

Linear algebra

3/10/11

15

We shall turn our attention to solving linear systems of equations $Ax = b$

$$A \in \mathbb{R}^{m \times n}$$

$$x \in \mathbb{R}^n$$

$$b \in \mathbb{R}^m.$$

We already saw examples of methods that required the solution of a linear system as part of the overall algorithm, e.g. the Vandermonde system for interpolation, which was a square system ($m=n$).

Another category of methods that leads to rectangular systems with $m > n$ is least square methods. They answer questions of the form:

→ What is the best n -order polynomial we can use to approximate (not interpolate) $(m+1)$ data points (where $m > n$).

→ More generally, find the solution that most closely satisfies m equations, in the presence of n ($n < m$) unknowns.

All these algorithms need to be conscious

3/10/11

26

• about error, and there are at least 3 sources for it.

→ Some algorithms are "imperfect" in the sense that they require several iterations to generate a good quality approximation. Thus, intermediate results are subject to error

→ Sometimes, it is not possible to find an "ideal" solution, e.g. because we have more equations than unknowns. In this case, not all equations will be satisfied exactly, and we need a notion of the "error" incurred in not satisfying certain equations fully.

→ Inputs to an algorithm are often corrupted by noise, roundoff error, etc. For example, instead of solving an "intended" system $Ax = b$ we may be solving $A^*x = b^*$ where the entries in A^* & b^* have been subject to noise and inaccuracy. It is important to know how those translate to errors in determining x .

(Vector and Matrix) Norms :

3/10/11

Norms are valuable tools in arguing about the extent and magnitude of error. We will introduce some concepts that we will use broadly later on.

Definition A vector norm is a function from \mathbb{R}^n to \mathbb{R} , with a certain number of properties. If $\underline{x} \in \mathbb{R}^n$, we symbolize its norm by $\|\underline{x}\|$. The defining properties of a norm are:

$$(i) \quad \|\underline{x}\| \geq 0 \quad \text{for all } \underline{x} \in \mathbb{R}^n$$

$$\text{also } \|\underline{x}\| = 0 \quad \text{iff } \underline{x} = \underline{0}$$

$$(ii) \quad \|\alpha \underline{x}\| = |\alpha| \cdot \|\underline{x}\| \quad \text{for all } \begin{array}{l} \alpha \in \mathbb{R} \\ \underline{x} \in \mathbb{R}^n \end{array}$$

$$(iii) \quad \|\underline{x} + \underline{y}\| \leq \|\underline{x}\| + \|\underline{y}\| \quad \text{for all } \underline{x}, \underline{y} \in \mathbb{R}^n$$

~~Note~~. Note that the properties above do not determine a unique form of a "norm" function, in fact many different valid norms exist. Typically, we will use subscripts ($\|\cdot\|_a, \|\cdot\|_b$) to denote different types of norms.

Vector norms - Why are they needed?

When dealing e.g. with the solution of a nonlinear equation $f(x) = 0$, the error $e = x_{\text{approx}} - x_{\text{exact}}$ is a single number, thus the absolute value $|e|$ gives us a good idea of the "extent" of error.

When solving a system of linear equations $A\underline{x} = \underline{b}$, the exact solution $\underline{x}_{\text{exact}}$ as well as any approximation $\underline{x}_{\text{approx}}$ are vectors, and the error:

$$\underline{e} = \underline{x}_{\text{approx}} - \underline{x}_{\text{exact}}$$

is a vector, too. It is not as straightforward to assess the "magnitude" of such a vector-valued error.

e.g. Consider $\underline{e}_1, \underline{e}_2 \in \mathbb{R}^{1000}$, and

$$\underline{e}_1 = \begin{pmatrix} 0.1 \\ 0.1 \\ 0.1 \\ \vdots \\ 0.1 \end{pmatrix} \quad \underline{e}_2 = \begin{pmatrix} 100 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Which one is worse? \underline{e}_1 has a modest amount of

error, distributed over all components. In \underline{e}_2 , all but one component are exact, but one of them has a huge discrepancy.

Exactly how we quantify and assess the extent of error is application-dependent. Vector norms are alternative ways to measure this magnitude, and different norms would be appropriate for different tasks.

Def. A vector norm is a function from $\mathbb{R}^n \rightarrow \mathbb{R}^+$, which maps a vector \underline{v} to the real number $\|\underline{v}\|$. This symbol must satisfy the properties:

- (i) $\|\underline{x}\| \geq 0$ for all $\underline{x} \in \mathbb{R}^n$. Also $\|\underline{x}\| = 0$ iff $\underline{x} = 0$
- (ii) $\|\alpha \underline{x}\| = |\alpha| \cdot \|\underline{x}\| \quad \forall \underline{x} \in \mathbb{R}^n, \alpha \in \mathbb{R}$
- (iii) $\|\underline{x} + \underline{y}\| \leq \|\underline{x}\| + \|\underline{y}\| \quad \forall \underline{x}, \underline{y} \in \mathbb{R}^n$ (triangle inequality).

