

# Math 307: Problems for section 1.1

1. Use Gaussian elimination to find the solution(s) to  $Ax = \mathbf{b}$  where

$$(a) \quad A = \begin{bmatrix} 1 & 2 & 3 & 4 \\ -1 & 2 & -3 & 4 \\ 5 & 6 & 7 & 8 \\ -5 & 6 & -7 & 8 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad (b) \quad A = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 3 \\ 1 \\ 0 \\ 1 \end{bmatrix}.$$

The process of Gaussian elimination is not unique — there is more than one way to do it.

We create the augmented matrix  $C = [A|\mathbf{b}]$  and perform Gaussian elimination on  $C$

$$(a) \quad \begin{bmatrix} 1 & 2 & 3 & 4 & 1 \\ -1 & 2 & -3 & 4 & 1 \\ 5 & 6 & 7 & 8 & 1 \\ -5 & 6 & -7 & 8 & 1 \end{bmatrix} \mapsto \begin{bmatrix} 1 & 2 & 3 & 4 & 1 \\ 0 & 4 & 0 & 8 & 2 \\ 5 & 6 & 7 & 8 & 1 \\ -5 & 6 & -7 & 8 & 1 \end{bmatrix} \mapsto \begin{bmatrix} 1 & 2 & 3 & 4 & 1 \\ 0 & 4 & 0 & 8 & 2 \\ 5 & 6 & 7 & 8 & 1 \\ 0 & 12 & 0 & 16 & 2 \end{bmatrix} \mapsto \begin{bmatrix} 1 & 2 & 3 & 4 & 1 \\ 0 & 4 & 0 & 8 & 2 \\ 0 & -4 & -8 & -12 & -4 \\ 0 & 12 & 0 & 16 & 2 \end{bmatrix} \mapsto \begin{bmatrix} 1 & 2 & 3 & 4 & 1 \\ 0 & 4 & 0 & 8 & 2 \\ 0 & 0 & -8 & -4 & -2 \\ 0 & 12 & 0 & 16 & 2 \end{bmatrix} \mapsto \begin{bmatrix} 1 & 2 & 3 & 4 & 1 \\ 0 & 4 & 0 & 8 & 2 \\ 0 & 0 & -8 & -4 & -2 \\ 0 & 0 & 0 & -8 & -4 \end{bmatrix}$$

$$(b) \quad \begin{bmatrix} 1 & 1 & 1 & 1 & 3 \\ 1 & 1 & -1 & -1 & 1 \\ 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \end{bmatrix} \mapsto \begin{bmatrix} 1 & 1 & 1 & 1 & 3 \\ 0 & 0 & -2 & -2 & -2 \\ 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \end{bmatrix} \mapsto \begin{bmatrix} 1 & 1 & 1 & 1 & 3 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \end{bmatrix} \mapsto \begin{bmatrix} 1 & 1 & 1 & 1 & 3 \\ 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \end{bmatrix} \mapsto \begin{bmatrix} 1 & 1 & 1 & 1 & 3 \\ 0 & -2 & -1 & -1 & -3 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \end{bmatrix} \mapsto \begin{bmatrix} 1 & 1 & 1 & 1 & 3 \\ 0 & -2 & -1 & -1 & -3 \\ 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

The solution to the system will be the same regardless of how the Gaussian elimination is done:

$$(a) \quad \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 0 \\ -1/2 \\ 0 \\ 1/2 \end{bmatrix} \quad (b) \quad \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \end{bmatrix} + s \begin{bmatrix} 0 \\ 0 \\ -1 \\ 1 \end{bmatrix}$$

2. Use MATLAB/Octave to find the solution(s) to  $Ax = b$  where

$$(a) \quad A = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{bmatrix} \quad b = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad (b) \quad A = \begin{bmatrix} 1 & 0 & 3 & 2 & -4 \\ 2 & 1 & 6 & 5 & 0 \\ -1 & 1 & -3 & -1 & 1 \end{bmatrix}, \quad b = \begin{bmatrix} 4 \\ 7 \\ -5 \end{bmatrix}.$$

(a)

```
> A=[1 1 1 1; 1 1 -1 -1; 1 -1 0 0; 0 0 1 -1];
> b=[1 1 1 1]';
> C=[A b];
> rref(C)
ans =
1 0 0 0 1
0 1 0 0 0
0 0 1 0 0.5
0 0 0 1 -0.5
```

and so the answer is

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0.5 \\ -0.5 \end{bmatrix}$$

We could have instead done

```
> A=[1 1 1 1; 1 1 -1 -1; 1 -1 0 0; 0 0 1 -1];
> b=[1 1 1 1]';
> A\b
ans =
1
0
0.5
-0.5
```

and found the same answer — be careful though (see next part).

