

Iterative methods for solving least-squares

Let's return now to our fundamental (regularized) least-squares estimators. When we are using a finite-dimensional basis, we have two forms for the estimator:

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A} + \delta \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \delta \mathbf{I})^{-1} \mathbf{y}.$$

Likewise, when we are using a kernel, the estimate map is parameterized by

$$\hat{\boldsymbol{\alpha}} = (\mathbf{K} + \delta \mathbf{I})^{-1} \mathbf{y}.$$

The first expression above requires solving an $N \times N$ system of equations, while the second require solving an $M \times M$ system (remember that M is the number of data points). The costs involved in solving these systems can be significant; if \mathbf{A} is $M \times N$, then constructing $\mathbf{A}^T \mathbf{A}$ can cost $O(MN^2)$ computations, and solving the $N \times N$ system $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{y}$ costs $O(N^3)$ computations. (Note that for $M \geq N$, the cost of constructing the matrix can actually exceed the cost to solve the system — but if we have the “kernel trick”, the cost of constructing the matrix can be greatly reduced.)

This cost can be prohibitive for even moderately large M and N . But learning and inverse problems with large M and N are common in the modern world. For example, a typical 3D MRI scan will try to reconstruct a $128 \times 128 \times 128$ cube of voxels from about 5 million non-uniformly spaced samples in the spatial Fourier domain. In this case, the matrix \mathbf{A} , which models the MRI machine, is $M = 5 \cdot 10^6$ by $N = 2.1 \cdot 10^6$.

With those values, MN^2 is huge ($\sim 10^{19}$); even storing the matrix $\mathbf{A}^T \mathbf{A}$ in memory would require terabytes of RAM.

In this section, we will overview two iterative methods — **steepest descent** and **conjugate gradients** — that solve symmetric positive-definite systems of equations

$$\mathbf{H}\mathbf{x} = \mathbf{b}$$

as an optimization program and then solve it by an iterative descent method. Each iteration is simple, requiring one or two applications of \mathbf{H} .

If \mathbf{H} is well-conditioned, then these methods can converge in very few iterations (especially conjugate gradients). This makes the cost of solving a least-squares problem dramatically smaller — about the cost of a few hundred applications of \mathbf{H} .

Moreover, we will not need to construct \mathbf{H} explicitly. All we need is a “black box” which takes a vector \mathbf{x} and returns $\mathbf{H}\mathbf{x}$. This is especially useful if it takes $\ll O(N^2)$ operations to apply \mathbf{H} .

In the MRI example above, it takes about one second to apply $\mathbf{H} = \mathbf{A}^T \mathbf{A}$, and the conjugate gradients method converges in about 50 iterations, meaning that the problem can be solved in less than a minute. Also, the storage requirement is on the order of $O(M + N)$, rather than $O(MN)$.

Recasting sym+def linear equations as an optimization program

We want to solve systems like

$$\underbrace{\mathbf{A}^T \mathbf{A}}_{\mathbf{H}} \mathbf{x} = \underbrace{\mathbf{A}^T \mathbf{y}}_{\mathbf{b}} \quad \text{or} \quad \underbrace{(\mathbf{K} + \delta \mathbf{I})}_{\mathbf{H}} \alpha = \underbrace{\mathbf{y}}_{\mathbf{b}}.$$

Note that in both cases above, \mathbf{H} is symmetric and positive definite (assuming in the former case that \mathbf{A} has full column rank). This means that the following quadratic program is **convex**:

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} - \mathbf{x}^T \mathbf{b}. \quad (1)$$

The consequence of this is that a necessary and sufficient condition for $\hat{\mathbf{x}}$ to be the minimizer of (1) is the gradient (evaluated at $\hat{\mathbf{x}}$) to be zero:

$$\nabla_{\mathbf{x}} \left(\frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} - \mathbf{x}^T \mathbf{b} \right) \Big|_{\mathbf{x}=\hat{\mathbf{x}}} = \mathbf{0}.$$

Since

$$\begin{aligned} \nabla_{\mathbf{x}} \left(\frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} - \mathbf{x}^T \mathbf{b} \right) &= \frac{1}{2} \nabla_{\mathbf{x}} (\mathbf{x}^T \mathbf{H} \mathbf{x}) - \nabla_{\mathbf{x}} (\mathbf{x}^T \mathbf{b}) \\ &= \mathbf{H} \mathbf{x} - \mathbf{b}, \end{aligned}$$

this means that $\hat{\mathbf{x}}$ is the solution to (1) if and only if

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

Steepest descent

Say you have an unconstrained optimization program

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x})$$

where $f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}$ is convex. One simple way to solve this program is to simply “roll downhill”. If we are sitting at a point

\mathbf{x}_0 , then $f(\cdot)$ decreases the fastest if we move in the direction of the negative gradient $-\nabla f(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_0}$.

