

On the Constructive Approach to Programming : The Case of "Partial Choleski Factorisation"  
 (A Tool for Static Condensation in Structural Analysis)

Alain BOSSAVIT  
 Bertrand MEYER

Electricité de France, Direction des Etudes et Recherches, Service IMA  
 1 avenue du Général de Gaulle 92141 CLAMART FRANCE

I - INTRODUCTION

From our experience in solving PDE's for engineering problems, we feel confident that the following statements will generally be agreed upon :

- There is a noticeable gap between *algorithms* presented in textbooks, and their implementation as *programs*.
- Too much unproductive work is done by structural engineers, physicists, numerical analysts, etc., to code again and again close variants of standard algorithms.

The two points are closely related, we think. The fact that, in spite of the existence and availability of popular finite-element codes, in spite of the broad marketing of well documented subroutine libraries, so many people continue to code their own versions of standard methods, is not to be blamed on human weaknesses only. It is also due to a tendency to reorganize old methods into new combinations in order to answer new questions. Thus, slight modifications of existing implementations are always required. But the general obscurity and poor readability of programs make these adaptations harder than they should be, and would be, if the gap we mentioned could be bridged.

Our contention is that this gap can be reduced, thanks to the advances in Computer Science in the recent years. We shall concentrate here on algorithms in linear algebra, relying on results obtained in the domain of proof-oriented program construction techniques, after the work of Hoare (2) and Dijkstra (1). In this approach, methods which were originally used for proving properties of existing programs are applied instead to the top-down design of new ones ; programs and their proofs will be developed in parallel.

The example we discuss is related to structural analysis by the finite-element method. It is often necessary to perform a "static condensation", which is nothing else than elimination of selected variables in a linear system. Though not new (3), the idea is currently in the development stage as far as software is concerned (4), and our implementation is (we hope) novel.

To avoid mixing all difficulties, we shall first state the basic concepts on a toy example (the square root), then proceed with the formal specification of the static condensation problem, the design of the program, and a brief account of the FORTRAN implementation.

II - PROGRAMS AS PREDICATE TRANSFORMERS - BASIC CONCEPTS

Underlined words below denote the few concepts we need. Consider the following text :

program square root (input  $a$  : REAL ; output  $x$  : REAL)

```
{P :  $a > 0$ }
  find  $x$ 
  {Q :  $x^2 = a$  and  $x > 0$ }
```

This is the text of a tautologically correct "program". The notation should be self-explanatory : sentences surrounded by braces { and } are treated as comments ; most often, they will be predicates, i.e. properties which must be verified by the variables of the program ( $x$  and  $a$  in our example). Texts such as find  $x$  are statements ; successive statements will be separated by semicolons. Proper indentation will help display the logical structure. A construct of the form {P} A {Q} is used to express that  $Q$  is true after execution of  $A$  if  $P$  was true before.

The approach known as top-down design calls now for refinements of the statement find  $x$ , which must eventually be expressed in terms of elementary actions. A seemingly very productive heuristics toward that end is uncoupling (or embedding) : let us introduce another variable  $y$ , also of type REAL, and a new predicate  $I$  which is a weaker form of the intended conclusion  $Q$

(I)  $xy = a$  and  $x > 0$  and  $y > 0$

so that  $Q = I$  and  $C$ , where  $C$  is

(C)  $x = y$  (within a prescribed margin of accuracy).

A possible way to proceed is to start with such a state of the variables where  $I$  is true, and, while keeping it true, to try to reach the "goal"  $C$ . One natural way to do this is to look for a program of the form

```
{P}
  establish I ;
  while not C do
    bring  $x$  and  $y$  closer to each other ;
    restore I
```

{C and I}

which we can assert is correct (i.e. will bring the variables from a state where the initial predicate is verified to a state where the final predicate is verified) provided the loop terminates properly.

This is shown by noting that  $I$  is a loop invariant, i.e.

$$(2) \{I \text{ and } \text{not } C\} A \{I\}$$

where  $A$  stands for the two statements of the loop body. Hoare's axioms for the while loop imply that, if the loop terminates, we can deduce from (2) that

$$(3) \{I\} \text{ while } \text{not } C \text{ do } A \{I \text{ and } C\}$$

which ensures the validity of program (1).

