

Lecture 8

Error Bounds and Residuals

In addition to being a reliable indicator of near singularity, the condition number also provides a quantitative bound for the error in the compounded solution to a linear system, as we will now see.

Let \mathbf{x} be the solution to $\mathbf{Ax} = \mathbf{b}$, and let $\hat{\mathbf{x}}$ be the solution to $\mathbf{A}\hat{\mathbf{x}} = \mathbf{b} + \Delta\mathbf{b}$ with a perturbed right hand side. If $\Delta\mathbf{x} = \hat{\mathbf{x}} - \mathbf{x}$, then

$$\begin{aligned}\mathbf{b} + \Delta\mathbf{b} &= \mathbf{A}\hat{\mathbf{x}} = \mathbf{A}(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{Ax} + \mathbf{A}\Delta\mathbf{x} \\ \text{i.e. } \Delta\mathbf{b} &= \mathbf{A}\Delta\mathbf{x} \Rightarrow \Delta\mathbf{x} = \mathbf{A}^{-1}\Delta\mathbf{b} \leq \|\mathbf{A}^{-1}\| \cdot \|\Delta\mathbf{b}\|\end{aligned}$$

Which leads to the bound

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \mathit{cond}(\mathbf{A}) \frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|}$$

for possible relative change in solution due to relative change in RSH \mathbf{b} .

Thus, the condition number of the matrix is an “amplification factor” that bounds the maximum relative change in the solution due to a given relative change in the RHS vector.

A similar result holds for relative change in the entries of the matrix \mathbf{A} . If $\mathbf{Ax} = \mathbf{b}$ and $(\mathbf{A} + \mathbf{E})\hat{\mathbf{x}} = \mathbf{b}$, then

$$\Delta\mathbf{x} = \hat{\mathbf{x}} - \mathbf{x} = \mathbf{A}^{-1}(\mathbf{Ax} - \mathbf{b}) = -\mathbf{A}^{-1}\mathbf{E}\hat{\mathbf{x}}$$

Taking norm, we get

$$\|\Delta\mathbf{x}\| \leq \|\mathbf{A}^{-1}\| \cdot \|\mathbf{E}\| \cdot \|\hat{\mathbf{x}}\| \text{ which gives the bound } \frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \mathit{cond}(\mathbf{A}) \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}$$

Result: the condition number times the relative change in the problem data bound the relative change in the solution.

One simple way of interpreting the results is that the computed solution loses about $\log_{10}(\mathit{cond}(\mathbf{A}))$ decimal digits of accuracy relative to the accuracy of the input. The example in the previous lecture, for instance has a condition number greater than 10^4 , so we would expect no correct digits in the solution to a linear system with the matrix unless the input data are accurate to more than four decimal digits and the solution is compared using arithmetic with more than four decimal digits of precision.

As a quantitative measure of sensitivity, the matrix condition number plays the same role for the problem of solving linear systems – and yields the same type of relationship between forward and backward error – as the general notion of condition number.

Note: The matrix condition number is never less than 1

Residuals

One way to verify a solution to an equation is to substitute it into the equation and see how closely left and right sides match.

Residual vector of approximate solution $\hat{\mathbf{x}}$ to linear system $\mathbf{Ax} = \mathbf{b}$ defined by

$$\mathbf{r} = \mathbf{b} - \mathbf{A}\hat{\mathbf{x}}$$

In theory, if \mathbf{A} is nonsingular, then $\|\Delta\mathbf{x}\| = \|\hat{\mathbf{x}} - \mathbf{x}\| = \mathbf{0}$, if and only if $\|\mathbf{r}\| = \mathbf{0}$. In practice, however these quantities are not necessarily small simultaneously.

Note that if the equation $\mathbf{Ax} = \mathbf{b}$ is multiplied by an arbitrary non-zero constant, then the solution is unaffected, but the residual is multiplied by same factor. Thus, the residual can be made arbitrarily large or small, depending on the scaling of the problem, and hence size of the residual is meaning less, unless it is considered relative to the size of the problem data and the solution. Thus the relative residual for the approximated solution $\hat{\mathbf{x}}$ is defined to be $\|\mathbf{r}\|/(\|\mathbf{A}\| \cdot \|\hat{\mathbf{x}}\|)$.

