

Farag S. Fattouh, B.Sc. (Hons. El.), B.Sc. (Hons. Phys.) (Cairo U.)

TECHNIQUE FOR THE PERIODIC COMPUTATIONAL RESPONSE · 17 OF & LARGE NONLINEAR CIRCUITS 15 ð, by Farag S. Fattouh, B.Sc. (Hons. Elect.), B.Sc. (Hons. Phys.), (Cairo U.) A thesis submitted to the Faculty of Graduate Studies and Research in partial fulfillment of the requirements for the degree of Master of Engineering. Department of Electrical Engineering McGill University, Montreal, Canada. December, 1977. Farag S. Fattouh 1977



E.

and the second se

The gradient algorithm for computing the steady-state, periodic response of circuits has been derived on the basis of a generalized tableau representation of the network equations. In contrast to the state variable formulation which was recently reported, the present one lends itself to straightforward implementation in modern transient analysis programs for large scale, nonlinear circuits, which make use of sparse matrix techniques. The algorithm has been implemented in one such program, SCAMPER, and details of the program as well as results of tests with several circuits are presented. The rate of convergence depends on the optimization subroutine

and is sensitive to scaling of both the variable as well as the gradient

vectors.

ABSTRACT

L'algorithme du gradient pour le calcul de la réponse périodique d'un circuit a été derivé à partit d'une representation généralisée en tableau des équations du circuit. En contraste avec la formulation selon les variables d'état qui a été récemment publiée, ce rapport presente une méthode plus aisement applicable avec les programmes moderne four l'analyse transitoire des réseaux non-linéaire et de grande dimension, utilisant des techniques de matrices creuses. On a implanté l'algorithme sur le programme SCAMPER, les détails du programme aussi bien que les résultats des tests sur plusieurs circuits sont presentés. La convergence des résultats dépend des sousprogrammes d'optimisation et est sensible aux transformations linéaires des variables aussi bien que du vecteur gradient.

RESUME

ii

ACKNOWLEDGEMENTS

The author wishes to express his gratitude to Professor N. Rumin for his continual encouragement and helpful criticism during the course of this work. Professor M.L. Blostein was the source of many valuable suggestions as well as of considerable assistance during the period of familiarization with the intricacies of SCAMPER, for which the author is deeply indebted.

Thanks are also extended to Mr. D.J. Millar who also helped to reduce the author's period of familiarization with SCAMPER, and to Dr. M.A. Nakhla and Professors C. Paige and A. Buckley who very kindly provided the optimization subroutines which were used in this work as well as many useful suggestions.

Special thanks are due to Mrs. P. Hyland for her skillful

The work was supported in part by the National Research

	• •	• .		
		1 1		
		1	TABLE OF CONTENTS	9 1
	,			س
	1	6		Page
•		,	<i>i</i>	<u>1090</u> ,
		(
•	ABSTRACT		' ' i J	i
	-			
	· · ·	1	· ·	
	RESUME			ii
			· · · · ·	•
	ACKNOWLEDG	FMENTIC		444
	ACKNOWIEDG.	EMENIS	· ·	111
				•
	TABLE OF C	ONTENTS		' iv
r	1		,	~
,	CHAPTER	I	INTRODUCTION	1
,	ļ ,			
	CUNDERD	TT	DEDTOTC SMEADY SMART ANALYSTS AT CODIMINS	7
	CHAPIER	ΥT	PERIODIC STERDI-STATE ANALIŞIS ALGORLIAMS	/
	ۍ بم د	• •		•
	,	2.1	Newton Method	9
	× ·	2.1.1	Computation of the Jacobian J	1.
		- \	Using the Sensitivity Circuit Approach	1.12
,		2.1.2	Computation of the Jacobian J	
			Using the Adjoint Network Approach	\19
		2.2	The Extrapolation Methods	25
		2.2.1	The e-Algorithm	26
	b	2.3	Gradient Method	31
		2.4	Comparison of Methods	34
	ı.	2.4.1	Computation Cost and Convergence Criteria	35
٠	γ.	242	Network Size and the Numerical Stability Problem	41
	1		successive and the numerical beability ribbles	
		-		
¢1	OUNDERD	***		• 1
	CHAPTER	TTT	A GENERALIZED APPROACH TO THE	4 m 1
1	•	-	GRADIENT METHOD	43
		(- -	• /	
		3.1	Introduction	43
`` •	-	3.2 /*	Derivation of the Gradient Equations Based	1
/ /	0	1	on a Generalized Adjoint Variational System	44
	1 61 1	3.3 -	Choice of the Objective Function 🏼 🌢	54
	1	3.4	Detailed Computations of the Gradient and	· · ·
		,	the Initial Conditions of the Adjoint System	56
		3.5	A Computational Technique	61
		x		
	1		e l	
	CHAPTER	IV	NUMERICAL CONSIDERATIONS AND RESULTS	66
1			/	
	1	4 1	Introduction	66
	,	A 2		'ch (
		4.4	Dearries Dearsles	20
	•	4.5	NUMELTOAT VESATOR	12

• , , , ,

; J

Ű

4

iv

Page

		· · .	¥ ,	, • <i>i</i>	
	CHAPTER	v	PROGRAMMING		-
			, · · · · ·	-	
r	t		PART A	· •	
,	· •	1			
• -	4	4		~ ~ ~	
7	•	•	DC AND TRANSLENT ANALYSIS PROGRAM SCAMPER	84	
-	. i .	•		4	
	Ø	5.1	Acceptable Element Types	85	
	• •	5.2	Order of Setting the Equations and Variables	86	6
	۳ ۳	5.3	Solution Procedures and Numerical Techniques	[*] 87 [*]	
	**	5 4	Code Ceneration	80	•
		· E E		101	
	,	· 5.5	Sparse, Marrix Technique	91	
	,	5.6	Pivoting Procedure	93	
	,	5.7	Numerical Considerations	96	
ø	,	5.7.1	Truncation Error "	97	
	1	5.7.2	Nonlinear Error	`97 '	
	, t , t	5.7.3	Step Size Control "	98	
			<u> </u>		
	, s *s - s		, מיזימגמ	,	
			PARI D		
	r .				•
			STEADY-STATE ANALYSIS ROUTINE	· 100	
) /*	Ser .				• ??
	,	5.8	Code Generation for the Adjoint System	100	· ·
		5.8.1	Transposition of the Jacobian J	103	•
	1	5.8.2	E-Code Generation	105	•' ,
	•	5.8.3	F- and S-Code Generation and Divoting	106	•
	5	5 9	Data Management	107	
		5.jo	Some Reports of the Backward Tatespetian	107	
łí		5.10	Some Aspects of the Backward Integration	· • • • •	•
₽Į			or the Adjoint System	109	
	< .	5.10.1	Setting the Initial Conditions of the		
			Adjoint Systems	111	· · ·
		.5.10.2	Truncation Error Computation	114	
*	· .	5.10.3	Interpolation and Step Size Control	115	
		5.11	Resetting the Initial Conditions		• •
	4 .		of the Original System	118	
	•	-			
	, *	ډ		-	ゆ
		* ***	CONCLUCTOR AND DECONDENSION MEDICA	100	•
	CHAFTER	VI.	CONCLUSIONS AND RECOMMENDATIONS	, 120 -	
	、			۹. s ^{pi}	-
	APPENDIX	Т	ACCEPTABILITY OF DOUBLE DEPENDENCY		1
	,		REACTIVE ELEMENT	124	1
	1	۰, ۲		• • • • • •	
	ì			, '	
1	APPENDIX	II ·	VERIFICATION OF ACCEPTABILITY		
	1		OF DOUBLE DEPENDENCY	128	
	`	4		120	
	n	-		dy ~	•
·	REFERENCE	s ·		138	
4			, ,		
				·	
	•			f	
		,			
•			, ,	•	~
				f	
	ι. ·		· · · · · · · · · · · · · · · · · · ·	, <i>,</i>	

Œ

CHAPTER I

INTRODUCTION

Many important classes of electronic circuits, such as power supplies, large signal amplifiers, oscillators, et cetera, operate, in the periodic steady-state mode and yet, because of their nonlinearity, the designer is usually forced to analyze them using time domain, transient analysis computer programs. The periodic steady state solution can be found by the continuous integration, ("brute force" method) of the system equations until all the transient components die away. However, such circuits are, more often than not, lightly damped, in which case this method becomes prohibitively expensive because the integration process has to be continued for a very large number of periods. Recently, this problem has received considerable attention which has yielded several procedures for computing the steady-state response of nonlinear circuits at a cost which is considerably less than that for the "brute force" method.

In 1972, Aprille and Trick [1] proposed a method for the steady-state analysis of nonlinear circuits with periodic input which is based on the Newton algorithm and which they implemented using a state variable formulation of the network equations. The Newton step is used iteratively to adjust the initial conditions at some time t_0 until the difference between these and the conditions at time $t = t_0 + T$, where T is the period, is negligible. An extension of this method to the oscillatory (autonomous) case was given in [2], where the period T is

considered as an unknown together with the initial conditions. The restriction of using the state variable representation of the network equations is inconvenient because it prevents the implementation of the method in transient analysis programs which are based on the tableau representation and which, consequently, can handle large problems. Generalization of the Newton method without this restriction were reported in [3], [4], and [5]. Director [5] used the concept of the "adjoint networks" which he had introduced in [6], while Trick et al [3], [4] based their approach on the "sensitivity circuits" concept. In fact, both the adjoint networks and the sensitivity circuits concepts can be applied to either the state variable or the tableau form, which If the network has N reactive elements, both ever is more suitable. techniques - the sensitivity circuits or the adjoint networks - involve the integration of N adjoint systems of equations in addition to the original network equations over one full period. For networks of moderate order N, the Newton method works guite well. But for high order systems (say N > 50), the computational cost may be prohibitive, and numerical difficulties may be encountered due to necessity of solving systems of large and full matrices.

A second direction for solving the periodic steady-state problem was initiated by Nakhla and Branin [10], [11], using an optimization-like approach, called the gradient method. A performance function ϕ , is defined which reflects the state of the system relative to its steady-state. The first derivatives (gradient) of the

objective function with respect to the initial conditions of the state variables x_0 , in addition to both x_0 and ϕ are passed to an optimization routine. The new values of x_0 returned by the optimization routine are used to update the initial conditions of the system. This process is repeated until the objective function ϕ reaches a specified minimum which is considered to correspond to the steady-state. The effort required in the gradient method for each iteration is the integration of only two systems of equations over one period, the original system in the forward direction, and the corresponding adjoint variational system in the backward direction. The gradient method is characterized by its low cost of integration per iteration step as well as by its freedom from potential numerical stability problems. Nakhla and Branin [10] implemented the gradient method using the state variable formulation, although they did also discuss a more general approach to their method. However, it was not particularly suitable for application with analysis programs based on the tableau representation of the network-It should be mentioned here that a similar gradient method equations. was used by Director [7 - 9] in optimizing nonlinear circuits. [8], [9] Current and Director showed that the optimization problem can be extended to include that of determining the periodic steady-state as In other words, the variable space can be extended to include well. not only the optimization parameters but also the initial conditions of the network.

.

