

# Math 170A Lecture Notes (Professor: Lei Huang)

Isabelle Mills

February 17, 2024

## Week 1 Notes (1/8 - 1/12/2024)

For this class, we shall define the number of flops an algorithm takes as the number of individual  $+$ ,  $-$ ,  $\times$ ,  $/$ , and  $\sqrt{\quad}$  operations on real numbers used in the algorithm.

For example: taking the inner product of two vectors  $\vec{u} = (u_1, u_2, \dots, u_n)$  and  $\vec{v} = (v_1, v_2, \dots, v_n)$  requires  $n$  multiplications and  $(n - 1)$  additions. So we'll say it has a flops count of  $2n - 1$ .

Technically, the word "flop" stands for floating (point) operation. Based on that knowledge, hopefully it is easier to guess what is and is not a flop. For instance, observe the code written below for taking an inner product of two  $n$ -vectors.

```
1  P = 0
2  for i = 1:n
3      P = P + v(i) * w(i)
4  end
```

Neither incrementing `i` nor initializing any other variables are counted towards the flop number. Because the code does  $n$  additions and multiplications between floating point numbers, we say this function has  $2n$  flops.

Here is how we formally define Big- $\mathcal{O}$  Notation:

For a sequence  $a_n$ , we define  $a_n = \mathcal{O}(b_n)$  if there exists real constants  $C, N \geq 0$  such that for  $n \geq N$ ,  $a_n \leq Cb_n$ .

Example problem:

```
1  function x = lowertriangsolve(L, b)
2      N = size(L);
3      n=size(b,1);
4      x=b;
5
6      for i=1:N(1):
7          for j=1:i-1
8              x(i) = x(i) - L(i,j)*x(j);
9          end
10
11         if L(i, j) == 0
12             error('matrix is singular')
13         end
14         x(i) = x(i)/L(i,i);
15     end
16 end
```

Inside the main for loop (lines 6-15):  
Line 8 has  $2(i - 1)$  flops.  
Line 14 has an additional flop.

Thus, the total number of flops is:

$$\sum_{i=1}^n (2i - 1) = 2 \left( \sum_{i=1}^n i \right) - n$$

This then simplifies to:

$$2 \left( \frac{n(n+1)}{2} \right) - n = 1n^2$$

So  $x$  has  $\mathcal{O}(n^2)$  flops as  $n$  goes to infinity.

## Lecture 4: 1/17/2024

Row operations used in Gaussian elimination can be represented via matrix multiplication as demonstrated below.

Action	Matrix representation of the operation
$\begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} \\ a_{2,1} & a_{2,2} & a_{2,3} & a_{2,4} \\ a_{3,1} & a_{3,2} & a_{3,3} & a_{3,4} \\ a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} \end{bmatrix} \xrightarrow{R_2 \leftrightarrow R_3} \begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} \\ a_{3,1} & a_{3,2} & a_{3,3} & a_{3,4} \\ a_{2,1} & a_{2,2} & a_{2,3} & a_{2,4} \\ a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$
$\begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} \\ a_{2,1} & a_{2,2} & a_{2,3} & a_{2,4} \\ a_{3,1} & a_{3,2} & a_{3,3} & a_{3,4} \\ a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} \end{bmatrix} \xrightarrow{\frac{2}{3}R_2 \rightarrow R_2} \begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} \\ \frac{2}{3}a_{2,1} & \frac{2}{3}a_{2,2} & \frac{2}{3}a_{2,3} & \frac{2}{3}a_{2,4} \\ a_{3,1} & a_{3,2} & a_{3,3} & a_{3,4} \\ a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & \frac{2}{3} & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$
$\begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} \\ a_{2,1} & a_{2,2} & a_{2,3} & a_{2,4} \\ a_{3,1} & a_{3,2} & a_{3,3} & a_{3,4} \\ a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} \end{bmatrix} \xrightarrow{R_3 + \frac{2}{3}R_2 \rightarrow R_3} \begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} \\ a_{2,1} & a_{2,2} & a_{2,3} & a_{2,4} \\ a_{3,1} + \frac{2}{3}a_{2,1} & a_{3,2} + \frac{2}{3}a_{2,2} & a_{3,3} + \frac{2}{3}a_{2,3} & a_{3,4} + \frac{2}{3}a_{2,4} \\ a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & \frac{2}{3} & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$

The importance of representing elementary row operations as matrices is that we can multiply these representations together to compose rows operations. Thus, these representations are central to many matrix decompositions.

For example, we can represent turning a matrix into row echelon form as follows:

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -\frac{2}{9} & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -2 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ -3 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 4 & 5 \\ 6 & 3 & 4 \\ 4 & 6 & 2 \end{bmatrix} = \begin{bmatrix} 2 & 4 & 5 \\ 0 & -9 & -11 \\ 0 & 0 & \frac{-50}{9} \end{bmatrix}$$

Here are some more observations about row operation matrices:

- Row scaling operations are represented by diagonal matrices and thus are also both lower and upper triangular.
- For  $i < j$ , adding a multiple of the  $i$ th row to the  $j$ th row is represented by a lower triangular matrix.
- For  $i > j$ , adding a multiple of the  $i$ th row to the  $j$ th row is represented by an upper triangular matrix.
- Row swaps are not represented as triangular matrices. However, they are permutation matrices.

Now note that in the normal algorithm for Gaussian elimination, assuming we never need to swap rows, all row operations will be such that their matrix representation is lower triangular.

Thus, given an invertible square matrix  $\mathbf{A}$ , we can represent doing Gaussian elimination on  $\mathbf{A}$  by the equation:  $\mathbf{L}_n \mathbf{L}_{n-1} \cdots \mathbf{L}_2 \mathbf{L}_1 \mathbf{A} = \mathbf{U}$  where  $\mathbf{L}_i$  is a lower triangular matrix and  $\mathbf{U}$  is an upper triangular matrix. Then, because the product of two lower triangular matrices is also lower triangular, we can multiply all the lower triangular matrices together to get an equation of the form  $\mathbf{L}\mathbf{A} = \mathbf{U}$ . Finally, we multiply both sides of the equation on the left by  $\mathbf{L}^{-1}$  to get that  $\mathbf{A} = \mathbf{L}^{-1}\mathbf{U}$ . And as the inverse of a lower triangular matrix is also lower triangular, we know that we have decomposed  $\mathbf{A}$  into the product of a lower triangular matrix and an upper triangular matrix. This algorithm is called LU decomposition.

As for why we would want to use LU decomposition, we can look at the number of flops different matrix operations require.

Assume we are given an invertible  $n \times n$  matrix  $\mathbf{A}$  and two vectors  $\vec{x}, \vec{b} \in \mathbb{R}^n$ .

Firstly note that row reduction takes approximately  $\frac{2}{3}n^3$  flops while the back substitution algorithm takes approximately  $n^2$  flops. Thus, for larger values of  $n$ , solving the matrix equation  $\mathbf{A}\vec{x} = \vec{b}$  takes around  $\frac{2}{3}n^3$  flops.

Now let's compare this to the number of flops the best algorithm we can make to do LU decomposition takes.

Assuming we don't need to do any row swaps, then the only elementary row operations we need to do are adding scaled rows to other rows. This is where our first optimization comes into play. The inverse of a row addition is merely a row subtraction. For example:

$$\begin{bmatrix} 1 & 0 & 0 \\ a & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ -a & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Thus, we can fairly straightforwardly get an expression of the form  $\mathbf{A} = \mathbf{L}_1 \mathbf{L}_2 \cdots \mathbf{L}_k \mathbf{U}$  by multiplying the inverse of each elementary row operation matrix to both sides of the equation representing the actions we took to reduce  $\mathbf{A}$  to  $\mathbf{U}$ .

Now, we apply another two shortcuts to speed up our calculations.

Firstly, observe that: 
$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & a & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & b & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & b & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & a & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & a & 1 & 0 \\ 0 & b & 0 & 1 \end{bmatrix}$$

There's nothing special about that column or that particular matrix size. In fact, this should make sense when you consider that those two elementary row operations don't effect each other.

Next, observe that: 
$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ a & 1 & 0 & 0 \\ b & 0 & 1 & 0 \\ c & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & d & 1 & 0 \\ 0 & e & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & f & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ a & 1 & 0 & 0 \\ b & d & 1 & 0 \\ c & e & f & 1 \end{bmatrix}$$

To intuit why this works, think about how the matrix product has every row affect the rows underneath it before it itself is affected by any rows above it. So, it is specifically the initial state of every row that is effecting every other row in the matrix product.

The end result of these observations is that the  $(i, j)$ th element of  $\mathbf{L}$  where  $i > j$  is just going to be the negative of whatever coefficient one multiplied a copy of row  $i$  by before then adding that to row  $j$  as part of doing the row reduction algorithm. But now note that means that every element in  $\mathbf{L}$  can be directly extracted from the calculations done to find  $\mathbf{U}$ . So finding  $\mathbf{U}$  takes approximately  $\frac{2}{3}n^3$  flops and finding  $\mathbf{L}$  takes no additional flops.

Once, we've found  $\mathbf{LU}$ , we can then use back substitution twice to solve any matrix vector equation  $\mathbf{A}\vec{x} = \mathbf{LU}\vec{x} = \vec{b}$ . This will take  $2n^2$  flops.

Technically, this means that if you are trying to solve a single matrix vector equation  $\mathbf{A}\vec{x} = \vec{b}$  by first decomposing  $\mathbf{A}$ , then it will take longer than if you just solved it directly. However, if you have many matrix vector equations involving  $\mathbf{A}$ , then you can work way faster by decomposing  $\mathbf{A}$ . This is because once you decompose  $\mathbf{A}$  once, you can just store  $\mathbf{L}$  and  $\mathbf{U}$  for later use. Thus, every subsequent matrix vector equation involving  $\mathbf{A}$  can be done in approximately  $2n^2$  flops instead of taking approximately  $\frac{2}{3}n^3$  flops.

