

so  $K$  is a square root of  $M^*M$ . There is only one positive semidefinite square root of any positive semidefinite matrix, so the  $K = Q\Sigma Q^*$  we specified is the unique choice.

If  $M$ , and hence  $K$ , is invertible the matrix  $U = MK^{-1}$  is also unique.

Likewise,  $K' = \sqrt{M^*M}$  is always unique.

When  $M$  is normal we have  $K = K'$  and so  $M = UK = KU$ .

**Any complex matrix  $M$**  can be written as  $M = UK = K'U$  where  $U$  is **unitary** and  $K$  and  $K'$  are **positive semidefinite**.

$K = \sqrt{MM^*}$  and  $K' = \sqrt{M^*M}$  so the positive semidefiniteness condition makes them unique.

These are called **Polar Decompositions of  $M$** .

Note that if  $M$  is normal then  $K = K'$  and  $M = UK = KU$ .

If  $M$  (and hence  $K$ ) is invertible  $U$  is unique.

## 67. Operator Norm and Error in Matrix Equations

### Operator Norm.

A **norm** on any real vector space  $V$  is a function

$$\|\cdot\|: V \rightarrow [0, \infty)$$

with the following properties.

- (1)  $\|\mathbf{v}\| = 0$  if and only if  $\mathbf{v} = 0$ .
- (2)  $\|c\mathbf{v}\| = |c|\|\mathbf{v}\|$  for all real  $c$  and any vector  $\mathbf{v}$ .
- (3)  $\|\mathbf{v} + \mathbf{w}\| \leq \|\mathbf{v}\| + \|\mathbf{w}\|$  for any pair of vectors  $\mathbf{v}$  and  $\mathbf{w}$ .

We defined norm earlier for inner product spaces, and created the norm using the inner product. But norms can arise elsewhere too.

67.1. **Exercise.** Prove that for any norm

$$|\|\mathbf{v}\| - \|\mathbf{w}\|| \leq \|\mathbf{v} - \mathbf{w}\|.$$

If  $M \in \mathbb{M}_{m \times n}(\mathbb{R})$ , a real  $m \times n$  matrix, we know that  $M$  has a singular value decomposition. Specifically, there is an orthogonal  $m \times m$  matrix  $P$  and orthogonal  $n \times n$  matrix  $Q$  and diagonal  $\Sigma$  so that

$$M = P\Sigma Q.$$

The diagonal entries of  $\Sigma$ , the numbers  $\sigma_1, \dots, \sigma_n$ , are the singular values. They form non-negative and decreasing sequence.

We define  $\|M\|_{op} = \sigma_1$ , the first (and largest) of these singular values. This is called the **operator norm of  $M$** .

**67.2. Exercise.** Suppose  $c_1 > c_2 > \dots > c_k > 0$ . Show that  $c_1$  is the maximum of all sums of the form  $c_i b^i$  where the  $b^i$  are non-negative and add to 1.

We can use the exercise to create an alternative definition of  $\|M\|_{op}$ , which does not mention the SVD decomposition. The claim is that

$$\|M\|_{op} = \text{maximum value of } \{ \|M\mathbf{u}\| \mid \mathbf{u} \text{ is a unit vector} \}.$$

Since  $P$  and  $Q$  in the equation  $M = P\Sigma Q$  are orthogonal, they don't change the length of vectors. So

$$\|M\mathbf{v}\| = \|P\Sigma Q\mathbf{v}\| = \|\Sigma Q\mathbf{v}\| = \|\Sigma\mathbf{w}\|$$

for vector  $\mathbf{w} = Q^{-1}\mathbf{v} = Q^T\mathbf{v}$ , which has the same length as  $\mathbf{v}$ .

$$\|M\mathbf{v}\|^2 = \|\Sigma\mathbf{w}\|^2 = (\Sigma\mathbf{w}) \cdot (\Sigma\mathbf{w}) = \sum_{i=1}^n (\sigma_i)^2 (w^i)^2.$$

Since  $\sum_{i=1}^n (w^i)^2 = 1$ , if we are looking for the maximum value which can be obtained it must be that the only nonzero coefficients among the  $w^i$  correspond to the singular value  $\sigma_1$ , in which case the sum itself is  $\sigma_1$ , as claimed.