A third direction for determining the steady-state of periodically forced nonlinear systems was given by Skelboe [16] and called the extrapolation method. ℓ_{i} This approach is based on the ε - algorithm by Wynn [15] and is, in fact, a sequence-to-sequence transformation where the ε - algorithm step is used to modify the initial conditions in a nonlinear sequence of operations. The algorithm does not require the computation of first derivatives (such as the sensitivity matrix for the Newton method, or the gradient vector for the gradient method), but, in place, it requires the integration of the system over many periods. The advantage of the extrapolation method is its quadratic rate of convergence (or even higher). But there are many restrictions on the method such as the assumption that the system is mildly nonlinear and the requirement that the order of the system be not too Jarge [16] .

As mentioned above, the gradient method introduced in [10], [11] was based on the state variable approach. In the work reported here, a more general formulation is developed for use in transient analysis programs based on tableau representation of network equations. The detailed implementation of the resultant algorithm in one such program SCAMPER is given as well as the numerical results too of some test problems solved on the new program.

In Chapter II, a survey of the above three methods for the periodic steady-state analysis (Newton method, gradient method, and extrapolation method) are presented. For the Newton method, the implementa-

tion is discussed according to two different approaches, the sensitivity circuits [4], and the adjoint networks [5]. Also, for both the Newton and the extrapolation methods, executable algorithms suited to the general representation of network equations are presented. In the case of the gradient method the approach is the same as that given by Nakhla [11]. The chapter ends with a discussion of the advantages and disadvantages of these methods.

The gradient method based on the general tableau representation is developed in Chapter III. Lagrangian multipliers are used in a manner similar to that employed by Director [8] to define the adjoint system of equations. A fairly general representation of the network nonlinearities is assumed. The method is applied to the dc and transient analysis program SCAMPER on the basis of a particular Choice of the objective function ϕ

In Chapter IV, some considerations related to the problem of scaling are presented, and the numerical results of some test problems are given. It should be noted that the same examples were studied by many authors including Trick [1], [3], Director [5], Skelboe [16], and Nakhla [11]. Since the algorithm presented in Chapter III is essentially a generalization of Nakhla's method and because the same optimization routine [22] was used, the results of both studies are compared where possible and are shown to be in good agreement.

Chapter V. is concerned with the programming problem. In. Part A, a review of the existing dc and transient analysis program SCAMPER [19], [20] is given. The key features and procedures are discussed, in particular, code generation, sparse matrix and pivoting techniques, and error and step size control. In Part B, the main features of the additional programming required for implementing the steady-state algorithm presented in Chapter III, are discussed. These features are : (1) relation between the original system and its adjoint, (2) the procedures for generating the necessary code, segments for integrating the adjoint system, (3) the problem of data manipulation (retrieval from disk, interpolation), (4), setting the initial conditions of both the original and adjoint systems, and (5) error and step size control.

6

The main results of the present work are summarized in Chapter VI, which concludes with a brief discussion of further improvements to the present method as well as of topics for work in this area.

CHAPTER II

PERIODÍC STEADY-STÀTE ANALYSIS ALGORITHMS

The periodic steady state of an electronic network can be analyzed either in the time domain or in the phase (frequency) domain. If the network is almost linear, with relatively few nonlinear elements, and the harmonic content is not too large, the phase domain is preferable [25]. But, when the network contains a significant number of nonlinear elements with sharp nonlinearities, the time domain is more suitable. Since it is this latter class of circuits which we are interested in, the following discussion will concentrate on time domain algorithms.

The objective here will be to examine several algorithms for the computation of the periodic steady state, from the point of view of their suitability for inclusion in an existing dc and transient analysis program based on the sparse tableau approach. Furthermore, although we are primarily interested in the nonautonomous case, ease of extension to the autonomous network is very desirable. Finally, an approach which would permit the subsequent inclusion of circuit optimization evidently merits particularly serious consideration.

The most recent algorithms are of the iterative type, and can be divided into two groups:

(i) The self sustained algorithms, i.e., those which do not need external steps, such as optimization. The Newton Methods [1 - 5] and the extrapolation methods [16], are members of this group.

The optimization - like algorithms, i.e., those that involve an external optimization step to achieve the circuit analysis. The best known method in this group is the gradient method [10], [11].

(ii)

ļ

In order to simplify the discussion, the reduced (state variable) form of equations will be used, but more general representations will be introduced where necessary.

Assume that we have a system of nonlinear differential equations given by,

 $\dot{q} = f(q, t)$

where q is the differentiated variable vector, and f is a \mathbb{R}^n vector and has a centinuous, first order partial derivative w.r. to q for $t < 0, \infty >$. Furthermore, we assume that the system has a unique periodic solution trajectory of period T, i.e.2

(2.1)

q (t) = q (t + n T) ; n = 1, 2, 3, ...

We seek a set of initial conditions $q(t_0)$ at $t = t_0$, which will put the system into its periodic steady-state, i.e.,

, *6*

 $\dot{q}(t_0) = q(t_0 + f)$ (2.2)

Following is a discussion of three steady-state analysis methods, with particular emphasis on the algorithms and computational techniques.

2,1 Newton Method

The application of the Newton method for the periodic steadystate analysis was first introduced by Aprille and Trick [1]. Several improvements of the computation techniques as well as extensions of the area of application are given in [2 - 5]. In this section, the Newton iteration will be stated, and two different computational techniques will be given, one using an approach based on sensitivity circuits [4], and the other on adjoint networks [5].

For iteration number i, assume q_0^i is the initial condition vector. It is required to adjust q_0 , until the difference between two succeeding iterations initial conditions, $|\Delta q_0^i| = |q_0^i| q_0^{i+1}|$, is less than an acceptable lower limit, at which the system will be considered in steady-state. The Newton iteration used to adjust the initial condition is given by:

16

1+

¢'s

where q_0 is an R^n initial condition véctor, $q(q_0, T)$ is the solution after one full period beyond the initial point t_0 , and i is the iteration number.

Equation (2.3) can be set in a more suitable form:

$$J^{i} \Delta q_{0}^{i} = E^{i}$$

where

~~ ~

$$J^{i} = I - D^{i} = I - \frac{\partial q^{i} (q_{0}, T)}{\partial q_{0}^{i}}, \qquad (2.5a)$$
$$\Delta q_{0}^{i} = q_{0}^{i} - q_{0}^{i+1}, \qquad (2.5b)$$

$$E^{i} = q^{i} (q_{0}, T) - q_{0}^{i}$$
 (2.5c)

The term
$$D^{i} = \frac{\partial q^{i}(q_{0}, T)}{\partial q_{0}^{i}}$$
 is called the sensitivity matrix, as it

gives the variation of the final state at, $t = t_0 + T$, w.r. to a perturbation in the initial conditions at, $t = t_0$.

Hence, if the Jacobian J and the right hand side vector E can be determined, equation (2.4) can be solved to give $4 q_0^i = q_0^i - q_0^{i+1}$. The iteration will be terminated if $|\Delta q_0^i|$ is less than a specified limit. The following are the main steps in the Newton algorithm:

10

(2, 4)

	2	. ~	, *	11
** * * * ** * *		· ·		, e
Algorithm 2.1			· * *	مبير
	. (.	· · · ·	< ' u	ŝ
(1) Start with an	initial guess q_0^{o}	say, by integr	ating	
for a few peri	iods).	٢		-
(2) $i = 0$.		· • • • •	· .	
(3) Integrate furt	ther one complete pe	eriod and determ	line	
q ⁱ (q ₀ , T) ; 1	hence determine E	•	\$	
· ·	· · · · · · · · · · · · · · · · · · ·	, , , , , , , , , , , , , , , , , , ,	۶ ۱۰۰۰	d'e -
(4) Compute J	(may be done simulta	-	leter~	
mining E ⁺).	*3 43	1	•.	
(5) Solve equation	n (2.4) (may be don	eusing ĹU fa	actoriza-	
tion) for Δ	q_0^{i} and thence det	ermine q ₀ ⁱ⁺¹ .	· ·	
(6) Is $ \Delta q_0^i $ les	s than certain acce	ptable limit?		
If yes, STO	P (the steady stat	e solution is fo	ound)	1
If not:	/	-	~	l, ,
i = i+l.		ور بر	/ 5/} 	o
IF i is gre	ater than a certain	limit, STOP.	~ · · \	•
If not, put	$\left \begin{array}{c} q_0^{i+1} \\ q_0 \end{array} \right $ as the new	initial conditi	ons, and	ي د
go to Step	3.	· /	,	
т г + ъ	be order of the syst	em n (number	of storage al	emente)
is small, this algo	prithm seems to be v	ery attractive	as it has rem	arkably

å,

(say $n \ge 50$), the drawbacks will outweigh the convergence properties.

good convergence criteria. But in the case of high order systems

Firstly, the cost of computing the Jacobian, J. will be very high, since it involves the analysis of n circuits (adjoint networks, or sensitiwity circuits). Secondly, the solution of equation (2.4) involves full matrices of order n which, in addition to slowing down the execution speed, has some critical numerical stability problems. Thirdly, it will not be possible to make use of the sparse matrix techniques used in most recent programs, for solving equation (2.4).

12

The heart of the above stated algorithm is the computation of the Jacobian J, which can be achieved using one of two methods: the sensitivity circuits approach [4], or the Adjoint Network method [6]. These will now be briefly discussed.

2.1.1

Computation of the Jacobian J Using the Sensitivity Circuit Approach

From equation (2.5a) -

 $J^{i} = I - D^{i} = I - \frac{\partial q^{i} (q_{0}, T)}{\partial q_{0}^{i}}$

, Put 🦾

$$D^{i} = \frac{\partial q^{i} (q_{0}, T)}{\partial q_{0}^{i}} = [D_{1}^{i}, D_{2}^{i}, \dots, D_{k}^{i}, \dots, D_{n}^{i}] (2.6)$$

where

is the kth column of the sensitivity matrix D^{i}

To show how equation (2.7) can be computed, it is necessary to examine the sensitivity circuit [4], and thence to show how the computation of each column of the sensitivity matrix D_k will correspond to the solution of a sensitivity circuit, after one full period with specified initial conditions.

Consider circuit equations defined in tableau form by

$$A i = 0$$
 KCL
 $A^{t}v = E$ KVL

The branch constitutive (B.C.) relations - self or mutually dependent - can be set in a general form.

For a non-reactive branch,

 $q_{c_1} = q_{c_1} (v_m)$

For a capacitive branch,

(2.10a)

(2.11a)

(2.8a)

(2.7)

(2.12a)

$$i_{c_{j}} = \frac{d_{c_{j}}}{d_{t}}, \text{ and } v_{m}(0) = v_{m0}$$

For an inductive branch

$$\mathbf{F}_{k_{j}} = \mathbf{F}_{k_{j}} (\mathbf{i}_{m})$$

where

$$dF_{\ell_j}$$

$$w_{j} = \frac{j}{dt}$$
, and $i_m(0) = i_{m0}$

Finally, for independent source "

$$I_{s} = I_{s}(t) \quad (current source) \quad (2.13a)$$

$$E_{s} = E_{s}(t) \quad (voltage source) \quad (2.14a)$$

Assuming that the system has a unique solution, let us perturb the initial conditions of the differentiated variables (generally, we will call them q_0) one at a time. The result is obtained by differentiating equation (2.8a) through (2.14a) with respect to the kth differentiated variable q_{0_L}

$$A \frac{\partial i}{\partial q_0} = 0$$
(2.8b)
$$A^{t} \frac{\partial v}{\partial q_0} = 0$$
(2.9b)



A close examination of the corresponding pairs (a - b) of equations (2.8 - 2.14) will show that equations (2.8b) - (2.14b) describe a circuit having the following properties:

> Linear circuit (time-varying if the original circuit contains nonlinear elements) having the same topology as the original circuit;

(2.13b)

(2.14b)

Þ. .

the type of an element is the same as that of the corresponding element of the original circuit;

the value of an element

∂e_s

ad^{0be}

2.