## Lecture 5: 1/19/2024

A crucial assumption we made in the last lecture was that we never would have to do row swaps when doing row reduction. In this lecture, we'll now allow ourselves to do row swaps.

Firstly, here are two general reasons to want to do row swaps.

- Firstly, sometimes we have no choice.

$$\begin{bmatrix} 0 & 4 & 1 \\ 1 & 3 & 4 \\ 2 & 2 & 5 \end{bmatrix}$$

For this matrix, if we were to divide the second and third row by  $a_{1,1}$ , we'd be dividing them by 0. Thus, we clearly can't do that.

- Secondly, dividing by small numbers causes more roundoff errors. So, we can slow the accumulation of roundoff errors by swapping rows in order to prioritize dividing by larger elements.

So here's an algorithm called partial pivoting for taking a PLU decomposition of an invertible  $n \times n$  matrix  $\mathbf{A}$ .

When doing Gaussian elimination, for every new pivot of  $\mathbf{A}$ , first perform a row swap with the top non-reduced row and the non-reduced row with the largest would-be pivot element. Only, after that do you then add a scaled version of the pivot row to the other rows like you were doing before.

Now note that performing the same row swap twice in a row is equivalent to doing nothing. Thus, every matrix representing a row swap is its own inverse.

Additionally, as we already covered, the inverse of a row addition operation is just a row subtraction operation. Because of this, we can easily get an equation for  $\mathbf{A}$  as the product of many elementary triangular matrices and permutation matrices.

Let's denote  $\mathbf{L}_i$  to be the lower triangular matrix representing all row additions of row  $i$  to the rest of the matrix. In other words, every element below the main diagonal in  $\mathbf{L}_i$  will be nonzero except for the elements in column  $i$ . Additionally, let's denote  $\mathbf{P}_{i,j}$  to be the permutation matrix swapping row  $i$  and row  $j$ . Then, we can say that:

$$\mathbf{A} = \mathbf{P}_{1,k_1} \mathbf{L}_1 \mathbf{P}_{2,k_2} \mathbf{L}_2 \cdots \mathbf{P}_{n-1,k_{n-1}} \mathbf{L}_{n-1} \mathbf{U}$$

(Note that  $i < k_i \leq n$  since it doesn't make sense to swap a row that has already been dealt with.)

We can rewrite this as  $\mathbf{P}_{1,k_1} \mathbf{A} = \mathbf{L}_1 \mathbf{P}_{2,k_2} \mathbf{L}_2 \cdots \mathbf{P}_{n-1,k_{n-1}} \mathbf{L}_{n-1} \mathbf{U}$ . And now here is where the magic starts. If we multiply both sides of that equation on the left by  $\mathbf{P}_{2,k_2}'$ , then we can say that  $(\mathbf{P}_{2,k_2}' \mathbf{L}_1 \mathbf{P}_{2,k_2}) = \mathbf{L}_1'$  where  $\mathbf{L}_1'$  is the matrix which would have resulted if we had only applied the permutation  $\mathbf{P}_{2,k_2}$  to the elements off the main diagonal.

Now we can do the same process again to  $(\mathbf{L}'_1 \mathbf{L}_2) \mathbf{P}_{3,k_3}$ , multiplying both sides of our equation by  $\mathbf{P}_{3,k_3}$  and then saying that  $\mathbf{L}'_2 = (\mathbf{P}_{3,k_3} \mathbf{L}'_1 \mathbf{L}_2 \mathbf{P}_{3,k_3})$  is a lower triangular matrix whose only nonzero elements below the diagonal are in columns 1 and 2. Doing this for all remaining permutation matrices on the right side of the equation, we will eventually get an equation of the form:

$$\mathbf{P}_{n-1,k_{n-1}} \cdots \mathbf{P}_{2,k_2} \mathbf{P}_{1,k_1} \mathbf{A} = \mathbf{L} \mathbf{U}$$

Finally, by multiplying those permutation matrices together, we get an equation:

$$\mathbf{P} \mathbf{A} = \mathbf{L} \mathbf{U}$$

Here is why the product  $\mathbf{P}_{i+1,k_{i+1}} \mathbf{L}'_{i-1} \mathbf{L}_i \mathbf{P}_{i+1,k_{i+1}}$  was equivalent to only applying the row permutation underneath the diagonal of  $\mathbf{L}'_{i-1} \mathbf{L}_i$ .

Observe the following diagram of where the permutation matrices send the elements in row  $i+1$ , row  $k_{i+1}$ , column  $i+1$ , and column  $k_{i+1}$  of the matrix  $\mathbf{L}'_{i-1} \mathbf{L}_i$ :

$$\begin{array}{c} \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ a & b & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 \\ c & d & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \\ \downarrow \\ \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ c & d & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ a & b & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \end{array}$$

Theorem: Every invertible matrix has a PLU-decomposition.

## Lecture 6: 1/22/2024

A matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is said to be positive definite if:

1.  $\mathbf{A}^T = \mathbf{A}$  ( $\mathbf{A}$  is symmetric)
2.  $\vec{x}^T \mathbf{A} \vec{x} = \langle \vec{x}, \mathbf{A} \vec{x} \rangle > 0$  for all vectors  $\vec{x} \in \mathbb{R}^n$  with  $\vec{x} \neq 0$ .

$$\langle \vec{x}, \mathbf{A} \vec{x} \rangle = \sum_{i=1}^n \sum_{j=1}^n x_i a_{i,j} x_j$$

Lemma: If  $\mathbf{A}$  is positive definite, then  $\mathbf{A}$  is invertible.

Proof: (we proceed towards a contradiction)

Assume there exists a vector  $\vec{y} \in \mathbb{R}^n$  not equal to 0 such that  $\mathbf{A} \vec{y} = 0$ . Then  $\langle \vec{y}, \mathbf{A} \vec{x} \rangle = \langle \vec{y}, 0 \rangle = 0$ . So,  $\mathbf{A}$  cannot be positive definite.

Theorem: Let  $\mathbf{M} \in \mathbb{R}^{n \times n}$  be invertible. Then  $\mathbf{A} = \mathbf{M}^T \mathbf{M}$  is positive definite.

Proof:

1.  $\mathbf{A}^T = (\mathbf{M}^T \mathbf{M})^T = \mathbf{M}^T (\mathbf{M}^T)^T = \mathbf{M}^T \mathbf{M} = \mathbf{A}$ . So  $\mathbf{A}$  is symmetric.
2. Note that  $\vec{x}^T \mathbf{A} \vec{x} = \vec{x}^T \mathbf{M}^T \mathbf{M} \vec{x} = (\mathbf{M} \vec{x})^T \mathbf{M} \vec{x} = \langle \mathbf{M} \vec{x}, \mathbf{M} \vec{x} \rangle = |\mathbf{M} \vec{x}|$ .  
Now,  $|\mathbf{M} \vec{x}| > 0$  for all  $\mathbf{M} \vec{x} \neq 0$ . And because  $\mathbf{M}$  is invertible, we know the only  $\vec{x}$  such that  $\mathbf{M} \vec{x} = 0$  is  $\vec{x} = 0$ . So,  $\vec{x}^T \mathbf{A} \vec{x} > 0$  for all  $\vec{x} \neq 0$ .  
If  $\mathbf{M}$  is not invertible, then  $\mathbf{A}$  is positive semidefinite as  $|\mathbf{M} \vec{x}|$  cannot be less than 0 but we can find a nonzero  $\vec{x}$  such that  $|\mathbf{M} \vec{x}| = 0$ .

Theorem (Cholesky decomposition):

Let  $\mathbf{A}$  be positive definite. Then there exists an upper triangular matrix  $\mathbf{R}$  such that  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$ . We call  $\mathbf{R}$  the Cholesky factor and can calculate it as follows:

If such an  $\mathbf{R}$  exists, then we can arrive at the following formula for each element of  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$ :

$$a_{i,j} = \sum_{k=1}^{\min(i,j)} r_{k,i} r_{k,j}$$

Now as  $\mathbf{A}$  must be symmetric based on our formula, we can safely ignore the elements of  $\mathbf{A}$  below the main diagonal. So, assuming that  $j \geq i$ , we can rearrange terms to get the following equation including the element  $a_{i,j}$ :

$$r_{i,i} r_{i,j} = a_{i,j} - \sum_{k=1}^{i-1} r_{k,i} r_{k,j}$$



And with that, we now have a way of expressing any element  $r_{i,j}$  in terms of elements of  $\mathbf{R}$  which are in previous rows or previous columns of  $\mathbf{R}$ . So, we can inductively solve for the coefficients of  $\mathbf{R}$  as follows:

for  $i = 1, \dots, n$ :

$$r_{i,i} = \sqrt{a_{i,i} - \sum_{k=1}^{i-1} r_{k,i}^2}$$

for  $j = (i + 1), \dots, n$ :

$$r_{i,j} = \frac{1}{r_{i,i}} \left( a_{i,j} - \sum_{k=1}^{i-1} r_{k,i} r_{k,j} \right)$$

end

end

Now it's worth noting what can go wrong in the above algorithm.

Firstly, if  $a_{i,i} - \sum_{k=1}^{i-1} r_{k,i}^2 < 0$ , then  $r_{i,i} \notin \mathbb{R}$ . So,  $\mathbf{R}$  does not exist

Secondly, if  $r_{i,i} = 0$ , then you can't divide by it to isolate  $r_{i,j}$ . Thus, one of two possibilities will arise:

- If  $\left( a_{i,j} - \sum_{k=1}^{i-1} r_{k,i} r_{k,j} \right) = 0$ , then you can set  $r_{i,j}$  to anything and maintain equality. Thus, if a Cholesky factorization exists, it is not unique.
- If  $\left( a_{i,j} - \sum_{k=1}^{i-1} r_{k,i} r_{k,j} \right) \neq 0$ , then there is nothing you can set  $r_{i,j}$  to. So, no Cholesky factor of  $\mathbf{A}$  exists.