A third step in the top-down design will be to express more precisely how  $I$  is to be "established" (for instance by the two assignments  $x + a ; y + 1$ ), how  $x$  and  $y$  are to be "brought closer to each other" (one possible way is the assignment  $x + (x+y)/2$ ), and how  $I$  will then be "restored" ( $y + a/x$ ).

The program now obtained (a version of Newton's algorithm) is in a form very close to any current programming language, and we may stop here. But the process could be pursued further down if this proved necessary. For instance, if ordinary floating-point division were not available, we might have to refine the statement  $y + a/x$  into a loop so designed as to achieve the final predicate  $xy = a$ .

### III - STATIC CONDENSATION

We consider now the  $n \times n$  rigidity matrix of a structure, in block form

$$(4) \begin{vmatrix} A & B \\ B^t & C \end{vmatrix}$$

where the last blocks correspond to a subset  $E \subset [1, n]$  of the variables, containing only those which may interact with other parts of a higher-level structure, and are called "internal" for that reason. It is often useful to factorize such a matrix in the following way :

$$(5) \begin{vmatrix} A & B \\ B^t & C \end{vmatrix} = \begin{vmatrix} S & 0 \\ T & H \end{vmatrix} \begin{vmatrix} S^t & T^t \\ 0 & I \end{vmatrix}$$

The matrix  $H(H = C - TT^t)$ , which links displacements (or other generalized variables) of external type and related forces may be used in a higher-level assembly process and  $T$  and  $S$  allow one to go back to the elementary displacements once external variables are known. It seems interesting to allow "external variables" to include load parameters, Lagrange multipliers, right-hand members in general. Thus, a software tool capable of solving (5), which may be called a "condenser" (it is an extension of the classical "solver"), appears as an essential piece of any finite-element code where sub-structuring is to be implemented.

It is not practical, however, to number the variables in such a way that internal variables will always appear first as in (5) : identical sub-structures may have different sets of external variables (fig. 1), so that a different set of pointers to the same file must be used for each condensation. If internal variables are treated first (for example through a Cholesky routine),  $T$  being computed next, and then  $H$ , a lot of page faults will occur if the file is not entirely in core.

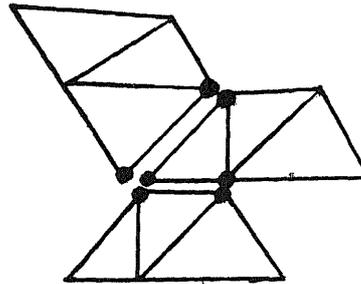


Figure 1 : Assembly of three identical substructures, with different sets of external variables

The problem to be solved may thus be stated as follows : given a  $n \times n$  symmetric matrix  $a$ , and a subset  $E$  of  $[1, n]$ , find a symmetric  $n \times n$  matrix  $s$  such that

$$I(n) \begin{cases} \text{if } 1 \leq j \leq i \leq n \text{ then} \\ \quad \text{if } j \notin E \text{ then} \\ \quad \quad \sum_{k \leq j} s_{ik} s_{jk} = a_{ij} , \\ \quad \text{if } j \in E \text{ and } i \notin E \text{ then} \\ \quad \quad \sum_{k \leq i} s_{ik} s_{jk} = a_{ij} , \\ \quad \text{if } j \in E \text{ and } i \in E \text{ then} \\ \quad \quad s_{ij} = a_{ij} - \sum_{k \leq n} s_{ik} s_{jk} \end{cases}$$

where  $\sum'$  stands for a summation where external indices are omitted. Entries  $s_{ij}$  correspond to  $S$  in (5) when  $i$  and  $j$  are both internal, to  $T$  when one of them is external, to  $H$  when both are external.

The above predicate, which we call  $I(n)$  for reasons which will be made clear below, is the final predicate of the program that is to be built. Obviously, it requires positive definitions of the internal-internal part of  $a$ , which will be assumed as initial predicate, although checked in practice by the program itself.