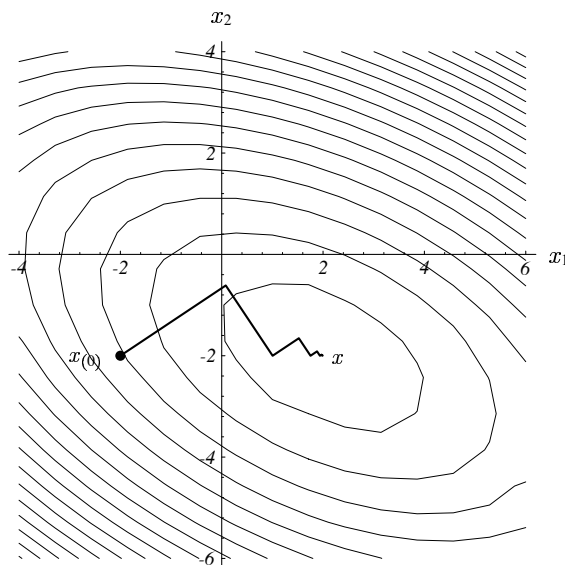
From a starting point \mathbf{x}_0 , we move to

$$\mathbf{x}_1 = \mathbf{x}_0 - \alpha_0 \nabla f(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_0}$$

then to

$$\begin{aligned} \mathbf{x}_2 &= \mathbf{x}_1 - \alpha_1 \nabla f(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_1} \\ &\vdots \\ \mathbf{x}_k &= \mathbf{x}_{k-1} - \alpha_{k-1} \nabla f(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_{k-1}}, \end{aligned}$$

where the $\alpha_0, \alpha_1, \dots$ are appropriately chosen **step sizes**.



(from Shewchuk, “... without the agonizing pain”)

For our particular optimization problem

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} - \mathbf{x}^T \mathbf{b},$$

we can explicitly compute both the gradient and the best choice of step size. The (negative) gradient is what we call the **residual**, the difference between \mathbf{b} and \mathbf{H} applied to the current iterate:

$$-\nabla \left(\frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} - \mathbf{x}^T \mathbf{b} \right) \Big|_{\mathbf{x}=\mathbf{x}_k} = \mathbf{b} - \mathbf{H} \mathbf{x}_k =: \mathbf{r}_k.$$

The steepest descent iteration can be written as

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{r}_k.$$

There is a nifty way to choose an optimal value for the step size α_k . We want to choose α_k so that $f(\mathbf{x}_{k+1})$ is as small as possible. It is not hard to show that along the ray $\mathbf{x}_k + \alpha \mathbf{r}_k$ for $\alpha \geq 0$,

$$f(\mathbf{x}_k + \alpha \mathbf{r}_k) \text{ is convex as a function of } \alpha.$$

Thus we can choose the value of α that makes the derivative of this function zero; we want

$$\frac{d}{d\alpha} f(\mathbf{x}_k + \alpha \mathbf{r}_k) = 0.$$

By the chain rule,

$$\begin{aligned} \frac{d}{d\alpha} f(\mathbf{x}_{k+1}) &= \nabla f(\mathbf{x}_{k+1})^T \frac{d}{d\alpha} \mathbf{x}_{k+1} \\ &= \nabla f(\mathbf{x}_{k+1})^T \mathbf{r}_k. \end{aligned}$$

So we need to choose α_k such that

$$\nabla f(\mathbf{x}_{k+1}) \perp \mathbf{r}_k,$$

or more concisely

$$\mathbf{r}_{k+1} \perp \mathbf{r}_k \quad (\mathbf{r}_{k+1}^T \mathbf{r}_k = 0).$$

