

Iterative methods for solving least squares

When \mathbf{A} has full column rank, our least squares estimate is

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

If \mathbf{A} is $M \times N$, then constructing $\mathbf{A}^T \mathbf{A}$ costs $O(MN^2)$ computations, and solving the $N \times N$ system $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{y}$, for example, by computing the inverse of $\mathbf{A}^T \mathbf{A}$, costs $O(N^3)$ computations. A similar complexity is involved in computing the SVD of \mathbf{A} , so that will not be any cheaper. (Note that for $M \geq N$, the cost of constructing the matrix actually exceeds the cost to solve the system.)

This cost can be prohibitive for even moderately large M and N . But least squares 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” (3D pixels) from about 5 million measurements. In this case, the matrix \mathbf{A} , which models the mapping from the 3D image \mathbf{x} to the set of measurements \mathbf{y} induced by the MRI machine, is $M \times N$ where $M = 5 \cdot 10^6$ and $N = 2.1 \cdot 10^6$.

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.

To address this we can consider approaches that return to our formulation of least squares as an optimization program and then solve it by an iterative descent method. Each iteration is simple, requiring one application of \mathbf{A} and one application of \mathbf{A}^T .

If $\mathbf{A}^T \mathbf{A}$ is “well-conditioned”, then these methods can converge in very few iterations. We will be more precise about this later, but

here “well-conditioned” roughly means that the ratio of the largest to the smallest singular values of \mathbf{A} is not too big. When this works, it can make the cost of solving a least squares problem dramatically smaller — about the cost of a few hundred applications of \mathbf{A} .

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

In the MRI example above, it turns out that \mathbf{A} has a special relationship to the Fourier transform, and because of this it takes about one second to apply $\mathbf{A}^T \mathbf{A}$, and a particular iterative method (the conjugate gradients method) converges in about 50 iterations, meaning that the problem can be solved in less than a minute.

To see how this approach works, recall that the least squares estimate is the solution to the optimization problem

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2.$$

Note that we can write this equivalently as

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} - 2\mathbf{x}^T \mathbf{A}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}.$$

We can ignore terms that do not depend on \mathbf{x} , and can also rescale the objective function by a constant (for convenience) to obtain

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

We have previously shown that a necessary and sufficient condition for $\hat{\mathbf{x}}$ to be the the minimizer of (1) is to satisfy

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{y}.$$

More generally, for any \mathbf{H} which is symmetric and positive definite and any vector \mathbf{b} , we can consider the optimization problem

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

and by the same argument we can show that $\hat{\mathbf{x}}$ is the solution to (2) if and only if

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

What remains is to show how we can actually solve an optimization problem of the form (2) without directly solving the system $\mathbf{H} \mathbf{x} = \mathbf{b}$. Here we will describe iterative methods — most prominently **gradient descent** — that do exactly this.

Gradient 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. We will give a more formal definition later, but for now let's just go with the very informal notion that convexity corresponds to a “bowl shape”. One simple way to solve this program is to simply “roll downhill”. If we are sitting at a point \mathbf{x}_0 , then as we mentioned previously in our review of multivariable calculus, f decreases the fastest if we move in the direction of the *negative gradient* $-\nabla f(\mathbf{x}_0)$, where we recall that this notation means the gradient of f with respect to \mathbf{x} evaluated at \mathbf{x}_0 .

Thus, suppose that from a starting point \mathbf{x}_0 , we take a step in the direction of the negative gradient, where the step size is controlled

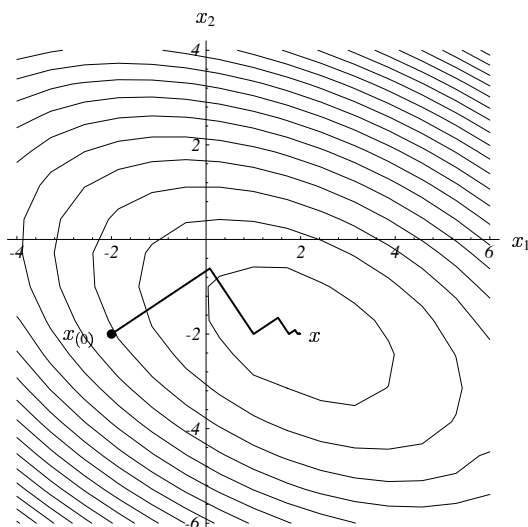
by a parameter α_0 :

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

We then repeat this process:

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

where, as before, 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 straightforward to compute (you've already done this on the homework):

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

Note that this is simply the difference between \mathbf{b} and \mathbf{H} applied to the current iterate \mathbf{x}_k . For this reason it is often called the **residual**, denoted by

$$\mathbf{r}_k := -\nabla f(\mathbf{x}_k) = \mathbf{b} - \mathbf{H}\mathbf{x}_k.$$

With this notation, the core gradient descent iteration can be written as

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

As mentioned above, in this problem 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 if we think of $f(\mathbf{x}_k + \alpha\mathbf{r}_k)$ as a function of α for $\alpha \geq 0$, then f is a (convex) quadratic function. Thus to find the α that minimizes $f(\mathbf{x}_{k+1})$, we can choose the value of α that makes the derivative of this function zero. Specifically, we want