We still have not verified that  $\|\cdot\|_{op}$  is actually a norm.

It is obvious that it satisfies items (1) and (2) of the requirements for a norm. It remains to verify (3), the triangle inequality.

Suppose  $K$  and  $M$  are  $m \times n$  matrices and  $\mathbf{v}$  is a unit vector for which  $\|K + M\|_{op} = \|(K + M)\mathbf{v}\|$ . Then we have

$$\|K + M\|_{op} = \|(K + M)\mathbf{v}\| \leq \|K\mathbf{v}\| + \|M\mathbf{v}\| \leq \|K\|_{op} + \|M\|_{op}$$

and the triangle inequality holds for operator norm.

Operator norm has a feature not possessed by many norms. It is called **sub-multiplicative** by virtue of the fact that

$$\|KM\|_{op} \leq \|K\|_{op}\|M\|_{op}$$

whenever  $K$  and  $M$  are compatible matrices for multiplication in this order. Since most vector spaces fail to possess a multiplication between members, there is no possibility of this inequality for most norms.

Assuming  $\|KM\|_{op} \neq 0$  select  $\mathbf{v}$  so that  $\|KM\|_{op} = \|KM\mathbf{v}\|$ . Then

$$\|KM\|_{op} = \|KM\mathbf{v}\| = \|M\mathbf{v}\| \left\| K \left( \frac{M\mathbf{v}}{\|M\mathbf{v}\|} \right) \right\| \leq \|K\|_{op} \|M\|_{op}.$$

*This implies that  $\mathbf{1} \leq \|A\|_{op} \|A^{-1}\|_{op}$  for invertible  $A$ .*

Suppose  $m \times n$  matrix  $B$  has entries  $b_{i,j}$  and rows  $R_i$ . Suppose every entry of  $B$  is smaller than  $\varepsilon$  in magnitude and  $\mathbf{w}$  is a unit vector. Let  $\mathbf{y} = B\mathbf{w}$ . Then

$$\begin{aligned} \|\mathbf{y}\|^2 &= \mathbf{y} \cdot \mathbf{y} = y^i y^i = \sum_{i=1}^m (R_i \mathbf{w})^2 = \sum_{i=1}^m (R_i^T \cdot \mathbf{w})^2 \leq \sum_{i=1}^m (\|R_i^T\| \|\mathbf{w}\|)^2 \\ &= \sum_{i=1}^m (\|R_i^T\|)^2 = \sum_{i=1}^m \sum_{j=1}^n b_{i,j}^2 < mn\varepsilon^2. \end{aligned}$$

This implies that  $\|B\|_{op}$  cannot exceed  $\varepsilon\sqrt{mn}$ . Together with the triangle inequality we have

$$\|A\|_{op} - \varepsilon\sqrt{mn} \leq \|A + B\|_{op} \leq \|A\|_{op} + \varepsilon\sqrt{mn}.$$

In other words, two matrices whose entries are all close have operator norms that are close too.

On the other hand, suppose a single entry of matrix  $B = C - A$  exceeds  $\delta$ . Suppose that entry is  $b_{i,j}$ . Then  $B\mathbf{e}_j$ , which is the  $j$ th column of  $B$ , has magnitude exceeding  $\delta$ . Since  $\mathbf{e}_j$  is a unit vector that means  $\|B\|_{op} > \delta$ .

