

Chapter 4

The major topics in Numerical Linear Algebra

1) Solution of the system $Ax=b$

- Direct methods
- Iterative methods

2) The Algebraic Eigenvalue problem

$$Av = \lambda v$$

3) Least-Squares problem.

Right inverse / Left inverse, inverse.

Suppose A is $m \times n$ and B is $n \times m \Rightarrow AB$ is $m \times m$ with $AB = I$. Then

B is called a right inverse of A

A is called a left inverse of B .

Ex.
$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ \alpha & \beta \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$
 holds for any α, β

$A_{2 \times 3} \quad B_{3 \times 2} \quad I_{2 \times 2}$

Note, $AB = I$ is possible only if $n \geq m$.

proof, suppose $n < m$

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \end{bmatrix} = \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix}$$

~~Then B has~~

Then # of rows of $B = n < m =$ # of columns of B

We know that there exists $x \in \mathbb{R}^m$ such that $x \neq 0$

$Bx = 0$. Then $ABx = Ix = x \neq 0$. ~~*~~

$$= A(\underbrace{Bx}_{=0}) = 0$$

$$\begin{bmatrix} A_1 & | & A_n \end{bmatrix} \begin{bmatrix} B_{kj} \\ | \\ B_{kj} \end{bmatrix} = B_{kj}$$

$$(x^T A) B = x$$

$$y^T B = x$$

$$\sum_{j=1}^n y_j B_j$$

B_j j -th row of

Remark This example shows that the right inverse is not unique, in general. The same is true of the left inverse.

Theorem A square matrix can have at most one right inverse. Also, it can have only one left inverse.

proof. In this case $m=n$

$AB = I$. want to show B is unique.

For any $x \in \mathbb{R}^n$, $ABx = A(Bx) = Ix = x$. let $y = Bx$

Now $Ay = \sum_{j=1}^n y_j A_j$ where A_j is the j -th column of A .

This shows that every x can be written as a linear combination of the columns of $A \Rightarrow$ the columns of A are linearly independent and form a basis for \mathbb{R}^n .

Now, let B_j be the j -th col. of B

$$AB_j = (AB)_j = I_j = e_j = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \leftarrow j$$

$$\Rightarrow \sum_{k=1}^n B_{kj} A_k = e_j \Rightarrow B_{kj} \text{ are uniquely determined.}$$

The proof of the uniqueness of the left inverse of a square matrix is similar. (Look at $x^T AB = x^T I = x^T$)

Theorem If A and B are square matrices such that $AB = I$, then $BA = I$

proof let $C = BA - I + B$. Then

$$AC = A(BA - I + B) = \underbrace{ABA}_I - A + \underbrace{AB}_I = A - A + I = I$$

$$\Rightarrow B = C \text{ by previous result}$$

$$= BA - I + B \Rightarrow BA = I \quad \checkmark$$

Remark This shows that if a square matrix has a one-sided inverse, then it has a two-sided inverse, or simply inverse, henceforth denoted by A^{-1} .

Elementary row operations / Elementary matrices

1. Interchange rows i and j
2. Multiply row i by a non zero number α
3. Multiply row j by α and add that to row i

we will show that if we apply a sequence of elementary operations to a system $Ax=b$, then the resulting system is equivalent to $Ax=b$.

Elementary matrices

To each type of elem. row operation we associate a square elementary matrix of the same type as follows

Elementary row op. \rightarrow corresponding Elementary matrix.

Ex. $n=3$. Interchange rows 1 and 3

$$\text{Elem. matrix} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \xleftarrow[\text{rows}]{\substack{\text{int.} \\ 1 \leftrightarrow 3}} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Ex. Add -2 row 3 to row 2

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & -2 \\ 0 & 0 & 1 \end{bmatrix} \longleftarrow \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Theorem. All 3 types of elem. matrices are invertible. Also

(i) Let M be an elementary matrix of type I

$$\text{Then } M^{-1} = M$$

(ii) Let M be an elementary matrix of type II corresponding to $\alpha \neq 0$. Then M^{-1} is the elementary matrix of type II with $\frac{1}{\alpha}$

(iii) Let M be an elementary matrix of type III corresp. to

$\alpha r_j + r_i \rightarrow r_i$
 M^{-1} is the elem. matrix of same type with $-\alpha r_j + r_i$

Lemma Consider the linear system $Ax = b$. (A $n \times n$)

Suppose M is any $n \times n$ invertible matrix. Then the system $(MA)x = Mb$ are equivalent.

proof

Suppose x satisfies $Ax = b$. Then multiply both sides by $M \Rightarrow M(Ax) = Mb$

$$(MA)x = Mb \quad (\text{associativity}).$$

Conversely, suppose x satisfies $(MA)x = Mb$. Then

Multiply both sides by $M^{-1} \Rightarrow M^{-1}((MA)x) = M^{-1}(Mb)$

$$\underbrace{(M^{-1}M)}_I Ax = \underbrace{M^{-1}M}_I b = b$$

Reduced row echelon form (rref)

Let A be an $m \times n$ matrix. zero row is a row which has only zeros.
nonzero " " a " " has at least one nonzero

We say that an $m \times n$ matrix is in rref if

- 1) All zero rows, if any, are stacked at the bottom.
- 2) The first nonzero of a nonzero row is 1 (leading one)
In a col. that contains a leading 1, all other elements are zero
- 3) The leading 1 of a nonzero row is to the right of the leading one, if any, of the preceding row.

Theorem Every $m \times n$ matrix can be reduced to rref by elementary row operations. Furthermore, the rref is the same independently of which sequence of elem. row ops, used. In other words the rref of a matrix is unique.

Defn. The rank of an $m \times n$ matrix is the number of nonzero rows in its rref.

Theorem. For an $m \times n$ matrix A

1. $\text{rank}(A) =$ number of linearly ind. rows of A .
 2. $\text{rank}(A) =$ " " " " columns of A .
- as a consequence $\text{rank}(A) \leq \min(m, n)$.
3. A square matrix is invertible iff $\text{rank}(A) = n$.
 4. $\text{rank}(A) = n$ iff $Ax = 0$ has only the trivial soln. $x = 0$.

Theorem For an $n \times n$ matrix the following are equivalent

1. A^{-1} exists, we say that A is nonsingular or invertible
2. $\det(A) \neq 0$
3. The rows of A form a basis of \mathbb{R}^n
4. The columns of A " " " "
5. As a map from $\mathbb{R}^n \rightarrow \mathbb{R}^n$ A is injective (one-to-one)
i.e. $Ax = Ay \Rightarrow x = y$
6. As a map from $\mathbb{R}^n \rightarrow \mathbb{R}^n$ A is surjective (onto)
i.e. Given any $b \in \mathbb{R}^n$, $\exists x$ such that $Ax = b$
7. $Ax = 0 \Rightarrow x = 0$ (The converse is always true!)
8. A is the product of elementary matrices
9. 0 is not an eigenvalue of A
10. The rref of A is the identity matrix

Gauss-Elimination: $-4.8-$
 $[A|b] \rightarrow [U|b']$

$A = LU$ and $PA = LU$ factorizations

Given a system $Ax = b$, we wish to solve it by reducing the system into an equivalent system which is "easier" to solve.

Consider the system $Ax = b$ where A is an $n \times n$ invertible matrix, hence has a unique solution x .

Suppose A can be factored as the product $A = LU$ where

- (i) L is lower triangular with $l_{ii} = 1, i = 1, \dots, n$
 - (ii) U is upper triangular.
- we will see later how to compute L and U

Ex. A invertible $\Rightarrow u_{ii} \neq 0, i = 1, \dots, n$

Then $Ax = b$ can be written as $LUx = b$
Step 0 Factor $A = LU$. (will do later) Cost $\frac{2}{3}n^3 + cn^2$

Step 1 we let $y = Ux$. Then $Ly = b$

$$\begin{bmatrix} l_{11} & 0 & \dots & 0 \\ l_{21} & l_{22} & 0 & \dots & 0 \\ & & \ddots & & \\ l_{n1} & l_{n2} & \dots & l_{nn} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$$

This system can be easily solved by Forward substitution

$y_1 = b_1 / l_{11}$ way or may not ignore the fact that $l_{ii} = 1$

For $i = 2, \dots, n$

$$y_i = \{b_i - \sum_{j=1}^{i-1} l_{ij} y_j\} / l_{ii}$$

Cost $\approx n(n-1)$ if $l_{ii} = 1$

Step 2 Once y has been calculated, x can be calculated from $Ux=y$ by Back substitution

$$\begin{bmatrix} u_{11} & u_{12} & \dots & u_{1n} \\ & u_{22} & \dots & u_{2n} \\ & & \ddots & | \\ & & & u_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ | \\ x_n \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ | \\ y_n \end{bmatrix}$$

$$x_n = y_n / u_{nn}$$

For $i = n-1, \dots, 1$

$$x_i = (y_i - \sum_{j=i+1}^n u_{ij} x_j) / u_{ii}$$

cost: n^2 ops.

