# Iterative Methods for Solving Least-squares

When $A$ has full column rank, our L.S. estimate is

$$\hat{x} = (A^T A)^{-1} A^T y$$

If $A$ is $M \times N$, then constructing $A^T A$ costs $O(MN^2)$ computations, and inverting $A^T A$ costs $O(N^3)$ computations. (Note that for $M > N$, the cost of constructing the matrix is actually larger than inverting it.)

This is a problem for even moderately large $M$ and $N$. But inverse problems with large $M$ & $N$ are common in the modern world.

Typical 3D MRI scan:

    Reconstruct a $128 \times 128 \times 128$ cube of voxels from $\approx 5 \cdot 10^6$ Fourier-domain samples

    $N \approx 2.1$ million , $M \approx 5$ million

It takes about .25 seconds for me to construct $A^T A$ and solve $A^T A x = A^T y$ when

$$M = 5000, \quad N = 1000$$

How long would it take (approx) for the example above? ( $M = 5,000,000$ $N = 2,000,000$ )

How much memory would you need to hold $A$ (assume double precision = 8 bytes per entry)?

In this section, we will overview two iterative methods — steepest descent & conjugate gradients — that reformulate

$$A^T A x = A^T y$$

as an optimization program. Each iteration is simple, and requires one application of $A$ & one application of $A^T$. If $A^T A$ is well conditioned, they can converge in very few iterations (especially CG). This makes the cost of solving this type of L.S. problem dramatically smaller — about the cost of a few to a couple hundred applications of $A$.

Moreover, we don't need to construct $A^T A$ or even $A$ explicitly, all we need is a "black box" which takes a vector $x$ and returns $Ax$. This is especially useful if it takes $\ll O(MN)$ time to apply $A$ or $A^T$

23

In the MRI example above, it takes about 30 seconds to apply $A$ or $A^T$ using a unequispaced FFT. CG (which we will learn about soon) converges in about 50 iterations $\Rightarrow$ the problem is solved in $\approx 50$ minutes. Also, the storage requirements are $O(M+N)$.

## Recasting as an optimization program

We want to solve

$$\underbrace{A^T A}\, x = \underbrace{A^T y}$$

or $\qquad H\, x = b \qquad\qquad H : N \times N$

Note $H$ is symmetric + def (if $A$ has full column rank).

Since $H$ is sym + def, then the solution $\hat{x}$ to

$$\min_{x \in \mathbb{R}^N} \tfrac{1}{2} x^T H x - x^T b \qquad\qquad (QP)$$

will satisfy $\quad H\hat{x} = b$.

24

Why? Well, if $H$ is sym+def, then (QP) is convex (strictly convex, actually). Thus a necessary and sufficient condition for $\hat{x}$ to be the minimizer is

$$\nabla_x \left( \tfrac{1}{2} x^T H x - x^T b \right) \Big|_{x=\hat{x}} = 0$$

$\uparrow$ gradient w.r.t. $x$

Since

$$\nabla_x \left( \tfrac{1}{2} x^T H x - x^T b \right) = \tfrac{1}{2} \nabla_x (x^T H x) - \nabla_x (x^T b)$$

$$= H x - b$$

the solution to (QP) must have

$$H \hat{x} = b.$$

25

# Steepest Descent

Say you have an unconstrained optimization program

$$\min_{x \in \mathbb{R}^N} f(x)$$

where $f(x): \mathbb{R}^N \to \mathbb{R}$ is convex. One way to solve this program is to simply "roll downhill". That is, from a starting point $x_0$ we move to

