})$, and $\mathbf{S} \in \mathbb{R}^{m \times n}$ be a sampling matrix with probabilities $p_i = \ell_i/r$, and $\mathbf{S}_{i*} = \mathbf{e}_j/\sqrt{mp_j}$ with $\Pr(j = i) = p_i$. If $m = O(r \log(r/\delta)/\epsilon^2)$, then \mathbf{S} is ϵ -subspace embedding of $\text{span}(\mathbf{A})$ with probability $1 - \delta$.

Proof: Let $\mathbf{U} \in \mathbb{R}^{n \times r}$ be orthonormal with $\text{span}(\mathbf{U}) = \text{span}(\mathbf{A})$.

For $k \in [m]$, let $\mathbf{X}_k = m\mathbf{U}^\top [\mathbf{S}_{k*}]^\top \mathbf{S}_{k*} \mathbf{U} - \mathbf{I}$, so

$$\frac{1}{m} \sum_k \mathbf{X}_k = \mathbf{U}^\top \mathbf{S}^\top \mathbf{S} \mathbf{U} - \mathbf{I},$$

and for ϵ -embedding, we need to bound its spectral norm.

Matrix Chernoff

Let \mathbf{X}_k for $k \in [m]$ be i.i.d copies of symmetric random $\mathbf{X} \in \mathbb{R}^{r \times r}$ with $\gamma, \sigma^2 > 0$, $\mathbb{E}[\mathbf{X}] = 0$, $\|\mathbf{X}\|_2 \leq \gamma$, and $\|\mathbb{E}[\mathbf{X}^2]\|_2 \leq \sigma^2$. Then for $\epsilon > 0$,

$$\Pr\left(\left\|\frac{1}{m} \sum_k \mathbf{X}_k\right\|_2 \geq \epsilon\right) \leq 2r \exp(-m\epsilon^2/(\sigma^2 + \gamma\epsilon/3)).$$

Apply to

$$\mathbf{X} = \frac{1}{p_j} [\mathbf{U}_{j*}]^\top \mathbf{U}_{j*} - \mathbf{I} \text{ with } \Pr(j = i) = p_i = \ell_i/r = \|\mathbf{U}_{i*}\|_2^2/r.$$

We have

$$\mathbb{E}[\mathbf{X}] =$$

$$\|\mathbf{X}\|_2 \leq$$

$$\mathbb{E}[\mathbf{X}^2] =$$

Matrix Chernoff

Let \mathbf{X}_k for $k \in [m]$ be i.i.d copies of symmetric random $\mathbf{X} \in \mathbb{R}^{r \times r}$ with $\gamma, \sigma^2 > 0$, $\mathbb{E}[\mathbf{X}] = 0$, $\|\mathbf{X}\|_2 \leq \gamma$, and $\|\mathbb{E}[\mathbf{X}^2]\|_2 \leq \sigma^2$. Then for $\epsilon > 0$,

$$\Pr\left(\left\|\frac{1}{m} \sum_k \mathbf{X}_k\right\|_2 \geq \epsilon\right) \leq 2r \exp(-m\epsilon^2/(\sigma^2 + \gamma\epsilon/3)).$$

Apply to

$$\mathbf{X} = \frac{1}{p_j} [\mathbf{U}_{j*}]^\top \mathbf{U}_{j*} - \mathbf{I} \text{ with } \Pr(j = i) = p_i = \ell_i/r = \|\mathbf{U}_{i*}\|_2^2/r.$$

We have

$$\mathbb{E}[\mathbf{X}] =$$

$$\|\mathbf{X}\|_2 \leq$$

$$\mathbb{E}[\mathbf{X}^2] =$$

so, $\|\mathbb{E}[\mathbf{X}^2]\|_2 \leq r - 1$.

Computing the leverage scores

- To compute the leverage scores exactly, we need U , i.e., compute the SVD of \mathbf{A} .
- Naive cost $O(nd^2)$.
- Can be approximately estimated in $O(nnz(\mathbf{A}) \log n + d^3)$ time.