3/22/2011 3

Some norms which can be proven to satisfy these properties, are: (Here $\underline{x} = (x_1, x_2, \dots, x_n)$)

1. The L_1 -norm (or 1-norm)

$$\|\underline{x}\|_1 = \sum_{i=1}^n |x_i|$$

2. The L_2 -norm (or 2-norm, or Euclidean norm)

$$\|\underline{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

3. The infinity norm (or max-norm)

$$\|\underline{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|$$

4. (Less common) L_p -norm

$$\|\underline{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}$$

It is relatively easy to show that these satisfy the defining properties of a norm. e.g. for $\|\cdot\|_1$: 3/22/11 L4

$$\bullet \|\underline{x}\|_1 = \sum_{i=1}^n |x_i| \geq 0$$

$$\bullet \text{ if } \underline{x} = \underline{0}, \text{ then } \|\underline{x}\|_1 = 0$$

$$\text{if } \|\underline{x}\|_1 = 0 \Rightarrow \sum_{i=1}^n |x_i| = 0 \Rightarrow |x_i| = 0 \forall i \Rightarrow \underline{x} = \underline{0}$$

$$\bullet \|\alpha \underline{x}\|_1 = \sum_{i=1}^n |\alpha x_i| = |\alpha| \sum_{i=1}^n |x_i| = |\alpha| \|\underline{x}\|_1$$

$$\bullet \|\underline{x} + \underline{y}\|_1 = \sum_{i=1}^n |x_i + y_i| \leq \sum_{i=1}^n |x_i| + \sum_{i=1}^n |y_i| = \|\underline{x}\|_1 + \|\underline{y}\|_1$$

Similar proofs can be given for $\|\cdot\|_\infty$ (just as easy),
 $\|\cdot\|_2$ (a bit more difficult) and $\|\cdot\|_p$ (rather complicated).

3/22/11 LS

We can actually define norms for (square) matrices, as well. A matrix norm is a function

$\| \cdot \| : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$ which satisfies:

- (i) $\|M\| \geq 0 \ \forall M \in \mathbb{R}^{n \times n}$. $\|M\| = 0$ iff $M = 0$
- (ii) $\|\alpha M\| = |\alpha| \|M\|$
- (iii) $\|M + N\| \leq \|M\| + \|N\|$
- (iv) $\|M \cdot N\| \leq \|M\| \cdot \|N\|$.

(Property (iv) is the one that has slightly different flavor than vector norms).

Although more types of matrix norms exist, one common category is that of matrix norms induced by vector norms.

Def. If $\| \cdot \|_*$ is a valid vector norm, its induced matrix norm is defined as

$$\|M\|_* = \max_{x \in \mathbb{R}^n, x \neq 0} \left\{ \frac{\|Mx\|_*}{\|x\|_*} \right\}$$

or equivalently:

$$\|M\|_* = \max_{\substack{x \in \mathbb{R}^n \\ \|x\|_* = 1}} \{ \|Mx\|_* \}$$

Note, again, that not all valid matrix norms are induced by vector norms. One notable example is the very commonly used Frobenius norm:

$$\|M\|_F = \sqrt{\sum_{i,j=1}^n M_{ij}^2}$$

We can easily show though that induced norms satisfy properties (i) through (iv). (i)-(iii) are rather trivial, eg:

$$\begin{aligned} \|M+N\| &= \max_{\underline{x} \neq 0} \frac{\|(M+N)\underline{x}\|}{\|\underline{x}\|} \leq \max_{\underline{x} \neq 0} \frac{\|M\underline{x}\| + \|N\underline{x}\|}{\|\underline{x}\|} \\ &= \max_{\underline{x} \neq 0} \frac{\|M\underline{x}\|}{\|\underline{x}\|} + \max_{\underline{x} \neq 0} \frac{\|N\underline{x}\|}{\|\underline{x}\|} = \|M\| + \|N\|. \end{aligned}$$

Property (iv) is slightly trickier to show.

