UNIVERSITY OF NEWCASTLE-UPON-TYNE CIVIL ENGINEERING DEPARTMENT

COMPUTER MODELLING OF WATER SUPPLY DISTRIBUTION NETWORKS USING THE GRADIENT METHOD.

ЪУ

Rubén O. Salgado-Castro

.

VOLUME TWO ******

APPENDICES

NEWCASTLE UNIVERSITY LIBRARY 088 22971 9 THESIS L3421

Thesis submitted in fulfilment of the requirements for the degree of Doctor of Philosophy in Civil Engineering.

November, 1988.

TABLE OF CONTENTS

VOLUME TWO ******

List	of	Figures	iii
List	of	Tables	v
List	of	main Symbols	vii

.

APPENDIX A

EFFICIENT SOLUTION OF LINEAR SYSTEMS OF EQUATIONS IN THE CONTEXT OF THE GRADIENT METHOD FOR THE ANALYSIS OF WATER SUPPLY DISTRIBUTION NETWORKS

A.1.	Intro	luction	A-1
	A.1.1.	Need for efficient linear solvers in water	
		distribution network analysis	A-1
	A.1.2.	Different classes of linear problems	A-2
•	A.1.3.	Direct methods and iterative methods for	
		solving determined systems of linear	
		equations	A-10
	. .		
A.Z.	Review	v of direct methods for the solution of	
	gener	al dense linear systems of equations	A-13
	A.2.1.	Gaussian elimination	A-13
	A.2.2.	Gauss-Jordan elimination	A-28
	A.2.3.	Matrix factorization methods	A-31
	A.2.4.	Comparison of the different algorithms for	
		the direct solution of general dense linear	
		systems	A-45
	A.2.5.	Error analysis in the solution of linear	
		systems	A-48
A.3.	Review	w of sparse direct methods for the solution of	
	linea	r systems of equations	A-57
	A.3.1.	Data structures	A-57
	A.3.2.	Fill-in	A-62
	A.3.3.	Sparse methods for banded matrices	A-69
	A.3.4.	The envelope ("skyline") method	A-71
	A.3.5.	Minimum degree algorithms	A-77
	A. 3. 6.	Quotient tree algorithms	A-89
	A.3.7	One-way dissection methods	A-104
	A 3 8	Nested dissection methods	A-109
	A 3 Q	Frontal and multifrontal methods	A_110
	Ω.υ.ν.	FIGHERI AND MUTCHTONERI MECHOUS	A-110
A. 4	. Review	w of iterative methods for the solution of	
	linea	r systems of equations	A-115
	A.4.1.	Introduction	A-115
	A.4.2.	Jacobi's method (method of simultaneous	
		displacements)	A-117

A.4.3.	Gauss-Seidel method (method of successive	
	displacements)	A-119
A.4.4.	Successive over (or under) relaxation method	A-120
A.4.5.	Relationship between the previous methods	A-120
A.4.6.	Convergence conditions for the previous	
	methods	A-124
A.4.7.	Standard conjugate gradient method for the	
	solution of linear systems of equations.	
	(Hestenes and Stiefel, 1952)	A-127
A.4.8.	Preconditioning	A-133
A.4.9.	Preconditioned (modified) conjugate gradient	
	method for the solution of linear systems	A-142

APPENDIX B

DERIVATION OF THE KRIGING ESTIMATOR EQUATIONS

B.1. Introd	luction	B-1
B.2. Statis	tical inference	B-5
B.3. Estima B.3.1. B.3.2.	Ation via Kriging Kriging under stationary conditions Relaxing the second-order stationarity	B-16 B-16
B.3.3.	condition: the intrinsic hypothesis Introducing uncertainty in the measurement	B-25
B.3.4.	data Universal Kriging and k-th order random	B-30
	functions (k-IRF)	B-32

APPENDIX C

DERIVATION OF THE BI-CUBIC SPLINES APPROXIMATION EQUATIONS

C.1.	Introducing cubic splines	C-1
C.3.	Data fitting using one-dimensional cubic splines	C-8
C.4.	Data fitting using bi-cubic B-splines	C-13
C.5.	The statistical problem in splines fitting	C-19

APPENDIX D

NUMERICAL RESULTS OF THE CALIBRATION EXERCISE D-1

FIGUE	?E		PAGE
Fig. Fig. Fig.	A.1. A.2. A.3.	Different linear systems of equations Arrow shaped linear system Cholesky factor of arrow shaped matrix	A-4 A-63
Fig. Fig	A.4.	(without reordering) Reordered arrow shaped linear system Cholesky factor of reordered arrow-shaped	A-63 A-63
Fig.	A. 6.	linear system Undirected graph corresponding to symmetric	A-64
Fig.	A.7.	matrix given by equation (98) Permutation matrix for interchanging rows 3	A-65
Fig.	A.8.	and 4 Reordered matrix [eq.(98)] and its	A-66
Fig.	A.9.	associated graph Diagonal storage scheme for banded matrices	A-66 A-69
Fig.	A. 10.	Minimum Dand ordering for arrow snaped linear system Two examples of Cuthill-McKee orderings for	A-70
Fig.	A. 12.	reducing bandwidth, from George (1981) Skyline representation	A-71 A-72
Fig.	A.13.	Graph for level structure example, from George and Liu (1981)	A-74
Fig.	A.14.	Level structure rooted at x5, for graph shown in Fig. A.13, from George and Liu	A 774
Fig.	A.15.	Algorithm for finding a pseudo-peripheral node, from George and Liu (1981)	A-76
Fig.	A.16.	Example of application of minimum degree algorithm, from George (1981)	A-81
Fig.	A.17.	Elimination graphs representing the process of creation of fill-in	A-82
Fig. Fig.	A. 18. A. 19.	The filled graph of F = L + L ¹ Filled (reordered) matrix	A-83 A-83
Fig.	A.20.	Sequence of quotient graphs Γ_i used for finding the reachable set of nodes in the elimination process	A-87
Fig.	A.21.	An example of a monotonely ordered tree (rooted at node 8)	A-96
Fig.	A.22.	Matrix corresponding to the monotonely ordered tree shown in Fig. A.21	A-97
Fig. Fig.	A.23. A.24.	Network example for quotient tree methods Quotient tree corresponding to the tree of	A-98
Fig.	A.25.	Matrix corresponding to the quotient tree	A-100
Fig.	A.26.	Level structure rooted at node "a" and its	A-101
Fig.	A.27.	Matrix corresponding to refined quotient tree of Fig. A.26. b)	A-102
Fig.	A.28.	Level structure rooted at node "t" and its corresponding quotient tree	A-103
Fig.	A.29.	Matrix corresponding to refined quotient tree of Fig. A.28. b)	A-103
rig.	A.30.	Rectangular grid partitioned with 2 dissectors	A-105

.

FIGUI	RE		PAGE
Fig.	A.31.	Example of a 40 nodes rectangular shaped	
Fig.	A.32.	graph, partitioned using one-way dissection Matrix structure corresponding to graph of	A-106
- Fia	A 33	Fig. A.31	A-107
r 1 g .	A. JJ.	graph, partitioned using nested dissection	A-109
Fig.	A.34.	Matrix structure corresponding to graph of Fig. A.33	A-109
Fig.	A.35.	Finite element definition and ordering for a frontal solution scheme, from Livesley	
Fia	A 36	(1983)	A-112
E.	A. 30.	solution scheme, from Livesley (1983)	A-113
Fig.	A.37.	Geometric interpretation of the convergence of the conjugate gradient algorithm for n=2 and n=3, from Gambolatti and Perdon (1984),	
Fie.	R 1	solving A h = b Estimated semi-variogram	A-134 B-9
Fig.	B.2.	Nugget effect and corresponding true semi-	ЪŪ
Fiø	ЪЗ	variogram	B-13
6.	D.J.	variogram	B-17
Fig.	C.1.	Representation of a B-spline	C-4
LTR.	C.2.	B-splines involved in the spline function for the interval $(: : : : :)$	C-5
Fig.	C.3.	Extended set of B-splines: interior and	
Fig.	C 4	exterior knots	C-6
0 •	V. I.	variables x and y, divided into panels R_{ij}	
		by knots λ_i i=1,2,h+1 and ν_j j=1,2,k+1	C-15

TABLE			PAGE
Table	A.1.	Comparison of the algorithms used in the direct solution of dense linear systems of	
Table	B.1.	equations Computation of the experimental semi-	A-46
Table	B.2.	variogram Basic analytical models for the experimental	B-8
Table	B.3.	semi-variogram Possible polynomial models for generalised	B-14
Table	D.1.	covariances, from Delhomme (1978) Summary of the comparison between true, initial and calibrated piezometric heads.	B-37
Table	D.2.	Example A	D-1
Table	D.3.	Example B Summary of the comparison between true, initial and calibrated piezometric heads	D-2
Table	D.4.	Example C Summary of the comparison between true, initial and calibrated piezometric heads.	D-3
Table	D.5.	Example C	D-4
Table	D.6.	Example E	D-5
Table	D.7.	Example F Summary of the performance of the calibrated heads Example A	D-6
Table	D.8.	Summary of the performance of the calibrated heads. Example B	ה- פ ח-8
Table	D.9.	Summary of the performance of the calibrated heads Example C	D-9
Table	D.10.	Summary of the performance of the calibrated heads. Example D	D-10
Table	D.11.	Summary of the performance of the calibrated heads. Example E	D-11
Table	D.12.	Summary of the performance of the calibrated heads. Example F	D-12
Table	D.13.	Summary of the comparison between true, initial and calibrated flows. Example A	– –– D–13
Table	D.14.	Summary of the comparison between true, initial and calibrated flows. Example B	D-14
Table	D.15.	Summary of the comparison between true, initial and calibrated flows. Example C	D-15
Table	D.16.	Summary of the comparison between true, initial and calibrated flows. Example D	D-16
Table	D.17.	Summary of the comparison between true, initial and calibrated flows. Example E	D-17
Table	D.18.	Summary of the comparison between true, initial and calibrated flows Example F	D-18
Table	D.19.	Summary of the performance of the calibrated flows. Example A	D-19

TABLE		PAGE
Table D.20.	Summary of the performance of the calibrated flows. Example B	D-20
Table D.21.	calibrated flows. Example C	D-21
Table D.22.	Summary of the performance of the	
Table D.23.	Summary of the performance of the	D-22
Table D.24.	calibrated flows. Example E	D-23
Table D 25	calibrated flows. Example F	D-24
Table D.23.	initial and calibrated C's. Example A	D-25
Table D.26.	Summary of the comparison between true, initial and calibrated C's. Example B	D-26
Table D.27.	Summary of the comparison between true,	
Table D.28.	Summary of the comparison between true,	D-27
Table D.29.	initial and calibrated C's. Example D Summary of the comparison between true.	D-28
m-11- D 20	initial and calibrated C's. Example E	D-29
Table D.30.	initial and calibrated C's. Example F	D-30
Table D.31.	Summary of the performance of the calibrated C's. Example A	D-31
Table D.32.	Summary of the performance of the	
Table D.33.	Summary of the performance of the	D-32
Table D 34	calibrated C's. Example C	D-33
	calibrated C's. Example D	D-34
Table D.35.	calibrated C's. Example E	D-35
Table D.36.	Summary of the performance of the calibrated C's Example F	D-36
	currented o bi brampto r iiiiiiiiiiiiiiiiiiiiii	<i>D</i> 00

.

LIST OF MAIN SYMBOLS

APPENDIX A

Symbol	Description
А, В	A capital letter denotes a <u>matrix</u> , particularly:
L	Denotes a lower triangular matrix.
U	Denotes an upper triangular matrix.
D	Denotes a diagonal matrix.
a, <u>A</u> , <u>a</u>	A lower case letter, or a capital (or lower case)
	<u>underlined</u> letter denotes a <u>vector</u> .
α, μ	A greek letter denotes a <u>scalar</u> .
ai	A lower case subindexed letter denotes a scalar,
	normally the i-th component of a vector.
a _{ij}	A lower case doubly subindexed letter denotes a scalar,
	normally the coefficient of a matrix A, located at the
	intersection of row "i" and column "j".
Ai	A capital subindexed letter denotes either a matrix in
	a sequence of matrices or a block-partitioned matrix.
(.) ^T	Denotes transposition, either of a vector or matrix.
A ⁻¹	Denotes matrix inversion.
(.)(k)	Refers to the k-th value of the variable (scalar,
	vector or matrix) within an iterative procedure.
x	Denotes the norm of a variable (vector or matrix).
x	Denotes the absolute value of a scalar.
λi	Denotes a particular eigenvalue of a matrix.
(A)	Spectral radius of?matrix A.
<x, y=""></x,>	Inner (scalar) product of vectors "x" and "y":
	$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = \mathbf{y}^T \mathbf{x} = \sum_{i=1}^n \mathbf{x}_i \mathbf{y}_i$

vii

Sym	bol		Descr	iption			_		
(<u>χ</u> ,	H)	State	variable	representing	the	piezometric	head	at	any
		point	$\underline{\mathbf{x}}$ of the	domain:					

 $(x, H) = (x_1, x_2, H)^T$ in a two-dimensional space.

- $Z(\underline{x}, H)$ Random function representing the value of the state variable H at a location \underline{x} (in general).
- $Z(\underline{x}_{O}, H)$ Random variable (at the particular location $\underline{x}=\underline{x}_{O}$. Simplified notation:

$$Z(x_i, H) \iff Zx_i \iff Z_i$$

E[Z(x_1 , H)] expected value of the random variable at $\underline{x}=\underline{x}_1$. $\Gamma(\underline{x}_1, \underline{x}_2)$ Semi-variogram:

$$\Gamma(x_1, x_2) \equiv \frac{1}{2} E[\{Z(x_1, H)Z(x_2, H)\}^2]$$

 $Cov(Z(x_1, H), Z(x_2, H))$ Covariance:

$$Cov(Z(x_1,H),Z(x_2,H)) = E[Z(x_1,H)Z(x_2,H)]$$

 $Var(Z(\underline{x},H))$ Variance of the random variable:

$$Var(Z(\underline{x},H)) = E[Z(\underline{x},H)^2] - \{E[Z(\underline{x},H)]\}^2 = \sigma^2$$

 χ_0^{i} Unknown weight of the Kriging estimator, where"i" refers to the measurement Y_i and "o" to the point being estimated χ_0 , i.e.:

$$Y_0^* = Y_0^*(x_0) = \sum_{i=1}^n x_0^i Y_i$$

APPENDIX C

Symbol	Description
One dimen	sional splines:
s(x)	General spline polynomial function.
×i	Set of knots associated to the spline.
h	Number of knots (χ_i , i=1,2,,h)
(x _r , f(x _r)) , r=1,2,,m set of data points.
$M_i(x)$	Cubic Basic (or fundamental) spline (B-spline).
$M_{n,i}(x)$	General B-spline of degree (n-1). For a cubic spline
	n=4, i.e.: $M_{4,i}(x) = M_i(x)$ (simplified notation).
ri	Linear (unknown) weights, combining the cubic B-splines
	into a general cubic spline polynomial:



 \underline{r} (h+4)*1) column vector of linear weights r_i .

A Matrix of observation equations (observation matrix), of order m*(h+4) in a cubic spline: $A_{r,i} = M_i(x)$.

A $\underline{\Gamma} = \underline{f}$ Observation equations, overdetermined system.

f (m*1) column vector with the value of the spline polynomial at each data point:

 $f_i = f(x_i) = s(x_i)$ i=1,2,...,m

 $A^{T}A\underline{\Gamma}=A^{T}\underline{f}$ Normal equations. Linear system of (h+4)*(h+4) equations (determined system).

Bi-cubic splines:

- s(x,y) General spline polynomial function.
- λ_i Set of knots associated to the independent variable x. μ_j Set of knots associated to the dependent variable y.

Symbol	Description									
h	Number of knots in the x-axis (λ_i , i=1,2,,h)									
k	Number of knots in the y-axis (μ_j , j=1,2,,k)									
R _{ij}	Set of (h+1)*(k+1) panels in which the x-y space is									
	sub-divided.									

 $(x_r, y_r, f(x_r, y_r))$, r=1,2,...,m set of data points.

- $M_i(x)$ Cubic B-spline related to the independent variable x.
- $N_j(y)$ Cubic B-spline related to the dependent variable y.
- r_{ij} Linear (unknown) weights, combining the cubic B-splines into a general cubic spline polynomial:

$$s(x,y) = \sum_{i=1}^{h+4} \sum_{j=1}^{k+4} r_{ij} M_i(x) N_j(y)$$

- $\underline{\Gamma}$ ((h+4)(k+4)*1) column vector of linear weights Γ_{ij} .
- A Matrix of observation equations (observation matrix), of order m*(h+4)(k+4).
- $A \underline{\Gamma} = \underline{f}$ Observation equations, overdetermined system.
- f (m*1) column vector with the value of the spline polynomial at each data point:

 $f_i = f(x_i) = s(x_i)$ i=1,2,...,m

 $A^{T}A\underline{\Gamma}=A^{T}\underline{f}$ Normal equations.Linear system of (h+4)(k+4)*(h+4)(k+4) equations (determined system).

APPENDIX A

EFFICIENT SOLUTION OF LINEAR SYSTEMS OF EQUATIONS IN THE CONTEXT OF THE GRADIENT METHOD FOR THE ANALYSIS OF WATER SUPPLY DISTRIBUTION NETWORKS

A.1. Introduction.

A.1.1. Need for efficient linear solvers in water distribution network analysis.

From a mathematical viewpoint, we have already seen in Chapter Two that the water distribution network analysis problem reduces to the solution of a set of non-linear simultaneous equations; different ways of assembling these equations lead to different network analysis algorithms.

Because a direct solution for simultaneous non-linear systems does not exist, the usual approach is to transform the problem into a sequence of linear problems, whose solution approximates step by step the solution of the original non-linear set.

As a result, for solving a non-linear system of equations we have to solve a number of linear systems and consequently, very efficient linear solvers are required in order to keep the computational resources needed under a reasonable limit.

Thus, network analysis algorithms rely heavily on the solution of linear systems of equations and, consequently, efficient linear solvers which can take advantage of the particular

A-1

features of the linear systems generated in network analysis are needed. In the context of the gradient network analysis method, we shall be especially interested in linear solvers for <u>symmetric</u> <u>positive definite systems</u>.

This Appendix presents a review of the methods currently available for the solution of linear systems of equations.

Because most of the real water distribution networks generate linear systems in the range of hundreds of unknowns, and since the matrices produced are sparse (i.e. with very few non-zeros per row), special consideration will be paid to those methods allowing the efficient storage and handling of matrices of this kind.

A.1.2. Different classes of linear problems.

In its most general format, the solution of a linear system can be stated as that of finding the "n" unknown values, x_i 's, in the following system of "m" linear equations:

 $a_{m1} x_1 + a_{m2} x_2 + a_{m3} x_3 + \ldots + a_{mn} x_n = b_m$ where the coefficients a_{ij} and the right hand side terms b_i are all known. These coefficients normally emerge from some physical property of the system under study.

In order to simplify the notation we shall use, whenever possible, capital letters to denote matrices, lower case letters

to represent vectors and greek letters to denote scalar numbers. In some instances, when the use of lower and upper case letters to differentiate vectors from matrices is not possible (or when it is not convenient), for the sake of clarity we shall denote vectors by an underlined letter (e.g.: P, Q, H, h, g are all vectors). Also notice that, in general, sub-indexed lower case letters will represent scalars. Thus, the linear system (1) can be expressed in matrix notation as:

$$A x = b \tag{2}$$

where:

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ \vdots & & \vdots & & \vdots \\ a_{m1} & a_{m2} & a_{m3} & \dots & a_{mn} \end{bmatrix}$$

A : is the <u>matrix of coefficients</u> and it has "m" rows and "n" columns, i.e. it is an (m x n) matrix.

 $x = (x_1, x_2, x_3, \dots, x_n)^T$

x : is the vector of unknowns, a column vector in the space \mathbb{R}^n ; this is a (n x 1) vector.

and

 $b = (b_1, b_2, b_3, \dots, b_m)^T$

b : is the <u>right hand side vector</u>, a column vector in the space \mathbb{R}^m ; this is a (m x 1) vector.

Depending on the relative values of "m" and "n", we may face different kinds of problems:

i) Underdetermined systems of equations:

In this particular case m < n. Clearly, in this case we have less equations than unknowns and the system has, in general, more than one "exact" solution, by which we mean that the equality in (1) or (2) holds.

A very simple example of an underdetermined system can be shown in a 2-dimensional plane, as in Fig. A.1.

The straight line L_1 in Fig. A.1 a) represents a "system" of just one equation in the unknowns x_1 and x_2 . The solution points of this system are simply all the infinite points on L_1 .

In underdetermined systems it is possible to find two cases:

a) Homogeneous case: where the right hand side (RHS) vector b is null (i. e. b = Q, or $b_i = 0$ i=1,2, ..., n).



linear equations.



Fig. A.1. Different linear systems of equations.

b) Non-homogeneous case: where some $b_i \neq 0$ (at least one of them must be non-zero).

In the previous example L_1 represents a non-homogeneous problem if and only if $c \neq 0$. If c = 0, we are in the homogeneous case and the straight line becomes $L_2 = a x_1 + b x_2 = 0$, where L_2 passes now through the origin, as shown in Fig. A.1. b)

In practice, it is not often that we need to solve an underdetermined system of equations, but one of these cases arises when, for example, we are trying to find an initial flow distribution in a pipe network, where we need to guarantee the fulfilment of the node balance; in this case, and using the notation defined in Chapter Two (Section 2.3.7), we need to solve the system:

$$A_{21} Q = q \tag{3}$$

where:

 $A_{21} = A_{12}^{T}$; with A_{12} the topological (connectivity) matrix. A_{12} is of order (NP x (NN-NS)), where NN is the number of nodes, NS is the number of sources and NP is the number of branches (pipes, valves, etc.). Then, A_{21} is an ((NN-NS) x NP) matrix.

and

Q: (NP x 1) vector of initial flows (unknown).

q: ((NN-NS) x 1) vector of nodal consumptions, which are given as data. This vector can be zero or non-zero, leading either to homogeneous or non-homogeneous problems.

We shall not study in more detail the solution of this kind of problem. O'Neil (1983) and George, Heath and Ng (1984), give a more detailed treatment of this problem.

ii) Overdetermined systems of equations.

In this particular case we have m > n, i.e. we have more equations than unknowns. No exact solution is possible and the equality in (1) and (2) does not hold.

the context of our small previous example, this can In be represented in a 2-dimensional space by a set of 3 straight lines L1, L2 and L3, where it is clear that no single point can satisfy simultaneously the 3 equations: point A represents the simultaneous fulfilment of L1 and L2, B that of L1 and L3 and C the corresponding to L2 and L3. See Fig. A.1. c)

The best we can do in an overdetermined problem is to find a solution that minimizes the difference between the RHS vector and the product A x. In other words we can minimize the residual vector: r

$$c = A x - b \tag{4}$$

A well-known technique for finding a solution to this problem is the "least square" algorithm, where the unknown "x" is chosen so that the square of the euclidean norm of "r" is minimized:

$$\min || r ||_{2}^{2} = r^{T}r = \sum_{i=1}^{m} r_{i}^{2} = \sum_{i=1}^{m} (\sum_{j=1}^{n} x_{j} - b_{j})^{2}$$
(5)

There is a wide variety of numerical methods available to solve this problem, but we shall not include them here; for details see for example Duff and Reid (1976), Golub and van Loan (1983).

In practice, as before, it is not common to have to solve overdetermined systems in civil engineering problems. Nevertheless, one example arises in estimation problems, when using weighted least squares estimators with less parameters to be estimated than the number of measurements available. With this same kind of problem we can get an underdetermined system if the number of measurements is less than the number of parameters to be estimated and also, we can get a determined system when both the number of measurements and parameters to be estimated In either underdetermined or overdetermined systems, coincide. because the matrices are not square matrices, their inverse does not exist and the concept of "pseudo-inverse" or "generalised inverse" is used instead; see Gelb (1974), Rao (1962) and Penrose (1955) for details.

In the context of an explicit calibration algorithm for water distribution networks [see Chapter Six], we have used a bi-cubic splines approach to approximate the piezometric head of unmeasured nodes, based on some measured heads. The estimation of the coefficients of the spline polynomial leads to an overdetermined system of equations, which is solved via the least-square technique (normal equations); see Appendix C for details on the cubic splines fitting.

iii) Determined systems of equations.

In this particular case m = n, then we have as many equations as unknowns. The matrix A is said to be square of order "n" (or "m").

In the context of our small example in a 2-dimensional space, the determined case corresponds to the problem of finding the intersection of two straight lines L1 and L2: i.e. the coordinates of point A: (x_1^*, x_2^*) , as shown in Fig. A.1. d)

The condition for the existence of a solution is that A must be <u>invertible</u>, this means that, from the mathematical point of view, one of the following statements must be true:

- a) The determinant of A is non-zero (i.e. A is non-singular).
- b) The rank of the matrix is equal to the order of the matrix (the rank is defined as the order of the largest submatrix of A with a non-vanishing determinant).

A simple representation of a non-invertible A matrix, arises when the straight lines of the previous figure are parallel. Numerical difficulties are found when both lines are "almost" parallel.

Most of the engineering problems involve the solution of determined systems of linear equations, and different methods of solution have been proposed, depending on the particular properties of the matrix A (symmetry, sparseness, positive or negative definiteness, etc.) and the computational resources available.

Since real engineering problems lead generally to positive definite systems, we shall concentrate our attention on methods for solving such problems, although different methods are available to solve systems of equations with indefinite matrices

A-8

[see Duff, (1980)].

The model problem for a determined linear system of equations is:

Solve: A x = b (6) where:

A : is a known non-singular square matrix of order "n",

b : a known (n x 1) column vector.

We need to find the vector $x = (x_1, x_2, ..., x_n)^T$, which satisfies equation (6).

The structure and numerical values of the matrix A and vector b are usually obtained from a mathematical model of the physical problem being studied.

From a strictly algebraic point of view, the solution of the linear system represented by equation (6) can be easily obtained if we know the inverse matrix of A, say A^{-1} , (A^{-1} is such that $A A^{-1} = A^{-1} A = I$: the identity matrix); thus, premultiplying equation (6) by A^{-1} , we get :

or

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

$$x = A^{-1} b$$
(7)

Although from the algebraic viewpoint the solution looks quite simple, the process of inverting a matrix (especially a large one, as is usual in many practical problems) is a very expensive process, demanding a significant amount of work (n^3 mathematical operations) and storage (n^2 storage positions, since the inverse matrix can be completely filled). As we shall see, looking at the most computationally efficient way of solving equation (6), the explicit inversion of A is not needed and a variety of methods can be used, according to the specific properties of A. If for some special reason the inverse of the matrix is needed, the efficient way of doing it is through the repeated application of our model problem (6), changing the right hand side vector, but maintaining the matrix A.

A.1.3, Direct methods and iterative methods for solving determined systems of linear equations.

From the numerical standpoint, a first sub-division of methods for solving a linear system of equations like (6) arises:

a) Direct methods: where an <u>exact</u> solution of the linear system is reached in a finite number of steps (or elementary arithmetic operations), provided that no rounding errors exist.

Usually direct methods are associated with algorithms like Gaussian or Gauss-Jordan elimination.

b) Iterative methods: where starting from an approximate initial solution (say x = Q, or a better one if available), an appropriate algorithm is applied in order to improve the initial solution and obtain an <u>approximate</u> solution of the linear system.

In this case, the number of steps and elementary arithmetic operations is <u>not finite</u> and it depends on various factors, e.g.: on the algorithm used, on how approximate the final solution required is, or on how close is the initial solution to the final

one.

In general, an approximate solution to the linear system at the k-th step of the iterative process will be denoted by $x^{(k)}$ and is such that:

$$\mathbf{x}^{(\mathbf{k})} \approx \mathbf{b}$$
 (8)

In this iterative context, we can define a residual vector "r" such that:

$$\mathbf{r} = \mathbf{b} - \mathbf{A} \mathbf{x}^{(\mathbf{k})} \tag{9}$$

and the iterative procedure for solving the linear system has the goal of producing a residual vector whose norm (or "length", roughly speaking) is smaller than some pre-specified value (accuracy). In other words, the iterative process will be stopped when the approximate solution $x^{(k)}$ produces a residual vector whose norm is smaller than a pre-specified accuracy (scalar), say ξ , i.e.

$$b - A x^{(k)} = r$$
 such that $||r|| < \epsilon$ (10)

 $x^{(k)}$ will be the solution corresponding to this level of accuracy. However, equation (10) relies on the definition of a norm for the residual vector, which can be any one of a variety of norms available in vector analysis, such as:

i) Euclidean norm :
$$||\mathbf{r}||_2 = \sqrt{\mathbf{r}^T \mathbf{r}}$$
 (11)

ii) Maximum norm :
$$||\mathbf{r}||_{\infty} = \max_{\substack{\mathbf{r} \\ 1 \le i \le n}} |\mathbf{r}_i|$$
 (12)

iii) Weighted p-norm :

$$||\mathbf{r}||_{p,w} = (\sum_{i=1}^{n} |\mathbf{r}_{i}|^{p} w_{i})^{(1/p)}$$
(13)

Iterative methods for the solution of linear systems are associated with algorithms like Jacobi, Gauss-Seidel, etc. We shall review most of these methods in some detail in the next sections.

In practice, and since computers do work with a finite word length, direct methods do not produce exact solutions in a finite number of arithmetic operations (effect of roundoff errors). As a result, in the computational sense, all the methods are in practice iterative. Nevertheless, we shall maintain the distinction between both approaches since it is important from the conceptual point of view.

The decision on which approach to be used depends on the characteristics of the matrix A and the computer resources available. Thus, it has been accepted for a long time that, for sparse matrices, iterative methods are the best ones, both from the storage and amount of work point of view; nowadays the availability of efficient direct solutions challenges this assertion. For dense matrices it has been accepted that direct methods are the best choice [Conte & de Boor, (1983)], although either big storage computers or special out-of-core data handling routines should be available for really big problems. A considerable amount of research has been carried out in the last decade in this field, resulting in new efficient methods becoming available, both in the direct and iterative fields.

Since water distribution network analysis using the gradient method produces symmetric linear systems of equations, we shall concentrate our search for efficient linear system solvers on these special systems.

A.2. Review of direct methods for the solution of general dense linear systems of equations.

A.2.1. Gaussian elimination.

i) Introduction.

The basic idea behind all kinds of elimination methods can be expressed in the following terms. Let us suppose that we have a linear system of equations such that:

$$T x = b \tag{14}$$

where:

T : is a triangular matrix of coefficients, i.e.: only the upper or lower part has non-zeros. Let us assume that we are dealing with an upper triangular matrix U:

$$\mathbf{T} = \mathbf{U} \tag{15}$$

then, the problem (14) can be represented as:

$$U x = b \tag{16}$$

The system (16) is then of the form:

It is easy to see that, from the last equation of the system (17), the solution for the last component of "x" can be obtained in a direct manner:

$$x_n = b_n / u_{nn}$$
(18)

Following the same idea, the $(n-1)^{th}$ component of x (i.e. x_{n-1}) can be computed directly from the $(n-1)^{th}$ equation of (17), provided that we have computed x_n from equation (18) previously:

$$x_{n-1} = (b_{n-1} - u_{n-1n} x_n) / u_{n-1n-1}$$
(19)

The key here is that, since we know x_n from equation (18), we can <u>eliminate</u> it from the rest of the system; this is done simply by sending the terms involving x_n to the right hand side vector. The same is valid for x_{n-1} , after computing it via equation (19).

Following the same procedure, we can express any unknown, say x_k , in terms of the previously eliminated unknowns x_{k+1} , x_{k+2} , ... x_n , via the following relationship:

$$x_{k} = (b_{k} - \sum_{j=k+1}^{n} u_{kj} x_{j}) / u_{kk}$$
 (20)

Equation (20) defines an algorithm to eliminate successively all the variables, producing the values of the "n" unknowns in the linear system represented by equation (16). This process is known as a <u>back-substitution</u> process. Back-substitution arises from the fact that the matrix of coefficients in equation (16) is upper-triangular. In the case of T = L, a lower triangular matrix, equation (14) becomes:

$$L x = b \tag{21}$$

and, instead of back-substitution, <u>forward substitution</u> can be implemented following the same pattern as before, but in the opposite direction, i.e. the order of the elimination starts with x_1 , followed by x_2 , x_3 , etc. up to x_n . The equation which defines the forward substitution algorithm becomes:

$$x_{k} = (b_{k} - \sum_{j=1}^{k-1} l_{jk} x_{k}) / l_{kk}$$
 (22)

٦

where:

lij : is the coefficient corresponding to row i and column j
 of L.

It can be easily demonstrated that the amount of work demanded by either backward or forward substitution is the same, and is equal to:

Number of subtractions =
$$n(n-1)/2$$

Number of divisions = n
Number of multiplications = $n(n-1)/2$ $n(n+1)/2$ (23)

Since normally we count the amount of work in terms of the number of multiplications/divisions (dropping the additions and subtractions), equation (23) leads to a simpler result which says that the amount of work for the substitution process is of the order of $n^2/2$.

Also, it is easy to see that the solution of any triangular linear system <u>exists</u>, provided that the diagonal elements (either u_{kk} or l_{kk}) are all non-zero, and that, in this case, the solution is <u>unique</u>.

In summary, we have shown so far that, if we have a linear system of equations where the matrix of coefficients is triangular, then the solution for "x" is easily obtainable via a substitution process.

Hence, to solve a linear system where A is a general matrix, all we need to do is to <u>transform it into a triangular matrix</u>, since we already know that the rest is simply a substitution step.

Gaussian elimination provides an algorithm for transforming a general matrix into a triangular one. This algorithm is based on the following theorem from linear algebra [Conte and de Boor (1983)]:

On solving the linear system Ax = b we are allowed to perform any of the following elementary operations:

- Multiplication of any equation by an arbitrary scalar.
- Add one equation to the multiple of any other.
- Interchange any two equations.

In so doing, we shall obtain a new system, say $A^* x = b^*$, which is equivalent to the original one (Ax = b); i.e. the solution of both systems is exactly the same (x).

Hence, we can use a combination of these elementary operations to transform our original system Ax = b into another one $A^* x = b^*$, adding the additional condition that A^* , the transformed matrix, must be triangular (upper or lower).

From now on, and for simplicity, we shall assume that we want to obtain an upper-triangular transformed matrix. It is straightforward to see that the same reasoning is valid in the case of a lower-triangular matrix.

To produce the upper triangular matrix from the original matrix, all we need to do is to proceed column by column (though a row by row approach is also possible), eliminating (i.e. leaving a zero in) all the positions under the diagonal elements. To do so, we are allowed to perform any combination of the elementary operations mentioned before. This process finishes with an upper triangular matrix.

An example helps to illustrate how Gaussian elimination works. Let us suppose we want to solve the linear system:

$$2 x_1 + 3 x_2 - x_3 = 5
4 x_1 + 8 x_2 - 3 x_3 = 3
-2 x_1 + 3 x_2 - x_3 = 1$$
(24)

The first step is to eliminate all those terms under the diagonal corresponding to the first row (i.e. under 2 in the upper left hand corner of the matrix, which is the so called "pivot element" at this stage); using a combination of elementary operations we get:

The second step is to eliminate the 6 under the diagonal of the second row, which leads to:

And we have got an upper triangular linear system, which can be easily solved by back-substitution, leading to the solution:

$$x_3 = 27$$

 $x_2 = 10$
 $x_1 = 1$

This is the basic algorithm for Gaussian elimination, which can be easily implemented in a computer. The only probable cause of problems can be the eventual presence of zeros in the current pivot element (at any stage of the elimination), in that case the procedure breaks-down. Moreover, even if the pivot is not exactly zero, but very small, the process can produce unacceptable errors leading to a completely wrong solution. The remedy to this stability problem is to determine a different strategy for the selection of the pivot element.

ii) Pivoting strategies.

As far as the stability of Gaussian elimination is concerned, when A is a general matrix, there are two main strategies for the selection of the pivot, in order to avoid a zero (or nearly zero) in the current pivot and to minimize the effect of rounding errors: scaled partial pivoting and total pivoting.

a) <u>Scaled partial pivoting</u>: under this strategy we are allowed to decide which equation (or row) will be picked up as pivotal equation. Let us assume that we are at the k-th step, then we shall choose the j-th equation as the new pivotal equation if and only if the j-th element in column k is the largest one (in absolute value), i.e.:

 $|w_{kj}|^2 |w_{ki}| \quad \forall i=k, k+1, ..., n i \neq j$

Once the new pivotal equation has been found, we perform a permutation between rows "k" and "j" and carry on with the elimination of the elements under this new pivot. Notice that if the original matrix is symmetric, this property is lost on

following this kind of pivoting.

In the previous example, represented by equation (24), partial pivoting, as described before, is performed in the following sequence:

Original system:

$$2 x_1 + 3 x_2 - x_3 = 5
4 x_1 + 8 x_2 - 3 x_3 = 3
-2 x_1 + 3 x_2 - x_3 = 1$$
(24)

The first pivot is found by looking at the first column and finding out that the second row has the biggest (absolute value) element. Then, the first and second rows are permuted, leading to:

The elimination of the terms under the pivot produces:

$$4 x_{1} + 8 x_{2} - 3 x_{3} = 3 - x_{2} + \frac{1}{2x_{3}} = \frac{7}{2} 7 x_{2} - \frac{5}{2x_{3}} = \frac{5}{2}$$
(28)

Now, the second pivot is chosen looking at the second column of the system (28); we have to choose between the second and third rows, and we find that the third row is the new pivotal equation. On permuting the second and third rows of (28), we get:

$$4 x_1 + 8 x_2 - 3 x_3 = 37 x_2 -5/2x_3 = 5/2- x_2 +1/2x_3 = 7/2$$
(29)

Eliminating the term under the pivot (7), we get:

$$4 x_1 + 8 x_2 - 3 x_3 = 37 x_2 -5/2x_3 = 5/22/14x_3 = 54/14$$
(30)

Hence, the solution obtained via back-substitution becomes:

$$x_3 = 27$$

 $x_2 = 10$
 $x_1 = 1$

as before.

Scaled partial pivoting is a powerful complement of the standard Gaussian elimination, which avoids most of the sources of trouble of the original method, at a very low expense.

Although the example does not show the full advantages of scaled partial pivoting in some "pathological" cases, such examples are presented in the literature (e.g. Conte and de Boor, (1983) chapter 4.3.)

b) <u>Total pivoting</u> (complete or full pivoting): in this case we are allowed to decide which equation (row) and which variable (column) will be picked-up as pivotal equation and pivot element. Thus, we have now two degrees of freedom for choosing the pivot, and the search of the pivot will be carried out looking both at the columns and rows to the left and under the last pivot. Obviously we shall choose that corresponding to the element with the maximum absolute value.

In our previous example, represented by equation (24):

$$2 x_1 + 3 x_2 - x_3 = 5
4 x_1 + 8 x_2 - 3 x_3 = 3
-2 x_1 + 3 x_2 - x_3 = 1$$
(24)

When choosing the first pivot, we find out that the maximum absolute value in the whole matrix is in the second row and second column (8), then, we permute both first and second rows and columns, obtaining:

$$8 x_{2} + 4 x_{1} - 3 x_{3} = 3$$

$$3 x_{2} + 2 x_{1} - x_{3} = 5$$

$$3 x_{2} - 2 x_{1} - x_{3} = 1$$
(31)

The elimination of the elements under the pivot leads to:

$$8 x_{2} + 4 x_{1} - 3 x_{3} = 3$$

$$1/2 x_{1} + 1/8 x_{3} = 31/8$$

$$-28 x_{1} + x_{3} = -1$$
(32)

We choose the second pivot looking at all the elements in the second and third rows and we find that the new pivot is -28 (third row, variable x_1), then we have to permute the second with the third rows:

$$8 x_{2} + 4 x_{1} - 3 x_{3} = 3$$

-28 x₁ + x₃ =-1
1/2 x₁ +1/8x₃ = 31/8 (33)

Hence, eliminating the element under the second pivot:

$$8 x_{2} + 4 x_{1} - 3 x_{3} = 3$$

-28 x₁ + x₃ =-1
8/3x₃ = 216/7 (34)

which produces by back-substitution:

$$x_3 = 27$$

 $x_2 = 10$
 $x_1 = 1$

as before.

Although total pivoting is recognised as more powerful than scaled partial pivoting, it is not frequently used, mainly because of the additional amount of work involved, which is relevant in large systems, since this involves a search through the whole matrix .

iii) Amount of work and storage required.

The amount of work involved by Gaussian elimination is:

-	Number	of	additions	:	(n ³ -n)/3		
-	Number	of	multiplications	:	(n ³ -n)/3	}	(35)
_	Number	of	divisions	:	n(n-1)/2		

Hence, we have the following number of multiplications/divisions:

$$n^{3}$$
 n^{2} $5n$
--- + --- (36)
3 2 6

٦

and, as far as the amount of work is concerned, the Gaussian elimination process is said to be of the order of $n^3/3$.

Thus, the complete solution of a linear system via Gaussian elimination plus back-substitution, requires the following amount of work:

Stage	Additions	Multiplications/divisions		
Factorization Substitution	$\frac{(n^3-n)/3}{n^2/2} - n/2$	$n^{3}/3 + n^{2}/2 - 5/6n$ $n^{2}/2 + n/2$		
Total	$n^3/3 + n^2/2 + 5/6n$	$n^3/3 + n^2 - n/3$		

As a result, we note that the solution of linear systems via a sequence of Gaussian elimination and back-substitution requires an amount of work which is, roughly speaking, of the order of $n^3/3$; the factorization process, i.e. the transformation of the original matrix into a triangular one is the most expensive item in the total cost.

As far as the cost of pivoting is concerned, scaled partial pivoting requires of the order of n^2 comparisons, while total pivoting needs of the order of $n^3/3$ comparisons [Golub and Van Loan (1983)].

Since the triangular matrix produced by Gaussian elimination can be completely filled with non-zeros, the amount of storage required is:

$$n(n+1)/2 = n^2/2 + n/2$$
 (37)

that is to say, from the storage point of view, Gaussian elimination, for general dense matrices, requires of the order of $n^2/2$ locations.

Because the explicit inversion of a matrix requires of the order of n^3 multiplications/divisions, it is now clear why Gaussian elimination with substitution is a more efficient method for solving linear systems. Indeed, we can produce the same results with just a third of the work. From the storage viewpoint, the solution of linear systems via Gaussian elimination requires only a half of the storage needed for an explicit inversion of the matrix. Unless it is specially justified, matrix inversion should be avoided.

We have to emphasise that, in the evaluation of the work and storage required presented so far, we have not taken into account that, in some cases, the special structure of the matrix (symmetry, sparseness) can be exploited, in order to reduce both the amount of work and the storage needed. Since most of the matrices produced in engineering applications have some special feature, and since in particular, the matrices generated in the context of water distribution analysis via the gradient method are always symmetric and positive definite, then we do have to take these properties into consideration, since we can get important savings in the amount of computation and storage

A-23

required.

iv) The product form of Gaussian elimination.

The whole process of Gaussian elimination can be understood 85 a transformation process, from a general matrix A into a triangular matrix. In fact, each step of Gaussian elimination can be formally expressed as a transformation, which in this case the form of the product of a matrix by another; takes in so doing, the whole elimination process becomes just a sequence of matrix operations. This concept is relevant from the practical point of view, since it can lead to important savings in theamount of work to be done, as will be shown in subsequent sections.

Gaussian elimination is equivalent to premultiplying the original system by a sequence of transformation matrices (T_i) . The structure of these matrices is such that, if we are eliminating the terms under the diagonal of the j-th column, the corresponding transformation matrix is of the following form:



where the j-th column has indeed the following form:

A-24

$$\begin{bmatrix} 0\\ \vdots\\ 0\\ 1\\ -a_{j+1,j}\\ \hline a_{j,j}\\ \hline -a_{j+2,j}\\ \hline a_{j,j}\\ \hline -a_{n,j}\\ \hline a_{j,j} \end{bmatrix} \int j-1 \text{ zeros}$$
(39)

Note that the values of the coefficients a_{ij} are those of the <u>current</u> matrix in the elimination process, rather than those of the original matrix .

Tł	nus, the ori	gin	al system is successively premultiplied	by	the
matr	cices T _i , i	=1,	2,,n, such that:		
Ori∉	inal system	1:	$A \mathbf{x} = \mathbf{b}$		
1^{st}	step	:	$T_1 A x = T_1 b$		
	or		$A_1 x = T_1 b$		
	with		$A_1 = T_1 A$		
2nd	step	:	$T_2 A_1 x = T_2 T_1 b$		
	or		$A_2 x = T_2 T_1 b$		
	with		$A_2 = T_2 A_1 = T_2 T_1 A$		
n-tł	¹ step	:	$T_n A_{n-1} x = T_n T_{n-1} \ldots T_2 T_1 b$	(40)
	or		$A_n \mathbf{x} = T_n T_{n-1} \ldots T_2 T_1 \mathbf{b}$	(41)
	with		$A_n = T_n T_{n-1} \dots T_2 T_1 A$	(42)

After the n-th step, A_n must be an upper triangular matrix, and the system (41) can be solved by back-substitution.

In the context of our previous small example [equation (24)]:
$$2 x_1 + 3 x_2 - x_3 = 5
4 x_1 + 8 x_2 - 3 x_3 = 3
-2 x_1 + 3 x_2 - x_3 = 1$$
(24)

The sequence of transformation matrices is:

For the first step:

$$T_{1} = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \longrightarrow A_{1} = T_{1} A = \begin{bmatrix} 2 & 3 & -1 \\ 0 & 2 & -1 \\ 0 & 6 & -2 \end{bmatrix}$$

For the second step:

$$T_{2} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -3 & 1 \end{bmatrix} \longrightarrow A_{2} = T_{2} T_{1} A = \begin{bmatrix} 2 & 3 & -1 \\ 0 & 2 & -1 \\ 0 & 0 & 1 \end{bmatrix}$$

Note that the matrices T_j are highly sparse and that they can be handled implicitly with a vector of length "n".

The practical importance of the product form of Gaussian elimination lies in the fact that, if we need to solve different linear systems with the same matrix of coefficients, for example:

$$A x = b$$
$$A y = c$$
$$A z = d$$

then, in this case, all we need to do is to transform the matrix A into:

 $A_n = T_n T_{n-1} \ldots T_2 T_1 A$

Ъ

С

d

and apply back-substitution to the transformed systems:

A-26

Hence, in this case, we can compute and store a transformation matrix:

$$T = T_n T_{n-1} \dots T_2 T_1$$
 (43)

with T_j as specified in (38) and (39), which can be used as many times as different right hand side vectors are given to us. As can be seen in our small example, the matrix T takes the form:

$$T = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 7 & -3 & 1 \end{bmatrix} = T_2 T_1 \text{ (note that } T_3 = I\text{)}$$

which is a lower triangular matrix with 1's in the diagonal.