3.

(i) for a linear element is the same as that of the soriginal;

(ii) for a nonlinear element is replaced, by the piece-

(iii) of an independent source is replaced by a zero source (i.e., short circuit for an independent

wise linearization at each time point;

voltage source, and open circuit for an independent

current source) ;

the "initial/conditions of the differentiated
variables are replaced by a source of unity value
if the differentiation is taken w.r. to this
variable, and with a zero valued source if not.

4.

17

(2.15)

This circuit is called the sensitivity circuit. It should be noted that, corresponding to each reactive element, there must be an equivalent sensitivity circuit. By solving the kth sensitivity circuit over one complete period starting at the initial point t_0 , the solution vector of the differentiated variables at the end of the period is just the kth column D_k of the sensitivity matrix D.

From a computational point of view, only one circuit (the original) is formed and the same coefficient Jacobian matrix, describing the linearized system of equations (in tableau form), is used but with a different forcing vector for the computation of each D_k . To show how this can be done, consider the general linearized tableau form

where B is the coefficient Jacobian matrix, w the algebraic variable vector, q the differential variable vector, and E is the forcing vector.

The tableau form for the kth sensitivity circuit is

$$\begin{bmatrix} \mathbf{u} \\ \mathbf{k} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \end{bmatrix} ; \lambda (\mathbf{0}) = \begin{bmatrix} \mathbf{e} \\ \mathbf{k} \end{bmatrix}$$

where e_k^{\langle} is the kth unity vector (i.e., vector of all zeros except the kth element value is unity), and u, λ^{\langle} are, respectively, the vectors of the algebraic and differentiated variables.

The above method can be implemented by means of the following algorithm for computing the sensitivity matrix D_{NXN} where J = I - D.

Algorithm 2.2

Start at an initial time point t₀;
 k = 0;
 form the coefficient Jacobian matrix (B) of the original circuit; save the nonlinear entries (if any);

solve the original system of equation's (2.15) and save the solution vector and the necessary variables ;

 $\mathbf{k} = \mathbf{k} + \mathbf{l}$

;

6.

7.

8.

5.

4.

load the Jacobian B (using the stored values of the nonlinear entries (if any))

If it is the first time point, put λ (0) = e_k . Solve the kth system of equations (2.16) and save the solution vector and the necessary variables (e.g., variables required to compute λ).

18

(2.16)

- 9. If k less than N go to step 5.
- 10. If $t < t_0 + T$, choose the next step size, and go to step 2.
- 11. Form the sensitivity matrix D, by column, where the kth column D_k is the solution vector of the differentiated variables of the kth system. (Then J = I - D).

2.1.2 Computation of the Jacobian J Using the Adjoint Network Approach

Transform equation (2.6) into the form

$$D^{i} = \frac{\partial q^{i}}{\partial q_{0}^{i}} = \begin{bmatrix} D_{1}^{i}, D_{2}^{i}, \dots, D_{k}^{i}, \dots, D_{N}^{i} \end{bmatrix}^{t} (2.17)$$

where

いたが

$$D_{k}^{i^{t}} = \left[\frac{\partial q_{k}^{i} (q_{0}, T)}{\partial q_{0}^{i}} \right]^{t} , \qquad (2.18)$$

is the kth row of the sensitivity matrix D.

Before describing an algorithm for computing .D, it is

necessary to review the adjoint network concept [6], and its application [5].

Consider an original network, with current and voltage variables i and v, respectively, and a corresponding "adjoint" network with current and voltage variables i and v, respectively. Let us assume the "adjoint" network to have the same topology as the original network. In that case, using Tellegen's theorem [27]

$$\int_{t_0}^{t_0+T} \sum_{j \in [v_j(t)]} [v_j(t)] \hat{i}_j(\tau) - i_j(t) \hat{v}_j(\tau)] dt = 0 \qquad (2.19)$$

where τ is the arbitrary time variable of the adjoint network.

If the original network variables are perturbed by Δv_j , Δi_j , Tellegen's theorem still holds for the perturbed states, i.e.

 $\int_{t_0} \sum_{j} [\Delta v_j | t) \hat{i}_j (\tau) - \Delta i_j (t) \hat{v}_j (\tau)] dt = 0 (2.20)$

Using the assumptions given in [5], the adjoint network can be extracted from the original network using the following properties:

The <u>time</u> variable of the adjoint system is $\tau = t_0 + t_f - t$, where the final time $t_f = t_0 + T$.

2.

1.

t₀+T

The <u>type</u> of an adjoint network element is the same as the corresponding one in the original network and both networks have the same <u>topology</u>. 20 ·

The value of an adjoint element is given by:

(i) for a linear element, its value;

3.

 (ii) for a nonlinear element, the piecewise linearization of the original element value. Thus the corresponding adjoint element is linear and time varying.

(iii) The independent source is replaced by a zero source (i.e., independent voltage source is replaced by a short circu the independent current source is by an open circuit).

It should be noted [5] that, for by diginal and adjoint networks each reactive element with initial on x_0 has to be replaced by a similar reactive element with z untial conditions accompanied by an auxiliary source (parallel current element or series voltage source for capacitive element. The value of that auxiliary source is equal to the initial condition of the corresponding branch.

Now there remains the problem of setting the initial conditions of the adjoint reactive elements. Suppose we are computing the kth row of the sensitivity matrix, D_k . There are two cases [5] to be considered.

(i)

ŕ,

If the kth element is a capacitor, the initial

condition of the kth adjoint network is

$$\hat{q}_{0_{k}} = \frac{1}{C_{k}}; e_{k}$$
 (2.21a)

where e_k is equal to (0, 0, ..., 1, ..., 0) with the one in the kth position. It can be shown that if the mth element of the kth adjoint network : is capacitive.



But if it is an inductor

$$\frac{\partial \mathbf{v}_{c_{\mathbf{k}}}}{\partial \mathbf{i}_{\ell_{0_{\mathbf{m}}}}} = -\mathbf{L}_{\mathbf{m}} \left(\mathbf{v}_{\ell_{\mathbf{m}}} \left(\mathbf{t}_{0} + \mathbf{T} \right) - \mathbf{v}_{\ell_{\mathbf{m}}} \left(\mathbf{t}_{0} \right) \right) \left| (2.23a)$$

(ii) If the [?]kth element is an inductor hence the initial condition of the kth adjoint system is

$$I_{0_k} = -\frac{1}{L_k} e_k$$

(2.21b)

If the mth element of the kth adjoint system is a

capacitor

$$\frac{\partial i_{\ell_{k}} (t_{0} + T)}{\frac{\partial v_{c0_{m}}}{}^{2} v_{c0_{m}}} = C_{m} (v_{c_{m}} (t_{0} + T) - v_{c_{m}} (t_{0}))$$
(2.22b)

But if it is an inductor

$$\frac{\partial i_{\ell_0} (t_0 + T)}{\frac{k}{2} \cdot t_0} = -L_m (\hat{i}_{\ell_m} (t_0 + T) - \hat{i}_{\ell_m} (t_0))$$
(2.23b)

Due to the restriction imposed on the adjoint time variable $\tau = t_0 + t_f - t$, the coefficient Jacobian of the adjoint networks is not computed at the same time point as the original network. So the adjoint computations have to be delayed until after the original network integration over the whole period $(t_0 t_0 t + T)$, and the necessary variables have to be stored. Then, after finishing the original system integration, the adjoint system can be integrated in the opposite direction, using the This is equivalent to integrating backward from stored information. $t = t_0 + T$ up to t_0 . From this discussion it should be clear that the adjoint network approach will slow down the 'execution speed relative to the sensitivity circuit approach. Also it may cause programming difficulties related to the retrieval of information from disk and the interpolation process as well.

A computational algorithm can be established for computing the sensitivity matrix D, using the adjoint network approach.

Algorithm 2.3

1,

2.

Start at an initial time t_0 . Integrate the original system up to a final time $t_f = t_0' + T$. At each (accepted) time point, store the necessary values of the nonlinear elements (if any).

Construct the coefficient Jacobian of the adjoint system and generate the part of the code that perform the LU factorization. Hence, for each adjoint system, generate the other parts of the code (such as that required for the forward and backward substitution) to be used for integrating each system, and save these codes.*

3,•

4.

5.

k = 0;

k = k + 1;

IT it is the first point, put the initial conditions according to equations (2.21a, b);

Load the Jacobian and the necessary variables.

Solve the system (by executing the generated code)

* This step is oriented for a specific time domain analysis program (SCAMPER) in which the integration is done using a generated machine code string. But it can be modified depending on the features of the program used in the analysis. and save the necessary variables (such as those needed for computing the derivatives $\hat{\mathbf{x}}$).

If k less than N go to step 4.

8.

9.

7.

If $t \le t_0$ go to step 9; If $t > t_0$, choose the next step size and go to step 3;

Formulate the matrix D, by row, using the solution vectors of the adjoint systems, where the kth row is that corresponding to the kth adjoint system.

2.2 The Extrapolation Methods

The extrapolation methods [16] are based on the

 ε - algorithm [15] which does not require the calculation of derivatives (such as the gradient in the optimization-like algorithms or the sensitivity matrix in the Newton methods). But, in place of this, for a system with N independent differentiated variables, 2M solution vectors ($M \le N$) must be determined over a time interval 2MT in order to carry the algorithm one step, where T is the period.

Originally, the ε - algorithm proposed by Wynn [15] was oriented to the solution of systems of linear differential equations. With faildly nonlinear systems, and near the equilibrium (steady-state

condition), the algorithm, hopefully, can be applied to find this solution with a quadratic convergence [12], [13].

2.2.1 The ε - Algorithm

The ε - algorithm is a nonlinear sequence-to-sequence transform. Assume that the system solution sequence is q_0, q_1, \ldots which may be produced by any integration process. Also assume vectors { ε_k^r } where k represents the transformation step and r is the sequence number at this step. If initial vectors are given by:

$$e_{-1}^{(r)} = 0$$
; $r = 1, 2, ... 2m$ (2.24a)

(2.24b)

The recursion formula

$$\varepsilon_{k+1}^{(r)} = |\varepsilon_{k-1}^{(r+1)} + (\varepsilon_{k}^{(r+1)} - \varepsilon_{k}^{(r)})^{-1}; k, r = 0, 1, \dots (2.25)$$

is known as the ε - Algorithm.

 $x = (x_1, x_2, \dots, x_N)$.