What is the advantage of LU factorization over applying Gauss Elimination to the augmented system $[A|b]$?

Suppose we want to solve systems $Ax=b$ with the same A but with different right-hand sides. Using Gauss Elimination, we have to repeat the costly part of reducing A to upper triangular form for every b .

We next consider the issue of factoring A into the product LU .

Defn. Let $A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ | & | & \ddots & | \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}$ be an $n \times n$ matrix.

For $k=1, \dots, n$, the k th leading principal minor of A is the $k \times k$ submatrix of A formed by the intersection of the first k rows of A and the first k columns of A .

Note that the $n \times n$ leading principal minor of A is A itself. Also, in some texts, the minor is actually the determinant of what we call here the minor.

Theorem. If all the leading principal minors of A are nonsingular, then A has an LU factorization. Furthermore, if $l_{ii} = 1$, then the factors L and U are unique.

Proof. Let $k=1$ the minor is $[a_{11}] \neq 0$ by assumption. So we can use it to introduce zeros in rows $2, \dots, n$ of the first column.

$$\begin{array}{c}
 \left[\begin{array}{ccc|ccc}
 a_{11} & a_{12} & \dots & a_{1n} & & \\
 a_{21} & a_{22} & \dots & a_{2n} & & \\
 \vdots & \vdots & & \vdots & & \\
 a_{n1} & a_{n2} & \dots & a_{nn} & &
 \end{array} \right] \rightarrow \left[\begin{array}{ccc|ccc}
 a_{11} & a_{12} & \dots & a_{1n} & & \\
 0 & a_{22}^{(2)} & \dots & a_{2n}^{(2)} & & \\
 \vdots & a_{32}^{(2)} & & \vdots & & \\
 \vdots & \vdots & & \vdots & & \\
 0 & a_{n2}^{(2)} & \dots & a_{nn}^{(2)} & &
 \end{array} \right] \\
 \\
 \begin{array}{c}
 A (= A^{(1)}) \xrightarrow{\quad} A^{(2)} \\
 \\
 A^{(2)} = M^{(1)} A^{(1)} \quad M^{(1)} = \left[\begin{array}{ccc|ccc}
 1 & 0 & \dots & 0 & & \\
 -\frac{a_{21}}{a_{11}} & 1 & & 0 & & \\
 \vdots & & & \vdots & & \\
 \vdots & & & \vdots & & \\
 -\frac{a_{n1}}{a_{11}} & & & 0 & & 1
 \end{array} \right]
 \end{array}
 \end{array}$$

We now show that $a_{22}^{(2)} \neq 0$. Indeed, we will show that all the leading principal minors of $A^{(2)}$ are invertible

$$A^{(2)} = M^{(1)} A^{(1)} = \begin{array}{c} \begin{array}{|cc|} \hline k & n-k \\ \hline M_{11}^{(1)} & 0 \\ \hline \end{array} \begin{array}{|cc|} \hline k & n-k \\ \hline A_{11}^{(1)} & A_{12}^{(1)} \\ \hline \end{array} \\ \begin{array}{|cc|} \hline n-k & n-k \\ \hline M_{21}^{(1)} & M_{22}^{(1)} \\ \hline \end{array} \begin{array}{|cc|} \hline k & n-k \\ \hline A_{21}^{(1)} & A_{22}^{(1)} \\ \hline \end{array} \end{array}$$

$$= \begin{array}{c} \begin{array}{|cc|} \hline k & n-k \\ \hline M_{11}^{(1)} A_{11}^{(1)} & M_{11}^{(1)} A_{12}^{(1)} \\ \hline \end{array} \\ \begin{array}{|cc|} \hline n-k & n-k \\ \hline M_{21}^{(1)} A_{11}^{(1)} & M_{21}^{(1)} A_{12}^{(1)} \\ \hline \end{array} \\ \begin{array}{|cc|} \hline k & n-k \\ \hline M_{22}^{(1)} A_{21}^{(1)} & M_{22}^{(1)} A_{22}^{(1)} \\ \hline \end{array} \end{array}$$

$$M_{11}^{(1)} = \begin{bmatrix} 1 & 0 \\ -\frac{a_{21}^{(1)}}{a_{11}^{(1)}} & 0 \\ -\frac{a_{k1}^{(1)}}{a_{11}^{(1)}} & 1 \end{bmatrix} \text{ is invertible for any } k,$$

$A_{11}^{(1)}$ is invertible, since A has that property
 \Rightarrow the product $M_{11}^{(1)} A_{11}^{(1)}$, which is the leading principal $k \times k$ submatrix of $A^{(2)}$ is also invertible.

In particular, for $k=2$, $\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} \\ 0 & a_{22}^{(2)} \end{bmatrix}$
 is invertible $\Leftrightarrow \det \neq 0$

$$a_{11}^{(1)} a_{22}^{(2)} \neq 0 \Rightarrow a_{22}^{(2)} \neq 0.$$

So, $a_{22}^{(2)}$ can be used to eliminate the elements below itself

-4.12

$$\Rightarrow A^{(3)} = M^{(2)} A^{(2)}$$

$$M^{(2)} = \begin{bmatrix} 1 & & & & & \\ & 1 & & & & \\ & -\frac{a_{21}^{(2)}}{a_{11}^{(2)}} & & & & \\ & & 1 & & & \\ & & & \ddots & & \\ & & & & 1 & \\ & -\frac{a_{n2}^{(2)}}{a_{22}^{(2)}} & & & & \\ & & & & & & 1 \end{bmatrix}, \quad A^{(3)} = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & & a_{2n}^{(2)} \\ \vdots & 0 & a_{33}^{(3)} & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & a_{n3}^{(3)} & a_{nn}^{(3)} \end{bmatrix}$$

As done before, we can show that all the leading principal minors of $A^{(3)}$ are invertible. In particular $a_{33}^{(3)} \neq 0$

The process can be continued:

$$A^{(m+1)} = M^{(m)} A^{(m)}, \quad m=1, \dots, n-1$$

until we arrive at

$$A^{(n)} = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ & a_{22}^{(2)} & & a_{2n}^{(2)} \\ & & \ddots & \vdots \\ 0 & & & a_{nn}^{(n)} \end{bmatrix} \equiv U$$

$$\begin{aligned} \Rightarrow U = A^{(n)} &= M^{(n-1)} A^{(n-1)} = M^{(n-1)} M^{(n-2)} A^{(n-2)} \\ &= \dots = \underbrace{M^{(n-1)} M^{(n-2)} \dots M^{(1)}}_M A^{(1)} \\ &\equiv MA \end{aligned}$$