Another immediate application of the product form of Gaussian elimination is matrix inversion, since this problem can be formulated as the solution of "n" linear systems of the form:

$$A x^{(i)} = e^{(i)}$$
 (44)
where:

$$x(i)$$
: are the columns vectors of the matrix A^{-1} .

 $e^{(i)}$: are vectors of the form $(0,0,\ldots,1,\ldots,0)^{T}$, i.e. null vectors with a one added in the i-th position.

Hence, we need to solve "n" times the same linear system, with "n" different right hand side vectors; the corresponding "n" solution vectors are assembled by columns to produce the inverse matrix:

$$A^{-1} = [x^{(1)} + x^{(2)} + \dots + x^{(n)}]$$
(45)

v) Gaussian elimination as a LU decomposition.

Gaussian elimination can also be interpreted as an equivalent triangular decomposition of the matrix A, since when multiplying T A we are actually getting an upper triangular matrix, say U :

$$T A = U$$
(46)

$$A = T^{-1} U \tag{47}$$

where T^{-1} is a lower triangular matrix: L , such that:

$$L = T^{-1}$$
(48)

and then

then,

$$A = L U$$
(49)

in terms of our small example:

$$\mathbf{L} = \mathbf{T}^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & 3 & 1 \end{bmatrix}$$

and

$$U = \begin{bmatrix} 2 & 3 & -1 \\ 0 & 2 & -1 \\ 0 & 0 & 1 \end{bmatrix}$$

Hence, Gaussian elimination is closely related to a family of factorization methods, which shall be reviewed in subsequent sections.

A.2.2. Gauss-Jordan elimination.

i) Introduction.

Gauss-Jordan elimination is similar to Gaussian elimination, the difference being that instead of producing a triangular matrix, eliminating only the non-zeros below the diagonal, we can carry on the elimination both below and above the diagonal, performing as many elementary operations as needed, until we get a transformed matrix which is an identity matrix . ii) The product form of the inverse.

In a similar manner to Gaussian elimination, Gauss-Jordan elimination can be explained in terms of the product of a succession of transformation matrices T_j^* .

The structure of these transformation matrices is now the following:



the j-th column has the following form:

$$\begin{bmatrix} -\frac{a_{1,j}}{a_{j,j}} \\ -\frac{a_{2,j}}{a_{j,j}} \\ \frac{-a_{2,j}}{a_{j,j}} \\ \vdots \\ \vdots \\ \frac{1}{a_{j,j}} \\ \vdots \\ \vdots \\ \frac{-a_{n-1,j}}{a_{j,j}} \\ \frac{-a_{n,j}}{a_{j,j}} \end{bmatrix}$$
 (51)



Equations (50) and (51) are for Gauss-Jordan, as equations (38) and (39) are for Gaussian elimination.

Note that the values of the coefficients a_{ij} are those of the <u>current</u> matrix in the elimination process, rather than those of the original matrix.

Thus, the successive premultiplications are carried out in the following sequence:

Ori∉	ginal system	1:	$A \mathbf{x} = \mathbf{b}$	
1^{st}	step	:	$T_1^* A x = T_1^* b$	
	or		$A_1^* x = T_1^* b$	
	with		$A_1^* = T_1^* A$	
2nd	step	:	$T_2^* A_1^* x = T_2^* T_1^* b$	
	or		$A_2^* x = T_2^* T_1^* b$	
	with		$A_2^* = T_2^* A_1^* = T_2^* T_1^* A$	
n-tł	¹ step	:	$T_n^* A_{n-1}^* x = T_n^* T_{n-1}^* \dots T_2^* T_1^* b$	(52)
	or		$A_n^* x = T_n^* T_{n-1}^* \dots T_2^* T_1^* b$	(53)
	with		$A_n^* = T_n^* T_{n-1}^* \dots T_2^* T_1^* A$	(54)

after the n-th step, A_n^* must be an identity matrix and the solution is exactly the current right hand side vector, hence no back-substitution is needed, i.e.:

$$A_n^* = T_n^* T_{n-1}^* \dots T_2^* T_1^* A = I$$
 (55)

then:

$$A^{-1} = T_n^* T_{n-1}^* \dots T_2^* T_1^*$$
(56)

this is the so-called <u>product form of the inverse</u> of the matrix A.

Gauss-Jordan elimination provides us with another way to buildup the inverse of a matrix by simply multiplying "n"

transformation matrices of the form of T_j^* , which are very sparse.

The amount of work needed for Gauss-Jordan elimination is:

- Number of additions : $n^3/2 n/2$
- Number of multiplications : $n^3/2 + n^2 n/2$
- Number of divisions : n

Since the inverse of a matrix can be a completely filled matrix, the storage requirements are n² positions.

As before, in the case of Gaussian elimination, the usefulness of Gauss-Jordan elimination for the repeated solution of linear systems which differ only in the right hand side vector is quite clear, with no need for further explanation.

A.2.3. Matrix factorization methods.

i) Introduction.

We have already seen that Gaussian elimination is, in essence, equivalent to factorizing the matrix of coefficients into the product of a lower and an upper triangular matrix, although in the standard Gaussian algorithm both matrices are not explicitly determined.

In fact we can follow different strategies for obtaining this factorization, each one leading to a different direct method for the solution of linear systems of equations. We shall review some of these factorization (or decomposition) strategies in the following sections.

The common idea behind any factorization strategy is that instead of having to solve the original system:

$$A x = b \tag{57}$$

if we factorize A as:

$$A = L U$$
(58)

where:

L is a lower triangular matrix.

U is an upper triangular matrix.

then, because the original problem is now :

$$L U x = b$$
 (59)

if we define:

$$\mathbf{U} \mathbf{x} = \mathbf{y} \tag{60}$$

where y is an auxiliary vector, the solution of the original system for x, can be split up into two substitution stages:

a) Forward substitution:

solve L y = b (61)

determining "y", and

b) Backward substitution: with "y" already known
solve equation (60) U x = y

determining "x", the original unknown.

Of course, this two stage process can be effectively used for solving multiple right hand side problems, where we only need to factorize the original matrix at the beginning, the rest being just a succession of forward and backward substitutions.

ii) LU triangular factorization (Crout method).

In this case the lower-upper triangular factorization takes the following form:

$$A = L U$$
(62)

where:

A : is a general unsymmetric matrix of coefficients.

- L : is a lower triangular matrix whose diagonal elements are all 1.
- U : is an upper triangular matrix.

Under these conditions, L and U are <u>unique</u> factors.

From equation (63) it can be seen that a recurrent scheme can be devised to compute the elements of the factors L and U, from the original elements of A.

On multiplying the first row of L by the columns of U, and identifying it with the first row of A, we find that the first row of U is simply the same first row of A and no computation is needed, i.e.:

$$u_{1i} = a_{1i} \quad \forall i=1,2,...,n$$
 (64)

The following step is to determine the elements of the first column of L; in this case, multiplying the first column of L by the first row of U and identifying the result with the first column of A, we get:

 $l_{i1} u_{11} = a_{i1} \quad \forall i=2,3,\ldots,n$ (65) then,

$$l_{i1} = a_{i1}/u_{11} \quad \forall i=2,3,...,n$$
 (66)

At this stage we are able to compute the second row of U, by

multiplying the second row of L by the columns of U and identifying the result with the second row of A, to obtain:

 $l_{21} u_{1i} + u_{2i} = a_{2i} \quad \forall i=2,3,...,n$ (67) then,

 $u_{2i} = a_{2i} - l_{21} u_{1i} \quad \forall i=2,3,\ldots,n$ (68)

Following the same pattern, the whole set of elements of L and U can be computed. Since we need to compute the terms $1/u_{ii}$, the algorithm breaks down when $u_{ii} = 0$. Also, the stability of the method can be affected when a very small u_{ii} is found. As in Gaussian elimination, partial or complete pivoting can help to stabilise the algorithm.

The amount of work demanded by this L U factorization is: - Number of multiplications : $n^3/3 - n/3$ - Number of divisions : n - Number of additions : $n^3/3 - n^2/2 + n/6$

Having determined the factors L and U, the solution of the linear system requires a forward and a backward substitution which, as we already know each substitution requires n(n+1)/2 multiplications/divisions, giving a total of $n^3/3 + n^2 + 2/3n$ multiplications/divisions for the solution of the linear system (factorization plus substitutions), which is basically the same amount of work needed for the Gaussian elimination process. The storage required is n^2 positions.

iii) L L^T factorization.

In the particular case when A is <u>symmetric and positive</u> <u>definite</u>, it can be seen that Crout's method leads to a symmetric

decomposition:

$$\mathbf{A} = \mathbf{L} \, \mathbf{U} = \mathbf{L} \, \mathbf{L}^{\mathrm{T}} \tag{69}$$

where $U = L^T$ is the transpose of L.

The amount of work is reduced to almost one half ($n^3/6 + n^2/2$ multiplications/divisions in the factorization stage) and the storage requirements are also halved.

Of course, the solution of a linear system using this symmetric decomposition demands a forward and a backward substitution, as before:

- Solve L y = b for y
- Solve $L^T x = y$ for x

We shall return to this kind of factorization later on, when dealing with Cholesky factorizations.

iv) L D U factorization.

An extension of the LU factorization consists of decomposing ever further the upper triangular matrix defined in equation (63), into the product of a diagonal matrix (D) and an upper (unit diagonal) triangular matrix:

or:

$$A = L D U$$

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

$$= \begin{bmatrix} 1 \\ 1_{21} & 1 \\ \vdots & \vdots \\ 1_{n1} & 1_{n2} & \dots & 1 \end{bmatrix} \cdot \begin{bmatrix} d_{11} \\ d_{22} \\ \vdots \\ d_{nn} \end{bmatrix} \cdot \begin{bmatrix} 1 & u_{12} & \dots & u_{1n} \\ 1 & \dots & u_{2n} \\ \vdots \\ 1 & \vdots \\ 1 & 1 \end{bmatrix}$$
(71)

where:

- A : is a general non-symmetric matrix of coefficients.
- L : is a lower triangular matrix whose diagonal elements are all 1.
- D : diagonal matrix
- U : is an upper triangular matrix whose diagonal elements are all 1.

Thus, instead of solving the original system:

$$A x = b \tag{72}$$

we are solving:

$$L D U x = b$$
(73)

and in this case the solution involves a three stage process:

a) Forward substitution:

solve L z = b (74)

for z, an auxiliary vector

- b) Diagonal inversion:
 - solve D y = z (75)

for y, by computing directly :

$$\mathbf{y} = \mathbf{D}^{-1} \mathbf{z} \tag{76}$$

this stage demands only n divisions, due to the fact that D is diagonal.

c) Back-substitution:

solve U x = y (77)

for x, the original unknown vector.

For the factorization, the algorithm requires $n^3/3 + n^2/2 + n/6$ multiplications/divisions. The advantage of this kind of factorization over LU factorization, becomes apparent in the handling of storage and in the implementation of the algorithms.

More details can be found for example in Golub and Van Loan (1983), algorithm 5.1.1.

v) L D L^{T} factorization.

If the matrix of coefficients of the linear system is a nonsingular symmetric matrix, it can be proved [Golub and Van Loan (1983), theorem 5.1.2.] that in the L D U factorization the upper triangular matrix U is equal to the transpose of the lower triangular matrix, i.e.:

$$\mathbf{U} = \mathbf{L}^{\mathrm{T}} \tag{78}$$

and

$$\mathbf{A} = \mathbf{L} \, \mathbf{D} \, \mathbf{L}^{\mathrm{T}} \tag{79}$$

It is easy to see that the amount of work for the decomposition stage now reduces to almost half of that of the L D U decomposition ($n^3/6$ multiplications/divisions, approximately), i.e. it is almost the same amount of work needed for the L L^T factorization.

Of course, the solution of the linear system demands successive forward substitution, diagonal inversion and backward substitution, in a similar way as in the case of the previous L D U decomposition.

From the stability point of view the decomposition is not numerically stable and the condition of non singularity of A does not seem to be sufficient, a more restrictive condition is needed. An example provided by Gill, Murray and Wright (1981) is useful to illustrate this problem:

the matrix $A = \begin{bmatrix} \epsilon & 1 \\ 1 & \epsilon \end{bmatrix}$

has a L D L^{T} factorization given by:

$$L = \begin{bmatrix} 1 & 0 \\ 1/\epsilon & 1 \end{bmatrix}, D = \begin{bmatrix} \epsilon & 0 \\ 0 & \epsilon - 1/\epsilon \end{bmatrix} \text{ and } L^{T} = \begin{bmatrix} 1 & 1/\epsilon \\ 0 & 1 \end{bmatrix}$$

hence, for $\varepsilon << 1$ the elements of the factors can become very large.

vi) Complete Cholesky factorization.

If we have, as usually occurs in engineering problems, a symmetric positive definite matrix A, the L D L^{T} factorization previously seen has the property that all diagonal elements of D are strictly positive [see theorem 5.2.1, Golub and Van Loan (1983) for a proof] and then, two alternative decompositions can be developed:

a) L D L^T Cholesky complete decomposition.

The matrix of coefficients of the linear system of equations can be decomposed as:

$$A = L D L^{T}$$
(80)

that is to say:

$$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} = \\ = \begin{bmatrix} 1_{11} & & & \\ 1_{21} & 1_{22} & & \\ \vdots & \vdots & & \\ 1_{n1} & 1_{n2} & \dots & 1_{nn} \end{bmatrix} \cdot \begin{bmatrix} d_{11} & & & \\ d_{22} & & & \\ & & & d_{nn} \end{bmatrix} \cdot \begin{bmatrix} 1_{11} & 1_{12} & \dots & 1_{1n} \\ 1_{22} & \dots & 1_{2n} \\ & & & & 1_{nn} \end{bmatrix}$$
(81)

where:

 $a_{ij} = a_{ji}$, since the matrix A is symmetric, and $d_{ii} > 0 \forall i=1,2,...,n$

An algorithm for the computation of the elements of L and D can be developed in a constructive fashion, considering first the product of L D as:

$$L^* = L D \tag{82}$$

then,

$$L^{*} = \begin{bmatrix} 1_{11} & d_{11} & 1_{22} & d_{22} \\ 1_{21} & d_{11} & 1_{32} & d_{22} & 1_{33} & d_{33} \\ \vdots & \vdots & \vdots & \vdots \\ 1_{n1} & d_{11} & 1_{n2} & d_{22} & 1_{n3} & d_{33} & \dots & 1_{nn} & d_{nn} \end{bmatrix}$$
(83)

Thus, A = L D $L^T = L^* L^T$, and it takes the form:

To determine the elements l_{ij} and d_{ij} we have to multiply the elements of L* by those of L^T and identify the results with the corresponding elements of A. In so doing, we find that we have one degree of freedom at our disposal, for example, we can impose the condition that all the diagonal elements of L* are 1:

$$l_{ii} d_{ii} = 1 \quad \forall i=1,2,..,n$$
 (85)

and then the algorithm becomes: Step 1 : Compute for j=i

$$l_{ji} = a_{ij} - \sum_{k=1}^{i-1} l_{jk} l_{jk} d_{kk} \quad j=i,(i+1),...,n \quad (86)$$

(at the beginning j=i=1 and $l_{ji}=a_{ij}$, then $l_{11}=a_{11}$) Step 2 : Compute

> $d_{ii} = (l_{ii})^{-1}$ (87) [this is due to equation (85)] the algorithm is stable at this stage, since if A is positive definite then $l_{ii} \neq 0 \quad \forall i$

Step 3 : Compute

 l_{ji} for j > i, using equation (86).

The amount of work required for the decomposition stage is of the order of $n^3/6 + n^2$ multiplications/divisions.

A couple of examples will illustrate the main features of the algorithm. First of all, an example with a symmetric matrix which is not positive definite will be considered, in this case we know in advance that the algorithm is not applicable:

Solve
$$\begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 0 \\ 3 & 0 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 6 \\ 3 \\ 6 \end{bmatrix}$$

determinant of A =-18, then the matrix is negative. the factorization algorithm gives: $l_{11} = 1$ $d_{11} = 1$ $l_{21} = 2$ $l_{31} = 3$ $l_{22} = 1 - l_{11}^2 d_{11} = 1 - 1 = 0$

An example where the factorization is possible is:

 $d_{22} = 1/0$ not defined !!, so the factorization is not possible.

Solve
$$\begin{bmatrix} 2 & 2 & 1 \\ 2 & 5 & 2 \\ 1 & 2 & 1 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = b$$

determinant of A = +1, then the matrix is positive definite the factorization algorithm gives:

$$\begin{bmatrix} 2 & 2 & 1 \\ 2 & 5 & 2 \\ 1 & 2 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 \\ 2 & 3 & 0 \\ 1 & 1 & 1/6 \end{bmatrix} \cdot \begin{bmatrix} 1/2 & 0 & 0 \\ 0 & 1/3 & 0 \\ 0 & 0 & 6 \end{bmatrix} \cdot \begin{bmatrix} 2 & 2 & 1 \\ 0 & 3 & 1 \\ 0 & 0 & 1/6 \end{bmatrix}$$

$$A = L \qquad D \qquad L^{T}$$
Note that d_{ii} > 0 \ Y i ; also l_{ij} > 0 \ Y i, j.

b) R R^T Cholesky complete factorization (square root factorization).

Using the fact that if A is symmetric and positive definite then D is also positive, then the factorization:

 $A = L D L^{T}$

can be expressed as:

$$A = L D^{1/2} D^{1/2} L^{T}$$
 (88)

where:

$$D^{1/2} = \begin{bmatrix} d_{11}^{1/2} \\ d_{22}^{1/2} \\ \vdots \\ d_{nn}^{1/2} \end{bmatrix}$$
(89)

The matrix $D^{1/2}$ exists since D is positive, the square roots in equation (89) being the reason for the restriction of positiveness in D.

Then, A can be expressed as:

 $R^{T} = D^{1/2} L^{T}$

$$A = R R^{T}$$
(90)

where:

 $R = L D^{1/2}$

and

which is exactly the same L L^T factorization we obtained before.

The algorithm can be derived from the fact that equation (90) is of the form:

$$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} = \begin{bmatrix} r_{11} & & & \\ r_{21} & r_{22} & & & \\ \vdots & \vdots & & \vdots \\ r_{n1} & r_{n2} & \dots & r_{nn} \end{bmatrix} \cdot \begin{bmatrix} r_{11} & r_{12} \dots & r_{1n} \\ r_{22} \dots & r_{2n} \\ & & \vdots \\ & & & r_{nn} \end{bmatrix}$$
(91)

Then, multiplying R R^T and identifying with A term by term, we get (note that $r_{ij} = r_{ji}$, because A is symmetric):

 $a_{11} = r_{11}^{2} \qquad -----> \qquad r_{11} = (a_{11})^{1/2}$ $a_{12} = r_{11} r_{12} \qquad -----> \qquad r_{12} = a_{12} / r_{11}$ $a_{1n} = r_{11} r_{1n} \qquad -----> \qquad r_{1n} = a_{1n} / r_{11}$ $a_{22} = r_{12}^{2} + r_{22}^{2} ----> \qquad r_{22} = (a_{22} - r_{12}^{2})^{1/2}$ etc.

This is a row-oriented version of the Cholesky factorization and, of course, a column-oriented algorithm can be derived following the same pattern; see for example Golub and van Loan (1983), algorithm 5.2.1., which is efficient since the elements a_{ij} are overwritten by r_{ij} (i 2 j).

The algorithm requires $n^3/6$ multiplications/divisions and additions, and n square roots. If A is not positive definite, the terms under the square root become negative and the algorithm breaks down. Our previous non-positive small example helps to illustrate this point:

Solve :

A-42

$$\begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 0 \\ 3 & 0 & 3 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 6 \\ 3 \\ 6 \end{bmatrix}$$

the determinant of A =-18, then the matrix is negative. the factorization algorithm gives:

$$r_{11} = (1)^{1/2} ----> r_{11} = 1$$

$$r_{11} r_{12} = 2 ----> r_{12} = 2$$

$$r_{12}^2 + r_{22}^2 = 1 ----> r_{22} = (1 - 2^2)^{1/2} \text{ not defined }!!$$

The main advantage of this version of Cholesky factorization is the fact that no permutations are needed to ensure stability, like in Gaussian elimination with pivoting. This is due to the step in the algorithm, where the diagonal elements are computed:

which implies that the terms r_{ik} cannot grow indefinitely (like in the example of the L D L^T factorization) and their maximum limit is indeed $(a_{kk})^{\frac{1}{2}}$. This is <u>always true</u>, even if a very small <u>pivot is found during the factorization</u>. The implications of this will become clearer when dealing with direct sparse methods, because it will allow us to concentrate on reducing the creation of fill-in, without having to worry about stability.

vii) Orthogonalization methods.

A completely different alternative to Gaussian elimination methods is the use of orthogonal factorization. A number of orthogonal methods have been proposed, all of them exploiting the fact that, in an orthogonal matrix, say Q, the product $Q^{T}Q$ is the

identity matrix; i.e. the inverse of an orthogonal matrix is its transpose.

As a result, if we want to solve our model problem A = b and if we are able to factorize A such that:

$$A = Q U \tag{92}$$

where:

Q is orthogonal

U is an upper triangular matrix

then, the linear system becomes:

$$Q U x = b \tag{93}$$

and premultiplying by Q^{T} , we get:

$$\mathbf{U} \mathbf{x} = \mathbf{Q}^{\mathrm{T}} \mathbf{b} \tag{94}$$

where a back substitution determines the unknown "x".

There are thus as many methods as factorizations that can be imagined.

The main advantage of all these methods is that they are especially well suited for ill-conditioned problems, where the matrix A becomes singular or nearly singular.

A method of this kind is the "singular value decomposition", briefly described by Press et al. (1986), where a FORTRAN implementation can be found. A more detailed description is provided by Golub and Van Loan (1983).

Unfortunately, these methods are much more costly than Gaussian elimination [Reid et al. (1986) reported at least twice the arithmetic cost of Gaussian elimination]. On top of that, the methods are not efficient in handling sparse matrices, since the orthogonal factors get a great deal of fill-in. Because of these reasons, orthogonal methods are not widely used, their use being restricted to solve ill-conditioned problems.

A.2.4. Comparison of the different algorithms for the direct solution of general dense linear systems.

We have briefly reviewed most of methods for solving general dense linear systems of equations, the main algorithm being Gaussian elimination, while all the rest can be understood as variations of that method.

Because the main objective in this section has been to introduce the basic ideas behind the solution of linear systems of equations, for the particular case of general dense matrices, we have avoided a number of details and we shall maintain the same policy in this comparison. This is due to the fact that most of the methods we are interested in are described in the following sections, since sparsity is a feature that must be explicitly considered when solving efficiently linear systems arising from water distribution network analysis.

Because all the methods described so far can be programmed to overwrite the original matrix of coefficients, the storage ceases to be relevant in comparative terms. Hence, the amount of work and stability remain as the only relevant factors of the comparison.

Table A.1 summarises the amount of work and storage requirements of the different algorithms reviewed. From a

Algorithm	Matrix	Anount of wo Multiplications/divisions	≯rk required Additions/subtractions	Storage required Positions (#)	Notes
Substitution: back or forward	Triangular	$\frac{n^2 + n}{2}$	$\frac{n^2 - n}{2 - 2}$		
Gaussian elimination	Unsymmetric	$\frac{n^3 + n^2 - 5}{3 + 6}$ n	$\frac{n^3}{3} - \frac{n}{3}$	<u>n</u> ² 2	Pivoting and permutations needed for stability .
Gauss elimin. + substitution	Unsymmetric	$\frac{n^3}{3} + n^2 - \frac{1}{3}n$	$\frac{n^3}{3} + \frac{n^2}{2} - \frac{5}{5}n$	<u>n</u> ² 2	Idea
Gauss-Jordan elimination	Unsymmetric	$\frac{n^3 + n^2 - 1}{2}$ n	$\frac{n^3}{2} - \frac{n}{2}$	n ²	Idem
L U factorization	Unsymmetric	$\frac{\underline{n}^3 - \underline{i}}{3} n$	$\frac{\underline{n}^3 - \underline{n}^2 + \underline{n}}{3 2 6}$	n ²	
L U factoriz.+ back+forw.subst.	Unsymmetric	$\frac{n^3 + n^2 + 2}{3}$ n	$\frac{n^3 + n^2 - 5}{3}$ n 3 2 6	n ²	
L L ^T factorization	Symmetric Pos.Defin.	$\frac{1}{6}\frac{n^3+n^2}{2}$	$\frac{2}{6}\frac{n^3}{4} + \frac{n^2}{4}$	<u>n</u> ² 2	Good from the stability point of view.
LL ^T factoriz.+ back+forw.subst.	Symmetric Pos.Defin.	$\frac{1}{2} \frac{n^3}{6} + 3n^2 + n$	<u>* n³ +5n² - n</u> 6 4	<u>n</u> ² 2	*********
L D U factorization	Unsymmetric	$\frac{n^3 + n^2 + 1}{3 + 2}$ n	$\frac{n^3}{3} - \frac{n^2}{2} + \frac{1}{6}n$	n ²	
LDU factoriz.+ back.+drag.inv+ forw.substitut.	Unsymmetric	$\frac{n^3}{3} + 3n^2 + 13 n$ 3 2 6	$\frac{n^3}{3} + \frac{n^2}{2} - \frac{5}{5}n$	n ²	
L D L ^T factorization	Symmetric Pos.Defin.	$\frac{1}{6} \frac{n^3 + n^2 + n}{2}$	$\frac{1}{2} \frac{n^3}{6} + 3n^2 - \frac{1}{2}n$ 6 4 2	$\frac{1}{2}\frac{n^2}{2}$	
LDL ^T factoriz.+ back.+diag.inv+ forw.substitut.	Symmetric Pos.Defin.	<u>* n</u> ³ + 2n ² + <u>5</u> n 6 2	$\frac{1}{6}\frac{n^3}{4} + 7n^2 - \frac{3}{2}n$	$\frac{1}{2}\frac{n^2}{2}$	
Explicit inversion	Unsymmetric	n ³ -1	$n^3 - 2n^2 + n$	n ²	
Notes: (#) It does not c	consider the	fact that most of the algor	rithms can overwrite the o	riginal matrix.	

Table A.1. Comparison of the algorithms used in the direct solution of dense linear systems of equations.

comparative point of view we can say that:

a) The explicit inversion of the matrix should be avoided when solving a linear system, unless it is explicitly needed for other purposes; in the latter case the product form of the inverse should be used since only requires of the order of $n^3/2$ multiplications and divisions.

b) When solving non-symmetric linear systems, Gaussian elimination either with partial or full pivoting can be the best choice, with Gauss-Jordan elimination being more expensive than Gaussian elimination. L D L^{T} factorization presents some stability problems and should be used with due care.

c) When solving multiple right hand side linear systems, a factorization of the matrix A is imperative, since this allows the solution for the different right hand sides with additional forward and backward substitutions.

d) For symmetric positive definite systems, the complete Cholesky factorization (square root method) seems to be the best choice, both from the stability and amount of work viewpoints. Symmetry must be exploited, since this reduces both storage and amount of work by 50 %.

e) For ill-conditioned problems try a method based on orthogonalization.

i) Introduction.

Although we have defined "direct methods" as those where an <u>exact</u> solution is reached in a <u>finite</u> number of elementary mathematical operations, in practice and due to the fact that all computers have a fixed word length, which determines its maximum accuracy, exact solutions are never reached. Then, <u>there are no</u> <u>exact solutions</u> just approximate solutions.

Two main sources of error can be identified when solving a linear system of equations like:

$$A \mathbf{x} = \mathbf{b} \tag{95}$$

a) Both the coefficients of the matrix A and the right had side vector b are computed based on some physical properties of the problem under study. As a result, both A and b are subject to errors in their estimation. This means that in practice we are dealing with a system like:

$$(A + \delta A) x = (b + \delta b)$$
(96)

where:

- SA : error in the estimation of the matrix of coefficients, which is another matrix of the same order of A.
- 6b : error in the estimation of b, which is another vector of the same order of b.

b) The second source of error appears even if we have exact estimates for A and b, and is due to rounding errors in the process of determining the vector x (Gaussian elimination or whatever method we are using).

Depending on the sensitivity of the solution for the unknowns x, either with respect to small changes in the coefficients of A and b or with respect to the solution process, we may be dealing with a well-conditioned problem (if the sensitivity is small) or with an ill-conditioned problem (if the system is highly sensitive to these changes). In the next sections we shall define more precisely the terms well-conditioned and ill-conditioned. In general, ill-conditioned will refer to systems which are close to being singular.

The main aim of this section is to find ways and means to assess if the computed solutions are accurate and, eventually, how to improve them.

ii) Effect of roundoff errors.

Let us assume that our estimation of A and b is exact. Then, we know that, because of roundoff errors, any computed solution is approximate, irrespective of whether it has been computed via a direct or an iterative method. The question to be resolved is how accurate is that computed solution.

Let \hat{x} be the approximate computed solution of the system (95), then we can define the <u>error</u> of our computed solution as:

$$\mathbf{e} = \mathbf{x} - \mathbf{\hat{x}} \tag{97}$$

Such an error cannot be computed directly from equation (97), since we do not know the value of the exact solution x.

If we premultiply our computed solution by the matrix A, we shall obtain a vector (i.e. $A \ x$), which is not exactly the original right hand side vector b; the difference between both will be referred to as the <u>residual</u> r:

$$\mathbf{r} = \mathbf{b} - \mathbf{A} \, \mathbf{\hat{x}} \tag{98}$$

The residual can be directly calculated from the computed solution x, and from the right hand side vector b (which is known).

Then, the residual tells us, in an indirect way, how far is the computed solution from the exact solution.

The error and the residual are closely related. In simple terms, it can be seen that if the residual is zero, then it means that our computed solution is exact and the error is zero as well. But the question still remains open in terms of whether a small residual is synonymous with a small error or not.

iii) An algorithm to compute the error: iterative improvement (iterative refinement).

To find a direct relationship between the error and the residual, we have to introduce equation (95) into equation (98):

$$\mathbf{r} = \mathbf{b} - \mathbf{A} \,\mathbf{\hat{x}} = \mathbf{A} \,\mathbf{x} - \mathbf{A} \,\mathbf{\hat{x}} \tag{99}$$

then, factorizing by A:

$$\mathbf{r} = \mathbf{A} \left(\mathbf{x} - \mathbf{\hat{x}} \right) \tag{100}$$

which, by equation (97) becomes:

$$\mathbf{r} = \mathbf{A} \mathbf{e} \tag{101}$$

Hence, once we know the residual (r) produced by the computed value (\hat{x}) , we can compute the error (e) by simply solving the following linear system:

$$A e = r \tag{102}$$

which is basically the same original system, but now with a different right hand side vector.

As we saw in previous sections, the solution of the linear system (102) can be obtained very easily if we have stored the original factorization of A from the solution for x and, in so doing, implies (n^2+n) multiplications/divisions plus (n^2-n) additions. Having calculated the error e, a new estimate of the vector x can be computed by:

$$x_{new} = x + e \tag{103}$$

The process can stop here if ||e||/||x|| is smaller than a prespecified accuracy, otherwise a new residual can be computed assigning x_{new} to x and repeating the previous sequence from equation (98) up to equation (103).

This is the <u>iterative improvement</u> procedure and it should converge to the exact solution, provided that the matrix A is well-conditioned. A high rate of convergence in the iterative improvement procedure is a good indication of a well-conditioned system.

In the previous computations, since the residual (r) becomes very small, its computation normally involves the use of doubleprecision arithmetic.

iv) Estimating limits for the error.

Returning to the question of the relationship between residual and error, equation (102) can be written as:

$$e = A^{-1} r \tag{104}$$

This relationship can be understood in the sense that e is proportional to r, with a proportionality constant represented by A^{-1} ; unfortunately, we do not know A^{-1} . Furthermore, we do not have a way for measuring the "size" of a matrix, if this size is going to be used as a proportionality factor. The latter problem is solved in a similar way as in the case of defining the "size" of a vector, i.e. using the concept of "norm", but this time applied to a matrix:

$$|| A || = \max \frac{|| A x ||}{|| x ||}$$
(105)

where the maximum is considered over all non-zero (nx1) vectors and A is a (nxn) matrix.

Again, as in the case of vector norms, we can have different matrix norms: euclidean norm, maximum norm, etc.

On using the following property of norms:

and considering equation (104) we get:

$$|| e || \leq ||A^{-1}|| || r ||$$
(106)

but, on applying the same to equation (101), we get:

 $|| r || \leq || A || || e ||$ (107)

which can be reordered as:

$$\frac{|| \mathbf{r} ||}{|| \mathbf{A} ||} \leq || \mathbf{e} || \qquad (108)$$

Hence, upon combining equations (108) and (106) we get:

$$\frac{|| r ||}{|| A ||} \leq || e || \leq || A^{-1} || || r ||$$
(109)

From equation (109) we can obtain lower and upper bounds for the relative error (||e||/||x||) in terms of the relative residual (||r||/||b||) :

$$\frac{|| b ||}{|| A || || x || || b || \leq \frac{|| e ||}{|| x ||} \leq \frac{|| A^{-1} || || b || || r ||}{|| x ||} (110)$$

But, from equation (109), for the particular case when x=0, where e = x and r = Ax = b, we get:

$$\frac{|| b ||}{|| A ||} \leq || x || \leq || A^{-1} || || b ||.$$
(111)

This means that, from the second part of (111):

$$\frac{|| \times ||}{|| b ||} \leq ||A^{-1}||$$
(112)

and, from the first part of (111):

$$\frac{|| \mathbf{b} ||}{|| \mathbf{x} ||} \leq || \mathbf{A} || \qquad (113)$$

Introducing equations (112) and (113) into equation (110) gives:

$$\frac{1}{||A|| ||A^{-1}|| ||b||} \leq \frac{||e||}{||x||} \leq ||A^{-1}|| ||A|| \frac{||r||}{||b||}$$
(114)

Equation (114) defines the lower and upper limits for the relative error ||e||/||x||.

The quantity $||A|| ||A^{-1}||$ is defined as the <u>condition number</u> of the matrix A:

Cond (A) =
$$||A|| ||A^{-1}||$$
 (115)

The condition number depends on the norm being used (euclidean, maximum, etc.), but it is always greater than 1.

For the particular case when Cond (A) = 1, equation (114) becomes:

$$\frac{|| r ||}{|| b ||} \leq \frac{|| e ||}{|| x ||} \leq \frac{|| r ||}{|| b ||}$$
(116)

which means that, in this case, the relative error is approximately the same as the relative residual. But, the greater the condition number, the less the information brought by the relative residual on the relative error.

In practice, the condition number is difficult to obtain, because the inverse matrix is usually not available.

The simplest estimate of ||e|| is the scalar $||\hat{e}||$, where \hat{e} is the solution of equation (102): A e = r.

An ill-conditioned system has a "large" condition number, where the value of "large" must be defined according to the particular computer being used. A large number is normally greater than 10^7 in single precision and greater than 10^{16} in double precision.

v) Effect of the errors in the matrix of coefficients.

Let us assume that the right hand side vector has been estimated accurately, then we are solving:

$$(A + \delta A) x = b$$
 (117)

where:

SA : matrix containing the error in the estimation of A. Let:

$$A = A + \delta A \tag{118}$$

be the estimated A matrix, and X be the computed solution of equation (117), then:

$$A \hat{\mathbf{x}} = \mathbf{b} \tag{119}$$

From equation (95):

$$x = A^{-1} b$$
 (120)

and introducing equation (119), we get:

$$\mathbf{x} = \mathbf{A}^{-1} \mathbf{A} \mathbf{\hat{x}} \tag{121}$$

but

or

$$\mathbf{A} = \mathbf{A} + \mathbf{A} - \mathbf{A} \tag{122}$$

and equation (121) becomes:

 $x = A^{-1} (A + A - A) \hat{x}$ (123)

 $x = \hat{x} + A^{-1} (A - A) \hat{x}$ (124)

recalling that $\delta A = A - A$, from (118), and introducing it into equation (124), we get, after reordering:

$$\mathbf{x} - \mathbf{\hat{x}} = \mathbf{A}^{-1} \left(\mathbf{\delta} \mathbf{A} \right) \mathbf{\hat{x}} \,. \tag{125}$$

Then, using the norm property: $||A B|| \leq ||A|| ||B||$, we get:

$$| \mathbf{x} - \hat{\mathbf{x}} || \leq || A^{-1} || || \delta A || || \hat{\mathbf{x}} ||$$
 (126)

but,

$$|| A^{-1}|| || \delta A || || \Re || = || A^{-1}|| || A|| \frac{|| \delta A ||}{|| A ||} || \Re ||$$

and so, equation (126) becomes:

$$\frac{|| \mathbf{x} - \hat{\mathbf{x}}||}{|| \hat{\mathbf{x}}||} \leq || \mathbf{A}^{-1} || || \mathbf{A} || \frac{|| \delta \mathbf{A} ||}{|| \mathbf{A} ||}$$
(127)

that is to say:

$$\frac{|| \times - \hat{\times} ||}{|| \hat{\times} ||} \leq \text{Cond} (A) \frac{|| \delta A ||}{|| A ||}$$
(128)

Equation (128) establishes that the greater the condition number of A, the greater the effect of the relative errors of the coefficients of A in the relative accuracy of the solution.

In addition, from equation (128), we know that if our estimation of the matrix of coefficients is accurate just up to the s-th decimal, and if Cond (A) \approx 10^t then, the maximum possible relative accuracy for x is 10^{t-s}; there is no point in trying to get a better solution because of the limited quality of the original data.

vi) Concluding remarks.

We have seen that computer solutions of linear systems of equations are in essence approximate, independent of whether the solution has been produced via a direct or an iterative algorithm.

The relative error depends on the condition number of the matrix of coefficients but, because condition numbers are difficult to estimate, a good and simpler approximation to the error is the norm $||\hat{e}||$, where \hat{e} is the solution of A e = r.

approximate solution & can be improved in its accuracy An via iterative improvement algorithm, which is not the too computationally expensive if the factorization of A has been There is a limit in the accuracy of \hat{x} in previously saved. the iterative improvement stage, since we cannot expect an estimate X more accurate than that allowed by the quality of the data.

A.3. Review of sparse direct methods for the solution of linear systems of equations.

A.3.1. Data structures.

As far as many engineering problems are concerned, the corresponding matrices generated in their solution are such that, even though their order may be in the range of hundreds or thousands, the number of non-zero elements is relatively small in comparison with the total number of elements. This is also the case in water distribution network analysis, where typically the matrices generated by the different algorithms have at most 5 - 6non-zero elements per row, irrespective of the size of the network. In practice, this means that for a network of, say, 1,000 nodes, i.e. producing a matrix with 1,000,000 elements, only 5,000-6,000 of these elements are actually non-zero. A11 these kinds of matrices are known as sparse matrices, where the relatively few non-zeros are scattered throughout the matrix.

Obviously, important savings of computer storage and execution time can be obtained by explicitly taking into account the sparsity of the matrix, since it is rather inefficient to store zeros and carry out additions or multiplications by zero. The

A-57

larger the size of the network, the bigger the relative savings of computer resources. In addition to that, some sparse matrices have special structures, like symmetric matrices, banded matrices or tridiagonal matrices; by explicitly taking into account the special structure of a matrix we can reduce even further the requirements of storage and computer time.

To take full advantage of the presence of a great deal of zeros in the sparse matrix, we need to store and operate only on the non-zeros and, to do so, we need some efficient way for keeping track of the position of all the non-zeros within the full matrix and, of course, of its numerical value. Roughly speaking, this means that we need to store the "coordinates" corresponding to the position of each non-zero within the matrix, which is normally dealt with by "cheap storage", say Integer*2 in a FORTRAN code, while the numerical value is kept in more lengthy storage, say Real*4 or Real*8 in a FORTRAN program.

There are a number of schemes for accessing the non-zeros in a sparse matrix, with a very different complexity/efficiency ratio. Simple data structures lead to relatively low efficiency schemes, while more complex data structures lead to the more efficient schemes.

Duff et al. (1986), Pissanetzky (1984) and George and Liu (1981), amongst others, give a detailed description of different data structures for sparse matrices; we are going to describe only some of them, especially those used in this Appendix. Since the matrices generated in water distribution network analysis are

A~58

symmetric positive definite matrices, we shall look carefully in that direction, although most of the topics covered in the next sections are applicable to more general matrices as well.

a) Coordinate scheme.

In this case we use two integer arrays to store the row and column indices of each non-zero, and a real number to store its non-zero value. The length of all these arrays is the same and equal to the total number of non-zeros within the matrix, say NZ. This data structure allows us to store the non-zeros in any order, and it can be be used straight-away for the case of nonsymmetric matrices.

The retrieval on the non-zero elements from the data structure is quite simple; the following example illustrates this point:

<u>Example</u>: Let a lower triangular matrix L of order N = 7 and 16 non-zeros (NZ = 16) be as follows:

then, the coordinate storage scheme, storing the elements by rows, would be: Row index: IROW IROW = (1, 2, 2, 3, 4, 4, 4, 5, 5, 5, 6, 6, 6, 7, 7, 7) Column index: ICOL ICOL = (1, 1, 2, 3, 1, 2, 4, 3, 4, 5, 3, 5, 6, 5, 6, 7)

Non-zero values array: VALUE

VALUE= (1₁₁, 1₂₁, 1₂₂, 1₃₃, 1₄₁, 1₄₂, 1₄₄, 1₅₃, 1₅₄, 1₅₅, 1₆₃, 1₆₅, 1₆₆, 1₇₅, 1₇₆, 1₇₇)

The length of each array is 16 (= NZ).

b) Implicit storage scheme.

A more elaborate scheme arises from observing that some further storage reduction can be achieved from the previous scheme. This can be done by avoiding the repetition of the row indices (a column oriented implicit scheme is also possible) and storing them in an implicit indexing scheme. The previous row index vector of length NZ is replaced by a new row index vector of length N+1, where N is the order of the matrix; the new row index points to the position of the first non-zero element in each row, which in our previous example means:

New row index: IROWI

IROWI= (1, 2, 4, 5, 8, 11, 14, 17)ICOL = (1, 1, 2, 3, 1, 2, 4, 3, 4, 5, 3, 5, 6, 5, 6, 7)VALUE= $(1_{11}, 1_{21}, 1_{22}, 1_{33}, 1_{41}, 1_{42}, 1_{44}, 1_{53}, 1_{54}, 1_{55}, 1_{63}, 1_{65}, 1_{66}, 1_{75}, 1_{76}, 1_{77})$

The column and non-zero array are the same as before.

Notice that since the difference between successive elements of the vector IROWI gives the number of non-zeros per row, we then need the element of order N+1 for completeness. Also note that IROWI(N+1)-1 corresponds to the amount of non-zero elements (NZ). The column-wise implicit scheme is straightforward. A further improvement of this scheme can be obtained by storing the diagonal non-zeros in a separate array, thus reducing the length of the column index array to NZ-N, and keeping the rest of the storage unchanged.

c) Compressed scheme.

One of the most efficient schemes is due to Sherman (1975) which takes advantage of the fact that, in the previous row-wise implicit scheme, the column index can be compressed even further by adding a new integer index of length N (integer array NZSUB). In the case of a column-wise implicit scheme it is the row index that can be compressed further.

In our previous example, presented by George (1981), this scheme is used for storing a lower triangular matrix by columns (the row-wise version is straightforward) :

Diagonal vector: DIAG DIAG = $(1_{11}, 1_{22}, 1_{33}, 1_{44}, 1_{55}, 1_{66}, 1_{77})$ Non-diagonal non-zero vector: LNZ LNZ = $(1_{21}, 1_{41}, 1_{42}, 1_{53}, 1_{63}, 1_{54}, 1_{65}, 1_{75}, 1_{76})$ XLNZ = (1, 3, 4, 6, 7, 9, 10)NZSUB = (2, 4, 5, 6, 5, 6, 7)XNZSUB = (1, 2, 3, 5, 6, 7, 8)

The compressed scheme splits the non-zero values into two real arrays, one for the diagonal elements DIAG (neither row nor column indices are needed for these elements) and another LNZ for the non-diagonal non-zeros.
The array XLNZ gives the amount of non-zeros per column (excluding the diagonal) and NZSUB gives the corresponding row index of these non-zeros using XNZSUB as a pointer to the first non-zero in NZSUB, as illustrated by the arrows in the previous structure. The extra savings produced by this scheme are due to the fact that in the array NZSUB some elements play a double (or even greater !!) role, like the 4 and 7 in the second and seventh position of NZSUB respectively.

The advantages of this more complex scheme become apparent in large systems, of the order of hundreds or thousands of equations.

<u>A.3.2. Fill-in</u>.

On using a direct method for the solution of a general linear system, we already know that the elimination process consists of transforming the original matrix into an upper (or lower) triangular matrix, which is subsequently solved via a backward (or forward) substitution. The fact that the original matrix is sparse does not generally imply that the transformed triangular matrix has the same sparsity pattern; in fact, it may not be sparse at all. Due to the elementary operations performed between the rows (or columns) of the original matrix, new nonzeros can be produced in places where a zero element existed. The newly created non-zeros are referred to as <u>fill-in</u> (or simply fill).

One of the worst cases of fill-in that can be expected is shown in Figures A.2 and A.3. The original system has the arrow shaped matrix shown in Fig. A.2.

If we factorize A into the Cholesky factors $L L^T$, then L has the structure shown in Fig. A.3, which corresponds to a lower triangular fully populated (dense) matrix.

Fig. A.2. Arrow shaped linear system.

Fig. A.3. Cholesky factor of arrow shaped matrix (without reordering).

However, if we permute rows and columns 1 and 5 in the original matrix A, producing the reordered matrix A^* shown in Fig. A.4,

Fig. A.4. Reordered arrow shaped linear system.

The corresponding new Cholesky factor L^* will have the sparsity pattern shown in Figure A.5.

Fig. A.5. Cholesky factor of reordered arrow-shaped linear system.

Comparing A^* with its factor L^* , we can see that no fill-in has been produced.

This extreme example illustrates very clearly the main feature on which most of the sparse techniques are based, namely, the fact that fill occurs only because of the way in which the original matrix has been set up; consequently, <u>fill can be</u> <u>reduced and minimized via a sensible reordering</u> of the original system. Thus, one of the first tasks to be tackled in sparse matrix handling is the search for an efficient ordering scheme that can reduce the amount of fill.

Formally, the reordering process can be understood as applying some permutation matrix P to the original linear system. That is to say, instead of solving:

$$A \mathbf{x} = \mathbf{b} \tag{95}$$

we solve:

$$(PAP^{T})Px = Pb$$
 (96)

This is because a permutation matrix P is also an orthogonal matrix, i.e.:

$$P P^{T} = P^{T} P = I$$
 (identity matrix) (97)

In this context, finding an adequate reordering strategy can be equivalent to find the matrix P which minimizes the fill in the matrix P A P^{T} .