Now unfortunately, this class is not interested in proving when and when not each of the above problems will arise. So, for now just know that:

1.  $\mathbf{A}$  is positive definite if and only if  $\mathbf{R}$  exists and is unique.
2.  $\mathbf{A}$  is positive semidefinite if and only if  $\mathbf{R}$  exists but is not invertible and thus also not unique.
3.  $\mathbf{A}$  is not positive definite or semidefinite if and only if  $\mathbf{R}$  does not exist.

The second theorem covered in lecture today let's us easily prove one direction in each of the above implications. So, the challenging task would be to prove the other direction, and to do that you would need to show that for any positive definite or semidefinite matrix  $\mathbf{A}$ , you can always use the above algorithm to find a matrix  $\mathbf{R}$ .

It takes approximately  $\frac{1}{3}n^3$  flops to do Cholesky decomposition as  $n$  gets larger. This is notably half of what LU decomposition takes.

## Lecture 7: 1/24/2024

A norm of a vector  $\vec{x} \in \mathbb{R}^n$  is a real number  $\|\vec{x}\|$  that is assigned to  $\vec{x}$  satisfying:

- $\vec{x} \neq 0 \implies \|\vec{x}\| > 0$  whereas  $\|\vec{0}\| = 0$ .
- $\|c\vec{x}\| = |c|\|\vec{x}\|$  for all  $c \in \mathbb{R}$ .
- $\|\vec{x} + \vec{y}\| \leq \|\vec{x}\| + \|\vec{y}\|$  for all  $\vec{x}, \vec{y} \in \mathbb{R}^n$

Some important vector norms:

1. Vector  $p$ -norm: For an integer  $p \geq 1$ , we define  $\|\vec{x}\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}$
  2. Infinity Norm:  $\|\vec{x}\|_\infty = \max \{|x_i| \mid 1 \leq i \leq n\}$
- 

A matrix norm assigns a real value  $\|\mathbf{A}\|$  to a matrix  $\mathbf{A}$  satisfying:

- $\mathbf{A} \neq \mathbf{0} \implies \|\mathbf{A}\| > 0$  whereas  $\|\mathbf{0}\| = 0$ .
- $\|c\mathbf{A}\| = |c|\|\mathbf{A}\|$  for all  $c \in \mathbb{R}$ .
- $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$  for all  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$

Some important matrix norms:

1. Frobenius Norm:  $\|\mathbf{A}\| = \left( \sum_{i,j=1}^n a_{i,j}^2 \right)^{\frac{1}{2}}$

Assuming  $\mathbf{A}$  is an  $m \times n$  matrix, this norm is equivalent to stringing out  $\mathbf{A}$ 's rows or columns to form an  $mn$  element vector and then taking the 2-norm of the resulting vector.

$$\left\| \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \right\| = \|(a, b, c, d, e, f, g, h, i)\|$$

2. Induced Norm: Let  $\|\cdot\|$  be a vector norm on  $\mathbb{R}^m$  and  $\mathbb{R}^n$ . Then the matrix norm induced by  $\|\cdot\|$  is:

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

The induced norm can be thought of as measuring the maximum stretch which a linear function  $L(\vec{x}) = \mathbf{A}\vec{x}$  applies to  $\vec{x}$  relative to  $\vec{x}$ 's starting length.

Importantly by the properties of a norm:

$$\frac{\|\mathbf{A}\vec{x}\|}{\|\vec{x}\|} = \left\| \frac{1}{\|\vec{x}\|} \mathbf{A}\vec{x} \right\| = \left\| \mathbf{A} \frac{\vec{x}}{\|\vec{x}\|} \right\| = \|\mathbf{A}\hat{x}\|$$

Thus, it's also common to see induced norms defined as:  $\|\mathbf{A}\| = \max_{\|\vec{x}\|=1} \|\mathbf{A}\vec{x}\|$

Some important formulas: (let  $\mathbf{A}$  be an  $m \times n$  matrix)

- $\|\mathbf{A}\|_1 = \max_{1 \leq j \leq n} \sum_{k=1}^m |a_{k,j}|$

To prove this, first note that:

$$\|\mathbf{A}\vec{x}\|_1 = \sum_{i=1}^m \left| \sum_{j=1}^n a_{i,j} x_j \right| \leq \sum_{i=1}^m \sum_{j=1}^n |a_{i,j}| |x_j| = \sum_{j=1}^n \left( |x_j| \sum_{i=1}^m |a_{i,j}| \right)$$

Now lets replace each  $a_{i,j}$  with the max element  $a_{i,k}$  in the  $i$ th row. Thus:

$$\sum_{j=1}^n \left( |x_j| \sum_{i=1}^m |a_{i,j}| \right) \leq \sum_{j=1}^n \left( |x_j| \max_{1 \leq k \leq m} \sum_{i=1}^m |a_{i,k}| \right)$$

And now we can separate the summands to get that:

$$\|\mathbf{A}\vec{x}\|_1 \leq \|\vec{x}\|_1 \max_{1 \leq k \leq m} \left( \sum_{i=1}^m |a_{i,k}| \right)$$

Now  $\|\mathbf{A}\|_1$  is the max value of  $\|\mathbf{A}\vec{x}\|_1$  when  $\|\vec{x}\|_1 = 1$ . So clearly:

$$\|\mathbf{A}\|_1 \leq \max_{1 \leq k \leq m} \left( \sum_{i=1}^m |a_{i,k}| \right)$$

Finally, to show equality note that the  $k$ th. standard unit basis vector  $e_k$  satisfies both that  $\|e_k\|_1 = 1$  and that:

$$\|\mathbf{A}e_k\| = \max_{1 \leq i \leq m} \left( \sum_{i=1}^m |a_{i,k}| \right)$$

.

- $\|\mathbf{A}\|_\infty = \max_{1 \leq i \leq m} \sum_{k=1}^n |a_{i,k}|$

This one is more obvious. The  $\vec{x}$  which maximizes  $\|\mathbf{A}\vec{x}\|_\infty$  while  $\|\vec{x}\|_\infty = 1$  is the vector whose elements are 1 and  $-1$  such that the  $k$ th element of  $\vec{x}$  has the same sign as the  $k$ th element of the row of  $\mathbf{A}$  whose elements have the largest sum of absolute values.

- $\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}$

( $\lambda_{\max}(\mathbf{M})$  refers to the largest eigenvalue of  $\mathbf{M}$ .)

Neither the lecture notes or the author of our textbook Börgers at this point gives a proof of this. However, at the very least we know from the proposition on page 8 that  $\mathbf{A}^T \mathbf{A}$  is positive definite or positive semidefinite. So, all of its eigenvalues are greater than or equal to zero, meaning that our distance formula will always give a real value.

## Lecture 8: 1/26/2024

Here is proof that induced norms are matrix norms.

### 1. Positivity:

If  $\|\mathbf{A}\| = 0$ , then  $\|\mathbf{A} \vec{x}\| = 0$  for all nonzero  $\vec{x}$ . This means that  $\mathbf{A} \vec{x} = \vec{0}$  for all nonzero  $\vec{x}$ . However, note that  $\mathbf{A} e_k$  equals the  $k$ th column of  $\mathbf{A}$ . So, all columns of  $\mathbf{A}$  must be zero. Hence,  $\|\mathbf{A}\| = 0 \implies \mathbf{A} = \mathbf{0}$

### 2. Multiplication by a constant:

$$\|c\mathbf{A}\| = \max_{\vec{x} \neq 0} \frac{\|c\mathbf{A} \vec{x}\|}{\|\vec{x}\|} = \max_{\vec{x} \neq 0} \frac{|c| \|\mathbf{A} \vec{x}\|}{\|\vec{x}\|} = |c| \left( \max_{\vec{x} \neq 0} \frac{\|\mathbf{A} \vec{x}\|}{\|\vec{x}\|} \right) = |c| \|\mathbf{A}\|$$

### 3. Triangle inequality:

$$\|\mathbf{A} + \mathbf{B}\| = \max_{\vec{x} \neq 0} \frac{\|(\mathbf{A} + \mathbf{B}) \vec{x}\|}{\|\vec{x}\|} = \max_{\vec{x} \neq 0} \frac{\|\mathbf{A} \vec{x} + \mathbf{B} \vec{x}\|}{\|\vec{x}\|}$$

Now, by the triangle inequality of the vector norm, we have that:

$$\frac{\|\mathbf{A} \vec{x} + \mathbf{B} \vec{x}\|}{\|\vec{x}\|} \leq \frac{\|\mathbf{A} \vec{x}\| + \|\mathbf{B} \vec{x}\|}{\|\vec{x}\|} = \frac{\|\mathbf{A} \vec{x}\|}{\|\vec{x}\|} + \frac{\|\mathbf{B} \vec{x}\|}{\|\vec{x}\|}$$

So,  $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$ .

Here are two more properties of induced norms:

### 4. $\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$ This is called submultiplicativity.

$$\|\mathbf{AB}\| = \max_{\vec{x} \neq 0} \frac{\|(\mathbf{AB}) \vec{x}\|}{\|\vec{x}\|} = \max_{\vec{x} \neq 0} \left( \frac{\|\mathbf{AB} \vec{x}\|}{\|\mathbf{B} \vec{x}\|} \frac{\|\mathbf{B} \vec{x}\|}{\|\vec{x}\|} \right) \text{ if } \mathbf{B} \vec{x} \neq 0.$$

Now this is obviously less than or equal to  $\max_{\mathbf{B} \vec{x} \neq 0} \left( \frac{\|\mathbf{AB} \vec{x}\|}{\|\mathbf{B} \vec{x}\|} \right) \max_{\vec{x} \neq 0} \left( \frac{\|\mathbf{B} \vec{x}\|}{\|\vec{x}\|} \right)$

Finally defining  $\vec{y} = \mathbf{B} \vec{x}$ , we get that:

$$\|\mathbf{AB}\| \leq \max_{\vec{y} \neq 0} \left( \frac{\|\mathbf{A} \vec{y}\|}{\|\vec{y}\|} \right) \max_{\vec{x} \neq 0} \left( \frac{\|\mathbf{B} \vec{x}\|}{\|\vec{x}\|} \right) = \|\mathbf{A}\| \|\mathbf{B}\|$$

5.  $\|\mathbf{A}\vec{x}\| \leq \|\mathbf{A}\| \|\vec{x}\|$  for all  $\vec{x} \in \mathbb{R}^n$ .