Algorithm:

Given $\mathbf{A} \in \mathbb{R}^{n \times d}$, a subspace ϵ -embedding $\mathbf{S}_1 \in \mathbb{R}^{m \times n}$ for \mathbf{A} , and a JL matrix $\mathbf{S}_2 \in \mathbb{R}^{d \times m'}$ so that $\|\mathbf{x}^\top \mathbf{S}_2\| = (1 \pm \epsilon)\|\mathbf{x}\|$ for n vectors, so $m' = O(\log(n)/\epsilon^2)$, then:

- 1 $\mathbf{W} = \mathbf{S}_1 \mathbf{A}$; // compute sketch
- 2 $[\mathbf{Q}, \mathbf{R}] = qr(\mathbf{W})$; // change of basis
- 3 $\mathbf{Z} = \mathbf{A}(\mathbf{R}^{-1} \mathbf{S}_2)$; // sketch of $\mathbf{A} \mathbf{R}^{-1}$
- 4 return $\|\mathbf{Z}_{i*}\|_2^2$ for $i \in [n]$

Correctness

- \mathbf{AR}^{-1} has singular values in $[1 - \epsilon, 1 + \epsilon]$.
For all \mathbf{x} , $\|\mathbf{AR}^{-1}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{S}_1\mathbf{AR}^{-1}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{Q}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{x}\|$
- Let \mathbf{U} be orthonormal with $\text{span}(\mathbf{U}) = \text{span}(\mathbf{A})$.
- \mathbf{AR}^{-1} is like \mathbf{U} .

Correctness

- \mathbf{AR}^{-1} has singular values in $[1 - \epsilon, 1 + \epsilon]$.
For all \mathbf{x} , $\|\mathbf{AR}^{-1}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{S}_1\mathbf{AR}^{-1}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{Q}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{x}\|$
- Let \mathbf{U} be orthonormal with $\text{span}(\mathbf{U}) = \text{span}(\mathbf{A})$.
- \mathbf{AR}^{-1} is like \mathbf{U} .
- Pick \mathbf{T} such that $\mathbf{AR}^{-1}\mathbf{T} = \mathbf{U}$.
- \mathbf{T} has singular values $(1 \pm \epsilon)$.

Correctness

- \mathbf{AR}^{-1} has singular values in $[1 - \epsilon, 1 + \epsilon]$.

For all \mathbf{x} , $\|\mathbf{AR}^{-1}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{S}_1\mathbf{AR}^{-1}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{Q}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{x}\|$

- Let \mathbf{U} be orthonormal with $\text{span}(\mathbf{U}) = \text{span}(\mathbf{A})$.

- \mathbf{AR}^{-1} is like \mathbf{U} .

- Pick \mathbf{T} such that $\mathbf{AR}^{-1}\mathbf{T} = \mathbf{U}$.

- \mathbf{T} has singular values $(1 \pm \epsilon)$.

For all \mathbf{x} ,

$$\|\mathbf{T}\mathbf{x}\| = \|\mathbf{Q}\mathbf{T}\mathbf{x}\| = \|\mathbf{S}_1\mathbf{AR}^{-1}\mathbf{T}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{AR}^{-1}\mathbf{T}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{U}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{x}\|$$

- Then \mathbf{T}^{-1} has singular values $(1 \pm 2\epsilon)$ for $\epsilon < 1/2$.

Correctness

- \mathbf{AR}^{-1} has singular values in $[1 - \epsilon, 1 + \epsilon]$.

For all \mathbf{x} , $\|\mathbf{AR}^{-1}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{S}_1\mathbf{AR}^{-1}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{Q}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{x}\|$

- Let \mathbf{U} be orthonormal with $\text{span}(\mathbf{U}) = \text{span}(\mathbf{A})$.

- \mathbf{AR}^{-1} is like \mathbf{U} .