Graph theory provides a good representation of the reordering process; for example, for the symmetric matrix given by equation (98):

$$A = \begin{vmatrix} 1 & 7 & & 8 \\ 7 & 2 & 9 & 10 \\ & 9 & 3 & 11 \\ & 10 & 4 \\ & & 11 & 5 & 12 \\ 8 & & & 12 & 6 \end{vmatrix}$$
(98)

the associated graph representation is shown in Fig. A.6. A graph associated with a matrix is the set $G^A = (X, E)$, where X is the set of nodes (one node per row or column) and E is the set of edges (one per each non-diagonal non-zero).



$$X = \{1, 2, 3, 4, 5, 6\}$$

E = {(6, 1), (1, 2), (2, 4), (2, 3),
(3, 5), (5, 6)}

Fig. A.6. Undirected graph corresponding to symmetric matrix given by equation (98).

Although graph theory concepts are very useful to understand some aspects of sparse matrices, we shall try to keep its use at a minimum level, to avoid further complexities. Any definition will be given immediately before it is needed; for details see either Wilson (1985), Deo (1974), Harary (1972) or Harris (1970). Note that for non-symmetric matrices, directed graphs (i.e. graphs whose edges have a direction associated with them) are needed to account for the non-symmetric connectivity.

To permute the 3^{rd} and 4^{th} rows and columns of A is equivalent to transforming A into (P A P^T), with the permutation matrix shown in Fig. A.7. which is an identity matrix with its 3^{rd} and 4^{th} rows interchanged.

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

Fig. A.7. Permutation matrix for interchanging rows 3 and 4.

On performing the pre and post-multiplication of A by P and P^{T} , we get the matrix and associated graph shown in Fig. A.8.



Thus, in graph terms, we can see that reordering does not mean a change in the <u>structure</u> of the matrix graph, but only changes in its <u>labelling</u>; in the previous example this swap of labels happens between nodes 3 and 4.

From the previous example it is also evident that if A is symmetric, then P A P^{T} is symmetric as well; this is due to the fact that the premultiplication by P interchanges the rows, while the postmultiplication by P^{T} swaps the columns.

Another important consideration apparent from the last example is that <u>for finding a "sensible" reordering</u> (or re-labelling in the graph of the original system), we do not actually need its <u>numerical values</u>. <u>but only its structure</u>. This is a very important feature, since it means that we can split the Gaussian elimination problem into a sequence of independent stages. <u>This</u> is true for symmetric positive definite matrices only. Usually the whole process is separated into the following stages:

i) ANALYSIS stage: this finds the minimum fill-in ordering and sets up the corresponding data structures.

ii) FACTORIZATION stage: where the numerical part of the factorization of A is carried out.

iii) SOLUTION stage: where the solution of the system takes place, for a given right hand side vector.

This three-stage process allows us to make more efficient the overall solution process, since if, for example, we need to solve a linear system for different right hand side vectors, then we only need to repeat the last SOLUTION stage, provided that we have stored the information from the ANALYSIS/FACTORIZATION stages. At the same time, if only the numerical values of the matrix are changed (but not its structure), then we need to repeat the FACTORIZATION and SOLUTION stages. This maximizes the

performance of the Gaussian elimination process, since it is quite usual that the ANALYSIS needs to be performed only once, while the remaining stages can be repeated as many times as needed, using the same ANALYSIS results.

As pointed out by many researchers in this area: Irons (1970), Reid (1977) and George (1981), the separation of the whole process into 3 independent stages is only suitable for symmetric positive definite matrices, where no risk of instability occurs. That is to say, in symmetric positive definite matrices we do not need to pivot for stability and then, we can concentrate our attention in reducing the fill-in. For more general sparse matrices, permutations are needed for stability reasons and then, the data structures are dependent on the numerical values and not only on the structure (or connectivity) of the matrix. This is clearly a conceptual problem, with very important practical consequences.

As a result, we have seen that the amount of fill produced in the Gaussian elimination process depends on the way we set up and reorder the original matrix. The problem is, as far as reducing the fill is concerned, to find an adequate reordering scheme that minimizes the amount of fill (if such a reordering exists). It must be emphasised that normally minimizing the amount of fill not only leads to a reduction in the storage requirements, but also to a reduction in the amount of computational work to be done.

A.3.3. Sparse methods for banded matrices.

Historically, banded matrices such as those arising from many differential equation problems, have been the starting point of the development of sparse matrix techniques. Even if the original linear system does not emerge in a banded form naturally, the aim has been to reorder it in such a way that the reordered system could have all its non-zeros "clustered" around the diagonal.

The key observation in sparse banded methods is that when factorizing a banded matrix, the eventual fill is concentrated within the band; thus, all the zeros outside the band can be ignored (or "exploited" in the sense of sparse matrix handling), and only the elements within the band are stored, irrespective of the fact that some of them may be zeros.

The storage scheme used for dealing with banded matrices has not been explained before, but simply reshapes the sub-diagonals of the banded matrix into vertical columns as shown in Fig. A.9.

a ₁₁ symmetric a ₂₁ a ₂₂ a ₃₃ a ₄₃ a ₄₂ a ₄₃ a ₄₂ a ₄₃	$=====> \begin{bmatrix} - & - & a_{11} \\ - & a_{21} & a_{22} \\ 0 & 0 & a_{33} \\ a_{42} & a_{43} & a_{44} \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ 0 & 0 & a_{nn} \end{bmatrix}$
<>	<>
half band width	half band width

Fig. A.9. Diagonal storage scheme for banded matrices.

The main advantage of this approach is its simplicity, both in terms of the data structures and in terms of programming. Of course, it is evident from Fig. A.9., that it is a sub-optimal solution, in the sense that some fill-in can still be exploited with very little extra effort. This is done in the next method.

Additionally, the lack of optimality of banded methods is due to other reasons, as can be shown using the same arrow shaped matrix already seen in Fig. A.2, where a band minimization approach produces the matrix (P A P^{T}) shown in Fig. A.10.

Fig. A.10. Minimum band ordering for arrow shaped linear system.

Clearly, the ordering shown in Fig. A.10 is poorer than that shown in Fig. A.4, where the permutation of rows and columns 1 and 5 produced no fill-in at all in the Cholesky factorization stage.

Cuthill and McKee (1969) developed an algorithm for producing minimum bandwidth orderings. The algorithm needs a starting or "root" node (r), and the original graph is relabelled starting from "r" and renumbering first those nodes with the minimum number of edges connected to them (i.e. minimum degree nodes). The example shown in Fig. A.11, taken from George (1981) illustrates how the algorithm works.



Fig. A.11. Two examples of Cuthill-McKee orderings for reducing bandwidth, from George (1981).

Two important features of the algorithm become apparent from this example:

a) The ordering (or labelling) depends on the selected root node.

b) Even for the same root node, different orderings can be obtained when the neighbouring unnumbered nodes have the same degree (like when numbering nodes 2 and 3 in Fig. A.11).

The classical example of the arrow shaped matrix showed quite clearly that minimizing the bandwidth has nothing to do with minimizing the fill-in (and nothing to do with minimizing the amount of work to be done). As a result, the method is suboptimal and because the extra effort needed for a further improvement is small, the band methods are not widely used today, except for problems where the banded structure comes up naturally (e.g. tridiagonal matrices and similar).

A.3.4. The envelope ("skyline") method.

The key observation leading to the envelope method is that fill-in actually occurs only in those positions between the first

non-zero in each row and the diagonal of the matrix, that is to say, within the "skyline"-shaped bordering line which surrounds all the non-diagonal non-zeros in the matrix. Figure A.12 shows this feature, by shading the location of possible fill-in.



Fig. A.12. Skyline representation.

In this example, fill in <u>may</u> only occur in the fourth and sixth rows, but never outside of the envelope (shaded area).

The number of non-zero elements inside the envelope is usually referred to as the <u>profile</u> of the matrix; thus, in the previous symmetric example [equation (98)], the profile is equal to 7.

In this context, orderings minimizing the profile of a matrix lead to low fill, and algorithms producing minimum profile are efficient from the storage point of view. One of the most widely used profile reduction algorithms was proposed by George (1974), who found that by simply reversing the ordering produced by the Cuthill-McKee algorithm eventually reduces the profile (in fact it does not make it worse), and the algorithm has been called "Reverse Cuthill-McKee" algorithm.

As pointed out in the case of the standard Cuthill-McKee algorithm, the ordering depends on the selection of the root node. Obviously the same happens in the Reverse Cuthill-McKee (or RCM) algorithm. Hence, we need to find a root node useful for our purposes.

It has been empirically found that good root nodes are those lying on the periphery of the graph (peripheral nodes); these are nodes such that the amount of edges between them and any other node in the graph (which is defined as the "eccentricity" of the node) is maximum.

George and Liu (1981) proposed the following algorithm for finding peripheral nodes; because the algorithm does not guarantee peripheral nodes, the nodes found by it are referred to as "pseudo-peripheral" nodes.

The algorithm relies on the concept of <u>rooted level</u> structure, which is explained immediately. Following the George and Liu (1981) definition, a level structure rooted at node "x" is a partitioning $\int (x)$ of the set of nodes of the matrix graph, such that:

 $f(x) = \{ L_0(x), L_1(x), \dots, L_{l(x)}(x) \}$

where:

x : is the root node of the level structure.

l(x) : is the eccentricity of node "x".

Rooted at "x", the following levels are defined, each one containing the adjacent nodes not previously included in a level: $L_0(x) = \{x\}$ $L_1(x) = \{Nodes adjacent to L_0(x)\}$: : $L_i(x) = \{Nodes adjacent to L_{i-1}(x), except those in L_{i-2}(x)\}$ $\forall i=2,3, ..., 1(x)$ Figure A.14 is an example of a rooted level structure, corresponding to the matrix whose graph is shown in Fig. A.13

In the level structure shown in Fig. 14 $\Gamma(x_5)$ the eccentricity of x_5 is $l_{(x5)} = 3$ and

$$\underline{\Gamma}$$
 (x₅) = { L₀(x₅), L₁(x₅), L₂(x₅), L₃(x₅) }



Fig. A.13. Graph for level structure example, from George and Liu (1981).



b) Level structure:



Fig. A.14. Level structure rooted at x_5 , for graph shown in Fig. A.13, from George and Liu (1981).

The algorithm proposed by George and Liu (1981) for finding pseudo-peripheral nodes is reported to be based on a previous one, proposed by Gibbs et al. (1976) and its steps are:

a) Pick up an arbitrary root node "r".

b) Generate the level structure rooted at node "r".

c) Choose a node in the last level $L_{l(r)}(r)$ with the minimum possible degree, say node "x".

d) Generate another level structure now rooted at "x".

If the eccentricity of "x" is greater than the eccentricity of "r", then assign r <--x and repeat the process from step c), otherwise:

e) "x" is a pseudo-peripheral node and stop.

For the previous example, the process of finding a pseudoperipheral node can be represented as in Fig. A.15.

Notice that in the previous example, node x_1 can be labelled as a pseudo-peripheral node as well.

On using the previous ideas, the following is the sequence for the Reverse Cuthill-McKee algorithm:

- Step 1) Determine a pseudo-peripheral node "r" and label it as x_1 . This can be done via the algorithm already introduced or another one.
- Step 2) Loop : label the unlabelled neighbours of x_i in an increasing order of degree, ¥ i=1,2, ..., N. This is the standard Cuthill-McKee algorithm.
- Step 3) Reverse the labelling: producing a new ordering A-75

(labelling), say:

 $y_1, y_2, ..., y_N$ such that $y_i = x_{(N-i+1)} \forall i=1, 2, ..., N$

As we have already pointed out, step 1 can be replaced by any other algorithm or, indeed, the root node can be specified by the user. In the context of water distribution networks which are fed gravitationally, a good candidate can be the node corresponding to the upstream reservoir in the system. In more general networks, with supply coming from gravitational and pumped sources, the previous choice might not be the best one.

First try: steps a) and b), $r = x_3$



eccentricity = 2



Third try: $r \leftarrow -x$ and back to steps c) and d)



Fig. A.15. Algorithm for finding a pseudo-peripheral node, from George and Liu (1981).

We have implemented the envelope method, with the RCM ordering, using the subroutines published by George and Liu (1981). We used implementation for solving the linear systems generated in this water distribution analysis. We would like to emphasise here that one of the main drawbacks found, deals with the amount of storage needed for the non-zeros out of the diagonal (array ENV). The real storage required is not predictable before running the program and an upper bound has to be defined on a trial-and-error basis using the program with problems of different sizes (although this procedure does not guarantee an upper limit); we found, using synthetic networks of up to 5000 nodes (\$ 5000 equations), that the maximum length of the array ENV tends to Also, the execution time obtained with this method is of N1.5 the order of 3 times the corresponding time of the efficient methods. For this reason we do not elaborate more on this algorithm.

A.3.5. Minimum degree algorithms.

One of the main shortcomings of the previous envelope method was its sub-optimality in terms of storage; in fact a great deal of the storage for the non-zeros enveloped within the skyline was indeed occupied by zeros, and the idea is now try to <u>exploit all</u> <u>the zeros</u>. The price to be paid is obviously further complications in data structures, algorithms and programming.

Algorithms based on the minimum degree concept have been used since Markowitz (1957), in linear programming applications, and subsequently by Tinney and Walker (1967), Ogbuobiri, Tinney and

Walker (1970) and Ogbuobiri (1970), in power networks.

Zollenkopf (1970) used a minimum degree approach for pivot selection within a sparse factorization method for symmetric positive definite matrices in power systems. Pivoting is apparently needed in this case because the author considered also the case when the matrix is asymmetrical in the element values (but remaining symmetric in its structure).

Shamir (1973) applied the concepts of minimum degree to water distribution networks, though he used a mixed strategy, pivoting both for size (partial pivoting) and density (to reduce fillin).

Nowadays it is understood that, for symmetric positive definite linear systems, pivoting is not necessary to ensure stability. As we said earlier, this means that in symmetric positive definite matrices, we can concentrate in the reduction of fill-in exclusively, without risking instabilities. Thus, the ANALYSIS stage does not include the numerical values of the non-zeros, but only their structure. Recall that for general sparse matrices this is not valid, and pivoting for stability is essential.

As we saw earlier, pivoting for stability meant that during the Gaussian elimination process, the next pivot element was chosen from the row which had the biggest diagonal absolute value (partial pivoting). We also mentioned full pivoting as another alternative.

Minimum degree consists basically in choosing the next pivot element, at each stage of the elimination sequence, from the row

```
A-78
```

with the fewest amount of non-zero elements.

It has to be said that when pursued simultaneously, pivoting for density and pivoting for stability are rather conflicting objectives, but this is not our problem in symmetric positive definite matrices.

An example of the use of the minimum degree algorithm is presented in Figure A.16, which has been taken from George (1981).

In its traditional implementation, which stores the Cholesky factors <u>explicitly</u>, the minimum degree algorithm generates a great deal of fill, requiring important quantities of storage; unfortunately, to make things worst, it happens that the amount of required storage is unpredictable until the factorization is completed. Thus, when using this implementation the user does not know how much storage will be needed, and the only way to tackle this problem is via experimentation.

George and Liu (1981) published a solution to this problem, using an implicit storage scheme, which has a predictable storage at the ANALYSIS stage. Thus, the computer program can determine, before the numerical part of the solution is started, what is the ofamount storage required. Α so-called fast hybrid implementation, which lies in between the explicit and implicit storage scheme, was compared by George and Liu (1981), chapter 9, and was found to be one of the best available from the point of view of the execution time; that comparison was carried out with problems of up to 2233 equations arising from finite element

problems.

Since we have used this implementation in our comparisons in the context of water distribution networks, we shall spend the rest of this section explaining, in the simplest possible terms, how this scheme works; for a more detailed explanation, see George and Liu (1981, chapter 5). The implementation relies heavily on graph theory concepts, especially in what is called "quotient graph", which is also the base for other forthcoming methods reviewed in this section.

Before explaining what is the fast hybrid implementation of George and Liu, we need to explain what is the implicit scheme. To do so, we shall introduce graphs into this problem, in order to visualise the process of creation of fill-in. The same problem of Gaussian elimination shown in Fig. A.16, can be represented using graphs as in Figure A.17.

Figure A.17 is self-explanatory and shows how the elimination sequence, using the minimum degree algorithm, takes place. The order of elimination is then:

(a, c, g, i, f, e, b, h, d)

In graph terms, fill-in, i.e. the creation of new darker edges in Fig. A.17, can be explained as follows: let us consider for example the elimination of "a" in G^{O} (see Fig. A.17), since "d" is "reachable" from "b" through "a". Then, in order to maintain the same "reachability" of the original graph G^{O} , after the elimination of "a", we need to add a new edge which accomplishes this purpose, otherwise the "reachability" or connectivity

р р bcdefghi a defghi Ъ С * * р a * * * b ⊕ * b * * * * С p * d e **(**) С * * * d * ж ж * f е ж * * ж ж f g ¥ * ж * g * h * * * * h * i * * * ж * * i * * * Ao A_1 р р р f f d bhi f b d е ghi е b h d е * * * * \oplus * Ð d * * \oplus \oplus * b * е ⊕ * * p f d * e * * * * * * *⊕⊕ * f Ð e f * \oplus * * * * * b * * Ð b *⊕ († († \oplus * h * Đ * * * * g * * × h * d * $\oplus \oplus$ * p * i * * h * * * p * * i * * * A_2 A_4 A₃ p p p p h d h d e b h d ъ d * * * * рb ∗ \oplus \oplus рh * 🕀 * рd р е Ð * + + d ⊕¥ Ъ * h ¥ (\mathbf{f}) \oplus (+)* d h * Ð É Ð d * * A5 A₆ A7 A₈

Notes: *****

- (1) "p" marks the new pivot.
- (2) (+) denotes fill-in.
- (3): the pivots are always swapped with the first row (and column), and then eliminated from the following graph.
- Fig. A.16. Example of application of minimum degree algorithm, from George (1981).



information contained in G^0 will be lost. On using the same idea throughout the whole elimination process, the new graph corresponding to the filled matrix of the original linear system becomes that shown in Fig. A.18.



Note: darker lines indicate that fill-in will occur in the original matrix, when the factorization is carried out following a minimum degree algorithm.

Fig. A.18. The filled graph of $F = L + L^{T}$.

The corresponding filled (reordered) matrix is shown in Fig. A.19, where an "plus" (+) sign denotes fill, which corresponds to the darker lines in Fig. A.18.

	8	С	g	i	f	е	Ъ	h	d
•	۲.						÷		ц. Т
a -	Ι T						Ť		*
е		*			*		*		
g			*					*	*
i				*	*			*	
f]	*		*	*	*	Ð	Ð	
e	1				*	*	×	Ť	*
b	*	*			Ð	*	*	(+)	\oplus
h	[*	*	(Ŧ)	*	(\mathbf{f})	¥	- Ā
d	*		*			*	Ť	\oplus	×

Fig. A.19. Filled (reordered) matrix.

Notice that G^F (Fig. A.18) has been found using only the information on the matrix structure, the numerical values not having been used at all.

According to Fig. A.17, the whole elimination can be represented as a succession of graphs: G^0 , G^1 , G^2 , ..., G^8 . This is an <u>explicit</u> representation.

The <u>implicit</u> scheme arises on formalising the "reachability" concept used to explain the development of Fig. A.17. Following George and Liu (1981), for a graph G = (X, E) and S being a subset of the nodal set X and taking $x \notin S$, then x is said to be reachable from another node "y" if there is a path through S connecting x and y (i.e. a series of nodes in between them). That path is formally written as:

(y,
$$\mu_1$$
, μ_2 , ..., μ_k , x / from y to x such that $\mu_i \in S$, $1 \le i \le k$
with $k \ge 0$)

Note that a neighbouring node of y not in S is also reachable from y (with k=0).

The set of all the reachable nodes from y through S is defined as:

Reach(y,S) = { $x \notin S$ such that x is reachable from y through S }

The practical application of the reach operator, as defined above, becomes clear if we define S as a variable set in the elimination process represented by Fig. A.17, which changes at each stage of the elimination process according to:

 $S_i = \{ x_1, x_2, \dots, x_i \}$

where x_i represent the nodes in the order in which they are eliminated.

Hence, for the elimination process shown in Fig. A.17 we have:

Reach(a,S _o)	= {	b, d }	,	$S_{o} = \{ \emptyset \}$	
$Reach(c, S_1)$	= {	b, f }	,	$S_1 = \{a\}$	
$Reach(g, S_2)$	= {	h, d }	,	$S_{2}^{-} = \{ a, c \}$	
Reach (i, S_3)	= {	f, h }	,	$S_3 = \{ a, c, g \}$	
$Reach(f, S_4)$	= {	e, b, h }	,	$S_4 = \{a, c, g, i\}$	(99)
$Reach(e, S_5)$	= {	b, h,d }	,	$S_5 = \{a, c, g, i, f\}$	/
$Reach(b, S_6)$	= {	h, d }	,	$S_6 = \{a, c, g, i, f, e\}$	
$Reach(h, S_7)$	= {	d }	,	$S_7 = \{a, c, g, i, f, e, b\}$	
Reach(d, S_8)	= {	Ø }	,	$S_8 = \{a, c, g, i, f, e, b, h\}$	ŀ

Then, on comparing the results of equation (99) with Fig. A.19 we can see that the Reach operator applied successively on a, c, g, ... etc., and S_i being the set of previously eliminated nodes, gives the column indices (or row indices because of the symmetry) corresponding to all the non-zeros in the Cholesky factorization in the subset of the uneliminated nodes (i.e. after the diagonal in Fig. A.19).

With this representation and with the Reach operator already defined, we have an <u>implicit scheme</u>; the sequence of graphs is no longer needed and all we have to do is to find a way of handling the Reach operator in the computer.

As we can intuitively assess, the main shortcoming of the implicit scheme is the fact that the amount of work involved in determining the reachable sets can become too large.

The fast hybrid solution proposed by George and Liu (1981) is a compromise between the explicit and implicit schemes. It arises by observing that because the computation of $\operatorname{Reach}(x_i, S_i)$ depends on the length of the paths between the node being eliminated and that in its reachable set (through S_i), the longer the paths, the more work needed to find the reachable nodes. The solution is then to reduce the length of these paths. On the

other hand, because we are interested only in finding the reachable nodes from the set of uneliminated nodes, we do not need the already eliminated nodes. In the explicit approach (Fig. A.17) we have actually deleted the eliminated nodes from the sequence of elimination graphs. In the formal implicit model, using the Reach operator, we can reduce the length of the path by "glueing" the neighbouring nodes already eliminated, creating what might be called "supernodes", like those shown in Fig. A.20, which is a graphical representation of the process of finding the reachable sets through the subsets S_i , where each S_i is formed by the already eliminated nodes. Each graph of Fig. A.20 is a <u>quotient graph</u> (or a "supergraph", i.e. a graph consisting of "supernodes").

In so doing, the paths between pivots and reachable nodes through S_i have a length of 2 (when it is through S_i) or 1 (when the reachable node is the neighbour of the pivot and it does not belong to S_i). We have then reduced the path lengths up to a maximum of 2.

The relationship between the fast hybrid approach using quotient graphs (Fig. A.20) and the explicit elimination scheme (Fig. A.17) becomes quite evident. It can also be deduced from Fig. A.20 [proved by a theorem in George and Liu (1981)] that the sequence of quotient graphs can be implemented using the same storage locations used for storing the original matrix. Then, the results of the ANALYSIS stage do not need more storage than that required for storing the original matrix.



Fig. A.20. Sequence of quotient graphs Γ_i used for finding the reachable set of nodes in the elimination process.

In summary, on applying the fast hybrid scheme of George and Liu (1981), we are able to determine the sequence of reachable sets [equation (99)] for each elimination step, which formally gives the structure of the Cholesky factorization $L L^T$ of the original matrix. This means that we get the amount of fill-in and the locations within the factorized matrix where the fill-in takes place, which allows us to set up the data structures needed for the following stages (FACTORIZATION/SOLUTION). All this can be done, when dealing with a symmetric positive definite matrix, without using the numerical values of the non-zeros . The search of the reachable sets is now meant to be efficient in terms of storage and speed, due to the introduction of "supernodes" and quotient graphs.

The minimum degree concepts give us the order in which the nodes are eliminated (e.g.: a, c, g, ..., h, d in our example) while the hybrid scheme gives the structure of the fill-in which such ordering produces.

additional improvement has been used in the implementation An provided by George and Liu (1981), and it is geared to reducing the time spent on searching for a minimum degree node (see Fig. A.17). The idea is, because we usually find more than one node with minimum degree (as c, g and i in Fig. A.17, during the first and second stages of the elimination), then we can use this information in order to eliminate these nodes as well, without having to carry out another minimum degree search in the next elimination stage. The enhancement used by George and Liu (1981) is based on that idea, and these nodes are said to be indistinguishable with respect to elimination. For more details see George and Liu (1981), chapter 5.

Despite the long reputation of the minimum degree algorithm and the much proclaimed advantages of the George and Liu (1981) implementation, we found in our testing that it is not as

successful either from the storage or from the execution time point of view. Moreover, we could not solve networks larger than 900 nodes within a reasonable execution time. So far, we have not been able to find an explanation for this behaviour; possible causes are the particular computer implementation used and the fact that the connectivity and "shape" of our test networks is too different from that used by George and Liu (1981).

A.3.6. Quotient tree algorithms.

These algorithms take advantage of the fact that an efficient way to deal with the sparseness of a linear system is to partition the matrix in blocks, some of which may eventually be dense matrices, but of a much more reduced order than the original system.

The advantages of this block-partitioning approach can be shown, for example, when instead of solving our model problem:

$$A x = b \tag{100}$$

where:

A : is a symmetric positive definite (n x n) matrix. we solve an equivalent 2x2 block partitioned system:

$$\begin{bmatrix} A_{11} & | & A_{12} \\ A_{21} & | & A_{22} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$
(101)

where:

 A_{11} : a (p x p) submatrix. A_{22} : a (q x q) submatrix, such that p + q = n. A_{21} : a (q x p) submatrix, such that $A_{21} = A_{12}^{T}$. x_1 , x_2 : partitioned vector of unknowns.

b₁, b₂ : partitioned right hand side vector.

Then, the corresponding Cholesky factorization of $A = L L^T$, can be expressed as a block-partitioned matrix as well:

$$L = \begin{bmatrix} L_{B} & | & 0 \\ \frac{1}{W^{T}} & | & L_{C} \end{bmatrix}$$
(102)

where:

 L_B : is the Cholesky factor of A₁₁, a (p x q) submatrix. L_C : is the Cholesky factor of A₂₂, a (q x p) submatrix. W : is a (q x q) submatrix and W^T its transpose.

The Cholesky factors can be computed algebraically by imposing the equivalence $A = L L^T$ and using equations (101) and (102):

$$\begin{bmatrix} L_{B} & | & O \\ W^{T} & | & L_{C} \end{bmatrix} \cdot \begin{bmatrix} L_{B}^{T} & | & W \\ 0 & | & L_{C}^{T} \end{bmatrix} = \begin{bmatrix} A_{11} & | & A_{12} \\ A_{21} & | & A_{22} \end{bmatrix}.$$
 (103)

Performing the multiplication on the left hand side of (103) we get:

$$\begin{bmatrix} -\frac{L_{B} \ L_{B}^{T} \ H_{B}^{T} \ H_{B}^{T} \ H_{B}^{T} \ H_{C}^{T} \ W^{T} \ W + L_{C} \ L_{C}^{T} \end{bmatrix} = \begin{bmatrix} A_{11} \ H_{12} \ A_{12} \ A_{21} \ H_{22} \end{bmatrix}$$
(104)

Then, identifying block by block in (104), we obtain: a) Upper left hand side block:

$$A_{11} = L_B \quad L_B^T \tag{105}$$

which simply states that L_B is the Cholesky factor of A_{11} . In other words, we can obtain L_B (or L_B^T) from A_{11} via Cholesky factorization.

b) Upper right hand side block:

$$L_B W = A_{12}$$
 (106)

which implies that having obtained L_B from (105), we can then obtain W by solving "q" linear systems of the form:

$$L_{B} W(*,i) = A_{12}(*,i) \quad \forall i=1,2,..,q$$
 (107)

where:

W(*,i) : denotes the i-th column (vector) of matrix W.

 $A_{12}(*,i)$: denotes the i-th column (vector) of matrix A_{12} . c) Lower right hand side block:

$$W^{T}W + L_{C}L_{C}^{T} = A_{22}$$
 (108)

which, on defining the auxiliary matrix C:

$$C = L_C L_C^T$$
(109)

becomes, after reordering :

$$C = A_{22} - W^T W$$
(110)

Hence, after having obtained W from step b) we can compute the auxiliar matrix C from (110) and, with C, we can get L_C which is the Cholesky factor of C in equation (109).

Thus, in this three stage process we are obtaining the Cholesky factors L and L^{T} by computing their corresponding submatrices L_{B} , W and L_{C} from A₁₁, A₁₂ and A₂₂.

With the Cholesky factorization of A, the solution stage of our original problem can be carried out by replacing our original linear system (using equation 103) by:

$$\begin{bmatrix} L_{B} & : & O \\ \frac{1}{W^{T}} & : & L_{C} \end{bmatrix} \cdot \begin{bmatrix} L_{B}^{T} & : & W \\ \frac{1}{O} & : & L_{C}^{T} \end{bmatrix} \cdot \begin{bmatrix} x_{1} \\ \frac{1}{x_{2}} \end{bmatrix} = \begin{bmatrix} b_{1} \\ \frac{1}{b_{2}} \end{bmatrix}$$
(111)

and then, because we already know L_B , L_C and W^T , its solution becomes just a series of substitution processes; indeed on defining the auxiliar vector:

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{\mathbf{B}}^{\mathrm{T}} & | & \mathbf{W} \\ -\mathbf{D} & | & \mathbf{L}_{\mathbf{C}}^{\mathrm{T}} \end{bmatrix} \cdot \begin{bmatrix} \mathbf{x}_1 \\ -\mathbf{x}_2 \end{bmatrix}$$
(112)

equation (111) becomes:

$$\begin{bmatrix} L_{B} & | & 0 \\ W^{T} & | & L_{C} \end{bmatrix} \cdot \begin{bmatrix} y_{1} \\ -y_{2} \end{bmatrix} = \begin{bmatrix} b_{1} \\ -b_{2} \end{bmatrix}$$
(113)

and the auxiliary vector $y = (y_1 + y_2)^T$ can be computed by the following steps:

i) Solving the upper part of (113):

$$L_B y_1 = b_1$$
 (114)

we obtain y_1 , by forward substitution.

ii) Compute, from the lower part of (113):

$$W^T y_1 + L_C y_2 = b_2$$
 (115)

which reordered gives:

$$L_{C} y_{2} = b_{2} - W^{T} y_{1}$$
 (116)

iii) Having y₁ from step i), y₂ can again be computed by a forward substitution from equation (116), which completes the computation of the auxiliary vector "y".

Now, having computed the vector y, we replace it into equation (112) and solve this last equation for x:

$$\begin{bmatrix} \mathbf{L}_{\mathbf{B}}^{\mathrm{T}} & | & \mathbf{W} \\ \hline \mathbf{0} & | & \mathbf{L}_{\mathbf{C}}^{\mathrm{T}} \end{bmatrix} \cdot \begin{bmatrix} \mathbf{x}_{1} \\ \hline \mathbf{x}_{2} \end{bmatrix} = \begin{bmatrix} \mathbf{y}_{1} \\ \hline \mathbf{y}_{2} \end{bmatrix}$$
(117)

which can be done via the following steps:

iv) Solving the lower part of (117):

$$L_{C}^{T} \mathbf{x}_{2} = \mathbf{y}_{2} \tag{118}$$

we obtain x_2 , by back substitution.

v) Compute the upper part of (117):

$$L_{B}^{T} x_{1} + W x_{2} = y_{1}$$
(119)

which reordered gives:

$$\mathbf{L}_{\mathbf{B}}^{\mathrm{T}} \mathbf{x}_{1} = \mathbf{y}_{1} - \mathbf{W} \mathbf{x}_{2} \tag{120}$$

vi) Having x_2 from step v), x_1 can be obtained by back substitution from (120).

What we have done so far is a bit of algebra, but the full advantage of the block partitioning approach has not become clear. The point is that, in the solution stage [steps i) to vi)], we have used the matrices L_B , L_C and W <u>explicitly</u>, and, because W is normally a full (dense) matrix, a great saving of storage can be achieved via a sensible way of handling this matrix. In fact, the ideal approach would be not to store it but to handle it in an implicit way; this can be possible by recalling that W can be written from equation (106) as:

$$W = L_{\rm B}^{-1} A_{12} \tag{121}$$

i.e. :

$$W^{\rm T} = A_{12}^{\rm T} L_{\rm B}^{-\rm T}$$
 (122)

This means that whenever we need to multiply W or W^{T} by a vector, like in equations (116) and (120), we can do it without including W or W^{T} in an explicit way; in fact, using equation (122), equation (116) becomes:

 $L_{C} y_{2} = b_{2} - (A_{12}^{T} L_{B}^{-T}) y_{1}$ (123) and, using (121), equation (120) becomes:

$$L_B^T x_1 = y_1 - (L_B^{-1} A_{12}) x_2$$
 (124)

We can also use the implicit form of W and W^T [equations (121) and (122)] to compute the product W^T W in equation (110):

 $W^{T} W = \{ A_{12}^{T} [L_{B}^{-T} (L_{B}^{-1} A_{12})] \}$ (125)

where the computations are carried out following the parenthesis sequence in (125) from the inside out, which is referred to as an asymmetric block factorization. The method then relies heavily on back and forward substitution routines, which normally are used in the SOLUTION stage.

With the implicit scheme, the factorization and solution stages require L_B, L_C and A₁₂, instead of L_B, L_C and W. Because A₁₂ is much sparser than W, equations (123), (124) and (125) are preferred instead of (116) and (120) and in the evaluation of the product W^T W. Now the full advantage of a block partitioning, with an implicit handling of W becomes clearer: we do not need to store the Cholesky factor L explicitly, we only need its diagonal components L_B and L_C plus the off-diagonal block of the original matrix, i.e. A₁₂.

It can be shown that the combination of block partitioning and asymmetric factorization produces not only storage savings, but also reduces the amount of arithmetic operations.

The process can be extended from a (2x2) block system to a more general one of, say, (rxr) blocks; in that case George and Liu (1981), section 6.2.2. showed that the off-diagonal submatrices of the Cholesky factor L can be computed implicitly as:

 $L_{ij} = A_{ij} L_{jj}^{-T} = (L_{jj}^{-1} A_{ji})^{T}$ i, j=1,2, ..., r that is to say, we only need to store the diagonal block components of L (i.e. L_{jj}) and the off-diagonal blocks of the original matrix (i.e. A_{ij} or A_{ji}), which is the same conclusion obtained in the case of a (2x2) block partitioning.

Obviously the greater the number of blocks (r), the greater the efficiency of this scheme.

Having illustrated with a small example the convenience of using the block partitioning approach, the remaining problem is to find some way of producing "good" block partitioning.

First of all, we have to define what we mean by a good partitioning. There is a close link between block partitioning and the concept of quotient graph, which was introduced previously, where the idea of "supernode" and, by extension, that of "supergraph" were introduced. In order to formalise that relationship, we have to return to these graph theory concepts.

Recall that a quotient graph is nothing more than a graph formed by "supernodes", which are simply a subset of nodes "glued" together; formally, if the original linear system has a graph defined by G = (X, E), and if we have a partitioning of X, say P, such that:

$$P = \{ Y_1, Y_2, \dots, Y_r \}$$
(126)

a quotient graph of G with respect to P, denoted as G/P, is the graph (P, E^*) where:

{ Y_i , Y_j } $\in E^*$ if and only if Adj(Y_i) $\cap Y_j \neq \phi$ (127)

Loosely speaking, equation (127) simply says that a link $\{Y_i, Y_j\}$ exists if and only if two supernodes Y_i and Y_j are adjacent.

A particular kind of connected graph in which there are no cycles (loops), is called a tree, say T = (X, E) and, as a result, every pair of nodes is connected by a unique path. Of course, we can also have "supertrees", or quotient graphs without loops. A rooted tree is a tree where every single node emanates from a root node, determining an ancestor-descendent structure. If we label (or number) the nodes of a rooted tree in such a way that every node is numbered <u>before</u> its ancestor, we get a tree which is known as a <u>monotonely ordered tree</u>. Fig. A.21 is an example of a monotonely ordered tree rooted at node 8.



Fig. A.21. An example of a monotonely ordered tree (rooted at node 8).

The matrix corresponding to the monotonely ordered tree shown in Fig. A.21 is as follows:



Fig. A.22. Matrix corresponding to the monotonely ordered tree shown in Fig. A.21.

As can be easily seen from Fig. A.22, the matrix of a monotonely ordered tree has a very interesting property indeed: it does not suffer fill-in during the factorization. Because of the labelling system used, where each node is numbered before its we are generating a matrix structure which ancestor, shows a series of arrow shaped arrangements pointing into the diagonal shaded "arrows" in Fig. A.22). (see

In other words, the numbering system adopted allows all the nodes to cluster themselves optimally, i.e. minimizing the <u>profile</u> of the matrix and thus reducing the amount of fill-in to the very minimum (i.e. no fill-in at all).

The results of the previous paragraph look quite interesting, but unfortunately in our water distribution networks we do not get trees (in general); we do get trees in sewage and drainage networks. Because loops are almost always found in water supply networks, we have to adapt this tree-oriented labelling system to a looped network (or graph) and the way of doing it is via the use of supernodes; we shall define the supernodes in such a way that the resulting quotient graph is a monotonely ordered
quotient tree ("supertree"). Thus the origin of the name of these methods has become clear.

To do so, we use again the concept of level structure, more precisely a rooted level structure.

First, note that we are looking for a tree with as many branches as possible; for example, if we use the matrix whose tree is shown in Fig. A.23, its level structure rooted at node "a" is that given in (128).



Fig. A.23. Network example for quotient tree methods.

$$L_{0} = \{ a \}$$

$$L_{1} = \{ b, f \}$$

$$L_{2} = \{ c, g, k \}$$

$$L_{3} = \{ d, h, 1, m \}$$

$$L_{4} = \{ e, i, n, s \}$$

$$L_{5} = \{ j, o, t \}$$

$$L_{6} = \{ p, q \}$$

$$L_{7} = \{ r \}$$
(128)

in this case, somebody would like to use the level structure itself as a partitioning, e.g.:

 $Y_{1} = \{ a \}$ $Y_{2} = \{ b, f \}$ $Y_{3} = \{ c, g, k \}$ $Y_{4} = \{ d, h, 1, m \}$ $Y_{5} = \{ e, i, n, s \}$ $Y_{6} = \{ j, o, t \}$ $Y_{7} = \{ p, q \}$ $Y_{8} = \{ r \}$ (129)

with the corresponding quotient tree shown in Fig. A.24. This quotient tree (Fig. A.24) generates a block-tridiagonal matrix, since each supernode is connected with only two others in a chained structure. The pattern of the generated block tridiagonal matrix is shown in Fig. A.25.

We can easily see that in this case no fill-in will occur, which is in fact due to the dense characteristics of the non-zero blocks of the matrix. The approach is clearly non-optimal, since if we compare the level structure (for example at level 5) with Fig. A.23, nodes "o" and "t" are quite far from each other and a different partitioning criterion is needed. An alternative is presented in Fig. A.26 where, for the same network of Fig. A.23, the partition is not the level structure itself, but a refinement has been used; in this case, we no longer glue nodes like "o" and "t", because they are in completely different branches. The new partitioning is produced by the following heuristic algorithm:

Let B_j be the graph defined by:

$$B_{j} = G (\bigcup_{i=j}^{l} L_{i})$$
(130)

which is referred to as the <u>section graph</u>. This graph is simply the union of all the nodes below the j-th level structure, including the nodes in L_j. The j-th level is refined in subsets Y such that:



.

Fig. A.24. Quotient tree corresponding to the tree of Fig. A.23.

*	*	*																	
*	*	*	*	*	*														Г
*	*	*	*	*	*					İ –				İ -					
	*	*	*	*	*	*	*	*	*								1		Γ
	*	*	*	*	*	*	*	*	*					l					
	*	*	*	*	*	*	*	*	*					[
			*	*	*	*	*	*	*	*	*	*	*						Γ
			*	*	*	*	*	*	*	*	*	*	*	ļ					
			*	*	*	*	*	*	*	*	*	*	*						
			*	*	*	*	*	*	*	*	*	*	*						
						*	*	*	*	*	*	*	*	*	*	*	Γ		Γ
						*	*	*	*	*	*	*	*	*	*	*			
						*	*	*	*	*	*	*	*	*	*	*			
						*	*	*	*	*	*	*	*	*	*	*			
										*	*	*	*	*	*	*	*	*	Γ
										*	*	*	*	*	*	*	*	*	
										*	*	*	*	*	*	*	*	*	
														*	*	*	*	*	5
														*	_*	*	*	*	
																	*	*	[;

rpqjoteinsdhlmcgkbfa

Fig. A.25. Matrix corresponding to the quotient tree of Fig. A.24.



a) Level structure
 b) Refined quotient tree
 Fig. A.26. Level structure rooted at node "a" and its corresponding quotient tree.

{ Y = L_i \cap C, G(C) results in a connected component of B_i } (131)

In other words, at a certain level "j", the nodes in that level are kept together (forming a unique supernode), if and only if the corresponding section graph spanned by these nodes is connected; otherwise, the nodes are re-grouped into separate supernodes, like in the case of nodes d, h, l and m at level L₃.

We can see immediately the advantage of this new partitioning approach, in the block tridiagonal matrix (Fig. A.25), we had 174 non-zeros and now in the refined quotient tree rooted at node "a" we have only 112 non-zeros. This means obviously less storage and less arithmetic operations.

As we saw before in the case of the quotient minimum degree algorithm, the level structures are dependent on the root node. In the level structure used in Fig. A.26 we chose "a" as the root

node just by chance, but we have to ask ourselves what is the best choice (if any) for such a root node. Since a partitioning with a maximum number of members will produce a matrix with the minimum number of non-zeros, it has been proposed that the concept of pseudo-peripheral nodes should be used to generate a root [George and Liu (1981), chapter 6]. Following this approach, if we generate the level structure rooted at node "t", and the corresponding refined partitioning of those levels using (131), see Fig. A.28, we get the matrix shown in Fig. A.29, which has only 108 non-zeros.

	r	- 1					T		1			1									· • •
r	×	ĸ	*			_										ł				.	
р	>	k	*		*	*															
q				*	*	*		_													
j			*	*	*	*	*	*								-					
ο			*	*	*	*	*	*							_						
е					*	*	*	*	*	*								-			
i					*	*	*	*	*	*											
d			-				*	*	*	*						*	*	*			
h							*	*	*	*						*	*	*			
t			-					_			*	*									
S								-			*	*		*	*						
n													*	*	*						
1												*	*	*	*	*	*	*			
m												*	*	*	*	*	*	*			
c								-	*	*				*	*	*	*	*	*	¥	
g									*	*]		*	*	*	*	*	*	*	
k		1							*	*				*	*	*	*	*	*	*	
Ъ																*	*	*	*	*	*
f		_	_													*	*	*	*	*	*
a		1						i											*	*	*
	L	1					l		<u> </u>							і м					

rpqjoeidhtsnlmcgkbfa

Fig. A.27. Matrix corresponding to refined quotient tree of Fig. A.26. b)



a) Level structure

b) Refined quotient tree

Fig. A.28. Level structure rooted at node "t" and its corresponding quotient tree.

*	*																		_
*	*			*	*	*													
	T	*		*	*	*								_					
			*	*	*	*													
	*	*	*	*	*	*	*	*											
	*	*	*	*	*	*	*	*											
	*	*	*	*	*	*	*	*				1							
				*	*	*	*	*	*	*		-							
				*	*	*	*	*	*	*				l					
							*	*	*	*	*	*							
							*	*	*	*	*	*							
	Γ								*	*	*	*	*	*					
									*	*	*	*	*	*					
											*	*	*	*	*	*			
											*	*	*	*	*	*			
	\square					-							*	*	*	*	*		
											l		*	*	*	*	*		
															*	*	*	*	
	П												-				*	*	*
													-	-				*	*

rqpedjocibhagflknmst

Fig. A.29. Matrix corresponding to refined quotient tree of Fig. A.28. b)

In our water network problems, we believe that this method is dependent on the shape of the networks under study; clearly in regular squared or round-shaped networks we obtain block tridiagonal matrices, whereas in branched networks we tend to obtain the sort of arrow shaped matrices we obtained at the beginning of this section (Fig. A.22).

The usefulness of some of these ideas for sewage and drainage networks looks very attractive and should be explored.

A.3.7. One-way dissection methods.

These methods were originally developed for solving linear systems arising in finite element applications. The background is again the quotient trees, and the one-way dissection algorithms are just another way of producing monotonely ordered supertree partitionings, as an alternative to the partitionings produced by the application of (131).

In essence, all dissection methods seek the division of the graph (and the associated matrix) into different members, in such a manner that the dissected graph produces a matrix with some desirable properties. Of course, these properties are related to low fill-in and the minimum number of arithmetic operations.

The idea of one-way dissection is to split up the network graph, by identifying some subset of nodes of the graph, which will be referred as "dissectors", and numbering each separated member sequentially, following a pre-specified direction. In

A-104

Fig. A.30, $\sigma = 2$ dissectors were used and 2 $\sigma + 1$ members were obtained (each dissector being a member itself).



Fig. A.30. Rectangular grid partitioned with 2 dissectors.

Each non-dissector member is numbered following an horizontal left-to right and top-down sequence (see arrows in Fig. A.30), while the dissectors are numbered after all the non-dissectors members have been numbered, following a top-down left-to-right sequence (see Fig. A.30). The same result can be obtained with a right-to-left, down-top strategy.

As a result of this particular numbering system the reordered matrix has a very special form. Fig. A.31 shows an example of a 40 nodes graph, renumbered according to this method, and Fig. A.32 presents the corresponding matrix structure induced by this numbering.

From Fig. A.32 we can see that the labelling system introduced by the one-dissection scheme produces $(2 \sigma + 1)$ blocks ; $(\sigma+1)$ blocks of size (ms x ms) and σ blocks of size (m x m). All the possible fill in the Cholesky factors is concentrated in the lower part of the matrices (see shaded areas in Fig. A.32).



Fig. A.31. Example of a 40 nodes rectangular shaped graph, partitioned using one-way dissection.

The selection of σ presents some practical problems, since it has to be an integer; non-dissectors members of different size can be defined, in order to cope with some cases when 1 and the chosen σ do not produce a constant $\delta = (1-\sigma)/(\sigma+1)$. A sensible selection of σ allows us to get an optimum matrix structure, which either minimizes fill-in or arithmetic operations.