To explain this, remember that  $\|\mathbf{A}\|$  is the maximum value that  $\frac{\|\mathbf{A}\vec{x}\|}{\|\vec{x}\|}$  could be.

## Sensitivity of $\mathbf{A}\vec{x} = \vec{b}$ with respect to perturbation

Due to rounding errors and limited precision, we know that  $\vec{b}$  and  $\mathbf{A}$  will typically not be perfectly accurate when doing calculations. We refer to introducing this inaccuracy as perturbing  $\vec{b}$  and  $\mathbf{A}$ . Now as seen in the demonstration below, this will effect the accuracy of our calculations. So, it would be nice to be able to predict and quantify that loss in accuracy.

Demonstration of perturbing  $\vec{b}$ :

$$\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 0 \end{bmatrix} \text{ gives a solution of } \vec{x} = \begin{bmatrix} 2 \\ 0 \end{bmatrix}. \text{ Meanwhile,}$$

$$\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 0.001 \end{bmatrix} \text{ gives a solution of } \vec{x}' = \begin{bmatrix} 1.999 \\ 0.001 \end{bmatrix}.$$

The error in  $\vec{b}$  is  $\|\vec{b} - \vec{b}'\|_2 = 0.001$ , whereas the corresponding error in  $\vec{x}$  is  $\|\vec{x} - \vec{x}'\|_2 = \sqrt{2}(0.001) \approx 0.0014$

For now, let's address perturbing  $\vec{b}$ ...

Assume  $\mathbf{A}$  is an invertible matrix.

Formally, we shall write:

$$\begin{array}{ll} \text{Original equation:} & \mathbf{A}\vec{x} = \vec{b} \\ \text{New equation:} & \mathbf{A}\vec{x}' = \vec{b}' \end{array} \quad \text{where } \vec{b}' = \vec{b} + \vec{\delta b} \text{ and } \vec{x}' = \vec{x} + \vec{\delta x}$$

Then for any vector norm and its induced matrix norm:  $\frac{\|\vec{\delta x}\|}{\|\vec{x}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\vec{\delta b}\|}{\|\vec{b}\|}$ .

Proof:

Firstly,  $\mathbf{A}\vec{x}' = \vec{b}' \implies \mathbf{A}\vec{x} + \mathbf{A}\vec{\delta x} = \vec{b} + \vec{\delta b}$ . Now since  $\mathbf{A}\vec{x} = \vec{b}$ , it must be the case that  $\mathbf{A}\vec{\delta x} = \vec{\delta b}$ . So,  $\vec{\delta x} = \mathbf{A}^{-1}\vec{\delta b}$ . From this, we now get that for any vector norm and corresponding induced matrix norm, we have that

$$\|\vec{\delta x}\| = \|\mathbf{A}^{-1}\vec{\delta b}\| \leq \|\mathbf{A}^{-1}\| \|\vec{\delta b}\|.$$

Secondly, using the same norms as before, we have that:

$\vec{b} = \mathbf{A} \vec{x} \implies \|\vec{b}\| = \|\mathbf{A} \vec{x}\| \leq \|\mathbf{A}\| \|\vec{x}\|$ . Thus dividing both sides by  $\|\vec{b}\| \|\vec{x}\|$ , we get that:

$$\frac{1}{\|\vec{x}\|} \leq \|\mathbf{A}\| \frac{1}{\|\vec{b}\|}$$

Thus in conclusion: 
$$\frac{\|\vec{\delta x}\|}{\|\vec{x}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\vec{\delta b}\|}{\|\vec{b}\|}$$

We call  $K(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$  the condition number of  $\mathbf{A}$  under some induced norm (indicated as a subscript to  $K$ ). Clearly, it is useful as it gives us an upper bound to the relative error of our solution  $\vec{x}$  given the relative error in  $\vec{b}$ .

- If  $K(\mathbf{A})$  is small (close to 1), then  $\mathbf{A}$  is called well conditioned.
- If  $K(\mathbf{A})$  is large, then  $\mathbf{A}$  is called ill conditioned.

## Lecture 9: 1/29/2024

Now let's consider perturbing  $\mathbf{A}$ . As before, we shall assume that  $\mathbf{A}$  is invertible.

Formally, we shall write:

Original equation:  $\mathbf{A} \vec{x} = \vec{b}$       where  $\mathbf{A}' = \mathbf{A} + \delta \mathbf{A}$  and  $\vec{x}' = \vec{x} + \vec{\delta x}$   
 New equation:  $\mathbf{A}' \vec{x}' = \vec{b}$

Then for any vector norm and its induced matrix norm: 
$$\frac{\|\vec{\delta x}\|}{\|\vec{x}'\|} \leq K(\mathbf{A}) \frac{\|\delta \mathbf{A}\|}{\|\mathbf{A}\|}.$$

Proof:

We can rewrite  $\mathbf{A}' \vec{x}' = \vec{b}$  as  $\mathbf{A} \vec{x} + \mathbf{A} \vec{\delta x} + \delta \mathbf{A} (\vec{x} + \vec{\delta x}) = \vec{b}$ . Then, by canceling out  $\vec{b}$  and  $\mathbf{A} \vec{x}$ , we get the formula:  $\mathbf{A} \vec{\delta x} + \delta \mathbf{A} \vec{x}' = 0$ . And, further rearranging of this yields:

$$\vec{\delta x} = -\mathbf{A}^{-1}(\delta \mathbf{A}) \vec{x}'$$

So by the submultiplicativity property of matrix norms:

$$\|\vec{\delta x}\| = \|\mathbf{A}^{-1}(\delta \mathbf{A}) \vec{x}'\| \leq \|\mathbf{A}^{-1}\| \|(\delta \mathbf{A}) \vec{x}'\| \leq \|\mathbf{A}^{-1}\| \|(\delta \mathbf{A})\| \|\vec{x}'\|$$

Dividing both sides by our above inequality by  $\|\vec{x}'\|$  and multiplying the greater side by  $1 = \frac{\|\mathbf{A}\|}{\|\mathbf{A}\|}$ , we finally get that:

$$\frac{\|\vec{\delta x}\|}{\|\vec{x}'\|} \leq \|\mathbf{A}^{-1}\| \|\mathbf{A}\| \frac{\|\delta \mathbf{A}\|}{\|\mathbf{A}\|} = K(\mathbf{A}) \frac{\|\delta \mathbf{A}\|}{\|\mathbf{A}\|}$$

Finally, let's consider perturbing  $\mathbf{A}$  and  $\vec{b}$ . Once again, we will assume  $\mathbf{A}^{-1}$  exists.

Formally, we shall write:

$$\begin{array}{ll} \text{Original equation:} & \mathbf{A} \vec{x} = \vec{b} \quad \text{where } \mathbf{A}' = \mathbf{A} + \delta\mathbf{A}, \vec{b}' = \vec{b} + \delta\vec{b}, \\ \text{New equation:} & \mathbf{A}' \vec{x}' = \vec{b}' \quad \text{and } \vec{x}' = \vec{x} + \delta\vec{x} \end{array}$$

Then for any vector norm and its induced matrix norm:

$$\frac{\|\delta\vec{x}\|}{\|\vec{x}'\|} \leq K(A) \left( \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} + \frac{\|\delta\vec{b}\|}{\|\vec{b}'\|} + \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \frac{\|\delta\vec{b}\|}{\|\vec{b}'\|} \right)$$

Proof:

Step 1:

$\mathbf{A}' \vec{x}' = \vec{b}'$  can be rewritten as:  $\mathbf{A} \vec{x} + \mathbf{A} \delta\vec{x} + \delta\mathbf{A}(\vec{x} + \delta\vec{x}) = \vec{b} + \delta\vec{b}$ . Then as  $\mathbf{A} \vec{x} = \vec{b}$ , we cancel those two terms to get that  $\mathbf{A} \delta\vec{x} + \delta\mathbf{A}(\vec{x}') = \delta\vec{b}$ . And so, a little rearranging yields:

$$\delta\vec{x} = \mathbf{A}^{-1} (\delta\vec{b} - \delta\mathbf{A} \vec{x}')$$

Next, using properties of the induced matrix norm, we get that:

$$\begin{aligned} \|\delta\vec{x}\| &= \|\mathbf{A}^{-1} (\delta\vec{b} - \delta\mathbf{A} \vec{x}')\| \leq \|\mathbf{A}^{-1}\| \|\delta\vec{b} - \delta\mathbf{A} \vec{x}'\| \\ &\leq \|\mathbf{A}^{-1}\| (\|\delta\vec{b}\| + \|\delta\mathbf{A} \vec{x}'\|) \\ &\leq \|\mathbf{A}^{-1}\| (\|\delta\vec{b}\| + \|\delta\mathbf{A}\| \|\vec{x}'\|) \end{aligned}$$

Therefore, by dividing both sides of our inequality by  $\|\vec{x}'\|$  and multiplying the greater side by  $1 = \frac{\|\mathbf{A}\|}{\|\mathbf{A}\|}$ , we finally get that:

$$\frac{\|\delta\vec{x}\|}{\|\vec{x}'\|} \leq K(A) \left( \frac{\|\delta\vec{b}\|}{\|\mathbf{A}\| \|\vec{x}'\|} + \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \right)$$

Step 2:

$$\vec{b}' = \mathbf{A}' \vec{x}' \implies \|\vec{b}'\| = \|(\mathbf{A}') \vec{x}'\| \leq \|\mathbf{A} + \delta\mathbf{A}\| \|\vec{x}'\|$$

So dividing by  $\|\vec{b}'\| \|\vec{x}'\|$ , we get that:

$$\frac{1}{\|\vec{x}'\|} \leq \frac{\|\mathbf{A} + \delta\mathbf{A}\|}{\|\vec{b}'\|} \leq \frac{\|\mathbf{A}\| + \|\delta\mathbf{A}\|}{\|\vec{b}'\|}$$

Step 3:

By combining the inequalities in steps 1 and 2, we know:

$$\frac{\|\vec{\delta x}\|}{\|\vec{x}'\|} \leq K(A) \left( \frac{\|\vec{\delta b}\|}{\|\mathbf{A}\|} \left( \frac{\|\mathbf{A}\| + \|\delta \mathbf{A}\|}{\|\vec{b}'\|} \right) + \frac{\|\delta \mathbf{A}\|}{\|\mathbf{A}\|} \right)$$

We can then manipulate that into the inequality claimed above.