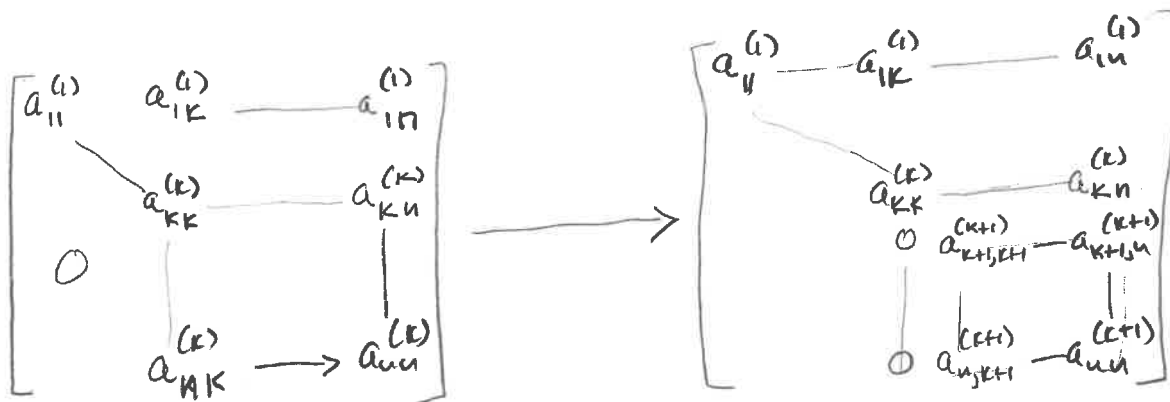
$$\Rightarrow \underbrace{\tilde{L}^{-1} L}_{\text{lower}} = \underbrace{\tilde{U} U^{-1}}_{\text{upper}} \Rightarrow \text{both are diagonal and diag are } = 1$$

with diag = 1

i.e. $\tilde{L}^{-1} L = \tilde{U} U^{-1} = I \Rightarrow \tilde{L} = L, \tilde{U} = U \checkmark$

work estimate for $A = LU$

Consider the step $A^{(k)} \rightarrow A^{(k+1)}$



$$a_{ij}^{(k+1)} = a_{ij}^{(k)} - \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} \cdot a_{kj}^{(k)}, \quad \begin{matrix} i = k+1, \dots, n \\ j = k+1, \dots, n \end{matrix}$$

$n-k$ multipliers $-\frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} \quad i = k+1, \dots, n$

Then for each pair $(i,j) \quad 2(n-k)^2$ ops.

Hence, cost of $A^{(k)} \rightarrow A^{(k+1)}$ is $(n-k) + 2(n-k)^2$ ops.

Total for LU factorization is

$$\sum_{k=1}^{n-1} [(n-k) + 2(n-k)^2] = \sum_{k=1}^{n-1} (n-k) + 2 \sum_{k=1}^{n-1} (n-k)^2$$

$$\sum_{k=1}^{n-1} (n-k) = 1 + 2 + \dots + (n-1) = \frac{(n-1)n}{2}$$

$$\sum_{k=1}^{n-1} (n-k)^2 = 1^2 + 2^2 + \dots + (n-1)^2 = \frac{(n-1)n(2n-1)}{6}$$

$$\begin{aligned} & \text{---4.15=} \\ \Rightarrow \text{ Total cost of } A=LU \text{ is } & \frac{(n-1)n}{2} + \frac{(n-1)n(2n-1)}{3} \\ & = \frac{(n-1)n}{6} \left[\frac{3 + 4n - 2}{4n+1} \right] = \frac{2}{3}n^3 - \frac{n^2}{2} - \frac{n}{6} = \boxed{\frac{2}{3}n^3 + O(n^2)}. \end{aligned}$$

Algorithm: $A=LU$

for $k=1, \dots, n-1$

 for $i=k+1, \dots, n$

$a_{ik} = a_{ik}/a_{kk}$ new column of L

 for $j=k+1, \dots, n$

$a_{ij} = a_{ij} - a_{ik}a_{kj}$

 }

 }

At the end, A contains
 L without the diagonals ($=1$) in its strictly lower
 triangular part
 U in its upper triangular part.

Theorem $A = LDL^T$ factorization. Assume as in the

$A = LU$ Theorem that all the leading principal minors of A are invertible. Assume in addition that A is symmetric. Then there exist a lower triangular matrix with unit diagonal elements and a diagonal matrix D with nonzero diagonal elements such that

$$A = LDL^T.$$

Moreover, L and D are unique.

proof.

As in $A = LU$,

$$\tilde{A}^{(2)} \equiv M_1 A^{(1)} = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{n2}^{(2)} & \dots & a_{nn}^{(2)} \end{bmatrix}$$

We can use $a_{11}^{(1)}$ to eliminate $a_{12}^{(1)} \dots a_{1n}^{(1)}$. Actually this can be done by the same multipliers in M_1 . Indeed,

$$M_1 A^{(1)} M_1^T = \begin{bmatrix} a_{11}^{(1)} & 0 & \dots & 0 \\ 0 & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{n2}^{(2)} & \dots & a_{nn}^{(2)} \end{bmatrix} \equiv A^{(2)}$$

First, the elements in the $(n-1) \times (n-1)$ submatrix $B^{(2)} = \begin{bmatrix} a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ \vdots & \ddots & \vdots \\ a_{n2}^{(2)} & \dots & a_{nn}^{(2)} \end{bmatrix}$ are not changed during the column elimination since the elements in rows 2, ..., n of column 1 of $\tilde{A}^{(2)}$ are all zero.

Furthermore, $A^{(2)} \equiv M_1 A^{(1)} M_1^T$ is symmetric. Indeed

$$(A^{(2)})^T = (M_1 A^{(1)} M_1^T)^T = (M_1^T)^T (A^{(1)})^T M_1^T = M_1 A^{(1)} M_1^T = A^{(2)}$$

Since $A^{(1)}$ is symmetric.

4/7-

Also, it is easy to show that all the leading principal minors of $A^{(2)}$ are invertible $\Rightarrow a_{22}^{(2)} \neq 0$.

The row and column elimination process can be continued without fail \Rightarrow

$$A^{(2)} = M_1 A^{(1)} M_1^T$$

$$A^{(3)} = M_2 A^{(2)} M_2^T$$

|

$$D \equiv A^{(n)} = M_{n-1} A^{(n-1)} M_{n-1}^T$$

$$= M_{n-1} M_{n-2} A^{(n-2)} M_{n-2}^T M_{n-1}^T$$

$$= \underbrace{M_{n-1} \dots M_1}_{=A} A^{(1)} \underbrace{M_1^T \dots M_{n-1}^T}_{=A^T}$$

$$\Rightarrow A = M_1^{-1} \dots M_{n-1}^{-1} D (M_{n-1}^T)^{-1} \dots (M_1^T)^{-1}$$

$$= M_1^{-1} \dots M_{n-1}^{-1} D (M_{n-1}^{-1})^T \dots (M_1^{-1})^T$$

$$= \underbrace{M_1^{-1} \dots M_{n-1}^{-1}}_L D \underbrace{(M_1^{-1})^T \dots (M_{n-1}^{-1})^T}_L$$

M_i : lower tri with unit diagonals $i=1, \dots, n-1$

$M_i^{-1} \Rightarrow$ " " " " " " $i=1, \dots, n-1$

$\Rightarrow M_1^{-1} \dots M_{n-1}^{-1}$ is lower triangular with unit elements.

Important Remark. The submatrix $B^{(2)}$ is not affected by the column elimination step, hence is symmetric.

So in practice, we don't need to do column elimination at all.

Also, during the row elimination process, that precedes the formation of $B^{(2)}$, we only create the lower triangular part of $B^{(2)}$ and all the subsequent submatrices.

Algorithm $A = LDL^T$

For $k=1, \dots, n-1$

$$d_k = a_{kk}$$

for $i=k+1, \dots, n$

$$l_{ik} = a_{ik} / a_{kk}$$

for $j=k+1, \dots, i$

$$a_{ij} = a_{ij} - l_{ik} a_{jk}$$

by symmetry $a_{kj} = a_{jk}$

}

}

$$d_k = a_{k+1,k+1}$$

}

-4.19-

Direct Factorization methods: Doolittle, Crout, Choleski

Ex.

$$A = \begin{bmatrix} 6 & 3 & 2 \\ 3 & 2 & 3/2 \\ 2 & 3/2 & 6/5 \end{bmatrix}$$

All the leading principal minors are invertible, hence there is an LU factorization

we write

$$\begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{bmatrix} = \begin{bmatrix} 6 & 3 & 2 \\ 3 & 2 & 3/2 \\ 2 & 3/2 & 6/5 \end{bmatrix}$$

