

Contents

1	Simultaneous linear equations	2
1	Theory of linear algebraic equations	2

mm FORTH M 9 Linear Algebra Linear Algebra :ontonts

§1 Simultaneous linear equations 551 Theory of linear algebraic equations
 §§2 Eigenvalue problems 52 Solving linear equations §§1 Cramer's rule §§2
 Pivotal Elimination 553 Testing §§4 Implementing pivotal elimination §§5 The
 program 556 Timing 53 Matrix inversion §§1 Linear transformations 5552 Ma-
 trix multiplication §§3 Matrix inversion §§4 Why invert matrices, anyway? §§5
 An example

54 LU decomposition 214 214 215 218 218 221 225 227 227 242 242 243 244
 245 245 246 213 214

51 Simultaneous linear equations

§§1 Theory oi linear algebraic equations

ChapterB-LhearNgebra Sewncronr'

Two common problems in scientific programming are the numerical solu-
 tion of simultaneous linear algebraic equation and computing the inverse of
 a given square matrix. This is sue". an important subject As an exercise in
 FORTH program development we shall now write programs to solve linear
 equations and invert matrices.

Simultaneous linear equations

We begin with a dose of mathematics (linear algebra), and then develop
 programs. Several actual programming session are reproduced *in toto* –
 mistakes and all– to give the flavor of program development and debugging
 1n FORTH.

Theory oi linear algebraic equations

We begin by stating the problem. Given the equations

$$\sum_{n=0}^{N-1} A_{mn} x_n = b_m . \quad (1)$$

– more compactly written $\vec{A} \cdot \vec{x} = \vec{b}$ – with known coefficient matrix A_{mn} and
 known inhomogeneous term b_m : under what conditions can we find unique
 values of x_n that simultaneoust satisfy the N equations 1?

The theory of simultaneous linear equations tells us that if not all the b_m 's are
 0, the necessary and sufficient condition for solvability is that the determinant
 of the matrix \vec{A} should not be 0. Contrariwise, if $\det(\vec{A}) = 0$, a solution with
 $\vec{x} \neq 0$ can be found when $\vec{b} = 0$.

Eigenvalue problems

Many physical systems can be represented by systems of linear equations.

Masses on springs, pendula, electrical circuits, structures, and molecules are examples. Such systems often can oscillate sinusoidally. If the amplitude of oscillation remains bounded, such motions are called stable. Conversely, sometimes the motions of physical systems are unbounded the amplitude of any small disturbance will increase exponentially with time. An example is a pencil balanced on its point. Exponentially growing motions are called -for obvious reasons - unstable.

Clearly it can be vital to know whether a system is stable or unstable. If stable, we want to know its possible frequencies of free oscillation; whereas for unstable systems we want to know how rapidly disturbances increase in magnitude. Both these problems can be expressed as the question: do linear equations of the form

$$Kx' = A x \quad (2)$$

have solutions? Here A is generally a complex number, called the eigenvalue (or characteristic value) of Eq. 2, and x' is often called the mass matrix. Frequently K is the unit matrix I ,

$$I, m = n \quad I'''' = 0, m = n, \quad (3)$$

but in any case, I must be positive-definite (we define this below). A non-trivial solution, x is 0, of the equation

$$Kx'' = 0 \quad (4)$$

exists if and only if $\lambda = 0$. This fact is useful in solving eigenvalue problems such as Eq. 2 above.

2.

eJuanvNoueiasa-Almm.

216 CMpterO-LklearNgabi-a Sclenwlcronnd

The secular equation (or determinantal equation)

$$\det(K - \lambda I) = 0 \quad (5)$$

is a polynomial of degree N in λ , hence has N roots (either real or complex). When $\lambda = \lambda_i$, these roots are called the eigenvalues (or 'characteristic values') of the matrix K .

Eigenvalue problems arising in physical contexts usually involve a restricted class of matrices, called real-symmetric or Hermitian

(after the French mathematician Hermite) matrices, for which.

$A_{ij} = \overline{A_{ji}}$. (The superscript \cdot denotes complex conjugation - \cdot :

see Ch. 7.) All the eigenvalues of Hermitian matrices are real numbers. How do we know? We simply consider Eq. 3 and its complex conjugate:

$$2\text{AM } x_m = 112 \text{ pm } x_n ; (6a)$$

$2 x:A'm = 1' 2 x: p_i ; (6b)$ Equation 6b can be rewritten (using the fact that K and flare. Hermitian)

$2 \text{MIA}'' \dots = 1' 2 x; P_m (6b')$ Multiply 6b' by x_m and 63 by x_i ; , sum both over m and i and subtract: this gives

$$0 = (i f i) 2 x; p_{i,m} x_m. (7)$$

However, as noted above, p is positive-definite, Le.

3. This follows from the fundamental theorem of algebra: a polynomial equation, $p(z): a_0 + a_1 z + a_2 z^2 + \dots + a_n z^n = 0$, of degree n (in a complex variable z) has exactly n solutions ,

WM''!

QuailWM 217

$$l t - p - X I Z x_{i,m} p_{i,m} x_m = 0 (8)$$

l'' for any non-zero vector x . Thus, Eq. 7 o A. E A, that is, l is real.

In vibration problems, the eigenvalue A usually stands for the square of the (angular) vibration frequency: ω^2 (oz. Thus, a positive eigenvalue). corresponds to a (double) real value, ω , of the angular frequency. Real frequencies correspond to sinusoidal vibration with time-dependence $\sin(\omega t)$ or $\cos(\omega t)$.

Conversely, a negative eigenvalue corresponds to an imaginary frequency, $i\omega$, and hence to a solution that grows exponentially in time, as

$$\sin(l' t) = i \sinh(\omega t) (9)$$

$$\cos(l' t) = \cosh(\omega t)$$

There are many techniques for finding eigenvalues of matrices. If only the largest few are needed, the simplest method is iteration:

make an initial guess $x^{(0)}$ and let

$$A x^{(0)} =$$

$$x^{(1)} =$$

Assuming the largest eigenvalue is unique, the sequence of vectors

$x^{(n)}$, $n = 1, 2, \dots$, is guaranteed to converge to the vector corresponding to that eigenvalue, usually after just a few iterations.