*So two matrices  $A$  and  $C$  are “close” in the sense that  $\|C - A\|_{op}$  is small if and only if they are “close” in the sense that every entry in  $A$  is close to the corresponding entry in  $C$ .*

### Condition Number and Error.

If  $A$  is an invertible matrix with operator norm we define the **condition number** of  $A$  to be

$$\mathbf{K}_A = \|A\|_{op} \|A^{-1}\|_{op}$$

Any of the standard mathematical software packages will produce this condition number.

**67.3. Exercise.** Suppose  $A$  is invertible with non-increasing positive singular values  $\sigma_1, \sigma_2, \dots, \sigma_n$ . Show that  $\mathbf{K}_A = \frac{\sigma_1}{\sigma_n}$ .

Condition numbers cannot be smaller than 1, but a very large condition number indicates that the matrix has features that interfere with the accuracy of a computation in situations where finite-precision arithmetic (usually to 16 places) rather than “perfect arithmetic” is being performed, and suggests that any calculation involving the matrix should be used with caution. Such matrices are called **ill-conditioned**.

We suppose given a matrix equation  $A\mathbf{x} = \mathbf{b}$  for invertible matrix  $A$ . The solution is given by  $\mathbf{x} = A^{-1}\mathbf{b}$  in terms of the target  $\mathbf{b}$ .

It may be that entries of coefficient matrix  $A$  are known exactly but there is variation  $\delta\mathbf{b}$  in knowledge of the target vector. We don't know  $\delta\mathbf{b}$  exactly, but we do have an idea of the range of possible variation: the maximum possible magnitude of the entries of  $\delta\mathbf{b}$ .

We suppose  $\mathbf{y} = \mathbf{x} + \delta\mathbf{x}$  where  $A\mathbf{y} = \mathbf{c} = \mathbf{b} + \delta\mathbf{b}$ .

So  $\mathbf{y}$  is the solution for erroneous target  $\mathbf{c}$ , and so probably differs from the solution  $\mathbf{x}$  we really would like to know for the true target  $\mathbf{b}$ .

But by how much could it differ?

We calculate that

$$A\mathbf{y} = A(\mathbf{x} + \delta\mathbf{x}) = \mathbf{c} = \mathbf{b} + \delta\mathbf{b} \implies A\delta\mathbf{x} = \delta\mathbf{b} \implies \delta\mathbf{x} = A^{-1}\delta\mathbf{b}.$$

So  $\|\delta\mathbf{x}\| \leq \|A^{-1}\|_{op} \|\delta\mathbf{b}\|$ . And also  $A\mathbf{y} = \mathbf{c}$  so  $\|A\|_{op} \|\mathbf{y}\| \geq \|\mathbf{c}\|$ .

Then we have

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{y}\|} \leq \|A\|_{op} \|A^{-1}\|_{op} \frac{\|\delta\mathbf{b}\|}{\|\mathbf{c}\|} = K_A \frac{\|\delta\mathbf{b}\|}{\|\mathbf{c}\|}$$

which is an expression for the worst-case relative “error” of calculated solution  $\mathbf{y}$  in terms of the relative uncertainty of the “target” vector  $\mathbf{c}$  and the condition number of exactly known matrix  $A$ .

We now suppose that  $\mathbf{b}$  is known exactly but matrix  $A$  has variation  $\delta A$ . We calculate our solution using slightly erroneous but invertible matrix  $C = A + \delta A$ .

We don't actually know  $\delta A$  but we do understand our ability to do the measurements that produced the entries of  $C$ , so we have an estimate of the maximum size of the entries of  $\delta A$ , and we assume  $\delta A$  to be small enough so that  $C + \delta B$  is invertible whenever  $\|\delta B\|_{op}$  does not exceed the maximum possible value of  $\|\delta A\|_{op}$ . We have matrix  $C$  in hand and know it is invertible, and with this supposition the (unknown but exact) matrix  $A$  is too.