So let's do this

$$\begin{aligned} & \mathbf{r}_{k+1}^T \mathbf{r}_k = 0 \\ \Rightarrow & (\mathbf{b} - \mathbf{H}\mathbf{x}_{k+1})^T \mathbf{r}_k = 0 \\ \Rightarrow & (\mathbf{b} - \mathbf{H}(\mathbf{x}_k + \alpha_k \mathbf{r}_k))^T \mathbf{r}_k = 0 \\ \Rightarrow & (\mathbf{b} - \mathbf{H}\mathbf{x}_k)^T \mathbf{r}_k - \alpha_k \mathbf{r}_k^T \mathbf{H}\mathbf{r}_k = 0 \\ \Rightarrow & \mathbf{r}_k^T \mathbf{r}_k - \alpha_k \mathbf{r}_k^T \mathbf{H}\mathbf{r}_k = 0 \end{aligned}$$

and so the optimal step size is

$$\alpha_k = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T \mathbf{H}\mathbf{r}_k}.$$

The steepest descent algorithm performs this iteration until $\|\mathbf{H}\mathbf{x}_k - \mathbf{b}\|_2$ is below some tolerance δ :

Steepest Descent, version 1

Initialize: \mathbf{x}_0 = some guess, $k = 0$, $\mathbf{r}_0 = \mathbf{b} - \mathbf{H}\mathbf{x}_0$.

while not converged, $\|\mathbf{r}_k\|_2 \geq \delta$ **do**

$$\alpha_k = \mathbf{r}_k^T \mathbf{r}_k / \mathbf{r}_k^T \mathbf{H}\mathbf{r}_k$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{r}_k$$

$$\mathbf{r}_{k+1} = \mathbf{b} - \mathbf{H}\mathbf{x}_{k+1}$$

$$k = k + 1$$

end while

There is a nice trick that can save us one of two applications of \mathbf{H} needed in each iteration above. Notice that

$$\begin{aligned}\mathbf{r}_{k+1} &= \mathbf{b} - \mathbf{H}\mathbf{x}_{k+1} = \mathbf{b} - \mathbf{H}(\mathbf{x}_k + \alpha_k \mathbf{r}_k) \\ &= \mathbf{r}_k - \alpha_k \mathbf{H}\mathbf{r}_k.\end{aligned}$$

So we can save an application of \mathbf{H} by updating the residual rather than recomputing it at each stage.

Steepest Descent, more efficient version 2

Initialize: $\mathbf{x}_0 = \text{some guess}$, $k = 0$, $\mathbf{r}_0 = \mathbf{b} - \mathbf{H}\mathbf{x}_0$.

while not converged, $\|\mathbf{r}_k\|_2 \geq \delta$ **do**

$$\mathbf{q} = \mathbf{H}\mathbf{r}_k$$

$$\alpha_k = \mathbf{r}_k^T \mathbf{r}_k / \mathbf{r}_k^T \mathbf{q}$$

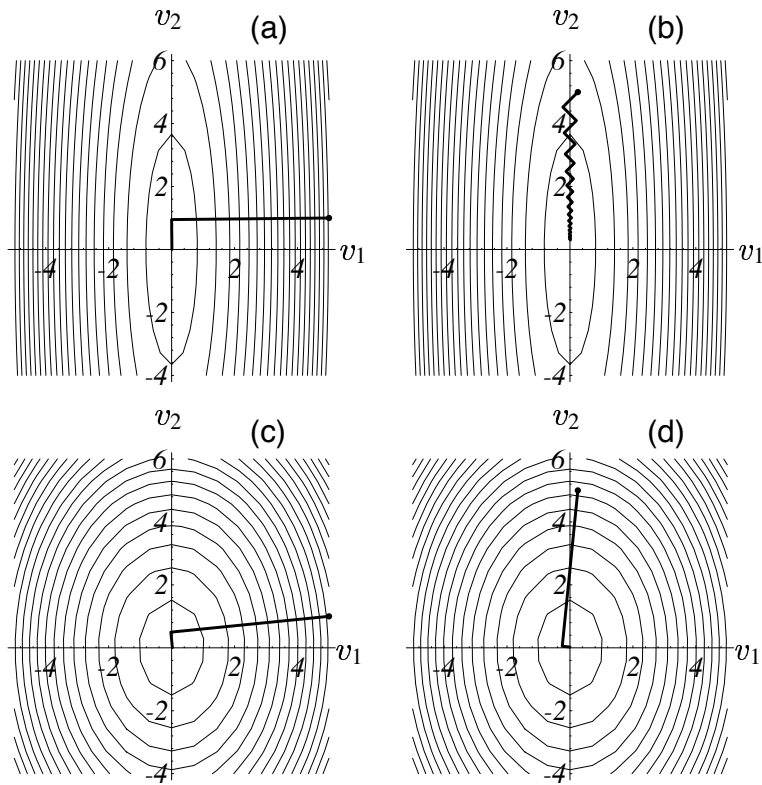
$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{r}_k$$