- Pick \mathbf{T} such that $\mathbf{AR}^{-1}\mathbf{T} = \mathbf{U}$.

- \mathbf{T} has singular values $(1 \pm \epsilon)$.

For all \mathbf{x} ,

$$\|\mathbf{T}\mathbf{x}\| = \|\mathbf{Q}\mathbf{T}\mathbf{x}\| = \|\mathbf{S}_1\mathbf{AR}^{-1}\mathbf{T}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{AR}^{-1}\mathbf{T}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{U}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{x}\|$$

- Then \mathbf{T}^{-1} has singular values $(1 \pm 2\epsilon)$ for $\epsilon < 1/2$.

- Hence, our output $\|\mathbf{e}_i^\top \mathbf{AR}^{-1}\mathbf{S}_2\|^2 = (1 \pm O(\epsilon))\|\mathbf{e}_i^\top \mathbf{U}\|^2$.

- \mathbf{AR}^{-1} has singular values in $[1 - \epsilon, 1 + \epsilon]$.

For all \mathbf{x} , $\|\mathbf{AR}^{-1}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{S}_1\mathbf{AR}^{-1}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{Q}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{x}\|$

- Let \mathbf{U} be orthonormal with $\text{span}(\mathbf{U}) = \text{span}(\mathbf{A})$.

- \mathbf{AR}^{-1} is like \mathbf{U} .

- Pick \mathbf{T} such that $\mathbf{AR}^{-1}\mathbf{T} = \mathbf{U}$.

- \mathbf{T} has singular values $(1 \pm \epsilon)$.

For all \mathbf{x} ,

$$\|\mathbf{T}\mathbf{x}\| = \|\mathbf{Q}\mathbf{T}\mathbf{x}\| = \|\mathbf{S}_1\mathbf{AR}^{-1}\mathbf{T}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{AR}^{-1}\mathbf{T}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{U}\mathbf{x}\| = (1 \pm \epsilon)\|\mathbf{x}\|$$

- Then \mathbf{T}^{-1} has singular values $(1 \pm 2\epsilon)$ for $\epsilon < 1/2$.

- Hence, our output $\|\mathbf{e}_i^\top \mathbf{AR}^{-1}\mathbf{S}_2\|^2 = (1 \pm O(\epsilon))\|\mathbf{e}_i^\top \mathbf{U}\|^2$.

$$\begin{aligned}\|\mathbf{e}_i^\top \mathbf{AR}^{-1}\mathbf{S}_2\|^2 &= (1 \pm \epsilon)\|\mathbf{e}_i^\top \mathbf{AR}^{-1}\|^2 = (1 \pm \epsilon)\|\mathbf{e}_i^\top \mathbf{U}\mathbf{T}^{-1}\|^2 \\ &= (1 \pm \epsilon)(1 \pm 2\epsilon)\|\mathbf{e}_i^\top \mathbf{U}\|^2\end{aligned}$$

Computational cost

- 1 $\mathbf{W} = \mathbf{S}_1 \mathbf{A};$ // $O(\text{nnz}(\mathbf{A})s)$
- 2 $[\mathbf{Q}, \mathbf{R}] = \text{qr}(\mathbf{W});$ // $O(d^2 m)$
- 3 $\mathbf{Z} = \mathbf{A}(\mathbf{R}^{-1} \mathbf{S}_2);$ // $O(d^2 m' + \text{nnz}(\mathbf{A})m')$
- 4 return $\|\mathbf{Z}_{i*}\|_2^2$ for $i \in [n]$ // $O(nm')$

Computational cost

- 1 $\mathbf{W} = \mathbf{S}_1 \mathbf{A};$ // $O(nnz(\mathbf{A})s)$
- 2 $[\mathbf{Q}, \mathbf{R}] = qr(\mathbf{W});$ // $O(d^2m)$
- 3 $\mathbf{Z} = \mathbf{A}(\mathbf{R}^{-1}\mathbf{S}_2);$ // $O(d^2m' + nnz(\mathbf{A})m')$
- 4 return $\|\mathbf{Z}_{i*}\|_2^2$ for $i \in [n]$ // $O(nm')$

If \mathbf{A} is dense, we use SRHT and fast JL.