$L \quad U = A$

$$6 = a_{11} = \text{row 1 of } L \cdot \text{col. 1 of } U = (1, 0, 0) \begin{pmatrix} u_{11} \\ 0 \\ 0 \end{pmatrix} = u_{11} \Rightarrow u_{11} = 6$$

$$3 = a_{12} = \text{row 1 of } L \cdot \text{col. 2 of } U = (1, 0, 0) \begin{pmatrix} u_{12} \\ u_{22} \\ 0 \end{pmatrix} = u_{12} \Rightarrow u_{12} = 3$$

$$2 = a_{13} = \text{row 1 of } L \cdot \text{col. 3 of } U = (1, 0, 0) \begin{pmatrix} u_{13} \\ u_{23} \\ u_{33} \end{pmatrix} = u_{13} \Rightarrow u_{13} = 2$$

$$3 = a_{23} = \text{row 2 of } L \cdot \text{col. 1 of } U = (l_{21}, 1, 0) \begin{pmatrix} u_{11} \\ 0 \\ 0 \end{pmatrix} = l_{21} u_{11}$$
$$\Rightarrow l_{21} = 3/u_{11} = 1/2$$

$$2 = a_{22} = \text{row 2 of } L \cdot \text{col. 2 of } U = (l_{21}, 1, 0) \begin{pmatrix} u_{12} \\ u_{22} \\ 0 \end{pmatrix} = l_{21} u_{12} + u_{22}$$
$$\Rightarrow u_{22} = 2 - l_{21} u_{12} = 2 - \frac{1}{2} \cdot 3 = \frac{1}{2} = u_{22}$$

$$3/2 = a_{23} = \text{row 2 of } L \cdot \text{col. 3 of } U = (l_{21}, 1, 0) \begin{pmatrix} u_{13} \\ u_{23} \\ u_{33} \end{pmatrix} = l_{21} u_{13} + u_{23}$$

$$\Rightarrow u_{23} = 3/2 - l_{21} u_{13} = 3/2 - \frac{1}{2} \cdot 2 = 1/2 = u_{23}$$

$$2 = a_{31} = \text{row 3 of } L \cdot \text{col. 1 of } U = (l_{31}, l_{32}, 1) \begin{pmatrix} u_{11} \\ 0 \\ 0 \end{pmatrix} = l_{31} u_{11}$$

$$\Rightarrow l_{31} = 2/u_{11} = 1/3$$

$$a_{32} = 3/2 = \text{row 3 of } L \cdot \text{col. 2 of } U = (l_{31}, l_{32}, 1) \begin{pmatrix} u_{12} \\ u_{22} \\ 0 \end{pmatrix} = l_{31}u_{12} + l_{32}u_{22}$$

$$\Rightarrow l_{32} = (3/2 - l_{31}u_{12})/u_{22} = (3/2 - 1/3 \cdot 3)/1/2 = \boxed{1 = l_{32}}$$

$$a_{33} = 6/5 = \text{row 3 of } L \cdot \text{col. 3 of } U = (l_{31}, l_{32}, 1) \begin{pmatrix} u_{13} \\ u_{23} \\ u_{33} \end{pmatrix} = l_{31}u_{13} + l_{32}u_{23} + u_{33}$$

$$\Rightarrow u_{33} = 6/5 - l_{31}u_{13} - l_{32}u_{23} = 6/5 - 1/3 \cdot 2 - 1 \cdot 1 = \frac{1}{30}$$

This is Doolittle's algorithm: Designed so that the diagonals of L are equal to 1.

Croft Algorithm is similar but we now want U to have ones along its diagonal.

Obviously these two algorithms give the same result as Gaussian elimination. However they offer a great variety in the sequence of operations. These can be attractive from the point of view of parallelism and other techniques to reduce CPU time.

One disadvantage of these two algorithms is the requirement of all the leading principal minors to be invertible, which is a strong condition. They can be modified to include row or column interchanges but the algorithms become more complicated.

One important variant is Choleski's algorithm which is widely used when it is known beforehand that the matrix is symmetric and positive definite.

Def. A real symmetric matrix is said to be positive definite if $x^T A x > 0 \quad \forall x \neq 0.$

Ex. $A = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$

Theorem let A be a real symmetric matrix.
The following statements are equivalent.

- (a) A is positive definite.
- (b) All leading principal minors have positive determinants.
- (c) All eigenvalues are positive.
- (d) A can be factored into $A = LL^T$ where L is a lower triangular matrix with nonzero diagonal elements. In particular, there is a factorization with $L_{ii} > 0$, in which case we have uniqueness.
- (e) $A = LDL^T$ with $L_{ii} = 1$ and $d_{ii} > 0$.

proof. (only (a) \Leftrightarrow (d)) we will assume (a) \Leftrightarrow (b) \Leftrightarrow (c).

(d) \Rightarrow (a). Suppose $A = LL^T$, $L_{ii} \neq 0$. For any $x \neq 0$,

$$x^T A x = x^T L L^T x = (L^T x)^T (L^T x) = y^T y \text{ with } y = L^T x$$

Now L^T is an upper triangular matrix with nonzero diagonals, hence it is invertible $\Rightarrow y \neq 0$ since $x \neq 0$.

Hence

$$y^T y = \sum_{i=1}^n y_i^2 > 0 \quad \checkmark$$

(a) \Rightarrow (d), (e) A is symmetric, and all its leading principal minors are invertible (have pos. determinant). Hence, A has an LDLT factorization $A = LDL^T$. ($L_{ii} = 1$)
We will show that $D_{ii} > 0, i = 1, \dots, n$.

Since L is invertible so is L^T . For each $i, i = 1, \dots, n$

$\exists y_i \neq 0$ such that

$$L^T y_i = e_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \leftarrow i$$

$$0 < y_i^T A y_i = y_i^T L D L^T y_i = (L^T y_i)^T D (L^T y_i) = e_i^T D e_i = D_{ii} \quad \checkmark$$

\uparrow
 A is spd, $y_i \neq 0$.

Since $D_{ii} > 0$, $\sqrt{D_{ii}}$ is well-defined, so let $\tilde{D} =$

$$D = \tilde{D}\tilde{D} = \tilde{D}\tilde{D}^T$$

$$\begin{bmatrix} \sqrt{D_{11}} & & & \\ & \sqrt{D_{22}} & & \\ & & \ddots & \\ & & & \sqrt{D_{nn}} \end{bmatrix}$$

From

$$A = LDL^T = L\tilde{D}\tilde{D}^TL^T = (L\tilde{D})(L\tilde{D})^T$$

(c) \Rightarrow (a) Exercise

(d) \Rightarrow (a) Exercise

Suppose $A = LL^T$, $L_{ii} > 0$. We will show uniqueness

Suppose $A = \tilde{L}\tilde{L}^T$, $\tilde{L}_{ii} > 0$. We will show $\tilde{L} = L$.

$$LL^T = \tilde{L}\tilde{L}^T \Rightarrow L = \tilde{L}\tilde{L}^TL^{-T} \Rightarrow \tilde{L}^{-1}L = \tilde{L}^TL^{-T}$$

i.e. $\tilde{L}^{-1}L = (L^{-1}\tilde{L})^T$

$\tilde{L}^{-1}L$ is lower triangular $(\tilde{L}^{-1}L)_{ii} = L_{ii}/\tilde{L}_{ii} > 0$

$L^{-1}\tilde{L}$ " " " with $(L^{-1}\tilde{L})_{ii} = \tilde{L}_{ii}/L_{ii} > 0$

\Downarrow $(L^{-1}\tilde{L})^T$ is upper triangular with $(L^{-1}\tilde{L})^T_{ii} = \tilde{L}_{ii}/L_{ii} > 0$

lower = upper \Rightarrow both diagonal

Also,

$$L_{ii}/\tilde{L}_{ii} = \tilde{L}_{ii}/L_{ii} \Rightarrow L_{ii}^2 = \tilde{L}_{ii}^2 \Rightarrow L_{ii} = \tilde{L}_{ii}$$

$$\Rightarrow \tilde{L}^{-1}L = I \Rightarrow \tilde{L} = L \quad \checkmark$$

Theorem Suppose A is a real, symmetric positive definite matrix. Then

(i) $a_{ii} > 0$

(ii) $a_{ij}^2 < a_{ii}a_{jj}$ for all i, j , $i \neq j$.