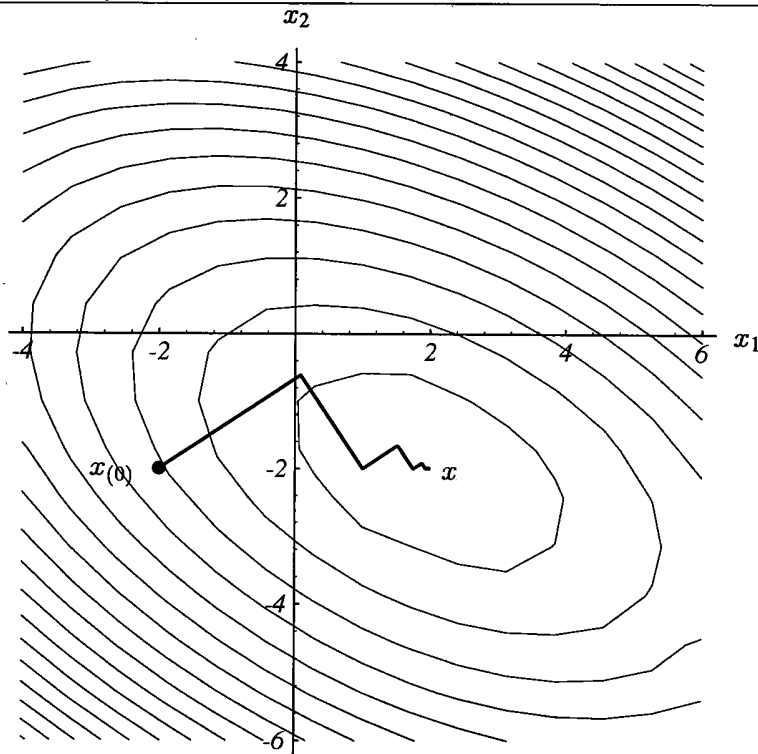$$X_1 = X_0 - \alpha_0 \cdot \nabla f(x) \big|_{x=x_0}$$

then to

$$X_2 = X_1 - \alpha_1 \nabla f(x) \big|_{x=x_1}$$
$$\vdots$$
$$X_K = X_{K-1} - \alpha_{K-1} \nabla f(x) \big|_{x=x_{K-1}}$$

for some appropriate $\alpha_1, \alpha_2, \ldots ; \alpha_K > 0$. At each iteration we are moving in the direction $-\nabla f(x_K)$; this is the direction of "steepest descent".

26

For our problem
$$\min_{x} \quad x^T H x - x^T b$$
the gradient is simply the residual
$$-\nabla \left( x^T H x - x^T b \right)\Big|_{x=x_k} = b - H x_k =: r_k$$
and so the iteration becomes
$$x_{k+1} = x_k + \alpha_k r_k$$

27

There is a nifty way to choose an optimal value for the $\alpha_K$. We want to choose $\alpha_K$ so that $f(x_{K+1})$ is as small as possible. Along indexed by $\alpha_K$ the ray $x_K + \alpha_K r_K$, $f(x_K + \alpha_K r_K)$ is again convex (in $\alpha_K$) so we want

$$\frac{d}{d\alpha_K} f(x_K + \alpha_K r_K) = 0$$

By the chain rule

$$\frac{d}{d\alpha_K} f(x_{K+1}) = \nabla f(x_{K+1})^T \frac{d}{d\alpha} x_{K+1}$$

$$= \nabla f(x_{K+1})^T r_K$$

So we need to choose $\alpha_K$ such that

$$\nabla f(x_{K+1}) \perp r_K$$

or more concisely

$$r_{K+1} \perp r_K \qquad \left( r_{K+1}^T r_K = 0 \right)$$

28

So let's do this:

$$r_{k+1}^T \, r_k = 0$$

$$\Rightarrow \quad (b - Hx_{k+1})^T r_k = 0$$

$$\Rightarrow \quad \left(b - H(x_k + \alpha_k r_k)\right)^T r_k = 0$$

$$\Rightarrow \quad (b - Hx_k)^T r_k - \alpha_k \, r_k^T H r_k = 0$$

$$\Rightarrow \quad r_k^T r_k - \alpha_k \, r_k^T H r_k = 0$$

$$\Rightarrow \quad \alpha_k = \frac{r_k^T r_k}{r_k^T H r_k}$$

So the steepest descent algorithm is

- Initialize $x_k$ = some guess, $k = 0$
- while not converged

$$\left[\begin{array}{l} r_k = b - Hx_k \\[4pt] \alpha_k = r_k^T r_k / r_k^T H r_k \\[4pt] x_{k+1} = x_k + \alpha_k r_k \\[4pt] k = k+1 \end{array}\right.$$