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{q}$$

$$k = k + 1$$

end while

The effectiveness of SD depends critically on how \mathbf{H} is conditioned and the starting point. Consider the two examples on the next page.



(from Shewchuk, “... without the agonizing pain”)

When the conditioning of \mathbf{H} is poor and we choose a bad starting point, convergence can take many iterations even in simple cases.

The method of conjugate gradients (CG)

An **excellent** resource for the material in this section is the manuscript: 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 can be inefficient because it can move in essentially the same direction many times.

CG avoids this by ensuring that each step is orthogonal (in an appropriate inner product) to all of the previous steps that have been taken. Miraculously, this can be done with very little overhead.

Suppose for a moment that we pre-determine N step directions $\mathbf{d}_0, \dots, \mathbf{d}_{N-1}$ that are orthogonal (but not necessarily normalized), $\mathbf{d}_j^T \mathbf{d}_i = 0$ for $i \neq j$. This means that $\{\mathbf{d}_k / \|\mathbf{d}_k\|_2, k = 0, \dots, N-1\}$ is an orthobasis for \mathbb{R}^N . Then given a starting point \mathbf{x}_0 , the initial error $\mathbf{e}_0 = \mathbf{x}_0 - \hat{\mathbf{x}}$, where $\hat{\mathbf{x}}$ is the solution that satisfies $\mathbf{H}\hat{\mathbf{x}} = \mathbf{b}$, can be expanded as

$$\mathbf{e}_0 = \sum_{\ell=0}^{N-1} c_\ell \frac{\mathbf{d}_\ell}{\|\mathbf{d}_\ell\|_2}, \quad \text{where} \quad c_\ell = \frac{\mathbf{d}_\ell^T \mathbf{e}_0}{\|\mathbf{d}_\ell\|_2^2}. \quad (2)$$

Given step sizes $\alpha_0, \alpha_1, \dots$, the error after the k th step is

$$\begin{aligned}
\mathbf{e}_k &= \mathbf{x}_k - \hat{\mathbf{x}} \\
&= \mathbf{x}_{k-1} + \alpha_{k-1} \mathbf{d}_{k-1} - \hat{\mathbf{x}} = \mathbf{e}_{k-1} + \alpha_{k-1} \mathbf{d}_{k-1} \\
&= \mathbf{e}_{k-2} + \alpha_{k-1} \mathbf{d}_{k-1} + \alpha_{k-2} \mathbf{d}_{k-2} \\
&\vdots \\
&= \mathbf{e}_0 + \sum_{\ell=0}^{k-1} \alpha_{\ell} \mathbf{d}_{\ell}.
\end{aligned}$$

Thus, if we choose the step sizes α_k carefully, then we can pick-off a component in (2) at every step. In particular, if we choose

$$\alpha_k = -\frac{c_k}{\|\mathbf{d}_k\|_2} = \frac{-\mathbf{d}_k^T \mathbf{e}_0}{\|\mathbf{d}_k\|_2^2}, \quad (3)$$

then we have

$$\mathbf{e}_k = \sum_{\ell=k}^{N-1} c_{\ell} \frac{\mathbf{d}_{\ell}}{\|\mathbf{d}_{\ell}\|_2}, \quad \text{and} \quad \|\mathbf{e}_k\|_2^2 = \sum_{\ell=k}^{N-1} |c_{\ell}|^2.$$

So we see that as k increases, there are fewer and fewer terms in the sum above, steadily decreasing the error until

$$\mathbf{e}_N = \mathbf{0}.$$

The argument above works for any set of orthogonal step directions $\{\mathbf{d}_k\}$. It would be beautiful, except that we **do not know** the initial error $\mathbf{e}_0 = \mathbf{x}_0 - \hat{\mathbf{x}}$. (If we did, we would have a solution in one step: just subtract \mathbf{e}_0 from \mathbf{x}_0 !) Thus there is now way we can compute the stepsizes in (3).