(b)

```
> A=[1 0 3 2 -4; 2 1 6 5 0; -1 1 -3 -1 1];
> b = [4 7 -5]';
> C=[A b];
> rref(C)
ans =
1 0 3 2 0 4
0 1 0 1 0 -1
0 0 0 0 1 0
```

The solution is

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} 4 \\ -1 \\ 0 \\ 0 \\ 0 \end{bmatrix} + s \begin{bmatrix} -2 \\ -1 \\ 0 \\ 1 \\ 0 \end{bmatrix} + t \begin{bmatrix} -3 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

If we had instead used

```
> A\b
```

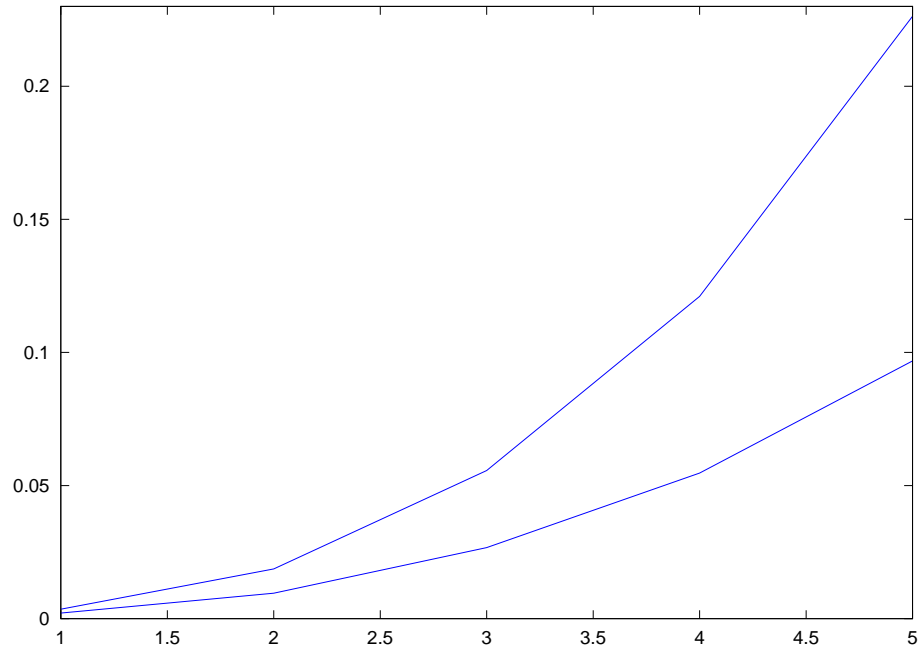
then MATLAB/Octave would have given only one of the many solutions.

3. Compute the time it takes to solve  $Ax = b$  using  $A\b$  on your computer for five or more random problems of increasing size. Make a plot of time vs. size. Repeat using the method  $A^{-1}*b$  and plot on the same graph. (If your calculation is taking too long, it can be interrupted by typing `<ctrl> c.`)

```
> A=rand(100,100); b=rand(100,1);
> tic(); A\b; toc();
Elapsed time is 0.00211289 seconds.
> tic(); A^(-1)*b; toc();
Elapsed time is 0.00353203 seconds.
> A=rand(200,200); b=rand(200,1);
> tic(); A\b; toc();
Elapsed time is 0.00957103 seconds.
tic(); A^(-1)*b; toc();
Elapsed time is 0.0186941 seconds.
> A=rand(300,300); b=rand(300,1);
> tic(); A\b; toc();
Elapsed time is 0.0267109 seconds.
tic(); A^(-1)*b; toc();
Elapsed time is 0.0556671 seconds.
> A=rand(400,400); b=rand(400,1);
> tic(); A\b; toc();
Elapsed time is 0.054698 seconds.
> tic(); A^(-1)*b; toc();
Elapsed time is 0.12107 seconds.
> A=rand(500,500); b=rand(500,1);
> tic(); A\b; toc();
Elapsed time is 0.0968079 seconds.
> tic(); A^(-1)*b; toc();
Elapsed time is 0.226295 seconds.

X=[0.00211289,0.00957103,0.0267109,0.054698,0.0968079];
Y=[0.00353203,0.0186941,0.0556671,0.12107,0.226295];
plot(X)
hold on
plot(Y)
print -depsc hmk1.1plot.eps
```

This last command produces the following eps file.