If \mathbf{A} is sparse, we can use OSNAP.

Total cost is :

$$O(nnz(\mathbf{A})(m' + s) + d^2(m + m')) = O((nnz(\mathbf{A}) \log n + d^3 \log d)/\epsilon^2).$$

Further Reading:

Drineas, Petros, et al. "Fast approximation of matrix coherence and statistical leverage." The Journal of Machine Learning Research 13.1 (2012): 3475-3506.

Preconditioning for least squares

- Solving least squares regression exactly requires $O(nd^2 + d^3)$ cost.
- Using sketching or sampling : $O((nnz(\mathbf{A}) \log n + d^3 \log d)/\epsilon)$.
- However, we only get an approximate solution:

$$\|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{b}\|_2 \leq (1 + \epsilon)\|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|_2$$

- For *machine precision regression*, we need reduce the dependence on ϵ to logarithmic.
- With iterative methods, such as the general class of Krylov or conjugate-gradient type algorithms :

$$\frac{\|\mathbf{A}(\mathbf{x}^{(m)} - \mathbf{x}^*)\|_2^2}{\|\mathbf{A}(\mathbf{x}^{(0)} - \mathbf{x}^*)\|_2^2} \leq 2 \left(\frac{\sqrt{\kappa(\mathbf{A}^\top \mathbf{A})} - 1}{\sqrt{\kappa(\mathbf{A}^\top \mathbf{A})} + 1} \right)^m .$$

So, need $m = O(\kappa(\mathbf{A}) \log(1/\epsilon))$ to get an ϵ error.

Preconditioning for least squares

- Pre-conditioning reduces the number of iterations needed for a given accuracy.
- Find a non-singular matrix \mathbf{R} , such that $\kappa((\mathbf{A}\mathbf{R}^{-1})^\top \mathbf{A}\mathbf{R}^{-1})$ is small.
- Applying CG method to $\mathbf{A}\mathbf{R}^{-1}$ would converge quickly.

Preconditioning for least squares

- Pre-conditioning reduces the number of iterations needed for a given accuracy.
- Find a non-singular matrix \mathbf{R} , such that $\kappa((\mathbf{A}\mathbf{R}^{-1})^\top \mathbf{A}\mathbf{R}^{-1})$ is small.
- Applying CG method to $\mathbf{A}\mathbf{R}^{-1}$ would converge quickly.
- Idea is similar to approximate leverage scores computation.

Apply a (sparse) subspace embedding matrix \mathbf{S} to \mathbf{A} .

Compute \mathbf{R} as $[\mathbf{Q}, \mathbf{R}] = qr(\mathbf{S}\mathbf{A})$.

We know that $\mathbf{A}\mathbf{R}^{-1}$ has singular values in $[1 - \epsilon_0, 1 + \epsilon_0]$ (almost orthonormal).

$$\kappa(\mathbf{A}\mathbf{R}^{-1}) \leq \frac{1 + \epsilon_0}{1 - \epsilon_0}.$$

After m iterations of CG, we have: $\|\mathbf{A}\mathbf{R}^{-1}(\mathbf{x}^{(m)} - \mathbf{x}^*)\|^2 \leq 2\epsilon_0^m \|\mathbf{x}^*\|^2$

Iterative Refinement

Given $\mathbf{A} \in \mathbb{R}^{n \times d}$, $\mathbf{b} \in \mathbb{R}^n$, and a subspace ϵ_0 -embedding $\mathbf{S} \in \mathbb{R}^{m \times n}$ for \mathbf{A} ,