We shall add the following requirement : assuming that  $a$  and  $s$  are stored in "symmetric mode", i.e. line after line, in one-dimensional arrays, the processing should be as "sequential" as possible. This will guide the search for a suitable loop-invariant.

Only Cholesky factorization is considered below ; using the LDL<sup>t</sup> variant would not be more difficult.

### IV - TOP-DOWN DESIGN OF THE SELECTIVE FACTORIZATION

We embed the final predicate  $I(n)$  in a family of predicates  $I(l)$ , obtained by substitution of  $l$  for  $n$  in  $I(n)$  ;  $I(0)$  is trivially true.

Choosing the "uncoupling"

(6)  $I(n) = I(l)$  and  $(l=n)$

then leads to the following program :

```

program selective_Cholesky (input a : SYMMETRIC_MATRIX,
                           E : subset of INTEGER;
                           output s : SYMMETRIC_MATRIX)

```

```

variable l : INTEGER ;
l ← 0 ; {I(0)}
while l < n do
  l ← l+1 ;
  restore I(l)
{I(n)}

```

For the second refinement step, we look for a sequence A of actions such that

(7)  $\{I(l-1)\} A \{I(l)\}$

It is sufficient to find A such that :

(8) 
$$A \left\{ \begin{array}{l} 1 \leq j \leq i \leq l-1 \text{ and} \\ i \in E \text{ and } j \in E \Rightarrow s_{ij} = \sum_{k < j} s_{ik} s_{jk} \end{array} \right\}$$

$\{J(l)\}$

where  $J(l)$  is defined as

$J(l)$  
$$\left\{ \begin{array}{l} \text{if } 1 \leq j \leq l \text{ then} \\ \quad \text{if } j \notin E \text{ then } \sum_{k < j} s_{ik} s_{jk} = a_{lj} , \\ \quad \text{if } j \in E \text{ and } l \notin E \text{ then} \\ \quad \quad \sum_{k < l} s_{lk} s_{jk} = a_{lj} \\ \quad \text{if } j \in E \text{ and } l \in E \text{ then} \\ \quad \quad \text{for all } i > j \text{ and } i \in E \\ \quad \quad \quad s_{ij} = a_{ij} - \sum_{k < l} s_{ik} s_{jk} \end{array} \right.$$

In order to solve problem (8), we again use the "uncoupling" strategy. Let  $l$ , an *INTEGER*, and  $y$ , a *REAL*, be new variables. We can express  $J(l)$  as :

(9)  $J(l) = J(l, r, y)$  and  $r=l$  and  $y=s_{ll}$  and  $K(l)$

where  $J(l, r, y)$  is :

$J(l, r, y)$  
$$\left\{ \begin{array}{l} \text{if } i < j \leq r \text{ and } r \leq l \text{ then} \\ \quad \text{if } j \notin E \text{ then } \sum_{k < l} s_{lk} s_{jk} = a_{lj} , \\ \quad \text{if } j \in E \text{ and } l \notin E \text{ then} \\ \quad \quad \sum_{k < l} s_{lk} s_{jk} + y s_{jl} = a_{lj} \end{array} \right.$$

and  $K(l)$  as :

$K(l)$  
$$\left\{ \begin{array}{l} \text{if } l \notin E \text{ and } 1 \leq j \leq i \leq l \text{ and } i \notin E \\ \text{and } j \notin E \\ \text{then} \\ \quad s_{ij} = a_{ij} - \sum_{k < l} s_{ik} s_{jk} \end{array} \right.$$

Examination of (9) suggests that action A should be constructed as the composition of a *while* loop admitting  $J(l, r, y)$  as an invariant, and suitable actions which will ensure  $y = s_{ll}$  and  $K(l)$ . In other words, we shall look for an A of the following form :

