

Geometric view of $Ax=b$. Let $a_1^T \dots a_n^T$ be rows of A :

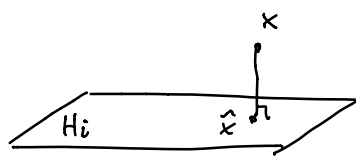
$$\begin{pmatrix} -a_1^T \\ \vdots \\ -a_n^T \end{pmatrix} \text{ So } Ax=b \Leftrightarrow a_i^T x = b_i \text{ for all } i.$$

$a_i^T x = b_i$ is the equation defining a hyperplane in \mathbb{R}^n .

Some people call it an affine subspace.

Kaczmarz: if there is an i s.t. $a_i^T x \neq b_i$,

move x to the closest \hat{x} s.t. $a_i^T \hat{x} = b_i$.



$$H_i = \{ \gamma : a_i^T \gamma = b_i \}$$

Obtain $\hat{x} = x - \gamma a_i$, for some γ

Solve for γ by $a_i^T \hat{x} = b_i \Rightarrow a_i^T (x - \gamma a_i) = b_i$

$$\Rightarrow \gamma a_i^T a_i = a_i^T x - b_i \Rightarrow \gamma = \frac{a_i^T x - b_i}{\|a_i\|^2}$$

Claim (not necessary, but useful)

$x - \gamma a_i$ is closest point to x on H_i .

$$\begin{aligned} \text{proof 1: } \text{dist}(x, x - \gamma a_i) &= \|\gamma a_i\|_2 = \frac{|a_i^T x - b_i|}{\|a_i\|_2^2} \cdot \|a_i\|_2 \\ &= \frac{|a_i^T x - b_i|}{\|a_i\|_2} \quad (1) \end{aligned}$$

whereas if $x - \delta \in H_i$,

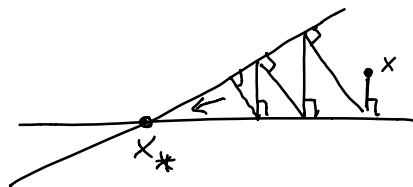
$$a_i^T (x - \delta) = b_i \Rightarrow a_i^T x - b_i = a_i^T \delta$$

$$\Rightarrow |a_i^T x - b_i| \leq \|a_i\|_2 \|\delta\|_2$$

$$\Rightarrow \frac{|a_i^T x - b_i|}{\|a_i\|_2} \leq \|\delta\|_2$$

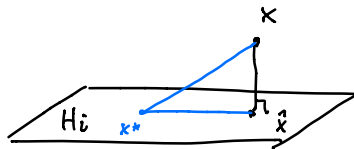
proof 2: $a_i / \|a_i\|$ is unit normal to H_i

kaczmarz. Do this repeatedly.



Let x_* be solution to $Ax_* = b$.

How much does $\|x - x_*\|_2$ decrease?



Claim $\|x - x_*\|_2^2 = \|x_* - \hat{x}\|_2^2 + \|x - \hat{x}\|_2^2$

by Pythagorean theorem because $x_* - \hat{x} \perp x - \hat{x}$

Here is an algebraic proof: $x - \hat{x} = \gamma a_i$,

whereas $\gamma a_i^T (x_* - \hat{x}) = \gamma (b_i - b_i) = 0$

So, (i) above implies

$$\|\hat{x} - x_*\|_2^2 = \|x - x_*\|_2^2 - \frac{(a_i^T x - b_i)^2}{\|a_i\|_2^2}$$

Leads to question: how choose i .

Originally, people tried $1, 2, \dots, n$,

Strohmer & Vershynin '06

suggested a non-uniform random choice

pick i with probability proportional to $\|a_i\|_2^2$

Leads to a very clean analysis.

Recall $\|A\|_F^2 = \sum_i \|a_i\|_2^2$.

So, set prob of i to $P_i = \frac{\|a_i\|_2^2}{\|A\|_F^2}$

$$\begin{aligned} \text{Then } \mathbb{E} \left[\|x - x_*\|_2^2 - \|\hat{x} - x_*\|_2^2 \right] &= \sum_i P_i \frac{(a_i^T x - b_i)^2}{\|a_i\|_2^2} = \sum_i \frac{\|a_i\|_2^2}{\|A\|_F^2} \frac{(a_i^T x - a_i^T x_*)^2}{\|a_i\|_2^2} \\ &= \frac{1}{\|A\|_F^2} \sum_i (a_i^T x - a_i^T x_*)^2 = \frac{1}{\|A\|_F^2} \|A(x - x_*)\|_2^2 \\ &\geq \frac{\sigma_n(A)}{\|A\|_F^2} \|x - x_*\|_2^2 = \frac{1}{\|A\|_F^2 \|A^{-1}\|_2^2} \|x - x_*\|_2^2 \end{aligned}$$

Implies $\mathbb{E} [\| \hat{X} - X_* \|_2^2] \leq \left(1 - \frac{1}{\|A\|_F^2 \|A^{-1}\|_2^2} \right) \|X - X_*\|_2^2$

If start at $x_0 = 0$, and x_t is t th iterate,

$$\mathbb{E} [\|x_t - X_*\|_2^2] \leq \left(1 - \frac{1}{\|A\|_F^2 \|A^{-1}\|_2^2} \right)^t \|X_*\|_2^2$$

Time: when pick i , need to compute $a_i^T x$
and subtract αa_i from x ,
which takes time $\Theta(n_i)$,
where $n_i = \#$ non-zero entries in a_i .