 $\epsilon_0^{(r)} = q_r; r = 0, 1, ... 2m$

There are two forms of the $\varepsilon \neq algorithm$, depending on the definition of the inverse of

The scalar ε - algorithm is based on the primitive inverse definition

 $x^{-1} = (x_1^{-1}, x_2^{-1}, \dots, x_N^{-1})$

so that each component is treated independently.

- The vector ε - algorithm involves the Euclidean norm definition of the inverse,

$$x^{-1} = x / ||x||_2^2$$
; $||x||_2^2 = \sum_{i=1}^{N} x_i^2$

The following discussion will be restricted to the vector ε - algorithm because the organization of the computation is easier and, moreover, there is some evidence that this form is more

stable [12].

ſ

Į.

金武院

It was shown by McLeod [14] that if a sequence $\{q_{\chi}\}$ satisfies the linear recursion form

 $\sum_{i=0}^{m} \alpha_{i} q_{r+i} = (\sum_{i=0}^{m} \alpha_{i}) q ; r = 0, 1, \dots$ $i=0 \qquad i=0$

where $\{\alpha\}$ is a real vector, and q is a constant vector, then

$$\frac{\varepsilon(\mathbf{r})}{2\mathbf{m}} = \frac{\mathbf{q}}{\mathbf{r}}; \ \mathbf{r} = 0, 1, \dots, \text{ if } \sum_{i=0}^{m} \alpha_{i} \neq 0$$

(2.26)

(2.27a)

$$\epsilon_{2m}^{(r)} = 0, r = 0, 1, ..., \text{ if } \sum_{i=0}^{m} \alpha_i = 0$$
(2.27b)

In the case of a nonlinear-system, where the solution sequence $\{q_k\}$ is produced, say, by the contraction mapping

$$q_{r+1} = F(q_r)$$
,

or

(

the ε -algorithm can be applied if [12], [13]

$$q_{r+1} - q' = F'(q)(q_r - q) + 0||q_r - q||_2^2$$
 (2.28)

where this notation means the last term has norm of the order $||q_{1} - q||_{2}^{2}$.

Under this condition, the convergence of the sequence q_r to the fixed point q of the mapping F is at least quadratic.

It was stated at the beginning that an integration over 2M periods is required for each iteration, hence, the determination and minimization of M is of great importance. This point has been studied by Skelboe [16] who examined closely each circuit and, in particular, the relative values of the circuit time constants and the period(s) of the forcing function(s). The maximum value of M is the order of the system N i.e., the number of independent reactive elements. One way of reducing M is by the proper choice of the starting time t_0 . This point can be selected such that some of the faster transients will have

died out by $t = t_0$
29 From the previous discussion, an executable procedure, , based on the ' ε - algorithm can be stated: Algorithm 2.4 Initialization step: $\tilde{k} = -1$ (i) Define (m), the starting point (t_0) and the (ii) acceptable error limit (δ). Integrate over $t \in (0, t_0)$ and set: (iii) $z = q(t_0)$ Integration step: k = k + j**(i**) $q_0 = z_{k-1}$ (ii) Integrate over $t \in (t_0, t_0 + 2 mT)$ and set: (iii) $q_{r} = q (t_{0} + r T)$, r = 1, 2,• • • 2m 3. algorithm step: ε-3 (i) put: $\varepsilon_{-1}^{(r)} = 0$, $\varepsilon_{0}^{(r)} = q_{r}$, r = 0, 1, ... 2mand

$$\varepsilon_{j+1}^{(r)} = \varepsilon_{j+1}^{(r+1)} + (\varepsilon_{j}^{(r+1)} - \varepsilon_{j}^{(r)})^{-1}, \quad r = 0,$$

 $j = j + L$

(iii) If j < 2m - 1, go to step (3 - ii).

4. Test of convergence step:

i)
$$z_k = \varepsilon_{2m}^{(0)}$$

(ii) If
$$||z_k - z_{k-1}||_2^2 \le \delta$$
 the steady-state

solution is found by using z_k as the initial condition at $t = t_0$. STOP.

If not go to step 2.

The computations in step (3) involve a tradeoff between execution time and storage requirements. If the execution time is to be minimal, at least 2 x (2M) vectors (each having N entries) must reside in the high speed memory, i.e., storage allocation of at least 4MN entries is required, which is four times bigger than that required for the Néwton algorithm. If this storage area is to be minimized, then 2M retrievals of data from disk memory are required, which will slow down the execution speed.

1, ... 2m-1

2.3 Gradient Method

The gradient method for the periodic steady-state analysis of nonlinear circuits was first proposed by Nakhla and Branin [10], [11]. Their implementation was based on the state variable formulation of the Although one of the objectives of this work was to circuit equations. make use of the gradient method in conjunction with the more general and⁷ convenient tableau representation of the circuit, the succeeding discussion will be based on the notationally simpler state variable approach as given in [11]. The detailed analysis of a more general tableau representation is left for a later Chapter (III).

Let the state variable equations describing the circuit be given by:

> $\dot{q} = f(q, t)$ (2.29)

where f and the differentiated variables $q \in R^{N}$. Given that the forcing function vis periodic and of period T, it is required to find a set of initial conditions $q(t_0)' = q_0$ such that

 $q(t_0) = q(t_0 + kT)$, $k = 1, 2, \dots$ For convenience k = 1. The deviation of the circuit from its periodic steady-state is expressed in terms of a performance function Φ which is defined by

 $\Phi = \left[\delta \left(\mathbf{q}_{0} \right) \right]^{\mathsf{L}} \left[\delta \left(\mathbf{q}_{0} \right) \right] ,$

(2.31)

(2.30)

where δ (q₀) is the discrepancy vector given by

$$\delta (q_0) = q (t_0) - q (t_0 + T) ,$$
 (2.32)

and the superscript t indicates transposition. Thus the objective is to find a set of q_0 which will minimize ϕ_0

The nonlinear programming technique used to minimize Φ requires the first derivatives (gradients) $\partial \Phi / \partial q_0$. The crux of the method is to demonstrate an efficient scheme for computing the gradient $G = \partial \Phi / \partial q_0$. It has been shown [11] that the gradient can be computed by only integrating two systems of equations over one complete period: the original system of equations from t = 0 to t = Tand a second system, which will be referred to as the adjoint variational system, from t = T to t = 0. From theorem 2.1 [11] and equations (2.29) to (2.32) it can be proved that the gradient G has the simple form

$$G = \frac{\partial \Phi}{\partial q_0} = 2 \left[\lambda \left(\delta, 0 \right) - \delta \left(q_0 \right) \right]$$
(2.33)

where λ ($\delta,$ t) $% \lambda$ is the solution of the adjoint variational system

$$\dot{\lambda} = -\left[\frac{\partial f}{\partial q}\right]^{t} \left[\lambda\right] \qquad (2.34)$$

$$q (q_0, t)$$

with initial conditions

 λ (δ , T) = δ (q_0)

From equations (2.33). to (2.35) it is evident that the adjoint system of equations has to be integrated backwards from t = Tto t = 0. This requires knowledge of the initial conditions λ (δ , T) = δ (q₀) and the trajectory of the coefficients of the Jacobian matrix $\frac{\partial}{\partial} \frac{f}{q}$ over the whole period. The initial conditions are given simply by the difference between the solutions of the original system at t = 0 and t = T. The Jacobian $\partial f / \partial q$ can be formed from the forward integration of the original system using one of two methods, namely, by either saving the nonlinear entries of $\partial f / \partial q$ directly; or by saving the necessary variables from which $\partial f / \partial q$ can be evaluated.

A computational algorithm can be written on the basis of the above discussion [10], [11].

Algorithm 2.5

- 1. Choose a first estimate of q_0 , perhaps by integrating equation (2.29) until a predefined time t_0 .
- 2. Using q as initial conditions, integrate equation (2.29) for one full period, saving the necessary variables to compute the Jacobian $\partial f / \partial q$.

3. At $t = t_0 + T$ compute the discrepancy vector δ (q₀)

33

(2.35)

defined by equation (2.32). Hence compute the performance function Φ defined by equation (2.31). If ϕ is less than or equal to an acceptable limit, the steady-state solution is found. STOP.

Otherwise, use $\delta(q_0)$ as initial conditions of the adjoint system defined by equation (2.34), and integrate backwards over one full period from t = T to t = 0 using the saved values from the forward integration (in reverse order) to compute the Jacobian.

Compute the gradient G defined by equation (2.33).

Pass the arguments q_0 , gradient G and the objective function Φ to the optimization routine Go to step (2) using the new arguments q_0 returned by the optimization routine as initial conditions.

2.4 Comparison of Methods

4.

5.

6.

7.

Three possible methods for computing the periodic steadystate response of nonlinear circuits have been reviewed. In order to

arrive at a decision as to which method is the most suitable, it is necessary to define the class of problems which has to be solved. The present work was undertaken with the objective of finding a method which would not only be used to compute economically the periodic steady-state response of large, stiff and nonlinear systems, but which could also be easily extended to include circuit optimization and sensitivity analysis. In the following sections the three methods are compared with respect to those characteristics which would have the greatest influence on the final choice.

2.4.1 Computation Cost and Convergence Criteria

Let us first consider the convergence problem. For the gradient method (GM), the rate of convergence depends essentially on (1) the stiffness of the network equations, (2) the initial state of the network relative to the steady-state conditions, and $\frac{1}{2}$ (3) the optimization routine used for minimizing the objective function Φ . With the optimization routine based on Fletcher's [22] variable metric method, the GM has been demonstrated by Nakhla [11] to converge quadratically near the equilibrium (steady-state) conditions. When the network initial conditions are quite far from the steady-state, the method still converges but at a slower rate. In the case of the Newton method (NM), a comparison of results reported by Trick [3] and Nakhla [11] shows that the. NM is more likely to converge faster

than the GM but with a narrower range of convergence. The convergence of NM depends on (1) the stiffness of the Network equations, and (2) the state of the network relative to the steady-state [1]. For the extrapolation method (EM), the few examples reported by Skelboe [16] seem to indicate that the rate of convergence of the EM is at least quadratic with a range as wide as that of the GM... In comparison to the Newton method [1], [3], the EM [16] seems to be even faster. The convergence of the EM depends on (1) the stiffness of the network equation, and (2) the mildness of the network nonlinearities.

From the above discussion and in the light of the published results for a relatively small number of modest size problems, it appears that, within the range of problems where the three methods GM, NM, and EM are applicable, the extrapolation method converges faster than the two other methods while the gradient method (with the variable metric optimization routine) is the slowest. But the number of iterations required for convergence by each of these three methods differ γ by ratios which are not so far from unity.

Let us turn now to the problem of computation cost.

 $\mathbf{J}, \mathbf{X} = \mathbf{E}$

where J is the Jacobian matrix, and E is the forcing vector.

Suppose that a LU factorization method is to be used for solving the network equations at each integration step, where the number of steps over one full period is n_s^* . The major portion of the computation cost per step can be deduced by considering the costs of (1) computing and loading the Jacobian (J), (2) computing the forcing vector (E), (3) factoring the matrix (LU), and (4) the forward and backward substitution (FB).