First, a lemma:

Lemma If $\|\cdot\|$ is a matrix norm induced by a vector norm $\|\cdot\|$, then:

$$\|Ax\| \leq \|A\| \cdot \|x\|$$

Proof: Since $\|A\| = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|}$, we have that

for an arbitrary $y \in \mathbb{R}^m$: ($y \neq 0$)

$$\|A\| = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} \geq \frac{\|Ay\|}{\|y\|} \Rightarrow$$

$\Rightarrow \|Ay\| \leq \|A\| \|y\|$. This holds for $y \neq 0$,

but we can see it is also true for $y = \underline{0}$

Prop. (iv)

$$\|MN\| = \max_{\|x\|=1} \|MNx\| \leq \max_{\|x\|=1} \|M\| \|Nx\| =$$

$$= \|M\| \cdot \max_{\|x\|=1} \|Nx\| = \|M\| \cdot \|N\| \Rightarrow$$

$$\|MN\| \leq \|M\| \cdot \|N\|.$$

Although the definition of an induced norm allowed us to prove certain properties, it does not necessarily provide a convenient formula for evaluating the matrix norm.

Fortunately, such formulas do exist for the L_1 and L_∞ induced matrix norms. Given here without proof:

$$\|A\|_1 = \max_j \sum_{i=1}^n |A_{ij}| \quad (\text{max. absolute column sum})$$

$$\|A\|_\infty = \max_i \sum_{j=1}^n |A_{ij}| \quad (\text{max. absolute row sum})$$

($\|\cdot\|_2$ is much more complicated!)

Where do these vector/matrix norms come handy?

Useful properties of matrix & vector norms

3/24/11 L

We previously saw that

$$\|Ax\| \leq \|A\| \cdot \|x\| \quad (1)$$

for any matrix A , and any vector x (of dimensions $m \times m$ and $m \times 1$, respectively).

Note that, when writing an expression such as (1), the matrix norm $\|A\|$ is understood to be the inferred norm from the vector norm used in $\|Ax\|$ and $\|x\|$. Thus

$$\|Ax\|_1 \leq \|A\|_1 \cdot \|x\|_1$$

and

$$\|Ax\|_\infty \leq \|A\|_\infty \cdot \|x\|_\infty$$

are both valid, but we cannot mix and match, e.g.:

$$\cancel{\|Ax\|_\infty \leq \|A\|_2 \cdot \|x\|_1} \rightarrow \text{NOT CORRECT}$$

When solving a linear system $A\underline{x} = \underline{b}$, computer algorithms are only providing an approximation (\underline{x}_{app}) to the exact solution (\underline{x}_{ex}). This is due to factors such as finite precision, roundoff errors or even imperfect solution algorithms. In either case, we have an error (error vector, in fact) defined as

$$\underline{e} = \underline{x}_{ap} - \underline{x}_{ex}$$

Naturally, we would like to have an understanding of the magnitude of this error (e.g. some appropriate norm $\|\underline{e}\|$). The problem is that we do not know the exact, pristine solution \underline{x}_{ex} !

One remedy is offered via the residual vector defined as:

$$\underline{r} = \underline{b} - A\underline{x}_{app}$$

The vector \underline{r} is something we can compute practically since it involves only known quantities ($\underline{b}, A, \underline{x}_{app}$).

Furthermore, we have:

$$\begin{aligned}
 \underline{r} &= \underline{b} - A \underline{x}_{\text{app}} \\
 &= A \underline{x}_{\text{ex}} - A \underline{x}_{\text{app}} \\
 &= -A (\underline{x}_{\text{app}} - \underline{x}_{\text{ex}}) \\
 &= -A \underline{e} \quad \Rightarrow \quad \underline{r} = -A \underline{e} \\
 &\quad \underline{e} = -A^{-1} \underline{r}
 \end{aligned}$$

The last equation links the error with the residual.

Furthermore, we can write

$$\|\underline{e}\| = \|A^{-1} \underline{r}\| \leq \|A^{-1}\| \|\underline{r}\|$$

This equation provides a bound for the error, as a function of $\|A^{-1}\|$ and the norm of the computable vector \underline{r} ! Note that:

→ We can obtain this estimate without knowing the exact solution, but

→ We need $\|A^{-1}\|$ and generally, computing A^{-1} is just as difficult (if not more) than finding $\underline{x}_{\text{ex}}$. However there are special cases where an estimate of $\|A^{-1}\|$ can be obtained.

A different source of error :

Sometimes, the right-hand-side (b) of $Ax = b$ has errors that make it deviate from its intended value. For example in the Vandermonde matrix method for polynomial interpolation, b contains the samples $(y_1 = f(x_1), y_2, \dots, y_n)$ where $y_i = f(x_i)$. An error in a measuring device supposed to sample $f(x)$ could lead to erroneous readings y_i^* instead of y_i . In general, measuring inaccuracies can lead to the right-hand-side vector \underline{b} being mis-represented as $\underline{b}^* (\neq b)$.

In this case, instead of the intended solution $x = A^{-1}b$ we in fact compute $x^* = A^{-1}b^*$.

How important is the error $e = x^* - x$ that is caused by this misrepresentation of b ?