But the argument above works not only for any orthobasis, but also for any orthobasis using any (valid) inner product. The key innovation in CG is to adaptively choose the step directions \mathbf{d}_k and step sizes α_k so that the steps are orthogonal in the \mathbf{H} inner product:

$$\langle \mathbf{d}_i, \mathbf{d}_j \rangle_{\mathbf{H}} = \mathbf{d}_j^{\mathrm{T}} \mathbf{H} \mathbf{d}_i.$$

It is easy to verify that if \mathbf{H} is sym+def, then $\langle \cdot, \cdot \rangle_{\mathbf{H}}$ is a valid inner product.

So again, suppose that we start at \mathbf{x}_0 with initial error \mathbf{e}_0 . If $\mathbf{d}_0, \dots, \mathbf{d}_{N-1}$ are \mathbf{H} -orthogonal vectors, then

$$\mathbf{e}_0 = \sum_{\ell=0}^{N-1} c_{\ell} \frac{\mathbf{d}_{\ell}}{\|\mathbf{d}_{\ell}\|_{\mathbf{H}}},$$

where $\|\mathbf{d}_{\ell}\|_{\mathbf{H}}^2 = \mathbf{d}_{\ell}^{\mathrm{T}} \mathbf{H} \mathbf{d}_{\ell}$, and

$$c_{\ell} = \frac{\langle \mathbf{e}_0, \mathbf{d}_{\ell} \rangle_{\mathbf{H}}}{\|\mathbf{d}_{\ell}\|_{\mathbf{H}}} = \frac{\mathbf{d}_{\ell}^{\mathrm{T}} \mathbf{H} \mathbf{e}_0}{\|\mathbf{d}_{\ell}\|_{\mathbf{H}}}.$$

As we will show below, the iterations below produce a set of \mathbf{H} -orthogonal step directions $\{\mathbf{d}_k\}$ with step sizes $\alpha_k = -c_k / \|\mathbf{d}_k\|_{\mathbf{H}}$.

Conjugate Gradients

Initialize: $\mathbf{x}_0 = \text{some guess}$

$$\mathbf{r}_0 = \mathbf{b} - \mathbf{H}\mathbf{x}_0$$

$$\mathbf{d}_0 = \mathbf{r}_0$$

for $k = 0$ to $N - 1$ **do**

$$\alpha_k = \mathbf{r}_k^T \mathbf{r}_k / \mathbf{d}_k^T \mathbf{H} \mathbf{d}_k$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$$

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{H} \mathbf{d}_k$$

$$\beta_{k+1} = \mathbf{r}_{k+1}^T \mathbf{r}_{k+1} / \mathbf{r}_k^T \mathbf{r}_k$$

$$\mathbf{d}_{k+1} = \mathbf{r}_{k+1} + \beta_{k+1} \mathbf{d}_k$$

end for

We will start our analysis of this iteration by establishing the following two facts:

Fact 1: $\langle \mathbf{r}_{k+1}, \mathbf{r}_\ell \rangle = \mathbf{r}_\ell^T \mathbf{r}_{k+1} = 0$ for $\ell = 0, \dots, k$.

That is, the **residual is orthogonal** to all previous residuals.

Fact 2: $\langle \mathbf{d}_{k+1}, \mathbf{d}_\ell \rangle_H = \mathbf{d}_\ell^T \mathbf{H} \mathbf{d}_{k+1} = 0$ for $\ell = 0, \dots, k$.

That is, the **direction is \mathbf{H} -orthogonal** to all previous directions.