To relate the error to the residual, we observe that

$$\|\Delta\mathbf{x}\| = \|\hat{\mathbf{x}} - \mathbf{x}\| = \|\mathbf{A}^{-1}(\mathbf{A}\hat{\mathbf{x}} - \mathbf{b})\| = \|\mathbf{A}^{-1}\mathbf{r}\| \leq \|\mathbf{A}^{-1}\| \cdot \|\mathbf{r}\|$$

Dividing both sides by $\|\hat{\mathbf{x}}\|$ and using the definition of $\mathit{cond}(\mathbf{A})$, we then have

$$\frac{\|\Delta\mathbf{x}\|}{\|\hat{\mathbf{x}}\|} \leq \mathit{cond}(\mathbf{A}) \frac{\|\mathbf{r}\|}{\|\mathbf{A}\| \cdot \|\hat{\mathbf{x}}\|}$$

Thus, a small relative residual implies a small relative error in the solution when, and only when, \mathbf{A} is well conditioned.

To see the implications of large residual, on the other hand, if computed solution $\hat{\mathbf{x}}$ exactly satisfies

$$(\mathbf{A} + \mathbf{E})\hat{\mathbf{x}} = \mathbf{b}$$

$$\text{Then } \|\mathbf{r}\| = \|\mathbf{b} - \mathbf{A}\hat{\mathbf{x}}\| = \|\mathbf{E}\hat{\mathbf{x}}\| \leq \|\mathbf{E}\| \cdot \|\hat{\mathbf{x}}\|$$

Which gives

$$\frac{\|\mathbf{r}\|}{\|\mathbf{A}\| \cdot \|\hat{\mathbf{x}}\|} \leq \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}$$

So large relative residual implies large backward error in the matrix and algorithm used to compute solution is unstable. Another way of saying this is that a stable algorithm will invariably produce a solution with small relative residual, irrespective of the conditioning of the problem, and hence a small residual by itself, sheds little light on the quality of the approximate solution.

Example: Small residual

$$A\mathbf{x} = \begin{bmatrix} 0.913 & 0.659 \\ 0.457 & 0.330 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0.245 \\ 0.127 \end{bmatrix} = \mathbf{b}$$

Consider two approximate solutions

$$\hat{\mathbf{x}}_1 = \begin{bmatrix} -0.0827 \\ 0.5 \end{bmatrix} \text{ and } \hat{\mathbf{x}}_2 = \begin{bmatrix} 0.999 \\ -1.001 \end{bmatrix}$$

The norms of their respective residuals are

$$\|\mathbf{r}_1\|_1 = 2.1 \times 10^{-4} \text{ and } \|\mathbf{r}_2\|_1 = 2.4 \times 10^{-2}$$

so which is the better solution?

We are tempted to say $\hat{\mathbf{x}}_1$ because of its much smaller residual. But the exact solution to this system is $\mathbf{x} = [1, -1]^T$, as is easily confirmed, so $\hat{\mathbf{x}}_2$ is actually much more accurate than $\hat{\mathbf{x}}_1$. The reason for this surprising behavior is that the matrix A is ill-conditioned as we saw previously, and because of its large condition number, a small residual does not imply a small error in the solution.

Iterative refinement

Given approximate solution \mathbf{x}_0 to linear system $A\mathbf{x} = \mathbf{b}$, compute residual

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

Now solve linear system $A\mathbf{z}_0 = \mathbf{r}_0$ and take $\mathbf{x}_1 = \mathbf{x}_0 + \mathbf{z}_0$ as new and “better” approximate solution, since

$$A\mathbf{x}_1 = A(\mathbf{x}_0 + \mathbf{z}_0) = A\mathbf{x}_0 + A\mathbf{z}_0 = (\mathbf{b} - \mathbf{r}_0) + \mathbf{r}_0 = \mathbf{b}$$

Process can be repeated to refine solution successively until convergence, potentially producing solution accurate to full machine precision.

Iterative refinement requires double storage, since both original matrix and LU factorization required

Due to cancellation, residual usually must be computed with higher precision for iterative refinement to produce meaningful improvement.

For these reasons, iterative improvement often impractical to use routinely, but can still be useful in some circumstances.