

1 THE CONDITION NUMBER AS A FUNDAMENTAL MATRIX QUANTITY

Here, we show that the condition number of a matrix determines how quickly various simple iterative methods will converge on symmetric positive $1 - \frac{1}{\kappa(\mathbf{A})}$ definite linear systems. Thus, throughout these notes, we will assume that \mathbf{A} is symmetric positive definite.

1.1 RICHARDSON

Recall the Richardson iteration for $\mathbf{Ax} = \mathbf{b}$:

$$\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{Ax}^{(k)} \quad \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \omega \mathbf{r}^{(k)}.$$

We can write this in terms of the gradient for the quadratic problem:

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{Ax} - \mathbf{x}^T \mathbf{b} \quad \text{with gradient} \quad \mathbf{g}(\mathbf{x}) = \mathbf{Ax} - \mathbf{b} = -\mathbf{r}(\mathbf{x})$$

which gives

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \omega \mathbf{g}^{(k)}.$$

Now consider the error vector

$$\mathbf{e}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}.$$

The evolution of the error is determined by

$$\mathbf{e}^{(k+1)} = \mathbf{x}^{(k+1)} - \mathbf{x} = \mathbf{x}^{(k)} + \omega \mathbf{r}^{(k)} - \mathbf{x} = \mathbf{x}^{(k)} + \omega \mathbf{b} - \omega \mathbf{Ax}^{(k)} - \mathbf{x}.$$

But note that $\omega \mathbf{b} = \omega \mathbf{Ax}$ for the true solution \mathbf{x} . Hence

$$\mathbf{e}^{(k+1)} = \mathbf{x}^{(k)} + \omega \mathbf{Ax} - \omega \mathbf{Ax}^{(k)} - \mathbf{x} = (\mathbf{I} - \omega \mathbf{A})(\mathbf{x}^{(k)} - \mathbf{x})$$

or simply

$$\mathbf{e}^{(k+1)} = (\mathbf{I} - \omega \mathbf{A})\mathbf{e}^{(k)} = (\mathbf{I} - \omega \mathbf{A})^k \mathbf{e}^{(0)}.$$

This converges quickly if we can make the spectral radius $\rho(\mathbf{I} - \omega \mathbf{A})$ small. For a symmetric positive definite matrix, there is an easy way to do this. The derivation is not particularly interesting. The choice is:

$$\omega = \frac{2}{\lambda_1 + \lambda_n}$$

where λ_1 and λ_n are the smallest and largest eigenvalues of \mathbf{A} respectively.¹ For this choice we have

$$\rho(\mathbf{I} - \omega \mathbf{A}) = \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}}.$$

There is no condition number yet, but it's hiding inside this formula! For a symmetric positive definite system, we have $\kappa(\mathbf{A}) = \frac{\lambda_{\max}}{\lambda_{\min}}$, and so we can adjust this expression to include this ratio:

$$\rho(\mathbf{I} - \omega \mathbf{A}) = \frac{\frac{1}{\lambda_{\min}} \lambda_{\max} - \lambda_{\min}}{\frac{1}{\lambda_{\min}} \lambda_{\max} + \lambda_{\min}} = \frac{\kappa(\mathbf{A}) - 1}{\kappa(\mathbf{A}) + 1} \leq \frac{\kappa(\mathbf{A}) - 1}{\kappa(\mathbf{A})} \leq 1 - \frac{1}{\kappa(\mathbf{A})}.$$

So the asymptotic error in Richardson on a symmetric positive definite system goes to 0 at a rate $1 - \frac{1}{\kappa(\mathbf{A})}$. Formally,

$$\|\mathbf{e}^{(k)}\| \leq \left(1 - \frac{1}{\kappa}\right)^k \|\mathbf{e}^{(0)}\|.$$

Learning objectives

1. See how the matrix condition number arises in terms of converge of Richardson with the optimal parameter.

2. See how the matrix condition number arises in terms of converge of Steepest Descent.

3. Compare the rates of convergence.

Up to this point, for Richardson, this is entirely general and actually has not used the assumption that \mathbf{A} is symmetric positive definite.