-4023-

Choleski's method consists on applying the direct factorization principle to $A = LL^T$, $L_{ii} > 0$

$$LL^T = A$$

$$\Rightarrow \begin{bmatrix} l_{11} & 0 & \dots & 0 \\ l_{21} & l_{22} & & 0 \\ \vdots & & \ddots & \\ l_{n1} & l_{n2} & \dots & l_{nn} \end{bmatrix} \begin{bmatrix} l_{11} & l_{21} & \dots & l_{n1} \\ 0 & l_{22} & & l_{n2} \\ \vdots & & \ddots & \\ 0 & \dots & 0 & l_{nn} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}$$

Algorithm (Choleski)

$$l_{11} = \sqrt{a_{11}}$$

For $j = 2, \dots, n$

$$\left. \begin{array}{l} l_{j1} = a_{j1} / l_{11} \\ \vdots \end{array} \right\}$$

for $i = 2, \dots, n-1$

$$l_{ii} = \left(a_{ii} - \sum_{k=1}^{i-1} l_{ik}^2 \right)^{1/2}$$

for $j = i+1, \dots, n$

$$l_{ji} = \left(a_{ji} - \sum_{k=1}^{i-1} l_{jk} l_{ik} \right) / l_{ii}$$

$\left. \begin{array}{l} \vdots \\ \vdots \end{array} \right\}$

$$l_{nn} = \left(a_{nn} - \sum_{k=1}^{n-1} l_{nk}^2 \right)^{1/2}$$

Pivoting strategies

The condition that allows for the $A=LU$, $A=LDL^T$ and LL^T factorizations, namely that all the leading principal minors be positive is quite restrictive. Indeed if we let

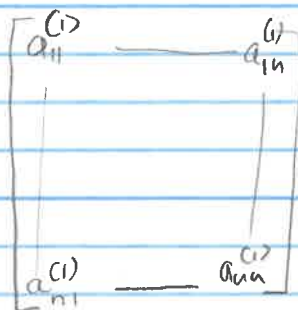
$$A = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$$

Then $A=LU$ cannot even "take off", so to speak. On the other hand this matrix is invertible and can be still factored if we consider more general factorizations. For instance

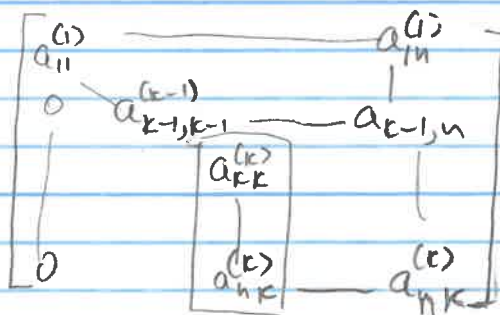
$$A^T = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \equiv LU, \quad L=I, \quad U=A^T.$$

In general, the $A=LU$ will fail whenever we encounter a zero pivot: $a_{kk}^{(k)} = 0$.

One popular way to continue the elimination process is to include row interchanges or pivoting. Indeed, if A is invertible, then it is easy to show that elimination can be concluded successfully if row interchanges are allowed.



if A is invertible, then there must be a nonzero in the first column, say $a_{p1}^{(1)} \neq 0$. we interchange rows 1 and p .



at the k -th stage, there must be a nonzero among

$$a_{kk}^{(k)}, \dots, a_{nk}^{(k)}.$$

-4.25-

There are other reasons, such as stability and control of roundoff errors to do row interchanges even if $a_{kk}^{(k)} \neq 0$.

Ex. Consider the linear system

$$\begin{bmatrix} .003000 & 59.14 \\ 5.291 & -6.130 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 59.17 \\ 46.78 \end{bmatrix}$$

which has exact solution $x=10, y=1$.

A has an LU factorization as can be verified.

However, suppose we do elimination using 4-decimal digit arithmetic with rounding.

$$m_{21} = -5.291 / .003 = -1763.\bar{6} \rightarrow -1764$$

$$\Rightarrow A^{(2)} = \begin{bmatrix} .003000 & 59.14 & | & 59.17 \\ 0 & fl(a_{22}) & | & fl(y^{(2)}) \end{bmatrix}$$

$$-6.130 + (-1764)(59.14) \\ = -104322.96 \rightarrow -104300$$

$$-6.130 - 104300 = -104306.130 \rightarrow -104300 = fl(a_{22}^{(2)})$$

$$46.78 - (1764)(59.17) \\ = 104375.88 \rightarrow 104400$$

$$46.78 - 104400 = -104353.22 \rightarrow -104400 = fl(y^{(2)})$$

$$fl(y^{(2)}) = \frac{-104400}{-104300} = 1.000958 \dots \rightarrow 1.001$$

$$fl(x) = fl\left(\frac{59.17 - (59.14)(1.001)}{.003}\right) = fl\left(\frac{-0.03}{.003}\right) = -10$$

Now we redo the elimination process after interchanging the two rows:

$$\begin{array}{ccc|ccc} 5.291 & -6.130 & 46.78 & 5.291 & -6.130 & 46.78 \\ .00300 & 59.14 & 59.17 & 0 & 59.14 & 59.14 \end{array} \rightarrow$$

$$fl(m_{21}) = fl\left(\frac{.003}{5.291}\right) = fl(5.67000567001 \times 10^{-4}) = -5.67 \times 10^{-4}$$

$$fl(59.14 + (-5.67 \times 10^{-4})(-6.130))$$

$$fl(5.67 \times 10^{-4} \times 6.130) = fl(34.7571 \times 10^{-4}) = 3.476 \times 10^{-3}$$

$$fl(59.14 + 3.476 \times 10^{-3}) = fl(59.14347571) = 59.14$$

$$fl(59.17 + (-5.67 \times 10^{-4}) \times (46.78))$$

$$fl(-5.67 \times 10^{-4} \times 46.78) = fl(.02652426) = -.02652$$

$$fl(59.17 - .02652) = fl(59.14348) = 59.14$$

$$\Rightarrow fl(y) = fl(59.14/59.14) = 1$$

$$\Rightarrow fl(x) = fl\left(\frac{46.78 + 6.130 \times 1}{5.291}\right)$$

$$fl(46.78 + 6.130) = fl(52.91) = 52.91$$

$$fl(x) = fl(52.91/5.291) = fl(10) = 10 \checkmark$$

One cause of the large errors seen in the previous calculation is the large multiplier m_{21} .

Theorem. Let A be an ^{invertible} $n \times n$ matrix. Then

There exists a permutation matrix P , a lower triangular matrix L with unit diagonal elements and an invertible upper triangular matrix U such that

$$PA = LU.$$

Furthermore L and U are unique.

proof.

$$\text{Let } A = A^{(1)} = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ a_{21}^{(1)} & a_{22}^{(1)} & \dots & a_{2n}^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1}^{(1)} & a_{n2}^{(1)} & \dots & a_{nn}^{(1)} \end{bmatrix}$$

Since A is invertible, there must be a nonzero element, say $a_{m1}^{(1)}$ in the first column.

Note if $a_{11}^{(1)}$ is nonzero, we may use it to eliminate

the elements below it in the first column. On the other hand, there might be other reasons, such as stability, control of roundoff errors etc. to make a row interchange.

$$P_{m_1} A^{(1)} = \begin{bmatrix} a_{m1}^{(1)} & a_{m2}^{(1)} & \dots & a_{mn}^{(1)} \\ \text{All other rows} \\ a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \end{bmatrix} \leftarrow m$$

$1 \leq m_1 \leq n$

are ex. before

-31-

$$\Rightarrow M_{n-1} \tilde{M}_{n-2} \tilde{M}_{n-3} \dots \tilde{M}_1 \underbrace{P_{n-1, n-1} \dots P_{1, 1}}_P A = U$$

} at most
n-2

$$\Rightarrow PA = \underbrace{\tilde{M}_1^{-1} \tilde{M}_2^{-1} \dots \tilde{M}_{n-2}^{-1} \tilde{M}_{n-1}^{-1}}_L U$$

uniqueness

Suppose $PA = LU = \tilde{L}\tilde{U}$.

we want to show $\tilde{L} = L$ and $\tilde{U} = U$.

$$LU = \tilde{L}\tilde{U} \Rightarrow \tilde{L}^{-1}L = \tilde{U}U^{-1}$$

$\tilde{L}^{-1}L$ is lower triangular with diagonal elements = 1

$\tilde{U}U^{-1}$ is upper triangular.