4. Compute the 1,2 and infinity norms for the following vectors. Is it always the same norm that is biggest? (The 2-norm is the standard Euclidean norm)

$$\mathbf{a} = \begin{bmatrix} 1 \\ 2 \\ 0 \\ 1 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 1 \\ -1 \\ 0 \\ 1 \end{bmatrix} \quad \mathbf{c} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$\|\mathbf{a}\|_1 = 4, \|\mathbf{a}\|_2 = \sqrt{6} \sim 2.4495, \|\mathbf{a}\|_\infty = 2$$

$$\|\mathbf{b}\|_1 = 3, \|\mathbf{b}\|_2 = \sqrt{3} \sim 1.7321, \|\mathbf{b}\|_\infty = 1$$

$$\|\mathbf{c}\|_1 = 3, \|\mathbf{c}\|_2 = \sqrt{3} \sim 1.7321, \|\mathbf{c}\|_\infty = 1$$

The 1-norm is always the biggest and the infinity norm the smallest

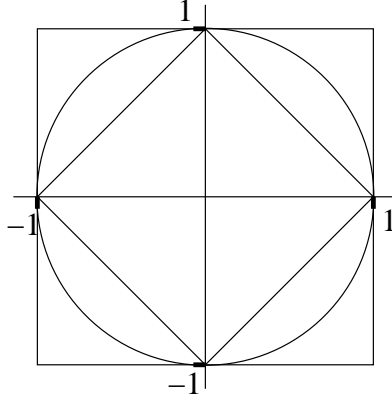
5. Find four vectors in two dimensions whose 1, 2 and infinity norms the same.

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

6. Suppose you were doing a problem each entry  $v_i$  in a vector  $[v_1, v_2, \dots, v_n]$  is positive and represents the yearly production of one of  $n$  factories. Which of the three norms we introduced have natural interpretations in this context.

The 1-norm is the total production of all the factories, and the infinity norm is the production of the most productive factory.

7. Draw a picture of the “unit circle” for the 1,2, and infinity norms in two dimensions. By “unit circle” we mean the set of all vectors whose norm is equal to one.



The unit circle for the 1-norm is the diamond, the unit circle for the 2-norm is the circle, the unit circle for the infinity norm is the square.

8. Recall that the Euclidean distance between two vectors  $\mathbf{v}$  and  $\mathbf{w}$  is  $\|\mathbf{v} - \mathbf{w}\|$  where we use the standard norm. If we use the 1-norm or  $\infty$ -norm in this formula, we obtain different distance functions. Consider vectors whose entries are either 0 or 1 (like  $[0, 1, 1, 0, 1]$ ). Describe in words the meaning of the 1-distance and the  $\infty$ -distance between two such vectors.

The 1-distance is the number of positions where the vectors have different entries, while the  $\infty$ -distance is either 1, if the vectors are different, or 0, if the vectors are the same.

9. The  $p$ -norm of a vector  $\mathbf{v} = [v_1, v_2, \dots, v_n]^T$  for  $1 \leq p \leq \infty$  is defined to be

$$\|\mathbf{v}\|_p = \left( \sum_{i=1}^n |v_i|^p \right)^{1/p}$$

Guess the MATLAB/Octave syntax to compute this norm and check that you are right. Which  $p$  corresponds to the standard (Euclidean) norm? What is the limit of the  $p$ -norm of a vector as  $p$  tends to infinity?

The MATLAB/Octave syntax for the  $p$  norm of the vector  $\mathbf{X}$  is `norm(X,p)`. The standard Euclidean norm corresponds to  $p = 2$ . The limit of the  $p$ -norm of a vector as  $p$  tends to infinity is the  $\infty$ -norm.