We establish these two facts by induction. We start at $k = 1$:

1. $\langle \mathbf{r}_1, \mathbf{r}_0 \rangle = \mathbf{r}_0^T \mathbf{r}_1 = 0$, since

$$\begin{aligned}\mathbf{r}_1 &= \mathbf{r}_0 - \frac{\mathbf{r}_0^T \mathbf{r}_0}{\mathbf{r}_0^T \mathbf{H} \mathbf{r}_0} \mathbf{H} \mathbf{r}_0 \\ \Rightarrow \mathbf{r}_0^T \mathbf{r}_1 &= \mathbf{r}_0^T \mathbf{r}_0 - \mathbf{r}_0^T \mathbf{r}_0 \frac{\mathbf{r}_0^T \mathbf{H} \mathbf{r}_0}{\mathbf{r}_0^T \mathbf{H} \mathbf{r}_0} \\ &= 0.\end{aligned}$$

2. $\langle \mathbf{d}_1, \mathbf{d}_0 \rangle_H = \mathbf{d}_0^T \mathbf{H} \mathbf{d}_1 = 0$, since

$$\begin{aligned}\mathbf{r}_1 &= \mathbf{r}_0 - \alpha_0 \mathbf{H} \mathbf{d}_0 \\ \Rightarrow \mathbf{r}_1^T \mathbf{r}_1 &= \mathbf{r}_1^T \mathbf{r}_0 - \alpha_0 \mathbf{r}_1^T \mathbf{H} \mathbf{d}_0 \\ \Rightarrow \mathbf{r}_1^T \mathbf{H} \mathbf{d}_0 &= -\frac{1}{\alpha_0} \mathbf{r}_1^T \mathbf{r}_1,\end{aligned}$$

since $\mathbf{r}_1^T \mathbf{r}_0 = 0$. Also,

$$\begin{aligned}\mathbf{d}_1 &= \mathbf{r}_1 + \frac{\mathbf{r}_1^T \mathbf{r}_1}{\mathbf{r}_0^T \mathbf{r}_0} \mathbf{d}_0 \\ \Rightarrow \mathbf{d}_0^T \mathbf{H} \mathbf{d}_1 &= \mathbf{d}_0^T \mathbf{H} \mathbf{r}_1 + \frac{\mathbf{r}_1^T \mathbf{r}_1}{\mathbf{r}_0^T \mathbf{r}_0} \mathbf{d}_0^T \mathbf{H} \mathbf{d}_0 \\ &= \frac{\mathbf{r}_1^T \mathbf{r}_1}{\mathbf{r}_0^T \mathbf{r}_0} \mathbf{d}_0^T \mathbf{H} \mathbf{d}_0 - \frac{\mathbf{r}_1^T \mathbf{r}_1}{\mathbf{r}_0^T \mathbf{r}_0} \mathbf{d}_0^T \mathbf{H} \mathbf{d}_0 \\ &= 0.\end{aligned}$$

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

$$\begin{aligned}\langle \mathbf{r}_j, \mathbf{r}_\ell \rangle &= \mathbf{r}_\ell^T \mathbf{r}_j = 0 \quad \forall j, \ell \leq k, \\ \langle \mathbf{d}_j, \mathbf{d}_\ell \rangle_H &= \mathbf{d}_\ell^T \mathbf{H} \mathbf{d}_j = 0 \quad \forall j, \ell \leq k.\end{aligned}$$

Then we will also have the following:

1. $\langle \mathbf{r}_{k+1}, \mathbf{r}_\ell \rangle = \mathbf{r}_\ell^\top \mathbf{r}_{k+1} = 0$ for all $\ell \leq k$.

To see this, notice that

$$\mathbf{r}_\ell^\top \mathbf{H} \mathbf{d}_k = (\mathbf{d}_\ell - \beta_\ell \mathbf{d}_{\ell-1})^\top \mathbf{H} \mathbf{d}_k \quad (4)$$

$$= \begin{cases} \mathbf{d}_k^\top \mathbf{H} \mathbf{d}_k & \ell = k \\ 0 & \ell < k, \end{cases} \quad (5)$$

where the second step follows directly from the fact that $\langle \mathbf{d}_k, \mathbf{d}_\ell \rangle_{\mathbf{H}} = 0$ for $\ell < k$. As a result

$$\mathbf{r}_\ell^\top \mathbf{r}_{k+1} = \mathbf{r}_\ell^\top \mathbf{r}_k - \frac{\mathbf{r}_k^\top \mathbf{r}_k}{\mathbf{d}_k^\top \mathbf{H} \mathbf{d}_k} \mathbf{r}_\ell^\top \mathbf{H} \mathbf{d}_k = 0 \quad \text{for all } \ell \leq k. \quad (6)$$

2. $\langle \mathbf{d}_{k+1}, \mathbf{d}_\ell \rangle_{\mathbf{H}} = \mathbf{d}_\ell^\top \mathbf{H} \mathbf{d}_{k+1} = 0$ for all $\ell \leq k$.