Now, $\tilde{L}^{-1}L = \tilde{U}U^{-1}$ implies "lower" = "upper"

\Rightarrow both must be diagonal. since $\tilde{L}^{-1}L$ has

diagonals = 1, $\tilde{L}^{-1}L = I \Rightarrow \tilde{L} = L$

$$\parallel$$

$$\tilde{U}U^{-1} = I \Rightarrow \tilde{U} = U. \quad \square$$

```

#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <ctype.h>
#include <math.h>

void factor(double *, int *, int);
void solve(double *, double *, int *, int);

int main() {
    double *b, *a;
    int i, dim, *piv;

    dim = 3;
    a = malloc(dim*dim*sizeof(double));
    b = malloc(dim*sizeof(double));
    piv = malloc(dim*sizeof(int));
    a[0] = 1.; a[1] = 5.; a[2] = 6.;
    a[3] = 2.; a[4] = 0.; a[5] = 4.;
    a[6] = 4.; a[7] = 2.; a[8] = 3.;
    b[0] = 33.; b[1] = 30.; b[2] = 21.;

    factor(a, piv, dim);
    solve(a, b, piv, dim);

    for(i=0;i<dim;i++) {
        printf(" piv[%d] = %d, b[%d] = %f\n", i, piv[i], i, b[i]);
    }
} /* end of main */

void factor(double *a, int *piv, int dim) {
    /******
    /* This function performs the PA=LU factorization of a matrix A
    /* The factors L and U are stored in A in the usual manner and
    /* piv returns the pivoting information
    /******
    int i, j, k, m;
    double t;

    for(i=0;i<dim;i++) piv[i] = i;
    for(k=0;k<dim-1;k++) {
        /* find pivot row for k-th column of matrix *a
        m=k;
        for(i=k+1;i<dim;i++) {
            if ( fabs(a[i*dim+k]) > fabs(a[m*dim+k]) ) m=i;
        }
        i = piv[k];
        piv[k] = piv[m];
        piv[m] = i;

        /* if m.ne.k, interchange rows k and m of L_k and U_k
        if(k!=m) {
            for(j=0;j<dim;j++) {
                t = a[m*dim+j];
                a[m*dim+j] = a[k*dim+j];
                a[k*dim+j] = t;
            }
        }
        /* compute multipliers. They will be stored in L_k
        for(i=k+1;i<dim;i++) {
            a[i*dim+k] = a[i*dim+k] / a[k*dim+k];
        }
        /* eliminate...
        for(i=k+1;i<dim;i++) {
            for(j=k+1;j<dim;j++) {
                a[i*dim+j] = a[i*dim+j] - a[i*dim+k] * a[k*dim+j];
            }
        }
        printf(" from factor, %d, %d, %d\n", piv[0], piv[1], piv[2]);
        /* end of function factor
    void solve(double *a, double *b, int *piv, int dim) {
        /******
        /* Computes the solution of LU x = P*F b. On output b contains x
        /* L and U are stored in the matrix a in the usual manner.
        /******
        int i, j, n;
        double sum, *temp;

        temp = malloc(dim*sizeof(double));
        for(i=0;i<dim;i++) temp[i] = b[piv[i]];

        /* forward solve
        for(i=1;i<dim;i++) {
            sum = 0.;
            for(j=0;j<i;j++) {
                sum += a[i*dim+j] * temp[j];
            }
            temp[i] = temp[i] - sum;
        }
        /* backward solve
        n = dim - 1;
        b[n] = temp[n] / a[n*dim+n];
        for(i=n-1;i>=0;i--) {
            sum = 0.;
            for(j=i+1;j<dim;j++) {
                sum += a[i*dim+j] * b[j];
            }
            b[i] = (temp[i] - sum) / a[i*dim+i];
        }
        free(temp);
    } /* end of function gen_solve
}

```

Vector and matrix norms

Defn. Let V be a vector space. A norm on V is a function, usually denoted by $\|\cdot\|$, $\|\cdot\|: V \rightarrow \mathbb{R}$ that has the following properties

- (i) $\|v\| \geq 0 \quad \forall v \in V$ and $\|v\| = 0 \Rightarrow v = 0$ "positivity"
- (ii) $\|\alpha v\| = |\alpha| \|v\|, \forall \alpha \in \mathbb{R}, \forall v \in V$ "homogeneity"
- (iii) $\|v+w\| \leq \|v\| + \|w\| \quad \forall v, w \in V$ "triangle inequality"

It follows from (iii) that $\|v-w\| \geq |\|v\| - \|w\||$
This is also called the triangle inequality.

The following are examples of norms on \mathbb{R}^n .

Ex. Euclidean norm

$$\|x\|_2 = \left(\sum_{i=1}^n |x_i|^2 \right)^{1/2} = \sqrt{x^T x}$$

Ex. l_1 -norm

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

Ex. l_∞ or sup or max norm

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

The above 3 norms are special cases ($p=2, 1, \infty$) of the l_p -norms: $1 \leq p < \infty$

$$\|x\|_p = \begin{cases} \left(\sum_{i=1}^n |x_i|^p \right)^{1/p} & 1 \leq p < \infty \\ \max_{1 \leq i \leq n} |x_i| & p = \infty \end{cases}$$

Exercise show that $\|\cdot\|_p$ satisfies properties

(i), (ii), (iii) above.

Norms are used to measure the "size" of vectors. They also provide a way to define "distance" between vectors. Indeed $\|v-w\|$ is viewed as the distance between vectors v and w .

Given a norm $\|\cdot\|$, a real number $p \geq 0$ and a vector v , the set

$$B_p(v) = \{w \in V, \|w-v\| \leq p\}$$

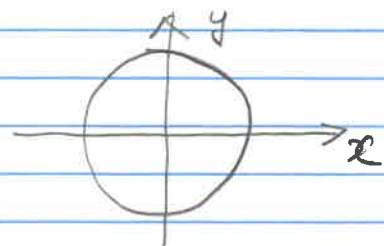
is the closed ball centered at v of radius p

i.e. it is the set of all vectors in V within a distance p of v . In particular

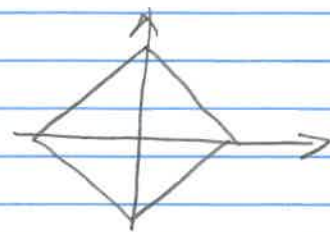
and $B_1(0) = \{v \in V, \|v\| \leq 1\}$ is the closed unit ball

$B_1(0) = \{v \in V, \|v\| = 1\}$ is the unit sphere

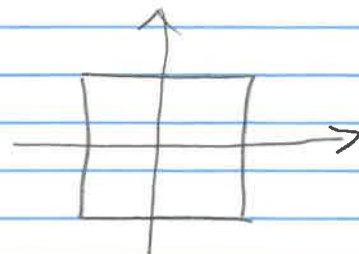
Ex. Unit sphere in \mathbb{R}^2
with respect to $\|\cdot\|_2$



Ex. Unit "sphere" in \mathbb{R}^2
with respect to $\|\cdot\|_1$



Ex. Unit "sphere" in \mathbb{R}^2
with respect to $\|\cdot\|_\infty$



Cauchy-Schwarz inequality in \mathbb{R}^n

$$|x^T y| \leq \|x\|_2 \|y\|_2 \quad \forall x, y \in \mathbb{R}^n.$$

This is a special case of Hölder's inequality

$$|x^T y| \leq \|x\|_p \|y\|_q \quad \frac{1}{p} + \frac{1}{q} = 1.$$

Ex. $p=2 \Rightarrow$ Cauchy-Schwarz,

Ex. $p=4 \Rightarrow |x^T y| \leq \|x\|_4 \|y\|_{4/3}$
 $q=4/3$ since $\frac{1}{4} + \frac{1}{4/3} = 1.$

Ex. $p=1, q=\infty$

$$|x^T y| \leq \|x\|_1 \|y\|_\infty.$$

Defn. A matrix norm on $\mathbb{R}^{n \times n}$ is a real-valued function $\|\cdot\|$ satisfying the following properties

- (i) $\|A\| \geq 0$ and $\|A\| = 0$ only if $A = 0$
- (ii) $\|\alpha A\| = |\alpha| \|A\|$, $\alpha \in \mathbb{R}$, $A \in \mathbb{R}^{n \times n}$
- (iii) $\|A+B\| \leq \|A\| + \|B\|$
- (iv) $\|AB\| \leq \|A\| \|B\|.$

Note that the first 3 properties are those of vector norms whereas (iv) is special to matrices and is useful in measuring the norm of the product of 2 or more matrices.