In fact, because the storage requirements can be expressed as a function of σ (or δ), it is possible to explicitly find the σ^* which minimizes the storage. The same thing can be done to minimize the amount of work in the Cholesky factorization and in its solution stage. Because storage and amount of work are conflicting objectives, we either choose one the other, thus "tuning" the algorithm for one of them, which usually is the storage. According to the testing carried out by George and Liu (1981), their <u>one-way dissection implementation was found to be the least demanding on storage</u> algorithm, when compared with the envelope method, the refined quotient tree method and the nested dissection method.

A-106



Fig. A.32. Matrix structure corresponding to graph of Fig. A.31.

We implemented the subroutines published by George and Liu (1981), with the gradient method for the solution of water distribution networks, and we can confirm the results of George and Liu. Furthermore, we extended the comparison to include the MA27 package of the Harwell Subroutine Library and we also found that <u>one-way dissection was still the best of the direct methods</u>. from the storage viewpoint.

A.3.8. Nested dissection methods.

Nested dissection is a further development of one-way dissection and it can be explained as follows.

When applying one-way dissection to the matrix whose graph is Fig. A.31 , we only considered vertical dissectors; actually, there is nothing to prevent us splitting the graph even further by using horizontal dissectors as well, as shown in Fig. A.33.

The result of applying nested dissection to the graph shown in Fig. A.33, as far as the corresponding matrix is concerned, is shown in Fig. A.34. On comparing the structures generated by one-way dissection and nested dissection (see Figs. A.32 and A.34), we can find that one of the main effects of nested dissection is to reduce the size of the diagonal submatrices, by increasing its number and its density; this is highly desirable, since we do have to store these diagonal submatrices.

We have used the nested dissection subroutines published by George and Liu and obtained the <u>fastest execution time</u> of all the direct methods we have reviewed so far, though in storage the algorithm is not as good as the one-way dissection.

According to George and Liu (1981) and George (1981), their algorithm is not as efficient as other implementations, but it is simpler and suitable for tri-dimensional meshes (in 3-D finite



Fig. A.33. Example of a 40 nodes rectangular-shaped graph, partitioned using nested dissection.



Fig. A.34. Matrix structure corresponding to graph of Fig. A.33.

A-109

element problems, for example), whereas some other implementations are restricted to 2-D problems (planar meshes).

Nested dissection seems to be better for problems where the network has a squared or rounded shape, while one-way dissection is more appealing for problems generating a strongly rectangularshaped graph.

George (1981) warned about the reliability of this routines: "there are no results guaranteeing the quality of its outputs"; in our testing with water distribution networks, we have found full agreement between the results obtained with this method and others.

A.3.9. Frontal and multifrontal methods.

The frontal method derives from the search for efficient methods for solving the huge linear systems generated in the application of the finite element method. Although the ideas behind the method had been used for a long time, the name and its formal presentation are due to Irons (1970).

The frontal method aims primarily at the reduction of storage, since the amount of operations is basically the same as in other Gaussian elimination-based algorithms, though the order in which these operations are carried out is different.

The foundations of the frontal method are strongly tied up with the way in which finite element matrices are produced. In the finite element method (FEM) a problem is discretized over its

continuous domain, aiming at evaluating the unknowns only at certain points (nodes) within the domain, rather than at all its (infinite) points. The discretization usually generates a mesh of nodes, which are joined by edges, forming basic geometric figures (typically triangles and/or rectangles). Each one of the elementary figures constitutes a "finite element" or simply "element".

The global matrix (or "stiffness" matrix, as a reminiscent from the structural analysis problems where FEM was originally developed) is actually assembled from each one of the elementary matrices relating the value of the unknown variable at the finite element nodes with the boundary conditions and with the corresponding physically meaningful parameters. This is usually formalised as:

$$A = \Sigma B (k)$$
(132)

 $\underline{b} = \Sigma \underline{c} (k) \tag{133}$

where:

and

A : global "stiffness" matrix.

- B(k): is the matrix coming from the k-th finite element, with its order equal to the number of nodes in the k-th element.
- $\underline{c}^{(k)}$: right hand side vector corresponding to the k-th finite element.
- b : right hand side vector of the global system.

The key observation here is that the finite element definition, i.e. its shape, size and nodes, provides the adequate basis for producing a desirable structure for the Gaussian elimination

process. In all the previous methods (like quotient tree and dissection algorithms), a great deal of effort was spent in producing the adequate orderings. In FEM, this structure comes up "naturally", from the finite element specification. In fact, the elimination order in FEM is determined by the sequence in which the global matrix A is assembled.

Another basic observation leading to the development of the frontal method, deals with the fact that we do not have to wait until the global matrix is fully assembled (and indeed <u>stored</u>) to start the elimination process. In fact all we need for eliminating a row and column, say "1", is to wait until its coefficients in the global matrix have been completely added; at that point, the elimination of the variable "1" can be carried out, and because we no longer need it, the corresponding row can be stored out of core until the substitution stage requires it.

The following example, from Livesley (1983), illustrates the elimination sequence. Fig. A.35 shows the element definition and ordering, while Fig. A.36 shows the corresponding elimination sequence.



Fig. A.35. Finite element definition and ordering for a frontal solution scheme, from Livesley (1983).



Fig. A.36. The elimination sequence in a frontal solution scheme from Livesley (1983).

The "front" advances across the mesh following the assembly sequence, which defines the elimination order. The frontal nodes (or variables) lie in the border between the already assembled elements and the un-assembled ones. The rows and columns corresponding to the frontal nodes are usually kept in a full (dense) matrix known as the "frontal matrix".

A-113

From Fig. A.36 it is clear that a node cannot be eliminated until the contributions from all its related elements have been assembled, like node 5, which has to wait until element D is assembled.

From the previous example it is also clear that the front-width (i.e. maximum number of frontal nodes) depends on the order in which the finite elements are assembled, which is reflected by its numbering. This leads to the conclusion that an efficient frontal method implementation has to include some way of minimizing the front width. This can be done using some of the ideas described previously for ordering sparse matrices.

In the previous example, the sequence for assembling the global matrix A ,using the notation of equation (132), is defined by: $A = [[[[B(A)] + B(B)] + B(C)] + B(D)] + B(E)] + \dots (134)$ where the assembling takes place from the inner brackets outwards. This sequence corresponds with the elimination sequence shown in Fig. A.36.

Equation (134) can indeed be assembled in many different ways, for example:

 $A = \{ [B(A) + B(B)] + [B(C) + B(D)] \} + \{ [B(E)] + ... \}$ (135)

where, again, the assembling is carried out following the outwards direction of bracketing. This implies that we can actually assemble A in such a way as to produce more than one "front", each one with its corresponding frontal matrix, which can be processed and kept aside temporarily. This is an extension of the frontal method, leading to the "multifrontal" method,

where a sensible selection of the fronts allows extra savings in terms of number of operations.

The complexity of the computer implementation of the frontal methods is considerable. We have used the routine MA27 of the Harwell Library, which consists of 15 subroutines with 2920 Fortran records (cards). As an additional feature, MA27 allows the solution of indefinite sparse symmetric linear systems, because some pivoting is performed within the frontal matrix. The results of the comparison between the multifrontal method implemented in MA27 and other methods, for the solution of water distribution networks, is presented in Chapter Five.

Reid (1981), describes the main features of frontal methods, while Duff and Reid (1982_a, 1982_b) describe the main characteristics of the MA27 routine. Duff, Erisman and Reid (1986, Chapter 10) give an up to date description of frontal and multifrontal methods for finite element problems.

It has to be emphasised that, even though frontal and multifrontal methods have been developed for and from FEM applications, their use is open to non-FEM problems.

A.4. Review of iterative methods for the solution of linear systems of equations.

A.4,1. Introduction.

When solving a linear system of equations there are a number of reasons why, sometimes, an iterative (approximate) approach can be most appropriate:

a) For example, if we are solving a non-linear system of equations via an iterative scheme, which builds up a sequence of linear systems, there is no point in getting an "exact" solution in the first iterations of the non-linear algorithm, since we know beforehand that this is not the true solution of our original non-linear system, and all we want is an approximate In those cases, considerable effort can be spared with one. an approximate/iterative solution of the first successive linear systems, provided that the iterative solution is cheaper than a direct solution.

b) In engineering practice, real problems usually imply linear systems of the order of hundreds or thousands and, in these cases a direct solution can become very expensive from the amount of work and storage point of views. This has been the main argument in favour of iterative methods for quite a long time, but nowadays the development of fast sparse direct algorithms is challenging this argument.

c) The availability of fast and accurate algorithms for the iterative solution of linear systems of equations, on the one hand, and the widespread use of cheaper microcomputers in reallife problems, on the other hand, makes the use of iterative methods for the efficient solution of linear systems very attractive indeed.

Thus, in the past 10-20 years, direct methods have been associated with big computers, with massive storage capacity and iterative methods have been linked to mini and microcomputers.

```
A-116
```

Though this situation is changing very rapidly, it is still considered like a kind of rule of thumb by a number of people.

In essence, when storage is a critical constraint, an iterative method seems to be the right choice for the solution of linear systems, at least to start with.

We shall review the most widely used iterative methods, in an increasing degree of complexity and sophistication. Starting with some traditional methods (like Jacobi, Gauss-Seidel, Successive over relaxation, etc.), a unified view of all these methods is presented, since all of them can be studied as particular cases of a single general method. Finally, we shall review a conceptually different approach to solving linear systems, based on optimisation methods (conjugate gradients). We shall not include iterative methods like "alternating directions" and "block iterative methods", which are applied especially to solve linear systems arising from differential equations (see Ortega and Poole (1981) and Stoer and Bulirsh (1980) for details on these methods).

Basically, all iterative methods aim at producing a sequence of solution vectors such as:

 $\underline{x}^{(o)}, \underline{x}^{(1)}, \underline{x}^{(2)}, \ldots, \underline{x}^{(k)}$

which converge towards the desired solution. The different methods have particular ways of producing this sequence.

A.4.2. Jacobi's method (method of simultaneous displacements).

Let us consider our model problem:

Solve A x = b (136)

that is to say:

a _{n1}	x ₁	+	a _{n2}	x ₂	+	a _{n3}	x ₃	+	••••	+	a _{nn}	×n	=	b _n		
:			:			:					:			:		
^a 31	\mathbf{x}_1	+	a ₃₂	×2	+	a 33	х ₃	÷	••••	+	a _{3n}	×n	Ξ	Ъ _З	(137	')
^a 21	\mathbf{x}_1	+	^a 22	x ₂	+	^a 23	x ₃	+	••••	+	^a 2n	×n	=	b ₂		
a ₁₁	\mathbf{x}_1	+	^a 12	×2	+	a 13	x3	+	• • • •	+	a _{1n}	×n	=	^b 1		

If we assume that $a_{ii} \neq 0 \forall i=1,2,..,n$, then we can apply the same idea of the simple iteration method (in one variable) to an n-dimensional space, and calculate χ , in an iterative fashion, by reordering equation (137) as:

 $a_{11} x_1 = b_1 - (a_{12} x_2 + a_{13} x_3 + \dots + a_{1n} x_n)$ $a_{22} x_2 = b_2 - (a_{21} x_1 + a_{23} x_3 + \dots + a_{2n} x_n)$ $a_{33} x_3 = b_3 - (a_{31} x_1 + a_{32} x_2 + \dots + a_{3n} x_n)$ $\vdots \qquad \vdots \qquad \vdots \qquad \vdots$ $a_{nn} x_n = b_n - (a_{n1} x_2 + a_{n2} x_2 + \dots + a_{nn-1} x_{n-1})$ (138)

from which the values of the unknowns can be expressed in a compact form as:

$$x_{i}^{(k+1)} = \frac{1}{----} (b_{i} - \sum_{\substack{j \neq i \\ j \neq i}} a_{ij} x_{j}^{(k)}) \quad \forall i=1,..,n \quad (139)$$

where (k+1) and (k) refer to the iteration index.

Then, Jacobi's algorithm can be expressed as: Step 1. Starting with k=0 and an initial solution $\underline{x}(0)$, which can be zero or a better solution, if available.

- Step 2. Compute $\underline{x}^{(k+1)}$ using equation (139)
- Step 3. Compute $E = || \underline{x}^{(k+1)} \underline{x}^{(k)} ||$ or

$$E = \frac{|| x^{(k+1)} - x^{(k)}||}{||x^{(k)}||}$$

or any other index that can be useful as a detention (convergence) criterion.

and, if $E > \epsilon$: a given small termination value, then go to step 4, otherwise stop; $\chi^{(k+1)}$ is the result for that level of accuracy.

Step 4. If (k+1) > M : a prescribed maximum number of iterations, then stop, because an accurate solution has not been found in a reasonable number of iterations. Otherwise [(k+1) ≤ M], increase k to k+1 and go to step 2 again.

The convergence conditions for this algorithm are given in a subsequent section.

As it can be seen from equation (139), when we are computing the value of $x_i^{(k+1)}$, the right hand side uses $x_j^{(k)}$ from the previous iteration. A better approximation can be obtained using the last value (of the current iteration) for those x_j whose index j is smaller than i [i.e. for those values x_j which have been already calculated in the (k+1) iteration]. This idea leads to the next method: Gauss-Seidel method.

A.4.3. Gauss-Seidel method (method of successive displacements).

This is the extension of Jacobi's method using the updated values of x_j in the computation of the following x_i for i > j, then equation (139) is replaced by:

$$x_{i}^{(k+1)} = \frac{1}{---} \begin{bmatrix} b_{i} - \sum_{i j \neq i} x_{j}^{(k+1)} - \sum_{i j \neq i} x_{j}^{(k)} \end{bmatrix} \quad \forall i=1,..,n \quad (140)$$

and the algorithm described for Jacobi's method remains the same, replacing equation (139) by (140).

No extra effort is needed for the Gauss-Seidel method, in comparison with Jacobi's method, since all we have to do is to consider the last updated values of the unknowns, rather than those of the previous iteration. As a result, there is no reason for using Jacobi's method instead of Gauss-Seidel.

A.4.4. Successive over (or under) relaxation method.

It has been found that further improvements in the rate of convergence can be achieved by varying the magnitude of the step $[x^{(k+1)}-x^{(k)}]$ taken in the Gauss-Seidel method.

In fact, if we define the new value found via equation (140) as an <u>intermediate</u> value, say x_i , then we can compute the final new value as:

$$x_i^{(k+1)} = x_i^{(k)} + w \{ \hat{x}_i - x_i^{(k)} \}$$
 (141)

"w" is called a relaxation parameter, more precisely, we have over-relaxation if w>1 and under-relaxation if w<1. Of course, in the case w = 1, we are simply applying the original Gauss-Seidel method. We shall discuss in the next sections how to choose "w" to get the maximum acceleration in the iterative procedure.

A.4.5. Relationship between the previous methods.

If we express our model problem A = b as:

[B + (A - B)] x = b (142)

where B is an arbitrary matrix, then:

B x + (A - B) x = b (143)

Using the same idea behind simple (Picard) iteration, we can solve (143) via the following iterative scheme:

$$B x^{(k+1)} + (A - B) x^{(k)} = b$$
(144)

where the superscript refers to the iteration counter.

Since we have the freedom to choose B, we can assume that we select a non-singular matrix B; thus, the unknowns in (144) can be computed as:

or

$$x^{(k+1)} = B^{-1} [b - (A - B) x^{(k)}]$$
(145)

$$x^{(k+1)} = B^{-1} b - B^{-1} A x^{(k)} + B^{-1} B x^{(k)}$$
(146)

then, reordering and considering $B^{-1} B = I$:

or

$$x^{(k+1)} = x^{(k)} - B^{-1} A x^{(k)} + B^{-1} b$$
 (147)

$$x^{(k+1)} = [I - B^{-1} A] x^{(k)} + B^{-1} b$$
 (148)

Equation (148) represents a common matrix formulation for all the iterative methods already seen.

Although interest in equation (148) is mainly theoretical, since for programming purposes we must consider the subindexed expressions (139), (140) and (141), the fact is that from equation (148) we can generate all the previous methods by choosing an adequate B matrix. The key factor for the successfulness of any method is that B has to be easily invertible, and that B must be as similar to A as possible. Equation (148) can be rearranged in a more general expression:

where:

$$x^{(k+1)} = M x^{(k)} + d$$
 (149)
 $M = [I - B^{-1} A]$
 $d = B^{-1} b$

We shall look at the iterative methods from the perspective of equation (149), especially the conditions for their convergence.

On decomposing matrix A as:

$$A = -E + D - F$$
 (150)

where: r

$$E = -\begin{bmatrix} 0 & & \\ a_{21} & 0 & & \\ \vdots & \vdots & \\ a_{n1} & a_{n2} & \dots & a_{nn-1} & 0 \end{bmatrix}$$
(151)

٦

i.e. a lower triangular matrix with the same corresponding coefficients of A.

$$D = \begin{bmatrix} a_{11} & & \\ & a_{22} & & \\ & & & a_{nn} \end{bmatrix}$$
(152)

i.e. a diagonal matrix with the same diagonal coefficients of A.

$$F = -\begin{bmatrix} 0 & a_{12} & \dots & a_{1n} \\ & & 0 & a_{n-1n} \\ & & & 0 \end{bmatrix}$$
(153)

i.e. an upper triangular matrix with the same corresponding coefficients of A.

Then, assuming $a_{ii}\neq 0$ $\forall i = 1,..,n$ (i.e. D non-singular):

a) The <u>Jacobi</u> method can be obtained from equation (148) by choosing B such that:

$$B = D$$
; (154)

in so doing:

$$M = I - B^{-1} A = I - D^{-1} (-E + D - F)$$

= I - D^{-1} (-E - F) - D^{-1} D
$$M = D^{-1} (E + F)$$
(155)

b) The <u>Gauss-Seidel</u> method can be obtained from equation (148) by choosing B as:

$$B = D - E$$
; (156)

in that case:

$$M = I - B^{-1} A = I - (D - E)^{-1}(D - E - F)$$

= I - [(D - E)^{-1}(D - E) - (D - E)^{-1}F]
= I - [I - (D - E)^{-1} F]
= (D - E)^{-1} F (157)

c) The <u>iterative improvement</u> method for the correction of a solution of the linear system, which has been previously obtained by a factorization method, can also be expressed as a member of this family of iterative methods.

Suppose we have decomposed the matrix of coefficients A in a LU fashion. If we acknowledge that, because of rounding errors during the factorization stage the equality A = LU does not hold, we have instead:

$$A \approx L U \tag{158}$$

Let us take:

then

$$B = L U$$
(159)
$$B^{-1} = (L U)^{-1} = U^{-1} L^{-1}$$

and

$$M = I - U^{-1} L^{-1} A$$
 (160)

d) The <u>successive over-relaxation</u> method can be obtained from equation (148) by choosing B such that:

B = (1/w) (D - w E) (161)

where now B and M depend on the relaxation factor (scalar) "w". Then

 $M = I - B^{-1} A = I - w(D - w E)^{-1}(-E + D - F)$

$$M = (D - wE)^{-1} [(D - wE) - w(- E + D - F)]$$

= (D - wE)^{-1} [(1-w)D + wF] (162)

which, for the particular case w = 1, gives the Gauss-Seidel method (equation 157).

A.4.6. Convergence conditions for the previous methods.

There is a very close relationship between the convergence behaviour of a linear system of equations and the way their eigenvalues are distributed. It is known (Duff, 1985), that slow convergence is associated with evenly distributed eigenvalues, while the presence of clustered eigenvalues implies fast convergence. So, the eigenvalues are the key to the study of the convergence of the iterative methods we are interested in.

Unfortunately, the process for obtaining the eigenvalues of a matrix is quite costly from the computational viewpoint, involving some iterative procedure. Nevertheless, a number of computer codes are available to compute the eigenvalues.

We shall summarise the results of applying some linear algebra concepts to the conditions for convergence of the iterative methods already seen. The details of most of these concepts, fully proven, can be found for example in Stoer and Bulirsh (1980).

We shall not involve ourselves in the mathematics, but we will introduce only the main concepts. A complete review of the computational procedures available to determine eigenvalues and eigenvectors can be found in Ruhe (1977).

A-124

The eigenvalue problem is usually defined in terms of finding a set of scalars λ_i for i=1,...,n such that:

$$A x = \chi B x$$
 $\forall x \neq Q$ (163)
where B is usually taken as the identity matrix, i.e.:

$$A x = \chi x \qquad \forall x \neq Q \qquad (164)$$

If A is a (nxn) matrix, there are "n" values satisfying this condition, some of them real or complex numbers. Nevertheless, if A is symmetric, its eigenvalues are all real; furthermore, if A is positive definite, its eigenvalues are all positive.

Let λ_1 , λ_2 , ..., λ_n be the eigenvalues of A, the matrix of coefficients of the linear system (where some of them can be repeated); then the spectral radius of the matrix A is defined as:

$$\ell(A) = \max\{|\chi_i|, i=1,...,n, \text{ where } \lambda_i \text{ is an eigenvalue of } A\}$$
 (165)

Because the spectral radius represents the scatter of the eigenvalues, it can be used to find some information about the convergence rate of an iterative method. In order to avoid the explicit computation of the eigenvalues, it is possible to relate the spectral radius and the norm of the matrix. On applying the norm operator to equation (164), it is easy to prove that:

$$|\lambda| \leq ||A|| \tag{166}$$

which is valid for <u>any norm</u>. Then, from the definition of spectral radius, we have:

$$\ell(A) \leq || A || \tag{167}$$

This means that any norm of A is an upper limit for its spectral radius. This result is useful, since it avoids the

explicit computation of the spectral radius and eigenvalues, when trying to determine the convergence properties of a method.

The basic (necessary and sufficient) condition for the convergence of any iterative algorithm, represented by its general formulation:

$$x^{(k+1)} = [I - B^{-1} A] x^{(k)} + B^{-1} b$$
 (148)

or

$$x^{(k+1)} = M x^{(k)} + d$$
 (149)

is :

$$((I-B^{-1}A) < 1)$$
 (150)

or simply:

f(M) < 1 (151)

On applying the result of equation (167), we find that a <u>sufficient</u> condition for convergence of the algorithm represented by (148) is simply:

$$|| I - B^{-1} A || < 1$$
 (152)

In general, the more distant from 1 the spectral radius or the norm of M is, the faster the convergence of the method. Equation (152) provides a quick means for getting a first check if a method is convergent or not, since this only implies the evaluation of any norm, say for example the maximum norm.

The results of applying these concepts to the iterative methods seen so far, can be summarised as follows:

- a) <u>Jacobi's</u> method: in general converges when the matrix A is strictly diagonally dominant ($|a_{ii}| > |a_{ik}|$, $\forall k \neq i$ and $\forall i$).
- b) <u>Gauss-Seidel</u> method: always converges when either

* A is strictly diagonally dominant, or

* A is a positive definite matrix.

In addition, if $J=D^{-1}E+D^{-1}F$ is non-negative, the Gauss-Seidel method converges faster than Jacobi's method.

c) <u>Successive over-relaxation</u>: this method always converges when A is positive definite or when the relaxation parameter "w" is such that:

0 < w < 2

In practice, the relaxation parameter must be found through a trial and error procedure, within its range of convergence.

The successive over-relaxation method converges faster than Gauss-Seidel, provided that $D^{-1}E$ and $D^{-1}F$ are non-negative matrices.

In general, although Jacobi and Gauss-Seidel methods may converge, their rate of convergence can be quite slow and, due to rounding errors, the whole process can be spoiled.

A.4.7. Standard conjugate gradient method for the solution of linear systems of equations. (Hestenes and Stiefel, 1952).

When A is a (nxn) symmetric and positive definite matrix, the solution of the model problem Ax = b can be understood as a minimization problem.

The main idea behind this method comes from the fact that the quadratic function:

$$Q = (1/2) \underline{z}^{T} A \underline{z} - \underline{z}^{T} \underline{b}$$
(153)

has a minimum which is exactly the solution of the model linear problem A z = b. Then, instead of solving the original problem,

we solve:

minimize (1/2)
$$\underline{\mathbf{x}}^{\mathrm{T}} \mathbf{A} \underline{\mathbf{x}} - \underline{\mathbf{x}}^{\mathrm{T}} \underline{\mathbf{b}}$$
 (154)

The method was developed by Hestenes and Stiefel (1952), and is basically an unconstrained minimization algorithm, where the successive direction vectors, are conjugate with respect to the matrix A and parallel to the error vectors (difference between true and computed \underline{x} 's).

We shall introduce the algorithm first and review its mathematical background afterwards.

We use $\langle \underline{x}, \underline{y} \rangle$ to denote the inner (scalar) product of two vectors:

$$\langle \underline{\mathbf{x}}, \underline{\mathbf{y}} \rangle \equiv \underline{\mathbf{x}}^{\mathrm{T}} \underline{\mathbf{y}} = \underline{\mathbf{y}}^{\mathrm{T}} \underline{\mathbf{x}}$$

The standard conjugate gradient method consists in:

- 1. a) Determine an initial solution vector $\underline{x}^{(0)}$.
 - b) Compute the residual vector: $\underline{r}(o) = \underline{b} A \underline{x}(o)$
 - c) Initialise the counter k=0 and take the direction p(0)=r(0)

2. a) Compute the auxiliary vector $\mu(k) = A \rho(k)$.

- b) Compute the scalar $\alpha_k = \langle \underline{r}(k), \underline{p}(k) \rangle / \langle \underline{p}(k), \underline{\mu}(k) \rangle$
- c) Compute the new vector $\underline{x}^{(k+1)} = \underline{x}^{(k)} + \alpha_k \underline{p}^{(k)}$
- d) Compute the new residual $\underline{r}^{(k+1)} = \underline{r}^{(k)} \alpha_k \underline{\mu}^{(k)}$
- e) Check for convergence: is $||r^{(k+1)}|| \le \epsilon$?

```
If yes, stop.
```

If not, continue

€ is the user specified accuracy.

3. a) Compute the scalar $\beta_k = \langle \underline{r}^{(k+1)}, \underline{\mu}^{(k)} \rangle / \langle \underline{p}^{(k)}, \underline{\mu}^{(k)} \rangle$

- b) Compute $\underline{p}^{(k+1)} = \underline{r}^{(k+1)} \beta_k \underline{p}^{(k)}$
- c) Increase k by 1 and go to step 2.a)

Thus, the algorithms builds up a sequence of vectors:

 $\underline{x}^{(0)}, \underline{x}^{(1)}, \underline{x}^{(2)}, \ldots, \underline{x}^{(n)}$

which converges to the solution of the linear system A $\underline{x} = \underline{b}$

As we shall see in the next paragraphs, it is implicit in the algorithm that consecutive directions $p^{(k)}$ and $p^{(k+1)}$ are Aconjugate. This property guarantees that the solution can be reached in at most "n" steps, where "n" is the order of Α, provided that the calculations are carried out in exact arithmetic. Then, at least in theory, this method fits into the category of the direct methods but, in practice, since rounding errors are present in all calculations, its behaviour corresponds to an iterative method.

Now we shall review the mathematical concepts behind the algorithm:

In the first step of the algorithm an initial solution vector $\underline{x}^{(0)}$ has to be provided; usually a null vector is the easiest starting point, but if a better approximation is available, the algorithm will find the final solution quicker.

Then, the algorithm moves the solution vector from $\underline{x}^{(k)}$ to $\underline{x}^{(k+1)}$ following the direction $\underline{p}^{(k)}$. The algorithm finds the best step size in the direction $\underline{p}^{(k)}$ via a one-dimensional optimisation of the step length; this is done by minimizing the objective function in that direction:

minimize $(1/2)(\underline{x}^{(k)}+\alpha\underline{p}^{(k)})^{T}A(\underline{x}^{(k)}+\alpha\underline{p}^{(k)}) - (\underline{x}^{(k)}+\alpha\underline{p}^{(k)})^{T}\underline{b}$ (155) α On differentiating equation (155) to determine the optimum value of α , say α_k , we get:

$$A (\underline{x}^{(k)} + \alpha_{k} \underline{p}^{(k)}) - \underline{b} = \underline{0}$$

then,

or

$$A \underline{x}^{(k)} + \alpha_k A \underline{p}^{(k)} = \underline{b}$$

$$\alpha_k A \underline{p}^{(k)} = \underline{b} - A \underline{x}^{(k)} = \underline{r}^{(k)}$$
(156)

where $\underline{r}^{(k)}$ is the residual vector at the k-th iteration.

Premultiplying equation (156) by $\underline{p}^{(k)}$ (using the inner product) and recalling that $\langle \underline{x}, \underline{y} \rangle = \langle \underline{y}, \underline{x} \rangle$, we can determine α_k :

$$\alpha_k \langle \underline{p}^{(k)}, A \underline{p}^{(k)} \rangle = \langle \underline{r}^{(k)}, \underline{p}^{(k)} \rangle$$

finally,

$$\alpha_{\mathbf{k}} = \frac{\langle \mathbf{r}^{(\mathbf{k})}, \mathbf{p}^{(\mathbf{k})} \rangle}{\langle \mathbf{p}^{(\mathbf{k})}, \mathbf{A}, \mathbf{p}^{(\mathbf{k})} \rangle}$$
(157)

which is the scalar we compute in step 2.b) of the algorithm.

In the computer implementation of the algorithm we do not need explicitly the matrix A, only its product by the vector $\underline{p}^{(k)}$, which is kept in an auxiliary vector, say $\underline{\mu}^{(k)}$.

Now we are able to update the vector of unknowns:

$$\underline{\mathbf{x}}^{(k+1)} = \underline{\mathbf{x}}^{(k)} + \alpha_{k} \underline{\mathbf{p}}^{(k)}$$
(158)

and the new residual:

$$\underline{\mathbf{r}}^{(k+1)} = \underline{\mathbf{b}} - \mathbf{A} \, \underline{\mathbf{x}}^{(k+1)} \tag{159}$$

which can be simplified in order to avoid the product A $\underline{x}^{(k+1)}$ because, on introducing $\underline{x}^{(k+1)}$ from equation (158) into equation (159), we get:

$$\underline{r}^{(k+1)} = \underline{b} - A \{ \underline{x}^{(k)} + \alpha_k \underline{p}^{(k)} \}$$

or

$$r^{(k+1)} = b - A x^{(k)} - \alpha_k A p^{(k)}$$

then

or

$$\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - \alpha_k \mathbf{A} \mathbf{p}^{(k)}$$
(160)

This means that we can overwrite the residual vector $\underline{\mathbf{r}}^{(k)}$ with its new value $\underline{\mathbf{r}}^{(k+1)}$ by just subtracting $\alpha_k \ \underline{\mu}^{(k)}$ [since we stored the product A $\underline{\mathbf{p}}^{(k)}$ in the auxiliary vector $\underline{\mu}^{(k)}$]. In so doing, we have avoided the explicit computation of the product Ax^(k+1) in equation (159). Equation (160) corresponds to step 2.d) of the algorithm.

If convergence is not achieved in step 2.e), we need to generate a new direction for the improvement of the solution; because we are looking for a new direction A-conjugate with the previous one; this means that we need:

$$\langle \underline{p}^{(k+1)}, A \underline{p}^{(k)} \rangle = 0$$
 (161)

We shall find the new direction as a linear combination of the residual and the previous direction, i.e.:

$$p^{(k+1)} = r^{(k+1)} - \beta_k p^{(k)}$$

Imposing the A-conjugate condition [equation (161)], we get:

$$\langle \underline{\mathbf{r}}^{(k+1)} - \beta_{k} \underline{\mathbf{p}}^{(k)}, A \underline{\mathbf{p}}^{(k)} \rangle = 0$$

$$\langle \underline{\mathbf{r}}^{(k+1)}, A \underline{\mathbf{p}}^{(k)} \rangle - \beta_{k} \langle \underline{\mathbf{p}}^{(k)}, A \underline{\mathbf{p}}^{(k)} \rangle = 0$$

which gives us the value of β_k we are looking for:

$$\beta_{k} = \frac{\langle \underline{r}^{(k+1)}, A \underline{p}^{(k)} \rangle}{\langle \underline{p}^{(k)}, A \underline{p}^{(k)} \rangle}$$
(162)

where again we can use the stored auxiliary vector $\mu(k)$ instead of A p(k).

It is at this stage that rounding errors affect the finite

termination property of the conjugate gradient method, since with exact arithmetic the vectors $\underline{p}(k)$ and $\underline{p}(k+1)$ are exactly A-conjugate (or A-orthogonal), but in the presence of rounding errors orthogonality is lost.

The amount of work needed at each step of the algorithm is mainly due to the product A $\underline{p}^{(k)}$, which makes the method very attractive from the computational point of view; the rest of the computations are mainly scalar operations. Only 4 vectors need to be stored: $\underline{x}^{(k)}$, $\underline{r}^{(k)}$ [which can overwrite the original right hand side vector <u>b</u>], $\underline{p}^{(k)}$ and $\underline{\mu}^{(k)} = A \underline{p}^{(k)}$.

These features have gained the conjugate gradient method a reputation for large sparse linear systems in reduced storage computers. In practice, the method is not recommended for dense (non-sparse) matrices.

The conjugate gradient method was developed in 1952, when it was originally regarded as a direct method, but it was abandoned shortly afterwards, due to convergence problems created by the lack of orthogonality. These problems and a comparison of different versions of the algorithm were analysed by Reid (1971), and interest in the method, now regarded as an iterative algorithm was renewed. Of the several possible versions of the conjugate gradient algorithm we have presented one of the most widely used ones.

Gambolatti and Perdon (1984) compared the conjugate gradient method with other iterative methods and offered a geometrical interpretation of the algorithm, which is summarised in Fig.

A.37. for 2 and 3-dimensional problems.

Parlett (1980) has shown the mathematical relationship between conjugate gradient methods and the Lanczos algorithm for solving eigenvalue problems, the conjugate gradient being a particular case of the Lanczos algorithm. We have used a Lanczos implementation for solving the linear systems in the early stages of the development of the gradient method for pipe distribution network analysis [Lewis (1977)]; as expected, for normallyconditioned problems, the convergence of the preconditioned conjugate gradient was faster than the Lanczos implementation by a factor of about two.

Nevertheless, the Lanczos algorithm is expected to behave better than the conjugate gradient for ill-conditioned problems. For more details on the Lanczos method see, for example, Scott (1981).

A.4.8, Preconditioning.

Any direct or iterative method for the solution of linear systems can be improved in its convergence rate via a technique known as "preconditioning".

In an earlier section (The effect of rounding errors in the solution of linear systems), we defined the condition number of a matrix as the norm product $||A|| \cdot ||A^{-1}||$, and we associated this number with the stability of the solution of the linear system. Matrices with low condition numbers (about 1) are said to be "well-conditioned" matrices, whereas those with larger condition

A-133


Fig. A.37. Geometric interpretation of the convergence of the conjugate gradient algorithm for n=2 and n=3, from Gambolatti and Perdon (1984), solving A h = b

numbers (say, 10⁴ or even greater) are considered as "illconditioned". We may get convergence problems with illconditioned matrices.

The main idea behind preconditioning is to transform an illconditioned problem into a well-conditioned one, via a suitable modification of the original system. Algebraically, this is done by a pre and/or post-multiplication of the original system of equations by some suitable matrices, which are referred to as preconditioning matrices (or simply preconditioners).

Clearly, the aim is to produce a new linear system, with a better behaviour from the rate of convergence viewpoint. A complete review of preconditioning methods, applied both to direct and iterative methods, can be found in Evans (1983).

We shall concentrate on the application of preconditioning to iterative methods.

The concept of preconditioning can be formalised as follows: instead of solving our original model problem A x = b, solve either

or

$$(MA) x = (Mb)$$
 (163)

 $(N A N^{T})(N^{-T} x) = (N b)$ (164)

where M and N are square preconditioning matrices.

The form of the preconditioning given by equation (163) seems simpler, but the main advantage of (164) is that, if A is symmetric, the new matrix of coefficients ($N A N^{T}$) remains symmetric, although in that case we need to determine a

transformed unknown variable, say $y = N^{-T} x$. Apart from our main objective, in the sense that the new system has to be better conditioned than the original one, we look for a new system which may be eventually easier to solve as well.

There are as many preconditioning methods as matrices M and N can be imagined. From an ideal point of view, the best choice for M in equation (163) would be $M = A^{-1}$, since then the solution for the vector x could be obtained in one single step. In general, with an adequate selection of M and N, the preconditioned problem will converge faster than the original one but, unfortunately, there is no simple rule to determine a general purpose preconditioning scheme.

Some possible preconditioning matrices are: a) Diagonal matrices:

For example M= diag(m_{11} , m_{22} , ..., m_{nn}) with $m_{ii} = 1/a_{ii}$. This kind of preconditioning matrix is more effective in highly diagonally dominant systems, since in this case (MA) is close to the identity matrix.

b) Complete factorizations of A:

Different kinds of factorizations of A can be used as preconditioners, since in general they are easier to invert than the original matrices. Thus, L U, L L^T, L D U, L D L^T, or Cholesky complete factorizations can be used as preconditioners [see section A.2.3.], because all of them produce easily invertible matrix factors.

A-136

c) Incomplete factorizations of A:

Since complete factorizations lead to loose the sparseness of the original system (because of fill-in), this kind of factorization does not seem adequate when dealing with large sparse linear systems; instead, an incomplete factorization is usually the correct answer to this problem. In an incomplete factorization we look for factors which have a similar sparseness pattern in relation to the original matrix. That is to say, we are now dealing with factorizations of the form:

or
$$A \approx L_s U_s$$

 $A \approx L_s D U_s$

where L_s and U_s correspond to sparse matrix factors, whose product approximates the original matrix A. Note that now we are just asking for an approximation to A.

In order to produce an incomplete factorization, there are at least two strategies:

i) Fixed pattern scheme: in this case the original matrix is decomposed in factors that have <u>exactly</u> the same sparsity pattern as A.

Kershaw (1978), used successfully the incomplete Cholesky factorization given by:

$$A = L D L^{T} + E$$
(165)

where A = symmetric and positive definite matrix.

- L = lower triangular matrix with the same sparsity as A.
- D = diagonal matrix
- E = error matrix containing all the non-zero elements which lie in the positions where A has zero values.

Then, the following approximation holds:

$$\mathbf{A} \approx \mathbf{L} \mathbf{D} \mathbf{L}^{\mathrm{T}}$$
(166)

Of course, in the case of symmetric positive definite systems, the same idea can be applied to the Cholesky "square root" factorization:

$$A = L L^{T} + E$$
(167)

with the approximation:

$$A \approx L L^{T}$$
(168)

The algorithm to perform this last factorization can be derived easily, from equation (168), expressed term by term:

$$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{22} & \dots & a_{21} \\ \vdots \\ \text{symmetric} & a_{nn} \end{bmatrix} \approx \begin{bmatrix} l_{11} & & & \\ l_{21} & l_{22} & & \\ \vdots & \vdots \\ l_{n1} & l_{n2} & \dots & l_{nn} \end{bmatrix} \cdot \begin{bmatrix} l_{11} & l_{21} & \dots & l_{n1} \\ \vdots & l_{22} & \dots & l_{n2} \\ \vdots & & & l_{nn} \end{bmatrix}$$
(169)

The first row of L^T is determined multiplying the first row of L by each column of L^T and comparing element by element with the coefficients of A:

$$a_{11} = l_{11} l_{11} \longrightarrow l_{11} = \sqrt{a_{11}}$$
(170)

$$a_{12} = l_{11} l_{21} \longrightarrow l_{21} = a_{12}/l_{11}$$

$$a_{11} = l_{11} l_{11} \longrightarrow l_{11} = a_{11}/l_{11} \quad \forall i=2,3,..,n$$
(171)
Equations (170) and (171) give us the first row of the matrix

Equations (170) and (171) give us the first row of the matrix L^{T} .

The diagonal elements of L^T are determined by multiplying row "i" of L and column "i" of L^T :

 $a_{ii} = l_{i1} l_{i1} + l_{i2} l_{i2} + \dots + l_{ii} l_{ii}$ then

$$l_{ii} = \int_{a_{ii}}^{i - 1} \frac{\sum_{k=1}^{i-1} l_{ik}^2}{k = 1}$$
(172)

The i-th row of L^T is determined by multiplying the i-th row of L by each column of L^T :

$$a_{ij} = l_{i1} l_{j1} + l_{i2} l_{j2} + \dots + l_{i(i-1)} l_{j(i-1)}$$

then
 $i-1$

$$l_{ji} = (a_{ij} - \sum_{k=1}^{r-1} l_{ik} l_{jk}) / l_{ii} \qquad \forall j = (i+1), (i+2), ..., n \qquad (173)$$

4

The main advantage of this decomposition lies in the fact that, when a sparse data structure is used to store A, then we only need to store the numerical part of L^{T} , since the same data structure (row and column indices) describe both A and L^{T} . Of course, the symmetry of A implies that only half of the matrix is actually stored, i.e. A is also stored like a triangular matrix, and the symmetry is handled implicitly. As a result, great economy in storage is achieved, although an auxiliary (integer) vector is required to scan L^{T} by columns, when a row-wise data structure is being used.

A couple of examples are shown to illustrate some features of this factorization scheme:

Example a) *******

-

$$A = \begin{bmatrix} 8.3 & -3.2 & 0.0 & -2.7 \\ -3.2 & 13.7 & -5.6 & 0.0 \\ 0.0 & -5.6 & 10.5 & -2.2 \\ -2.7 & 0.0 & -2.2 & 9.2 \end{bmatrix}$$

A-139

$$L^{T} = \begin{bmatrix} 2.881 & -1.111 & 0.0 & -0.937 \\ 3.531 & -1.586 & 0.0 \\ 2.826 & -0.779 \\ 2.778 \end{bmatrix}$$
$$L L^{T} = \begin{bmatrix} 8.299 & -3.2 & 0.0 & -2.699 \\ -3.2 & 13.7 & -5.6 & 1.041 \\ 0.0 & -5.6 & 10.499 & -2.200 \\ -2.699 & 1.041 & -2.200 & 9.199 \end{bmatrix}$$

We can see that the only apparent discrepancy seems to be the element in the second row and fourth column of L L^{T} which is 1.041 instead of 0.0.

Example b) *******

A	=	16.0 0.0 -2.0 -1.0	0.0 9.0 -2.0 -1.0 0.0	0.0 -2.0 10.0 1.0 0.0	-2.0 -1.0 1.0 20.0 1.0	$ \begin{bmatrix} -1.0 \\ 0.0 \\ 0.0 \\ 1.0 \\ 25.0 \end{bmatrix} $
LT	=	4.0	0.0 3.0	0.0 -0.67 3.091	-0.5 -0.33 0.252 4.424	-0.25 0.0 0.0 0.198 4.990
L L ^T	=	16.0 0.0 -2.0 -1.0	0.0 9.0 -1.999 -0.999 0.0	0.0 -1.999 9.999 0.999 0.0	-2.0 -0.999 0.999 19.999 1.0	$ \begin{array}{c} -1.0\\ 0.0\\ 0.0\\ 1.0\\ 25.0 \end{array} $

Now the agreement between A and L L^T seems to be better than in example a).

These two examples "suggest" some of the most relevant features of this sort of factorization: they behave better in the case of large matrices, and also when a strong diagonally dominant matrix is factorized.

<u>ii) Relative magnitude scheme</u>: now the criterion for keeping the factors with a similar sparseness changes and, instead of their relative position with respect to the original non-zeros, some terms of the factorization are "rejected" or "thrown away" because their magnitude is less than a pre-established amount.

Meijerink and Van der Vorst (1977) proposed a decomposition along this line, for symmetric M-matrices (i.e. when all the non-diagonal coefficients are less or equal to zero).

Munksgaard (1980) proposed an incomplete L D L^{T} factorization of this type, for symmetric positive definite matrices (without extra requirements on its structure).

Ajiz and Jennings (1984), following the same idea, presented a $L L^T$ factorization, though some adjustments are recommended to the diagonal values in order to maintain positive definiteness.

The advantage of this proposition is its flexibility, because it allows a wide set of possible incomplete factorizations to be obtained, ranging from a complete Cholesky decomposition up to one where all off-diagonal terms are rejected. The disadvantage is that it is not possible to predict beforehand the fill-in created, and so it is quite possible to run out of storage during the factorization, unless a dense factor L has been dimensioned, in which case it is obvious that a complete factorization should have been used. Less evident as a disadvantage is the fact that this approach needs a more complex data structure, because the original data structure does not represent the incomplete factor.

A-141

Mainly because of storage constraints we prefer the incomplete factorization proposed by Kershaw (1978), though in some cases, if storage is less restrictive, the approach of Ajiz and Jennings (1984) can produce faster execution times.

Chen and Tewarson (1986) have shown how incomplete factorization can be applied to solve some banded shifted coefficient matrices.

A.4.9. Preconditioned (modified) conjugate gradient method for the solution of linear systems.

Although there are as many preconditioned conjugate gradient methods as possible preconditioning matrices, we have followed a scheme proposed by Ajiz and Jennings (1984), but with an incomplete Cholesky factorization based on the "square root" decomposition [see equations (170) to (173)], instead of the incomplete factorization proposed by Ajiz and Jennings (based on the magnitude of the non-zero values). This is exclusively in order to keep the storage needs under control.

The subroutine MA31 of the Harwell Library implements a preconditioned conjugate gradient algorithm, based on the incomplete factorization proposed by Munksgaard (1980).

Given the incomplete Cholesky factorization:

1

$$A \approx L L^T$$
, (174)

the "standard" conjugate gradient algorithm is applied to the following preconditioned linear system:

$$(L^{-1} A L^{-T}) (L^{T} \underline{x}) = (L^{-1} \underline{b})$$
 (175)

which is equivalent to solving the system:

$$(L^{-1} A L^{-T}) y = \underline{b}*$$
 (176)

where:

($L^{-1} \wedge L^{-T}$) is the new (preconditioned) symmetric matrix $\chi = L^T \chi$ is the new transformed unknown $b^* = L^{-1} b$ is the new transformed right hand side vector.

Thus, the preconditioned conjugate gradient method is, loosely speaking, no more than the standard conjugate gradient applied over a transformed (preconditioned) linear system. When the solution of the preconditioned system has been obtained for the vector $\underline{\mathbf{y}}$, a back substitution is needed to compute the original unknown $\underline{\mathbf{x}}$.