Now, for the gradient method it has been shown that an integration of two systems (the original system in the forward direction, and the adjoint system in the backward direction) over one full period is required for each iteration. If we assume that the cost of integrating the adjoint system is the same as that for, the original one, and that the number of iterations required for convergence is $n_{\rm G}$, then, the integration cost of the gradient method is

* In some time-domain analysis programs, such as NDP [17], the integration process proceeds as follows: at each time point, the system is solved with a predetermined step size h. If the system does not converge, a Newton iteration is applied repeatedly (with the same step size) k times until it converge. If we consider the number of time steps equal to n_t , and the average number of Newton iterations per step is k, then the number of integration steps n_s will be given by $n_s = n_t k$.

In some other programs, such as SCAMPER [19], [20], the integration process is as follows: at each time point a step size h is chosen and the system of equations is solved. If it does not converge, the step size is reduced until convergence occurs. In this case, the number of integration steps n will be the number of solution trials (both accepted or not accepted).

 $GM = 2 \cdot n_{g} \cdot \eta_{g}$. (J + E + LU + FB)

38

(2.36)

(2.37)

(2.38)

For the Newton method, in addition to the original system, an integration of N adjoint systems (N sensitivity circuits [4] simultaneously with the original system, or N adjoint networks [5] in the reverse (backward) direction) must be carried out, where N is the number of reactive elements in the network. Consider the case [4], where the integration of the N sensitivity circuits is carried out simultaneously with integrating the original system. The integration of the N sensitivity circuits can be carried out using the same Jacobian J and LU factorization used for integrating the original system. Hence the computation cost of the Newton method can be approximately given by:

 $NM = n \cdot \eta_{N} \cdot [J + LU_{v} + N \cdot (E + FB)]$

where n_{N} is the number of iterations required for convergence.

For the extrapolation method, it has been shown that an integration of one system (the original) has to be carried out over 2M periods per iteration where $M \le N$. Accordingly, the computation

cost of the extrapolation method will be given by

 $EM = 2 \cdot n_s \cdot M \cdot \eta_E \cdot (J + E + LU + FB)$

where η_{E} is the number of iterations required for convergence.

From equations (2.36) - (2.38) the relative cost of the three methods can be given by

 $GM : NM : EM = 2 \cdot n'_G \cdot (J + E + LU + FB) :'$

 $\eta_{N} \cdot [J + LU + N \cdot (E + FB)] : 2 \cdot M \cdot \eta_{E} \cdot (J + E + LU + FB)$ (2.39)

In order to simplify this equation, let us consider a practical situation. As mentioned before, we are interested in moderate and large scale nonlinear networks, with quite a large number of reactive elements N (N \geq 20). Suppose that the network equations are to be written in the tableau form with the number of equations equal to NE (which is proportional to the number of branches in the network), Consequently, the Jacobian J will be sparse. For sparse matrices with quite large dimension NE, Norin and Pattle [26] reported that LU and FB requires O (NE) operations, with LU = 0.4 O(NE), FB = 0.5 O (NE), approximately. In general we can write

> $LU = \alpha \cdot NE$ LU

 $\mathbf{FB} = \alpha_{\mathbf{FB}} \cdot \mathbf{NE}$

(2.40)

The cost of computing and loading the Jacobian J and the forcing vector E depends on (1) the number of equations NE , (2) number and types of nonlinear elements, and (3) the number of reactive elements. For simplicity, let us assume that these two costs, J and E, depends linearly on the number of equations NE such that

$$J = \alpha_{J} \cdot NE$$

$$E = \alpha_{E} \cdot NE$$
(2.41)

2

It has been found experimentally* that the cost of J or E is less than that for LU or FB (e.g., $\alpha_E \approx 0.3 (\alpha_{LU} + \alpha_{FB})$). If we normalize the units of the computation cost such that $\alpha_J + \alpha_E + \frac{1}{\alpha_{LU}} + \alpha_{FB} = 1$, then equation (2.39) can be replaced by

: NM : EM = 2 η_{G} : η_{N} [α_{J} + α_{J} + α_{LU} + N (α_{E} + α_{FB})] : 2 . M . η_{E}

* A network, with the number of branches NB = 78 of which 40 were nonlinear, and the number of tableau equations NE = 199, was analysed using SCAMPER [19], [20] with the steady-state algorithm included (see Chapter V). Both J_{\parallel} and E are computed in part using Fortran statements and in part using machine instructions (code), while LU and FB are done using a generated code only. The length of the codes for this problem was as follows:

LU + FB = 7083, E = 1498, J = 146.

The computation of J was mainly by Fortran statements, while the computation of E was mainly using the generated code (E code). From a close examination of the network we expect that the total cost of E will be approximately 50% higher than the number given above. Also we expect the cost of computing the Jacobian to be less than or equal to the cost of computing E. According to these estimates, a rough cost comparison can be given as follows:

E, $J \approx 1.5 - (1498 / 7083)$ (LU + FB) E, $J \approx 0.3$ (LU + FB).

i.e.,

40

(2.42)

Assuming M is of the order of N ($M \le N$), it appears that for networks with quite a large number of reactive elements (N > 20), the gradient method yields the lowest computational cost provided that η_G , η_N and η_E are approximately equal.

2.4.2 Network Size and the Numerical Stability Problem

For the Newton method, the sensitivity matrix $\partial q (q_0, T) / \partial q_0$ is, in general, a full matrix. When the order of the system N becomes large (say N > 50), difficulties may be encountered because : (1) a solution of a large scale system with a full Jacobian matrix may present serious numerical stability problems as well as poor accuracy due to round-off error ; (2) the storage requirements for the Jacobian is of the order of N² which may be prohibitively large.

In the case of the extrapolation method, the numerical properties of the method are less apparent due to the nonlinearity of the ε - algorithm step used with this method. But it was reported by Brezinski [12] that the round-off error may cause severe numerical problems especially for stiff systems. Furthermore, it was shown previously that storage of order N² is required to minimize the execution time.

For the gradient method, the optimization step does not require solution of simultaneous set of equations (even with the variable metric method [22], the only matrix operations involved are addition or subtraction). Accordingly, there will be none of the numerical problems which may be encountered with the other methods. What concerns the storage (hence the order of the system N) problem, if a variable metric algorithm is used for the optimization, then storage of the order of N^2 is required just as for the other two methods. But if a conjugate gradients method could be used for the optimization, then storage of only order N will be required.

In conclusion, it can be said that the gradient method does not have the numerical stability problems such as those that may be encountered with both the Newton and the extrapolation methods for large scale systems. Furthermore, the three methods suffer from the same storage problem except that the gradient method would not have this problem if a conjugate gradients method (which does not involve matrices) could be used for the optimization.

CHAPTER III

A GENERALIZED APPROACH TO THE GRADIENT METHOD

3.1 Introduction

In Section 2.3 a gradient method for determining the periodic steady-state response, based on the state variable approach [10], [11], was outlined. Although an extension to the general algebraic-differential formulation of the system equations was also given [10], [11], it is still far from being directly implementable within most third generation time domain analysis programs which are based on the tableau formulation of the Network equations to take advantage of highly sophisticated sparse matrix techniques. In this chapter a more general and convenient formulation of the gradient method will be derived which can be implemented directly within most of the recent time domain analysis programs.

The derivation involves the use of the adjoint system concept. In Director's [6] treatment of the adjoint network, the parametric representation of branch types involved one dependency on a circuit variable as well as a set of design parameters, and the derivation of the adjoint network was carried out type-by-type. In the present approach to the derivation of the adjoint variational system, the equations describing the network are set in a general form, but with an assumed upper limit of two on the number of dependencies (linear or nonlinear). During the course of the derivation it is shown that dependencies within these limits are acceptable for representing the adjoint system, with the implicit assumption that the dependency can be represented as an element in the original network.

The approach in the present work is similar to that of Director [8] but the equation formulation is more general. After 'a` general adjoint system as well as a gradient form are developed, two useful objective functions, reflecting the state of the network relative to the steady-state are considered. One of these, based on normalizing the contribution of each reactive element into the weighted, square of voltage, is selected. On the basis of this choice a more specific approach is taken, which leads to convenient forms of the equations for computing the gradient, as well as for the initial conditions of the adjoint system. The chapter ends with a computational algorithm suited to application within a general time domain analysis program such as SCAMPER [19], [20].

3.2 Derivation of the Gradient Equations Based on a Generalized Adjoint Variational System

In order to minimize an objective function ϕ , which reflects how far a circuit is from the periodic steady-state, it is necessary to have an efficient procedure for computing the gradient $G = \partial \phi / \partial q_0^*$ where q_0^* is the initial condition vector of the differentiated variables representing the system. To this end, we formulate the original system of equations in the general tableau form.

45

h(u, q, t) = 0; Algebraic equations (3.la)

 $g(u, q, q, q_0, t) = 0$; Differential equations (3.1b)

where

is an algebraic variables vector , u

 $\int \int \left(\frac{u}{q} \right)^{\circ} = E(t)$

- is a differentiated variables vector , q
 - = $q(t_0)$ is the initial condition vector, q _

 $= q'(q_0, t)$.

In recent time domain, circuit analysis programs the solution of the system of equations (3.1) at any time point is achieved by linearizing the nonlinear terms at a point (the last accepted point) or a predicted one), and replacing the differentiated terms q by the slope of a polynomial of order $R \leq 6$ [28]. The linearized equations corresponding to (3.1) can be expressed in the form

(3.2)

where J is the coefficient Jacobian matrix, and E (t) is the forcing vector.

Before giving a detailed representation of J, it is important to specify limits on the acceptable range of branch types and We shall use the following equation to describe an eledependencies. ment or a part of a controlled (coupled) branch;

$$x_{i} = \alpha (x_{i}) \cdot x_{j}$$
(3.

where x may be an algebraic (u) or differentiated (q) variable. In this way branches ranging from linear self-dependent to nonlinear with two dependencies will be considered within the scope of our analysis*.

Returning to the linearization of the system of equations, equation (3.1b) can be set in the form

$$Mu + A \frac{dq}{dt} = 0$$

where M is a matrix of 0 or +1 entries,

and

 $\neg A = [a_{ij}]$, where

 $a_{ij} = 0, \pm 1$ for linear elements,

= a_{j} (q_j)₂ for a nonlinear single dependency,

(self or mutual) ,

= a (x), r \neq j for nonlinear double dependency.

From equations (3.1a) and (3.4) the linearized equations of the original system can be expressed in the form

* For more details about the different types of branches under consideration, see Part A of Chapter V.