- 1 $m = O(\log(1/\epsilon))$
- 2 $\mathbf{W} = \mathbf{S}\mathbf{A}$;
- 3 $[\mathbf{Q}, \mathbf{R}] = qr(\mathbf{W})$;
- 4 $\mathbf{x}^{(0)} \leftarrow \mathbf{0}$;
- 5 for $j = 0, 1, \dots, m$:

$$\mathbf{x}^{(j+1)} \leftarrow \mathbf{x}^{(j)} + (\mathbf{R}^\top)^{-1} \mathbf{A}^\top (\mathbf{b} - \mathbf{A}\mathbf{R}^{-1} \mathbf{x}^{(j)})$$

- 6 return $\mathbf{R}^{-1} \mathbf{x}^{(m+1)}$

Cost:

For SRHT or OSNAP: $O(nnz(\mathbf{A}) \log(n/\epsilon) + d^3 \log^2 d + d^2 \log(1/\epsilon))$

For Countsketch: $O((nnz(\mathbf{A}) + d^4) \log(1/\epsilon))$.

Sketch based preconditioning

Let $\mathbf{x}^{(j+1)} \leftarrow \mathbf{x}^{(j)} + (\mathbf{R}^\top)^{-1} \mathbf{A}^\top (\mathbf{b} - \mathbf{A} \mathbf{R}^{-1} \mathbf{x}^{(j)})$.

We have

$$\begin{aligned} \mathbf{A} \mathbf{R}^{-1} (\mathbf{x}^{(j+1)} - \mathbf{x}^*) &= \mathbf{A} \mathbf{R}^{-1} \left(\mathbf{x}^{(j)} + (\mathbf{R}^\top)^{-1} \mathbf{A}^\top (\mathbf{b} - \mathbf{A} \mathbf{R}^{-1} \mathbf{x}^{(j)}) - \mathbf{x}^* \right) \\ &= \\ &= \end{aligned}$$

Sketch based preconditioning

Let $\mathbf{x}^{(j+1)} \leftarrow \mathbf{x}^{(j)} + (\mathbf{R}^\top)^{-1} \mathbf{A}^\top (\mathbf{b} - \mathbf{A} \mathbf{R}^{-1} \mathbf{x}^{(j)})$.

We have

$$\begin{aligned} \mathbf{A} \mathbf{R}^{-1} (\mathbf{x}^{(j+1)} - \mathbf{x}^*) &= \mathbf{A} \mathbf{R}^{-1} \left(\mathbf{x}^{(j)} + (\mathbf{R}^\top)^{-1} \mathbf{A}^\top (\mathbf{b} - \mathbf{A} \mathbf{R}^{-1} \mathbf{x}^{(j)}) - \mathbf{x}^* \right) \\ &= \\ &= \end{aligned}$$

where $\mathbf{A} \mathbf{R}^{-1} = \mathbf{U} \Sigma \mathbf{V}^\top$. We know $\mathbf{A} \mathbf{R}^{-1}$ has singular values in $[1 - \epsilon_0, 1 + \epsilon_0]$.

So, diagonal entries of $\Sigma - \Sigma^3$ are at most $\sigma_i (1 - (1 - \epsilon_0)^2) \leq 3\sigma_i \epsilon_0$ for $\epsilon_0 \leq 1$. Hence,

$$\|\mathbf{A} \mathbf{R}^{-1} (\mathbf{x}^{(m+1)} - \mathbf{x}^*)\| \leq 3\epsilon_0 \|\mathbf{A} \mathbf{R}^{-1} (\mathbf{x}^{(m)} - \mathbf{x}^*)\|$$

and by choosing $\epsilon_0 = 1/2$, say, $O(\log(1/\epsilon))$ iterations suffice to attain ϵ relative error.

Further Reading

- Avron, Haim, Petar Maymounkov, and Sivan Toledo. “Blendenpik: Supercharging LAPACK’s least-squares solver.” *SIAM Journal on Scientific Computing* 32.3 (2010): 1217-1236.
- Clarkson, Kenneth L., and David P. Woodruff. “Low-rank approximation and regression in input sparsity time.” *Journal of the ACM (JACM)* 63.6 (2017): 1-45.

Questions?