

## Notes for 2016-09-02

### 1 Notions of error

The art of numerics is finding an approximation with a fast algorithm, a form that is easy to analyze, and an error bound. Given a task, we want to engineer an approximation that is good enough, and that composes well with other approximations. To make these goals precise, we need to define types of errors and error propagation, and some associated notation – which is the point of this lecture.

#### 1.1 Absolute and relative error

Suppose  $\hat{x}$  is an approximation to  $x$ . The *absolute error* is

$$e_{\text{abs}} = |\hat{x} - x|.$$

Absolute error has the same dimensions as  $x$ , and can be misleading without some context. An error of one meter per second is dramatic if  $x$  is my walking pace; if  $x$  is the speed of light, it is a very small error.

The *relative error* is a measure with a more natural sense of scale:

$$e_{\text{rel}} = \frac{|\hat{x} - x|}{|x|}.$$

Relative error is familiar in everyday life: when someone talks about an error of a few percent, or says that a given measurement is good to three significant figures, she is describing a relative error.

We sometimes estimate the relative error in approximating  $x$  by  $\hat{x}$  using the relative error in approximating  $\hat{x}$  by  $x$ :

$$\hat{e}_{\text{rel}} = \frac{|\hat{x} - x|}{|\hat{x}|}.$$

As long as  $\hat{e}_{\text{rel}} < 1$ , a little algebra gives that

$$\frac{\hat{e}_{\text{rel}}}{1 + \hat{e}_{\text{rel}}} \leq e_{\text{rel}} \leq \frac{\hat{e}_{\text{rel}}}{1 - \hat{e}_{\text{rel}}}.$$

If we know  $\hat{e}_{\text{rel}}$  is much less than one, then it is a good estimate for  $e_{\text{rel}}$ . If  $\hat{e}_{\text{rel}}$  is not much less than one, we know that  $\hat{x}$  is a poor approximation to  $x$ . Either way,  $\hat{e}_{\text{rel}}$  is often just as useful as  $e_{\text{rel}}$ , and may be easier to estimate.

Relative error makes no sense for  $x = 0$ , and may be too pessimistic when the property of  $x$  we care about is “small enough.” A natural intermediate between absolute and relative errors is the mixed error

$$e_{\text{mixed}} = \frac{|\hat{x} - x|}{|x| + \tau}$$

where  $\tau$  is some natural scale factor associated with  $x$ .

## 1.2 Errors beyond scalars

Absolute and relative error make sense for vectors as well as scalars. If  $\|\cdot\|$  is a vector norm and  $\hat{x}$  and  $x$  are vectors, then the (normwise) absolute and relative errors are

$$e_{\text{abs}} = \|\hat{x} - x\|, \quad e_{\text{rel}} = \frac{\|\hat{x} - x\|}{\|x\|}.$$

We might also consider the componentwise absolute or relative errors

$$e_{\text{abs},i} = |\hat{x}_i - x_i| \quad e_{\text{rel},i} = \frac{|\hat{x}_i - x_i|}{|x_i|}.$$

The two concepts are related: the maximum componentwise relative error can be computed as a normwise error in a norm defined in terms of the solution vector:

$$\max_i e_{\text{rel},i} = \|\|\hat{x} - x\|\|$$

where  $\|\|z\|\| = \|\text{diag}(x)^{-1}z\|$ . More generally, absolute error makes sense whenever we can measure distances between the truth and the approximation; and relative error makes sense whenever we can additionally measure the size of the truth. However, there are often many possible notions of distance and size; and different ways to measure give different notions of absolute and relative error. In practice, this deserves some care.

### 1.3 Dimensions and scaling

The first step in analyzing many application problems is *nondimensionalization*: combining constants in the problem to obtain a small number of dimensionless constants. Examples include the aspect ratio of a rectangle, the Reynolds number in fluid mechanics<sup>1</sup>, and so forth. There are three big reasons to nondimensionalize:

- Typically, the physics of a problem only really depends on dimensionless constants, of which there may be fewer than the number of dimensional constants. This is important for parameter studies, for example.
- For multi-dimensional problems in which the unknowns have different units, it is hard to judge an approximation error as “small” or “large,” even with a (normwise) relative error estimate. But one can usually tell what is large or small in a non-dimensionalized problem.
- Many physical problems have dimensionless parameters much less than one or much greater than one, and we can approximate the physics in these limits. Often when dimensionless constants are huge or tiny and asymptotic approximations work well, naive numerical methods work poorly. Hence, nondimensionalization helps us choose how to analyze our problems — and a purely numerical approach may be silly.

## 2 Forward and backward error

We often approximate a function  $f$  by another function  $\hat{f}$ . For a particular  $x$ , the *forward* (absolute) error is

$$|\hat{f}(x) - f(x)|.$$

In words, forward error is the function *output*. Sometimes, though, we can think of a slightly wrong *input*:

$$\hat{f}(x) = f(\hat{x}).$$

In this case,  $|x - \hat{x}|$  is called the *backward* error. An algorithm that always has small backward error is *backward stable*.