If all the eigenvalues are wanted, then the only choice is to solve the secular equation 5 for all N roots.

52 Solving linear equations

The test-and-development cycle in FORTH is short compared with most languages. It is usually easy to create a working program that subsequently can be tuned for speed, again much

4. Fordaritywenowonitthevcaor''- 'anddyad''0'symbdsfi'omvectorsandmauiwe
emvuouoim-me.

218

ChapterQ-LlnaarNgebra SclentlflcFORmi

more rapidly than with other languages. For me this is the chief benefit of the language.

"I Cramer's rule

Cramer's rule is a constructive method for solving linear equa

tions by computing determinants. It is completely impractical as a computer algorithm because it requires $O(N!)$ steps to solve N linear equations, whereas pivotal elimination (that we look at below) requires $O(N^3)$ steps, a much smaller number. Nevertheless Cramer's rule is of theoretical interest because it is a closed form solution.

Consider a square $N \times N$ matrix A . Pretend for the moment we know how to compute the determinant of an $(N-1) \times (N-1)$ matrix. The determinant of A is defined to be

$$\det(A) = \sum_{j=1}^N A_{1j} \det(A_{1j})$$

where

the $\det(A_{1j})$'s are called co-factors of the matrix elements A_{1j} and are in fact determinants of specially selected $(N-1) \times (N-1)$ sub-matrices of A .

The sub-matrices are chosen by striking out of A the n 'th column and m 'th row (leaving an $(N-1) \times (N-1)$ matrix). To illustrate, consider the 3×3 matrix

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}$$

and

produce the co-factor of A_{12} :

$$C_{12} = (-1)^{1+2} \det \begin{pmatrix} a_{11} & a_{13} \\ a_{31} & a_{33} \end{pmatrix}$$

$$C_{12} = - \det \begin{pmatrix} a_{11} & a_{13} \\ a_{31} & a_{33} \end{pmatrix}$$

We also attach the factor $(-1)^{m+n}$ to the determinant of the submatrix when we compute a_{mn} .

K's-DO

A determinant changes sign when any two rows or any two

columns are interchanged. Thus, a determinant with two identical rows or columns is exactly zero. What would happen to Eq. 11 if instead of putting A_{ij} in the sum we put A_{ji} where $k \neq m$? By inspection we realize that this is the same as evaluating a determinant in which two rows are the same, hence we get zero. Thus Eq. 11 can be rewritten more generally

$$\det A = \sum_{\sigma \in S_N} \text{sgn}(\sigma) \prod_{i=1}^N A_{i, \sigma(i)} \quad (13)$$

This feature of determinants lets us solve the linear equation $Ax = b$ by construction: try

$$x_i = \frac{\det A_i}{\det A}$$

$$x_i = \frac{\det A_i}{\det A} \quad (14)$$

We see from Eq. 13 that Eq. 14 solves the equation. Equation 14 also makes clear why the solution cannot be found if $\det(A) = 0$.

A determinant also vanishes when a row is a linear combination of any of the other rows. Suppose row 0 can be written

$$R_0 = \sum_{i=1}^N c_i R_i$$

The only number equal to its negative is zero.

Of course, the determinant is zero.

Chapter 9 Linear Algebra Scientific FORTH

$N-1$ or $E-2$ fins amt;

$$m=i$$

that is, the O 'th equation can be derived from the other $N-1$ equations, hence it contains no new information. We do not really have N equations for N unknowns, but at most $N-1$ equations. The N unknowns therefore cannot be completely specified, and the determinant tells us this by vanishing.

As an example, we now use Cramer's rule to evaluate the determinant of Eq. 11. We will write

$$\det A = \sum_{\sigma \in S_N} \text{sgn}(\sigma) \prod_{i=1}^N A_{i, \sigma(i)}$$

$$4 \text{ mini-v}$$

The determinant of a 1×1 matrix is just the matrix element, hence

$$\det A = A_{11}$$

$$320: 1$$

$\det A = 2 \cdot 6 + (1) \cdot 4 \cdot 1 = 8 \text{ am} = (18 \cdot 4) = 14 \cdot 320 = (3 \cdot 2) = 1 \cdot A \cdot 11 = 1 \cdot 8 + 0 \cdot (14) + 5 \cdot 1 = 13$. How many operations does it take to evaluate a determinant? We

see that a determinant of order N requires N determinants of order $N-1$ to be evaluated, as well as N multiplications and $N-1$

WWW-1

Glaciers-Wm 221

additions. If the addition time plus the multiplication time is r , then

$T_N = N(r + T_m)$.

It is easy to see the solution to this is

- Nire.

" 1 $T_N = N(r + T_m)$ "an

In other words, the time required to solve N linear equations by Cramer's rule increases so rapidly with N as to render the method thoroughly impractical.

§2 Pivotal Elimination

The algorithm we shall use for solving linear equations is elimination, just as we were taught in high school algebra. However we modify it to take into account the experience of 40 years's solving linear equations on digital computers (not to mention a previous 20 years's worth on mechanical calculators!), to minimize the buildup of round-off error and consequent loss of precision.

The necessary additional step involves pivoting selecting the largest element in a given column to normalize all the other

6. Professors always say this, hee, heel See R. Sedgewick, Algorithm (Addison-Wesley Publishing Company, Reading, MA 1983).

5"3.1

$N!$ means $N \times (N-1) \times (N-2) \times \dots \times 2 \times 1$. The base of the "natural" logarithm is $e = 2.7182818\dots$. A computer stores a number with finite precision say 6-7 decimal places with 32-bit floating-

point numbers. This is enough for many purposes, up to daily in science and engineering where the data are rarely measured to better than 1 number, about 10% in magnitude, are multiplied. Their product is of order 10% and is known to six significant figures. Now add it to a third number of order unity: 10%. Later, a fourth number also of order unity is subtracted from the sum; the result will be a number of order 10% but now known only to two significant figures. Matrix arithmetic is full of multiplications and additions. The lesson is clear to minimize the inevitable) kind of mass of data with which we must cope. They may be reduced to products and sums as easily as possible.