To see this suppose the largest component of the vector  $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$  is  $x_i$  so that  $\|\mathbf{x}\|_\infty = |x_i|$ . Then

$\|\mathbf{x}\|_p = \left( \sum_{j=1}^n |x_j|^p \right)^{1/p} = \left( |x_j|^p \sum_{j=1}^n \lambda_j^p \right)^{1/p} = |x_j| \left( \sum_{j=1}^n \lambda_j^p \right)^{1/p}$  where  $0 \leq \lambda_j = |x_j|/|x_i| \leq 1$  and  $\lambda_i = 1$ . We will show that  $\lim_{p \rightarrow \infty} \log \|\mathbf{x}\|_p = \log \|\mathbf{x}\|_\infty$ . From the formula above

$$\log \|\mathbf{x}\|_p = \log \|\mathbf{x}\|_\infty + \frac{1}{p} \log(\lambda_1 + \dots + \lambda_n)$$

Since  $1 \leq \lambda_1 + \dots + \lambda_n \leq n$  we have  $0 \leq \log(\lambda_1 + \dots + \lambda_n) \leq \log n$  and the second term on the right tends to zero as  $p$  tends to infinity.

10. Show that for any square matrix  $A$  (with real entries),  $\|A\|_{HS}^2 = \text{tr}(A^T A)$ .

If  $A = [A_{ij}]$  then  $A^T = [A_{ij}^T] = [A_{ji}]$ . Here  $i$  and  $j$  run from 1 to  $n$ . The formula for matrix multiplication gives the  $i, j$  entry of  $A^T A$  as

$$(A^T A)_{ij} = \sum_{k=1}^n A_{ik}^T A_{kj} = \sum_{k=1}^n A_{ki} A_{kj}$$

The trace of a matrix is the sum of the diagonal elements. These are the elements with indices  $i, i$ . Thus

$$\text{tr} A^T A = \sum_{i=1}^n (A^T A)_{ii} = \sum_{i=1}^n \sum_{k=1}^n A_{ki} A_{ki} = \sum_{i=1}^n \sum_{k=1}^n A_{ki}^2 = \|A\|_{HS}^2$$

11. **Guess whether each of the following statements about  $n \times n$  matrices  $A$  is true or false by testing them on a few random matrices. Bonus: prove that your guess is correct.**

(a)  $\|A^2\| = \|A\|^2$

False: for example

```
>A=[1 2; 0 4];
> norm(A^2)
ans = 18.875
> norm(A)^2
ans = 20.208
```

(b)  $\|A^2\| \leq \|A\|^2$

True: To prove it note that for any  $\mathbf{x} \neq 0$ ,  $\|A^2 \mathbf{x}\| = \|\mathbf{A}(\mathbf{A}\mathbf{x})\| \leq \|\mathbf{A}\| \|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{A}\| \|\mathbf{x}\| = \|\mathbf{A}\|^2 \|\mathbf{x}\|$ . This implies that  $\|A^2 \mathbf{x}\| / \|\mathbf{x}\| \leq \|\mathbf{A}\|^2$ . Since this is true for every  $\mathbf{x}$  is also must be true for the maximum. This shows that  $\|A^2\| \leq \|A\|^2$ .

(c)  $\|A^T A\| = \|A\|^2$

True: but proving it is a bit tricky (and may use material that you are not yet familiar with). Here are two different ways to do it.

$$\begin{aligned} \|A\|^2 &= \left( \max_{\mathbf{x}: \mathbf{x} \neq \mathbf{0}} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|} \right)^2 = \max_{\mathbf{x}: \mathbf{x} \neq \mathbf{0}} \frac{\|A\mathbf{x}\|^2}{\|\mathbf{x}\|^2} \\ &= \max_{\mathbf{x}: \mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^T A^T A \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \quad \text{using the fact that } \|\mathbf{v}\|^2 = \mathbf{v}^T \mathbf{v} \end{aligned}$$

But  $A^T A$  is a real symmetric positive definite matrix and therefore has an orthonormal basis of eigenvectors  $\mathbf{e}_1, \dots, \mathbf{e}_n$  and positive eigenvalues  $\lambda_1, \dots, \lambda_n$ . Writing  $\mathbf{x} = \alpha_1 \mathbf{e}_1 + \dots \alpha_n \mathbf{e}_n$  we have

$$\|A\|^2 = \max_{\alpha: \alpha \neq \mathbf{0}} \frac{\sum_{i=1}^n \alpha_i^2 \lambda_i}{\sum_{i=1}^n \alpha_i^2} = \max_i \lambda_i$$

We also have that

$$\|A^T A\| = \max_{\alpha: \alpha \neq \mathbf{0}} \frac{\sum_{i=1}^n \alpha_i \lambda_i}{\sum_{i=1}^n \alpha_i} = \max_i \lambda_i.$$

Comparing the two results we have our proof.