46

3)

(3.4)

. In general, the objective function Φ , which reflects the deviation from the periodic steady-state, can be defined in an integral form

$$\Phi = \Phi (u, q, q_0, t_0) = \int_{t_0}^{t_0+T} \phi (u, q, q_0, t) dt \qquad (3.6)$$

If we choose Φ to be specifically the sum of the weighted squares of the differentiated variable differences, i.e.,

$$\Phi = [W (q (T) - q_0)]^{t} [q (T) - q_0] , \qquad (3.7)$$

where W is a constant weighting coefficient vector, then the integral form for Φ can be written as follows:

$$\Phi = \int_{0}^{T} \phi \, dt = \int_{0}^{T} \left[(q(t) - q_0)^{t} \frac{dq}{dt} \right] dt , \qquad (3.8)$$

where to

is taken to be zero for simplicity of notation. $\label{eq:simplicity}$

Since the variables u and q must satisfy the network equations (3.5), by using Lagrangian multipliers [8] λ_1 and λ_2 , equation (3.8) for Φ can be rewritten as

ິ 48

$$\Phi = \int_{0}^{T} [\phi (u, q, q_{0}, t) + \lambda_{1}^{t} h (u, q, t) + (\lambda_{2}^{t} M u) + \lambda_{2}^{t} A \frac{d q}{d t}] d t - \int_{0}^{T} (\lambda_{1}^{t}, \overline{\lambda_{2}^{t}}) E (t) dt$$
(3.9)

where λ_1 and λ_2 are assumed to be independent of the initial conditions q_0 . 0

The gradient G will then be

「流行学习」

まとうしていていた 他が発行

ſ

\If

$$G = \frac{d}{d} \frac{\phi}{q_0} = \frac{d}{d} \frac{T}{q_0} \int_{0}^{T} [\phi + \lambda_1^{t} h + \lambda_2^{t} M u + \lambda_2^{t} A \frac{d}{d} \frac{q}{d}] dt$$

$$= G T_1 + G T_2 + G T_3 + G T_4 . \qquad (3.10) \neq$$

$$G T_1 = \frac{d}{d} \frac{T}{q_0} \int_{0}^{T} \frac{\phi}{\phi} dt = \int_{0}^{T} [\frac{\partial}{\partial} \frac{\phi}{u} \frac{\partial}{\partial} \frac{u}{q_0} + \frac{\partial}{\partial} \frac{\phi}{q} \frac{\partial}{\partial} \frac{q}{q_0} + \frac{\partial}{\partial} \frac{\phi}{q_0}] dt$$
we use for ϕ the form given by (3.8) then

ţ

$$GT_{1} = \int_{0}^{T} \left[W \frac{d}{d} \frac{d}{t}\right]^{t} dt + \int_{0}^{T} \left[W \left(q \left(t\right) - q_{0}\right)\right]^{t} \cdot \frac{d}{dt} \left(\frac{\partial}{\partial} \frac{q}{q_{0}}\right) dt$$

which, after some straightforward manipulation, can be simplified to:

$$G T_{l_{T}} = [W (q (T) - q_{0})]^{t} \cdot \frac{\partial q (T)}{\partial q_{0}} - [W (q (T) - q_{0})]^{t}$$
(3.11a)

The second term G_{2}^{T} can be expanded as follows

$$T_{2} = \frac{d}{dq_{0}} \int_{0}^{T} \lambda_{1}^{t} h d t = \int_{0}^{T} \lambda_{1}^{t} \left[\frac{\partial h}{\partial u} + \frac{\partial u}{\partial q_{0}} + \frac{\partial h}{\partial q} + \frac{\partial q}{\partial q_{0}} + \frac{\partial h}{\partial q_{0}} \right] dt.$$

But since h is independent of q_0 , therefore, $\partial h / \partial q = 0$, whence $G \cdot T_2$ becomes

$$G T_{2} = \int_{0}^{T} \lambda_{1}^{t} \left[\frac{\partial h}{\partial u} \frac{\partial u}{\partial q_{0}} + \frac{\partial h}{\partial q} \frac{\partial q}{\partial q_{0}} \right] dt \qquad (3.11b)$$

The third term GT_3 is

からないというであるというない

G

$$T_{3} = \frac{d}{dq_{0}} \int_{0}^{T} \lambda_{2}^{t} M u d t = \int_{0}^{T} \lambda_{2}^{t} M \frac{\partial u}{\partial q_{0}} dt \qquad (3.11c)$$

In examining the fourth term G_4 , which depends on A, we will have to consider the acceptability of different types of nonlinear dependencies. We therefore proceed to reformulate G_4 component-wise, i.e., the inner products will be explicitly shown.

$$G T_{4} = \frac{d}{d q_{0}} \int_{0}^{T} [\lambda_{2}^{t} A \frac{d q}{d t}] d t = (I_{1}, I_{2}, ..., I_{k}, ..., I_{N}) \quad (3.12)$$

where I_k is the kth entry of GT_4 , k = 1, 2, ... N.

$$\begin{split} \mathbf{I}_{\mathbf{k}} &= \frac{\mathbf{d}}{\mathbf{d}} \frac{\mathbf{q}}{\mathbf{q}_{\mathbf{k}}} \int_{0}^{T} \left(\sum_{i=1}^{N} \mathbf{\lambda} \frac{\mathbf{d}}{\mathbf{d}_{t} \mathbf{t}} \right) \, \mathbf{dt} \\ &= \frac{\mathbf{d}}{\mathbf{d}} \frac{\mathbf{q}}{\mathbf{q}_{\mathbf{k}}} \int_{0}^{T} \left(\sum_{i=1}^{N} \sum_{j=1}^{N} \lambda_{2_{i}} \mathbf{a}_{ij} \left(\mathbf{x}_{i} \right) \frac{\mathbf{d}}{\mathbf{d}t} \frac{\mathbf{d}}{\mathbf{d}t} \frac{\mathbf{d}}{\mathbf{q}} \frac{\mathbf{d}}{\mathbf{q}_{\mathbf{k}}} \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{N} \left(\int_{0}^{T} \left(\lambda_{2_{i}} \mathbf{a}_{ij} \left(\mathbf{x}_{i} \right) \frac{\mathbf{d}}{\mathbf{d}t} \frac{\mathbf{d}}{\mathbf{d}t} \frac{\mathbf{d}}{\mathbf{d}q_{\mathbf{k}}} \frac{\mathbf{d}}{\mathbf{q}_{\mathbf{k}}} \frac{\mathbf{d}}{\mathbf{q}_{\mathbf{k}}} \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{N} \left(\sum_{i=1}^{N} \left(\lambda_{2_{i}} \frac{\mathbf{d}}{\mathbf{d}t} \mathbf{d}t \right) \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{I}_{\mathbf{k}1} + \mathbf{I}_{\mathbf{k}2} \right) , \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\lambda_{2_{i}} \mathbf{a}_{ij} \left(\mathbf{x}_{i} \right) \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\lambda_{2_{i}} \mathbf{a}_{ij} \left(\mathbf{x}_{i} \right) \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\lambda_{2_{i}} \mathbf{a}_{ij} \left(\mathbf{x}_{i} \right) \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\lambda_{2_{i}} \mathbf{a}_{ij} \left(\mathbf{x}_{i} \right) \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\lambda_{2_{i}} \mathbf{a}_{ij} \left(\mathbf{x}_{i} \right) \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\lambda_{2_{i}} \mathbf{a}_{ij} \left(\mathbf{x}_{i} \right) \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\lambda_{2_{i}} \mathbf{a}_{ij} \left(\mathbf{x}_{i} \right) \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\mathbf{d}_{i} \mathbf{d}_{i} \mathbf{d}_{i} \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\mathbf{d}_{i} \mathbf{d}_{i} \mathbf{d}_{i} \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\mathbf{d}_{i} \mathbf{d}_{i} \mathbf{d}_{i} \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\mathbf{d}_{i} \mathbf{d}_{i} \mathbf{d}_{i} \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\mathbf{d}_{i} \mathbf{d}_{i} \mathbf{d}_{i} \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\mathbf{d}_{i} \mathbf{d}_{i} \mathbf{d}_{i} \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\mathbf{d}_{i} \mathbf{d}_{i} \mathbf{d}_{i} \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\mathbf{d}_{i} \mathbf{d}_{i} \mathbf{d}_{i} \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\mathbf{d}_{i} \mathbf{d}_{i} \mathbf{d}_{i} \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\mathbf{d}_{i} \mathbf{d}_{i} \mathbf{d}_{i} \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\mathbf{d}_{i} \mathbf{d}_{i} \mathbf{d}_{i} \right) \, \mathbf{dt} \\ &= \sum_{i=1}^{N} \sum_{j=1}^{T} \left(\mathbf{d}_{i} \mathbf{d}_{i} \mathbf{d}_{i} \right) \, \mathbf{dt} \\ &= \sum_{i=1}$$

Ţ

C.

ſ

 $I_{k} = \sum_{i=1}^{N} \sum_{j=1}^{N} [\lambda_{2_{i}}^{x} a_{ij}(x_{r}) \frac{\partial q_{j}}{\partial q_{0_{k}}}]_{0}^{T}$ $- \int_{0}^{T} \frac{d \lambda_{2_{i}}}{d t} a_{ij}(x_{r}) \frac{d q_{j}}{d q_{0_{k}}} dt$ $+ |\int_{0}^{T} \lambda_{2_{1}} \frac{\partial a_{ij}(x_{r})}{\partial x_{r}} (\frac{\partial x_{r}}{\partial q_{0_{k}}} \frac{\partial q_{j}}{\partial t} - \frac{\partial x_{r}}{\partial t} \frac{\partial q_{j}}{\partial q_{0_{k}}}) dt]$ (3.14)

51

(3.11d)

In the case of a single dependency, i.e., $x_r \rightarrow q_r$ and $r \rightarrow j$, the term $\frac{\partial x_r}{\partial q_0} = \frac{\partial q_j}{\partial t} = \frac{\partial x_r}{\partial t} = \frac{\partial q_j}{\partial q_0}$

is identically zero. For a double dependency, where $x_r \neq q_j$, it is shown in Appendix I that this term is identically zero for a class of systems with a unique periodic steady-state solution. Hence, using this fact, and equations (3.12) and (3.14)

$$\mathbf{F} \mathbf{T}_{4} = \lambda_{2}^{t} (\mathbf{T}) \mathbf{A} (\mathbf{T}) \frac{\partial \mathbf{q} (\mathbf{T})}{\partial \mathbf{q}_{0}} - \lambda_{2}^{t} (\mathbf{0}) \mathbf{A} (\mathbf{0})$$
$$- \int_{0}^{T} \left[\frac{d \lambda_{2}^{t}}{d t} \mathbf{A} \frac{\partial \mathbf{q}}{\partial \mathbf{q}_{0}} \right] dt$$

From equations (3.11a) - (3.11d), the gradient can take the form

$$G = \int_{0}^{T} \left[\lambda_{1}^{t} \frac{\partial}{\partial u} + \lambda_{2}^{t} M\right] \frac{\partial}{\partial q} \frac{u}{q_{0}} dt$$

$$+ \int_{0}^{T} \left[\lambda_{1}^{t} \frac{\partial}{\partial q} - \frac{d}{dt} \frac{\lambda_{2}^{t}}{dt} A\right] \frac{\partial}{\partial q} \frac{q}{q_{0}} dt$$

$$+ \left[\lambda_{2}^{t} (T) A (T) - (W (q_{0} - q (T)))^{t}\right] \frac{\partial}{\partial q_{0}} \frac{q}{q_{0}}$$

$$+ \left[W (q_{0} - q (T))\right]^{t} - \lambda_{0}^{t} (0) A (0) \qquad (3.15)$$

Since the Lagrangian multipliers can be chosen arbitrarily, they can, in particular, be specified in such a way as to greatly simplify the form of the gradient equation, e.g., take λ_1 and λ_2 such that both the first and second terms in (3.15) are identically zero.