222

Chapter 9 - Linear Algebra Sclermflc FORTH

elements. This will be clearer with a concrete illustration rather than further description: Consider the 3 X 3 system of equations:

$105x + 0y + 324z = 4$ (15) $116x + 2y + 2z = 2$

We check that the determinant is #0; in fact $\det(A) = 13$. The first step in solving these equations is to transpose rows in A and in b to bring the largest element in the first column to the/loo position.

Natl-a;

e The x's are not relabeled by this transposition. ,

e We choose the row (n=1 second row) with the largest (in absolute value) element Ano because we are eventually going to divide by it, and want to minimize the accumulation of roundoff error in the floating point arithmetic.

Transposition gives 3 2 4 xo 4 1 0 5 x1 = 0 (16) 1 1 6 x2 2

Now divide row 0 by the neono (in this case, 3) to get 1 2/3 4'3 xo V3 1 0 5
x1 = 0 (17) 1 1 6 X2 2

Subtract row 0 times Ana from rows with n i 0 1 2/3 4': 10 V: o 2/, 11/3 x, =
-Va (13) 0 1'3 W: x, 26

Since —An— i —A21— we do not bother to switch rows 1 and 2, but divide row 1 by An = -2/3, getting

cranes-mm 223

l '3 V: to V3 0 1 -1V2 xi - 2 (19) 0 V) IV! 12 V3

We now multiply row 1 by A'' a 1/3 and subtract it from row 2, and also divide through by A'' to get

173 V3 10 V3 0 1 -u/2 x; 8 2 (20) 0 0 1 x; 0

The resulting transformed system of equations now has 0's to the left and below the principal diagonal, and 1's along the diagonal. Its solution is almost trivial, as can be seen by actually writing out the equations:

1Xo+2/3X— +V3X2=4/3 (21.0) 0x0 +1x1+(11/2)X2 = 2 (21.1) Oxo+0x1+1x2=0
(21.2)

That is, from Eq. 21.2, x; = 0. We can back-substitute this in 21.1, then solve for x, to get

x, = 2 (11/2)-(0) = 2

and similarly, from 21.0 we find xo=V32/3(2)+V3(0)=0. We test to, see whether this is the correct solution by direct trial: 1-O+0-2+5-O=0 3-0+2-2+4-0=4 1-O+2-2+6-O=2

his works we have indeed found the solution.

cmvliouieieaa-Alrlgium.

224

Chapter 9 - Unear Algebra Scientific FORTH

How much time is needed to perform pivotal elimination? We concentrate on the terms that dominate as $N \rightarrow \infty$.

The pivot has to be found once for each row; this takes $N-k$ comparisons for the k 'th row. Thus we make $\sum_{k=0}^{N-1} (N-k)$ comparisons of 2 real numbers. (For complex matrices we compare the squared moduli of 2 complex numbers, requiring two multiplications and an addition for each modulus.)

We have to divide the k 'th (pivot) row by the pivot, at a point in the calculation when the row contains $N-k$ elements that have not yet been reduced to 0. We have to do this for $k = 0, 1, \dots, N-1$,

requiring $\sum_{k=0}^{N-1} (N-k)$ divisions.

The back-substitution requires 0 steps for row N , 1 multiplication and 1 addition for row $N-1$, 2 each for it, etc. That is, it requires

$$\sum_{k=0}^{N-1} (N-k) = N(N+1)/2$$

multiplications and additions.

The really time-consuming step is multiplying the k 'th row by A''_{kj} , $j \geq k$, and subtracting it from row j . Each such step requires $N-k$ multiplications and subtractions, for $j=k+1$ to $N-1$, or $\sum_{k=0}^{N-1} (N-k)(N-k-1)/2$ multiplications and subtractions. This has to be repeated for $k = 0$ to $N-2$, giving approximately $N^3/3$ multiplications and subtractions.

In other words,

the leading contribution to the time is $\frac{1}{3}N^3$, which is a lot better than $N!$ as with Cramer's rule.

When we optimize for speed, only the innermost loop requires

$O(N^2)$ operations; the $O(N^3)$ operations need careful tuning; the $O(N^2)$

operations comparing floating point numbers, dividing by the

pivot, and back-substituting need not be optimized because for large N they are overshadowed by the innermost loop.

9.

The exception to this general rule, where more complete optimization would pay, would be an application that requires solving many sets of equations of relatively small order.

MM!

"3

cleansmm 225

Mia

Since we have worked out a specific 3×3 set of linear equations, we might as well use it for testing and debugging. Begin with some test words that do the following jobs:

e Create matrix and vector (inhomogeneous term) e Initialize them to the values in the example

e Display the matrix at any stage of the calculation e Display inhomogeneous term at any stage

```
\ Display words
V{ or M{ stands torwhatanarrayputsonstack
U: G. F. x. ;
..M (M{-- ) FINIT DUP
D.LEN 0 DO CR DUP
D.LEN 0 DO DUP (--M{ M{ )
J}{I}} DUP>R G@ G.
LOOP
LOOP DROP ;
```

```
..v (V{~) DUP D.LEN
0 DO on DUP 1}{0} G@ G. LOOPDROP:
```

We define a word to put ASCII numbers into matrices:
:GEFF# BL TEXT PAD S->F;

This word takes the next string¹⁰ from the input stream (set off by ASCII blank, BL) and uses the PIS/FORTH word S ¿ F to convert the string to floating point format and put it on the 87stack.

10. The word 'IEX' inputs a counted string using the FORTH-D standard word WOIID, and places it "PAD" his is why PAD appears in the definition of Gfl-FlA definition of EXT mid-it be: 'IEX' (delimiter-) WORD DUP C@ 1+ uns wxrcuovs; Note that all words prefaced with C are byte operations by FORTH convention.

OJilthNoUetm-MlIdusleseived.