From the computational point of view, we do not need to store the transformed matrix $(L^{-1} \wedge L^{-T})$ explicitly, because all we need to do with this matrix is to multiply it by the conjugate direction vector $\mathbf{p}^{(k)}$, to form the auxiliary vector:

$$\mu(k) = (L^{-1} A L^{-T}) \rho(k)$$

which is similar to that computed in step 2.a of the standard conjugate gradient method. As a result, we can build up the auxiliary vector $\underline{\mu}$ in a three stage process:

a) Back substitution: solve $L^T \underline{v}(k) = \underline{p}(k)$ for $\underline{v}(k)$ b) Premultiplication: compute $\underline{w}(k) = A \underline{v}(k)$ c) Forward substitution: solve $L \underline{\mu}(k) = \underline{w}(k)$ for $\underline{\mu}(k)$ (177)

Then, the preconditioned conjugate gradient algorithm can be carried out in the following steps (compare with the standard method in section 6.7):

- 0. Compute the corrected right hand side vector $\underline{b} * = L^{-1} \underline{b}$ which is a simple premultiplication.
- 1. a) Determine an initial solution vector $\mathbf{y}^{(o)}$.
 - $(\underline{y}^{(o)}=\underline{0} \text{ or a better approximation , if available,}$ see Note 1 at the end of the algorithm).
 - b) Compute the residual vector: $\underline{r}^{(0)} = \underline{b} * (L-1AL-T) \underline{y}^{(0)}$ $(\underline{r}^{(0)} = \underline{b} * \text{ if we started with } \underline{y}^{(0)} = \underline{0}, \text{ see Note 1})$
 - c) Initialise the counter k=0 and take the direction p(o)=r(o)
- 2. a) Compute the auxiliary vector $\underline{\mu}(k)$ with the three stage process given in equation (177)
 - b) Compute the scalar $\alpha_k = \langle \underline{r}(k), \underline{r}(k) \rangle / \langle \underline{p}(k), \underline{\mu}(k) \rangle$ (see Note 2).
 - c) Compute the new vector $y^{(k+1)} = y^{(k)} + \alpha_k p^{(k)}$
 - d) Compute the new residual $\underline{r}^{(k+1)} = \underline{r}^{(k)} \alpha_k \underline{\mu}^{(k)}$
 - e) Check for convergence: is ||r^(k+1)|| ≤ ∈ ?
 If yes, go to step 4.
 If not, continue
 ∈ is the user specified accuracy.
- 3. a) Compute the scalar $\beta_k = \langle \underline{r}^{(k+1)}, \underline{r}^{(k+1)} \rangle / \langle \underline{r}^{(k)}, \underline{r}^{(k)} \rangle$ (see Note 3).
 - b) Compute $\underline{p}^{(k+1)} = \underline{r}^{(k+1)} + \beta_k \underline{p}^{(k)}$ (see Note 4).
 - c) Increase k by 1 and go to step 2.a)

4. Solve $L^T \underline{x} = \underline{y}$ for \underline{x} via a back substitution process and stop, because \underline{x} is the solution of the original system (see Note 5).

Notes:

=====

1. When solving iteratively a system of non-linear equations, via a sequence of linear systems, it appears logical to keep the solution of the k-th linear system (say $\chi^{(k)}$), and use it as an initial solution vector for the (k+1)-th linear system.

We have applied this approach in the context of the gradient method for the water distribution analysis problem but, when solving systems of more than 165 unknowns we reached a break even point, balancing the time saved starting with the vector χ from the previous (non-linear) iteration with the additional work involved in computing the new matrix (L⁻¹ A L^{-T}) via the three stage process (177). In this case, we recommend a start with the initial solution $\chi^{(0)}=Q$, as in step 1.a), i.e. $r^{(0)}=b^*$

Gambolatti (1980) suggested that further improvement in the convergence rate of the preconditioned conjugate gradient can be obtained via the pre-computation of the initial solution by the Newton iterative algorithm, but we have not experimented with his recommendation.

2. There is an alternative formulation of the algorithm, for step 2.b), which is equivalent to that used in the standard conjugate gradient algorithm:

Compute $\alpha_k = \langle \underline{r}^{(k)}, \underline{p}^{(k)} \rangle / \langle \underline{p}^{(k)}, \underline{\mu}^{(k)} \rangle$

3. Here there is an alternative formulation for step 3.a), equivalent to that used in the standard method:

Compute $\beta_k = \langle \underline{r}^{(k+1)}, \underline{\mu}^{(k)} \rangle / \langle \underline{p}^{(k)}, \underline{\mu}^{(k)} \rangle$

4. If the alternative β_k of Note 3 has been used, the new direction vector p(k+1) in step 3.b) must be computed with:

$$p(k+1) = r(k+1) - \beta_{k} p(k)$$

i.e. only a change of sign.

5. Obviously, because the preconditioned conjugate gradient solved the transformed problem for the auxiliary unknown $\underline{\mathbf{y}} = \mathbf{L}^T \underline{\mathbf{x}}$, we need to find the value of the original unknown $\underline{\mathbf{x}}$ via a back substitution process, which is done in step 4.

We shall not repeat here the reasons for the computation of the optimal parameters α^k and β^k , because they remain the same as those explained in the case of the standard conjugate gradient.

Gambolatti and Perdon (1984) interpreted the conjugate gradient method in geometric terms. The standard conjugate gradient involves the minimization of a functional which represents a hyper-ellipsoid in an n-dimensional space; the more "spherical" this hyper-ellipsoid the faster is the convergence of the algorithm. Thus, the preconditioning can be interpreted as an effort to turn the hyper-ellipsoid into a sphere.

The development of conjugate gradient methods for the solution linear systems of equations is still a very active research of Sartoretto (1984) compared the preconditioned conjugate field. gradient with other iterative techniques used in microcomputers, using problems coming from the finite element analysis of demonstrating thestructures. Some tests goodness of preconditioning are presented in Jackson and Robinson (1985). Samuelsson, Wiberg and Bernspâng (1986) compared some preconditioned conjugate methods with direct methods for structural mechanics problems.

APPENDIX B

DERIVATION OF THE KRIGING ESTIMATOR EQUATIONS

B.1. Introduction.

The estimation technique known as Kriging owes its name to D. G. Krige, a South African geologist who, in the early fifties, approached the problem of estimating the ore content in a deposit, based on limited sampling.

G. Matheron and co-workers in the School of Mines of Paris at Fontainebleu (France), in the sixties and seventies, extended the ideas of Krige, introducing the terms "geostatistics" and "regionalized variables". Geostatistics refers to the application of statistic concepts to natural phenomena, whereas regionalized variable stands for variables describing phenomena which are spatially and, eventually, temporally distributed, with an underlying structure, the structural characteristics being a consequence of the genesis of the phenomena.

So far, geostatistics has been applied to a vast variety of fields, some of them well apart from the original mining engineering: the petroleum industry, geophysics, geotechnical engineering, forestry, geochemistry, pollution problems, etc. [Huijbregts (1973)]. In water resources engineering, geostatistics has been used to study the spatial variation of rainfall, to interpolate and to determine average values of rainfall depths within a basin [Delhomme, (1978)], to the

identification of parameters in groundwater flow [Hoeksema and Kitanidis (1983), (1984) and (1985)] and other problems .

We shall review the statistical and mathematical background of Kriging, looking at its application to the problem of piezometric head estimation in the context of water supply distribution networks. Most of the concepts presented here can be found, with much more detail, in the publications by de Marsily (1986) and Journel and Huijbregts (1978), but they are reviewed here for completeness.

The geostatistical process is usually carried out in a twostage procedure:

a) <u>Representation of the model structure</u>: which accounts for the construction of an experimental variogram, based on the field data, and the determination of the best analytical model corresponding to that experimental variogram. This is also known as the "statistical inference" or the "structural analysis" stage.

b) <u>The estimation of the unmeasured states (Kriging)</u>: based on the previously determined structure, Kriging provides the best linear unbiased estimates of the unknown states.

Both stages involve completely different activities and require different techniques. In the first stage (statistical inference), we are dealing with an analytical process, trying to specify in statistical terms (through the variogram) the structure of the phenomenon under study; this is clearly an activity which demands

a great deal of judgement, even though the computer can be used in an interactive way to ease the calculations involved. The second stage (Kriging itself) is the numerical part of the process and is usually carried out with the help of the computer, since it basically involves the assembly and solution of linear systems of equations. The success of the overall process relies heavily on the ability of the modeller to incorporate the . structure of the phenomenon into the variogram model and, to a lesser extent, into the Kriging itself.

We shall review each stage with some more detail, but first some notation needs to be introduced.

Because we are dealing with a water supply distribution network, which is physically spread under the ground level, the geometrical space over which the problem is defined is a twodimensional space. It can be argued that the alternative is a one-dimensional space, since we are dealing with pipes, which are basically one-dimensional objects, but we believe the twodimensional representation is more appropriate because we are (in general) concerned with looped networks, in which case the heads and flows are not dependent on the properties of one isolated pipe, but depend on the interaction of the different elements of the network, which do not necessarily need to be located in a one-dimensional arrangement. In fact, looped water distribution networks are essentially two-dimensional arrangements of pipes and nodes, whereas open (branched) water distribution networks are easily handled in a one-dimensional fashion.

Nevertheless, we restrict ourselves to the study of the piezometric heads over a discrete set of points in the twodimensional space, these points are the nodes of the water distribution network. Occasionally, we may also consider other points, located in the pipes joining those nodes, but in either case we are dealing with a discrete phenomenon rather than a continuous one. As a result, head measurements are available at some points of the space and head estimates (based on those measurements) are required at unmeasured points.

Let $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)^T$ be the two-dimensional vector representing the location, projected into an horizontal plane, of the nodes and any point within the network.

Since we are interested in the description of the piezometric head plane (surface actually) we introduce the state variable H, representing the piezometric head at any point of the network, thus:

 $(\mathbf{x}, \mathbf{H}) = (\mathbf{x}_1, \mathbf{x}_2, \mathbf{H})^{\mathrm{T}}$

represents a point in the piezometric plane at coordinates (x_1, x_2) with a piezometric head of value H.

Due to the fact that we are modelling the phenomenon as a stochastic process, instead of a deterministic one, the previous definition is better represented as a <u>random function</u>, say Z(x,H), whereas $Z(x_0,H)$ denotes a <u>random variable</u> at the particular location x_0 and $Z(x,H_1)$ denotes a particular <u>realisation</u> of the random function.

In order to simplify the notation of the random variables, sometimes we may drop the state variable (H) or, even simpler, we may use sub-indices to represent the value of a random variable at a particular location; thus, the following expressions are all equivalent:

 $Z(\underline{x}_i, H) \iff Z(\underline{x}_i) \iff Z_i$

B.2. Statistical inference.

In the geostatistical sense, the fact that $Z(\underline{x}_1,H)$ is a "regionalized variable" (i.e. it has an underlying structure or "spatial distribution"), means that there is a correlation between $Z(\underline{x}_1,H)$ and another variable $Z(\underline{x}_2,H)$, separated by a distance $\underline{d} = \underline{x}_1 - \underline{x}_2$.

The variability of the regionalized variables is characterised by a statistical function referred to as the semi-variogram, which is defined as:

Semi-variogram: $\Gamma(\underline{x}_1, \underline{x}_2) \equiv \frac{1}{2} \mathbb{E} \left[\{Z(\underline{x}_1, H) - Z(\underline{x}_2, H)\}^2 \right]$ (1) where E denotes mathematical expectation.

Some authors [de Marsily (1986), for example], refer to the semi-variogram as simply the variogram, the only difference being the scalar factor ½ applied to the expectation in the previous definition. We shall try to keep the present notation (semivariogram) since it does not lead to any confusion, although we may refer to it occasionally as simply the variogram.

Hence, in theory, the semi-variogram depends on \underline{x}_1 and \underline{x}_2 , and its estimation requires the availability of a great number of realisations: [$Z(\underline{x}_1, H_1)$, $Z(\underline{x}_2, H_1)$], [$Z(\underline{x}_1, H_2)$, $Z(\underline{x}_2, H_2)$],, [$Z(\underline{x}_1, H_r)$, $Z(\underline{x}_2, H_r)$] of the regionalized variable.

In practice, these numerous realisations are not available and we only have <u>one</u> realisation at different points of the space; physically, this corresponds to a set of head measurements at the different nodes in the network.

order to be able to estimate the semi-variogram from the In measured data some simplifying assumptions are made. Firstly, the random process (or random function) is assumed to be ergodic, which allows us to use the information from the unique realisation available to determine the properties of the ensemble all possible realisations. Secondly, assumptions with of different degree of stringency are made with respect to the behaviour of the statistics of different order (mean, variance and higher order moments) associated with the random process. The most stringent assumption used is that of <u>stationarity</u>, which means that all the statistics, including higher order moments, are assumed to be constant. Because for many practical problems hypothesis is too demanding and far from the real this characteristics of the random process under study, a more relaxed known as weak stationarity (or second-order assumption stationarity) is often used, whereby only the mean and variance assumed to be invariant. Further relaxation in theare stationarity of the random process occurs when the weak stationarity is not required for the random process itself, but

for its <u>increments</u>, which is the <u>intrinsic hypothesis</u>. We shall follow the usual procedure of deriving the Kriging equations for different stationarity conditions, starting with the most stringent one (stationarity) and we shall relax it successively to second-order stationarity and second-order stationarity for the increments (intrinsic hypothesis). However, the question concerning which conditions are applicable to a particular problem has to be answered from the data available and from a physical understanding of the process being studied.

On assuming that the intrinsic hypothesis holds (the weakest assumption so far), it is possible to estimate the semi-variogram from the measurement data. Indeed, the intrinsic hypothesis formally establishes that:

$$E[Z(x_1, H) - Z(x_2, H)] = 0$$
 (2)

$$\operatorname{var}[Z(\underline{x}_1, H) - Z(\underline{x}_2, H)] = 2 \Gamma(\underline{d})$$
(3)

where $\underline{\mathbf{d}} = \underline{\mathbf{x}}_1 - \underline{\mathbf{x}}_2$.

With these hypotheses, the semi-variogram $\Gamma(\underline{d})$ can be computed, because from (3):

$$\Gamma(\underline{d}) = \frac{1}{2} \operatorname{var}[Z(\underline{x}_1, H) - Z(\underline{x}_2, H)]$$
(4)

which can be expanded to: $\Gamma(\underline{d}) = \frac{1}{2} E[\{Z(\underline{x}_1, H) - Z(\underline{x}_2, H)\}^2] - \{E[Z(\underline{x}_1, H) - Z(\underline{x}_2, H)]\}^2$ (5)

The last term of (5) vanishes because of the stationarity assumption on the mean of the increments [equation (2)], hence:

$$\Gamma(\underline{d}) = \frac{1}{2} \mathbb{E}[\{Z(\underline{x}_1, H) - Z(\underline{x}_2, H)\}^2]$$
(6)

Equation (6) allows us to estimate the semi-variogram from the data, as the arithmetic mean of the squared differences between two pairs of measurements $[Z(x_1, H), Z(x_2, H)]$:

$$\Gamma^{*}(\underline{d}) = \frac{1}{2} \left\{ \begin{array}{c} 1 & N(\underline{d}) \\ ---- & \Sigma & [Z(\underline{x}_{1}, H) - Z(\underline{x}_{2}, H)]^{2} \right\}$$
(7)
$$N(\underline{d}) \quad i=1$$

where $N(\underline{d})$ represents the number of pairs of data at distance "d".

For computational purposes we define classes of distance between pairs of measurements, for example:

Class	Between distances (m)	
$\begin{array}{c} C_1\\ C_2\\ C_3\\ C_4\\ \vdots\end{array}$	(0 , 100] (100, 200] (200, 300] (300, 400] : :	

and, for each class, we simply average the squared differences: $\frac{1}{2}$ { [$Z(\underline{x}_1, H) - Z(\underline{x}_2, H)$]² }, for all \underline{x}_1 and \underline{x}_2 belonging to this particular class. The process can be carried out in a tabular format [de Marsily (1986)], as presented in Table B.1.

Table B.1. Computation of the experimental semi-variogram.

Class	c ₁	C ₂	Сз	• • •	$C_{\mathbf{k}}$	•••
Number of pairs: N(d)	N ₁	N2	N3		Nk	•••
Average distance: d	d_1	d2	dg	• • •	\mathtt{d}_k	• • •
Average of $\{1/(2N(d))\}[Z_{i}-Z_{j}]^{2}$	a1	a2	ag	• • •	ak	(*)
(*) $a_{k} = \frac{1}{2} \frac{N_{k}}{N_{k}} [Z_{i} - Z_{j}]^{2}$						

The estimated semi-variogram $\Gamma^*(\underline{d})$ is the resulting plot with the average distance as abscissa and the corresponding average of the squared differences as ordinate, as shown in Fig. B.1.

The estimated semi-variogram is then a broken line obtained by joining the successive points produced by the previous computation, i.e. we get as many points as distance classes.

Some details on the semi-variogram become relevant:

* Note that for a set of "n" measurements, ,we get n(n-1)/2pairs; each one has to be classified according to its distance (d) and allocated into the corresponding distance class (C_i). This is without considering the direction, in which case the number of pairs is reduced, since we only have to include the pairs belonging to a certain direction; this implies that a sectorization of the geometrical space is needed for computational purposes.



Fig. B.1. Estimated semi-variogram.

* Because the number of pairs decreases with the distance, it means that the uncertainty in the estimated semi-variogram increases with the distance. Because of that, sometimes it can be convenient to discard the pairs too distant from each other; this implies the use of a "moving neighbourhood", made up of all the points surrounding the point to be Kriged, which are contained within a circle of a pre-specified radius, thus reducing the uncertainty in the semi-variogram at longer distances. Also, some improvement in the estimation of the semi-variogram in the shorter distances can be obtained by defining the distance classes with a variable length, taking into consideration the number of pairs included.

* Some data points (outliers) can induce great errors in the estimated semi-variogram, and it might be necessary to eliminate some of them in order to stabilise the numerical computations.

The practicalities in the semi-variogram estimation are key factors in its robustness and representativeness, as a "summary" of the spatial correlation between the regionalized variables; Armstrong (1984) discusses a number of the common causes of problems in estimating the semi-variogram and also proposes the corresponding corrective measures.

Some relevant features of the semi-variogram need to be stressed:

a) <u>Continuity</u>: Due to the spatial continuity of the regionalized variable, the semi-variogram is also continuous. For a particular direction in the distance vector <u>d</u>, and because of

the increments vanish at a zero distance, the semi-variogram starts, in theory, at the origin $[\Gamma(0)=0]$ and it increases with the distance (absolute value of the distance). The way in which the semi-variogram increases with the distance is an indicator of the degree of spatial correlation of the regionalized variable.

Range and sill: we restrict ourselves to a particular b) direction in the distance vector, and if we carry on increasing distance (modulus), we may reach a point where thethe correlation no longer exists; the distance corresponding to this situation, say "r", is referred to as the "range" of the semivariogram, while its associated value, $\Gamma(r)$ say, is known as the "sill". Changing the direction in the distance vector normally implies that different ranges and sills can be obtained for the same regionalized variable; in the three-dimensional case, in mining or groundwater for example, we may have two completely different values for the range and sill, depending on whether we are moving in an horizontal or a vertical direction. In water supply distribution networks, we may expect different semivariograms in the case, for example, of long networks where the flow pattern in the longitudinal direction is very different with respect to the flow distribution in the transverse direction. The study of the effect of anisotropies in the regionalized variable on the semi-variogram is important in the statistical inference process, in order to obtain a robust variogram.

c) <u>Existence of the sill</u>: As we implied before, the presence of a range and sill in the semi-variogram <u>may</u> be possible. In fact, there are some cases where the semi-variogram does not reach a

sill for any value of the distance; this implies that the regionalized variable has an infinite variance, and only the introduction of additional assumptions (intrinsic hypothesis) can allow us to study the behaviour of such a phenomenon. We might expect the lack of a sill in the case of water supply distribution networks with a steep piezometric plane.

d) <u>Behaviour near the origin</u>: In theory, the value of the semivariogram at the origin (distance=0) is null, but this may not become clear from the measurement data. First of all, because we only compute the value of points of the semi-variogram centred at the average distance of each class, we may not get points close enough to the origin (depending on how we have chosen the class length). This may produce an apparent "jump" in the semivariogram near the origin, as shown in Fig. B.2., this is known as the "nugget effect".

The second practical issue leading to the nugget effect is related to the sampling procedure itself, and it explains the name of the effect; in the mining field, if a sample contains a nugget, i.e. a small piece of mineral in its natural state (gold is the typical case), the ore content will be very high, whereas a very close sample, taken just outside the vein of high mineral content, will show a very low ore content (if any at all). This may explain the erratic behaviour of the semi-variogram at the origin, in some cases. We do not expect this kind of behaviour in the semi-variogram of water supply distribution systems, and we assume that a continuous variogram, passing through the shall origin is the best representation of the phenomenon under study.



(a) The nugget effect





(b) The true semi-variogram

Fig. B.2. Nugget effect and corresponding true semi-variogram

Measurement errors can produce a small nugget effect, and we have to be aware of this in the event that the data show some discontinuity near the origin.

Once the experimental semi-variogram is available from the data, the next step is to fit an analytical expression to it, in order to set up the theoretical basis for the Kriging equations. However, not any analytical function can be used to represent the estimated semi-variogram, due to the main features of the semivariogram which have been reviewed in the previous paragraphs. In addition, for reasons which will become apparent later on, only positive definite functions can be accepted as suitable models. In practice, only a few analytical models are commonly used; Table B.2 summarises these common models.

All the previous basic models, except the monomial one, are models with a sill, although only the spherical and cubic models reach the sill in a finite distance, given by the parameter B_0 ; the exponential and Gaussian models only reach the sill asymptotically. Another difference between all the basic models is concerned with their behaviour near the origin, where different concavities are obtained in different models.

Model	Analytical expression for $\Gamma(d)$: $\Gamma(d) = A_0 + B_0$.[function(d, C ₀)]	Conditions
Monomial in d ^{CO}	B _o d ^{Co}	C _O <2 A _O =0
Cabonical	$B_0[1.5(d/C_0) - 0.5(d/C_0)^3]$	d <c<sub>o ; A_o=0</c<sub>
Spherical	B _o	$d > C_0$; $A_0 = 0$
Exponential	B _o [1 - e ^(-d/Co)]	A _o =0
Gaussian	B _o {1-e[-(d/Co) ²]}	A ₀ =0
0.1.	$\begin{array}{r} B_{o}[7(d/C_{o})-8.75(d/C_{o})^{3}+3.5(d/C_{o})^{5} - \\ - 0.75(d/C_{o})^{7}] \end{array}$	d <c<sub>o A_o=0</c<sub>
CUDIC	Bo	d>C _o ; A _o =0

Table B.2. Basic analytical models for the experimental semivariogram.

Having determined the experimental variogram on the one hand, and having a set of available models on the other hand, the aim is to find out which model (or a combination of them) fits the experimental semi-variogram best. This process can be greatly simplified by a computer program which can be run interactively by the user to produce the best fit model in an interactive way; eventually, if the computer program is designed to compute an index of the goodness of fit (i.e. squared sum of residuals or any other index), the process can be automated to give the best model, leaving to the user the decision to accept it, reject it or modify it. Huijbregts (1975) warns about the automatic fitting of variogram models without performing the physical (structural) interpretation of them, the main shortcoming of this being the loss of insight into the behaviour of the regionalized variable. An important feature of these programs is the ability to display (on the screen or in a printout) both graphically the experimental variogram and data and the fitted analytical model, since such a plot can be very helpful when accepting or rejecting the model.

For a more detailed exposition on statistical inference, in the geostatistical context, see Journel and Huijbregts (1978), Chapter III: Structural Analysis, where some FORTRAN program listings for this purpose are included. Also, the paper by Armstrong (1984) is relevant to this subject.

Perhaps the most important feature in the inference of the semi-variogram, is that everything we have said so far is valid only when the intrinsic hypothesis holds [which is expressed

through equations (2) and (3)]. If a drift or trend is detected (i.e.: $E[Z(\underline{x}+\underline{d})-Z(\underline{x})]=m(\underline{d})\neq 0$), then the semi-variogram cannot be estimated from the data and a more complex method is required, such as that known as "intrinsic random functions" [see Delfiner and Delhomme (1973) and de Marsily (1986)]. We shall review this case later on.

B.3. Estimation via Kriging.

In order to derive the Kriging estimation equations, we shall start from some ideal situation, assuming stationarity in the random function; then we shall see how a non stationary process can be dealt with and also how measurement errors can be handled within Kriging. Finally, the most general version of Kriging, known as "universal Kriging" will be outlined. In this case the mean is not constant and the drift is explicitly handled using the concept of an order-k intrinsic random function.

B.3.1. Kriging under stationary conditions.

A stochastic process is said to be second-order stationary or weakly stationary if its mean is a constant value and its covariance depends only on its lag; in our case, this means that

$$E [Z(\underline{x}, H)] = m$$
(8)

and

$$Cov(\underline{x}_1, \underline{x}_2) = E [\{Z(\underline{x}_1, H) - m\} \{Z(\underline{x}_2, H) - m\}] = C(\underline{d})$$
(9)

where $\underline{d} = \underline{x}_1 - \underline{x}_2$: distance (vector) between points \underline{x}_1 and \underline{x}_2 ; we are primarily concerned with the magnitude of this vector rather than its direction. In these conditions, we can expand the expectation in (9) to give:

$$C(d) = E[Z(x_1, H)Z(x_2, H)] - mE[Z(x_1, H) - mE[Z(x_2, H)] + m^2$$

then, on introducing the constant mean condition, we obtain:

$$C(\underline{d}) = E[Z(\underline{x}_1, H)Z(\underline{x}_2, H)] - m^2 - m^2 + m^2$$

which reduces to

$$C(\underline{d}) = E[Z(\underline{x}_1, H)Z(\underline{x}_2, H)] - m^2$$
 (10)

The covariance, under stationary conditions, is related to the semi-variogram [see de Marsily, (1986)] through the relationship:

$$\Gamma(\underline{d}) = C(\underline{Q}) - C(\underline{d}) \tag{11}$$

and the relationship can be seen graphically in Fig. B.3, where the semi-variogram ordinates are simply the differences between the horizontal line through C(O) and the covariance curve (shaded area).





Thus, because we already know how to estimate the semivariogram from the measurement data, we can assume that the covariance is actually known. Additionally, we can also assume that the mean "m" is also known and, under these conditions, we can introduce a change of variable, to obtain the following zero mean process (dropping the state H in the notation from now on):

$$Y(\underline{x}) = Z(\underline{x}) - m \tag{12}$$

i.e.: we are dealing with a zero mean process: E[Y(x)]=0.

Now we can proceed to compute an estimate (Y_0^*) of the unknown variable Y_0 at a point χ_0 , based on the available measurements at the other points Y_i :

$$Y_i = Y(x_i)$$
, $i = 1, ..., n$ (13)

The Kriging estimate Y_0^* , corresponding to an unmeasured point χ_0 , is then computed as a linear weighted combination of the measurements:

$$Y_0^* = Y^*(x_0) = \sum_{i=1}^n x_0^i Y_i$$
 (14)

where the λ_0^{i} are the unknown weights of the Kriging estimator, the superscript "i" refers to the measurement Y_i , whereas the sub-index "o" refers to the point being estimated \underline{x}_0 . We are dealing basically with a point estimator, where equation (14) is applied one point at a time, using the same set of measurements; we then get a different set of weights for each point to be estimated. Y_0^* is then an estimate of the true value of the variable Y_0 , which is not known. Kriging estimates an optimum Y_0^* in the sense that it minimizes the square of the errors: $Y_0^* - Y_0$; because Y_0 is not known, we have to formalise the optimality condition stating that the <u>expected value</u> of the square of the errors is the quantity to be minimized:

min E [
$$(Y_0^* - Y_0)^2$$
] (15)

In other words, we are minimizing the <u>variance</u> of the errors of the estimation; the expectation is taken at a fixed point, one point \underline{x}_0 at a time, over all the possible realisations of the variable $Y(\underline{x}_0, H)$.

On developing expression (15) and introducing the Kriging estimator from equation (14), we get: $E [(Y_0^* - Y_0)^2] = E [(\sum_{i} \chi_0^i Y_i - Y_0)^2]$ $E[(Y_0^* - Y_0)^2] = E[(\sum_{i} \chi_0^i Y_i)(\sum_{j} \chi_0^j Y_j)] - 2 E[(\sum_{i} \chi_0^i Y_i Y_0)] + i + E[Y_0^2]$ $E [(Y_0^* - Y_0)^2] = \sum_{i} \sum_{j} \chi_0^i \chi_0^j E[Y_i Y_j] - 2 \sum_{i} \chi_0^i E[Y_i Y_0] + i + E[Y_0^2]$ (16)

On reordering equation (10), we can compute the expected value of the product $Y_i Y_j$, assuming that the covariance and mean are known:

$$E[Y_{i} Y_{j}] = C(x_{i} - x_{j}) + m^{2}$$
(17)

but, because of the variable change (12), Y is a process with zero mean, hence:

$$E[Y_i \ Y_j] = C(\underline{x}_i - \underline{x}_j) \tag{18}$$

and, in particular, for $x_i = x_j$: $E[Y_0^2] = C(Q)$ (19)

which is the dispersion variance of Y or var(Y).

Introducing the results of equations (18) and (19) into equation (16), we obtain:

$$E [(Y_{0}^{*} - Y_{0})^{2}] = \sum_{i j} \sum_{i j} \lambda_{0}^{i} \lambda_{0}^{j} C[x_{i} - x_{j}] - 2 \sum_{i j} \lambda_{0}^{i} C[x_{i} - x_{0}] + C(0)$$
(20)

which is a quadratic function of the weights λ_0^{i} . According to our optimality criterion [equation 15], this function ought to be minimized, the variables being the weights λ_0^{i} . The necessary conditions for a minimum of (20) are given by the following set of equations in the unknowns λ_0^{i} :

$$\frac{\partial}{\partial x_0^{i}} = [(Y_0^* - Y_0)^2] = 0 \quad i=1, ..., n \quad (21)$$

i.e.

$$\frac{\partial}{\partial x_0^{i}} E[(Y_0^* - Y_0)^2] = 2 \sum_{j} x_0^{j} C[\underline{x}_1 - \underline{x}_j] - 2 C[\underline{x}_1 - \underline{x}_0] = 0$$

$$i = 1, \dots, n$$
(22)

hence:

$$\sum_{j=1}^{n} x_{0}^{j} C[x_{i}-x_{j}] = C[x_{i}-x_{0}] \qquad i = 1, ..., n$$
 (23)

A sufficient condition for a unique solution of (23) is that the matrix of covariances (C) has to be positive definite, and that is the reason why we wanted positive definite functions for the semi-variogram model.

Some important features of the Kriging estimator become apparent from equation (23), namely:

The weights χ_0^{j} do not depend explicitly on the value of the i) measurements Zi. The weights determined with (23) depend on the structure of the regionalized variable, expressed through the semi-variogram or covariance matrix C. This means that the same set of weights can be used with different sets of measurements for the same location, provided that the information brought up by the new measurements does not alter the structure of the problem (semi-variogram). This rather striking property of Kriging can be used advantageously in real-time estimation, thesemi-variogram being updated off-line, if necessary, or a bank of different time-dependent semi-variograms can be made available based on past measurements.

ii) The estimation considers both the distances between the point being estimated and the measurement points $[\underline{x}_i - \underline{x}_0]$ and the distances between the measurement points $[\underline{x}_i - \underline{x}_j]$. The former means that Kriging gives a higher weight to data closer to the point being estimated and the latter means that Kriging discriminates redundant information [see Delfiner and Delhomme (1973; Fig. 1) for an example of an estimator which does not discriminate redundant information].

iii) If we set $\underline{x}_0 = \underline{x}_j$, i.e.: \underline{x}_0 coincides with one measurement point \underline{x}_j , the solution of (23) gives $\lambda_0^{j=1}$ and $\lambda_0^{k=0} \forall k \neq j$, which produces $\underline{Y}_0^* = \underline{Y}(\underline{x}_j)$, i.e.: the estimate coincides with the measurement. This is why Kriging is said to be an <u>exact</u> <u>interpolator</u>, in other words: the estimates "replicate" exactly the measurements.

iv) Armstrong (1984) showed that, even in reduced order cases (4x4 systems), the Kriging linear system can be <u>ill-conditioned</u>. In that case, small changes in the variogram produce such changes in the matrix of coefficients, which may lead to significant changes in the weights (λ_0^{i}) . Armstrong recommended that the computation of the condition number of the Kriging matrix be incorporated as a standard step in the estimation process, so that some warning can be given to the user in the presence of an ill-conditioned system [see Appendix A for details of ill-conditioned linear systems of equations].

v) The right hand side of (23) is the only term dependent on the point to be estimated (\underline{x}_0); this means that for estimating Y_0^* for different locations, we only need to solve equation (23) for different right hand side vectors. Thus, efficient solutions of the linear system (23) can be obtained via Gaussian elimination with a triangular decomposition, which can be carried out only <u>once</u>, stored and used repeatedly for each different point (\underline{x}_0) to be estimated.

Appendix A includes a review of the theory and performance of the most widely used linear solvers; all of them are applicable

to solve the linear systems generated by Kriging, due to the symmetry and positive definiteness of the system.

Having computed the weights λ_0^i , we go back to equation (14) to compute the Kriging estimate Y_0^* corresponding to the particular location \underline{x}_0 , i.e.:

$$Y_{O}^{*} = Y^{*}(X_{O}) = \sum_{i=1}^{n} X_{O}^{i} Y_{i}$$

Now, the last problem to be solved is the determination of the estimation variance. Since we do not know Y_0 , the true value of the variable at \underline{x}_0 , we cannot compute the error $Y_0^*-Y_0$ directly, only its variance:

$$var(Y_0^* - Y_0) = E[(Y_0^* - Y_0)^2] - \{ E[Y_0^* - Y_0] \}^2$$
(24)

Introducing equation (14):

$$E[Y_{O}^{*} - Y_{O}] = E[\sum_{i} \chi_{O}^{i} Y_{i}] - E[Y_{O}]$$

or

$$E[Y_0^* - Y_0] = \sum_{i} \chi_0^i E[Y_i] - E[Y_0]$$

but, because of the change of variable (equation 12), in general $E[Y_j]=0 \quad \forall j$, then:

$$E[Y_0^* - Y_0] = 0$$
 (25)

As a result, equation (24) simplifies to:

$$var(Y_0^* - Y_0) = E[(Y_0^* - Y_0)^2]$$
(26)

Note that the right hand side of (26) is simply the quadratic function we had to minimize [equation (20)]. But we already know
the right hand side in equation (20) after determining the weights λ_0^{i} by solving the linear system (23); therefore:

$$\operatorname{var}(Y_0^* - Y_0) = \sum_{i j} \sum_{j \in \mathcal{X}_0} \lambda_0^j C(x_i - x_j) - 2 \sum_{i j} \lambda_0^i C(x_i - x_0) + C(0)$$

Introducing (23) into the above equation gives:

$$\operatorname{var}(Y_0^* - Y_0) = \sum_{i} \chi_0^{i} C(\chi_i - \chi_0) - 2 \sum_{i} \chi_0^{i} C(\chi_i - \chi_0) + C(0)$$

which, on simplifying, yields:

$$\operatorname{var}(Y_0^* - Y_0) = - \sum_{i} \chi_0^i C(\chi_i - \chi_0) + C(0)$$

and recalling that C(0) = var(Y):

$$var(Y_0^*-Y_0) = var(Y) - \sum_{i} \chi_0^{i} C(x_i-x_0)$$
 (27)

٦

In words, (27) means that the estimation variance is less than the dispersion variance of the random function; only when $\underline{x}_0 = \underline{x}_j$ is the variance of the estimation exactly equal to the measurement variance. The optimum combination of measurements with a given variance produces then, via Kriging, an estimate which, on average, has a minimum mean square error.

Finally, having solved the problem in terms of the transformed variable Y, we need to go back to the original variable Z:

$$Z = Y + m$$

Then,

$$Z_{O}^{*} = Z^{*}(X_{O}) = m + \sum_{i=1}^{n} X_{O}^{i} (Z_{i} - m)$$
 (28)

and the Kriging variance (σ_0^2) becomes:

$$\sigma_0^2 = \operatorname{var}(Z_0^* - Z_0) = \operatorname{var}(Z) - \sum_i \chi_0^i C(x_i - x_0)$$
 (29)

The second term on the right hand side of equation (29) is known as the <u>variance reduction</u>, i.e.: the variance reduction brought about by the optimal combination of measurements.

It must be stressed that, strictly speaking, the variance given by (27) and (29) is the estimation variance, or the variance of the <u>estimator</u>, rather than the variance of the <u>estimates</u>. In other words, it is a sort of "average" of the variance of the estimates.

B.3.2. Relaxing the second-order stationarity condition: the intrinsic hypothesis.

The second-order stationarity conditions used to derive the previous equations apply to stochastic processes where the mean of the random function is known and where the covariance is finite; then the variogram tends to a sill. The stationarity conditions may not represent the real conditions in some applications where the mean may not be known and the variance may not be finite; to overcome these problems a less stringent set of conditions is introduced, which is referred to as the "intrinsic hypothesis".

In the intrinsic hypothesis the process satisfies the following conditions:

$$E [Z(\underline{x}_1, H) - Z(\underline{x}_2, H)] = m$$

where m is an unknown constant.

and

var
$$[Z(\underline{x}_1, H) - Z(\underline{x}_2, H)] = 2 \Gamma(\underline{d})$$

i.e.: the variance is a function of \underline{d} , the distance vector given by $\underline{d} = \underline{x}_1 - \underline{x}_2$, and not of \underline{x}_1 and \underline{x}_2 .

Under these conditions we shall define the following Kriging estimator:

$$Z_{o}^{*} = Z^{*}(X_{o}) = \sum_{i=1}^{n} X_{o}^{i} Z_{i}$$
 (30)

where, as before, λ_0^{i} is a set of weights which has to be determined.

For the moment we shall assume that the mean is constant, though unknown, i.e.:

$$E[Z_0^*] = E[Z_0] = m$$
 (31)

Since we are looking for an unbiased estimator, this means that:

$$E[Z_0^* - Z_0] = 0$$

where:

$$Z_0$$
 is the true value of the variable Z at \underline{x}_0 , which is unknown.

Introducing the estimator (30) into (31), we get

$$E[\sum_{i=1}^{n} \chi_{o}^{i} Z_{i}] = E[Z_{o}] = m$$

 \mathbf{or}

$$\sum_{i=1}^{n} \chi_0^i E[Z_i] = m$$

i.e.:

$$\sum_{i=1}^{n} \chi_{o}^{i} m = m$$

which implies that:

$$\sum_{i=1}^{n} \chi_0^{i} = 1$$
(32)

Now, as before in the second-order stationarity case, we impose the condition of minimum variance on the estimator, i.e. we need:

minimum E [
$$(Z_0^* - Z_0)^2$$
]. (33)

Introducing the estimator (30) into (33) gives:

$$E[(Z_{o}^{*} - Z_{o})^{2}] = E[(\Sigma_{i} \lambda_{o}^{i} Z_{i} - Z_{o})^{2}]$$
(34)

But now, because of (32), we can have the identity:

$$Z_{o} = \sum_{i=1}^{n} \chi_{o}^{i} Z_{o}$$
(35)

and introducing (35) into (34):

$$E [(Z_0^* - Z_0)^2] = E[(\sum_{i} \chi_0^i Z_i - \sum_{i} \chi_0^i Z_0)^2]$$

and factorizing:

$$E [(Z_0^* - Z_0)^2] = E[\{ \sum_{i} \chi_0^i (Z_i - Z_0)\}^2]$$

or

$$E[(Z_0^* - Z_0)^2] = E[\sum_{i} \chi_0^i (Z_i - Z_0) \sum_{j} \chi_0^j (Z_j - Z_0)]$$

which leads to

$$E[(Z_{0}^{*} - Z_{0})^{2}] = \sum_{i j} \sum_{\lambda_{0}^{i}} \lambda_{0}^{j} E[(Z_{i} - Z_{0})(Z_{j} - Z_{0})].$$
(36)

Recalling the definition of the semi-variogram in equation (1):

$$\Gamma(\underline{x}_{i},\underline{x}_{j}) \equiv \frac{1}{2} \mathbb{E} \left[\{ \mathbb{Z}(\underline{x}_{i},\mathbb{H}) - \mathbb{Z}(\underline{x}_{j},\mathbb{H}) \}^{2} \right]$$

or simply

$$\Gamma(\underline{x}_{i} - \underline{x}_{j}) = \frac{1}{2} \mathbb{E} \left[(Z_{i} - Z_{j})^{2} \right]$$
(37)

which can be re-arranged (adding and subtracting $\rm Z_{O}$) as

 $\Gamma(x_i - x_j) = \frac{1}{2} E [\{(Z_i - Z_o) - (Z_j - Z_o)\}^2]$

On developing the squared binomial:

$$\Gamma(\underline{x}_{i}-\underline{x}_{j}) = \frac{1}{2} \mathbb{E}[(Z_{i}-Z_{o})^{2}] - \mathbb{E}[(Z_{i}-Z_{o})(Z_{j}-Z_{o})] + \frac{1}{2} \mathbb{E}[(Z_{j}-Z_{o})^{2}]$$
(38)

and using the semi-variogram definition [equation (37)] again, in the first and third terms on the right side of (38), we get

$$\Gamma(\underline{x}_{i}-\underline{x}_{j}) = \Gamma(\underline{x}_{i} - \underline{x}_{o}) - \mathbb{E}[(Z_{i}-Z_{o})(Z_{j}-Z_{o})] + \Gamma(\underline{x}_{j} - \underline{x}_{o})$$

which allows us to compute the expectation needed in equation (36), i.e.

$$E[(Z_i-Z_o)(Z_j-Z_o)] = -\Gamma(\underline{x}_i-\underline{x}_j) + \Gamma(\underline{x}_i-\underline{x}_o) + \Gamma(\underline{x}_j-\underline{x}_o)$$

which, when introduced into (36) gives:

$$E[(Z_{o}^{*}-Z_{o})^{2}] = -\sum_{i j} \chi_{o}^{i} \chi_{o}^{j} \Gamma(\underline{x_{i}}-\underline{x_{j}}) + \sum_{i j} \chi_{o}^{i} \chi_{o}^{j} \Gamma(\underline{x_{i}}-\underline{x_{o}}) + \sum_{i j} \chi_{o}^{i} \chi_{o}^{j} \Gamma(\underline{x_{j}}-\underline{x_{o}})$$

$$+ \sum_{i j} \chi_{o}^{i} \chi_{o}^{j} \Gamma(\underline{x_{j}}-\underline{x_{o}})$$

$$(39)$$

This can be simplified, because we can factorize by $\Sigma \chi_0^{i}$ or $\Sigma \chi_0^{j}$ in the last two terms of (39), both factors being exactly 1 according to (32); in addition

$$\sum \chi_{0}^{i} \Gamma(\underline{x}_{i}-\underline{x}_{0}) = \sum \chi_{0}^{j} \Gamma(\underline{x}_{j}-\underline{x}_{0})$$
whence (39) reduces to
$$E[(Z_{0}^{*}-Z_{0})^{2}] = -\sum_{i j} \chi_{0}^{i} \chi_{0}^{j} \Gamma(\underline{x}_{i}-\underline{x}_{j}) + 2\sum_{i} \chi_{0}^{i} \Gamma(\underline{x}_{i}-\underline{x}_{0}) \quad (40)$$

This is the quadratic form corresponding to equation (20) in the second-order stationary case. Equation (40) has to be minimized, subject to the equality constraint represented by equation (32). The minimization can be carried out via the Lagrange multipliers technique, minimizing (without restrictions) the expanded objective function:

$$L(\chi_{0}^{i}, \mu) = \frac{1}{2} E[(Z_{0}^{*} - Z_{0})^{2}] - \mu [\sum_{i} \chi_{0}^{i} - 1]$$
(41)

where the factor $\frac{1}{2}$ in the first term and the negative sign in the second have been added for convenience, but do not alter the main objective. The coefficient " μ " is the (unknown) Lagrange multiplier, which has to be determined, together with the set of weights λ_0^{i} ; to do so we have to make the partial derivatives of (41), with respect to μ and λ_0^{i} (i=1,...,n), equal to zero, which leads to a set of (n+1) linear equations in (n+1) unknowns:

$$\begin{bmatrix} \Sigma & \chi_0^j & \Gamma(\chi_i - \chi_j) + \mu = \Gamma(\chi_i - \chi_0) & i = 1, \dots, n \\ j & & \\ \Sigma & \chi_0^i = 1 \\ i & & \\ \end{bmatrix}$$
(42)

On introducing the notation:

$$r_{ij} = r(x_i - x_j)$$

the linear system represented by (42) can be expanded as follows:

$$\begin{bmatrix} 0 & \Gamma_{12} & \Gamma_{13} & \cdots & \Gamma_{1n} & 1 \\ \Gamma_{12} & 0 & \Gamma_{23} & \cdots & \Gamma_{1n} & 1 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ \Gamma_{n1} & \Gamma_{n2} & \Gamma_{n3} & \cdots & 0 & 1 \\ 1 & 1 & 1 & \cdots & 1 & 0 \end{bmatrix} * \begin{bmatrix} \chi_0^1 \\ \chi_0^2 \\ \vdots \\ \vdots \\ \chi_0^n \\ \mu \end{bmatrix} = \begin{bmatrix} \Gamma_{10} \\ \Gamma_{20} \\ \vdots \\ \vdots \\ \Gamma_{n0} \\ 1 \end{bmatrix}$$
(43)

•The solution of (43) exists and is unique provided that the matrix of coefficients is positive definite. The matrix is symmetric and it only needs to be solved once, since it is just the right hand side vector that depends on the point (\underline{x}_0) being estimated.

We now proceed to compute the variance of the error of the estimator, which is:

$$var(Z_0^*-Z_0) = E[(Z_0^*-Z_0)^2] - \{ E[Z_0^*-Z_0] \}^2$$
(44)

which, on introducing the condition of unbiased estimator [equation (29)], reduces to :

which has already been computed in (40). In addition, we can use the results of the linear system (42) to get:

$$\operatorname{var}(Z_{O}^{*}-Z_{O}) = - \sum_{i} (\lambda_{O}^{i} \Gamma(\underline{x}_{i}-\underline{x}_{O})-\mu) + 2\sum_{i} \lambda_{O}^{i} \Gamma(\underline{x}_{i}-\underline{x}_{O})$$

or

$$\sigma_{0}^{2} = \operatorname{var}(Z_{0}^{*}-Z_{0}) = \sum_{i} \chi_{0}^{i} \Gamma(\chi_{i}-\chi_{0}) + \mu$$
(45)

B.3.3. Introducing uncertainty in the measurement data.

In the Kriging estimation procedures already reviewed, we have implicitly assumed that the data were actually error-free. This is not the case in practical applications, where, as in water distribution networks, the data come from measurements, which do have an associated error, say ϵ_i , corresponding to each measurement Z_i.

To handle the measurement errors within Kriging, the following assumptions are made:

- The errors are random with zero mean: $E[\epsilon_i]=0$ $\forall i=1,2,..,n$ - The errors are uncorrelated to each other:

 $Cov(\epsilon_i, \epsilon_j) = 0 \quad \forall i \neq j$

- The errors are uncorrelated with the measurements:

$$Cov(\epsilon_i, Z(\underline{x})) = 0 \quad \forall i \quad \forall \underline{x}$$

- The error variance is known: σ_i^2 and $\sigma_i \neq \sigma_j$ ¥i≠j

Of course, in the event that some measurements are error-free we only need σ_k =0, where k is the index of that particular error-free measurement.