29

Notice that
$$r_{k+1} = b - Hx_{k+1} = b - H(x_k + \alpha_k r_k)$$
$$= r_k - \alpha_k H r_k$$

So we can save an application of $H$ using

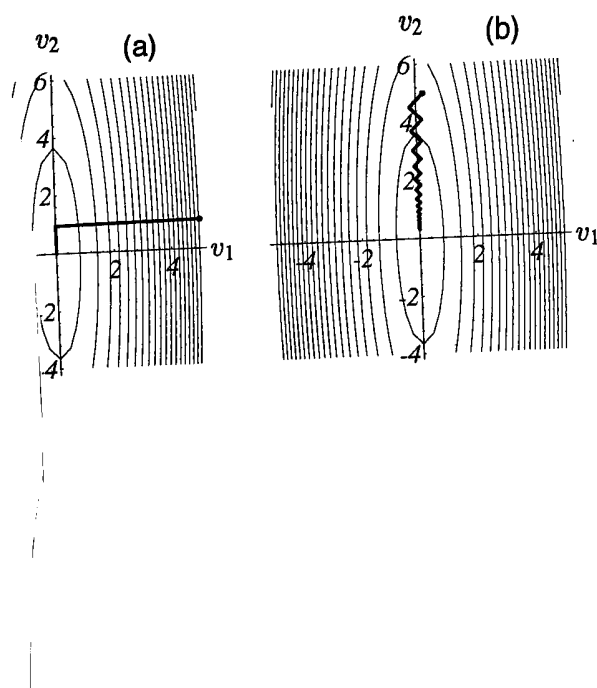Initialize $k = 0$
$$x_0 = \text{some guess}$$
$$r_0 = b - Hx_0$$

while not converged
$$q = H r_k$$
$$\alpha_k = r_k^T r_k / r_k^T q$$
$$x_{k+1} = x_k + \alpha_k r_k$$
$$r_{k+1} = r_k - \alpha_k q$$

30

The effectiveness of S.D. depends critically
on how $H'$ is conditioned and the
starting point

From Shewchuk:

# The conjugate gradient (CG) method

An excellent resource for the material in this section is

J. Shewchuk: "An introduction to the conjugate gradient method without the agonizing pain"

We can see from the example on the last page that steepest descent is sometimes inefficient because it can move in essentially the same direction many times.

CG avoids this by ensuring that each step is $\perp$ (in an appropriate inner product) to all of the previous steps that have been taken.

Consider the general iteration

$$X_{k+1} = X_k + \alpha_k \cdot d_k$$

↖ direction to move @ $k^{th}$ iteration

32

What we'd like is to choose $\alpha_K$ so that

$$e_{K+1} = X_{K+1} - X$$

the error at iteration $K+1$ is _orthogonal_ to the direction $d_K$. This would mean (thanks to the $\perp$-principle) that we are optimally aligned with $X$ (the true solution) along the direction $d_K$. Then if $d_0, d_1, \ldots, d_{N-1}$ is a set of $\perp$ directions we would only need to move in each exactly once.

We would like $\alpha_K$ s.t.

$$d_K \perp e_{K+1} \quad \text{i.e.} \quad \langle d_K, e_{K+1} \rangle = d_K^T e_{K+1} = 0$$

So we would need

$$d_K^T ( X_K + \alpha_K d_K - X ) = 0$$

$$\Rightarrow d_K^T ( e_K + \alpha_K d_K ) = 0$$

$$\Rightarrow \alpha_K = -\frac{d_K^T e_K}{d_K^T d_K}$$

This is fine, except we have no idea what $e_K$ is (if we did, we could solve the entire problem instantly).