Chapter 9 - Unear Algebra Scientlflc FORTH

GET-F# is used in the following:

```
: <DO-VEC> 0DO GET-F# DUP I0} G!
LOOP DROP;
```

```
:TAB->VEC DUP D.LEN <DO-VEC> ;
```

```
:TAB->MAT DUP D.LEN DUP * <DO-VEC> ;
```

The file EX.3 will be found in the accompanying program disk-ette. The explanatory notes below refer to EX.3.

Notes:

e The word `.` (emits everything up to a terminating) `.`

e HS/FORTH-I, because of its segmented dictionary, uses a word `TASK` to define a task-name, and `FORGET-TASK` to `FORGET` everything following the task-name. The FORTH word `FOR-GE` fails in HS/FORTH when it has to reach too deeply into the dictionary.

e Ex.3 uses the word `}row{` defined below.

This is what it looks like when we load EX.3:

```
FLOAD EX.3 Loading EX.3 ReadadimenslonSvectorand 3113 matrbltror-
natablethendispiaythem Now displaying vector. V.V CR 0.0000000 4.0000000
2.0000000 Now displayhg matrbc A .M CR 1.0000000 0.000(1)") SIXJOOCDO
3.0000000 2.0000000 4.0000000 1.0000000 1.0000000 6.0000000 say Exa FORGET-
TASKto graceitfly FORGET these words olt
```

WM"

Mil-Wm 227

"4 WWW

Now we have to define the FORTH words that will implement this algorithm. In pseudocode, we can express pivotal elimination as shown in Fig. 9-1 below.

ch

A B SOLVE (solutiontoAx=btinB) BEGIN

setn=0

findpivot(rowm i n)

swaprowsman

divide row n by pivot: A'M. =AMJAM

subtract: A'... =A..A...A'... increment n n N = UNTIL Back-substitute: x" =
b'n 2 Am... xm(12)

m;n

DONE

Fig. 9-1 Pseudocode for pivotal elimination

"5 The program

Whenever we write a complex program we are faced with the problem "Where do we begin?" FORTH emphasizes the con-

struction of components. and generally cares little which component we define first.

For example, we must swap two rows of a matrix. The most obvious way moves data:

```
row1 -> temp row2 -> row1
temp -> row2
```

Although data moves are relatively fast, up to $N^2/2$ swaps may be needed, each row containing N data; row-swapping is an $O(N^3)$ operation which must be avoided. Indirectly setting up a list of row pointers and exchanging those, uses $O(N)$ time.

OWVNObletmAlmm.

228

Chapter 9 - Linear Algebra Scientific FORTRAN

We anticipate computing the address of an element of A in a DO loop via the phrase

```
A(J, I, row, col)
```

Suppose we have an array of row pointers and a word `row` that looks them up. Then the (swapped) element would be

```
A(J, row(I),
```

To implement this syntax we define

```
:swapper (n - 1) CREATE 0 DO I, LOOP DOES OVER ++ @ ; this is a
defining word
```

We would define the array of row indices via

```
A D.LEN swapper row create array row
```

We have initialized the vector `row` by filling it with integers from 0 to $N-1$ while defining it.

Suppose we wanted to re-initialize the currently vectored row-5 index array: this is accomplished via

```
A D.LEN 'row refill .i
```

where

```
ramm, .. .m' .
```

```
: refill (n adr -) SWAP 0 DO I 2* DDUP + ! 2 +LOOP DROP;
```

11.

```
mm"... 'mnLk,fi-u_.. aims . Vannr
```

The HS/FORTH words VAR and IS define a data structure that can be changed, like a FORTH VARIABLE: 0 VAR X 3 IS X but has the run-time behavior of a CONSTANT: X . 3 olt.

```
:ADRS >R 2'OVER + SWAPR> 2' +;
(amn-a+2ma+2n)
OVAR ?SWAP \to keep track of swaps
```

```
:}SWAP (amn-)DDUP -
IF DDROP DROP \noswap - deanup
```

```
ELSE ADRS DDUP (.-121 2)
@SWAP@ (-'12[21[1I)
nor! SWAP 1
```

```
-1 ?SWAP XOR IS ?SWAP \ '' ?SWAP
THEN ;
```

Test this with a 3-dimensional row-index array:

```
3 swapper }row{
:TEST ODOI}row{ CR. LOOP;
3 TEST
```

```
0
```

```
1
```

```
2 ok
```

Now swap rows 1 and 2:

```
'}row{ 1 2 }SWAP 3TEST
0
```

```
2
```

```
10k
```

and back again:

```
'}row{ 1 2 }SWAP 3 TEST
0
```

```
1
```

```
20k
```

Next, we need a word to find the pivot element in a given column, searching from a given n. The steps in this procedure are shown in Fig. 9-2 on page 230 below.

We anticipate using generic operations *Git* defined in Chapter 5

to perform fetch, store and other useful manipulations on the matrix elements, without specifying their type until mn-time. It

```
:A{{ col# }}PIVOT
```

```
N col# 1+ D0
```

```
Yes IF
```

```
No 1 Put larger on fstack
```

```
set I.PIV = I
```

```
t__l THEN
```

```
. LOOP
```

```
; \ exit
```

```
M{{ | }}{ 00"" }}I > M{{ col# }}{ 00h" }}I 7
```

Fig. 9-2 Pseudocode for finding the pivot element

is thus useful to have a place to store the type (other than the second cell of the array data structure). We arrange this via

OVART

The rest of the definition is rendered self-explanatory by the vertical commenting style (a must for long12 words):

```
OVAR Length \Length of matrix
```

```

OVAR COL \currentcol#
OVAR I.PIV
OVAR a{{
: INITIALIZE (M{{ -- ::--)
IS a{{
a type@ IS T
a LEN@ IS Length ;

\ |.PIV is used to return the result
\ a{{ stores the adress of M{{

\ initialize T

\ length - > Length

```

12. While Brodie (11", p.180) quotes Charles Moore. FORTH's inventor, as offering l-line defini- tions as the goal in FORTH programming, sometimes this is just not possible.

```

i
a place to store data type
i

2))PIVOT (M{{ col-- ::--)

IS COL ZECOL IS LPN

INITIAU

fl: COL row COL}} (--aog.ofl[M{% oot,ool}}]t)

> 8 SF >F (87:--|1st.olt)

R> COL1+ \Sloop|infgLCOL+1

DO I .
a{{ I }row{ COL }} (- - seg.ofl[M{fII,col}}] t)
>FS GABS FS>F(87: |old.olt Ineweltl)
xoup F< \test' newett > olden
IF FSWAP IISI.PIV THEN FDROP

LOOP \ondloop

FDROP ; \deanupfstack

```

\ Usage: A{{ 2 "PIVOT

Mates:

0 We avoid filling up the parameter stack with addresses that have to be Eagle , by putting arguments in named storage locations (J). We anticipate setting all the parameters in the beginning using IN! 21". which can be extended later if necessary.