Under these assumptions, the impact of measurement errors on the Kriging equations is noticeable only in equation (43), the linear system of equations, where the first "n" zero diagonals are replaced by $-\sigma_1^2$; the restriction for the λ_0^{i} adding up to 1 and the computation of the Kriging variance remains unchanged. In summary, we get:

To determine the weights $\lambda_0{}^i$ and the Lagrange multiplier u:

$$\sum_{j} \lambda_{0}^{j} \Gamma(x_{1}-x_{j}) - \lambda_{0}^{i} \sigma_{1}^{2} + \mu = \Gamma(x_{1}-x_{0})$$

$$= 1, \dots, n$$

$$\sum_{i} \lambda_{0}^{i} = 1$$
(46)

or, using the notation

$$\Gamma_{ij} = \Gamma(\underline{x}_{i} - \underline{x}_{j})$$

$$\begin{bmatrix} -\sigma_{1}^{2} & \Gamma_{12} & \Gamma_{13} & \dots & \Gamma_{1n} & 1 \\ \Gamma_{12} & -\sigma_{2}^{2} & \Gamma_{23}^{2} & \dots & \Gamma_{1n} & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \Gamma_{n1} & \Gamma_{n2} & \Gamma_{n3} & \dots & -\sigma_{n}^{2} & 1 \\ 1 & 1 & 1 & \dots & 1 & 0 \end{bmatrix} * \begin{bmatrix} \underline{x}_{0}^{1} \\ \underline{x}_{0}^{2} \\ \vdots \\ \vdots \\ x_{0}^{n} \\ \mu \end{bmatrix} = \begin{bmatrix} \Gamma_{10} \\ \Gamma_{20} \\ \vdots \\ \vdots \\ \Gamma_{n0} \\ 1 \end{bmatrix}$$

$$(47)$$

To determine the Kriging variance:

$$\sigma_{0}^{2} = \operatorname{var}(Z_{0}^{*}-Z_{0}) = \sum_{i} \chi_{0}^{i} \Gamma(x_{i}-x_{0}) + \mu$$
 (48)

i.e. exactly the same as equation (45).

B.3.4. Universal Kriging and k-th order random functions (k-IRF).

In many practical applications the mean cannot be assumed constant, i.e.: we have $E[Z(\underline{x})] = m(\underline{x})$, a function of \underline{x} . Because of this, the real semi-variogram cannot be obtained from the raw measurement data. This is evidenced in:

$$\Gamma_{raw}(\underline{x}_1 - \underline{x}_2) = \Gamma_{real}(\underline{x}_1 - \underline{x}_2) + \frac{1}{2} \{ E[Z(\underline{x}_1) - Z(\underline{x}_2)] \}^2$$
(49)

Thus, the real semi-variogram coincides with that extracted from the raw data only when $E[Z(x_1)]=E[Z(x_2)]=m=constant$. If a drift exists, i.e.: $m=m(\underline{d})$, then equation (49) says that the semi-variogram estimated from the raw data will be a <u>biased</u> estimate of the real semi-variogram.

Nevertheless, some authors [Delhomme, (1979)] accept the possibility that, when the drift or trend of the mean is "small", the intrinsic hypothesis can still give relatively accurate results for short distances. We shall concentrate on the most general case, where that approximation is not acceptable, i.e. in the case where a strong trend is present; in that case, a generalisation of the intrinsic hypothesis will allow us to tackle the problem within the framework of Kriging.

In fact, in the intrinsic hypothesis, what we did was to assume that if the original process Z(x) was not stationary, then the (or increments) were indeed stationary. The differences generalisation we are looking at does not stop there; if this difference is still not stationary, then why not take the difference again, and again and again, until a stationary process found. The successive differences are then removing the nonis stationarity. The idea is not new; in fact it is widely used in stochastic processes, particularly in detrending time series data.

Then, from now on, we are going to assume that $Z(\underline{x})$ has been differenced as many times as needed in order to produce a stationary process. Thus, $Z(\underline{x})$ represents the differenced stationary process and not the original one.

We shall follow the same sequence as before, firstly we address the problem of how to infer the semi-variogram from the raw data and, secondly, how to carry out the estimation itself via universal (or "generalised") Kriging.

B-33

Following Delhomme (1978) and de Marsily (1986), both based on previous work by Matheron, and for a random function Z, we shall

refer to $\sum_{i=0}^{n} \chi^{i} Z_{i}$ as a generalised increment of order k if the

weights λ^{i} satisfy the condition:

$$\sum_{i=0}^{n} \lambda^{i} f^{1}(\underline{x}_{i}) = 0$$
(50)

for all the monomial functions $f^1(\underline{x}_i)$ of degree less or equal than "k"; e.g. in a two-dimensional geometrical space, where $\underline{x}=(\underline{x}_1,\underline{x}_2)^T$, this implies that the maximum degree must be $\underline{x}_1^{p*\underline{x}_2^q}$ with $0 \leq p+q \leq k$.

In other words, in a first-order generalised increment we are asking for the simultaneous fulfilment of:

$$\begin{array}{c}
\overset{n}{\Sigma} \chi^{i} = 0 \\
\overset{n}{\Sigma} \chi^{i} \times 1^{i} = 0 \\
\overset{n}{\Sigma} \chi^{i} \times 1^{i} = 0 \\
\overset{n}{\Sigma} \chi^{i} \times 2^{i} = 0
\end{array}$$
(51)

where the superscript "i" is just an index, not an exponent.

Similarly, in the case of a second-order generalised increment we shall be asking for the extended conditions:

$$\begin{bmatrix} n \\ \Sigma \\ i=0 \end{bmatrix}^{n} \chi^{i} = 0 \qquad \begin{bmatrix} n \\ \Sigma \\ i=0 \end{bmatrix}^{i} (x_{1}^{i})^{2} = 0 \\ i=0 \qquad i=0 \end{bmatrix}^{n} \begin{bmatrix} \chi^{i} \\ x_{2}^{i} \end{bmatrix}^{2} = 0 \\ i=0 \qquad i=0 \end{bmatrix}$$
(52)

Of course, the intrinsic hypothesis already seen in previous sections implies a generalised increment of order zero, i.e. only

$$\sum_{i=0}^{n} \chi^{i} = 0$$

is fulfilled.

The concept of the generalised covariance of order k is now introduced simply as the covariance of a generalised increment of order k. The generalised covariance is denoted by $K(\underline{d})$. For k=0, the generalised covariance of order zero is simply the semi-variogram, with a negative sign.

Thus, according to the theory of the intrinsic random functions [Matheron, (1973)], if

$$\sum_{i=0}^{n} \chi^{i} Z_{i}$$

is a generalised increment of order k, its variance is given by:

$$\operatorname{var}(\overset{n}{\Sigma}\chi^{i} Z_{i}) = \underset{i j}{\Sigma} \overset{\Sigma}{\chi}^{i} \chi^{j} K(\underline{x}_{i}-\underline{x}_{j})$$
(53)

As before, in the semi-variogram case, only some functions can be used to represent a generalised covariance; basically, they must guarantee that the variances of generalised increments must be always positive.

Matheron (1973) showed that appropriate models for the isotropic generalised covariance of order k can be represented by polynomials of order 2k+1; he also demonstrated that the even powers are redundant and, as a result, they do not need to be used in the polynomials.

Delfiner (1976) presented an analytical expression for these polynomials and also formalised the conditions for positive definiteness. Table B.3 presents some of the most widely used models for generalised covariances.

As can be seen from Table B.3, the generalised covariances are linear functions of the parameters A_i , which can then be computed via regression analysis or similar techniques.

Kitanidis (1983) evaluated different techniques for the estimation of the parameters of the generalised covariance models. He stressed the point that the determination of the parameters A_i of the generalised covariance model is a process involving considerable errors; he recommended a maximum likelihood technique and a minimum variance unbiased quadratic technique as adequate methods for fitting the unknown parameters to the generalised covariance.

B-36

Drift	Drift order k	Model for the generalised covariance
Constant	0	$K(d) = A_0 \delta + A_1 d $
Linear	1	$K(d) = A_0 \delta + A_1 d + A_3 d ^3$
Quadratic	2	$K(d) = A_0 \delta + A_1 d + A_3 d ^3 + A_5 d ^5$
Constraints:	A_0^{20} , $A_3^{2} - (A_3^{2})^{2} - (A_3^{2})^{2} - \delta^{2} = 1$ if	A ₁ ≤0 , A ₅ ≤0 and 10/3) $\sqrt{A_1 A_5}$ (in R ²) $\sqrt{10 A_1 A_5}$ (in R ³) d=0, \$=0 otherwise

Having produced and checked a valid model for the generalised covariance, we can proceed to determine the universal Kriging equations, which is done, as usual, imposing the unbiased condition on the estimator Z_0^* and the minimum variance on the residual $Z_0^*-Z_0$. The difference now is in the fact that the mean is no longer constant, but has a trend or drift:

$$\mathbf{E}[\mathbf{Z}(\underline{\mathbf{x}})] = \mathbf{m}(\underline{\mathbf{x}}) \tag{54}$$

Although the value of $m(\underline{x})$ is not known, we assume that it can be represented (at least locally) in terms of a linear combination of basis functions (usually polynomials), e.g.: in the two-dimensional case, where $\underline{x} = (x_1, x_2)^T$: $m(\mathbf{x}) = a_1 + a_2 x_1 + a_3 x_2 + a_4 x_1^2 + a_5 x_1 x_2 + a_6 x_2^2 + \dots$ (55)

or, in general :

$$m(\underline{x}) = \sum_{k=1}^{w} a_k p^k(\underline{x})$$
(56)

where :

- the coefficients a_i have to be determined by fitting, applying some knowledge of the physical behaviour of the variable Z.
- the polynomials $p^k(\underline{x})$ are usually monomials of first or second degree (linear or quadratic drift), for example in a twodimensional space:

* linear drift: $m(\underline{x}) = a_1 + a_2x_1 + a_3x_2$ * quadratic drift:

 $m(x) = a_1 + a_2x_1 + a_3x_2 + a_4x_1^2 + a_5x_1x_2 + a_6x_2^2$ * constant mean (no drift):

 $m(\underline{x}) = a_1$

As before, the Kriging estimator is :

$$Z_{o}^{*} = \sum_{i=1}^{n} \chi_{o}^{i} Z_{i}$$
(57)

The unbiased estimator must satisfy:

$$\mathbf{E}[\mathbf{Z}_{O}^{*}] = \mathbf{E}[\mathbf{Z}_{O}] \tag{58}$$

where Z_O is the true (unknown) parameter being estimated.

Introducing (54) and (57) into (58) gives:

$$E[\Sigma \chi_0^i Z_i] = E[Z_0] = m(\underline{x}_0)$$
(59)

and introducing (56) into (59):

$$\sum_{i=1}^{n} \chi_{o}^{i} \left\{ \sum_{k=1}^{w} p^{k}(\underline{x}_{i}) \right\} = \sum_{k=1}^{w} a_{k} p^{k}(\underline{x}_{o})$$

or, reordering the first term:

$$\sum_{k=1}^{w} a_{k} \left\{ \sum_{i=1}^{n} \lambda_{o}^{i} p^{k}(\underline{x}_{i}) \right\} = \sum_{k=1}^{w} a_{k} p^{k}(\underline{x}_{o})$$

which is satisfied when:

$$\sum_{i} \chi_{0}^{i} p^{k}(x_{i}) = p^{k}(x_{0}) \qquad k=1,...,m$$
(60)

This is a generalisation of the previous condition $\Sigma \chi_0^{1}=1$ which is obtained from (60) when no drift is assumed, i.e.: $m(\chi)=a_1$.

The rest of the Kriging equations are obtained from the minimum variance condition. For this purpose we could use the generalised covariance concept, provided that $\Sigma \ \lambda^i Z_i$ is a generalised increment of order k, which is true since:

$$\sum_{i=1}^{n} \chi_{0}^{i} Z_{i} = Z_{0}$$

can be re-written as:

$$\sum_{i=0}^{n} \chi_0^i Z_i = 0$$

with $\lambda_0^0 = -1$

and provided that the mean (drift) is expressed as a function of monomials of order "k" or lower.

Then, we can minimize the generalised covariance (equation 53) subject to the condition (60) by making zero all the partial derivatives with respect to λ_0^{i} of the Lagrangian function:

$$L(\lambda_0^{i}, \mu_k) = \sum_{i j} \sum_{k_0^{i} \lambda_0^{j}} K(x_i - x_j) + \mu_k [\sum_{i} \lambda_0^{i} p^k(x_i) - p^k(x_0)]$$
(61)

which leads to the following set of linear equations:

and $\begin{bmatrix} \sum_{j} \chi_{0}^{j} K(x_{1}-x_{j}) + \sum_{k} \mu_{k} p^{k}(x_{1}) = K(x_{1}-x_{0}) \\ j = 1, \dots, n \\ j = 1, \dots, n \\ j = 1, \dots, n \\ i & k = 1, \dots, m \end{bmatrix}$ (62)

This is a system of (n+m) equations in (n+m) unknowns: λ_0^i i=1,...,n and μ_k k=1,...,m, which can be expanded as:

$$\begin{bmatrix} K_{11} & K_{12} & K_{13} & \cdots & K_{1n} & 1 & P_{12} & \cdots & P_{1k} \\ K_{21} & K_{22} & K_{23} & \cdots & K_{1n} & 1 & P_{22} & \cdots & P_{2k} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ K_{n1} & K_{n2} & K_{n3} & \cdots & K_{nn} & 1 & P_{n2} & \cdots & P_{nk} \\ 1 & 1 & 1 & \vdots & \vdots & 1 & 0 & 0 & \cdots & 0 \\ P_{21} & P_{22} & P_{23} & \cdots & P_{2n} & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ P_{m1} & P_{m2} & P_{m3} & \cdots & P_{mn} & 0 & 0 & \cdots & 0 \end{bmatrix} * \begin{bmatrix} \chi_{0}^{1} \\ \chi_{0}^{2} \\ \vdots \\ \chi_{0}^{n} \\ -\mu_{1} \\ -\mu_{2} \\ \vdots \\ \vdots \\ -\mu_{m} \end{bmatrix} = \begin{bmatrix} K_{10} \\ K_{20} \\ \vdots \\ K_{n0} \\ 1 \\ P_{20} \\ \vdots \\ P_{m0} \end{bmatrix}$$
(63)

using the notation

and

$$K_{ij} = K(x_i - x_j)$$

$$p_{ij} = p^i(x_j)$$

Note that the system (63) is symmetric. Its size is (n+m)x(n+m) which is not particularly large in water resources applications, since is mainly dominated by "n", the amount of data points (piezometric head measurements in our case). Appendix A includes a review of efficient linear solvers for symmetric positive definite linear systems and, because of that, we do not go into more details here.

In the particular case of the intrinsic random functions of first order, the linear system (62) reduces to:

and

$$\sum_{j}^{\Sigma} \chi_{0}^{j} K(\underline{x}_{1} - \underline{x}_{j}) - \mu_{1} - \mu_{2} x_{1}^{i} - \mu_{3} x_{2}^{i} = K(\underline{x}_{1} - \underline{x}_{0})]$$

$$\lim_{\substack{i=1, \dots, n \\ \underline{x}_{1} = (x_{1}^{i}, x_{2}^{i})^{T}}$$

$$(64)$$

and in the case of intrinsic random functions of second order, the linear system becomes:

and

$$\begin{bmatrix} \Sigma & \chi_0 j & K(\underline{x}_1 - \underline{x}_j) & + \sum_{k=1}^{6} \mu_k & p^k(\underline{x}_1) & = & K(\underline{x}_1 - \underline{x}_0) \end{bmatrix} \\ \downarrow & \downarrow \\ \downarrow \\ \vdots & \downarrow \\ \vdots & \chi_0 i & p^k(\underline{x}_1) & = & p^k(\underline{x}_0) \end{bmatrix} \qquad k=1, \dots, m$$
(65)

with $p^1(x)$, ..., $p^6(x)$ being 6 polynomials in x_1 , x_2 as defined in equation (52).

To compute the corresponding Kriging variance we need to go back to the expression for the generalised covariance (equation 53) and to introduce $\Sigma \chi^{j} K(\underline{x_{i}}-\underline{x_{j}})$ from the linear system (62), obtaining:

$$\operatorname{var}(Z_0^* - Z_0) = K(0) + \sum_{k} \mu_k p^k(\underline{x}_0) - \sum_{i} \aleph_0^i K(\underline{x}_i - \underline{x}_0)$$
(66)

which, for the particular case of an intrinsic random function of first order reduces to:

$$\operatorname{var}(Z_{0}^{*}-Z_{0}) = K(0) + \mu_{1} + \mu_{2}x_{1}^{0} + \mu_{3}x_{2}^{0} - \sum_{i} \lambda_{0}^{i} K(\underline{x}_{i}-\underline{x}_{0})$$
(67)

Equations (62) and (66) are the equations for the universal Kriging and the Kriging variance respectively, in the most general case of intrinsic random function of order k. It is not too difficult to re-derive all the previous equations obtained by assuming stationarity or the intrinsic hypothesis, by simply considering them as particular cases of the most general intrinsic random function of order k.

APPENDIX C

DERIVATION OF THE BI-CUBIC SPLINES APPROXIMATION EQUATIONS

C.1. Introducing cubic splines.

In words, a cubic spline is a set of cubic polynomials joined end to end in a continuous and smooth manner.

Mathematically, and following Hayes and Halliday (1974), a polynomial function s(x) is a cubic spline with ordered interior knots λ_1 , λ_2 , ..., λ_h (i.e.: $\lambda_1 < \lambda_2 < ... < \lambda_h$) if it satisfies the following conditions:

i) In each interval defined by the knots, the function s(x) is a polynomial of degree less or equal than 3.

ii) The function s(x), its first and second derivatives [s'(x)]and s''(x), respectively] are continuous. Higher order derivatives do not need to comply with the continuity condition, in particular, a cubic spline will have third derivatives which are discontinuous at the knots $[\lambda_i, i=1, 2, ..., h]$.

The set of interior knots (another set of knots known as exterior knots will be defined later on) divides the domain over which the independent variable (x) is defined into h+1 intervals. A cubic spline is defined as a piecewise polynomial function (of degree less or equal than 3) on each one of these intervals and, provided that the specified knots are different (i.e. simple knots), the cubic spline will be twice-differentiable over its domain. When two coincident knots are specified, the

differentiability of the spline is reduced in one order, at the particular point where the coincident knots are located. Thus, by an adequate knot selection, we are allowed to specify discontinuities in the derivatives of the splines and, indeed, in the spline polynomial itself.

The one-dimensional splines fitting problem consists of finding the analytic expression for s(x) which best approximates a set of "m" data points { x_r , $f(x_r)$ }, $r=1,2, \ldots m$, for a given set of knots λ_i , $i=1,2, \ldots h$. "Best" here is used in the leastsquares sense. Note that we are not asking for the function s(x)to pass through the data points, which is the <u>interpolation</u> problem, but only to <u>approximate</u> them; "fitting" is used as an equivalent term for "approximation".

Cox (1987a) argues that in the case of noisy data, approximation methods should be used instead of interpolation. Besides, reliable computer software for splines interpolation, particularly for two-dimensional splines interpolation with scattered data, is not widely available [Anthony and Cox (1987 b)]. For these reasons, we shall apply a spline approximation approach to the piezometric head estimation problem.

From a mathematical point of view, a spline can be represented in various ways, but the B-splines representation (also known as "basic" or "fundamental" or "normalised" splines) has gained a reputation as being the most adequate and stable representation [see Cox (1972) and Cox (1982b) for a discussion on this subject].

Also, <u>cubic</u> splines, rather than splines of any other order, have became the most widely used splines in engineering applications, mainly due to their continuity properties and to the availability of proven and reliable software. The NAG subroutine library deals with cubic splines, while NPL's Data Approximation Subroutine Library deals with splines of any order.

Cubic B-splines are bell-shaped functions, which are defined to be non-zero only within four adjacent knot intervals (defined by five different knots λ_i), being zero everywhere else. Fig. C.1 shows the shape of a B-spline $M_i(x)$ with knots λ_{i-4} , λ_{i-3} , λ_{i-2} , λ_{i-1} and λ_i . Different basic splines can be defined, leading to different magnitudes for their peak point or area under the curve; this degree of freedom in the definition of a B-spline is used to scale-up the shape of the B-spline.

The following <u>explicit</u> representation of the B-spline, $M_i(x)$, has been shown to have some advantages over earlier divideddifference representations.

$$M_{i}(x) = \sum_{s=0}^{4} v_{is} (x - \lambda_{i-s})^{\frac{3}{4}}$$
(1)

with

$$v_{is} = 1 / \{ \pi (\lambda_{i-s} - \lambda_{i-r}) \}$$
(2)
$$r = 0 r \neq s$$

and where the notation $(x - \lambda_j)^2$ stands for:



Fig. C.1. Representation of a B-spline.

$$(x - \lambda_j)_{\downarrow}^3 = \begin{cases} (x - \lambda_j)^3 & \text{if } x > \lambda_j \\ 0 & \text{if } x \leq \lambda_j \end{cases}$$
(3)

i.e. the function $(x - \lambda_j)^2$ has non-zero non-negative values only from $x = \lambda_j$ onwards.

The weighting factors v_{is} in equation (2), have been defined such that the area under the cubic B-spline is exactly equal to 1/4. Note that the weights do not depend either on the data points or on the the points to be approximated, but only on the knots; they can therefore be pre-computed and kept in storage, before the actual approximation process is carried out.

The cubic B-splines previously defined are the "building blocks" of splines fitting. The idea is now to express the polynomial function s(x), describing the spline as a <u>linear</u> <u>combination</u> of the B-splines, having one B-spline associated with each knot. As a result, a cubic spline combines four Bsplines in an interval $(\chi_i - \chi_{i-1})$, i.e.: the B-splines $M_i(x)$, $M_{i+1}(x)$, $M_{i+2}(x)$ and $M_{i+3}(x)$, as shown in Fig. C.2.

Since normally we shall be interested in defining a spline for a limited range of the independent variable x, say $x \in [a,b]$, and due to the fact that for a convenient representation of the spline in any interval we want to combine 4 B-splines (as shown in Fig. C.2), and assuming that:

$$a < \lambda_1 < \lambda_2 < \ldots < \lambda_h < b$$

this implies that, for completeness, we need 8 extra knots, known as <u>exterior knots</u>, to be able to produce 4 B-splines in the first and last intervals: $[a, \lambda_1]$ and $[\lambda_h, b]$, respectively. These exterior knots are normally labelled as:

 $\lambda_{i-3} < \lambda_{i-2} < \lambda_{i-1} < \lambda_0 \leq a$

and

b $\leq \lambda_{h+1} < \lambda_{h+2} < \lambda_{h+3} < \lambda_{h+4}$

Thus, the extended set of B-splines looks like that shown in Fig. C.3.



Fig. C.2. B-splines involved in the spline function for the interval (λ_{i-1}, λ_i).

Cox (1972) demonstrated the conveniences of using a recursive relationship for the numerical representation of the B-spline $M_i(x)$, instead of the explicit representation of equation (1). The reasons for such a preference are mainly concerned with stability considerations, particularly when two or more knots are placed very close to of each other, in which case the explicit representation "blows up". This does not change most of what we have previously said, but only changes the way the B-splines are computed.

The recursive expression proposed by Cox (1972), and adopted by Hayes and Halliday (1974) for evaluating B-splines, both in the NPL and NAG libraries is:

$$M_{n,i}(x) = \frac{(x - \lambda_{i-n}) M_{n-1,i-1}(x) + (\lambda_{i} - x) M_{n-1,i}(x)}{(\lambda_{i} - \lambda_{i-n})}$$
(4)

where:

 $M_{n,i}(x)$: represents a B-spline of degree (n-1), given as a linear combination of B-splines of lower degree. Thus, n=4 for our cubic B-splines and $M_{4,i}(i)$ becomes the equivalent of our previous $M_i(x)$.



Fig. C.3. Extended set of B-splines: interior and exterior knots.

The evaluation of the B-splines using equation (4) starts with:

$$M_{1,i}(\mathbf{x}) = \begin{bmatrix} 1./(\lambda_i - \lambda_{i-1}) & \text{if } \mathbf{x} \in [\lambda_{i-1}, \lambda_i] \\ 0.0 & \text{if } \mathbf{x} \notin [\lambda_{i-1}, \lambda_i] \end{bmatrix}$$
(5)

The recursive definition for the B-splines given by equations (4) and (5) implies that the area under its curve is exactly 1/n, that is to say, for the cubic B-spline (n=4), the area is exactly the same as before: 1/4.

For a point, either a data point or a point where an approximate value is needed, lying in the interval (χ_{j-1}, χ_j) , the sequence for evaluating a cubic B-spline is given in the following array, where all B-splines not shown are considered as zero, and we are proceeding from left to right:

The recursive scheme for the evaluation of B-splines represented by equation (4) is stable, even in the case of coincident knots ($\chi_i = \chi_{i-1}$). Sometimes, it may be desirable to have coincident knots, in order to introduce discontinuities in the second derivatives, first derivatives or in the function itself.

C.3. Data fitting using one-dimensional cubic splines.

It has already been said that the idea is to combine the cubic B-splines <u>linearly</u>, in order to produce a general cubic spline polynomial function s(x), that is to say:

$$s(\mathbf{x}) = \sum_{i=1}^{h+4} \mathbf{M}_{i}(\mathbf{x})$$
(6)

where:

 r_i are the linear weights (to be determined) corresponding to the different B-splines for each interval (χ_{i-1}, χ_i)

The curve fitting problem consists of representing the set of data points (measurements for example): $\{x_r, f(x_r)\}, r=1,2,\ldots,m,$ in the form of equation (6). This implies that the coefficients r_i have to be determined as the solution of the following linear system of equations:

$$s(x_{1}) = \sum_{i=1}^{h+4} \Gamma_{i} M_{i}(x_{1}) = f(x_{1}) = f_{1}$$

$$s(x_{2}) = \sum_{i=1}^{h+4} \Gamma_{i} M_{i}(x_{2}) = f(x_{2}) = f_{2}$$

$$\vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots$$

$$s(x_{m}) = \sum_{i=1}^{h+4} \Gamma_{i} M_{i}(x_{m}) = f(x_{m}) = f_{m}$$
(7)

or, in a compact format:

$$\sum_{i=1}^{h+4} F_i M_i(x_r) = f(x_r) = f_r \quad r=1,2,... m$$
 (8)

which is clearly a linear system of "m" equations in the h+4 unknowns Γ_i . Because normally m > h+4, (7) is an overdetermined system of equations and an exact solution does not exist; a solution in the least-squares sense in sought. This system is known as the <u>observation equations</u>.

An even more compact matrix representation of such a system is:

$$A \underline{\Gamma} = \underline{f} \tag{9}$$

where:

- A is a rectangular mx(h+4) matrix with coefficients: $A_{ri} = M_i(x_r).$
- \underline{r} and \underline{f} are column vectors of lengths (h+4)x1 and mx1, respectively.

The structure of the matrix A has some important features:

i) Due to the fact that for every data point there are only 4 non-zero B-splines involved, and they are adjacent: $M_{j-3}(x)$, $M_{j-2}(x)$, $M_{j-1}(x)$ and $M_j(x)$ (if the data point lies in the interval $\lambda_{j-1} \langle x \langle \lambda_j \rangle$, this implies that A has only four non-zero elements per row and they are adjacent. In addition, if the data points are sorted in ascending order, according to their magnitudes, the first of the four non-zero of each row of A will never lie at the left of the first non-zero of the previous row. Two or more non-zeros will lie in the same column of A, if and only if they lie in the same knot interval, say $(\lambda_{j-1}, \lambda_j)$.

ii) This means that A has a special banded structure and, as a result, this represents some advantages as far as the storage and handling of the system is concerned.

To illustrate this point, some examples of the structure of the observation matrix A are as follows:

Example a:



the corresponding structure of the observation matrix is:

where "x" represents a non-zero value.

Example b:



the corresponding structure of the observation matrix is:

		Г								•
		X								
		x	х	х	х					
			х	х	х	х				ĺ
		1			х	х	х	х		
		}			х	х	х	х		
Α	=	1			х	х	х	х		
		}				\mathbf{x}	х	х	х	
						х	х	х	х	
									х	
		L								J

Since normally m > h+4, i.e. there are more measurements than unknowns, the linear system (9) is <u>overdetermined</u>. Clearly, a

unique solution is not possible, and the best we can expect is a solution in a statistical sense, such as to minimise some preestablished criterion; this is normally achieved by resolving the overdetermination in the least-squares sense, which leads to the following set of "normal" equations [premultiplying equation (9) by the transpose of the observation matrix A]:

$$\mathbf{A}^{\mathrm{T}} \mathbf{A} \mathbf{\underline{\Gamma}} = \mathbf{A}^{\mathrm{T}} \mathbf{\underline{f}} \tag{10}$$

where:

[A^T A] is a squared (h+4)x(h+4) banded symmetric matrix, with bandwidth (2n-1).

The system (10) is not a very large linear system; it does not have hundreds or thousands of equations, and its banded structure allows us to use efficient specialised linear solvers.

The numerical solution of (10) can be obtained via a direct method, i.e. a Cholesky factorization [see Appendix A for a review on the solution of linear systems of equations]. However, the system (10) can present instabilities which may require a more stable numerical solution, i.e. orthogonal factorization, Givens rotations or Householder transformations [see Golub and Van Loan (1983) for details on these methods]. The illconditioning of (10) becomes manifest particularly when m >>h+4. Iterative solutions (like the conjugate gradient method) of the normal equations have also been reported [Cox (1982a)].

A few words ought to be said with respect to the existence of the solution of the normal equations. The solution indeed exists, if and only if the data fulfil some conditions known as the

Schoenberg and Whitney conditions; for interpolation problems, these conditions establish a relationship between the location of data points and knots, which is normally expressed as:

For one-dimensional approximation problems, the Schoenberg and Whitney conditions are restricted to any subset of the data abscissae; if the conditions hold for any such subset, then the solution of the normal equations exists and is indeed unique.

In essence, the solution of equation (10) leads to the weights $\underline{\Gamma}$ corresponding to each inter-knot interval, from λ_1 onwards. The weights $\underline{\Gamma}$ represent the way in which the B-splines are linearly combined together to approximate the data set.

Having determined $\underline{\Gamma}$ from equation (10), they can be introduced in the expression for the general cubic spline:

$$s(x) = \sum_{i=1}^{h+4} \Gamma_i M_i(x)$$

to obtain a polynomial function that can be used to approximate any unmeasured point $x \in [a,b]$. For any of these points, say x, we need to evaluate the set of 4 B-splines $M_i(x)$ corresponding to x, using the recursive stable equation (4).

One of the critical aspects of spline fitting is concerned with the placement of the knots λ_i . They not only determine the number of unknowns in the normal equations, but they also influence the conditioning of the problem and the accuracy of the approximation. In principle, it would seem that adding more knots should lead to an improved approximation, but this is not always true, since too many knots may produce <u>overfitting</u> of the data set, which is manifested through unwanted oscillations of the fitted spline [see Cox, Harris and Jones (1987)]. The automatic placement of knots in splines approximation is a subject of ongoing research.

C.4. Data fitting using bi-cubic B-splines.

So far, the application of cubic splines for approximating data with only one independent variable has been reviewed. In Chapter 6 of this Thesis, the convenience and inconvenience of modelling the piezometric heads of the water distribution system 85 а continuous surface, "floating" above the network has been discussed; in following this line of thinking, the onedimensional cubic spline approach is clearly insufficient for describing the piezometric plane, and two independent variables (x,y) are needed.

The one-dimensional cubic B-splines concepts can be extended to the two-dimensional case in a relatively straightforward manner, leading to a bi-cubic spline (or doubly cubic spline). In the two-dimensional case the measurement data and any point in the space, will be represented by the triples: $\{x, y, f(x,y)\}$, where (x,y) caters for the horizontal coordinates of each point, and f(x,y) represents the value of the function (piezometric head in

our case) at that particular point (x,y).

As in the one-dimensional case, the basic splines (which have not been defined yet) will be associated with a set of knots. To introduce the knots, let us consider the plane defined by the independent variables (x, y), and let us define a rectangle R on it, such as:

R = { (x, y) such that $a \leq x \leq b$ and $c \leq y \leq d$ } where:

a is the minimum value of the independent variable x.

b is the maximum value of the independent variable x.

c is the minimum value of the independent variable y.

d is the maximum value of the independent variable y.

In this rectangle R, two different set of interior knots λ_i i=1,2, ... h+1 and μ_j j=1,2,... k+1 are introduced, with similar characteristics as in the one-dimensional case although here, for convenience, b = λ_{h+1} and d = μ_{k+1} . These two sets of interior knots divide the rectangle R into a set of (h+1)(k+1) panels R_{ij} , as shown in Fig. C.4.

For completeness, similar to the one-dimensional case, let us add exterior knots to the sets λ_i and μ_j such that:

 $\lambda_{-3} < \lambda_{-2} < \lambda_{-1} < \lambda_{0} = a$ $b = \lambda_{h+1} < \lambda_{h+2} < \lambda_{h+3} < \lambda_{h+4}$ $\mu_{-3} < \mu_{-2} < \mu_{-1} < \mu_{0} = c$ $d = \mu_{k+1} < \mu_{k+2} < \mu_{k+3} < \mu_{k+4}$

and



Fig. C.4. Rectangular subspace R for the independent variables x and y, divided into panels R_{ij} by knots λ_i i=1,2, ... h+1 and μ_j j=1,2, ... k+1.

As in the one-dimensional case, where the B-splines were defined as the basis functions (or "building blocks") for representing general cubic splines functions, we need to specify basis functions for the two-dimensional problem. To do so, the concept of B-splines can be extended to the two-dimensional space in the following way.

Let us define two sets of one-dimensional B-splines, one related to the independent variable "x" and to the set of knots λ_i , say $M_i(x)$, and the second related to the independent variable "y" and to the set of knots μ_j , say $N_j(y)$. Both one-dimensional B-splines [i.e. $M_i(x)$ and $N_j(y)$], have similar properties to those used in the previous section.

Let us think of the sets of B-splines $M_i(x)$ and $N_j(y)$ as twoone-dimensional structures; then, we can consider the set of all cross-products (or tensor products) of $M_i(x)$ with $N_j(y)$ as the for the bi-cubic spline. functions Although the basis mathematical concept of a tensor product is rather abstract, in this case it can be thought of simply as a summation of the set of all possible combinations of scalar products of the elements of the B-splines $M_i(x)$ and $N_j(y)$, at each point of the (x, y)space (or at each panel R_{ij} of the subspace R). For more details on tensor products see: de Boor (1978; chapter XVII).

The general bi-cubic spline polynomial function becomes, then, a linear combination of the new basis function $M_i(x)N_j(y)$, i.e.:

$$s(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{h+4} \sum_{j=1}^{h+4} \Gamma_{ij} M_i(\mathbf{x}) N_j(\mathbf{y})$$
(11)

where:

 r_{ij} is the new vector of "weights" (to be determined).

Note that, because the original basis functions $M_i(x)$ and $N_j(y)$ were non-zero over four adjacent intervals only (λ_{i-4} , λ_{i-3} , λ_{i-2} , λ_{i-1} , λ_i and μ_{j-4} , μ_{j-3} , μ_{j-2} , μ_{j-1} , μ_j , respectively), and because of the cross-product properties, the new basis function $M_i(x)N_j(y)$ will be non-zero over only 4x4 = 16 adjacent panels [see shaded area in Fig. C.4]. Similarly, as the originals $M_i(x)$ and $N_j(y)$ had only 4 non-zero basis functions at any point "x", the new basis functions $M_i(x)N_j(y)$ will have only 16 non-zero basis functions at any point (x,y). If this (x,y) point lies in

the panel R_{ij} , the non-zero $M_i(x)N_j(y)$ will be those within the shaded area of Fig. C.4.

Thus, with the help of vectorial cross-products (tensor products), the new bi-cubic B-splines have been introduced, based on the one-dimensional B-splines previously defined. The evaluation of each one-dimensional B-spline [$M_i(x)$ or $N_j(y)$] remains exactly as before, i.e. the stable method (recursive algorithm) proposed by Cox (1972) will be used [equations (4) and (5)], rather than the explicit form [equation(1)], based on the same stability reasons mentioned earlier.

Having the basis functions, the problem of approximating the data points $\{x_r, y_r, f(x_r, y_r)\}$, r=1,2,...m has to be tackled. The spline fitting problem now consists of determining the coefficients Γ_{ij} , as the least-squares solution of the following overdetermined linear observation system:

h+4 k+4

$$\Sigma \Sigma \Gamma_{ij} M_i(x_r) N_j(y_r) = f(x_r, y_r) = f_r$$
 (12)
i=1 j=1 r=1, 2, ... m

which, in matrix format, becomes:

$$A \underline{\Gamma} = \underline{f} \tag{13}$$

where:

A is a (mx(h+4)(k+4)) matrix.

- $\underline{\Gamma}$ is an (h+4)(k+4)x1 column vector of unknowns.
- <u>f</u> is a mx1 column vector, with the data values of the function being approximated.

As in the one-dimensional fitting problem, the structure of A is determined by the way the data are re-ordered; thus, if the data points are sorted according to the panel $R_{i,i}$ to which they
belong, allowing the index "i" to run faster than "j" (i.e. following the panels in Fig. C.4 from left to right and from downwards upwards), then the observation matrix has the following structure:

where:

 $A_{i,i}$ is a rectangular (rx(h+4)) sub-matrix.

r is the number of data points within the panels $R_{1,i}$, $R_{2,i}$, ..., $R_{h+1,i}$ (i.e. the number of data points in the i-th row of panels in Fig. C.4).

With this data arrangement, the vector $\underline{\Gamma}$ gets its elements ordered by i=1,2, ... h+4 for fixed "j", where j=1,2, ... k+4.

Note that if Γ_{ij} is associated with the panel R_{ij} , then the arrangement of $\underline{\Gamma}$ follows the panels in Fig. C.4, from left to right and from downwards upwards.

The problem now becomes that of solving the overdetermined linear system (13). In principle, we can proceed as in the onedimensional case, and form the corresponding normal equations by pre-multiplying (13) by A^{T} :

$$[\mathbf{A}^{\mathrm{T}} \mathbf{A}] \mathbf{\underline{\Gamma}} = \mathbf{A}^{\mathrm{T}} \mathbf{\underline{f}}$$
(14)

which is now a (h+4)(k+4)x(h+4)(k+4) linear system of equations.

Unfortunately, for bivariate approximation problems with scattered data, there is no equivalent of the Schoenberg and Whitney conditions, to guarantee the existence and uniqueness of the solution. To make things worse, it is often the case that the observation matrix A is rank-deficient, usually due to the fact that it may happen that some panels R_{i,i} remain without data points within them, which implies that the least squares solution of the observation system has infinite solutions. A unique solution is produced via the addition of a new condition, usually in the sense that, of all the infinite possible solutions, a "minimum norm" solution is sought. This forces us to use more sophisticated linear solvers, able to cope with rank deficiency in a stable manner: singular value decomposition and orthogonal transformations are cited as safe methods to solve the problem [Cox (1986), Cox (1987a)]. These methods belong to a very specialised field of numerical linear algebra and, because of that, we do not present more details here.

C.5. The statistical problem in splines fitting.

From a statistical point of view, the values of the "weights", $\underline{\Gamma}$, determined through the solution of the normal equations (10) or (14), for the one-dimensional and two-dimensional cases, respectively, can be interpreted in the case of noisy measurements as the expected values of the weights.

Always in the statistical sense of the problem, we may not only be interested in the expected values of the weights, but also in the corresponding dispersion statistics (as the variance and covariance), in order to estimate the error associated with the approximates computed by the cubic-splines at unmeasured points.

This is particularly true when the measurements are subject to errors, where we want to know the impact of those errors on the spline estimates.

The solution to the latter problem can be obtained in a straightforward manner, when solving the least squares estimation problem, as we did through the normal equations. In fact, all the information required is already contained in the inverse of the matrix $A^{T}A$ [see Walpole and Myers (1985), section 10.4].

Under the assumption that the measurement errors are independent, with zero mean and variance σ^2 , the matrix $\sigma^2[A^TA]^{-1}$ represents the variance-covariance matrix, and displays directly the variances of the weights in the main diagonal, and the covariances in the off-diagonal coefficients. Thus, if:

$$C = [A^{T}A]^{-1} = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1r} \\ c_{21} & c_{22} & \dots & c_{2r} \\ \vdots & \vdots & & \vdots \\ c_{r1} & c_{r2} & \dots & c_{rr} \end{bmatrix}$$
(15)

then, the variances are given by:

$$\sigma^{2}_{\Gamma i} = \operatorname{var}(\Gamma_{i}) = \sigma^{2} c_{ii} \quad \forall i=1,2,\ldots,r \quad (16)$$

and the covariances by:

$$\sigma_{\text{rirj}} = \operatorname{cov}(\Gamma_i, \Gamma_j) = \sigma^2 c_{ij} \quad \forall i \neq j$$
(17)

The problem now is basically a numerical one, namely: how to obtain the coefficients of the matrix $C = [A^TA]^{-1}$ in (15), in an efficient way. In practice, this means without having to <u>compute</u> explicitly the inverse of A^TA .

The solution can be almost direct, and very efficient from the computational point of view, provided that the normal equations have been solved using a direct Gaussian elimination method, with a Cholesky factorization [see Appendix A for details]. This is actually the main argument in favour of direct solutions when solving the normal equations, since if an iterative scheme (conjugate gradient) is used, the statistical information becomes unavailable [see Cox (1982a)].

Let us assume that we do have the Cholesky factorization of the normal system of equations, say:

$$A^{T}A = R^{T}R$$
(18)

where:

and

R is an upper triangular matrix and R^{T} is a lower triangular matrix.

Hence, following Cox (1986 and 1987a), to compute the covariance cov (Γ_i , Γ_j) we need to specify two (rx1) auxiliary column vectors h_1 and h_2 , such that:

$$\mathbf{h}_1 = (0, 0, \dots, 0, 1, 0, 0, \dots 0)^{\mathrm{T}}$$
(19)

$$\mathbf{h}_2 = (0, 0, \dots, 0, 0, 1, 0, \dots, 0)^{\mathrm{T}}$$
(20)

where the "1" is in the i-th and j-th position, respectively.

Then, the covariance cov (Γ_i, Γ_j) can be computed as:

cov
$$(\Gamma_{i}, \Gamma_{i}) = \sigma^{2} h_{1}^{T} [A^{T}A]^{-1}h_{2}$$
 (21)

on introducing the Cholesky factorization (18):

cov
$$(\Gamma_{i}, \Gamma_{j}) = \sigma^{2} \underline{h}_{1}^{T} [R^{T}R]^{-1}\underline{h}_{2}$$
 (22)

which can be re-written as:

cov
$$(\Gamma_{i}, \Gamma_{j}) = \sigma^{2} (\mathbb{R}^{-T} h_{1})^{T} (\mathbb{R}^{-T} h_{2})$$
 (23)

Upon defining the auxiliary vectors \underline{u}_1 and \underline{u}_2 such that:

$$\mathbb{R}^{\mathrm{T}}\mathbf{u}_{1} = \mathbf{h}_{1} \tag{24}$$

and

$$R^{T} \underline{u}_{2} = \underline{h}_{2} \tag{25}$$

and noting that \underline{u}_1 and \underline{u}_2 can be computed via a forward substitution process, since \mathbb{R}^T is a lower triangular matrix (the Cholesky factor) and \underline{h}_1 and \underline{h}_2 are known vectors, [once Γ_i and Γ_j have been chosen, (equations 19 and 20)], then equation (23) can be re-written as:

$$\operatorname{cov} (\Gamma_{i}, \Gamma_{j}) = \sigma^{2} \underline{u}_{1}^{T} \underline{u}_{2}$$
(26)

which gives the covariance.

To compute the variance $[var(\Gamma_i)]$, the process is similar, though simpler, because now $\underline{h}_1 = \underline{h}_2 = \underline{h}$, such that:

$$\mathbf{h} = (0, 0, \dots, 0, 0, 1, 0, \dots, 0)^{\mathrm{T}}$$
(27)

with the "1" in the i-th position. Then, on applying:

$$\operatorname{var}(\Gamma_{i}) = \sigma^{2} h^{T} (A^{T}A)^{-1} h \qquad (28)$$

and following the same procedure as before, we need to solve:

$$\mathbb{R}^{\mathrm{T}} \mathbf{u} = \mathbf{h} \tag{29}$$

and compute the variance as:

$$\operatorname{var}(\Gamma_{i}) = \sigma^{2} \underline{u}^{\mathrm{T}} \underline{u}$$
(30)

This is the efficient way to compute the variance and covariances, since it only requires two forward substitutions plus one inner product, to compute each covariance, and only one forward substitution plus an inner product for each variance. The explicit inversion of $[A^{T}A]$ has been avoided.

The variance can be computed for Γ_i i=1,2,...,r, where "r" is the number of weights, i.e.: r=h+4 and r=(h+4)(k+4), for the onedimensional and two-dimensional cubic splines, respectively.

Upon having the variance, and because the spline polynomial is just a linear combination of weighted B-splines, we can compute the respective variance of the splines estimates using the definition of the spline, and the following property of the variances of two independent random variables X and Y:

$$var(aX+bY) = a^2 var(X) + b^2 var(Y)$$
(31)

Thus, for the one-dimensional spline case:

$$var(s(x)) = var(\sum_{i=1}^{h+4} M_i(x) \Gamma_i) = \sum_{i=1}^{h+4} M_i(x) \}^2 var(\Gamma_i)$$
(32)

while, for the bi-cubic splines, a similar formula is obtained.