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

By the chain rule,

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

In summary, we need to choose α_k such that

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

Let's do exactly 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 gradient descent algorithm performs this iteration until $\|\mathbf{H}\mathbf{x}_k - \mathbf{b}\|_2$ is below some tolerance ϵ :

Gradient Descent, version 1

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

while $\|\mathbf{r}_k\|_2 \geq \epsilon$ (not converged) **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

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

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

Gradient Descent, more efficient version 2

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

while $\|\mathbf{r}_k\|_2 \geq \epsilon$ (not converged) **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

One small caveat with this approach: you will note that we are updating \mathbf{r}_k in a recursive fashion. This has the consequence that over time small numerical errors can accumulate in \mathbf{r}_k . This can result in a situation where the algorithm seemingly fails to converge (because $\|\mathbf{r}_k\|_2$ remains large), despite actually making good progress. The typical way to handle this is to “reset” \mathbf{r}_k by computing $\mathbf{r}_k = \mathbf{b} - \mathbf{H}\mathbf{x}_k$ every so often (e.g., once every 50 iterations). Exactly how often you will want to do this depends on the relative costs associated with potentially running more unnecessary iterations versus that of the extra matrix multiply that this “resetting” step requires.

Example: Swarm robotics

Suppose that we have N robots with positions $\mathbf{p}^{(1)}, \mathbf{p}^{(2)}, \dots, \mathbf{p}^{(N)}$, where each $\mathbf{p}^{(n)}$ is a vector in \mathbb{R}^D , with $D = 2$ or 3 , depending on the application. Suppose that we want these robots to meet at the same location. We do not care where this is, we simply want the robots to all converge to the same point. We can pose this as the solution to a convex optimization problem. Specifically, set

$$\mathbf{x} = \begin{bmatrix} \mathbf{p}^{(1)} \\ \mathbf{p}^{(2)} \\ \vdots \\ \mathbf{p}^{(N)} \end{bmatrix},$$

so that $\mathbf{x} \in \mathbb{R}^{ND}$. Next, for each robot we define a neighborhood or a set of indices \mathcal{N}_n corresponding to the robots to which robot n can measure its distance. In other words, if $m \in \mathcal{N}_n$, robot n can compute $\|\mathbf{p}^{(n)} - \mathbf{p}^{(m)}\|_2$. We will assume for the sake of simplicity that these are symmetric in the sense that $m \in \mathcal{N}_n$ if and only if $n \in \mathcal{N}_m$. We would like all of these distances to be zero, so a natural objective function that we might want to minimize is

$$f(\mathbf{x}) = \sum_{n=1}^N \sum_{m \in \mathcal{N}_n} \|\mathbf{p}^{(n)} - \mathbf{p}^{(m)}\|_2^2.$$

We can compute the gradient of this function by noting that

$$\nabla_{\mathbf{p}^{(n)}} f(\mathbf{x}) = \sum_{m \in \mathcal{N}_n} 2(\mathbf{p}^{(n)} - \mathbf{p}^{(m)}) + \sum_{m: n \in \mathcal{N}_m} 2(\mathbf{p}^{(n)} - \mathbf{p}^{(m)}).$$

If we make the simplifying assumption that the neighborhoods are symmetric, so that $m \in \mathcal{N}_n$ if and only if $n \in \mathcal{N}_m$, then this simplifies

to

$$\nabla_{\mathbf{p}^{(n)}} f(\mathbf{x}) = 4 \sum_{m \in \mathcal{N}_n} (\mathbf{p}^{(n)} - \mathbf{p}^{(m)}).$$

Putting this all together, we can write

$$\nabla f(\mathbf{x}) = 4 \begin{bmatrix} \sum_{m \in \mathcal{N}_1} (\mathbf{p}^{(1)} - \mathbf{p}^{(m)}) \\ \sum_{m \in \mathcal{N}_2} (\mathbf{p}^{(2)} - \mathbf{p}^{(m)}) \\ \vdots \\ \sum_{m \in \mathcal{N}_N} (\mathbf{p}^{(N)} - \mathbf{p}^{(m)}) \end{bmatrix}.$$

In this case the update rule $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$ nicely de-couples so that the n^{th} robot has the update rule (ignoring the multiplicative factor of 4):

$$\mathbf{p}_{k+1}^{(n)} = \mathbf{p}_k^{(n)} - \alpha_k \sum_{m \in \mathcal{N}_n} (\mathbf{p}_k^{(n)} - \mathbf{p}_k^{(m)}).$$

This update rule plays a fundamental role in many swarm robotics problems and is known as the **consensus equation**. Note that the update for each robot depends only on *local* information (the difference between its own position and that of its neighbors), and hence each robot can compute its own update without any form of global coordination.

For a sufficiently small step size (as we will see next time), this algorithm is guaranteed to converge. Moreover, provided that the neighborhoods are fully connected (so that there is at least some indirect path between any pair of robots) then the global optimum of this problem will be for all robots to converge to the same point.