If you are wondering why the upper bounds we derived in this lecture are for  $\frac{\|\vec{\delta x}\|}{\|\vec{x}'\|}$  instead of  $\frac{\|\vec{\delta x}\|}{\|\vec{x}\|}$ , know that I am too. I asked the professor and he said its because the upper bounds for  $\frac{\|\vec{\delta x}\|}{\|\vec{x}\|}$  are more difficult. That said, our textbook does derive upper bounds for  $\frac{\|\vec{\delta x}\|}{\|\vec{x}\|}$  when perturbing  $\mathbf{A}$ .

Proposition: Let  $\|\cdot\|$  be an induced matrix norm. Then:

1.  $\|\mathbf{I}\| = 1$
2.  $K(\mathbf{A}) \geq 1$ .

Proof:

$$1. \|\mathbf{I}\| = \max_{\vec{x} \neq 0} \frac{\|\mathbf{I}\vec{x}\|}{\|\vec{x}\|} = \max_{\vec{x} \neq 0} \frac{\|\vec{x}\|}{\|\vec{x}\|} = 1$$

$$2. \mathbf{I} = \mathbf{A}\mathbf{A}^{-1} \implies 1 = \|\mathbf{I}\| = \|\mathbf{A}\mathbf{A}^{-1}\| \leq \|\mathbf{A}\|\|\mathbf{A}^{-1}\| = K(\mathbf{A})$$

## Lecture 10: 2/5/2024

A matrix  $\mathbf{Q}$  is orthogonal if  $\mathbf{Q}^{-1} = \mathbf{Q}^T$

If  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  is orthogonal, then:

1.  $\langle \mathbf{Q}\vec{x}, \mathbf{Q}\vec{y} \rangle = \vec{x}^T \mathbf{Q}^T \mathbf{Q} \vec{y} = \vec{x}^T \vec{y} = \langle \vec{x}, \vec{y} \rangle$
2.  $\|\mathbf{Q}\vec{x}\|_2 = \langle \mathbf{Q}\vec{x}, \mathbf{Q}\vec{x} \rangle = \langle \vec{x}, \vec{x} \rangle = \|\vec{x}\|_2$

Theorem (QR decomposition / full QR decomposition):

Suppose  $\mathbf{A} \in \mathbb{R}^{n \times m}$  where  $n \geq m$ . Then there exists an orthogonal matrix

$\mathbf{Q} \in \mathbb{R}^{n \times n}$  and a matrix  $\mathbf{R} = \begin{bmatrix} \hat{\mathbf{R}} \\ \mathbf{0} \end{bmatrix} \in \mathbb{R}^{n \times m}$  where  $\hat{\mathbf{R}}$  is an  $m \times m$  upper triangular matrix such that  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ .



Some uses of QR decompositions:

- Solving matrix-vector equations:

Using the QR decomposition of a matrix, you can rewrite an equation  $\mathbf{A}\vec{x} = \vec{b}$  as follows:

$$\mathbf{A}\vec{x} = \mathbf{Q}\mathbf{R}\vec{x} = \vec{b} \Rightarrow \mathbf{R}\vec{x} = \mathbf{Q}^{-1}\vec{b}$$

And now you can find  $\vec{x}$  using back substitution and matrix vector multiplication.

- The Least Square Problem (LSE): Suppose  $\mathbf{A} \in \mathbb{R}^{n \times m}$  where  $n \geq m$  and  $b \in \mathbb{R}^n$ . Then the least square problem is of the form:  $\min_{\vec{x} \in \mathbb{R}^m} (\|\vec{b} - \mathbf{A}\vec{x}\|_2)^2$ .

For  $n \geq m$ , if the matrix-vector equation  $\mathbf{A}\vec{x} = \vec{b}$  has a solution, then the solution of the least square problem is also the solution of that matrix vector equation.

However, if there is no solution to the matrix-vector equation, then the least square problem is asking us to find a vector  $\vec{x}$  that minimizes the distance between  $\mathbf{A}\vec{x}$  and  $\vec{b}$ .

**Step 1.** Find a QR-decomposition such that  $\mathbf{A} = \mathbf{Q}\mathbf{R}$

$$\begin{aligned} \text{Step 2. } (\|\vec{b} - \mathbf{A}\vec{x}\|_2)^2 &= (\|\vec{b} - \mathbf{Q}\mathbf{R}\vec{x}\|_2)^2 = (\|\mathbf{Q}^T(\vec{b} - \mathbf{Q}\mathbf{R}\vec{x})\|_2)^2 \\ &= (\|\mathbf{Q}^T\vec{b} - \mathbf{Q}^T\mathbf{Q}\mathbf{R}\vec{x}\|_2)^2 = (\|\mathbf{Q}^T\vec{b} - \mathbf{R}\vec{x}\|_2)^2 \end{aligned}$$

$$\begin{aligned} \text{Now denote: } \mathbf{Q}^T\vec{b} &= \begin{bmatrix} \vec{c} \\ \vec{d} \end{bmatrix} \text{ and } \mathbf{R}\vec{x} = \begin{bmatrix} \hat{\mathbf{R}}\vec{x} \\ \vec{0} \end{bmatrix} \\ \dots \text{where } \vec{c}, \hat{\mathbf{R}}\vec{x} &\in \mathbb{R}^m \text{ and } \vec{d} \in \mathbb{R}^{n-m} \end{aligned}$$

$$\begin{aligned} (\|\mathbf{Q}^T\vec{b} - \mathbf{R}\vec{x}\|_2)^2 &= \left( \left\| \begin{bmatrix} \vec{c} \\ \vec{d} \end{bmatrix} - \begin{bmatrix} \hat{\mathbf{R}}\vec{x} \\ \vec{0} \end{bmatrix} \right\|_2 \right)^2 = \left\| \begin{bmatrix} \vec{c} - \hat{\mathbf{R}}\vec{x} \\ \vec{d} \end{bmatrix} \right\|_2^2 \\ &= (\|\vec{c} - \hat{\mathbf{R}}\vec{x}\|_2)^2 + (\|\vec{d}\|_2)^2 \end{aligned}$$

$$\text{Step 3. Now } \min_{\vec{x} \in \mathbb{R}^m} (\|\vec{b} - \mathbf{A}\vec{x}\|_2)^2 = \left( \min_{\vec{x} \in \mathbb{R}^m} (\|\vec{c} - \hat{\mathbf{R}}\vec{x}\|_2)^2 \right) + (\|\vec{d}\|_2)^2$$

**Step 4.** If  $\mathbf{A}$  has full-rank, then  $\hat{\mathbf{R}} \in \mathbb{R}^{m \times m}$  is invertible. So:  $\hat{\mathbf{R}}\vec{x} = \vec{c}$  has a unique solution, meaning that  $\min_{\vec{x} \in \mathbb{R}^m} (\|\vec{b} - \mathbf{A}\vec{x}\|_2)^2 = (\|\vec{d}\|_2)^2$ . So we get the theorem on the next page:

Theorem: (For if  $\mathbf{A}$  has full rank...)

- The solution to the least squares problem is the solution to the linear system  $\hat{\mathbf{R}}\vec{x} = \vec{c}$  where  $\vec{c}$  is the first  $m$  values of  $\mathbf{Q}^T \vec{b}$ .
- Meanwhile, the optimal value is  $\min_{\vec{x} \in \mathbb{R}^m} (\|\vec{b} - \mathbf{A}\vec{x}\|_2)^2 = (\|\vec{d}\|_2)^2$  where  $\vec{d}$  is the remaining elements of  $\mathbf{Q}^T \vec{b}$ .

## Lecture 11: 2/7/2024

A QR decomposition of  $\begin{bmatrix} -1 & -1 & 1 \\ 1 & 3 & 3 \\ -1 & -1 & 5 \\ 1 & 3 & 7 \end{bmatrix}$  is:  $\mathbf{QR} = \begin{bmatrix} -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 2 & 4 & 2 \\ 0 & 2 & 8 \\ 0 & 0 & 4 \\ 0 & 0 & 0 \end{bmatrix}$

Meanwhile, a reduced QR decomposition of that same matrix would be:

$$\hat{\mathbf{Q}}\hat{\mathbf{R}} = \begin{bmatrix} -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 2 & 4 & 2 \\ 0 & 2 & 8 \\ 0 & 0 & 4 \end{bmatrix}$$

...where  $\hat{\mathbf{R}}$  is the  $m \times m$  upper triangular inside  $\mathbf{R}$  while  $\hat{\mathbf{Q}}$  is the matrix of the first  $m$  columns in  $\mathbf{Q}$ .

Two notes:

- $\mathbf{QR} = \hat{\mathbf{Q}}\hat{\mathbf{R}}$ .
- $\hat{\mathbf{Q}}$  is not orthogonal but it does have the property that  $\hat{\mathbf{Q}}^T \hat{\mathbf{Q}} = I_m$ . In other words, the columns of  $\hat{\mathbf{Q}}$  are an orthonormal set (a.k.a it's column-orthogonal).

Using the Gram-Schmit algorithm, for any  $\mathbf{A} \in \mathbb{R}^{n \times m}$  with full rank, we can calculate a column-orthogonal matrix  $\hat{\mathbf{Q}}$  and upper triangular matrix  $\hat{\mathbf{R}}$  such that  $\mathbf{A} = \hat{\mathbf{Q}}\hat{\mathbf{R}}$ .