Let us introduce some notation:

3/24/11

5

$$\text{Let } \underline{\delta b} := \underline{b^*} - \underline{b}$$

$$\underline{\delta x} := \underline{x^*} - \underline{x}$$

$$A \underline{x} = \underline{b}$$

$$A \underline{x^*} = \underline{b^*}$$

$$A(\underline{x^*} - \underline{x}) = \underline{b^*} - \underline{b}$$

$$A \underline{\delta x} = \underline{\delta b}$$

$$\underline{\delta x} = A^{-1} \underline{\delta b}$$

Taking norms:

$$\|\underline{\delta x}\| = \|A^{-1} \underline{\delta b}\| \leq \|A^{-1}\| \|\underline{\delta b}\| \quad (1)$$

Thus the error in the computed solution $\underline{\delta x}$ is proportional to the error in \underline{b} .

An even more relevant question is: How does the relative error $\frac{\|\underline{\delta x}\|}{\|\underline{x}\|} = \frac{\|\underline{x^*} - \underline{x}\|}{\|\underline{x}\|}$ compare to the

relative error in \underline{b} $\frac{\|\underline{\delta b}\|}{\|\underline{b}\|}$? This may be more

useful to know, since $\|\underline{\delta b}\|$ may be impossible to compute (if we don't know the real \underline{b} !).

For this, we write

3/24/11 16

$$Ax = b \Rightarrow \|b\| = \|Ax\| \leq \|A\| \cdot \|x\|$$

$$\Rightarrow \frac{1}{\|x\|} \leq \|A\| \cdot \frac{1}{\|b\|} \quad (2)$$

Multiplying (1) & (2) we get:

$$\frac{\|\delta x\|}{\|x\|} \leq \underbrace{\|A\| \cdot \|A^{-1}\|}_{\kappa(A)} \cdot \frac{\|\delta b\|}{\|b\|}$$

Thus the relative error in \underline{x} is bounded by a multiple of the relative error in \underline{b} ! The multiplicative constant $\kappa(A) = \|A\| \cdot \|A^{-1}\|$ is called the condition number of A , and is an important measure of the sensitivity of a linear system $Ax=b$ to being solved on a computer, in the presence of inaccurate values.

e.g. If the relative error $\frac{\|\delta b\|}{\|b\|}$ is 0.0001%, but $\kappa(A) = 100,000$ (could happen!), Then we could have up to a 10% error in the computed \underline{x} !

Why is this always relevant?

3/24/11

7

Simply, almost any b will have some, small relative error, due to the fact it is represented on a computer up to machine precision! The relative error will be at least as much as the machine epsilon due to roundoff!

$$\frac{\|\delta b\|_{\infty}}{\|b\|} \geq \epsilon \approx 10^{-7} \quad (\text{for floats}).$$

But, how bad can the condition number get?

Very bad, at times. For example:

Hilbert matrices $H_n \in \mathbb{R}^{n \times n}$ $(H_n)_{ij} = \frac{1}{i+j-1}$

$$H_5 = \begin{bmatrix} 1 & 1/2 & 1/3 & 1/4 & 1/5 \\ 1/2 & 1/3 & 1/4 & 1/5 & 1/6 \\ 1/3 & 1/4 & 1/5 & 1/6 & 1/7 \\ 1/4 & 1/5 & 1/6 & 1/7 & 1/8 \\ 1/5 & 1/6 & 1/7 & 1/8 & 1/9 \end{bmatrix}$$

$$\kappa_{\infty}(H_5) = \|H_5\|_{\infty} \cdot \|H_5^{-1}\|_{\infty} \approx 10^6 !$$

Thus, any attempt at solving $H_5 x = b$ would be subject to a relative error up to 10% just due to roundoff errors in b !

Another case: near-singular matrices

3/24/11

LE

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 6+\varepsilon \end{bmatrix}$$

As $\varepsilon \rightarrow 0$ A becomes singular (non-invertible)

In this case, $\kappa(A) \rightarrow \infty$.

What is the best case for $\kappa(A)$?

Lemma For any vector-induced matrix norm, we have $\|I\| = 1$.

Proof From definition:

$$\|I\| = \max_{x \neq 0} \frac{\|Ix\|}{\|x\|} = \max_{x \neq 0} \frac{\|x\|}{\|x\|}$$

Using property (iv) of matrix norms, we get:

$$I = A \cdot A^{-1} \Rightarrow 1 = \|I\| = \|A \cdot A^{-1}\| \leq \|A\| \cdot \|A^{-1}\|$$

Thus $\boxed{\kappa(A) \geq 1}$

The "best" conditioned matrices are of the form $A = c \cdot I$!
And have $\kappa(A) = 1$.