33

what we can do is choose the $\alpha_K$ such that $d_K$ is $H-\perp$ to $e_{K+1}$. That is,

$$\langle d_K, e_{K+1} \rangle_H = d_K^T H e_{K+1} = 0$$

( Recall that if $H$ is an $N \times N$ sym+def matrix, then $\langle x, y \rangle_H = x^T H y$ is a valid inner product on $\mathbb{R}^N$. )

We need

$$d_K^T H (x_K + \alpha_K d_K - x) = 0$$

$$\Rightarrow d_K^T H e_K + \alpha_K d_K^T H d_K = 0$$

$$\Rightarrow \alpha_K = -\frac{d_K^T H e_K}{d_K^T H d_K}$$

Notice that

$$H e_K = H(x_K - x) = H x_K - b = -r_K$$

which we definitely know, so we could take

$$\alpha_K = \frac{d_K^T r_K}{d_K^T H d_K}$$

residual

34

So given a set of $H-\perp$ search directions $d_0, d_1, \ldots, d_{N-1}$, $\langle d_i, d_j \rangle_H = 0 \quad i \neq j$

We could implement the following

Initialize $x_0 =$ some guess

for $K = 0$ to $N-1$

$$\begin{cases} r_K = b - H x_K \\[2mm] \alpha_K = d_K^T r_K \big/ d_K^T H d_K \\[2mm] x_{K+1} = x_K + \alpha_K d_K \end{cases}$$

This procedure would converge to the <u>exact</u> solution after the $N^{th}$ step (more on this later).

All we need now is a set of $H-\perp$ direction vectors $d_j$. The beauty of CG is that it generates these directions "on the fly" by running what is essentially Gram-Schmidt.

35

Here is the CG algorithm:

Initialize $\quad X_0 =$ some guess

$\qquad r_0 = b - H x_0$

$\qquad d_0 = r_0$

for $\quad K = 0$ up to $N - 1$

$$\alpha_K = r_K^T r_K \Big/ d_K^T H d_K$$

$$X_{K+1} = X_K + \alpha_K d_K$$

$$r_{K+1} = r_K - \alpha_K H d_K$$

$$\beta_{K+1} = r_{K+1}^T r_{K+1} \Big/ r_K^T r_K$$

$$d_{K+1} = r_{K+1} + \beta_{K+1} d_K$$

(We will show below that $r_K^T r_K = r_K^T d_K$, so the $\alpha_K$ really is the same as before.)

These choices of $\alpha_K$ and $\beta_{K+1}$ are ensuring two important things:

① $\langle r_\ell, r_{K+1} \rangle = 0$

$\quad$ for $\ell = 0, \cdots, j$

$\quad$ "the residuals are $\perp$"

② $\langle d_\ell, d_{K+1} \rangle_H = 0$

$\quad$ for $\ell = 0, \cdots, K$

$\quad$ "the directions are $H - \perp$"

36

We can establish these two facts by induction.
We start with the following:

① $\langle r_0, r_1 \rangle = r_0^T r_1 = 0$

Since
$$r_1 = r_0 - \frac{r_0^T r_0}{r_0^T H r_0} \cdot H r_0$$

$$\Rightarrow r_0^T r_1 = r_0^T r_0 - r_0^T r_0 \frac{r_0^T H r_0}{r_0^T H r_0}$$

$$= 0$$

② $\langle d_0, d_1 \rangle_H = d_0^T H d_1 = 0$

Since
$$r_1 = r_0 - \alpha_0 H d_0$$

$$\Rightarrow r_1^T r_1 = r_1^T r_0 - \alpha_0 r_1^T H d_0$$

$$\Rightarrow r_1^T H d_0 = -\frac{1}{\alpha_0} \cdot r_1^T r_1 \qquad (\text{since } r_1^T r_0 = 0)$$

and also
$$d_1 = r_1 + \frac{r_1^T r_1}{r_0^T r_0} \cdot d_0$$

$$\Rightarrow d_0^T H d_1 = d_0^T H r_1 + \frac{r_1^T r_1}{r_0^T r_0} \cdot d_0^T H d_0$$

$$= -\frac{r_1^T r_1}{r_0^T r_0} \cdot d_0^T H d_0 + \frac{r_1^T r_1}{r_0^T r_0} \cdot d_0^T H d_0 = 0.$$

37

Now at step $k+1$, suppose we have

$$\langle r_\ell, r_j \rangle = r_\ell^T r_j = 0 \quad \forall j, \ell \le k$$

$$\langle d_\ell, d_j \rangle_H = d_\ell^T H d_j = 0 \quad \forall j, \ell \le k$$