0 We have used a word from the COMPLEX arithmetic lexicon, namely XDUP(FOVER FOVER) to double-copy the contents

of the fstack, since the test F_i drops both arguments and leaves a flag.

0 There is little to be gained by optimizing (and if we did we should have had to avoid the generic Gx words because they cannot be optimized by the PIS/FORTH recursive-descent optimizer) because only 1/2 time is used, negligible (for large

compared with the [3 in the innermost oop.

The Gx operations are found in the file GLIBJ-TH.

To test the word "PIVOT we can add some lines to the test file:

oncn.(weammmmo.1.2)

on

```
G! .(A{{0}}PIVOT mv .) BL EMIT A{{0}}PNOT IPN.
on .(A{{1}}PNOT IPIV .) BL em A{{1}}PNOT IPN .
on .(A{{2}}PIVOT IPN .) BL am A{{2}}PNOT new.
```

one:

The result is the additional screen output

```
A{{ 0 }}PIVOT IPIV 1
A{{1}}PIVOT IPIV 1
A{{ 2 }}PIVOT IPIV 2
```

OR

From our pseudocode expression (see Fig. 9-1, p. 227) of the pivotal elimination algorithm we see that the next operation is to multiply a row by a constant. To do this we define

OVAR ISTEP

\ include phrase

T #BYTES DROP IS ISTEP

\ in INITIALIZE


```
: DO(ROW*X)LOOP (seg off Lfin I.beg - - :: x - - x)
DO DDUP I + DDUP T >FS G'NP TFS>
ISTEP /LOOP DDROP ;
```

\ G*NP means "(generic) multiply, no pop"

```
:}}ROW*X M{{ row - :: x x)
UNDER }row{ 0 A (- - r seg offt)
DROP ROT (- -seg off r)
ISTEP * Length ISTEP * SWAP \Iloop indices
DO(ROW*X)LOOP ; \Iloop
```

\Ex: A{{ 2 }ROW*X

Mates:

o While DO(ROW*X)LOOP and ??ROW*X clear the parameter stack in approved FORTH fashion, they leave the constant multiplier x on the fstack. That is, we anticipate next multiplying the corresponding row of the inhomogeneous term b (in the matrix equation $A \cdot x = b$) by the same constant x .

e We assume INITIALIZE (including INITJSTEP) has been invoked before "PIVOT" or "110 'X.

o Anticipating the (possible) need to optimize, we factor the loop right out of the word. We also come to the addresses fast by putting base addresses on the stack and then incrementing them by the cell-size, in bytes.

Subtracting row 1 (times a constant) from row I is quite similar to multiplying a row by a constant. The process can be broken down into the pseudocoded actions in Fig. 9-3 below.

```
no (K-JTOL)
x e@

M[I,K] e@ G'NP (:: - x M[1,K] 'X)
M{J,K} 6@ GR- (:: --XM[J,K]-X'M[I,K])
M{J.K} G1

LOOP
```

Fig. "Steps In nm-rwx

A first attempt might look as follows:

```
\auxiliarywords

:4DUP DOVER DOVER ;

(abcd--abcdabcd)

: + + (51.0152.02n 1.01 +n 52.02+n)
DUP>R + ROT R> + -ROT;

:}}R1-R2"X (r1r2-- :x--x)

> R > R

a R> row 0 DROP \M r10
aila@imioli one. \Miaoli
Length iSTEP " \ L'lSTEP (upper limit)

R> ISTEP * \ r2'ISTEP (lower limit)

DO 4DUP I ++ \duplicateaincbasadr

T >FS G'NP (2: xx*M[r2,k])
DDUP T >FS GR-
T >FS \ I result
ISTEP [LOOP DDROP DDROP ;
\ INITIALIZE is assumed
```

As with HROW'X we avoid the execution-speed penalty of calculating addresses within the loop by not using the phrase

However, we are doing a lot of unnecessary work inside the loop by making 5 choices based on datatype. There are also a lot , unnecessary moves to/from the ifstack. This is an example where speed has been sacrificed to compactness of code: one program solves equations of any ("scientific") data format REAL'4, REAL-8, COMPLEX'8 and COMPIJEX' 16. i

Conversely, we can both eliminate the extra work and accelerate? execution simply by defining a separate inner loop for each data; type, letting the choice take place once, outside the inner loopX This multiplies the needed code fourfold (or more-fold, if we? optimize).

I

The 4 inner loops are

```

: Re.LP (\$1.01 52.02 L'istep r2*istep - - :2 x - - x)
FS>F >R
DO 4DUP I + + R32@L F*NP
DDUP R32@L FR- R32!L
4 /LOOP R> F>FS ;

```

```

: DRe.LP ($1.01 $2.02 L'istep r2'istep - - z: x - - x)
FS>F >R
DO 4DUP 1 + + R64@L F*NP
DDUP R64@L FR- R64!L
8 /LOOP R> F>FS ;

```

```

: X.LP (51.01 52.02 L*istep r2*istep - - z: x - - x)
FS>F >R
DO 4DUP 1 + + CP@L X'NP
DDUP CP@L CPR- CP!L
8 /LOOP R> F>FS ;

```

```

: DX.LP (\$1.01 52.02 L'istep r2'istep - - :: x - - x)
FS>F >R
DO 4DUP I + + DCP@L X'NP
DDUP DCP@L CPR- DCP!L
16 /LOOP R> F>FS;

```

Not surprisingly, the loops contain some duplicate code. but this- ' is a small price to pay for the speed increase. A significant furthe r3

increase can be obtained easily. using the HS/FOR'TH recursive- descent optimizer to define these inner loops, or by redefining the loops in assembler (for ultimate speed).