APPENDIX D

NUMERICAL RESULTS OF THE CALIBRATION EXERCISE

Table D.1. Summary of the comparison between true, initial and calibrated piezometric heads. Example A.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	MAXIMUM RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]
1	TRUE INITIAL INTERPOL. KRIGING SPLINES	62.136954 62.040502 62.123656 62.343461 36.283353	1.780894 1.810364 1.766405 1.602880 403.289659	1.334501 1.345498 1.329062 1.266049 20.082073	0.130198 0.053762 0.643541 51.459823	0.096452 0.014019 0.211023 25.900618	0.001502 0.000439 0.088072 402.249599	0.998448 0.999786 1.003323 0.583926
II	TRUE INITIAL INTERPOL. KRIGING SPLINES	62.136954 62.040502 62.122005 62.356386 33.503773	1.780894 1.810364 1.768055 1.589679 617.090146	1.334501 1.345498 1.329682 1.260825 24.841299	0.130198 0.055034 0.657572 76.383997	0.076452 0.015274 0.220374 28.680198	0.001502 0.000431 0.090076 631.640732	0.998448 0.999759 1.003531 0.539192
111	TRUE INITIAL INTERPOL. KRIGING SPLINES	62.136954 62.040502 62.124011 62.354969 33.578977	1.780894 1.810364 1.765906 1.591085 617.007472	1.334501 1.345498 1.328874 1.261382 24.839635	0.130198 0.053375 0.655473 76.361830	0.096452 0.013766 0.219276 28.604994	0.001502 0.000383 0.089847 631.661541	0.978448 0.999792 1.003509 0.540403
1V	TRUE INITIAL INTERPOL. KRIGING SPLINES	62.136954 62.040502 62.109964 62.355181 33.560587	1.780894 1.810364 1.684000 1.590768 624.443795	1.334501 1.345498 1.297690 1.261257 24.988873	0.130198 0.195602 0.654835 76.545792	0.096452 0.055025 0.219255 28.623284	0.001502 0.005610 0.089969 639.937695	0.998448 0.999566 1.003512 0.540108
V	TRUE INITIAL INTERPOL. KRIGING SPLINES	62.136954 62.040502 62.110472 62.339564 34.046619	1.780894 1.810364 1.798262 1.606536 540.108914	1.334501 1.345498 1.340993 1.267492 23.240243	0.13019B 0.108172 0.646702 71.219991	0.096452 0.026649 0.209684 28.137352	0.001502 0.001488 0.086764 546.933422	0.998448 0.999574 1.003261 0.547929
NOTES: ===== [1] : MAXIMUM ABSOLUTE VALUE RESIDUAL = MAX.(ABS(TRUE-ESTIMATED)) [2] : AVERAGE OF THE RESIDUALS=(SUMMATION OF ABSOLUTE VALUES OF THE RESIDUALS)/N [3] : VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = { SSAVR - N t AVRtt2 } / (N-2) SSAVR : SUMMATION SQUARES ABSOLUTE VALUES RESIDUALS AVR : AVERAGE ABSOLUTE VALUE RESIDUALS N : NUMBER OF DATA POINTS [4] : RATIO AVERAGE ESTIMATED VALUES/AVERAGE TRUE VALUES								

,

Table D.2. Summary of the comparison between true, initial and calibrated piezometric heads. Example B.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	MAXIMUM RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATID AVERAGES EST/TRUE [4]	
I	TRUE INITIAL INTERPOL. KRIGING SPLINES	53.068512 52.717053 53.068335 53.932906 34.311754	35.740184 35.932102 35.250806 32.593471 425.290821	5.978309 5.994339 5.937239 5.709069 20.622580	0.500262 0.204632 2.741375 49.416708	0.351459 0.046582 0.865810 19.209650	0.023382 0.004287 1.712570 279.293286	0.993377 0.999997 1.016288 0.646556	
 II	TRUE INITIAL INTERPOL. KRIGING SPLINES	53.068512 52.717053 53.065770 53.938091 34.253451	35.740184 35.932102 35.261951 32.571324 425.983291	5.978309 5.994339 5.938177 5.707129 20.639363	0.500262 0.204608 2.747570 49.413179	0.351459 0.046917 0.869579 19.267953	0.023382 0.004126 1.716390 279.822898	0.993377 0.999948 1.016386 0.645457	
III	TRUE INITIAL INTERPOL. KRIGING SPLINES	53.068512 52.717053 53.069100 53.940904 34.399069	35.740184 35.932102 35.247428 32.559316 423.970079	5.978309 5.974339 5.936954 5.706077 20.590534	0.500262 0.204608 2.751158 49.423569	0.351459 0.046862 0.872392 19.122335	0.023382 0.004300 1.716841 278.393582	0.993377 1.000011 1.016439 0.648201	
IV	TRUE INITIAL INTERPOL. KRIGING SPLINES	53.068512 52.717053 53.102704 53.934059 35.306851	35.740184 35.932102 33.922002 32.588918 372.016585	5.978309 5.994339 5.824260 5.708670 19.287731	0.500262 0.749739 2.742666 45.259086	0.351459 0.178470 0.866494 18.214553	0.023382 0.071623 1.713824 225.682231	0.993377 1.000644 1.016310 0.665307	
v	TRUE INITIAL INTERPOL. KRIGING SPLINES	53.068512 52.717053 52.805732 53.767893 33.952638	35.740184 35.932102 36.798907 33.305292 422.206230	5.978309 5.994339 6.066210 5.771074 20.547658	0.500262 1.098131 2.539383 50.320178	0.351459 0.266632 0.887150 19.568768	0.023382 0.110035 1.332804 274.153415	0.993377 0.995048 1.013179 0.639789	
NOTES [1]: [2]: [3]: [4]:	INTERS STATES AND AND AND AND AND AND AND AND AND AND								

Table D.3. Summary of the comparison between true, initial and calibrated piezometric heads. Example C.

.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	MAXIMUN RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]
I	TRUE INITIAL INTERPOL. KRIGING SPLINES	62.136956 62.878059 62.201538 62.341660 33.516220	1.780891 0.938416 1.641184 1.604408 621.797560	1.334500 0.968719 1.281087 1.266652 24.935869	1.066049 0.243984 0.645151 76.730861	0.741103 0.082440 0.210333 28.667753	0.149333 0.009005 0.087045 636.829002	1.011927 1.001039 1.003294 0.539393
II	TRUE INITIAL INTERPOL. KRIGING SPLINES	62.136956 62.878059 62.199914 62.350918 33.473344	1.780891 0.938416 1.642798 1.595251 622.038356	1.334500 0.968719 1.281717 1.263032 24.940697	1.066049 0.240220 0.654034 76.725757	0.741103 0.082454 0.216252 28.710629	0.149333 0.008611 0.089106 637.057828	1.011927 1.001013 1.003443 0.538703
III	TRUE INITIAL INTERPOL. NRIGING SPLINES	62.136956 62.878059 62.202519 62.355235 33.602055	1.780891 0.938416 1.640165 1.590922 621.290876	1.334500 0.968719 1.280689 1.261317 24.925707	1.066049 0.245286 0.658027 76.734329	0.741103 0.082802 0.219529 28.581918	0.149333 0.009111 0.089857 636.377371	1.011927 1.001055 1.003513 0.540774
IV	TRUE INITIAL INTERPOL. NRIGING SPLINES	62.136756 62.698589 62.150291 62.350108 33.560637	1.780891 0.981271 1.651821 1.595858 624.445852	1.334500 0.990591 1.285232 1.263273 24.988915	0.900043 0.163188 0.649373 76.545904	0,561633 0.075048 0,215389 28.623336	0.142984 0.006495 0.089098 639.939976	1.009039 1.000215 1.003430 0.540108
۷	TRUE INITIAL INTERPOL. KRIGING SPLINES	62.136956 62.741539 62.186270 62.314798 34.226023	1.780891 1.102719 1.752223 1.631318 523.721954	1.334500 1.050104 1.323716 1.277231 22.884972	0.847173 0.208746 0.619053 70.013527	0.604583 0.056611 0.203224 27.957950	0.090009 0.005916 0.077470 528.936504	1.009730 1.000794 1.002862 0.550816
<pre>NOTES: H</pre>								

.

Table D.4. Summary of the comparison between true, initial and calibrated piezometric heads. Example D.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	MAXIMUM RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATID AVERAGES EST/TRUE [4]	
I	TRUE INITIAL INTERPOL. KRIGING SPLINES	23.469569 22.684429 23.391649 24.309381 23.644199	92.318030 93.275043 90.859003 154.818543 118.963877	9.608227 9.657901 9.531999 12.442610 10.907056	1.141310 0.868950 23.097968 13.247504	0.785140 0.082335 1.403829 0.818050	0.009442 0.036972 17.676265 4.345164	0.966546 0.996680 1.035783 1.007441	
II	TRUE INITIAL INTERPOL. KRIGING SPLINES	23.469569 22.684429 23.389110 24.327721 23.668704	92.318030 93.275043 90.875278 154.661081 118.754159	9.608227 9.657901 9.532853 12.436281 10.897438	1.141310 0.870788 23.103627 13.253927	0.785140 0.083621 1.385520 0.804035	0.009442 0.036733 17.701572 4.380930	0.966546 0.996572 1.036564 1.008485	
III	TRUE INITIAL INTEPPOL. ARIGING SPLINES	23.469569 22.684429 23.401347 24.313210 23.705721	92.318030 93.275043 90.809464 154.788988 118.721159	9.608227 9.657901 9.529400 12.441422 10.895924	1.141310 0.856382 23.098135 13.247997	0.785140 0.086661 1.412530 0.795535	0.009442 0.035715 17.727623 4.367497	0.966546 0.997093 1.035946 1.010062	
IV	TRUE INITIAL INTERPOL. KRIGING SPLINES	23.469569 22.684429 23.388163 24.323520 23.705645	92.318030 93.275043 90.417509 154.703687 118.719071	9.608227 9.657901 9.508812 12.437994 10.895828	1.141310 1.236070 23.097093 13.248164	0.785140 0.123363 1.412602 0.795328	0.009442 0.056530 17.745359 4.368942	0.966546 0.996531 1.036385 1.010059	
v	TRUE INITIAL INTERPOL. KRIGING SPLINES	23.469569 22.684429 23.195262 23.661852 22.795223	92.318030 93.275043 89.943616 155.801525 120.156063	9.608227 9.657901 9.483861 12.482048 10.961572	1.141310 1.785770 22.440190 12.332561	0.785140 0.274644 1.908507 1.492429	0.009442 0.184791 15.371171 3.159400	0.966546 0.988312 1.008193 0.971267	
NOTES: [1]: [2]: [3]: [4]:	IOTES: (1]: MAXIMUM ABSOLUTE VALUE RESIDUAL = MAX.(ABS(TRUE-ESTIMATED)) (2]: AVERAGE OF THE RESIDUALS=(SUMMATION OF ABSOLUTE VALUES OF THE RESIDUALS)/N (3]: VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = { SSAVR - N ‡ AVR‡‡2 } / (N-2) SSAVR : SUMMATION SQUARES ABSOLUTE VALUES RESIDUALS AVR : AVERAGE ABSOLUTE VALUE RESIDUALS N : NUMBER OF DATA POINTS (4]: RATIO AVERAGE ESTIMATED VALUES/AVERAGE TRUE VALUES								

Table D.5. Summary of the comparison between true, initial and calibrated piezometric heads. Example E.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN, DEV.	NAXIMUM RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]
I	TRUE INITIAL INTERPOL. KRIGING SPLINES	90.580945 90.565063 90.576717 91.090711 90.641002	6.435581 6.444084 6.438456 10.609632 8.703889	2.536845 2.538520 2.537411 3.257243 2.950235	0.098798 0.091366 5.828270 3.736779	0.019310 0.009282 0.580414 0.372636	0.000276 0.000274 1.460438 0.419807	0.999825 0.999953 1.005628 1.00066 3
II	TRUE INITIAL INTERPOL. KRIGING SPLINES	90.580945 90.565063 90.576821 91.099321 90.633813	6.435581 6.444084 6.438314 10.613703 8.681062	2.536845 2.538520 2.537383 3.257868 2.946364	0.098798 0.090545 5.833304 3.723417	0.019310 0.009197 0.580345 0.369745	0.000276 0.000275 1.472233 0.406742	0.999825 0.999954 1.005723 1.000584
III	TRUE INITIAL INTERPOL. KRIGING SPLINES	90.580945 90.565063 90.576989 91.119078 90.696569	6.435581 6.444084 6.438518 10.677210 8.683292	2.536845 2.538520 2.537423 3.267600 2.946743	0.098798 0.089223 5.965369 3.877208	0.019310 0.009100 0.586325 0.345689	0.000276 0.000273 1.521479 0.448671	0.999825 0.999956 1.005941 1.001276
IV	TRUE INITIAL INTERPOL. KRIGING SPLINES	90.580945 90.565063 90.581264 91.116472 90.679997	6.435581 6.444084 6.431398 10.675143 8.646093	2.536845 2.538520 2.536020 3.267284 2.940424	0.098798 0.128526 5.959070 3.792621	0.019310 0.013024 0.585652 0.342484	0.000276 0.000556 1.517881 0.429829	0.977825 1.000004 1.005912 1.001074
v	TRUE INITIAL INTERPOL. KRIGING SPLINES	90.580945 90.565063 90.560459 91.124170 90.686189	6.435581 6.444084 6.430524 10.692372 8.691012	2.536845 2.538520 2.535848 3.269919 2.948052	0.098798 0.121307 5.985603 3.869142	0.019310 0.022875 0.589721 0.349365	0.000276 0.000486 1.531876 0.444341	0.999825 0.999774 1.005997 1.001162
NOTES: [1] : [2] : [3] : [4] :	MAXIMUM ABSO AVERAGE OF TI VARIANCE OF SSAVR : 9 AVR : 1 N : 1 RATIO AVERAGI	LUTE VALUE HE RESIDUALS THE ABSOLUTI SUMMATION SU AVERAGE ABS NUMBER OF DI E ESTIMATED	RESIDUAL = S=(SUMMATION E VALUE OF T DUARES ABSOL DLUTE VALUE NTA POINTS VALUES/AVER	MAX.(ABS(TF OF ABSOLUT HE RESIDUAL UTE VALUES RESIDUALS AGE TRUE VA	IUE-ESTIMATE E VALUES OF S = (SSAVR RESIDUALS	D)) The Residu - N \$ Avri	IALS)/N I\$2 } / (N-2	?)

D-5

Table D.6. Summary of the comparison between true, initial and calibrated piezometric heads. Example F.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	MAXIMUM RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]
I	TRUE INITIAL INTERPOL. KRIGING SPLINES	68.544348 68.415341 68.528709 69.688396 55.582129	376.086468 370.880707 374.788283 388.812582 865.717269	19.392949 19.258263 19.359449 19.718331 29.423074	0.376937 0.290282 8.010388 36.994771	0.165889 0.041405 1.381892 13.295901	0.014416 0.003382 2.680300 207.565165	0.998118 0.999772 1.016691 0.810893
II	TRUE INITIAL INTERPOL. KRIGING SPLINES	68.544348 68.415341 68.518600 69.696173 55.349084	376.086468 370.880707 374.463647 389.074867 894.105688	19.392949 19.258263 19.351063 19.724981 29.901600	0.376937 0.306612 8.031336 36.911984	0.165889 0.050655 1.380769 13.841272	0.014416 0.003443 2.690190 211.172737	0.998118 0.999624 1.016804 0.807493
III	TRUE INITIAL INTERPOL. KRIGING SPLINES	68.544348 68.415341 68.509465 69.705559 58.515908	376.086468 370.880707 374.194520 390.055931 661.971743	19.392949 19.258263 19.344108 19.749834 25.728812	0.376937 0.329019 8.085776 38.447331	0.165889 0.058857 1.395836 10.777761	0.014416 0.003935 2.709800 139.034831	0.998118 0.999491 1.016941 0.853694
IV	TRUE INITIAL INTERPOL. KRIGING SPLINES	68.544348 68.415341 68.497198 69.700594 58.511196	376.096468 370.880707 373.863877 389.899959 661.797609	19.392949 19.258263 19.336077 19.745885 25.725427	0.376937 0.460254 8.076814 38.451710	0.165889 0.071618 1.393295 10.781083	0.014416 0.007033 2.701085 138.996665	0.998118 0.999312 1.016869 0.853625
V	TRUE INITIAL INTERPOL. KRIGING SPLINES	68.544348 68.415341 68.476801 69.674382 58.658442	376.086468 370.880707 372.963306 389.029010 646.886053	19.392949 19.258263 19.312258 19.723818 25.433955	0.376937 0.443014 7.999846 38.507682	0.165889 0.091037 1.381557 10.611937	0.014416 0.012138 2.654147 134.807513	0.998118 0.999015 1.016486 0.855774
NOTES: ===== [1] : MAXINUM ABSOLUTE VALUE RESIDUAL = MAX.(ABS(TRUE-ESTIMATED)) [2] : AVERAGE OF THE RESIDUALS=(SUMMATION OF ABSOLUTE VALUES OF THE RESIDUALS)/N [3] : VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = { SSAVR - N t AVRtt2 } / (N-2) SSAVR : SUMMATION SQUARES ABSOLUTE VALUES RESIDUALS AVR : AVERAGE ABSOLUTE VALUE RESIDUALS N : NUMBER OF DATA POINTS [4] : RATIO AVERAGE ESTIMATED VALUES/AVERAGE TRUE VALUES								

Table D.7. Summary of the performance of the calibrated heads. Example A.

		PERF	ORMANC	E I	NDEX
CASE	PROCEDURE	AVERAGE	VARIANCE M R	AXIMUM ESID.	VARIANCE RESID.
I	INTERPOL.	0.981	0.758	0.829	0.915
	KRIGING	-3.584	-35.488 -	23.431	-999.999
	SPLINES	-999.999	-999.999 -9	99.999	-999.999
II	INTERPOL.	0.976	0.810	0.821	0.918
	KRIGING	-4.176	-41.100 -	24.508	-999.999
	SPLINES	-999.999	-999.999 -9	99.999	-999.999
III	INTERPOL.	0.982	0.741	0.832	0.935
	KRIGING	-4.109	-40.483 -	24.346	-999.999
	SPLINES	-999.999	-999.999 -9	99.999	-999.999
IV	INTERPOL.	0.922	-9.810	-1.257	-12.950
	KRIGING	-4.119	-40.622 -	24.296	-999.999
	SPLINES	-999.999	-999.999 -9	99.999	-999.999
v	INTERPOL.	0.925	0.653	0.310	0.019
	KRIGING	-3.413	-34.004 -	23.672	-999.999
	SPLINES	-999.999	-999.999 -9	99.999	-999.999
NOTE:	A VALUE OF -	999.999	INDICATES TH	AT THE	INDEX
	IS ACTUALLY	LESS THAN	V OR EQUAL T	0 -999.	999

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.957 -3.880 -999.999	-1.370 -38.340 -999.999	0.307 -24.051 -999.999	-2.033 -999.999 -999.999

	FREQUENCY OF IMPROVEMENT (%)						
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES			
INTERPOL. KRIGING SPLINES	100.000 0.000 0.000	80.000 0.000 0.000	80.000 0.000 0.000	80.000 0.000 0.000			

			· · · · · · · · · · · · · · · · · · ·	
		PERF	ORMANCE I	NDEX
CASE	PROCEDURE	AVERAGE	VARIANCE MAXIMUM RESID.	VARIANCE RESID.
I	INTERPOL.	1.000	-5.502 0.833	0.966
	KRIGING	-5.049	-267.833 -29.029	-999.999
	SPLINES	-999.999	-999.999 -999.999	-999.999
II	INTERPOL.	1.000	-5.209 0.833	0.969
	KRIGING	-5.122	-271.631 -29.165	-999.999
	SPLINES	-999.999	-999.999 -999.999	-999.999
III	INTERPOL.	1.000	-5.592 0.833	0.966
	KRIGING	-5.161	-273.701 -29.244	-999.999
	SPLINES	-999.999	-999.999 -999.999	-999.999
IV	INTERPOL.	0.991	-88.752 -1.246	-8.383
	KRIGING	-5.065	-268.612 -29.057	-999.999
	SPLINES	-999.999	-999.999 -999.999	-999.999
V	INTERPOL.	0.441	-29.432 -3.819	-21.146
	KRIGING	-2.960	-159.964 -24.767	-999.999
	SPLINES	-999.999	-999.999 -999.999	-999.999
NOTE:	A VALUE OF -	-999.999	INDICATES THAT THE	INDEX
	IS ACTUALLY	LESS THA	N OR EQUAL TO -999.	999

Table D.8. Summary of the performance of the calibrated heads. Example B.

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.886 -4.671 -999.999	-26.898 -248.348 -999.999	-0.513 -28.252 -999.999	-5.326 -999.999 -999.999

	FREQUENC	CY OF IMPR	OVEMENT ((%)
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	100.000 0.000 0.000	0.000 0.000 0.000	60.000 0.000 0.000	60.000 0.000 0.000

The second second second second second second second second second second second second second second second se		and the second second second second second second second second second second second second second second second		
		PERF	ORMANCE I	NDEX
CASE	PROCEDURE	AVERAGE	VARIANCE MAXIMUM RESID.	VARIANCE RESID.
I	INTERPOL.	0.992	0.973 0.948	0.996
	KRIGING	0.924	0.956 0.634	0.660
	SPLINES	-999.999	-999.999 -999.999	-999.999
II	INTERPOL.	0.993	0.973 0.949	0.997
	KRIGING	0.917	0.951 0.624	0.644
	SPLINES	-999.999	-999.999 -999.999	-999.999
III	INTERPOL.	0.992	0.972 0.947	0.996
	KRIGING	0.913	0.949 0.619	0.638
	SPLINES	-999.999	-999.999 -999.999	-999.999
IV	INTERPOL.	0.999	0.974 0.967	0.998
	KRIGING	0.856	0.946 0.479	0.612
	SPLINES	-999.999	-999.999 -999.999	-999.999
v	INTERPOL.	0.993	0.998 0.939	0.996
	KRIGING	0.913	0.951 0.466	0.259
	SPLINES	-999.999	-999.999 -999.999	-999.999
NOTE:	A VALUE OF	-999.999	INDICATES THAT THE	INDEX
	IS ACTUALLY	LESS THA	N OR EQUAL TO -999.	999

Table D.9. Summary of the performance of the calibrated heads. Example C.

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.994 0.905 -999.999	0.978 0.951 -999.999	0.950 0.564 -999.999	0.997 0.563 -999.999

	FREQUENC	Y OF IMPR	OVEMENT ((%)
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	100.000 100.000 0.000	100.000 100.000 0.000	100.000 100.000 0.000	100.000 100.000 0.000

•

	<u>, , , , , , , , , , , , , , , , , , , </u>	PERFORMANCE INDEX
CASE	PROCEDURE	AVERAGE VARIANCE MAXIMUM VARIANCE RESID. RESID.
I	INTERPOL. KRIGING SPLINES	0.990 -1.324 0.420 -14.333 -0.144 -999.999 -408.581 -999.999 0.951 -774.217 -133.729 -999.999
II	INTERPOL. KRIGING SPLINES	0.989 -1.273 0.418 -14.135 -0.195 -999.999 -408.782 -999.999 0.936 -762.062 -133.860 -999.999
III	INTERPOL. KRIGING SPLINES	0.992 -1.485 0.437 -13.308 -0.155 -999.999 -408.587 -999.999 0.910 -760.159 -133.739 -999.999
IV	INTERPOL. KRIGING SPLINES	0.989 -2.944 -0.173 -34.845 -0.183 -999.999 -408.550 -999.999 0.910 -760.038 -133.742 -999.999
v	INTERPOL. KRIGING SPLINES	0.878 -5.156 -1.448 -382.031 0.940 -999.999 -385.586 -999.999 0.262 -845.138 -115.761 -999.999
NOTE:	A VALUE OF - IS ACTUALLY	999.999 INDICATES THAT THE INDEX LESS THAN OR EQUAL TO -999.999

Table	D.	10.	Summary	\mathbf{of}	the	performance	\mathbf{of}	the	calibrated	heads.
		F	Example I).						

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.968 0.053 0.794	-2.436 -999.999 -780.323	-0.069 -404.017 -130.166	-91.730 -999.999 -999.999

	FREQUENC	CY OF IMPR	OVEMENT (%)
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR. RES
INTERPOL. KRIGING SPLINES	100.000 20.000 100.000	0.000 0.000 0.000	60.000 0.000 0.000	0.000 0.000 0.000

		PERF	ORMAN	CE I	NDEX
CASE	PROCEDURE	AVERAGE	VARIANCE	MAXIMUM RESID.	VARIANCE RESID.
I	INTERPOL.	0.929	0.886	0.145	0.014
	KRIGING	-999.999	-999.999	-999.999	-999.999
	SPLINES	-13.299	-999.999	-999.999	-999.999
II	INTERPOL.	0.933	0.897	0.160	0.007
	KRIGING	-999.999	-999.999	-999.999	-999.999
	SPLINES	-10.081	-999.999	-999.999	-999.999
III	INTERPOL.	0.938	0.881	0.184	0.022
	KRIGING	-999.999	-999.999	-999.999	-999.999
	SPLINES	-52.001	-999.999	-999.999	-999.999
IV	INTERPOL.	1.000	0.758	-0.692	-3.058
	KRIGING	-999.999	-999.999	-999.999	-999.999
	SPLINES	-37.897	-999.999	-999.999	-999.999
v	INTERPOL.	-0.664	0.646	-0.508	-2.101
	KRIGING	-999.999	-999.999	-999.999	-999.999
	SPLINES	-42.912	-999.999	-999.999	-999.999
NOTE:	A VALUE OF	-999.999	INDICATES	THAT THE	INDEX
	IS ACTUALLY	LESS THA	N OR EQUAL	TO -999	.999

Table D.11. Summary of the performance of the calibrated heads. Example E.

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.627 -999.999 -31.238	0.813 -999.999 -999.999	-0.142 -999.999 -999.999	-1.023 -999.999 -999.999

	FREQUENC	Y OF IMPR	OVEMENT (%)
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	80.000 0.000 0.000	100.000 0.000 0.000	60.000 0.000 0.000	60.000 - 0.000 0.000

		PERF	ORMANCE I	NDEX
CASE	PROCEDURE	AVERAGE	VARIANCE MAXIMUM RESID.	VARIANCE RESID.
I	INTERPOL.	0.985	0.938 0.407	0.945
	KRIGING	-77.643	-4.976 -450.616	-999.999
	SPLINES	-999.999	-999.999 -999.999	-999.999
II	INTERPOL.	0.960	0.903 0.338	0.943
	KRIGING	-78.716	-5.225 -452.981	-999.999
	SPLINES	-999.999	-999.999 -999.999	-999.999
III	INTERPOL.	0.927	0.868 0.238	0.925
	KRIGING	-80.021	-6.201 -459.157	-999.999
	SPLINES	-999.999	-999.999 -999.999	-999.999
IV	INTERPOL.	0.866	0.821 -0.491	0.762
	KRIGING	-79.329	-6.041 -458.137	-999.999
	SPLINES	-999.999	-999.999 -999.999	-999.999
v	INTERPOL.	0.726	0.640 -0.413	0.291
	KRIGING	-75.728	-5.181 -449.428	-999.999
	SPLINES	-999.999	-999.999 -999.999	-999.999
NOTE:	A VALUE OF	-999.999	INDICATES THAT THE	INDEX
	IS ACTUALLY	LESS THA	N OR EQUAL TO -999	999

Table	D.	12.	Summary	of	the	performance	of	the	calibrated	heads.
			Example H	Γ.						

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.893 -78.288 -999.999	0.834 -5.525 -999.999	0.016 -454.064 -999.999	0.773 -999.999 -999.999

	FREQUENC	CY OF IMPR	OVEMENT ((%)
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	100.000 0.000 0.000	100.000 0.000 0.000	60.000 0.000 0.000	100.000 0.000 0.000

Table D.13. Summary of the comparison between true, initial and calibrated flows. Example A.

-

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	MAXIMUM RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]
	TRUE	2.856538	14.284065	3.779427				
	INITIAL	2.859369	14.272675	3.777919	0.047100	0.019569	0.000320	1.000991
I	INTERPOL.	2.859375	14.271362	3.777746	0.050200	0.018475	0.000272	1.000993
	KRIGING	2.817956	14.436289	3,799512	0.617300	0.169644	0.041886	0.986494
	SPLINES	2.880875	14.210631	3.769699	0.389500	0.137150	0.021196	1.008520
	TRUE	2.856538	14.284065	3,779427				
	INITIAL	2.859369	14.272675	3.777919	0.047100	0.019569	0.000320	1.000991
П	INTERPOL.	2.855188	14.280479	3.778952	0.039200	0.012700	0.000137	0.999527
	KRIGING	2.851169	14.420433	3.797424	0.196100	0.057406	0.004129	0.998121
	SPLINES	2.840825	14.138351	3.760100	0.392700	0.101812	0.019210	0.994499
	TRUE	2.856538	14.284065	3.779427				
	INITIAL	2.859369	14.272675	3.777919	0.047100	0.019569	0.000320	1.000991
Ш	INTERPOL.	2,858187	14.280074	3.778898	0.029200	0.010775	0.000092	1.000578
	KRIGING	2.856588	14.384611	3.792705	0.191700	0.063800	0.003917	1.000018
	SPLINES	2.877188	14,141629	3,760536	0.369200	0.090925	0.013911	1.007229
	TRUE	2.856538	14.284065	3.779427				
	INITIAL	2.859369	14.272675	3.777919	0.047100	0.019569	0.000320	1.000991
I۷	INTERPOL.	2.859187	14.282264	3.779188	0.160800	0.049100	0.002967	1.000928
	LRIGING	2.856431	14.387656	3.793106	0.190100	0.078156	0.004605	0.999963
	SPLINES	2.879231	14.147690	3.761342	0.280900	0.102869	0.009580	1.007944
	TRUE	2.856538	14.284065	3.779427				
	INITIAL	2.859369	14.272675	3.777919	0.047100	0.019569	0.000320	1.000991
۷	INTERPOL.	2.857250	14.282359	3.779201	0.037200	0.011712	0.000186	1.000249
	KRIGING	2.953675	14.927282	3.863584	0.493500	0.175388	0.023376	1.034005
	SPI INES	2.967669	14.557097	3.815376	0.750500	0.175556	0.049815	1.038904

Table D.14. Summary of the comparison between true, initial and calibrated flows. Example B.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	MAXIMUM RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]
I	TRUE INITIAL INTERPOL. KRIGING SPLINES	6.364231 6.370150 6.370125 6.364119 6.403319	59.767118 59.722451 59.721357 59.754468 59.590010	7.730920 7.728030 7.727959 7.730101 7.719457	0.120100 0.119800 0.129000 0.625300	0.052619 0.052544 0.054500 0.228337	0.002130 0.002133 0.001539 0.050065	1.000930 1.000926 0.999982 1.006142
II	TRUE INITIAL INTERPOL. KRIGING SPLINES	6.364231 6.370150 6.362062 6.362281 6.355156	59.767118 59.722451 59.781664 59.814570 59.281006	7.730920 7.728030 7.731850 7.733988 7.699416	0.120100 0.118800 0.142300 0.545200	0.052619 0.039594 0.051950 0.167837	0.002130 0.001774 0.001880 0.041532	1.000930 0.999659 0.999694 0.998574
III	TRUE INITIAL INTERPOL. NRIGING SPLINES	6.364231 6.370150 6.364687 6.364988 6.372306	59.767118 59.722451 59.778066 59.814052 59.203720	7.730920 7.728030 7.731628 7.733955 7.694395	0.120100 0.118800 0.074700 0.580600	0.052619 0.034981 0.040094 0.177437	0.002130 0.002042 0.001251 0.042499	1.000930 1.000072 1.000119 1.001269
IV	TRUE INITIAL INTERPOL. KRIGING SPLINES	6.364231 6.370150 6.365875 6.364444 6.381637	59.767118 59.722451 59.765616 59.832539 58.795146	7.730920 7.728030 7.730822 7.735150 7.667799	0.120100 0.119800 0.340500 1.008700	0.052619 0.036169 0.103550 0.327106	0.002130 0.001956 0.017282 0.129761	1.000930 1.000258 1.000033 1.002735
V	TRUE INITIAL INTERPOL. KRIGING SPLINES	6.364231 6.370150 6.360250 6.636481 6.681875	59.767118 59.722451 59.776458 63.835945 62.332049	7.730920 7.728030 7.731524 7.989740 7.895065	0.120100 0.495800 1.157000 3.847700	0.052619 0.185394 0.382650 0.910531	0.002130 0.024362 0.076651 1.257681	1.000930 0.999374 1.042778 1.049911
NOTES: ===== [1] : MAXIMUM ABSOLUTE VALUE RESIDUAL = MAX.(ABS(TRUE-ESTINATED)) [2] : AVERAGE OF THE RESIDUALS=(SUMMATION OF ABSOLUTE VALUES OF THE RESIDUALS)/N [3] : VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = (SSAVR - N t AVRtt2) / (N-2) SSAVR : SUMMATION SQUARES ABSOLUTE VALUES RESIDUALS AVR : AVERAGE ABSOLUTE VALUE RESIDUALS N : NUMBER OF DATA POINTS [4] : RATIO AVERAGE ESTIMATED VALUES/AVERAGE TRUE VALUES								

Table D.15. Summary of the comparison between true, initial and calibrated flows. Example C.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	MAXINUM RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]	
	TRUE	2.856538	14.284065	3.779427					
	INITIAL	2.858787	14.277518	3.778560	0.073300	0.031275	0.000477	1.000788	
I	INTERPOL.	2.858813	14.280935	3.779013	0.139800	0.058012	0.002114	1.000796	
	KRIGING	2.843750	14.320438	3,784235	0.204500	0.061863	0.004110	0.995523	
	SPLINES	2.883269	14.204004	3.768820	0.427900	0.146269	0.025934	1.009358	
	TRUE	2.856538	14.284065	3.779427					
	INITIAL	2.858787	14.277518	3.778560	0.073300	0.031275	0.000477	1.000788	
II	INTERPOL.	2.864813	14.226220	3.771766	0.138800	0.046013	0.001928	1.002897	
	KRIGING	2.853413	14.393494	3.793876	0.161600	0.041737	0.002391	0.998906	
	SPLINES	2.841006	14.132745	3.759354	0,378800	0.099606	0,018068	0.994563	
	TRUE	2.856538	14.284065	3.779427					
	INITIAL	2.858787	14.277518	3.778560	0.073300	0.031275	0.000477	1.000788	
Ш	INTERPOL.	2.857250	14.293493	3.780574	0.144200	0.051950	0.002589	1.000249	
	KRIGING	2,855688	14.386401	3.792941	0.195100	0.057787	0.003410	1.000053	
	SPLINES	2.880925	14.132468	3.759317	0.360100	0.093737	0.015501	1.008537	
-	TRUE	2.856538	14.284065	3.779427					
	INITIAL	2.865569	14.252602	3.775262	0.150500	0.062294	0.003243	1.003162	
IV	INTERPOL.	2,858187	14.284406	3.779472	0.067600	0.026787	0.000526	1.000578	
	KRIGING	2.855669	14.370472	3.790841	0.151000	0.067819	0.003061	0.999696	
	SPLINES	2.879238	14.147652	3.761337	0.280900	0.102875	0.009580	1.007947	
	TRUE	2.856538	14.284065	3.779427					
	INITIAL	2.971906	14.869644	3.856118	0.383700	0.161781	0.013103	1.040388	
۷	INTERPOL.	2.843750	14.330026	3.785502	0.326600	0.093087	0.008230	0.995523	
	KRIGING	2.950088	14.859979	3.854864	0.342900	0.155000	0.013841	1.032749	
	SPLINES	2.966525	14.538424	3.812929	0.733500	0.178225	0.045836	1.038504	
NOTES:									
=====									
[1]:	NAXIMUN ABSO	LUTE VALUE	RESIDUAL =	MAX. (ABS (T	RUE-ESTIMATE	ED)}			
[2]:	AVERAGE OF TH	HE RESIDUAL	S=(SUMMATIO	N OF ABSOLUT	TE VALUES OF	THE RESIDU	JALS)/N		
[3]:	VARIANCE OF	THE ABSOLUT	E VALUE OF	THE RESIDUAL	.S = { SSAVE	R - N I AVR	112 } / (N-3	2)	
	SSAVR : S	SUMMATION S	QUARES ABSO	LUTE VALUES	RESIDUALS				
	AVR :	AVERAGE ABS	OLUTE VALUE	RESIDUALS					
	N : 1	NUMBER OF D	ATA POINTS						
[4]:	[4] : RATID AVERAGE ESTIMATED VALUES/AVERAGE TRUE VALUES								

Table D.16. Summary of the comparison between true, initial and calibrated flows. Example D.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	NAXINUM RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]	
I	TRUE INITIAL INTERPOL. KRIGING SPLINES	5.000003 4.999995 5.000039 4.997993 4.999997	46.933771 46.941706 47.312672 62.404206 57.302259	6.850823 6.851402 6.878421 7.899633 7.569826	0.523500 3.840300 16.024600 10.014300	0.064099 0.294291 2.165528 2.081771	0.008231 0.307374 10.264412 5.146496	0.999998 1.000007 0.999978 0.999999	
II	TRUE INITIAL INTERPOL. KRIGING SPLINES	5.000003 4.999995 4.997800 5.061939 5.116812	46.933771 46.941706 47.360765 69.411966 57.417879	6.850823 6.851402 6.881916 8.331384 7.577459	0.523500 4.315300 21.050000 9.752800	0.064099 0.308802 2.454491 2.017170	0.008231 0.355305 15.383614 5.133866	0.999998 0.999559 1.012387 1.023362	
III	TRUE INITIAL INTERPOL. KRIGING SPLINES	5.000003 4.999995 5.005050 4.993333 5.074139	46.933771 46.941706 47.203915 72.943998 57.927176	6.850823 6.851402 6.870511 8.540726 7.610991	0.523500 3.332300 33.209000 12.566800	0.064099 0.220570 2.555598 2.059450	0.008231 0.224825 18.477858 5.654480	0.999998 1.001009 0.998666 1.014827	
IV	TRUE INITIAL INTERPOL. KRIGING SPLINES	5.000003 4.999995 5.002828 4.994833 5.073988	46.933771 46.941706 47.563499 74.154957 57.994937	6.850823 6.851402 6.896630 8.611327 7.615441	0.523500 3.905700 33.654300 12.544100	0.064099 0.437233 2.615462 2.067881	0.008231 0.312139 19.334366 5.680241	0.977998 1.000565 0.978966 1.014797	
v	TRUE INITIAL INTERPOL. IRIGING SPLINES	5.000003 4.999995 4.999900 5.072791 5.152630	46.933771 46.941706 51.326560 75.137247 59.058413	6.850823 6.851402 7.164256 8.668174 7.684947	0.523500 9.799800 34.481700 12.571500	0.064099 1.020781 2.632107 2.103636	0.008231 2.927971 19.451386 5.817014	0.999998 0.999979 1.014557 1.030525	
NOTES: [1]: [2]: [3]: [4]:	NOTES: ===== [1] : MAXIMUM ABSOLUTE VALUE RESIDUAL = MAX.(ABS(TRUE-ESTIMATED)) [2] : AVERAGE OF THE RESIDUALS=(SUMMATION OF ABSOLUTE VALUES OF THE RESIDUALS)/N [3] : VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = (SSAVR - N & AVR**2) / (N-2) SSAVR : SUMMATION SQUARES ABSOLUTE VALUES RESIDUALS AVR : AVERAGE ABSOLUTE VALUE RESIDUALS N : NUMBER OF DATA POINTS [4] : RATIO AVERAGE ESTIMATED VALUES/AVERAGE TRUE VALUES								

.

Table D.17. Summary of the comparison between true, initial and calibrated flows. Example E.

.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	NAXIMUN RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]
I	TRUE INITIAL INTERPOL. KRIGING SPLINES	-0.714126 -0.717209 -0.716406 -0.336327 -0.647950	21.495997 21.352192 21.372964 18.456918 23.764316	4.636378 4.620843 4.623090 4.296152 4.874866	0.267600 0.443100 4.446700 5.367700	0.044977 0.071807 1.235058 1.073869	0.002412 0.004806 1.213173 1.134717	1.004318 1.003193 0.470963 0.907333
II	TRUE INITIAL INTERPOL. KRIGING SPLINES	-0.714126 -0.717209 -0.706772 -0.318991 -0.657779	21.495997 21.352192 21.416409 18.719249 24.470366	4.636378 4.620843 4.627787 4.326575 4.946753	0.269600 0.365100 4.053400 9.343000	0.044977 0.058980 1.094189 1.092908	0.002412 0.004430 1.018508 1.740104	1.004318 0.989703 0.446687 0.921097
111	TRUE INITIAL INTERPOL. KRIGING SPLINES	-0.714126 -0.717209 -0.714139 -0.695466 -0.694970	21.495997 21.352192 21.494502 24.594553 24.164441	4.636378 4.620843 4.636216 4.959290 4.915734	0.267600 0.383100 9.004700 8.407400	0.044977 0.053401 1.113822 1.006849	0.002412 0.003729 2.452523 1.309989	1.00431B 1.000019 0.973870 0.973176
IV	TRUE INITIAL INTERPOL. KRIGING SPLINES	-0.714126 -0.717209 -0.713978 -0.695793 -0.699783	21.495997 21.352192 21.502128 24.834116 26.184839	4.636378 4.620843 4.637039 4.983384 5.117112	0.269600 0.491700 9.338200 13.391900	0.044977 0.074714 1.148873 1.251661	0.002412 0.008010 2.625475 2.556300	1.004318 0.999793 0.974329 0.979916
V	TRUE INITIAL INTERPOL. NRIGING SPLINES	-0.714126 -0.717209 -0.715056 -0.685712 -0.685941	21.495997 21.352192 21.522782 23.741603 24.260616	4.636378 4.620843 4.639265 4.872536 4.925507	0.269600 0.835100 7.792400 8.177600	0.044977 0.127837 0.989923 1.040444	0.002412 0.018792 1.780893 1.295714	1.004318 1.001302 0.960212 0.960533
<pre>NOTES: ===== [1] : NAXINUM ABSOLUTE VALUE RESIDUAL = MAX.(ABS(TRUE-ESTIMATED)) [2] : AVERAGE OF THE RESIDUALS=(SUMMATION OF ABSOLUTE VALUES OF THE RESIDUALS)/N [3] : VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = { SSAVR - N \$ AVR\$\$\$ [3] : VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = { SSAVR - N \$ AVR\$\$\$ [3] : VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = { SSAVR - N \$ AVR\$\$\$ [3] : VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = { SSAVR - N \$ AVR\$\$\$ [3] : VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = { SSAVR - N \$ AVR\$\$\$ [3] : VARIANCE OF THE ABSOLUTE VALUE RESIDUALS = { SSAVR - N \$ AVR\$\$\$ [4] : RATIO AVERAGE ESTIMATED VALUES/AVERAGE TRUE VALUES</pre>								

-

Table D.18. Summary of the comparison between true, initial and calibrated flows. Example F.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	NAXINUN RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]
I	TRUE INITIAL INTERPOL. KRIGING SPLINES	4.255524 4.274763 4.276528 4.269430 2.848872	31.234204 31.246031 31.246759 35.985537 43.911869	5.588757 5.589815 5.589880 5.998795 6.626603	0.838200 0.875700 9.800500 28.928800	0.084158 0.102532 1.573010 4.924292	0.009166 0.016960 2.849809 18.056179	1.004521 1.004936 1.003268 0.669453
II	TRUE INITIAL INTERPOL. KRIGING SPLINES	4.255524 4.274763 4.251584 4.256825 2.807373	31.234204 31.246031 31.305232 36.907174 32.011668	5.588757 5.589815 5.595108 6.075128 5.657885	0.838200 0.822700 11.222500 14.799500	0.084158 0.082948 1.557878 4.107767	0.009166 0.012840 3.487768 10.423112	1.004521 0.999074 1.000306 0.659701
111	TRUE INITIAL INTERPOL. KRIGING SPLINES	4.255524 4.274763 4.252107 4.281368 2.961984	31.234204 31.246031 31.164734 35.606763 37.289550	5.588757 5.589815 5.582538 5.967140 6.106517	0.838200 0.860700 7.750000 18.144300	0.084158 0.075950 1.369353 3.621890	0.009166 0.011884 2.605558 10.546634	1.004521 0.999197 1.006073 0.696033
IV	TRUE INITIAL INTERPOL. KRIGING SPLINES	4.255524 4.274763 4.252837 4.281750 2.960927	31.234204 31.246031 31.158797 35.596979 37.272008	5.588757 5.589815 5.582007 5.966320 6.105081	0.838200 0.865700 7.757700 18.150000	0.084158 0.059032 1.354954 3.628024	0.009166 0.009102 2.608005 10.512316	1.004521 0.999369 1.006163 0.695784
V	TRUE INITIAL INTERPOL. IRIGING SPLINES	4.255524 4.274763 4.248747 4.275504 2.961483	31.234204 31.246031 31.307175 36.298322 37.480801	5.588757 5.589815 5.595282 6.024809 6.122157	0.838200 1.871000 7.906800 18.416000	0.084158 0.195639 1.393187 3.614879	0.009166 0.109416 2.585588 10.534920	1.004521 0.998407 1.004695 0.695915
<pre>JTES: 1] : MAXIMUM ABSOLUTE VALUE RESIDUAL = MAX.(ABS(TRUE-ESTIMATED)) 2] : AVERAGE OF THE RESIDUALS=(SUMMATION OF ABSOLUTE VALUES OF THE RESIDUALS)/N 3] : VARIANCE OF THE ARSOLUTE VALUE OF THE RESIDUALS = (SSAVR - N & AVR**2) / (N-2) SSAVR : SUMMATION SQUARES ABSOLUTE VALUES RESIDUALS AVR : AVERAGE ABSOLUTE VALUE RESIDUALS WITH A AVERAGE ABSOLUTE VALUE RESIDUALS</pre>								
4]:R	N : NU ATIO AVERAGE	ESTIMATED	VALUES/AVERA	IGE TRUE VA	LUES			

		PERF	ORMAN	CE I	NDEX
CASE	PROCEDURE	AVERAGE	VARIANCE	MAXIMUM RESID.	VARIANCE RESID.
I	INTERPOL.	-0.004	-0.244	-0.136	0.278
	KRIGING -	184.733	-177.615	-170.771	-999.999
	SPLINES	-72.902	-40.567	-67.387	-999.999
II	INTERPOL.	0.773	0.901	0.307	0.817
	KRIGING	-2.597	-142.343	-16.335	-165.491
	SPLINES	-29.806	-162.665	-68.515	-999.999
III	INTERPOL.	0.661	0.877	0.616	0.917
	KRIGING	1.000	-76.926	-15.565	-148.833
	SPLINES	-52.206	-155.384	-60.444	-999.999
IV	INTERPOL.	0.124	0.975	-10.655	-84.968
	KRIGING	0.999	-81.717	-15.290	-206.090
	SPLINES	-63.255	-142.358	-34.568	-895.254
v	INTERPOL.	0.937	0.978	0.376	0.662
	KRIGING -	999.999	-999.999 -	108.782	-999.999
	SPLINES -	999.999	-573.619 -	252.898	-999.999
NOTE:	A VALUE OF -9	99.999 I	NDICATES T	HAT THE	INDEX
	IS ACTUALLY L	ESS THAN	OR EQUAL	TO -999.	999

Table D.19. Summary of the performance of the calibrated flows. Example A.