This follows from the expansion

$$\mathbf{d}_\ell^\top \mathbf{H} \mathbf{d}_{k+1} = \mathbf{d}_\ell^\top \mathbf{H} \mathbf{r}_{k+1} + \beta_{k+1} \mathbf{d}_\ell^\top \mathbf{H} \mathbf{d}_k.$$

Notice that

$$\begin{aligned} \mathbf{r}_i^\top \mathbf{r}_{k+1} &= \mathbf{r}_i^\top \mathbf{r}_k - \alpha_k \mathbf{r}_i^\top \mathbf{H} \mathbf{d}_k \\ \Rightarrow \mathbf{r}_i^\top \mathbf{H} \mathbf{d}_k &= \begin{cases} \frac{1}{\alpha_k} \mathbf{r}_k^\top \mathbf{r}_k & i = k \\ -\frac{1}{\alpha_k} \mathbf{r}_{k+1}^\top \mathbf{r}_{k+1} & i = k+1 \\ 0 & i < k. \end{cases} \end{aligned} \quad (7)$$

Then for $\ell = k$

$$\begin{aligned} \mathbf{d}_k^\top \mathbf{H} \mathbf{d}_{k+1} &= -\frac{1}{\alpha_k} \mathbf{r}_{k+1}^\top \mathbf{r}_{k+1} + \beta_{k+1} \mathbf{d}_k^\top \mathbf{H} \mathbf{d}_k \\ &= \frac{-\mathbf{r}_{k+1}^\top \mathbf{r}_{k+1}}{\mathbf{r}_k^\top \mathbf{r}_k} \mathbf{d}_k^\top \mathbf{H} \mathbf{d}_k + \frac{\mathbf{r}_{k+1}^\top \mathbf{r}_{k+1}}{\mathbf{r}_k^\top \mathbf{r}_k} \mathbf{d}_k^\top \mathbf{H} \mathbf{d}_k \\ &= 0. \end{aligned}$$

For $\ell < k$,

$$\mathbf{d}_\ell^\top \mathbf{H} \mathbf{d}_{k+1} = \mathbf{d}_\ell^\top \mathbf{H} \mathbf{r}_{k+1} + \beta_{k+1} \mathbf{d}_\ell^\top \mathbf{H} \mathbf{d}_k.$$

For the first term

$$\mathbf{d}_\ell^\top \mathbf{H} \mathbf{r}_{k+1} = 0,$$

since $\mathbf{H} \mathbf{d}_\ell = \alpha_\ell^{-1}(\mathbf{r}_\ell - \mathbf{r}_{\ell+1})$ and we have (6); for the second term

$$\beta_{k+1} \mathbf{d}_\ell^\top \mathbf{H} \mathbf{d}_k = 0,$$

since the $\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_k$ are \mathbf{H} -orthogonal already.

We have established that the direction \mathbf{d}_k that CG moves on iteration k is \mathbf{H} -orthogonal to all previous directions. Now let's look at the step sizes, where we want to establish that $\alpha_k = -c_k / \|\mathbf{d}_k\|_{\mathbf{H}} = -\mathbf{d}_k^\top \mathbf{H} \mathbf{e}_0 / \|\mathbf{d}_k\|_{\mathbf{H}}^2$. Start by noting (5) above, and recall that

$$\mathbf{r}_k = \mathbf{b} - \mathbf{H} \mathbf{x}_k = \mathbf{H}(\hat{\mathbf{x}} - \mathbf{x}_k) = -\mathbf{H} \mathbf{e}_k.$$