Now we can redefine }}R1-R2'X using vectored execution, via HS/FOR'I'H's CASE: . . . :CASE construct, or the equivalent high- level version given here:

```

: CASE: CREATE
DOES> OVER + + @ EXECUTE ;
: ;CASE [COMPILE] ; ; IMMEDIATE

```

```

CASE: DO(R1-R2"X)LOOP
Re.LP DRe.LP X.LP DX.LP ;CASE

```

```

2}}R1-R2'X (r1 r2-- ::x--x)

```

```

>R >R

```

```

2112a mm; 328: {main
Length I.STEP * \ ULSTEP (upper limit)
R> I.STEP ' \ r2'I.STEP (lower limit)

```

```
T DO(R1-R2'X)LOOP DDROP DDROP ;
\ We assume INITIALIZE has been invoked
```

Another word is needed to perform the same manipulation on the inhomogeneous term. Since this latter process runs in $O(N)$ time we need not concern ourselves unduly with efficiency.

```
:}V1V2*X (r1 r2-- ::x--)
SWAP >R >R
b{ 0.0 } DROP DDUP \V[0] V[0]
R> irowilSTEP' + >FS G'
R> row ISTEP " + DDUP >FS GR- FS> ;
```

Hats:

0 After }V1-VZ'X executes, the multiplier x is no longer needed, so we drop it here by using G' rather than $G'NP$.

We now combine the words }SWAP, }}PIVOT, }}ROW'X. }RI-RZ'X and }V1-VZ'X to implement the triangularization portion of pivotal elimination.

Since the Gaussian elimination method makes it very easy to compute the determinant of the matrix as we go, we might as well do so, especially as it is only an $O(N)$ process. By evaluating the determinants of simple cases, we realize the determinant is simply the product of all the pivot elements, multiplied by

(.1) swaps. Hence we need to keep track of swaps, as well as to multiply the determinant (at a given stage) by the next pivot element. The need to do these things has been anticipated in defining }SWAP above.

Computing the determinant also lets us test as we go, that the equations can be solved: if at any stage the determinant is zero, (at least) two of the equations must have been equivalent. Should it be necessary to test the condition of the equations, this too can be found as we proceed, by computing the determinant.

Here is a simple recipe for computing the determinant, with checking to be sure it does not vanish identically. We use the intelligent fstack (ifstack) defined in Ch. 7.

```
DCOMPLEX SCALAR DET \ room for any type
:INIT\_DER T 'DEI' I \set Type
TG=1 DEI' GI; \setdet=
```

%1.E-10 FCONSTANT CONDITION

```
:DETERMINANT (- - 2: x - - x)
DET >FS G*NP
?SWAP IF GNEGATE THEN
```

```
FS.DUP GABS FS>F CONDITION F<
ABOR'1" determinant too small!' DET FS> ;
```

13.

The condition of a system of linear equations refers to how accurately they can be solved. Equations can be hard to solve precisely if the inverse matrix A^{-1} has some large elements. Thus, the error vector for a calculated solution, $d = -b - A^{-1}szlc$ can be small, but the difference between the exact solution and the calculated solution can be large, since (see §9§3 below)

$x_{Bar-ct} = A^{-1}W_6 A$ A test for ill-condition is whether the determinant gets small compared with the precision of the arithmetic. See, 43.3., A. Ralston, A First Course in Numerical Analysis (McGraw-Hill Book Co., New York, 1965).

Now we define the word that triangularizes the matrix:

```
0VARb( \tokeepthesteokehort
: INITIAUZE (M{V{o- ::--})

1e

1Sa

3 .TYPE IS T

a D.LEN is Length

a D.TYPE #BYTES DROP ISISTEP

iN .DET ; \setdet=1

:}/PIVOT }row{) DROP DDUP T >FS G' TFS>;
(segofIt--::x--))

:TRIANGULAFIIZE (M{{ V{ - -)

INITIAUZE
Length 0 DO \ loop 1 - by rows
a{{1} PIVOT \findpivotincoil
' )row I 1.PlV }SWAP \ exchange rows
a{{ I }row{ I }} >FS \pivot->iistack
DETERMINANT \ mic det
1/6 a{{ i }}ROW'X \row i Ipivot
b{ 1 }/PIVOT \inhom. term /pivot
Length 11+ DO \ioop2-byrows
```

```

a l}row{J}} >FS \x->ifstack
a lJ}}R1-R2"X

\ row[i] = row[i]-row[j]'x
b{ I J }V1-V2*X
\ same for b{ and drop x
LOOP \ and loop 2
LOOP ; \end loop 1
\Usege: A{{ B{ TRIANGULARIZE

```

Now at last we can back-solve the triangularized equations to find

the unknown x 's. The word for this is BACK-SOLVE, defined as follows:

```

: }BACK-SOLVE (- - )
0 LENGTH 2- DO \ outer loop
T G=0
LENGTH I 1 + DO \ inner loop
a? J }row{ I }} > PS
b I }row{ } >FS G*G+\
LOOP \ inner loop
b{ 1 ){ } DROP DDUP T >FS \
GR- T FS > \
1 +LOOP ; \outer loop J

```

Putting the entire program together we have the linear equation solver given in the file SOLVEL. Examples of solving dense 3x3 and 4 X4 systems are included for testing purposes.

§§6 Timing

e should like to know how much time it will take to solve a

given system. (Of course it is also useful to know whether the solution is correct!) We time the solution of 4' sets of N equations, with 4 different values of N . The running time can be expressed as a cubic polynomial in N with undetermined coefficients:

$$TN = a_0 + a_1 N + a_2 N^2 + a_3 N^3 \quad (19)$$

Evaluating 19 for 4 different values of N , and measuring the 4 times TN_i , we produce 4 inhomogeneous linear equations in 4

unknowns: $a_i, i = 0, 1, 2, 3$. As luck would have it, we just hap-

pen to have on hand a linear equation solver, and are thus in a position to determine these coefficients numerically.