### Classic Gram-Schmit Algorithm:

We'll denote  $\mathbf{A} = [\vec{a}_1 \ \cdots \ \vec{a}_m]$  and  $\hat{\mathbf{Q}} = [\vec{q}_1 \ \cdots \ \vec{q}_m]$

For  $i = 1, \dots, m$ :

For  $k = 1, \dots, (i - 1)$ :

$$r_{k,i} = \langle \vec{a}_i, \vec{q}_k \rangle$$

end

$$r_{i,i} = \left\| \vec{a}_i - \sum_{k=1}^{i-1} (r_{k,i} \vec{q}_k) \right\|_2 \text{ and } \vec{q}_i = \frac{\vec{a}_i - \sum_{k=1}^{i-1} (r_{k,i} \vec{q}_k)}{r_{i,i}}$$

end

To remember this algorithm, note that:

- $\hat{\mathbf{Q}}\hat{\mathbf{R}} = \mathbf{A} \implies \hat{\mathbf{R}} = \hat{\mathbf{Q}}^T \mathbf{A}$ . So,  $r_{i,j}$  equals the inner product of the  $i$ th column of  $\hat{\mathbf{Q}}^T$  and the  $j$ th column of  $\mathbf{A}$ .
- If  $V_i$  is the space spanned by the first  $i$  columns of  $\hat{\mathbf{Q}}$ , then  $\vec{q}_{i+1}$  is equal to the normalized result of  $\vec{a}_{i+1} - \text{Proj}_{V_i}(\vec{a}_{i+1})$ . However, each  $\vec{q}_i$  is a unit vector. So, the projection formula is simplified.

### Modified Gram-Schmit Algorithm: (more stable for computers)

```

1  function [Q, R] = modifiedGS(A)
2      % Storing the number of columns in A
3      m = size(A, 2);
4
5      % Copy Q into A (Matlab doesn't seem to alias the matrices)
6      Q = A;
7
8      for i=1:m
9          R(i,i) = norm(Q(:,i));
10         Q(:,i) = Q(:,i) / R(i,i) %Q(:,i) * (1 / R(i,i))
11
12         for j = (i+1):m
13             R(i,j) = Q(:,i)' * Q(i,j); % A' is the transpose of A
14             Q(:,j) = Q(:,j) - (R(i,j) * Q(:,i));
15         end
16     end
17 end

```

The way that this modified algorithm differs from the original one is that instead of building a column of  $\hat{\mathbf{Q}}$  every iteration, it instead builds a row of  $\hat{\mathbf{R}}$  every iteration. This is much more numerically stable.

## Lecture 12: 2/9/2024

A matrix  $\mathbf{P} \in \mathbb{R}^{n \times n}$  is called a projector if  $\mathbf{P}^2 = \mathbf{P}$ .

Given  $\vec{u} \in \mathbb{R}^n$  with  $\|\vec{u}\|_2 = 1$ , then  $\mathbf{P} = \vec{u} \vec{u}^T$  is an orthogonal projector.

**Proof:**

- $\mathbf{P}^2 = (\vec{u} \vec{u}^T)(\vec{u} \vec{u}^T) = \vec{u} \langle \vec{u}, \vec{u} \rangle \vec{u}^T = 1 \cdot \vec{u} \vec{u}^T = \mathbf{P}$
- $\mathbf{P}^T = (\vec{u} \vec{u}^T)^T = (\vec{u}^T)^T \vec{u} = \vec{u} \vec{u}^T = \mathbf{P}$

Properties of  $\mathbf{P} = \vec{u} \vec{u}^T$ :

- $\mathbf{P} \vec{u} = \vec{u}$ .
- $\langle \vec{u}, \vec{v} \rangle = 0 \implies \mathbf{P} \vec{v} = 0$

Let  $\vec{u} \in \mathbb{R}^n$  with  $\|\vec{u}\|_2 = 1$ . Then the matrix  $\mathbf{Q} = \mathbf{I} - 2\vec{u} \vec{u}^T$  is called a (Householder) reflector.

If  $\|\vec{u}\|_2 \neq 1$ , you can always normalize it and construct a reflector:  $\mathbf{Q} = \mathbf{I} - \frac{2}{(\|\vec{u}\|_2)^2} \vec{u} \vec{u}^T$

Properties of  $\mathbf{Q} = \mathbf{I} - 2\vec{u} \vec{u}^T$  where  $\|\vec{u}\|_2 = 1$ :

- $\mathbf{Q} \vec{u} = -\vec{u}$ .

**Proof:**  $\mathbf{Q} \vec{u} = (\mathbf{I} - 2\vec{u} \vec{u}^T) \vec{u} = \vec{u} - 2\vec{u} = -\vec{u}$

- $\langle \vec{u}, \vec{v} \rangle = 0$

**Proof:**  $\mathbf{Q} \vec{v} = (\mathbf{I} - 2\vec{u} \vec{u}^T) \vec{v} = \vec{v} - 2\vec{u}(\vec{u}^T \vec{v}) = \vec{v} - 2\vec{u}(0) = \vec{v}$

- $\mathbf{Q}^T = \mathbf{Q}$

**Proof:**  $\mathbf{Q}^T = (\mathbf{I} - 2\vec{u} \vec{u}^T)^T = \mathbf{I}^T - 2(\vec{u} \vec{u}^T)^T = \mathbf{I} - 2\vec{u} \vec{u}^T = \mathbf{Q}$

- $\mathbf{Q}^2 = \mathbf{I}$

**Proof:**  $\mathbf{Q}^2 = (\mathbf{I} - 2\vec{u} \vec{u}^T)^2 = \mathbf{I} - 4\vec{u} \vec{u}^T + 4(\vec{u} \vec{u}^T)^2 = \mathbf{I}$

Basically,  $\mathbf{Q}$  reflects any vector  $\vec{x}$  about the hyperplane perpendicular to  $\vec{u}$ .

Suppose  $\vec{x} \in \mathbb{R}^n$  and  $\vec{y} = [\|\vec{x}\|_2 \ 0 \ \cdots \ 0]^T = \|\vec{x}\|_2 e_1$ . (Thus  $\|\vec{y}\| = \|\vec{x}\|$ )  
 Let  $\vec{v} = \vec{x} - \vec{y}$ , and  $\vec{u} = \frac{\vec{v}}{\|\vec{v}\|_2}$ . Then  $\mathbf{Q} = \mathbf{I} - 2\vec{u}\vec{u}^T$  satisfies:  $\mathbf{Q}\vec{x} = \vec{y}$ .

**Proof:**

$$\begin{aligned}
 \mathbf{Q}\vec{x} &= (\mathbf{I} - 2\vec{u}\vec{u}^T)\vec{x} = \vec{x} - 2\langle \vec{u}, \vec{x} \rangle \vec{u} \\
 &= \vec{x} - \frac{2}{(\|\vec{x} - \vec{y}\|_2)^2} ((\vec{x} - \vec{y})^T \vec{x})(\vec{x} - \vec{y}) \\
 &= \vec{x} - \frac{1}{(\|\vec{x} - \vec{y}\|_2)^2} (2(\|\vec{x}\|_2)^2 - 2\vec{y}^T \vec{x})(\vec{x} - \vec{y}) \\
 &= \vec{x} - \frac{1}{(\|\vec{x} - \vec{y}\|_2)^2} ((\|\vec{x}\|_2)^2 + (\|\vec{y}\|_2)^2 - 2\vec{y}^T \vec{x})(\vec{x} - \vec{y}) \\
 &= \vec{x} - \frac{1}{(\|\vec{x} - \vec{y}\|_2)^2} (\|\vec{x} - \vec{y}\|_2)^2 (\vec{x} - \vec{y}) \\
 &= \vec{x} - \vec{x} + \vec{y} = \vec{y}
 \end{aligned}$$


---

Now let's consider taking a full QR decomposition:

Consider  $\mathbf{A} = [\vec{a}_1 \ \vec{a}_2 \ \cdots \ \vec{a}_m] \in \mathbb{R}^{n \times m}$ .

**Step 1.** Compute a reflector  $\mathbf{Q}_1$  such that  $\mathbf{Q}_1 \vec{a}_1 = \|\vec{a}_1\|_2 e_1$ . Thus:

$$\mathbf{Q}_1 \mathbf{A} = [\mathbf{Q}_1 \vec{a}_1 \ \mathbf{Q}_1 \vec{a}_2 \ \cdots \ \mathbf{Q}_1 \vec{a}_m] = \begin{bmatrix} \|\vec{a}_1\|_2 & \vec{b}^T \\ \vec{0} & \tilde{\mathbf{A}}_2 \end{bmatrix}$$

So, now consider  $\tilde{\mathbf{A}}_2 = [\tilde{a}_2 \ \cdots \ \tilde{a}_m] \in \mathbb{R}^{(n-1) \times (m-1)}$

**Step 2.** Now consider a reflector  $\tilde{\mathbf{Q}}_2 \in \mathbb{R}^{(n-1) \times (m-1)}$  such that  $\tilde{\mathbf{Q}}_2 \tilde{a}_2 = \|\tilde{a}_2\|_2 e_1$ .

Then define:  $\mathbf{Q}_2 = \begin{bmatrix} \mathbf{I}_1 & \vec{0} \\ \vec{0} & \tilde{\mathbf{Q}}_2 \end{bmatrix}$

Therefore, we have that  $\mathbf{Q}_2 \mathbf{Q}_1 \mathbf{A} = \begin{bmatrix} \hat{\mathbf{R}} & \mathbf{B} \\ \vec{0} & \tilde{\mathbf{A}}_3 \end{bmatrix}$ , where  $\tilde{\mathbf{A}}_3 \in \mathbb{R}^{(n-2) \times (m-2)}$ ,

$\hat{\mathbf{R}}$  is  $2 \times 2$  upper triangular, and  $\mathbf{B}$  is some rectangular matrix.

Note that only the lowest row in  $\mathbf{B}$  is effected by this step.