An alternative proof. Let's write the dot product as  $\mathbf{x}^T \mathbf{y}$  (thinking of vectors as  $n \times 1$  matrices) and recall that  $\mathbf{y}^T \mathbf{x} = \|\mathbf{y}\| \|\mathbf{x}\| \cos(\theta)$  where  $\theta$  is the angle between the vectors. If we let  $\mathbf{y}$  run

through all unit vectors, then the left side is biggest when  $\theta = 0$  and  $\cos(\theta) = 1$ . This gives a formula for  $\|\mathbf{x}\|$ , namely

$$\|\mathbf{x}\| = \max_{\mathbf{y}: \|\mathbf{y}\|=1} \mathbf{y}^T \mathbf{x}$$

Now we find

$$\begin{aligned} \|A^T A \mathbf{x}\| &= \max_{\mathbf{y}: \|\mathbf{y}\|=1} \mathbf{y}^T A^T A \mathbf{x} \\ &= \max_{\mathbf{y}: \|\mathbf{y}\|=1} (\mathbf{A}\mathbf{y})^T A \mathbf{x} \\ &\leq \max_{\mathbf{y}: \|\mathbf{y}\|=1} \|\mathbf{A}\mathbf{y}\| \|A \mathbf{x}\| \\ &\leq \max_{\mathbf{y}: \|\mathbf{y}\|=1} \|A\| \|\mathbf{y}\| \|A\| \|\mathbf{x}\| \\ &= \|A\|^2 \|\mathbf{x}\| \end{aligned}$$

Since this is true for every  $\mathbf{x}$  it implies  $\|A^T A\| \leq \|A\|^2$ . To get the opposite inequality, start the same way, but this time we use that the maximum over all unit vectors  $\mathbf{y}$  is bigger than the value for the particular unit vector  $\mathbf{x}/\|\mathbf{x}\|$ .

$$\begin{aligned} \|A^T A \mathbf{x}\| &= \max_{\mathbf{y}: \|\mathbf{y}\|=1} \mathbf{y}^T A^T A \mathbf{x} \\ &= \max_{\mathbf{y}: \|\mathbf{y}\|=1} (\mathbf{A}\mathbf{y})^T A \mathbf{x} \\ &\geq \frac{(\mathbf{A}\mathbf{x})^T A \mathbf{x}}{\|\mathbf{x}\|} \\ &= \frac{\|A \mathbf{x}\|^2}{\|\mathbf{x}\|} \end{aligned}$$

Dividing by  $\|\mathbf{x}\|$  this gives

$$\frac{\|A^T A \mathbf{x}\|}{\|\mathbf{x}\|} \geq \frac{\|A \mathbf{x}\|^2}{\|\mathbf{x}\|^2}$$

Taking the maximum over  $\mathbf{x}$  and using that the maximum of the squares of positive number is the square of the maximum gives

$$\|A^T A\| \geq \|A\|^2$$

(d)  $\|A\| \leq \|A\|_{HS}$

True: again two possible ways to do it.

We need to show that  $\|A\|^2 = \|A^T A\| \leq \|A\|_{HS}^2 = \text{tr}(A^T A)$  (using results from previous questions). As in the proof for (c), we use the fact that  $A^T A$  is a real symmetric positive definite matrix with positive eigenvalues so that  $\|A^T A\| = \max_i \lambda_i$ . But from general properties of matrices we have that the trace of a matrix is equal to the sum of its eigenvalues  $\text{tr}(A^T A) = \sum_{i=1}^n \lambda_i$ . Since the eigenvalues for a real symmetric positive definite matrix are all positive, comparing these two results gives us our proof.

An alternative proof. Note that if  $A = [A_{ij}]$  then  $(A\mathbf{x})_i = \sum_j A_{ij} x_j$ . Using the Cauchy-Schwartz inequality this gives

$$(A\mathbf{x})_i^2 \leq \sum_j A_{ij}^2 \|\mathbf{x}\|^2.$$

Now summing over  $i$  gives

$$\|A\mathbf{x}\|^2 \leq \sum_i \sum_j A_{ij}^2 \|\mathbf{x}\|^2 = \|A\|_{HS}^2 \|\mathbf{x}\|^2$$

which proves the  $\|A\|^2 \leq \|A\|_{HS}^2$ .