Then we will also have the following:

① $\quad \langle r_\ell, r_{k+1} \rangle = r_\ell^T r_{k+1} = 0 \quad \forall \ell \le k$

To see this, first note that

$$r_\ell^T H d_k = (d_\ell - \beta_\ell d_{\ell-1})^T H d_k$$

$$= \begin{cases} d_k^T H d_k & \ell = k \\ 0 & \ell < k \end{cases} \quad \begin{array}{l} \text{since} \\ \langle d_\ell, d_k \rangle_H = 0 \\ \text{for } \ell < k \end{array}$$

As a result

$$r_\ell^T r_{k+1} = r_\ell^T r_k - \frac{r_k^T r_k}{d_k^T H d_k} r_\ell^T H d_k = 0 \quad \forall \ell \le k$$

38

② $\langle d_\ell, d_{K+1} \rangle_H = d_\ell^T H d_{K+1} = 0 \quad \forall \ell \leq K$

This follows from the expansion
$$d_\ell^T H d_{K+1} = d_\ell^T H r_{K+1} + \beta_{K+1} d_\ell^T H d_K$$

Notice that
$$r_i^T r_{K+1} = r_i^T r_K - \alpha_K r_i^T H d_K$$

$$\Rightarrow \quad r_i^T H d_K = \begin{cases} \frac{1}{\alpha_K} r_K^T r_K & i = K \\[2mm] -\frac{1}{\alpha_K} r_{K+1}^T r_{K+1} & i = K+1 \\[2mm] 0 & i < K \end{cases} \qquad (*)$$

Then for $\ell = K$
$$d_K^T H d_{K+1} = \frac{-1}{\alpha_K} r_{K+1}^T r_{K+1} + \beta_{K+1} d_K^T H d_K$$

$$= \underbrace{\frac{-r_{K+1}^T r_{K+1}}{r_K^T r_K} \cdot d_K^T H d_K} + \underbrace{\frac{r_{K+1}^T r_{K+1}}{r_K^T r_K} \cdot d_K^T H d_K}$$

$$= 0$$

39

and for $\ell < K$

$$d_\ell^T H d_{K+1} = \underbrace{d_\ell^T H r_{K+1}}_{\substack{=0 \text{ by } (\#) \\ \text{on last page}}} + \underbrace{\beta_{K+1} d_\ell^T H d_K}_{\substack{=0 \text{ since the } d_0, d_1, \ldots, d_K \\ \text{are } H\perp \text{ already}}}$$

So we have established that the direction $d_K$ that CG moves on iteration $K$ is $H-\perp$ to all previous directions.

Now we will establish that this means that CG will converge <u>exactly</u> in at most $N$ iterations.

Let the error at the $K^{th}$ iteration be

$$e_K = X_K - X$$

Notice that since $X_{K+1} = X_K + \alpha_K d_K$

$$\Rightarrow e_{K+1} = X_{K+1} - X$$
$$= e_K + \alpha_K d_K$$

40

Unrolling this expression, we get

$$e_K = e_0 + \sum_{j=0}^{K-1} \alpha_j d_j$$

Now remember that $d_0, \ldots, d_{N-1}$ are all H-$\perp$. Thus

$$\frac{d_0}{\|d_0\|_H}, \frac{d_1}{\|d_1\|_H}, \ldots, \frac{d_{N-1}}{\|d_{N-1}\|_H} \qquad \text{is an H-orthobasis for } \mathbb{R}^N$$

As such, we can write

$$e_0 = \sum_{j=0}^{N-1} \langle e_0, \frac{d_j}{\|d_j\|_H} \rangle_H \cdot d_j / \|d_j\|_H$$

$$= \sum_{j=0}^{N-1} \frac{\langle e_0, d_j \rangle_H}{\|d_j\|_H^2} \cdot d_j$$

$$= \sum_{j=0}^{N-1} \frac{\langle e_0 + \sum_{i=0}^{j-1} \alpha_i d_i, d_j \rangle_H}{\|d_j\|_H^2} \cdot d_j \qquad \left(\text{since } d_j \text{ are } H\text{-}\perp\right)$$

$$= \sum_{j=0}^{N-1} \frac{d_j^T H e_j}{d_j^T H d_j} \cdot d_j$$

$$= \sum_{j=0}^{N-1} \frac{-d_j^T r_j}{d_j^T H d_j} \cdot d_j \qquad \text{since } He_j = -r_j$$