For this timing and testing chore we need an exactly soluble dense system of linear equations, of arbitrary order. The simplest siqch system involves a rank-1 matrix, that is, a matrix of the form :

14.

We employ the standard notation that a vector u is a column and an adjoint vector V is a row. Their outer product, $u^T v$, is a matrix. Given a column v , we construct its adjoint (row) V by taking the complex conjugate of each element and placing it in the corresponding position in a

row-vector.

Gustavo WNW 23.

$A - I - u v^T$ (20) In terms of matrix elements,

$A_{ii} - u_i u_i^T$ (20')

The solution of the system $A x = b$ is simple: using the standard notation

$v^T x = \sum_i v_i x_i$ (21) we have $X = D^{-1} (V, X) U$. (22)

The coefficient (v, x) is just a (generally complex) number; we compute it from 22 via

$(V, X) = (V, I x) + (V, u) (V, X)$ (23) $O(1) - (v, b) (V, X) = 1 - (v, u)$. (24)

$x_i = 0, +$

Everything in 25 is determined, hence the solution is known in closed form.

An example that embodies the above idea is given in the file EX.20, where we make the special choices

Chapter 9 Unear Algebra Scientific FORTH

N even

GAO-s.s 01

We give the times only for the case of a highly optimized inner loop (written in assembler), for the type REAL4:

Table 9-1 Execution times for various N (9.54 MHz 8086 + 8087 machine)

N Time 20 1.20 50 7.75 75 19.6 100 38.8

The coefficients of the cubic polynomial extracted from the above are (in seconds, 9.54 MHz 8086 machine, machine-coded inner

loop)

$a_0 = 0.32727304$ $a_1 = -0.01084850$ $a_2 = 0.00241636$ $a_3 = 0.00001539$

from them we can extrapolate the time to solve a 350×350 real system to be about 16 minutes. On a 25 MHz 80386/80387 machine with 32-bit addressing implemented, the time should decrease five- to tenfold.

It is interesting to explore a bit further the extracted value of a_3 ,

. . . 1 which is 3- fined above as Re.LP and hand-coded in assembler for ultimate speed). Converting to clock cycles on an IBM-PC compatible (running at 9.54 MHz) we have an average of 440 clock cycles per traversal (of this loop). Recall the operations that are needed: a 32-bit memory fetch, a floating-point multiply, another 32-bit memory fetch, a floating-point subtraction, and a 32-bit store. The initial fetch-and-multiply can be compressed into a single co-processor operation, as can the fetch-and-subtract. The times in cycles for these basic operations are

the time needed to evaluate the innermost loop (de-

TaUe92ExecutiontImesotlnnennostiooperntions

Operation lime (cpu clock5) memory FMUL 133 memory FSUBR 128 mem-
ory FSTP 100

overhead 61

Total 422

15. thnlthinkthatsolvinglmxlflystems keyto myPh.D.researchin1966 took-
thebetter putofanhornonanlnnumnheseresuhsseeminaedible.

16. See WP.

§3 Matrix Inversion

§§1 Unear transformations

ChapterO-LlnaarAlgebra SciandflcFOfl

We see from Table 9-2 above that the time computed from t cpu and co-processor specifications (422 clocks) is close to measured time (440 clocks). The slight difference doub comes both from measurement error and from the fact th timing of CISC chips is not an exact art (for example, there a periodic interruptions for dynamic memory refreshment and f the system clock).

Finally, we confirm that little is to be gained by optimizing out loops. Suppose, 5.3., we could halve a; by clever programmin then we should cut the N =350 time from 16 to 13 minutes, the time for N = 1000 by some 20 minutes in 5 hours, or 6

We introduce this subject with a brief discourse on liners algebra of square matrices.

Suppose A is a square (NXN) matrix and x is an N-dimensiomi vector (a column, or NXI matrix). We can think of the symbolic operation

$$y = A \cdot x \quad (26)$$

as a linear transformation of the column x to a new column y. In terms of matrix elements and components,

$$Nl Ym = E Amnxn \quad (27)$$

n-o

The transformation is linear because if x is the sum of 2 columns (added component by component)

$$it = x'' + is), (23) \quad -$$

1

we can calculate A^{-1} either by first adding the two vectors and then transforming, written

$$A^{-1}x = A^{-1}(x_1 + x_2) \quad (29)$$

or we could transform first and then add the transformed vectors. The identity of the results is called the distributive law

$$A^{-1}x = A^{-1}(x_1 + x_2) \quad \text{"amuse"} \quad (30) \quad \text{of linear transformations.}$$

"2 Matrix m Now, suppose we had several square matrices, A, B, C, \dots . We could imagine performing successive linear transformations on a

$$\text{vector } x \text{ via } y = Ax \quad (31a) \quad z = By \quad (31b) \quad w = Cz \quad (31c)$$

These can conveniently be written

$$w = C[B(Ax)] \quad (32)$$

The concept of successive transformations leads to the idea of multiplying two matrices to obtain a third:

$$D = BA \quad E = CD \quad (33) \quad \text{In terms of matrix elements we have, for example } D_{ij} = \sum_k B_{ik} A_{kj}$$

The important point is that the (matrix) multiplications may be performed in any order, so long as the left-to-right ordering of the factors is maintained:

$$D = B(Ax) = (BA)x$$

$$E = C(Dx) = (CD)x$$

$$C(A^{-1}x) = (CA^{-1})x \quad (35)$$

Equation 35 is known as the associative law of matrix multiplication. Finally, we note that as hinted above the left-to-right order of factors in matrix multiplication is significant. That is, in general,

$$A(Bx) \neq (AB)x \quad (36)$$

We say that, unlike with ordinary or even complex arithmetic, matrix multiplication even of square matrices does not in general obey the commutative law.