**Step 3.** Repeat step 2. over and over again on  $\tilde{\mathbf{Q}}_k$  until you can't any more. At the end, you will get an equation of the form:

$$\mathbf{Q}_m \cdots \mathbf{Q}_2 \mathbf{Q}_1 \mathbf{A} = \begin{bmatrix} \hat{\mathbf{R}} \\ \vec{0} \end{bmatrix}$$

**Step 4.** Each  $Q_i$  is orthogonal by the properties of a reflector. And since the product of orthogonal matrices is orthogonal, we can thus coalesce the  $Q$  in order to get that  $\mathbf{A} = \mathbf{QR}$ .

## Lecture 13: 2/12/2024

Let  $\mathbf{A} \in \mathbb{R}^{n \times m}$  with  $\text{rank}(\mathbf{A}) = r$ . Then there exists orthogonal matrices  $\mathbf{U} \in \mathbb{R}^{n \times n}$  and  $\mathbf{V} \in \mathbb{R}^{m \times m}$ , as well as positive numbers  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  such that letting:

$$\Sigma = \begin{bmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_r & & \\ & & & \ddots & \\ & & & & 0 \end{bmatrix} \in \mathbb{R}^{n \times m},$$

we have that  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$ .

$\mathbf{U}\Sigma\mathbf{V}^T$  is called the SVD of  $\mathbf{A}$ . Also, each  $\sigma_i$  is called a singular value of  $\mathbf{A}$ . Specifically, SVD stands for "Singular Value Decomposition".

**Geometric Explanation:**

There exists an orthonormal basis  $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_m$  for  $\mathbb{R}^m$ , an orthonormal basis  $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_n$  for  $\mathbb{R}^n$ , and real numbers  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  such that:

$$\mathbf{A}\vec{v}_i = \begin{cases} \sigma_i \vec{u}_i & \text{if } i \in \{1, \dots, r\} \\ \vec{0} & \text{if } i \in \{r+1, \dots, m\} \end{cases} \quad \mathbf{A}^T \vec{u}_i = \begin{cases} \sigma_i \vec{v}_i & \text{if } i \in \{1, \dots, r\} \\ \vec{0} & \text{if } i \in \{r+1, \dots, n\} \end{cases}$$

To see this, let  $\mathbf{U} = [\vec{u}_1 \ \vec{u}_2 \ \dots \ \vec{u}_n]$ , and  $\mathbf{V} = [\vec{v}_1 \ \vec{v}_2 \ \dots \ \vec{v}_m]$ . Also let  $\Sigma$  be a matrix defined as seen above.

$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T \implies \mathbf{A}\mathbf{V} = \mathbf{U}\Sigma$ . Therefore, we have that:

$$[\mathbf{A}\vec{v}_1 \ \mathbf{A}\vec{v}_2 \ \dots \ \mathbf{A}\vec{v}_m] = \begin{bmatrix} \sigma_1 \vec{u}_1 & \sigma_2 \vec{u}_2 & \dots & \sigma_r \vec{u}_r & \vec{0} & \dots & \vec{0} \end{bmatrix}$$

Meanwhile  $\mathbf{A}^T = \mathbf{V}\Sigma^T\mathbf{U}^T$ . Thus,  $\mathbf{A}^T\mathbf{U} = \mathbf{V}\Sigma^T$ , which means that:

$$[\mathbf{A}^T \vec{u}_1 \ \mathbf{A}^T \vec{u}_2 \ \dots \ \mathbf{A}^T \vec{u}_m] = \begin{bmatrix} \sigma_1 \vec{v}_1 & \sigma_2 \vec{v}_2 & \dots & \sigma_r \vec{v}_r & \vec{0} & \dots & \vec{0} \end{bmatrix}$$

**Important fact:**

For any  $\vec{w} \in \mathbb{R}^m$ , we have that  $\vec{w} = \sum_{i=1}^m c_i \vec{v}_i$ . Therefore:

$$\mathbf{A}\vec{w} = \sum_{i=1}^m c_i (\mathbf{A}\vec{v}_i) = \sum_{i=1}^r c_i \sigma_i \vec{u}_i$$

Proposition:

(Note: In this class, the professor is writing  $R(\mathbf{A})$  to refer to the "range" [a.k.a column space] of  $\mathbf{A}$ . I don't personally like it but I'm also not going to rewrite  $R(\mathbf{A})$  every time he says it.)

$$R(\mathbf{A}) = \text{span}\{\vec{u}_1, \dots, \vec{u}_r\}$$

$$N(\mathbf{A}) = \text{span}\{\vec{v}_{r+1}, \dots, \vec{v}_m\}$$

$$R(\mathbf{A}^T) = \text{span}\{\vec{v}_1, \dots, \vec{v}_r\}$$

$$N(\mathbf{A}^T) = \text{span}\{\vec{u}_{r+1}, \dots, \vec{u}_n\}$$

$$R(\mathbf{A}^T) = N(\mathbf{A})^\perp$$

$$R(\mathbf{A}) = N(\mathbf{A}^T)^\perp$$

Similarly to QR decompositions, we can also find a condensed / reduced SVD of a matrix  $\mathbf{A}$ ...

Let  $\mathbf{A} \in \mathbb{R}^{n \times m}$  with  $\text{rank}(\mathbf{A}) = r$ . Then there exists  $\hat{\mathbf{U}} \in \mathbb{R}^{n \times r}$ ,  $\hat{\Sigma} \in \mathbb{R}^{r \times r}$ , and  $\hat{\mathbf{V}} \in \mathbb{R}^{m \times r}$  such that:

- $\hat{\mathbf{U}}^T \hat{\mathbf{U}} = \mathbf{I}_r$
- $\hat{\mathbf{V}}^T \hat{\mathbf{V}} = \mathbf{I}_r$
- $\hat{\Sigma}$  is a diagonal square matrix with decreasing entries as one goes down the diagonal.
- $\mathbf{A} = \hat{\mathbf{U}} \hat{\Sigma} \hat{\mathbf{V}}^T$

One last note for this lecture:

We won't be learning how to calculate SVDs of matrices yet...

That said, know that the singular values of  $\mathbf{A}$  are uniquely determined by  $\mathbf{A}$ .

## Lecture 14: 2/14/2024

Theorem: For  $\mathbf{A} \in \mathbb{R}^{n \times m}$ , we have that  $\|\mathbf{A}\|_2$  equals the largest singular value of  $\mathbf{A}$ .

Proof:

Let  $\mathbf{U}\Sigma\mathbf{V}^T$  be the SVD of  $\mathbf{A}$  where  $\mathbf{U} = [\vec{u}_1 \ \vec{u}_2 \ \dots \ \vec{u}_n]$ ,

$\mathbf{V} = [\vec{v}_1 \ \vec{v}_2 \ \dots \ \vec{v}_m]$ , and the singular values in  $\Sigma$  are  $\sigma_1, \dots, \sigma_r$ .

From last lecture, we know that  $\mathbf{A}\vec{v}_1 = \sigma_1\vec{u}_1$ . So:

$$\|\mathbf{A}\|_2 = \max_{\vec{x} \neq 0} \frac{\|\mathbf{A}\vec{x}\|_2}{\|\vec{x}\|_2} \geq \frac{\|\mathbf{A}\vec{v}_1\|_2}{\|\vec{v}_1\|_2} = \frac{\sigma_1\|\vec{u}_1\|}{\|\vec{v}_1\|} = \frac{\sigma_1(1)}{(1)} = \sigma_1$$

Meanwhile, note that  $\vec{x} \in \mathbb{R}^m$  can be expressed as a linear combination:  $c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_m\vec{v}_m$  because the columns of  $\mathbf{V}$  form a basis for  $\mathbb{R}^m$ . And, since  $\vec{v}_i \cdot \vec{v}_j = 0$  for all  $i \neq j$ , we know that:

$$(\|\vec{x}\|_2)^2 = \vec{x} \cdot \vec{x} = \sum_{i=1}^n c_i^2 (\vec{v}_i \cdot \vec{v}_i)$$

Also, because each  $\|\vec{v}_i\|_2 = 1$ , this further simplifies to:  $(\|\vec{x}\|_2)^2 = \sum_{i=1}^n c_i^2$

At the same time:

$$\mathbf{A}\vec{x} = \mathbf{A}(c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_m\vec{v}_m) = c_1\sigma_1\vec{u}_1 + c_2\sigma_2\vec{u}_2 + \dots + c_r\sigma_r\vec{u}_r.$$

And as each  $u_i$  is a unit vector and  $u_i \cdot u_j = 0$  for all  $i \neq j$ , we thus have that:

$$(\|\mathbf{A}\vec{x}\|_2)^2 = \sum_{i=1}^r c_i^2 \sigma_i^2$$

Then because  $\sigma_1$  is the largest singular value, we have that:

$$\sum_{i=1}^r c_i^2 \sigma_i^2 \leq \sum_{i=1}^r c_i^2 \sigma_1^2 = \sigma_1^2 \sum_{i=1}^r c_i^2 = \sigma_1^2 (\|\vec{x}\|_2)^2$$

Thus,  $(\|\mathbf{A}\vec{x}\|_2)^2 \leq \sigma_1^2 (\|\vec{x}\|_2)^2$ . And since  $\sigma_1$ ,  $\|\mathbf{A}\vec{x}\|$ ,  $\|\vec{x}\|$  are all positive, we have that  $\|\mathbf{A}\vec{x}\|_2 \leq \sigma_1 \|\vec{x}\|_2$ . Therefore:

$$\|\mathbf{A}\|_2 = \max_{\vec{x} \neq 0} \frac{\|\mathbf{A}\vec{x}\|_2}{\|\vec{x}\|_2} \leq \max_{\vec{x} \neq 0} \frac{\sigma_1 \|\vec{x}\|_2}{\|\vec{x}\|_2} = \sigma_1$$

Two notes:

- On page 12, it was mentioned that  $\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}$ . So, consider that a hint about how SVDs are calculated.
- $\mathbf{A}^T$  and  $\mathbf{A}$  have the same singular values. So  $\|\mathbf{A}\|_2 = \|\mathbf{A}^T\|_2$ .