 $\lambda_{1}^{t} \frac{\partial h}{\partial u} + \lambda_{2}^{t} M = 0$

ここのであるま

and

$$\lambda_1^{t} \frac{\partial h}{\partial a} - \frac{d \lambda_2^{t}}{d t} A$$

In matrix form, it can be replaced by

52

(3.16)

L.F.

٢¥

Furthermore, and to avoid the computation of the term $\partial q'(T) / \partial q_0$ (which is, in general, a full matrix), the initial conditions of λ_2 , can be chosen so that the third term of equation (3.15) is identically zero, i.e.,

$$\lambda_{2}^{t}$$
 (T) A (T) - [W (q₀ - q (T))]^t = 0 7 (3.17)

Accordingly, the gradient takes the simplified form

 $G = [W (q_0 - q(T))]^{t} - \lambda_2^{t} (0) A (0)$

Equations (3.16) - (3.18) are the results that we If we compare equations (3.5) and (3.16), the coefficient sought. Jacobian matrix of the latter is just the transpose of the Jacobian of (3.5) with a change in the sign of the time derivative term. (3.16) is called the adjoint variational system of the original System one described by (3.5). Relation /(3.17) defines the initial conditions of that adjoint system for computing G. Note that these conditions are given at t = T, which means that to satisfy equation (3.18), the integration of the adjoint system has to be carried out over one full period from t' = T to t = 0, i.e., in the backward direction. It should also be noted that, due to this reversal of the direction of integration, the sign of the differentiated terms in (3.16) will be restored to the original sign. Let us define a new time base τ , for the adjoint system where $\tau = T - t$, so that $d/dt = -d/d\tau$.

53

(3.18)

In that case (3.16) takes the form

$$\begin{bmatrix} \frac{\partial}{\partial u} & \frac{\partial}{\partial q} \\ M & A \frac{d}{d\tau} \end{bmatrix}^{t} \begin{bmatrix} \lambda_{1} \\ \lambda_{2} \end{bmatrix}_{t} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} .$$
(3.19)

In summary, (3.19) defines the adjoint system which,

when integrated over one full period, starting with initial conditions at $\tau = 0$ (i.e., t = T) given by equation (3.17), has a final solution at $\tau = T$ (i.e., t = 0) in which λ_2 defines the second term of the gradient equation (3.18).

3.3 Choice of the Objective Function Φ

そうちょう シート・シート からちょう こうちょうちょう ないない

If the objective function Φ is to be a meaningful measure of the circuit's deviation from the periodic steady-state, then the components of Φ , corresponding to the contributions of the individual differentiated variables, should, in general, be of the same order of magnitude. However, the representation of the branch variables (i, v, q, f) in the Kirchhoff's laws equations in SCAMPER depends on the type of dependency and on whether the branch is nonlinear. Thus, the current in a q-dependency (capacitive) branch is replaced by dq / dt if the element is linear, and by c dv / dt if it is nonlinear. Because the units of the differentiated variables differ, their values

may also differ by many orders of magnitude. Therefore, it is necessary to weight the contribution of each storage element to the performance function Φ .

One particularly convenient way of achieving this weighting is to convert all the branch variables to voltages when computing ϕ Thus, for example, the variable associated with a linear capacitive branch, namely q_i , would be converted to a voltage v_i by defining $v_i \equiv q_i / C_i$. An objective function which is based on this approach and which has the same form as equation (3.7) is

$$\phi = \sum_{i=1}^{N} \phi_{i} = \sum_{i=1}^{N} \omega_{i} \cdot \frac{1}{2} (v (T) - v_{0})_{i}^{2}$$
(3.20)

where N is the system order, and ω_{1} is an additional weighting factor i for reactive element i.

An alternative form of Φ , which also tends to reduce the spread in the contributions of the storage elements to the objective \tilde{z} function, is the sum over the stored energy, i.e.,

$$\sum_{i=1}^{N} \phi_{i} = \sum_{over i} \omega_{i} \cdot \frac{1}{2} (v (T) - v (0))_{i} (q (T) - q (0))_{i}$$
all C

$$\sum_{\text{over}} \omega_{j} \cdot \frac{1}{2} (i (T) - i (0)) (f (T) - f (0)) j$$

where f is the inductor flux.

all I

55

(3.21)

It can, in fact, be shown that (3.21) and (3.20) can be made equivalent by a suitable choice of the ω 's. Nevertheless, equation (3.20) is more convenient to use and was selected for the work reported here.

Detailed Computations of the Gradient and the Initial Conditions of the Adjoint System

The initial conditions λ (T) of the adjoint system are defined by equation (3.17), while the gradient G can be computed through equation (3.18). Both equations are given in matrix form, and appear to be difficult to handle. In order to see how easy it is, in fact, to compute λ_2 (T) and G, it is helpful to review the actual representation of the system variables and the order in which the equations are set in the dc and transient analysis program SCAMPER [19], [20] which was used for this work. (For more details see Part A of Chapter V).

Reactive Variable Representation:

d q d t

The current through a q-dependency (capacitive) element in Kirchhoff's current law equations is replaced by:

(i)

if the element is linear (self or mutually dependent), i.e., the differentiated variable is charge q.

(ii) $c \frac{d v}{d t}$ if the branch is nonlinear, where v is often a voltage, i.e., the differentiated variable is a voltage.

in Kirchhoff's voltage law is replaced by:

(i) $\frac{d f}{d t}$ if the element is linear.

(ii) $L\frac{di}{dt}$ if the element is nonlinear, where i is often a current.

The Order of the System Variables and Equations

「たいいいい

Consider the number of elements is NE, and the number of nodes NN. In the tableau representation there are (2NE + NN - 1)equations and variables. The equations are ordered in the following way: first, the branch Kirchhoff's voltage law equations KVL (1 to NE); second, the nodal Kirchhoff's current law equations KCL (NE + 1 to NE + NN - 1); third, the branch constitutive relations E.C. (NE + 1 to NE + NN - 1); third, the branch constitutive relations E.C. (NE + NN to 2NE + NN - 1). The variables appear in the order of : first, branch current or charge (1 to NE); second, branch voltage or flux (NE + 1 to 2 NE); third, node voltages (2NE + 1 to 2 NE + NN - 1).

Now, a close examination of equations (3.17) and (3.18), in the light of the above comments, reveals a very simple procedure for computing the different components of the gradient G of the performance

function defined by equation (3.20) as well as the initial conditions λ (T). The following relations for λ (T) and G can be proved to be equivalent to those of equations (3.17) and (3.18).

For a linear capacitor C,

$$\phi_{c_{\ell}} = \frac{1}{2} \cdot \omega_{c_{\ell}} \cdot \frac{1}{c_{\ell}^{2}} (q_{0} - q^{(T)})_{c_{\ell}}^{2} \qquad (3.21a)$$

The initial conditions λ (T) in the solution of the adjoint system

 $\lambda (T) - \lambda (T) = \omega_{c_{\ell}} \frac{1}{c_{\ell}^2} (q_0 - q(T))_{c_{\ell}}$ (3.22a) NE+N1 NE+N2 ℓc_{ℓ}

where N1 and N2 are the numbers of nodes to which the branch is

connected.

will be

The corresponding gradient G

$$G_{c_{\ell}} = C_{\ell} \frac{d \Phi}{d q_{0}} = \omega_{c_{\ell}} \frac{1}{C_{\ell}} (q_{0} - q (T))_{c_{\ell}} - C_{\ell} (\lambda^{\prime} (0) - \lambda (0))_{NE+N1}$$

$$NE+N2$$

$$(3.23a)$$

while the argument passed to the optimization routine x_{c_1}

$$\mathbf{x}_{\mathbf{c},\boldsymbol{\ell}} = \frac{1}{C_{\boldsymbol{\ell}}} \mathbf{q}_{\mathbf{0}_{\mathbf{c},\boldsymbol{\ell}}}$$

For a nonlinear capacitor C_n (•) :

58

(3.24a)

$$\begin{aligned} & \oint_{C_n} = \frac{1}{2} u_{C_n} (v_0 - v_1(T))_{C_n}^2 , \qquad (3.21b) \\ & \lambda_n(T) = \lambda_n(T) = u_{C_n} \frac{1}{C_n}(T) (v_0 - v_1(T))_{C_n} , \qquad (3.22b) \\ & \text{Where } C_n(T) \text{ is the value of } C_n \text{ at } t_1 = T + t_0 . \\ & G_{C_n} = \frac{4}{4} \frac{\phi}{v_0} = u_{C_n} (v_0 - v_1(T))_{C_n} - C_n(0) (\lambda_n(0) - \lambda_n(0)) \\ & \text{RE-NL NE+N2} \\ & G_{C_n} = \frac{4}{2} v_{C_n} (v_0 - v_n(T))_{C_n} - C_n(0) (\lambda_n(0) - \lambda_n(0)) \\ & \text{RE-NL NE+N2} \\ & (3.23b) \end{aligned}$$
where $C_n(0)$ is the value of C_n at $t = t_0$, $(3.24b)$
For a linear inductor L_k :
 $& \phi_{1_k} = \frac{1}{2} \cdot u_{1_k} \cdot (\frac{R}{4_k})^2 (t_0 - f_n(T))_{1_k}^2 , \qquad (3.21c)$
where R is a weighting resistance.
The initial condition $\lambda_n(T)$ in the solution of the adjoint system
 $& \lambda_{1_k}(T) = -u_{1_k} (\frac{R}{4_k})^2 (t_0 - f_n(T))_{1_k} , \qquad (3.22c)$
where i is the branch number $(1 \le i \le NE)$.

 $\|$

۰ ٥

، ۲.

ĩ

The corresponding gradient G

$$G_{l_{\ell}} = \frac{L_{\ell}}{R} \frac{d \phi}{d f_{0}} = \omega_{l_{\ell}} \left(\frac{R}{L_{\ell}}\right) \left(f_{0} - f(T)\right)_{l_{\ell}} - \frac{L_{\ell}}{R} \lambda_{i} \left(0\right)$$
(3.23*q*)

while the argument passed to the optimization routine X

$$x_{l_{\ell}} = \frac{R}{L_{\ell}} f_{0_{l_{\ell}}}$$
(3.24c)

ł

For a nonlinear inductor L_n (•) :

$$\phi_{1} = \frac{1}{2} \omega_{1} R^{2} (i_{0} - i (T))_{n}^{2} , \qquad (3.21d)$$

$$\lambda_{i}(T) = \omega_{1} \frac{R^{2}}{n} \frac{L_{i}(T)}{L_{n}(T)} (i_{0} - i(T))_{1}$$
(3.22d)

where L_{n} (T) is the value of L_{n} at $t = T + t_{0}$ $G_{1_{n}} = \frac{1}{R} \left(\frac{d}{d} \frac{\phi}{1_{0}} \right)_{n} = \omega_{1} R \left(i_{0} - i_{0} (T) \right)_{n} - \frac{L_{n}(0)}{R} \lambda_{i} (0)$ (3.23d)

where $L_n(0)$ is the value of L_n at $t = t_0$,

and

·(3.24d)

1

3.5 A Computational Technique

The proceeding formulation of the gradient method has been implemented in the dc and transient analysis program SCAMPER by means of the algorithm outkined below.