41

This seems like the appropriate place to mention that $d_j^T r_j = r_j^T r_j$, since

$$d_j = r_j + \beta_j d_{j-1}$$

$$= r_j + \beta_j r_{j-1} + \beta_j \beta_{j-1} d_{j-2}$$

$$= r_j + \beta_j r_{j-1} + \beta_j \beta_{j-1} r_{j-2} + \beta_j \beta_{j-1} \beta_{j-2} d_{j-3}$$

$$\text{etc}$$

$$= r_j + \sum_{i=0}^{j-1} \gamma_i r_i \qquad \text{for some scalars } \gamma_i$$

$$\Rightarrow r_j^T d_j = r_j^T r_j + \sum_{i=j-1} \gamma_i \underbrace{r_j^T r_i}_{0}$$

Thus

$$e_0 = \sum_{j=0}^{N-1} -\alpha_j \cdot d_j$$

and

$$e_K = e_0 + \sum_{j=0}^{K-1} \alpha_j d_j$$

$$= -\sum_{j=0}^{N-1} \alpha_j d_j + \sum_{j=0}^{K-1} \alpha_j d_j$$

$$= -\sum_{j=K}^{N-1} \alpha_j d_j = -\sum_{j=K}^{N-1} \alpha_j \cdot \|d_j\|_H \cdot \frac{d_j}{\|d_j\|_H}$$

42

By Parseval (for $\langle \cdot, \cdot \rangle_H$), we have

$$\| e_K \|_H^2 = \sum_{j=K}^{N-1} \alpha_j^2 / \| d_j \|_H^2$$

This is obviously monotonically decreasing with $K$, and

$$\| e_N \|_H^2 = 0$$

$\Rightarrow$ $\boxed{X_N = X = \text{true solution} \ !}$

Since each iteration of CG is a matrix-vector multiply — which is $O(N^2)$ — and we converge in $N$ iterations, CG solves $Hx = b$ in $O(N^3)$ computations in general, the same as with other solvers.

43

BUT, if H is specially structured so that it takes $\ll O(N^2)$ computations to apply, then CG takes advantage of this. The real cost is $N$ applications of H.

In addition, it is often the case that $\|e_K\|_H^2$ is acceptably small for relatively modest values of K. This is particularly true if H is well-conditioned.

Moral: CG can get an approximate (but still potentially very good) solution using _much_ less computation than solving the system directly.

It also significantly outperforms steepest descent.

44

# Convergence Guarantees

How many iterations do we need for steepest descent or CG to converge within a certain precision? There are "worst case" bounds that depend on the condition number $\mathbb{K}$ of $H$

$$\mathbb{K} = \frac{\lambda_{max}(H)}{\lambda_{min}(H)} = \frac{max\ eigenvalue}{min\ eigenvalue}$$

For steepest descent, we will have

$$\| e_K \|_H \leq \delta \cdot \| e_0 \|_H$$

in at most

$$K \leq \left\lceil \frac{1}{2} \mathbb{K} \cdot \log(2/\delta) \right\rceil$$

iterations.

For CG, we need at most

$$K \leq \left\lceil \frac{1}{2} \sqrt{\mathbb{K}} \cdot \log(2/\delta) \right\rceil$$

(these are "natural logs")

45

See Shewchuk for a nice derivation of these.

Say the condition number is $\underline{K} = 100$. How many iterations do you need to get 6 digits of precision $(\delta = 10^{-6})$ ?

$$\text{SD:} \quad \left\lceil \frac{1}{2} \cdot 100 \cdot \log(10^6) \right\rceil = 691$$

$$\text{CG:} \quad \left\lceil \frac{1}{2} \cdot 10 \cdot \log(2 \cdot 10^6) \right\rceil = 73 \quad .$$

Again, these are worst-case bounds, and performance is typically better.

46