---

<sup>1</sup>Or any of a dozen other named numbers in fluid mechanics. Fluid mechanics is a field that appreciates the power of dimensional analysis

A *condition number* is a tight constant relating relative output error to relative input error. For example, for the problem of evaluating a sufficiently nice function  $f(x)$  where  $x$  is the input and  $\hat{x} = x + h$  is a perturbed input (relative error  $|h|/|x|$ ), the condition number  $\kappa[f(x)]$  is the smallest constant such that

$$\frac{|f(x+h) - f(x)|}{|f(x)|} \leq \kappa[f(x)] \frac{|h|}{|x|} + o(|h|)$$

If  $f$  is differentiable, the condition number is

$$\kappa[f(x)] = \lim_{h \neq 0} \frac{|f(x+h) - f(x)|/|f(x)|}{|(x+h) - x|/|x|} = \frac{|f'(x)||x|}{|f(x)|}.$$

If  $f$  is Lipschitz in a neighborhood of  $x$  (locally Lipschitz), then

$$\kappa[f(x)] = \frac{M_{f(x)}|x|}{|f(x)|}.$$

where  $M_f$  is the smallest constant such that  $|f(x+h) - f(x)| \leq M_f|h| + o(|h|)$ . When the problem has no linear bound on the output error relative to the input error, we say the problem has an *infinite* condition number. An example is  $x^{1/3}$  at  $x = 0$ .

A problem with a small condition number is called *well-conditioned*; a problem with a large condition number is *ill-conditioned*. A backward stable algorithm applied to a well-conditioned problem has a small forward error.

### 3 Perturbing matrix problems

To make the previous discussion concrete, suppose I want  $y = Ax$ , but because of a small error in  $A$  (due to measurement errors or roundoff effects), I instead compute  $\hat{y} = (A + E)x$  where  $E$  is “small.” The expression for the *absolute* error is trivial:

$$\|\hat{y} - y\| = \|Ex\|.$$

But I usually care more about the *relative error*.

$$\frac{\|\hat{y} - y\|}{\|y\|} = \frac{\|Ex\|}{\|y\|}.$$

If we assume that  $A$  is invertible and that we are using consistent norms (which we will usually assume), then

$$\|Ex\| = \|EA^{-1}y\| \leq \|E\|\|A^{-1}\|\|y\|,$$

which gives us

$$\frac{\|\hat{y} - y\|}{\|y\|} \leq \|A\| \|A^{-1}\| \frac{\|E\|}{\|A\|} = \kappa(A) \frac{\|E\|}{\|A\|}.$$

That is, the relative error in the output is the relative error in the input multiplied by the condition number  $\kappa(A) = \|A\| \|A^{-1}\|$ . Technically, this is the condition number for the problem of matrix multiplication (or solving linear systems, as we will see) with respect to a particular (consistent) norm; different problems have different condition numbers. Nonetheless, it is common to call this “the” condition number of  $A$ .

For some problems, we are given more control over the structure of the error matrix  $E$ . For example, we might suppose that  $A$  is symmetric, and ask whether we can get a tighter bound if in addition to assuming a bound on  $\|E\|$ , we also assume  $E$  is symmetric. In this particular case, the answer is “no” — we have the same condition number either way, at least for the 2-norm or Frobenius norm<sup>2</sup>. In other cases, assuming a structure to the perturbation does indeed allow us to achieve tighter bounds.

As an example of a refined bound, we consider moving from condition numbers based on small norm-wise perturbations to condition numbers based on small *element-wise* perturbations. Suppose  $E$  is elementwise small relative to  $A$ , i.e.  $|E| \leq \epsilon |A|$ . Suppose also that we are dealing with a norm such that  $\|X\| \leq \| |X| \|$ , as is true of all the norms we have seen so far. Then

$$\frac{\|\hat{y} - y\|}{\|y\|} \leq \|EA^{-1}\| \leq \| |A| |A^{-1}| \| \epsilon.$$

The quantity  $\kappa_{\text{rel}}(A) = \| |A| |A^{-1}| \|$  is the *relative condition number*; it is closely related to the *Skeel condition number* which we will see in our discussion of linear systems<sup>3</sup>. Unlike the standard condition number, the relative condition number is invariant under column scaling of  $A$ ; that is  $\kappa_{\text{rel}}(AD) = \kappa_{\text{rel}}(A)$  where  $D$  is a nonsingular diagonal matrix.

What if, instead of perturbing  $A$ , we perturb  $x$ ? That is, if  $\hat{y} = A\hat{x}$  and  $y = Ax$ , what is the condition number relating  $\|\hat{y} - y\|/\|y\|$  to  $\|\hat{x} - x\|/\|x\|$ ? We note that

$$\|\hat{y} - y\| = \|A(\hat{x} - x)\| \leq \|A\| \|\hat{x} - x\|;$$

<sup>2</sup>This is left as an exercise for the student

<sup>3</sup>The Skeel condition number involves the two factors in the reverse order.

and

$$\|x\| = \|A^{-1}y\| \leq \|A^{-1}\| \|y\| \quad \implies \quad \|y\| \geq \|A^{-1}\|^{-1} \|x\|.$$

Put together, this implies

$$\frac{\|\hat{y} - y\|}{\|y\|} \leq \|A\| \|A^{-1}\| \frac{\|\hat{x} - x\|}{\|x\|}.$$

The same condition number appears again!