¹ The way to determine this quantity is to look at how $\mathbf{I} - \omega \mathbf{A}$ changes the eigenvalues of \mathbf{A} . This transform maps the region $[\lambda_1, \lambda_n]$ to the region $[1 - \omega \lambda_n, 1 - \omega \lambda_1]$. We now want to pick ω to minimize $\max(|1 - \omega \lambda_n|, |1 - \omega \lambda_1|)$. Note that when ω is small enough, then $1 - \omega \lambda_n$ and $1 - \omega \lambda_1$ are both positive and $1 - \omega \lambda_1$ determines the spectral radius, which decreases with ω . As ω increases, $1 - \omega \lambda_n$ goes negative first (assuming $\lambda_1 \neq \lambda_n$) and so at some point we have $|1 - \omega \lambda_n| = -1 + \omega \lambda_n = |1 - \omega \lambda_1| = 1 - \omega \lambda_1$, which gives $\omega = 2/(\lambda_1 + \lambda_n)$ as required. This equivalency point is minimizer as further increasing ω just results in a larger spectral radius.

1.2 STEEPEST DESCENT

We will now show that the steepest descent iteration converges for a symmetric positive definite system that also depends on $\kappa(\mathbf{A})$.² Recall that steepest descent uses a dynamic choice of ω , called α or γ , that minimizes the function

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{x}^T \mathbf{b}$$

at each step. The iteration is

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \gamma_k \mathbf{g}(\mathbf{x}_k) \quad \gamma_k = \frac{\mathbf{g}(\mathbf{x}_k)^T \mathbf{g}(\mathbf{x}_k)}{\mathbf{g}(\mathbf{x}_k)^T \mathbf{A} \mathbf{g}(\mathbf{x}_k)}.$$

We are going to tweak this setup slightly. Note that at a solution $\mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$ we have

$$f(\mathbf{A}^{-1} \mathbf{b}) = \frac{1}{2} \mathbf{b}^T \mathbf{A}^{-T} \mathbf{A} \mathbf{A}^{-1} \mathbf{b} - \mathbf{b}^T \mathbf{A}^{-T} \mathbf{b} = -\frac{1}{2} \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b}.$$

The strategy we are going to use is to study the rate $f(\mathbf{x}_k) \rightarrow -\frac{1}{2} \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b}$. But this is a slightly annoying constant to have around, so we just study the function

$$s(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{x}^T \mathbf{b} + \frac{1}{2} \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b}$$

instead. This function is just shifted by a constant, and so the gradient is unchanged. Now, we can study the rate at which $s(\mathbf{x}_k) \rightarrow 0$ instead, which makes life slightly easier.

This shifted function s is also nice for another reason. Let $n(\mathbf{x}) = \sqrt{\mathbf{x}^T \mathbf{A}^{-1} \mathbf{x}}$. Then we can show that $n(\mathbf{x})$ is a vector norm.³ Typically we write this as

$$\|\mathbf{x}\|_{\mathbf{A}^{-1}} = \sqrt{\mathbf{x}^T \mathbf{A}^{-1} \mathbf{x}}.$$

Using this norm, we can write

$$s(\mathbf{x}) = \frac{1}{2} \|\mathbf{A} \mathbf{x} - \mathbf{b}\|_{\mathbf{A}^{-1}}^2 = \frac{1}{2} \|\mathbf{g}(\mathbf{x})\|_{\mathbf{A}^{-1}}^2.$$

The goal is to show that $s(\mathbf{x}^{(k+1)}) \leq s(\mathbf{x}^{(k)})$ (constant less than 1). We have the following that allow us to do so

$$\mathbf{g}(\mathbf{x}_k) = \mathbf{g}_k \quad \gamma_k = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{A} \mathbf{g}_k}$$

$$\mathbf{x}^{(k+1)} = (\mathbf{I} - \gamma_k \mathbf{A}) \mathbf{g}_k$$

$$s(\mathbf{x}^{(k+1)}) = \frac{1}{2} \|(\mathbf{I} - \gamma_k \mathbf{A}) \mathbf{g}_k\|_{\mathbf{A}^{-1}}^2 = \frac{1}{2} \mathbf{g}_k^T \mathbf{A}^{-1} \mathbf{g}_k - \gamma_k \mathbf{g}_k^T \mathbf{g}_k + \frac{1}{2} \gamma_k^2 \mathbf{g}_k^T \mathbf{A} \mathbf{g}_k = s(\mathbf{x}_k) - \frac{1}{2} \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{A} \mathbf{g}_k}$$