Suppose, again, that  $\mathbf{x}$  is the solution we really want,  $A\mathbf{x} = \mathbf{b}$ , and  $\mathbf{y} = \mathbf{x} + \delta\mathbf{x}$  is the solution we *have*, with

$$C\mathbf{y} = (A + \delta A)\mathbf{y} = (A + \delta A)(\mathbf{x} + \delta\mathbf{x}) = \mathbf{b}.$$

$$\begin{aligned}
\text{Then } (A + \delta A)\mathbf{y} &= \mathbf{b} \\
\implies A\mathbf{x} + A\delta\mathbf{x} + \delta A\mathbf{y} &= \mathbf{b} \implies A\delta\mathbf{x} = -\delta A\mathbf{y} \\
\implies \delta\mathbf{x} &= A^{-1}\delta A\mathbf{y} \implies \|\delta\mathbf{x}\| \leq \|A^{-1}\|_{op}\|\delta A\|_{op}\|\mathbf{y}\|.
\end{aligned}$$

This now produces

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{y}\|} \leq \|A\|_{op}\|A^{-1}\|_{op} \frac{\|\delta A\|_{op}}{\|A\|_{op}} = K_A \frac{\|\delta A\|_{op}}{\|A\|_{op}}$$

expressing the solution error relative to the disturbed solution in terms of the condition number of matrix  $A$  and the relative matrix error.

The problem here is that in this situation we don't *have*  $A$ , so we can't calculate  $K_A$ . We do have  $C = A + \delta A$ , and so  $\|A\| = \|C - \delta A\| \geq \|C\| - \|\delta A\|$ .

$$\begin{aligned}
\frac{\|\delta\mathbf{x}\|}{\|\mathbf{y}\|} &\leq \|A^{-1}\|_{op}\|\delta A\|_{op} = \|C\|_{op}\|C^{-1}\|_{op} \frac{\|A\|_{op}\|A^{-1}\|_{op}}{\|C\|_{op}\|C^{-1}\|_{op}} \frac{\|\delta A\|_{op}}{\|A\|_{op}} \\
&= K_C \frac{K_A}{K_C} \frac{\|\delta A\|_{op}}{\|A\|_{op}} \leq \left(\frac{K_A}{K_C}\right) K_C \frac{\|\delta A\|_{op}}{\|C\|_{op} - \|\delta A\|_{op}}
\end{aligned}$$

The final expression gives the the relative solution error in terms of quantities actually in our possession, except for the factor  $\frac{K_A}{K_C}$ .

If  $\delta A$  is small  $K_C = \frac{\sigma_1}{\sigma_n}$  might be close to  $K_A$ , where the  $\sigma_i$  are the singular values for  $C$ . However if  $\sigma_n$  is close to zero even a small  $\delta A$  could move the smallest singular value for  $C$  relatively closer to zero and blow up the factor  $\frac{K_A}{K_C}$ .

So if the least singular value  $\sigma_n$  for  $C$  is large compared to  $\delta A$  the estimate

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{y}\|} \leq K_C \left( \frac{\|\delta A\|_{op}}{\|C\|_{op} - \|\delta A\|_{op}} \right),$$

is likely to be useful. Otherwise you should be very cautious about using the approximate solution  $\mathbf{y}$ .

I should emphasize here that *this does not mean the solution is necessarily bad or useless*. This is worst-case scenario. It is a warning flag.

In the exercise below we investigate a typical situation where  $\frac{K_A}{K_C}$  is near one so we *can* use this uncertainty estimate.

**67.4. Exercise.** Suppose  $\|C\|_{op} = 1.5$  and  $K_C = 1000$  and the uncertainty in the entries of  $C$  produces maximum possible  $\|\delta A\|_{op} = 10^{-4}$ . The least singular value of  $C$  is large compared to  $10^{-4}$ .

You solve  $C\mathbf{y} = \mathbf{b}$  and find  $\|\mathbf{y}\| = 2.3$ .

The top entry  $y^1$  of solution vector  $\mathbf{y}$  is 1.3423192. How many of those decimal places should you trust?