A 
$$\left\{ \begin{array}{l} r \leftarrow 0 ; y \leftarrow 1 ; \\ \text{while } r < l \text{ do} \\ \quad r \leftarrow r+1 ; \\ \quad \text{restore } J(l, r, y) ; \\ \text{ensure } y = s_{ll} \text{ while maintaining } J(l, l, y) ; \\ \text{ensure } K(l) \end{array} \right.$$

We could pursue the process further, and show how the last sub-actions of A can be systematically developed ; the approach of proof-directed top-down program design should be clear by now, however, and we shall give the final product without further justification. The program for selective Choleski factorization, expressed in our notation, appears on the next page.

program selective\_choleski

(input a : SYMMETRIC MATRIX,  
 e : FUNCTION (INTEGER → LOGICAL) ;  
 output s : SYMMETRIC MATRIX)

{input assertion :  
 the internal-internal part of a is  
 positive-definite}

variables

n, l, r, i, j : INTEGER,  
 y : REAL ;

n ← order (a) ; l ← 0 ;

{I(l)}

while l < n do {loop invariant : I(l)}

l ← l+1 ;

{restore I(l) : }

r ← 0 ; y ← 1 {J(l,r,y) is satisfied}

while r < l do {loop invariant :  
 J(l,r,y)}

r ← r+1 ;

{restore J(l,r,y) : }

if not e(v) then

if r < l then

$$s_{lr} ← (a_{lr} - \sum_{k<r} s_{lk} s_{rk}) / s_{rr}$$

else

$$s_{ll} ← \sqrt{a_{ll} - \sum_{k<l} s_{lk}^2}$$

{positive definiteness  
 guarantees that the square  
 root is defined}

elseif not e(l) then

$$s_{lr} ← (a_{lr} - \sum_{k<l} a_{lk} a_{rk}) / y ;$$

if not e(l) then

y ← s<sub>ll</sub> ;

for 0 < r < l while e(r) do

$$s_{lr} ← s_{lr} / y ;$$

for 1 ≤ j ≤ i ≤ l  
 while e(i) and e(j)  
 do

$$s_{ij} ← s_{ij} - s_{li} s_{lj}$$

{J(l,r,y)}

{I(l)}

V - A FORTRAN IMPLEMENTATION OF THE SELECTIVE  
 FACTORIZATION.

Included below is a FORTRAN program which is little more than a literal translation of the above routine. It embodies a few implementation decisions ; in particular :

- s and a have been stored as one-dimensional arrays, where the relative position of an element (i, j) is computed via a function ADDRESS (I, J). In the implementation below, it is assumed that s and a are full and fit in core. If skyline is considered, and/or if secondary storage is required, only ADDRESS has to be modified.

- the  $\sum'$  operator is effected through a call to the SCLPRD subroutine, which must be consistent with ADDRESS (pursuing the approach further, we could have defined these functions as belonging to the same "virtual machine", or "abstract data type").

VI - CONCLUSION

We hope to have shown that program design, all the way from the search for an algorithm to the writing and documenting of the final code, can proceed in a disciplined way not unlike the classical mathematical discourse.

We certainly do not mean to imply, however, that this is a smooth and easy process ; indeed, there are choices involved at each step, and sometimes there is no better way to resolve them than trial and error. Uncoupling (9), for instance, requires a good deal of insight, and we had to try a few blind alleys before we found it. Top-down proof-directed program design is certainly not a magical recipe for solving numerical problems ; what the method does bring, however, is a better control of the whole programming process, a better relationship between mathematical methods and their implementations, and, of course, better programs, easier to code, debug, document and modify.

REFERENCES

- (1) Dijkstra, E. W. : *A Discipline of Programming* ; Prentice-Hall, 1976.
- (2) Hoare, C. A. R. : *An Axiomatic Basis for Computer Programming* ; CACM, 12, 10, 1969, pp. 576-583.
- (3) Kron, G. : *Diakoptics, Piecewise Solution of Large Systems* ; London, McDonald, 1963.
- (4) Wilson, E. L., and Dovey, H. H. : *Solution or Reduction of Equilibrium Equations for Large Complex Structural Systems* ; *Advances in Engineering Software*, 1, 1, 1978, pp. 19-26.