Let's compare this to GD and its improvement CG
each step requires mult by A , which takes time $\sum_i n_i$

$$\text{and (GD)}: \|A(x - x_*)\|_2^2 \leq \left(1 - \frac{1}{\|A\|_2^2 \cdot \|A^{-1}\|_2^2} \right)^t \|x - x_*\|_2^2$$

$$\text{CG} \leq 2 \left(1 - \frac{2}{\|A\|_2 \|A^{-1}\|_2} \right)^{2t} \|x - x_*\|_2^2$$

So, CG needs fewer iterations, but each takes longer.

Sometimes Kaczmarz is a win.

(Note: most analyses of CG (conjugate gradients)
state bounds in terms of $\kappa(A^T A)$, whereas we
are using $\kappa(A)$)

Understanding perturbed matrices.

Consider $B = A + \varepsilon R$ where R has independent $\mathcal{N}(0,1)$ entries Normal / Gaussian

Recall x is $\mathcal{N}(0, \sigma^2)$ - normal, mean 0, variance σ^2
has density $p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-x^2/2\sigma^2}$

If x is $\mathcal{N}(0,1)$, εx is $\mathcal{N}(0, \varepsilon^2)$.

And if $x(1), \dots, x(n)$ are independent $\mathcal{N}(0,1)$,

then $t^T x = \sum_i t(i) x(i)$ is $\mathcal{N}(0, \|t\|_2^2)$

It's unlikely $\kappa(B)$ is big

$$\Pr \left[\kappa(A + \varepsilon R) \geq (\|A\|_2 + \varepsilon\sqrt{n}) / \lambda \right] \leq \sqrt{n} \lambda / \varepsilon$$

we prove with $n^{3/2}$

$\|A + \varepsilon R\|_2$ is approximately $\|A\|_2$, plus $\varepsilon\sqrt{n}$.

$\sigma_n(A + \varepsilon R)$ is unlikely to be small.

For $B = \begin{pmatrix} b_1 & \dots & b_n \\ \vdots & & \vdots \end{pmatrix}$ define

$$\text{height}(i) = \text{dist}(b_i, \text{span}(b_j : j \neq i))$$

$$= t^T b_i \text{ where } t \text{ is the unit normal to } \text{span}(b_j : j \neq i)$$

$$\text{i.e. } t^T b_j = 0 \text{ for } j \neq i, \quad \|t\| = 1$$

Geometry lem 1 $\exists i$ s.t. $\text{height}(i) \leq \sqrt{n} \sigma_n(B)$

proof let v be s.t. $\|v\|_2 = 1$ and $\|Bv\|_2 = \sigma_n$

Let i be s.t. $|v(i)| = \|v\|_\infty$, so $|v(i)| \geq 1/\sqrt{n}$.

Now $\left\| \sum_j v(j) b_j \right\| = \sigma_n$

$$\left\| \sum_j \frac{v(j)}{|v(i)|} b_j \right\| = \frac{\sigma_n}{|v(i)|} \leq \sqrt{n} \sigma_n$$

$$\left\| b_i - \sum_{j \neq i} \frac{v(j)}{|v(i)|} b_j \right\| \text{ and } \sum_{j \neq i} \frac{v(j)}{|v(i)|} b_j \in \text{span}(b_j : j \neq i)$$

Probability lem 2 for each i ,

$$\Pr[\text{height}(i) \leq \lambda] \leq \lambda/\varepsilon$$

proof of lem Sample b_j for $j \neq i$. Having fixed these,

let t be the unit normal to $\text{span}(b_j : j \neq i)$

so $\text{height}(i) = t^T b_i$

$$b_i = a_i + \varepsilon \tau_i \text{ so } \Pr[\text{height}(i) \leq \lambda] = \Pr[|t^T b_i| \leq \lambda]$$

$$= \Pr\left[\varepsilon t^T \tau_i \in [-t^T a_i - \lambda, -t^T a_i + \lambda] \right]$$

$$= \frac{1}{\sqrt{2\pi} \varepsilon} \int_{-t^T a_i - \lambda}^{-t^T a_i + \lambda} e^{-x^2/2\varepsilon^2} dx \leq \frac{2\lambda}{\sqrt{2\pi} \varepsilon} \leq \lambda/\varepsilon$$

Thm $\Pr[\sigma_n(B) \leq \lambda] \leq n^{3/2} \lambda / \varepsilon$

proof lem 1 \Rightarrow

$$\Pr[\sigma_n(B) \leq \lambda] \leq \Pr[\exists i \text{ s.t. } \text{height}(i) \leq \sqrt{n} \lambda]$$

$$\leq \sum_i \Pr[\text{height}(i) \leq \sqrt{n} \lambda]$$

$$\leq \sum_i \sqrt{n} \lambda / \varepsilon, \text{ by lem 2}$$

$$\leq n^{3/2} \lambda / \varepsilon$$

Note: Improved to $2.35 \sqrt{n} \lambda / \varepsilon$ in SST'06
and a tight $\sqrt{n} \lambda / \varepsilon$ is BKMS'19

Plot it?