Induced or natural matrix norms

These are matrix norms that are derived from vector norms. So if $\|\cdot\|$ is a vector norm, we define the corresponding natural or induced matrix norm by

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|}.$$

Theorem The natural matrix norm defined above satisfies all 4 properties of a matrix norm proof.

(i) Since $\|Ax\| \geq 0$ and $\|x\| \geq 0$ as vector norms it follows that $\|A\| \geq 0$. Now if $\|A\| = 0$

$$0 = \|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} \Rightarrow \|Ax\| = 0 \quad \forall x \Rightarrow A = 0.$$

$$(ii) \quad \|\alpha A\| = \sup_{x \neq 0} \frac{\|(\alpha A)x\|}{\|x\|} = \sup_{x \neq 0} \frac{|\alpha| \|Ax\|}{\|x\|}$$

$$= |\alpha| \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = |\alpha| \|A\|.$$

$$(iii) \quad \|A+B\| = \sup_{x \neq 0} \frac{\|(A+B)x\|}{\|x\|} = \sup_{x \neq 0} \frac{\|Ax+Bx\|}{\|x\|}$$

$$\leq \sup_{x \neq 0} \frac{\|Ax\| + \|Bx\|}{\|x\|} \leftarrow (ii) \text{ of vector norms}$$

$$\leq \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} + \sup_{x \neq 0} \frac{\|Bx\|}{\|x\|}$$

$$= \|A\| + \|B\|.$$

$$(iv) \quad \|AB\| = \sup_{x \neq 0} \frac{\|(AB)x\|}{\|x\|} = \sup_{x \neq 0} \frac{\|A(Bx)\|}{\|x\|}$$

Now $\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} \geq \frac{\|Ax\|}{\|x\|}$ for any $x \neq 0$
and hence $\|Ax\| \leq \|A\| \|x\|$. Using this,

$$\|AB\| \leq \sup_{x \neq 0} \frac{\|A\| \|Bx\|}{\|x\|} = \|A\| \sup_{x \neq 0} \frac{\|Bx\|}{\|x\|} = \|A\| \|B\| \quad \checkmark$$

There are matrix norms which are not natural norms. An example of such is the

Frobenius norm

$$\|A\|_F = \left\{ \sum_{i,j=1}^n |a_{ij}|^2 \right\}^{1/2}.$$

Defn. Condition number of a matrix. Given a matrix norm $\|\cdot\|$

$$\kappa(A) \equiv \begin{cases} \|A\| \|A^{-1}\| & \text{if } A \text{ is invertible} \\ \infty & \text{if } A \text{ is singular.} \end{cases}$$

If $\|\cdot\|$ is a natural norm, then from $I = AA^{-1}$

$$1 = \|I\| = \|AA^{-1}\| \leq \|A\| \|A^{-1}\| = \kappa_{\|\cdot\|}(A).$$

If $\kappa_{\|\cdot\|}(A) \gg 1$, we say A is ill-conditioned.

Lemma Let A be square matrix and $\|\cdot\|$ any natural matrix norm. Then if λ is an eigenvalue of A , we have

$$|\lambda| \leq \|A\|$$

proof

$$Av = \lambda v, v \neq 0.$$

If $\lambda = 0$, then $|0| \leq \|A\|$ since $\|A\| \geq 0$.

If $\lambda \neq 0$, then $\|\lambda v\| = \|Av\| \Rightarrow |\lambda| \|v\| = \|Av\| \leq \|A\| \|v\|$

divide by $\|v\| \Rightarrow |\lambda| \leq \|A\|$ ✓

Defn. The spectrum of a square matrix is the set of all its eigenvalues. $\sigma(A)$.

Defn. The spectral radius of a matrix is

$$\rho(A) \equiv \max_{\lambda \in \sigma(A)} |\lambda|, \text{ i.e. largest in modulus of the eigenvalues of } A.$$

Lemma If λ is an eigenvalue of A , then $1-\lambda$ is an eigenvalue of $I-A$.

proof.

$$Av = \lambda v \Rightarrow (I-A)v = v - Av = v - \lambda v = (1-\lambda)v. \quad \square$$

Lemma Suppose $\|E\| < 1$ for some matrix E and an induced matrix norm $\|\cdot\|$. Then $I-E$ is invertible and

$$\|(I-E)^{-1}\| \leq \frac{1}{1-\|E\|}.$$

proof.

$\|E\| < 1 \Rightarrow |\lambda| < 1$ for all eigenvalues of A .
Now the eigenvalues of A are $1-\lambda$, and $|\lambda| < 1$
 $\Rightarrow 1-\lambda \neq 0$ hence $I-E$ is invertible.

To get the bounds, write (formally)

$$(I-E)^{-1} = I + E + E^2 + \dots + E^n + \dots$$

We will show that the matrix series is absolutely convergent.

$$\text{let } S_n = I + E + \dots + E^n$$

$$\|S_n\| = \|I + \dots + E^n\|$$

$$\leq \|I\| + \|E\| + \|E^2\| + \dots + \|E^n\|$$

$$\leq 1 + \|E\| + \|E\|^2 + \dots + \|E\|^n$$

$$= \frac{1 - \|E\|^{n+1}}{1 - \|E\|}$$

As $n \rightarrow \infty$ $\|E\|^{n+1} \rightarrow 0$ since $\|E\| < 1$. ✓

Thus the right hand side exists as a matrix.

Also

$$\|(I-E)^{-1}\| \leq \lim_{n \rightarrow \infty} \frac{\|S_n\|}{1 - \|E\|} = \frac{1}{1 - \|E\|} \quad \square$$

Remark The same result holds for $I+E$. In this

Case $(I+E)^{-1} = I - E + E^2 - E^3 + \dots$

Lemma Suppose A, E are two matrices such that A is invertible

$\|A^{-1}E\| < 1$ where $\|\cdot\|$ some natural matrix norm.

Then

$A+E$ is invertible and

$$\|(A+E)^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}E\|}$$

proof

write $A+E = A(I+A^{-1}E)$. Since $\|A^{-1}E\| < 1$ by the previous lemma $I+A^{-1}E$ is invertible and

$$\|(I+A^{-1}E)^{-1}\| \leq \frac{1}{1 - \|A^{-1}E\|}$$

Now

$$(A+E)^{-1} = (A(I+A^{-1}E))^{-1} = (I+A^{-1}E)^{-1}A^{-1}$$

$$\Rightarrow \|(A+E)^{-1}\| \leq \|(I+A^{-1}E)^{-1}A^{-1}\|$$

$$\leq \|(I+A^{-1}E)^{-1}\| \|A^{-1}\| \quad \text{property (iv)}$$

$$\leq \frac{\|A^{-1}\|}{1 - \|A^{-1}E\|}$$

Remark If $\|EA^{-1}\| < 1$, then $A \pm E$ are invertible and

$$\|(A \pm E)^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|EA^{-1}\|}$$

Note It is possible to have $\|E\| < 1$ in one norm and $\|E\| > 1$ in another.

ex. $A = \begin{bmatrix} 0 & \frac{1}{4} \\ \frac{1}{2} & 0.51 \end{bmatrix}$, $\|A\|_{\infty} = \max\{\frac{1}{4}, 1.01\} = 1.01 > 1$

$$\|A\|_1 = \max\{\frac{1}{2}, 0.78\} = 0.78 < 1$$

However $\|A\|_1 < 1$ is sufficient to ensure that

$I \pm E$ are invertible and $\|(I \pm E)^{-1}\|_1 < \frac{1}{1 - 0.78} = \frac{1}{0.22}$

Rounding Error Analysis for Gaussian Elimination

In solving the linear system by Gaussian Elimination, roundoff errors are introduced at two stages

- 1) The input stage of A and b , i.e. machine representation of A and b
- 2) The elimination process. permutation does not cause roundoff errors, however, roundoff errors are introduced when computing $a_{ij}^{(k)} = fl(a_{ij}^{(k)} - m_{ik} a_{jk}^{(k)})$.
The multipliers $m_{ik} = fl(a_{ik}/a_{kk}^{(k)})$.

According to the IEEE standard (of accuracy)

$$fl(x) = x(1 + \delta) \quad |\delta| \leq u$$

$$fl(x \odot y) = (x \odot y)(1 + \delta), \quad |\delta| \leq u$$

$$\text{rel. Error} \leq \delta$$

u is the machine roundoff unit: $u = 2^{-24}$ in SP
 $u = 2^{-52}$ in DP.