§§3 Matrix Inversion With this introduction, what does it mean to invert a matrix? First of all, the concept can apply only to a square matrix. Given an $N \times N$ matrix A , we seek another $N \times N$ matrix A^{-1} with the property that

$$A'' - AaA - A'' = I \quad (37)$$

where I is the unit matrix defined in the beginning of this chapter (Eq. 4 1's on the main diagonal, 0's everywhere else). We have: implied in Eq. 37 that a matrix that is an inverse with respect to left-multiplication is also an inverse with respect to right-multiplication. Put another way, we imply that

$$(r')'' a A \quad (38)$$

The condition that - given A we can construct A^{-1} is the same.— as the condition that we should be able to solve the linear equation

$$A^{-1}b;$$

t namely, $\det(A) \neq 0$.

5'4 We invert matrices, We have developed a linear equation solver program already why should we be interested in a matrix inversion program?

Here is why: The time needed to solve a single system (that is, with a given inhomogeneous term) of linear equations is conditioned by the number of floating-point multiplications required:

about $143/111$): number needed to invert the matrix is about N^3 , roughly $3X$ as many, meaning roughly $3x$ as long to invert as to solve, if the matrix is large. Clearly there is no advantage to inverting unless we want to solve a number of equations with the same coefficient matrix but with different inhomogeneous terms. In this case, we can write

$$x = A^{-1}b$$

and just recalculate for each b . Clearly this breaks even relative to solving 3 sets of equations for 3 different b 's and is superior to re-solving for more than 3 b 's

s55 Anmmmpie

Let us now calculate the inverse of our 3×3 matrix from before. The equation is (let $c = A^{-1}$ be the inverse")

$$105 \ 000301 \ \text{can} \ 100 \ 3 \ 2 \ 4 \ C10 \ C'' \ C12 = 0 \ 1 \ 0 \ (40) \ 116 \ \text{CanCa} \ 001$$

18. Note: if C is a right-inverse, $AC = I$, it is also a left-inverse, $CA = I$.

OJUIUIVNanim-Alrldasraaarvad.

246

CMptarO-Urngebra ScilntlficFORT

It is easy to see that Eq. 40 is like 3 linear equations, the unknown being each column of the matrix C . The brute-force method I calculate A is then to work on the right-hand-side all at once we triangularize, and back-solve. It is easy to see that triangularization by pivotal elimination leads to

$\frac{1}{2} \frac{3}{3} V_s$ can cat can 0 V_3 0 0 1 $\frac{1}{2} C_{10}$ Ctr $C_{12} = '3': V_2$ 0 (41) 0 0 1 C_{20}
 C_{12} 022 V_{13} "V13 13 1

results of back-solving, column by column. This is N times a.

if

1 total time required for brute-force inversion is $\approx 5 N^3$. By keep

To construct the inverse we replace the right-hand side with thq

much work as back-solving a single set of equations, hence th

ing track of zero elements we could reduce this time by another elimination) of $O(N^3)$. However, the brute-force method is unsatisfactory for another reason: it takes twice the storage of more sophisticated algorithms. Modern matrix packages therefore "8:" LU decomposition for both linear equations and matrix inver

sion.

§4 LU decomposition

We now investigate the LU decomposition algorithm". Suppose a given matrix A could be rewritten

awam 1m 0 0 lmflm $A: 310311 = L'U = 1401110$ 0 11 (42) 0 0

then the solution of

19.

Sec. 2.3., WH. Press, B.P. Flannery, SA. Teukolsky and W.T. Vetterling. Numerical Recipe: (Cambridge University Press, Cambridge, 1986), p. 31H.

owe-mm 241

$(L-u)-nL-(u-x) -b$ (43) can be found in two steps: First. solve $Lty - b$ (44) for y : $U-x$ (45) via $looYo = bo$ $MOYO'I'ttIYI = 01$ (46) $130\% + 421Y1 + la$ ': $= D$: etc.

which can be solved successively by forward substitution. Next solve 45 successively (by back-substitution) for x :

$I'NtN-t XN1 = YNt I'M-2 N-2 XN-z + I'M2 N-1 xNi = YNz$ (47) etc.

The n 'th term of Eq. 46 requires n multiplications and n additions. Since we must sum n from 0 to $N-1$, we find $N(N-1)/2 \approx 4 N^2/2$ multiplications and additions to solve all of 46. Similarly, solving

47 requires about $N^2/2$ additions and multiplications. Thus, the dominant time in solving must be the time to decompose accord-

ing to Eq. 42.

The decomposition time is $\approx 5 N^3/3$, and the method for decomposing is described clearly in Numerical Recipes. The equations to be solved are

$N-1$

$[a_{lk}/h = A_{lk}] \quad (48)$

constituting $N^2 + N$ equations for N^2 unknowns. Thus we may arbitrarily choose $L_{11} = 1$.

$e_{jk} = A_{jk} - \sum_{l=1}^{j-1} L_{jl} A_{lk}$

Chapters 1-4: Unear Ngobra ScbndflcFORTH

The equations 48 are easy to solve if we do so in a sensible order. Clearly,

$1 \leq k \leq j$ (49) "knaos $k \leq j$

so we can divide up the work as follows: for each j , write

$L_{ml} = (A_{ml} - \sum_{k=1}^{j-1} L_{mk} A_{lk}) / e_{jj} \quad m = j+1, \dots, N$

(50)

$L_{1m} = (A_{1m} - \sum_{k=1}^{j-1} L_{1k} A_{lk}) / e_{jj}$, $m = j+1, j+2, \dots, N-1$. $\sum_{i=0}^{\infty} i=0$

Inspection of Eq. 50 makes clear that the terms on the right side are always computed before they are needed. We can store the computed elements A_{lk} and e_{jj} in place of the corresponding elements of the original matrix (on the diagonals we store 4).

since $A_{jj} = 1$ is known).

To limit roundoff error we again pivot, which amounts to permuting so the row with the largest diagonal element is the current one. Much of the code developed for the Gauss elimination method is applicable, as the file LU.FTH shows.