(e)  $\text{cond}(A) = \text{cond}(A^{-1})$

True: easy to prove since

$$\text{cond}(A^{-1}) = \|A^{-1}\| \|(A^{-1})^{-1}\| = \|A^{-1}\| \|A\| = \text{cond}(A)$$

(f)  $\text{cond}(A) \geq 1$

True: To prove it notice that we showed

$$\min_{\mathbf{x} \neq 0} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|} = \frac{1}{\|A^{-1}\|}$$

But

$$\|A\| = \max_{\mathbf{x} \neq 0} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|}$$

is bigger than the left side of the previous equality (since the max is bigger than the min). Hence

$$\|A\| \geq \frac{1}{\|A^{-1}\|}$$

which implies that

$$\text{cond}(A) = \|A^{-1}\| \|A\| \geq \frac{\|A^{-1}\|}{\|A^{-1}\|} = 1$$

12. For a diagonal matrix, the eigenvalues are equal to the diagonal entries. So the matrix norm is the absolute value of the largest eigenvalue. Show that this is not true for an arbitrary matrix. (In MATLAB/Octave `eig(A)` computes the eigenvalues of  $A$ ) For a matrix  $A$  with real entries, there is a relationship between the norm of  $A$  and the eigenvalues of  $A^T A$ . Can you guess what it is using MATLAB/Octave?

Consider the matrix  $\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ .

```
> A=[1 1; 0 1];
> eig(A)
ans =
```

```
    1
    1
> norm(A)
ans = 1.6180
```

It turns out that the norm is the square root of the largest eigenvalue of  $A^T A$ .

```
> eig(A'*A)
ans =

    0.38197
    2.61803
> sqrt(2.61803)
ans = 1.6180
```



13. Suppose  $A$  has a large condition number. This means that in the equation  $Ax = b$ , a small relative error in  $b$  may result in a large relative error in  $x$ . Is it possible, though, that for some choices of  $b$  and  $\Delta b$  the relative error of  $x$  is not large. Illustrate this using the matrix  $A = \begin{bmatrix} 1000 & 0 \\ 0 & 1 \end{bmatrix}$ .

If we take  $b = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$  and  $\Delta b = \begin{bmatrix} 0 \\ 0.1 \end{bmatrix}$  then we find that  $x = b$  and  $\Delta x = \Delta b$  so the relative error in  $b$  and  $x$  are the same, even though the matrix has condition number 1000.

14. A famous ill-conditioned matrix is the hilbert matrix whose  $i, j$  entry is  $1/(i+j)$ . In MATLAB/Octave the hilbert matrix of size  $n$  can be generated using `hilb(n)`. What is the condition number when  $n=5$ ,  $n=10$ ?

```
> cond(hilb(5))
ans = 4.7661e+05
> cond(hilb(10))
ans = 1.6025e+13
```

15. In “single precision” computer calculations, we cannot trust more than approximately the first 7 significant figures. Assuming that the relative error in the right-hand-side of the matrix equation  $Ax = b$  is  $\|\Delta b\|/\|b\| = 1.1921 \times 10^{-7}$ , give an upper bound on the relative error  $\|\Delta x\|/\|x\|$  in the solution of the equation for the following matrices  $A$ :

$$(a) \quad A = \begin{bmatrix} 0 & 2 & 4 \\ 3 & 1 & 1 \\ 2 & 0 & 1 \end{bmatrix} \quad (b) \quad A = \begin{bmatrix} 4 & 1.99 \\ 2.01 & 1 \end{bmatrix}$$

Interpret these bounds to say how many significant figures we can trust in the solution.

Remember from lectures that we found

$$\frac{\|\Delta x\|}{\|x\|} \leq \text{cond}(A) \frac{\|\Delta b\|}{\|b\|} = 1.1921 \times 10^{-7} \text{cond}(A).$$

For (a) we have from MATLAB/Octave that  $\text{cond}(A) \approx 8.1226$ . Therefore

$$\|\Delta x\|/\|x\| \leq 8.1226 \times 1.1921 \times 10^{-7} \approx 9.68 \times 10^{-7} \approx 10^{-6}.$$

So the error in  $\|x\|$  is roughly  $10^{-6}\|x\|$ . The first 6 digits can be trusted.

For (b) we have from MATLAB/Octave that  $\text{cond}(A) \approx 2.5 \times 10^5$ . Therefore

$$\|\Delta x\|/\|x\| \leq 2.5 \times 10^5 \times 1.1921 \times 10^{-7} \approx 0.029803.$$

1 digit can be trusted (the second digit may be off by 1).