At the first step, we have $\mathbf{d}_0 = \mathbf{r}_0$, and so

$$\alpha_0 = \frac{\mathbf{r}_0^\top \mathbf{r}_0}{\mathbf{d}_0^\top \mathbf{H} \mathbf{d}_0} = \frac{\mathbf{d}_0^\top \mathbf{r}_0}{\mathbf{d}_0^\top \mathbf{H} \mathbf{d}_0} = \frac{\mathbf{d}_0^\top \mathbf{H}(\hat{\mathbf{x}} - \mathbf{x}_0)}{\mathbf{d}_0^\top \mathbf{H} \mathbf{d}_0} = \frac{-\mathbf{d}_0^\top \mathbf{H} \mathbf{e}_0}{\mathbf{d}_0^\top \mathbf{H} \mathbf{d}_0}.$$

At subsequent steps, since

$$\mathbf{d}_k = \mathbf{r}_k + \sum_{i=0}^{k-1} \gamma_i \mathbf{r}_i \quad \text{for some } \gamma_i \in \mathbb{R},$$

by Fact 1, we have

$$\mathbf{r}_k^\top \mathbf{r}_k = \mathbf{d}_k^\top \mathbf{r}_k,$$

and so

$$\alpha_k = \frac{\mathbf{d}_k^T \mathbf{r}_k}{\mathbf{d}_k^T \mathbf{H} \mathbf{d}_k} = \frac{-\mathbf{d}_k^T \mathbf{H} \left(\mathbf{e}_0 + \sum_{\ell=0}^{k-1} \alpha_\ell \mathbf{d}_\ell \right)}{\mathbf{d}_k^T \mathbf{H} \mathbf{d}_k} = \frac{-\mathbf{d}_k^T \mathbf{H} \mathbf{e}_0}{\mathbf{d}_k^T \mathbf{H} \mathbf{d}_k}.$$

So finally, this means that for the method of conjugate gradients,

$$\mathbf{e}_k = \sum_{\ell=k}^{N-1} \left(\frac{\mathbf{d}_\ell^T \mathbf{r}_\ell}{\mathbf{d}_\ell^T \mathbf{H} \mathbf{d}_\ell} \right) \mathbf{d}_\ell, \quad \|\mathbf{e}_k\|_{\mathbf{H}}^2 = \sum_{\ell=k}^{N-1} \frac{|\mathbf{d}_\ell^T \mathbf{r}_\ell|^2}{\mathbf{d}_\ell^T \mathbf{H} \mathbf{d}_\ell}.$$

As k increases, the number of (positive) terms in the sum above gets smaller and smaller, until finally

$$\mathbf{e}_N = \mathbf{0}.$$

Thus **CG is guaranteed to converge exactly in N steps.**

Since each iteration of CG involves a vector-matrix multiply, each of which are $O(N^2)$, and we converge in $O(N)$ iterations, CG solves $\mathbf{H}\mathbf{x} = \mathbf{b}$ in $O(N^3)$ computations in general, the same as other solvers.

But there are two important things to realize:

1. If \mathbf{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 \mathbf{H} .
2. It is often the case that $\|\mathbf{e}_k\|_{\mathbf{H}}^2$ is acceptably small for relatively modest values of k . This is particularly true if \mathbf{H} is well-conditioned. Each iteration (application of \mathbf{H}) gets us closer, in a measurable way, to the solution.

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.

Convergence Guarantees

We can actually talk intelligently about how many iterations we need for steepest descent and CG to converge to within a certain precision. Here we present (but do not prove) two “worst case” bounds that depend on the condition number κ of \mathbf{H} :

$$\kappa = \frac{\lambda_{\max}(\mathbf{H})}{\lambda_{\min}(\mathbf{H})} = \frac{\text{max eigenvalue}}{\text{min eigenvalue}}.$$

For steepest descent, we will have

$$\|\mathbf{e}_k\|_{\mathbf{H}} \leq \delta \|\mathbf{e}_0\|_{\mathbf{H}}$$

in at most¹

$$k \leq \left\lceil \frac{1}{2} \kappa \log \left(\frac{1}{\delta} \right) \right\rceil$$

iterations.

For CG, we need at most

$$k \leq \left\lceil \frac{1}{2} \sqrt{\kappa} \log \left(\frac{2}{\delta} \right) \right\rceil.$$

There are nice derivations for both of these bounds in the Shewchuk manuscript mentioned at the beginning of these notes.

¹These are natural logarithms.

Example:

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

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

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

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