According to the IEEE standard

$$fl(A) = A + \delta A \quad \text{where } |(\delta A)_{ij}| \leq u \quad i, j = 1, \dots, n$$

$$fl(b) = b + \delta b \quad |(\delta b)_i| \leq u \quad i = 1, \dots, n.$$

Hence, instead of solving the system $Ax = b$, we are solving a perturbed system

$$\boxed{(A + \delta A)(x + \delta x) = b + \delta b} \quad (*)$$

Here $x + \delta x$ is the exact solution of the perturbed system. $(*)$

The question is: how large/small can the error δx be? The following result gives an upper bound for the relative error $\frac{\|\delta x\|}{\|x\|}$

Theorem, Consider the system $(A + \delta A)(x + \delta x) = b + \delta b$ and assume $\|A^{-1}\| \|\delta A\| < 1$. Then

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\kappa(A)}{1 - \|A^{-1}\| \|\delta A\|} \left\{ \frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right\}.$$

proof

$$\begin{aligned} (A + \delta A)(x + \delta x) = b + \delta b &\Rightarrow A x + A(\delta x) + (\delta A)x + (\delta A)(\delta x) = b + \delta b \\ &\Rightarrow (A + \delta A)(\delta x) = \delta b - (\delta A)x \end{aligned}$$

By a previous result, $\|A^{-1}(\delta A)\| < 1$ implies that

$A + \delta A$ is invertible. Hence

$$\delta x = (A + \delta A)^{-1} [\delta b - (\delta A)x]$$

Take norms

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

Divide by $\|x\|$

$$\begin{aligned} \frac{\|\delta x\|}{\|x\|} &\leq \frac{\|A^{-1}\|}{1 - \|A^{-1}(\delta A)\|} \left[\frac{\|\delta b\|}{\|x\|} + \|\delta A\| \right] \\ &= \frac{\kappa(A)}{1 - \|A^{-1}(\delta A)\|} \left[\frac{\|\delta b\|}{\|A\| \|x\|} + \frac{\|\delta A\|}{\|A\|} \right]. \end{aligned}$$

Finally, from $b = Ax$, we have $\|b\| = \|Ax\| \leq \|A\| \|x\|$.

$$\Rightarrow \frac{\|\delta x\|}{\|x\|} \leq \frac{\kappa(A)}{1 - \|A^{-1}(\delta A)\|} \left[\frac{\|\delta b\|}{\|b\|} + \frac{\|\delta A\|}{\|A\|} \right]. \quad \square$$

According to ^{the} IEEE standard, $\frac{\|b\|_\infty}{\|b\|_\infty} \leq u$ and $\frac{\|bA\|_\infty}{\|A\|_\infty} \leq u$.

Thus, the preceding error bound tells us that

if $\kappa(A)$ is not too large, we expect the relative error in x not to be large.

Ex.
let $A = \begin{bmatrix} 3.9999 & 1.00 \\ 4 & 1 \end{bmatrix}$.

clearly A is invertible since $\det(A) = -10^{-4} \neq 0$.

However, if we are using 4-decimal digit arithmetic with rounding then

$$fl(A) = \begin{bmatrix} 4 & 1 \\ 4 & 1 \end{bmatrix} \text{ which is singular.}$$

let us assume that we went past the input stage of A and b and now A and b contain machine numbers. what is the magnitude of the roundoff errors committed during the Gaussian elimination phase?
? Here by Gauss-elimination we mean

- (1) $PA = LU$ factorization
- (2) Forward substitution
- (3) Back substitution.

we have the following fundamental results

Thm 1 Suppose $fl(A) = A$. The computed factors L and U satisfy

$$PA = LU + E$$

where the "error matrix" E satisfies the bound

$$\|E\|_\infty \leq \frac{u^2 p u}{1-u}, \quad p = \max_{i,j,k} |a_{ij}^{(k)}|$$

Another way of interpreting this is to say that the computed factors L and U are not exact L, U factors of the perturbed matrix $PA - E$.

$\frac{57}{1x}$

$$A = \begin{bmatrix} 2.345 & 5.347 \\ 1.207 & 9.423 \end{bmatrix}$$

we don't do pivoting (P-I) and use 4-decimal digit arithmetic with rounding.

The computed factors L and U are

$$L = \begin{bmatrix} 1 & 0 \\ 0.5147 & 1 \end{bmatrix}$$

$$fl\left(\frac{1.207}{2.345}\right) = fl(0.514712) = 0.5147$$

$$U = \begin{bmatrix} 2.345 & 5.347 \\ 0 & 6.671 \end{bmatrix}$$

$$fl(9.423 - fl(0.5147 \times 5.347))$$

$$\frac{2.7521009}{2.752} \rightarrow 2.752$$

$$fl(9.423 - 2.752) = 6.671$$

Now multiply LU exactly

$$LU = \begin{bmatrix} 2.345 & 5.347 \\ 1.2069715 & 9.4231009 \end{bmatrix} = \begin{bmatrix} 2.345 & 5.347 \\ 1.207 & 9.423 \end{bmatrix} - \begin{bmatrix} 0 & 0 \\ 2.85 \times 10^{-5} & -1.009 \times 10^{-4} \end{bmatrix}$$

$PA - E$

u for this machine (arithmetic) is 5×10^{-4}

Theorem 2 let L, U be the computed factors as in the preceding result. Then the computed solution \hat{x} of $LUx = Pb$ using Forward Then Back substitution satisfies

$$(L + \delta L)(U + \delta U)\hat{x} = Pb$$

where the perturbation (error) matrices δL and δU satisfy

$$\|\delta L\|_{\infty} \leq \frac{n(n+1)}{2}(1.01)u \text{ and } \|\delta U\|_{\infty} \leq \frac{n(n+1)}{2}(1.01)\rho u. \square$$

Putting the above two Theorems together we have

Theorem 3 The solution \hat{x} computed by Gaussian Elimination with partial pivoting exactly satisfies

$$(A + \delta A)\hat{x} = b$$

where provided $n^2 u \leq 1$,

$$\|\delta A\|_{\infty} \leq [1.01n^3 + 4(n+1)^2]\rho u, \quad \rho = \max_{i,j,k} |a_{ij}^{(k)}|.$$

we now compare \hat{x} to the exact solution x of $Ax = b$

$$A\hat{x} + (\delta A)\hat{x} = b = Ax \Rightarrow A(\hat{x} - x) = -(\delta A)\hat{x}$$

$$\Rightarrow \hat{x} - x = -A^{-1}(\delta A)(A + \delta A)^{-1}b$$

$$= -A^{-1}(\delta A)(A + \delta A)^{-1}Ax$$

$$(A + \delta A)^{-1} = [A(I + A^{-1}(\delta A))]^{-1} = (I + A^{-1}(\delta A))^{-1}A^{-1}$$

$$\Rightarrow \hat{x} - x = -A^{-1}(\delta A)(I + A^{-1}(\delta A))^{-1}A^{-1}Ax$$

$$\Rightarrow \|\hat{x} - x\|_{\infty} \leq \|A^{-1}\|_{\infty} \|\delta A\|_{\infty} \|(I + A^{-1}(\delta A))^{-1}\|_{\infty} \|x\|_{\infty}$$

$$\Rightarrow \frac{\|\hat{x} - x\|_{\infty}}{\|x\|_{\infty}} \leq \|A^{-1}\|_{\infty} \|\delta A\|_{\infty} \frac{1}{1 - \|A^{-1}(\delta A)\|_{\infty}} \frac{\|A\|_{\infty}}{\|A\|_{\infty}}$$

$$\rightarrow \frac{\|\hat{x} - x\|_{\infty}}{\|x\|_{\infty}} \leq \frac{\kappa_{\infty}(A)}{1 - \|A^{-1}(SA)\|_{\infty}} \frac{[(1.01)^n + 4(n+1)^2] \rho}{\|A\|_{\infty}} \cdot u$$

The quantity $\frac{\rho}{\|A\|_{\infty}} = \frac{\max_{i,j,k} |a_{ij}^{(k)}|}{\|A\|_{\infty}}$ is called the

"growth factor". It can be as large as 2^{n-1} .

However, it is extremely rare for it to exceed 16.