Theorem: Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be an invertible matrix with singular values  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0$ . Then:

$$K_2(\mathbf{A}) = \|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2 = \frac{\sigma_1}{\sigma_n}$$

Proof:

Let  $\mathbf{U}\Sigma\mathbf{V}^T$  be the SVD of  $\mathbf{A}$ . Note that as  $\mathbf{A}$  is invertible (and square), we know that  $\Sigma$  is a diagonal square matrix which contains  $n$  strictly positive singular values. So,  $\Sigma$  is invertible.

Now note:  $\mathbf{A}^{-1} = (\mathbf{U}\Sigma\mathbf{V}^T)^{-1} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T$ . Thus, we see that an SVD of  $\mathbf{A}^{-1}$  is  $\mathbf{V}\Sigma^{-1}\mathbf{U}^T$ . Additionally,  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0$  implies that  $\frac{1}{\sigma_n} \geq \dots \geq \frac{1}{\sigma_2} \geq \frac{1}{\sigma_1} > 0$ . So the largest singular value in  $\Sigma^{-1}$  is  $\frac{1}{\sigma_n}$ . We therefore conclude that  $\|\mathbf{A}^{-1}\|_2 = \frac{1}{\sigma_n}$ .

Hence,  $\|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2 = \sigma_1 \frac{1}{\sigma_n}$ .



Problem: For  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and an integer  $k > 0$ , find the matrix  $\mathbf{X}$  with rank  $k$  that minimizes  $\|\mathbf{A} - \mathbf{X}\|_2$ .

To do this, first consider the following lemma:

Let  $\mathbf{A} \in \mathbb{R}^{n \times m}$  be a matrix with  $\text{rank}(\mathbf{A}) = r$ . Consider the singular values of  $\mathbf{A}$ :  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ , with their respective right and left singular vectors:  $\vec{v}_1, \dots, \vec{v}_r$  and  $\vec{u}_1, \dots, \vec{u}_r$ . Then, we can write that:

$$\mathbf{A} = \sum_{i=1}^r \sigma_i \vec{u}_i \vec{v}_i^T$$

Proof:

$$\begin{aligned} \mathbf{A} = \hat{\mathbf{U}} \hat{\Sigma} \hat{\mathbf{V}}^T &= \begin{bmatrix} \vec{u}_1 & \dots & \vec{u}_r \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{bmatrix} \begin{bmatrix} \vec{v}_1^T \\ \vdots \\ \vec{v}_r^T \end{bmatrix} \\ &= \begin{bmatrix} \sigma_1 \vec{u}_1 & \dots & \sigma_r \vec{u}_r \end{bmatrix} \begin{bmatrix} \vec{v}_1^T \\ \vdots \\ \vec{v}_r^T \end{bmatrix} \\ &= \sigma_1 \vec{u}_1 \vec{v}_1^T + \dots + \sigma_r \vec{u}_r \vec{v}_r^T \end{aligned}$$

In other words,  $\mathbf{A}$  is a linear combination of  $r$  many rank 1 matrices.

If  $\mathbf{A}$  has the SVD:  $\mathbf{A} = \sum_{i=1}^r \sigma_i \vec{u}_i \vec{v}_i^T$ , then the optimal matrix  $\mathbf{X}$  is  $\mathbf{A}_k = \sum_{i=1}^k \sigma_i \vec{u}_i \vec{v}_i^T$ .

Specifically,  $\mathbf{A}_k$  is the rank- $k$  approximation of  $\mathbf{A}$ .

Proposition:  $\|\mathbf{A} - \mathbf{A}_k\|_2 = \sigma_{k+1}$

Proof:  $\mathbf{A} - \mathbf{A}_k = \sum_{i=k+1}^r \sigma_i \vec{u}_i \vec{v}_i^T.$

Thus, the largest singular value of  $\mathbf{A} - \mathbf{A}_k$  is  $\sigma_{k+1}$ .

## Lecture 15: 2/16/2024

The least square problem revisited:

We previously covered a method using the QR decomposition of  $\mathbf{A}$ . But that required  $\mathbf{A}$  to be full rank. Now, here is a method using an SVD of  $\mathbf{A}$  which works even if  $\mathbf{A}$  is not full rank.

Let  $\mathbf{A} \in \mathbb{R}^{n \times m}$  where  $n \geq m$ .

Suppose  $\mathbf{A}$  has an SVD:  $\mathbf{U}\Sigma\mathbf{V}^T$ . Then:

$$(\|\vec{b} - \mathbf{A}\vec{x}\|_2)^2 = (\|\vec{b} - \mathbf{U}\Sigma\mathbf{V}^T\vec{x}\|_2)^2 = (\|\mathbf{U}^T\vec{b} - \Sigma\mathbf{V}^T\vec{x}\|_2)^2$$

As a reminder,  $\mathbf{U}$  is an orthogonal matrix and so  $\mathbf{U}\vec{x} \cdot \mathbf{U}\vec{y} = \vec{x} \cdot \vec{y}$ .

Thus, we get that:  $(\|\vec{b} - \mathbf{A}\vec{x}\|_2)^2 = \left( \left\| \begin{bmatrix} \vec{c} \\ \vec{d} \end{bmatrix} - \begin{bmatrix} \hat{\Sigma}\vec{y} \\ \vec{0} \end{bmatrix} \right\|_2 \right)^2$

...where  $\vec{y}$  is the vector containing the first  $r$  elements of  $\mathbf{V}^T\vec{x}$

with  $r = \text{rank}(\mathbf{A})$  and  $\mathbf{U}^T\vec{b} = \begin{bmatrix} \vec{c} \\ \vec{d} \end{bmatrix}$  where  $\vec{c} \in \mathbb{R}^r$ . Also,

$$\Sigma = \begin{bmatrix} \hat{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \text{ such that } \hat{\Sigma} \in \mathbb{R}^{r \times r}.$$

This can be rewritten as  $(\|\vec{b} - \mathbf{A}\vec{x}\|_2)^2 = (\|\vec{c} - \hat{\Sigma}\vec{y}\|_2)^2 + (\|\vec{d}\|_2)^2$ .  $\vec{d}$  is independent from  $\vec{x}$ . So we need only worry about minimizing  $(\|\vec{c} - \hat{\Sigma}\vec{y}\|_2)^2$  in order to minimize  $(\|\vec{b} - \mathbf{A}\vec{x}\|_2)^2$ . Luckily,  $\hat{\Sigma}$  is invertible since it is a diagonal matrix with only positive entries. So, we can solve for  $\vec{y}$  such that  $\vec{c} - \hat{\Sigma}\vec{y} = \vec{0}$

Therefore, any solution of the least squares problem of  $\mathbf{A}$  and  $\vec{b}$  satisfies that:  $\mathbf{V}^T\vec{x} = \begin{bmatrix} \vec{y} \\ \vec{z} \end{bmatrix}$  where  $\vec{z} \in \mathbb{R}^{m-r}$  is any vector.

Because  $\mathbf{V}^T$  is orthogonal, we can thus say that  $\vec{x} = \mathbf{V} \begin{bmatrix} \vec{y} \\ \vec{z} \end{bmatrix}$ . So the set of solutions to the least square problem is:

$$\left\{ \mathbf{V} \begin{bmatrix} \vec{y} \\ \vec{z} \end{bmatrix} \mid \vec{z} \in \mathbb{R}^{m-r} \right\}$$

Note: One important solution is the solution that minimizes  $\|\vec{x}\|_2$ . To get this, note that  $\|\mathbf{V}\vec{v}\|_2 = \|\vec{v}\|_2$  because  $\mathbf{V}$  is orthogonal. So, the solution to the least square problem that minimizes  $\|\vec{x}\|_2$  is  $\vec{x} = \mathbf{V} \begin{bmatrix} \vec{y} \\ \vec{0} \end{bmatrix}$ .

If  $\mathbf{A}$  is not a square matrix, there is no inverse. A pseudo inverse is a generalization that works for all matrices.

Note that if  $\mathbf{A}$  has the SVD:  $\mathbf{U}\Sigma\mathbf{V}^T$ , and is invertible, then  $\mathbf{A}^{-1} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T$  where:

$$\Sigma^{-1} = \begin{bmatrix} \frac{1}{\sigma_1} & & \\ & \ddots & \\ & & \frac{1}{\sigma_n} \end{bmatrix}$$

Importantly, no matter how  $\mathbf{A}$  isn't invertible (whether that be because it is rank deficient or non-square), the expression  $\mathbf{U}\Sigma\mathbf{V}^T$  is mostly well defined. The only problem that arises if  $\mathbf{A}$  is not invertible is that  $\Sigma$  is not invertible anymore either.

$$\text{Given } \Sigma = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{bmatrix} \in \mathbb{R}^{n \times m}, \text{ define } \Sigma^+ = \begin{bmatrix} \frac{1}{\sigma_1} & & \\ & \ddots & \\ & & \frac{1}{\sigma_r} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

Then, define  $\mathbf{A}^+ = \mathbf{V}\Sigma^+\mathbf{U}^T \in \mathbb{R}^{m \times n}$ .

Note that if  $\mathbf{A}$  is invertible, then  $\mathbf{A}^+ = \mathbf{A}^{-1}$ .

Theorem: If  $\mathbf{A} \in \mathbb{R}^{n \times m}$  where  $n \geq m$ , then the solution  $\vec{x}$  to the least square problem of  $\mathbf{A}$  and  $\vec{b}$  such that  $\|\vec{x}\|_2$  is minimized is given by  $\vec{x} = \mathbf{A}^+ \vec{b}$ .

Why: (recall how  $\vec{c}$ ,  $\vec{d}$ , and  $\vec{y}$  were defined on the previous page...)

- $\mathbf{U}^T \vec{b} = \begin{bmatrix} \vec{c} \\ \vec{d} \end{bmatrix}$
- $\Sigma^+ \begin{bmatrix} \vec{c} \\ \vec{d} \end{bmatrix} = \begin{bmatrix} \vec{y} \\ \vec{0} \end{bmatrix}$

## Lecture 16: 2/21/2024