$$s(\mathbf{x}^{(k+1)}) = s(\mathbf{x}_k) \left(1 - \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{A} \mathbf{g}_k s(\mathbf{x}_k)}\right) = s(\mathbf{x}_k) \left(1 - \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{A} \mathbf{g}_k \mathbf{g}_k^T \mathbf{A}^{-1} \mathbf{g}_k}\right)$$

because $s(\mathbf{x}_k) = \frac{1}{2} \mathbf{g}_k^T \mathbf{A}^{-1} \mathbf{g}_k$.

The key is that this quantity $\frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{A} \mathbf{g}_k \mathbf{g}_k^T \mathbf{A}^{-1} \mathbf{g}_k}$ is fairly close to a condition number. Let θ be the inverse quantity so

$$\theta = \frac{\mathbf{g}_k^T \mathbf{A} \mathbf{g}_k \mathbf{g}_k^T \mathbf{A}^{-1} \mathbf{g}_k}{(\mathbf{g}_k^T \mathbf{g}_k)^2}.$$

Then we have

$$\theta = \frac{\mathbf{g}_k^T \mathbf{A} \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{g}_k} \frac{\mathbf{g}_k^T \mathbf{A}^{-1} \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{g}_k} \leq \max_{\mathbf{g}} \left[\frac{\mathbf{g}_k^T \mathbf{A} \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{g}_k} \right] \max_{\mathbf{g}} \left[\frac{\mathbf{g}_k^T \mathbf{A}^{-1} \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{g}_k} \right] = \lambda_1 \frac{1}{\lambda_n}.$$

So we have

$$s(\mathbf{x}^{(k+1)}) = s(\mathbf{x}^{(k)}) \left(1 - \frac{1}{\kappa(\mathbf{A})}\right)$$

² This is a very slick proof that involves a number of interesting quantities; it's been designed over years to be clever and simple, so it's the sort of proof that would be hard to come up with yourself, so read through it a few times to see what is going on.

Recall that \mathbf{A} is symmetric and so $\mathbf{A}^{-T} = \mathbf{A}^{-1}$.

³ This is a good exercise. The only challenging step is the triangle inequality. The easy way to show this is to use the fact that $n(\mathbf{x}) = \|F^{-1} \mathbf{x}\|$ where F is the Cholesky factor of \mathbf{A}^{-1} (which exists because we have a symmetric positive definite \mathbf{A} and \mathbf{A}^{-1}).

which is exactly the same rate as Richardson. To improve this, we need a stronger bound on θ .

One such stronger bound is called the Kantorovich inequality. Let \mathbf{A} be a symmetric positive definite matrix and let \mathbf{v} be any vector with $\mathbf{v}^T \mathbf{v} = 1$, then

$$(\mathbf{v}^T \mathbf{A} \mathbf{v})(\mathbf{v}^T \mathbf{A}^{-1} \mathbf{v}) \leq \frac{(\lambda_1 + \lambda_n)^2}{4\lambda_1 \lambda_n}.$$

This gives us a better bound on θ and we get

$$1 - 1/\theta \leq \frac{(\lambda_1 - \lambda_n)^2}{(\lambda_1 + \lambda_n)^2} \leq \left(\frac{\kappa(\mathbf{A}) - 1}{\kappa(\mathbf{A}) + 1} \right)^2 \leq (1 - 1/\kappa(\mathbf{A}))^2.$$

This completes the proof.

1.3 COMPARISON

Which method is faster? For Richardson we have

$$\|\mathbf{e}^{(k)}\| \leq \left(1 - \frac{1}{\kappa}\right)^k \|\mathbf{e}^{(0)}\|$$

whereas for Steepest Descent we have

$$s(\mathbf{x}^{(k+1)}) \leq \left(1 - \frac{1}{\kappa}\right)^{2k} s(\mathbf{x}^{(0)}).$$

To conclude that steepest descent is faster requires a comparison of $s(\mathbf{x}^{(k)})$ and $\|\mathbf{e}^{(0)}\|$. We will continue to study this relationship.