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.498 -237.066 -243.633	0.697 -295.720 -214.918	-1.898 -65.349 -96.763	-16.459 -504.082 -979.050

	FREQUEN	CY OF IMPR	OVEMENT ((%)
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	80.000 40.000 0.000	80.000 0.000 0.000	60.000 0.000 0.000	80.000 0.000 0.000

	+ p - g	PERFORMAN	CE I	NDEX
CASE	PROCEDURE	AVERAGE VARIANCE	MAXIMUM RESID.	VARIANCE RESID.
I	INTERPOL. KRIGING SPLINES	$\begin{array}{rrrr} 0.008 & -0.050 \\ 1.000 & 0.920 \\ -42.610 & -14.722 \end{array}$	0.005 -0.154 -26.108	-0.003 0.478 -551.471
II	INTERPOL.	0.866 0.894	0.022	0.306
	KRIGING	0.891 -0.129	-0.404	0.221
	SPLINES	-1.351 -117.440	-19.608	-379.195
III	INTERPOL.	0.994 0.940	0.022	0.081
	KRIGING	0.984 -0.104	0.376	0.655
	SPLINES	-0.861 -158.095	-22.371	-397.106
IV	INTERPOL.	0.923 0.999	0.005	0.157
	KRIGING	0.999 -1.145	-7.038	-64.831
	SPLINES	-7.648 -472.515	-69.540	-999.999
V	INTERPOL.	0.548 0.956	-16.042	-129.818
	KRIGING -	-999.999 -999.999	-91.807	-999.999
	SPLINES -	-999.999 -999.999 -	999.999	-999.999
NOTE:	A VALUE OF -9	999.999 INDICATES T	'HAT THE	INDEX
	IS ACTUALLY I	LESS THAN OR EQUAL	TO -999.	999

Table D.20. Summary of the performance of the calibrated flows. Example B.

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.668 -199.225 -210.494	0.748 -200.091 -352.554	-3.198 -19.805 -227.525	-25.855 -212.695 -665.554

	FREQUENC	CY OF I	MPROVEMENT	(%)
PROCEDURE	AVERAGE	VARIAN	CE MAX.RES	5 VAR. RES
INTERPOL. KRIGING SPLINES	100.000 80.000 0.000	80.00 20.00 0.00	00 80.000 00 20.000 00 0.000	0 60.000 0 60.000 0 0.000

		PERFORMAN	CE I	NDEX
CASE	PROCEDURE	AVERAGE VARIANCE	MAXIMUM RESID.	VARIANCE RESID.
I	INTERPOL.	-0.023 0.771	-2.638	-18.641
	KRIGING	-31.332 -29.866	-6.784	-73.242
	SPLINES	-140.271 -148.540	-33.078	-999.999
II	INTERPOL.	-12.538 -77.063	-2.586	-15.337
	KRIGING	-0.931 -278.370	-3.860	-24.126
	SPLINES	-46.695 -533.205	-25.706	-999.999
III	INTERPOL.	0.900 -1.074	-2.870	-28.460
	KRIGING	0.996 -243.327	-6.084	-50.106
	SPLINES	-116.581 -535.163	-23.135	-999.999
IV	INTERPOL.	0.967 1.000	0.798	0.974
	KRIGING	0.991 -6.542	-0.007	0.109
	SPLINES	-5.318 -17.798	-2.484	-7.726
v	INTERPOL.	0.988 0.994	0.275	0.605
	KRIGING	0.342 0.033	0.201	-0.116
	SPLINES	0.091 0.811	-2.654	-11.237
NOTE:	A VALUE OF	999.999 INDICATES T LESS THAN OR EQUAL	HAT THE TO -999.	INDEX 999

Table D.21. Summary of the performance of the calibrated flows. Example C.

}	AVERAGE	INDEX
PROCEDURE	AVERAGE VARIANCE	MAX. RES VAR. RES
INTERPOL. KRIGING SPLINES	-1.941 -15.074 -5.987 -111.615 -61.755 -246.779	-1.404 -12.172 -3.307 -29.496 -17.411 -603.792

	FREQUENC	CY OF IMPH	ROVEMENT	(%)
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	60.000 60.000 20.000	60.000 20.000 20.000	40.000 20.000 0.000	40.000 20.000 0.000

Table D.22. Summary of the performance of the calibrated flows. Example D.

		PERF	ORMANCE I	NDEX
CASE	PROCEDURE	AVERAGE	VARIANCE MAXIMUM RESID.	VARIANCE RESID.
I	INTERPOL.	-19.250	-999.999 -52.814	-999.999
	KRIGING	-0.563	-999.999 -936.003	-999.999
	SPLINES	0.437	-999.999 -364.938	-999.999
II	INTERPOL.	-999.999	-999.999 -66.950	-999.999
	KRIGING	-999.999	-999.999 -999.999	-999.999
	SPLINES	-999.999	-999.999 -346.077	-999.999
III	INTERPOL.	-999.999	-999.999 -39.519	-745.078
	KRIGING	-999.999	-999.999 -999.999	-999.999
	SPLINES	-999.999	-999.999 -575.257	-999.999
IV	INTERPOL.	-999.999	-999.999 -54.663	-999.999
	KRIGING	-999.999	-999.999 -999.999	-999.999
	SPLINES	-999.999	-999.999 -582.368	-999.999
v	INTERPOL.	-164.766	-999.999 -349.430	-999.999
	KRIGING	-999.999	-999.999 -999.999	-999.999
	SPLINES	-999.999	-999.999 -575.688	-999.999
NOTE:	A VALUE OF -	999.999]	INDICATES THAT THE	INDEX
	IS ACTUALLY	LESS THAN	VOR EQUAL TO -999.	999

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR. RES
INTERPOL. KRIGING SPLINES	-636.803 -800.112 -799.912	-999.999 -999.999 -999.999	-112.675 -987.200 -488.865	-949.015 -999.999 -999.999

	FREQUENC	CY OF IMPR	OVEMENT	(%)
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.000 0.000 20.000	0.000 0.000 0.000	0.000 0.000 0.000	0.000 0.000 0.000

Table	D.	23.	Summary	of	the	performance	of	the	calibrated	flows.
			Example	Ε.						

		PERFORMANCE INDEX
CASE	PROCEDURE	AVERAGE VARIANCE MAXIMUM VARIANCE RESID. RESID.
I	INTERPOL. KRIGING SPLINES	0.453 0.268 -1.701 -2.970 -999.999 -445.618 -271.042 -999.999 -459.738 -247.806 -395.403 -999.999
II	INTERPOL. KRIGING SPLINES	-4.690 0.694 -0.834 -2.373 -999.999 -371.842 -225.047 -999.999 -333.037 -426.801 -999.999 -999.999
III	INTERPOL. KRIGING SPLINES	1.000 1.000 -1.019 -1.390 -35.633 -463.270 -999.999 -999.999 -37.607 -343.325 -971.949 -999.999
IV	INTERPOL. KRIGING SPLINES	0.998 0.998 -2.326 -10.028 -34.361 -537.835 -999.999 -999.999 -20.644 -999.999 -999.999 -999.999
v	INTERPOL. KRIGING SPLINES	0.909 0.965 -8.595 -59.700 -83.941 -242.848 -834.416 -999.999 -82.577 -368.592 -919.051 -999.999
NOTE:	A VALUE OF - IS ACTUALLY	-999.999 INDICATES THAT THE INDEX LESS THAN OR EQUAL TO -999.999

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR. RES
INTERPOL. KRIGING SPLINES	-0.266 -430.787 -186.721	0.785 -412.283 -477.304	-2.895 -666.101 -857.280	-15.292 -999.999 -999.999

	FREQUEN	CY OF IMPR	OVEMENT ((%)
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR. RES
INTERPOL. KRIGING SPLINES	80.000 0.000 0.000	100.000 0.000 0.000	0.000 0.000 0.000	0.000 0.000 0.000

•

Table D.24. Summary of the performance of the calibrated flows. Example F.

•

.

		PERFORMANCE INDEX
CASE	PROCEDURE	AVERAGE VARIANCE MAXIMUM VARIANCE RESID. RESID.
I	INTERPOL. KRIGING SPLINES	-0.192 -0.127 -0.091 -2.424 0.478 -999.999 -135.710 -999.999 -999.999 -999.999 -999.999 -999.999
II	INTERPOL. KRIGING SPLINES	0.958 -35.067 0.037 -0.962 0.995 -999.999 -178.260 -999.999 -999.999 -999.999 -310.745 -999.999
III	INTERPOL. KRIGING SPLINES	0.968 -33.502 -0.054 -0.681 -0.804 -999.999 -84.489 -999.999 -999.999 -999.999 -467.582 -999.999
IV	INTERPOL. KRIGING SPLINES	0.980 -39.651 -0.067 0.014 -0.858 -999.999 -84.659 -999.999 -999.999 -999.999 -467.876 -999.999
v	INTERPOL. KRIGING SPLINES	0.876 -37.067 -3.983 -141.496 -0.079 -999.999 -87.983 -999.999 -999.999 -999.999 -481.720 -999.999
NOTE:	A VALUE OF IS ACTUALLY	-999.999 INDICATES THAT THE INDEX LESS THAN OR EQUAL TO -999.999

.

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.718 -0.054 -999.999	-29.083 -999.999 -999.999	-0.832 -114.220 -545.584	-29.110 -999.999 -999.999

FREQUENCY OF IMPROVEMENT (%)							
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR. RES			
INTERPOL.	80.000	0.000	20.000	20.000			
SPLINES	40.000	0.000	0.000	0.000			

Table D.25. Summary of the comparison between true, initial and calibrated C's. Example A.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	MAXIMUM RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]
 I	TRUE INITIAL INTERPOL. KRIGING	122.500000 122.890938 119.379625 116.817375	566.666667 596.145498 453.853417 452.467591	23.804761 24.416091 21.303836 21.271284	6.221000 17.267000 17.267000	2.022688 5.442625 11.506250	2.243039 37.148887 21.444897	1.003191 0.974528 0.953570
	SPLINES	111.619937	418.452325	20.456107	17.267000	11.454188	12.881961	0.911183
II	TRUE INITIAL INTERPOL. KRIGING SPLINES	122.500000 122.890938 119.185125 116.401375 111.619937	566.666667 596.145498 459.350265 477.925369 418.452325	23.804761 24.416091 21.432458 21.861504 20.456107	6.221000 17.267000 17.267000 17.267000	2.022688 5.537625 11.809750 11.454188	2.243039 36.315976 15.638503 12.881961	1.003191 0.972940 0.950215 0.911183
111	TRUE INITIAL INTERPOL. KRIGING SPLINES	122.500000 122.890938 119.301563 116.396250 111.619937	566.666667 596.145498 456.534883 478.276939 418.452325	23.804761 24.416091 21.366677 21.869544 20.456107	6.221000 17.267000 17.267000 17.267000 17.267000	2.022688 5.457187 11.814875 11.454188	2.243039 36.825283 15.567570 12.881961	1.003191 0.973890 0.950173 0.911183
IV	TRUE INITIAL INTERPOL. KRIGING SPLINES	122.500000 123.736213 122.858750 117.225062 112.284687	566.666667 702.925795 788.387341 578.875420 494.513023	23.804761 26.512748 28.078236 24.059830 22.237649	19.671900 18.461000 22.372000 22.372000	6.396437 10.110000 12.627313 11.053438	22.428561 29.841839 36.986530 32.087600	1.010092 1.002929 0.956939 0.916610
V	TRUE INITIAL INTERPOL. KRIGING SPLINES	122.500000 122.890938 124.703875 116.555125 111.619937	566.6666657 596.145498 923.377391 467.769095 418.452325	23.804761 24.416091 30.387125 21.627970 20.456107	6.221000 18.461000 17.267000 17.267000	2.022688 8.890500 11.656000 11.454188	2.243039 37.744147 18.185262 12.981961	1.003191 1.017991 0.951470 0.911183
NOTES: ===== [1] : MAXIMUM ABSOLUTE VALUE RESIDUAL = MAX.(ABS(TRUE-ESTIMATED)) [2] : AVERAGE OF THE RESIDUALS=(SUMMATION OF ABSOLUTE VALUES OF THE RESIDUALS)/N [3] : VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = { SSAVR - N & AVR**2 } / (N-2) SSAVR : SUMMATION SQUARES ABSOLUTE VALUES RESIDUALS AVR - AUERAGE ABSOLUTE VALUE DESIDUALS								
[4];	N : RATIO AVERAG	NUMBER OF D DE ESTIMATED	ATA POINTS VALUES/AVE	RAGE TRUE V	ALUES			

Table D.26. Summary of the comparison between true, initial and calibrated C's. Example B.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	MAXINUH RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]	
I	TRUE INITIAL INTERPOL. NRIGING SPLINES	122.500000 122.890938 119.657187 122.622563 113.200125	565.666667 596.145498 424.691198 675.892952 465.662323	23.804761 24.416091 20.608037 25.997941 21.579210	6.221000 17.267000 18.461000 17.267000	2.022688 5.532688 12.730688 11.631000	2.243039 31.803443 14.750940 13.326076	1.003191 0.976793 1.001001 0.924083	
II	TRUE INITIAL INTERPOL. KRIGING SPLINES	122.50000 122.890938 118.330875 122.622563 113.200125	566.666667 576.145498 418.961833 675.892952 465.662323	23.804761 24.416091 20.468557 25.997941 21.579210	6.221000 17.267000 18.461000 17.267000	2.022688 6.111125 12.730688 11.631000	2.243039 36.077826 14.750940 13.326076	1.003191 0.965966 1.001001 0.924083	
III	TRUE INITIAL INTERPOL. KRIGING SPLINES	122.500000 122.890938 118.872063 122.622563 113.200125	566.666667 596.145498 436.324698 675.892952 465.662323	23.804761 24.416091 20.888387 25.997941 21.579210	6.221000 17.267000 18.461000 17.267000	2.022688 5.567937 12.730688 11.631000	2.243039 31.063965 14.750940 13.326076	1.003191 0.970384 1.001001 0.924083	
IV	TRUE INITIAL INTERPOL. KRIGING SPLINES	122.500000 123.736213 118.616813 127.067375 113.903062	566.666667 702.925795 360.368703 935.437515 551.062498	23.804761 26.512748 18.983380 30.584923 23.474720	19.671900 17.267000 32.639000 22.372000	6.396437 8.423312 13.667000 11.612562	22.428561 30.098341 65.855510 34.150292	1.010092 0.968301 1.037285 0.929821	
¥	TRUE INITIAL INTERPOL. KRIGING SPLINES	122.500000 122.890938 120.268812 124.253625 114.718687	566.666667 596.145498 574.629776 700.411627 490.040616	23.804761 24.416091 23.971437 26.465291 22.136861	6.221000 17.843000 18.461000 17.267000	2.022688 10.084187 12.791250 11.816563	2.243039 26.602814 14.790742 13.506217	1.003191 0.981786 1.014315 0.936479	
NOTES: ===== [1] : MAXIMUM ABSOLUTE VALUE RESIDUAL = MAX.(ABS(TRUE-ESTIMATED)) [2] : AVERAGE OF THE RESIDUALS=(SUMMATION OF ABSOLUTE VALUES OF THE RESIDUALS)/N [3] : VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = (SSAVR - N # AVR##2) / (N-2) SSAVR : SUMMATION SQUARES ABSOLUTE VALUES RESIDUALS AVR : AVERAGE ABSOLUTE VALUE RESIDUALS N : NUMBER OF DATA POINTS [4] : PATIO AVERAGE ESTIMATED VALUES/AVERAGE TRUE VALUES									

Table D.27. Summary of the comparison between true, initial and calibrated C's. Example C.

PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	NAXINUM	AVERAGE	VARIANCE	RATIO AVERAGES
				RES.[1]	RES.[2]	RES. [3]	EST/TRUE [4]
TDIIC	199 500000	 511 111117					
TNITIA	122.300000	300,000007	23.004/01	74 104000	21 710012	10 511001	1 201700
TNICODAL	14/.217V02	000,041000 900 301513	21,3313/0	34,004000 AL 120000	17 007470	40.364271	1.201/00
INTERFOL.	130.47/43/	700.301313	31.3V7/0/ 37 105173	40.02VVVV	13,77/438	138.883330	1.114203
COL THEC	13/ 4/2023	134.070323 LLZ 005880	27,103172 25 75044L	15 519000	19.9/2023	112.0/8718	1.1////8
3FLIME3	132.0//10/		23./30440	13.318000	10.3//18/	10.626710	1.084/12
TRUE	122.500000	566.666667	23.804761				
INITIAL	147.219062	860.341005	29.331570	34.684000	24.719062	40.564291	1.201788
INTERPOL.	134.497375	874,595216	29.573556	46.506000	11.997375	103.276303	1.097938
KRIGING	137.472688	734.686817	27.105107	35.930000	14.972688	112.107661	1.122226
SPLINES	132.877187	663.085449	25.750446	15,518000	10.377187	10.626910	1.084712
TRUE	122.500000	566.666667	23.804761				
INITIAL	147.219062	860.341005	29.331570	34.684000	24.719062	40.564291	1.201788
INTERPOL.	136.498625	980.540824	31.313598	46.634000	13.998625	158.953740	1.114274
LRIGING	137.380813	730.457963	27.026986	35.959000	14.880813	108.113889	1.121476
SPLINES	132.877187	663.085449	25.750446	15.518000	10.377187	10.626910	1.084712
TRUE	122.500000	566.666667	23,804761				
INITIAL	148.4835001	012.213738	31.815307	45.606000	25.983500	140.251862	1.212110
INTERPOL.	137.3849381	041.147562	32.266818	47.007000	14.884938	158,527745	1.121510
KRIGING	138.228250	799.584177	28.276920	38.470000	15,728250	163.225904	1.128394
SPLINES	134.533500	729.374341	27.006931	30.045000	12.033500	69.358936	1.098233
TDIIC	177 500000	 511 61117					
INITIAL	147.219062	860.341005	29.331570	34.684000	74.719062	40.564291	1 201788
INTERPOL	138.0623121	208.461234	34.762929	54.153000	15.562312	217.351680	1,127039
KRIGING	138.420125	757.052799	27.514592	41.867000	15.920125	151,195141	1.129960
SPL INES	133.010000	652.365272	25.541441	15.518000	10.510000	10.524881	1.085796
AXIMUM ABSO	LUTE VALUE	RESIDUAL =	NAX.(ABS(T	RUE-ESTINATI	ED)}		
VERAGE OF T	HE RESIDUAL	S=(SUMMATIC	N OF ABSOLU	TE VALUES OF	F THE RESID	UALS)/N	
ARIANCE OF	THE ARGOLUT	E VALUE OF	THE RESIDUAL	LS = { SSAVI	R - N I AVR	##2 } / {N-2	2)
SSAVR :	SUMMATION S	QUARES ABGO	UTE VALUES	RESIDUALS			
AVR :	AVERAGE ABS	OLUTE VALUE	RESIDUALS				
N :	NUMBER OF D	ATA POINTS					
	PROCEDURE TRUE INITIAL INTERPOL. KRIGING SPLINES TRUE INITIAL INTERPOL. KRIGING SPLINES TRUE INITIAL INTERPOL. KRIGING SPLINES TRUE INITIAL INTERPOL. KRIGING SPLINES TRUE INITIAL INTERPOL. KRIGING SPLINES TRUE INITIAL INTERPOL. KRIGING SPLINES TRUE INITIAL INTERPOL. KRIGING SPLINES TRUE INITIAL INTERPOL. KRIGING SPLINES TRUE INITIAL INTERPOL. KRIGING SPLINES TRUE INITIAL INTERPOL. KRIGING SPLINES TRUE INITIAL INTERPOL. KRIGING SPLINES	PROCEDURE AVERAGE TRUE 122.500000 INITIAL 147.219062 INTERPOL. 136.497437 KRIGING 137.472625 SPLINES 132.877187 TRUE 122.500000 INITIAL 147.219062 INTERPOL. 134.497375 KRIGING 137.472668 SPLINES 132.877187 TRUE 122.500000 INITIAL 147.219062 INTERPOL. 134.497375 KRIGING 137.472668 SPLINES 132.877187 TRUE 122.500000 INITIAL 147.219062 INTERPOL. 136.498625 KRIGING 137.380813 SPLINES 132.877187 TRUE 122.500000 INITIAL 148.4835001 INTERPOL. 137.3849381 KRIGING 138.228250 SPLINES 134.533500 TRUE 122.500000 INITIAL 147.219062 INTERPOL.	PROCEDURE AVERAGE VARIANCE TRUE 122.500000 566.666667 INITIAL 147.219062 860.341005 INTERPOL. 136.497437 980.301513 KRIGING 137.472625 734.690323 SPLINES 132.877187 663.085449 TRUE 122.500000 566.666667 INITIAL 147.219062 860.341005 INTERPOL. 134.497375 874.595216 KRIGING 137.472668 734.686817 SPLINES 132.877187 663.085449 TRUE 122.500000 566.666667 INITIAL 147.219062 860.341005 INTERPOL. 136.498625 980.540824 KRIGING 137.380813 730.457963 SPLINES 132.877187 663.085449 TRUE 122.500000 566.666667 INITIAL 147.219062 860.341005 INTERPOL. 137.3849381041.147562 KRIGING KRIGING 138.228250 799.584177 SPLINES <	PROCEDURE AVERAGE VARIANCE STAN.DEV. TRUE 122.500000 566.666667 23.804761 INITIAL 147.219062 860.341005 29.331570 INTERPOL. 136.497437 980.301513 31.309767 KRIGING 137.472625 734.690323 27.105172 SPLINES 132.877187 663.085449 25.750446 TRUE 122.500000 566.666667 23.804761 INITIAL 147.219062 860.341005 29.331570 INTERPOL. 134.497375 874.595216 29.573556 KRIGING 137.472608 734.686817 27.105107 SPLINES 132.877187 663.085449 25.750446 TRUE 122.500000 566.666667 23.804761 INITIAL 147.219062 860.341005 29.331570 INTERPOL. 136.492625 980.540824 31.313588 KRIGING 137.380913 730.457963 27.026986 SPLINES 132.877187 663.085449 25.750446 <	PROCEDURE AVERAGE VARIANCE STAN.DEV. MAXIMUM RES.[1] TRUE 122.500000 566.666667 23.804761 INITIAL 147.219062 860.341005 29.331570 34.664000 INTERPOL. 136.497437 980.301513 31.309767 46.620000 KRIGING 137.472625 734.690323 27.105172 35.814000 SPLINES 132.877187 663.085449 25.750446 15.518000 TRUE 122.500000 566.666667 23.804761 11111AL 147.219062 860.341005 29.331570 34.684000 INTERPOL. 134.497375 874.595216 29.573556 46.56600 KRIGING 137.472668 734.686817 27.105107 35.930000 SPLINES 132.877187 663.085449 25.750446 15.518000 TRUE 122.500000 566.666667 23.804761 11111AL 147.219062 860.341005 29.331570 34.684000 INTERPOL. 136.498625 980.54082 25.750446 15.518000 TRUE <td>PROCEDURE AVERAGE VARIANCE STAN.DEV. MAXIMUM RES.[1] AVERAGE RES.[1] TRUE 122.500000 566.66667 23.804761 INITIAL 147.219062 860.341005 29.331570 34.684000 24.719062 INTERPOL 156.497437 780.301513 31.309767 46.620000 13.997438 KRIGING 137.472625 734.690323 27.105172 35.814000 14.972625 SPLINES 132.877187 663.085449 25.750446 15.518000 10.377187 TRUE 122.500000 566.66667 23.804761 1 11.977375 TRUE 122.500000 566.66667 23.804761 1.977375 KRIGING 137.472688 734.68817 27.105107 35.930000 14.972688 SPLINES 132.877187 663.085449 25.750446 15.518000 10.377187 TRUE 122.500000 566.666667 23.804761 11.111AL 147.219062 860.341005 29.331570 34.684000 24.719062 INTERPOL.</td> <td>PROCEDURE AVERAGE VARIANCE STAN.DEV. MAXIMUM AVERAGE VARIANCE RES.[1] RES.[2] RES.[3] RES.[1] RES.[2] RES.[3] INUE 122.500000 566.666667 23.804761 RES.[1] RES.[2] RES.[3] INTERPOL. 136.497437 960.301513 31.30767 46.620000 13.997438 158.88550 KRIGING 137.472625 734.690323 27.105172 35.814000 14.972625 112.078918 SPLINES 132.877187 663.085449 25.750446 15.518000 10.377187 10.626910 TRUE 122.500000 566.666667 23.804761 11.111AL 147.219062 460.341005 29.331570 34.684000 24.719062 40.564291 INTERPOL. 134.497375 874.595216 29.573556 46.590000 14.972688 112.107661 SPLINES 132.877187 663.085449 25.750446 15.518000 13.97962 40.564291 INTERPOL. 136.492625 980.540624 31.315588</td>	PROCEDURE AVERAGE VARIANCE STAN.DEV. MAXIMUM RES.[1] AVERAGE RES.[1] TRUE 122.500000 566.66667 23.804761 INITIAL 147.219062 860.341005 29.331570 34.684000 24.719062 INTERPOL 156.497437 780.301513 31.309767 46.620000 13.997438 KRIGING 137.472625 734.690323 27.105172 35.814000 14.972625 SPLINES 132.877187 663.085449 25.750446 15.518000 10.377187 TRUE 122.500000 566.66667 23.804761 1 11.977375 TRUE 122.500000 566.66667 23.804761 1.977375 KRIGING 137.472688 734.68817 27.105107 35.930000 14.972688 SPLINES 132.877187 663.085449 25.750446 15.518000 10.377187 TRUE 122.500000 566.666667 23.804761 11.111AL 147.219062 860.341005 29.331570 34.684000 24.719062 INTERPOL.	PROCEDURE AVERAGE VARIANCE STAN.DEV. MAXIMUM AVERAGE VARIANCE RES.[1] RES.[2] RES.[3] RES.[1] RES.[2] RES.[3] INUE 122.500000 566.666667 23.804761 RES.[1] RES.[2] RES.[3] INTERPOL. 136.497437 960.301513 31.30767 46.620000 13.997438 158.88550 KRIGING 137.472625 734.690323 27.105172 35.814000 14.972625 112.078918 SPLINES 132.877187 663.085449 25.750446 15.518000 10.377187 10.626910 TRUE 122.500000 566.666667 23.804761 11.111AL 147.219062 460.341005 29.331570 34.684000 24.719062 40.564291 INTERPOL. 134.497375 874.595216 29.573556 46.590000 14.972688 112.107661 SPLINES 132.877187 663.085449 25.750446 15.518000 13.97962 40.564291 INTERPOL. 136.492625 980.540624 31.315588

Table D.28. Summary of the comparison between true, initial and calibrated C's. Example D.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	HAXINUM	AVERAGE	VARIANCE	RATIO AVERAGES
					RES.[1]	RES.[2]	RES. [3]	EST/TRUE [4]
	TRUE	140.000000	0.000000	0.000000				
	INITIAL	139.958833	7.499320	2.738489	8.522000	2.237533	2.480356	0.999706
I	INTERPOL.	139,192178	86.478615	9.299388	18.792000	7.526600	30.338138	0.994230
	KRIGING	136.527167	181.567801	13.474710	20.412000	13.389444	13.492355	0.975194
	SPLINES	137.622156	190.145060	13.789310	20.412000	13.640400	8.779885	0.983015
	TRUE	140.000000	0.000000	0.000000				
	INITIAL	139.958833	7.499320	2.738489	8.522000	2.237533	2.480356	0.999706
П	INTERPOL.	138.995589	95.519798	9.773423	20.379000	7.970922	32.827118	0.992826
	KRIGING	136.333611	184.074020	13.567388	20.412000	13.601756	11.615097	0.973812
	SPLINES	136.728044	180.292340	13.427298	20.412000	13.450600	9.179766	0.976629
	TRUE	140.000000	0.000000	0.000000				
	INITIAL	139,958833	7.499320	2.738489	8.522000	2.237533	2.480356	0.999706
Ш	INTERPOL.	139.758500	75.524815	8.690501	18.792000	6.708589	30.497248	0.998275
	KRIGING	136.327722	184.567057	13.585546	20.412000	13.657544	10.616752	0.973769
	SPL INES	135.940489	182.494425	13.509050	21.670000	13.761400	8.680518	0.971003
	TRUE	140.000000	0.000000	0.000000				
	INITIAL	139.941799	14.998452	3.872783	12.052300	3.164327	4.960670	0.999584
I۷	INTERPOL.	141.093283	139.495700	11.810872	20.018000	10,585194	28.183798	1.007809
	KRIGING	136.465333	193.623291	13.914859	23.067000	13.705056	17.406331	0.974752
	SPLINES	135.895606	189.202131	13.755077	24.847000	13.765272	15.688661	0.970683
	TRUE	140.000000	0.000000	0.000000				
	INITIAL	139.958833	7.499320	2.738489	8.522000	2.237533	2.480356	0.999706
٧	INTERPOL.	140.077217	158.886754	12.605029	21.670000	11.411817	28.092593	1.000552
	KRIGING	136.113194	184.420192	13.580140	20.412000	13.702972	10.852023	0.972237
	SPLINES	136.254539	183,084631	13.530877	20.412000	13.649317	9.902152	0.973247
OTES:								
====								
1] : M	AXIMUM ABSO	LUTE VALUE	RESIDUAL =	MAX. (ABS (T	RUE-ESTIMAT	ED)}		
2]:A	VERAGE OF T	HE RESIDUAL	S=(SUMMATIO	N OF ABSOLU	TE VALUES O	F THE RESID	JALS)/N	
3] : Y	ARIANCE OF	THE ABSOLUT	E VALUE OF	THE RESIDUA	LS = { SSAV	R – N 🖡 AVR	\$\$2 } / {N-3	2)
	SSAVR :	SUMMATION S	IQUARES ABSO	LUTE VALUES	RESIDUALS			
	AVR :	AVERAGE ABS	SOLUTE VALUE	RESIDUALS				
	N :	NUMBER OF D	ATA POINTS					
4] : R	ATIO AVERAG	E ESTIMATEI	VALUES/AVE	RAGE TRUE V	ALUES			

Table D.29. Summary of the comparison between true, initial and calibrated C's. Example E.

				and the second se				and the second se	
CAS	E PROCEDUR	E AVERAGE	VARIANCE	STAN.DEV.	. MAXIMUN RES.[1]	AVERAGE RES.[2]	VARJANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]	
I	TRUE INITIAL INTERPOL KRIGING SPLINES	140.000000 139.958833 139.959117 133.260300 134.701350	0.000000 7.499320 85.458171 167.756177 166.904531	0.000000 2.738489 9.244359 12.952072 12.919154	8.522000 19.986000 21.670000 20.412000	2.237533 7.114017 14.315844 13.577439	2,480356 34,762087 7,386423 9,815186	0.999706 0.999708 0.951859 0.962153	
II	TRUE INITIAL INTERPOL. KRIGING SPLINES	140.000000 139.958833 139.964678 133.278372 134.189333	0.000000 7.499320 83.243194 161.905021 163.253766	0.000000 2.738489 9.123771 12.724190 12.777080	8.522000 19.986000 21.670000 20.412000	2.237533 6.903667 14.057872 13.731356	2.480356 35.515988 8.658260 7.645472	0.999706 0.999748 0.951988 0.958495	
111	TRUE INITIAL INTERPOL. KRIGING SPLINES	140.000000 139.958833 140.051139 132.354694 134.730900	0.000000 7.499320 84.346600 147.510161 167.586697	0.000000 2.738489 9.184040 12.145376 12.945528	8.522000 19.986000 21.670000 21.670000	2.237533 7.002450 14.047183 13.562789	2.480356 35.237848 7.905840 10.587471	0.999706 1.000365 0.945391 0.962364	
IV	TRUE INITIAL INTERPOL. KRIGING SPLINES	140.000000 139.941799 139.628605 132.401378 134.281450	0.000000 14.998452 104.934688 156.283385 171.450543	0.000000 3.872783 10.243763 12.501335 13.093912	12.052300 19.986000 24.847000 24.847000	3.164327 8.321639 14.074833 13.686539	4.960670 35.635932 15.222409 16.056918	0.999584 0.997347 0.945724 0.959153	
v	TRUE INITIAL INTERPOL. KRIGING SPLINES	140.000000 139.958833 137.235400 132.666867 134.993917	0.000000 7.499320 118.755303 152.716219 171.974287	0.000000 2.738489 10.897491 12.357840 13.113897	8.522000 20.298000 21.670000 21.670000	2.237533 9.453600 14.022656 13.682317	2.480356 36.776641 9.108982 8.973663	0.999706 0.980253 0.947620 0.964242	
(OTES:	IOTES:								
 1]: MAXIMUM ABSOLUTE VALUE RESIDUAL = MAX.{ABS(TRUE-ESTIMATED)} 2]: AVERAGE OF THE RESIDUALS={SUMMATION OF ABSOLUTE VALUES OF THE RESIDUALS)/N 3]: VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = { SSAVR - N ‡ AVR‡‡2 } / (N-2) SSAVR : SUMMATION SQUARES ABSOLUTE VALUES RESIDUALS AVR : AVERAGE ABSOLUTE VALUE RESIDUALS N = NUMPER OF DATA POINTS 									
4]:R	ATID AVERAGE	ESTIMATED	VALUES/AVER	AGE TRUE VA	LUES				

Table D.30. Summary of the comparison between true, initial and calibrated C's. Example F.

CASE	PROCEDURE	AVERAGE	VARIANCE	STAN.DEV.	MAXIMUM RES.[1]	AVERAGE RES.[2]	VARIANCE RES. [3]	RATIO AVERAGES EST/TRUE [4]
I	TRUE INITIAL INTERPOL. KRIGING SPLINES	140.000000 139.969736 139.866701 132.113661 133.594437	0.000000 7.491670 52.512202 118.103721 165.135903	0.000000 2.737092 7.246530 10.867554 12.850522	8.522000 20.412000 21.670000 21.670000	2.222908 5.242322 12.927937 14.063609	2.537376 25.033985 12.632965 7.519413	0.999784 0.999048 0.943669 0.954246
II	TRUE INITIAL INTERPOL. KRIGING SPLINES	140.000000 139.969736 139.530891 132.308672 133.784931	0.000000 7.491670 54.774347 119.922215 162.353426	0.000000 2.737092 7.400969 10.950900 12.741798	8.522000 20.412000 21.670000 21.670000	2.222908 5.397684 12.835730 13.870092	2.537376 25.841654 13.792095 7.757154	0.999784 D.995549 0.945062 0.955607
III	TRUE INITIAL INTERPOL. KRIGING SPLINES	140.00000 139.969736 139.469500 133.351609 134.002477	0.000000 7.491670 50.787031 141.238452 167.204936	0.000000 2.737092 7.126502 11.884379 12.930775	8.522000 20.412000 21.670000 21.670000	2,222908 5,150810 13,129943 13,913879	2.537376 24.527662 12.374684 8.718445	0.999784 0.996211 0.952511 0.957161
IV	TRUE INITIAL INTERPOL. KRIGING SPLINES	140.000000 139.957221 139.767529 133.245954 134.011236	0.000000 14.983226 59.377350 144.269356 176.107339	0.000000 3.870817 7.705670 12.011218 13.270544	12.052300 19.255000 24.847000 24.847000	3.143651 5.887805 13.060161 14.000236	5.074737 24.707900 18.704550 15.127812	0.99969 4 0.998339 0.951757 0.957223
v	TRUE INITIAL INTERPOL. KRIGING SPLINES	140.000000 139.969736 140.410747 133.654310 134.123264	0.000000 7.491670 94.910647 140.366317 171.372956	0.000000 2.737092 9.742210 11.847629 13.090949	8.522000 20.412000 21.670000 21.670000	2.222908 7.831782 12.962885 14.042563	2. 537376 33.583106 11.928107 7.820383	0.999784 1.002934 0.954674 0.958023
NOTES: 1] : MAXIMUM ABSOLUTE VALUE RESIDUAL = MAX.(ABS(TRUE-ESTIMATED)) 2] : AVERAGE OF THE RESIDUALS=(SUMMATION OF ABSOLUTE VALUES OF THE RESIDUALS)/N 3] : VARIANCE OF THE ABSOLUTE VALUE OF THE RESIDUALS = { SSAVR - N t AVRtt2 } / (N-2) SSAVR : SUMMATION SQUARES ABSOLUTE VALUES RESIDUALS AVR : AVERAGE ABSOLUTE VALUE RESIDUALS N : NUMBER OF DATA POINTS 4] : RATIO AVERAGE ESTIMATED VALUES/AVERAGE TRUE VALUES								
		PERF	ORMAN	CE I	NDEX			
--------	--------------	-----------	-------------	-------------------	--------------------			
CASE	PROCEDURE	AVERAGE	VARIANCE	MAXIMUM RESID.	VARIANCE RESID.			
I	INTERPOL.	-62.709	-13.645	-6.704	-273.295			
	KRIGING	-210.664	-14.007	-6.704	-90.406			
	SPLINES	-773.546	-24.279	-6.704	-31.983			
II	INTERPOL.	-70.898	-12.253	-6.704	-261.133			
	KRIGING	-242.359	-8.062	-6.704	-47.609			
	SPLINES	-773.546	-24.279	-6.704	-31.983			
III	INTERPOL.	-65.936	-12.957	-6.704	-268.537			
	KRIGING	-242.769	-7.990	-6.704	-47.169			
	SPLINES	-773.546	-24.279	-6.704	-31.983			
IV	INTERPOL.	0.916	-1.648	0.119	-0.770			
	KRIGING	-17.207	0.992	-0.293	-1.719			
	SPLINES	-67.284	0.720	-0.293	-1.047			
v	INTERPOL.	-30.780	-145.424	-7.806	-282.156			
	KRIGING	-230.244	-10.255	-6.704	-64.730			
	SPLINES	-773.546	-24.279	-6.704	-31.983			
NOTE :	A VALUE OF -	-999.999	INDICATES 1	HAT THE	INDEX			
	IS ACTUALLY	LESS THAN	N OR EQUAL	TO -999.	999			

Table D.31. Summary of the performance of the calibrated C's. Example A.

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR.RES
INTERPOL. KRIGING SPLINES	-45.881 -188.649 -632.293	-37.185 -7.865 -19.279	-5.560 -5.422 -5.422	-217.178 -50.327 -25.796

	FREQUENCY OF IMPROVEMENT (%)			
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR. RES
INTERPOL. KRIGING SPLINES	20.000 0.000 0.000	0.000 20.000 20.000	20.000 0.000 0.000	0.000 0.000 0.000

······		<u> </u>			
		PERF	ORMAN	CE I	NDEX
CASE	PROCEDURE	AVERAGE	VARIANCE	MAXIMUM RESID.	VARIANCE RESID.
I	INTERPOL.	-51.879	-22.196	-6.704	-200.036
	KRIGING	0.902	-12.729	-7.806	-42.248
	SPLINES	-564.898	-10.740	-6.704	-34.296
II	INTERPOL.	-112.730	-24.106	-6.704	-257.706
	KRIGING	0.902	-12.729	-7.806	-42.248
	SPLINES	-564.898	-10.740	-6.704	-34.296
III	INTERPOL.	-85.120	-18.550	-6.704	-190.796
	KRIGING	0.902	-12.729	-7.806	-42.248
	SPLINES	-564.898	-10.740	-6.704	-34.296
IV	INTERPOL.	-8.867	-1.292	0.230	-0.801
	KRIGING	-12.650	-6.325	-1.753	-7.621
	SPLINES	-47.362	0.987	-0.293	-1.318
V .	INTERPOL.	-31.573	0.927	-7.227	-139.664
	KRIGING	-19.121	-19.584	-7.806	-42.482
	SPLINES	-395.178	-5.757	-6.704	-35.257
NOTE:	A VALUE OF - IS ACTUALLY	-999.999 I LESS THAN	INDICATES	THAT THE TO -999.	INDEX 999

Table D.32. Summary of the performance of the calibrated C's. Example B.

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	-58.034 -5.813 -427.447	-13.043 -12.819 -7.398	-5.422 -6.596 -5.422	-157.801 -35.369 -27.893

	FREQUENC	CY OF IMPR	OVEMENT ((%)
PROCEDURE	AVERAGE	VARIANCE	MAX.RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.000 60.000 0.000	20.000 0.000 20.000	20.000 0.000 0.000	0.000 0.000 0.000

	فيسون والمستعيرة المتحفظ المسور فالمتحد المعا			
		PERFORM	ANCE I	NDEX
CASE	PROCEDURE	AVERAGE VARIA	NCE MAXIMUM RESID.	VARIANCE RESID.
I	INTERPOL.	0.679 -0.	984 -0.807	-14.342
	KRIGING	0.633 0.	673 -0.066	-6.634
	SPLINES	0.824 0.	892 0.800	0.931
II	INTERPOL.	0.764 -0.	099 -0.798	-5.482
	KRIGING	0.633 0.	373 -0.073	-6.638
	SPLINES	0.824 0.	892 0.800	0.931
III	INTERPOL.	0.679 -0.9	986 -0.808	-14.355
	KRIGING	0.638 0.6	589 -0.075	-6.104
	SPLINES	0.824 0.1	392 0.800	0.931
IV	INTERPOL.	0.672 -0.1	134 -0.062	-0.278
	KRIGING	0.634 0.1	727 0.288	-0.354
	SPLINES	0.786 0.8	367 0.566	0.755
v	INTERPOL.	0.604 -3.	776 -1.438	-27.710
	KRIGING	0.585 0.5	580 -0.457	-12.893
	SPLINES	0.819 0.5	915 0.800	0.933
NOTE:	A VALUE OF -	999.999 INDICA	TES THAT THE	INDEX
	IS ACTUALLY	LESS THAN OR EQ	QUAL TO -999	.999
		AVERAGE	I	NDEX
	PROCEDURE	AVERAGE VARIAN	ICE MAX.RES	VAR. RES

Table D.33. Summary of the performance of the calibrated C's. Example C.

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.680 0.625 0.815	-1.196 0.668 0.892	-0.782 -0.077 0.753	-12.433 -6.525 0.896

	FREQUENC	Y OF IMPR	OVEMENT ((%)
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	100.000 100.000 100.000	0.000 100.000 100.000	0.000 20.000 100.000	0.000 0.000 100.000

Table D.34. Summary of the performance of the calibrated C's. Example D.

		PERF	ORMAN	CE I	NDEX
CASE	PROCEDURE	AVERAGE	VARIANCE	MAXIMUM RESID.	VARIANCE RESID.
I	INTERPOL.	-384.064	-131.976	-3.863	-148.606
	KRIGING	-999.999	-585.184	-4.737	-28.590
	SPLINES	-999.999	-641.875	-4.737	-11.530
II	INTERPOL.	-594.284	-161.234	-4.719	-174.161
	KRIGING	-999.999	-601.478	-4.737	-20.929
	SPLINES	-999.999	-576.977	-4.737	-12.697
III	INTERPOL.	-33.414	-100.423	-3.863	-150.180
	KRIGING	-999.999	-604.710	-4.737	-17.321
	SPLINES	-999.999	-591.182	-5.466	-11.248
IV	INTERPOL.	-351.862	-85.504	-1.759	-31.279
	KRIGING	-999.999	-165.657	-2.663	-11.312
	SPLINES	-999.999	-158.133	-3.250	-9.002
V	INTERPOL.	-2.518	-447.881	-5.466	-127.279
	KRIGING	-999.999	-603.746	-4.737	-18.142
	SPLINES	-999.999	-595.019	-4.737	-14.938
NOTE:	A VALUE OF - IS ACTUALLY	999.999 1 LESS THAN	INDICATES T	THAT THE TO -999.	INDEX 999

	AVERA	GE	I	NDEX
PROCEDURE	AVERAGE V	ARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	-273.228 - -999.999 - -999.999 -	185.404 512.155 512.637	-3.934 -4.322 -4.585	-126.301 -19.259 -11.883

	FREQUEN	CY OF IMPR	OVEMENT	(%)
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.000 0.000 0.000	0.000 0.000 0.000	0.000 0.000 0.000	0.000 0.000 0.000

		· · · · · · · · · · · · · · · · · · ·			
		PERF	ORMAN	CE I	NDEX
CASE	PROCEDURE	AVERAGE	VARIANCE	MAXIMUM RESID.	VARIANCE RESID.
I	INTERPOL.	0.014	-128.856	-4.500	-195.419
	KRIGING	-999.999	-499.395	-5.466	-7.868
	SPLINES	-999.999	-494.328	-4.737	-14.659
II	INTERPOL.	0.264	-122.212	-4.500	-204.031
	KRIGING	-999.999	-465.098	-5.466	-11.185
	SPLINES	-999.999	-472.896	-4.737	-8.501
III	INTERPOL.	-0.543	-125.500	-4.500	-200.832
	KRIGING	-999.999	-385.901	-5.466	-9.159
	SPLINES	-999.999	-498.385	-5.466	-17.220
IV	INTERPOL.	-39.720	-47.949	-1.750	-50.605
	KRIGING	-999.999	-107.576	-3.250	-8.416
	SPLINES	-999.999	-129.673	-3.250	-9.477
v	INTERPOL.	-999.999	-249.762	-4.673	-218.845
	KRIGING	-999.999	-413.693	-5.466	-12.487
	SPLINES	-999.999	-524.876	-5.466	-12.089
NOTE:	A VALUE OF - IS ACTUALLY	-999.999] LESS THAN	INDICATES T	HAT THE TO -999.	INDEX 999

Table	D.	35.	Summary	of	the	performance	of	the	calibrated	C,	s.
			Example	Ε.							

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	-207.997 -999.999 -999.999	-134.856 -374.333 -424.031	-3.985 -5.023 -4.731	-173.947 -9.823 -12.389

	FREQUEN	CY OF IMPR	OVEMENT ((%)
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	40.000 0.000 0.000	0.000 0.000 0.000	0.000 0.000 0.000	0.000 0.000 0.000

Table D.36. Summary of the performance of the calibrated C's. Example F.

CASE	PROCEDURE	PERFORMAN AVERAGE VARIANCE	CE INDEX MAXIMUM VARIANCE
			RESID. RESID.
I	INTERPOL. KRIGING	-18.400 -48.132 -999.999 -247.525	-4.737 -96.340 -5.466 -23.788
	SPLINES	-999.999 -484.876	-5.466 -7.782
 T T	INTERPOL.	-239.267 -52.456	-4.737 -102.722
11	SPLINES	-999.999 -255.257	-5.466 -20.545 -5.466 -8.346
	INTERPOL.	-306.269 -44.957	-4.737 -92.442 -5.466 -22.795
111	SPLINES	-999.999 -497.128	-5.466 -10.806
	INTERPOL.	-28.531 -14.705	
ΙV	SPLINES	-999.999 -137.148	-3.250 -7.886
	INTERPOL.	-183.203 -159.499	
v	SPLINES	-999.999 -522.272	-5.466 -21.099 -5.466 -8.499
NOTE:	A VALUE OF -	999.999 INDICATES T	HAT THE INDEX
	15 ACTUALLY	LESS INAN OK EQUAL	10 -898.998
	[т N D 型 V

	AVER	AGE	I	NDEX
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	-155.134 -999.999 -999.999	-63.950 -259.790 -422.013	-4.100 -5.023 -5.023	-97.677 -21.760 -8.664

	FREQUENC	CY OF IMPR	OVEMENT ((%)
PROCEDURE	AVERAGE	VARIANCE	MAX. RES	VAR. RES
INTERPOL. KRIGING SPLINES	0.000 0.000 0.000	0.000 0.000 0.000	0.000 0.000 0.000	0.000 0.000 0.000