Algorithm 3.1

1. Initialization Step.

(i) Define constants (minimum limit of Φ (Φ) and G (G) and maximum number of iterations NIT (NIT_{max})).

(ii) Integrate the original system of equations
 (3.5) until a predefined time t₀, yielding
 a suitable set of initial conditions for the steady-state analysis.

(iii) At $t' = t_0$ save the initial conditions of the reactive variables q, the solution vector (u, q) and the values of the nonlinear reactive elements (for computing matrix A (0)).

(iv) NIT =

いたいないないないであっているとうないないである

2. Forward Integration.

(i)

NIT =

NIT + 1

£ *

(ii) Integrate the original system (3.5) over one full period T from $t = t_0$. Save the necessary nonlinear entries of the Jacobian J as well as the step size at a sequence of time points.

(iii) At $t \approx t_0 + T$ compute q (T) and the values of nonlinear reactive elements (for computing A (T)). $-\frac{4}{2}$

- (iv) Compute the performance function I defined by
 equation (3.20), using equations (3.21a, b, c,
 and d) .
- (v) If $\Phi \leq \Phi$, a steady-state solution has been found. STOP.
- (vi) IF NIT > NIT \max . STOP. (vii) Compute the initial conditions λ_2 (T) for the adjoint system defined by equation (3.17) using equations (3.22a, b, c, and d). Also compute the first term of the gradient defined by equation (3.18) using equations (2.23a, b, c, and d).

3. Backward Integration.

With the initial conditions λ (T) computed in

(i)

step (2. vii), integrate the adjoint system
of equations (3.19) over one full period
using the stored values of the nonlinear en, tries of the Jacobian as well as the corresponding step sizes.

At the end of the period $\tau = T$ (i.e., $t = t_0$) use the solution of the adjoint system λ (0) and A (0) for completing the computation of the gradient (second term) using equations (3.23a, b, c, and d).

. Optimization.

(i)

(ii)

Scale the argument vector X according to equations (3.24a, b, c, and d).

(ii) Call the optimization routine.

(iii) Reset the initial conditions of the original
system (u, q) to the values saved fromstep (1. iii).

(iv)

いたがたいないないであるというない

Update the reactive element variables q with the values q_0 returned by the optimization routine.

· √` (v)

Remarks

(1)

The optimization routine returns a new set of initial conditions q_0 . Due to the nonlinearity of the system, and the change in q_0 (which may be large relative to the previous values of q_0) the system may fail completely to satisfy the nonlinear error criteria at the starting point of integrating the original system. A way to overcome this problem is to start the integration with values of q_0 which satisfy the system equations at t_0 . This is the reason for step (1. iii). In fact, the only values which it is necessary to save for this step are the components of (u, q) corresponding to the nonlinear branches.

(2)

Steps (2. iv) and (2. vii) can be computed simultaneously.

(3)

Some optimization techniques consider a minimum to have been found (which in our application means a steady-state solution is found) if the gradient
vector |G| (or its norm $||G||_2^2$) is less than a certain limit. Hence another termination criterion can be added in this case.

(4)

For more details about resetting the initial conditions as well as the integration process, see Part B of Chapter V.

CHAPTER IV

NUMERICAL CONSIDERATIONS AND RESULTS

4.1 Introduction

The periodic steady-state analysis algorithms 3.1 was implemented in the dc and transient analysis program SCAMPER. For the optimization step, two different optimization routines were made available to the user. These routines are Fletcher's [22] variable metric routine (VM), which requires full matrix storage, and an automatic restart, conjugate gradient routine (ARCG) by Powell [23] which requires less storage.

The convergence of the algorithm as well as the execution time are greatly influenced by the scaling of the different arguments. In Section 4.2 the scaling features are discussed briefly and some suggestions are made, based on our limited experience with the program, regarding the selection of a proper scheme for scaling the arguments.

Three examples are presented in Section 4.3. These were solved using both the VM and the ARCG optimization routines, however, since the VM routine gave better convergence, most of the detailed results are given for that routine and only some final results are shown for ARCG for comparison. An important result that comes out of these examples is that the time required for integrating the adjoint system of equations can be controlled and is usually less than half the time required for integrating the original system of equations describing the network.

Although these, same examples were also analyzed by Nakhla [11] using a similar gradient method and the same VM routine, a detailed comparison of results is, in general, not meaningful because of the different strategies used for the transient analysis. In our case, the transient analysis starts after the dc solution is found, while in Nakhla's program the transient analysis starts from zero initial conditions.

4.2 Scaling

「こうない」というない ない、ない、ないないないない、こういいいいいいいいいい

The results obtained for several problems demonstrate clearly that the optimization step of the algorithm is the most critical one. The available optimization routines, VM and ARCG, may require a large number of iterations to reach the region of quadratic convergence, depending on the starting point. The convergence rate as well as the execution time can be greatly affected by the choice of the scaling parameters, so that it is important to understand their respective functions in order to use them properly. The four kinds of scaling which are available are discussed below and some guide-lines are given for their proper use.

I - The equivalent resistance, RLS, for the inductive elements:

The initial conditions, q₀, corresponding to the dif-

ternally to be equivalent to current sources $I_{0_{\chi}}$, while those corresponding to the capacitive elements are made equivalent to voltage "sources $V_{0_{C}}$. The equivalent resistance RLS is introduced to scale the inductive variables by converting them to voltages $V_{0_{\chi}} = \text{RLS} \cdot I_{0_{\chi}}$. A choice of the value of RLS other than unity may be necessary in cases where there are order of magnitude differences between the V's and I's of the C's and L's, respectively. From the computational point of view, the equivalent resistance RLS scales the objective functions of the inductive elements ϕ_{χ} by (RLS)² and the variables passed to the optimization routine, X_{χ} , by RLS, while the gradient will be scaled indirectly by RLS and (1 / RLS). (See equations 3.23c and d).

II - Weighting of the objective function ϕ_i :

The weighting factor ω_i of the objective function ϕ_i given by equation (3.20) is specified via the parameter FCTR_i = ω_i . This kind of weighting is useful for controlling the influence of each storage element's deviation from the periodic steady-state on the conversion process. For example, the FCTR_k of ϕ_k corresponding to element k associated with a longer time constant, and yet probably having a relatively small initial difference between q_0 and q (T), may be increased to advantage w. r. t. the other FCTR_i.

This kind of weighting affects both the objective function , and therefore ϕ , and the gradient vector G .

III - Scaling the optimization variables X : .

SCLX (i) = SCLX, i = L, 2, ... N.

In the steady-state analysis of circuits, cases are often encountered where the values of the variables differ by orders of magnitude. This may result in poor convergence due to the search for a minimum over an elongated contour of the performance function Φ with respect to these variables. Convergence may be improved by scaling the variables to be of the same order of magnitude. Care must be taken in using this type of scaling as the relative magnitudes of the variables may change drastically as the search for a minimum Φ proceeds. Four options are available in using SCLX.

Individual scaling of each element. In this
case M values have to be provided. If
M < N, the first. M variables are scaled in
'the given order of their corresponding branches,
by the specified M values, while the remaining N - M variables will be scaled by unity.</pre>

(ii)

(i)

Common scaling for all C's and L's, respectively. In this case two values, SCLX and SCLX, must be provided where all the capacitive variables will be scaled by the first value, SCLX , and all the inductive variables by the

second, $SCLX_{l}$. This particular form of scaling is simple to use and yet quite important. If, for example, inductor currents are of the order of milliamps while the capacitor voltages are of the order of volts, the ratio $SCLX_{l}$ / $SCLX_{c}$ could be, say, 1000.

All the variables to be scaled in such a way as to cause their scaled values to be equal at the first steady-state iteration. In this case a parameter, XREF, must be supplied whereupon the variables will be scaled internally to be equal to XREF for the first iteration of the steady-state analysis.

(iv) No scaling at all.

(iii)

IV - Scaling the initial conditions of the adjoint system SCLD :

If this scaling is required all the initial conditions of the adjoint system, λ_2 (T), will be scaled equally in such a way as to cause the initial condition with the largest equivalent voltage to be equal to SCLD at the beginning of each integration of the adjoint system. The reason for introducing this type of scaling is that the gradient may not have to be computed as accurately during the

first few iterations as within the quadratic region of convergence. Hence, if the parameter SCLD is chosen to be small enough relative to the largest equivalent voltage of λ_2 (T) at the first iteration and large enough relative to that same component of λ_2 (T) at the last iteration (which is normally within the specified error of computation), the average time for integrating the adjoint system may be reduced. Furthermore, the higher accuracy during the last few iterations may improve the convergence.

On the basis of our admittedly limited experience with the program, the following guidelines are proposed for using the steadystate analysis effectively.

(i)

Perform a transient analysis over few periods. Examine the reactive (differentiated) variables as well as those corresponding to sharp nonlinearities, such as exponentials, over the last period and choose the starting point t_0 in such a way that (1) regions where the sign of the time derivatives is changing rapidly are avoided, (2) it does not correspond to the region of conduction of junctions shunted by reactive elements (if any), and (3) as many variables as possible are near their largest absolute values. Choose the equivalent resistance RLS, corresponding to the inductive elements (if any), in such a way as to cause the average equivalent voltage of the inductive elements to be comparable to the average voltage of the capacitive elements.

72

Weight each objective function ϕ_1 according to its expected influence on the transient response of the circuit (taking into account the effect of RLS).

Scale the arguments passed to the optimization routine in such a way as to make them all approximately equal (taking into account the effect of RLS) .

4.3 Numerical Results

Three examples are presented to illustrate the effectiveness of the proposed steady-state algorithm 3.1, and to demonstrate the importance of using scaling. Besides the scaling parameters introduced in the preceding section, the following notations will be used in these examples

(ii)⁻

(iii)

(iv)

TF, TB

the time required for one full period integration of the original and the adjoint system, respectively ;

NP

the equivalent total number of forward integrations, required for completing the steadystate analysis, i.e., NP = TT / TF, where TT is the total time to convergence.

Example 1

The importance of using a steady-state analysis algorithm rather than using continuous transient analysis ("brute force" method) is demonstrated by this simple example which was also analysed by Nakhla [11]. The clipper circuit shown in Figure 4.1 was analysed twice, once using the gradient method and once by the "brute force" In fact, with the dc source EC = 5 volts, the diode never technique. conducts so that the circuit is effectively a linear one. The circuit has two time constants T_1 and T_2 which, with the diode off, can be simply shown to be given by : $T_1 \approx C_1 (R_1 R_2 / (R_1 + R_2))$ and The ratio T_0 / T_1 is approximately 50, while $T_2 \approx C_1 (R_1 + R_2).$ the ratio of the longer time constant T_2 to the period of the input signal is approximately 100 . Hence, using the brute force method, the steady-state condition is not expected to be reached before a few hundred full periods of integration. " This is confirmed by the results



shown in the first two columns of Table 4.1. Note that with the steadystate algorithm, Φ is reduced to less than 10^{-9} in two iterations (equivalent to approximately 2.64 periods of forward integration) whereas 49 periods by the brute force method reduces Φ to only 0.234 x 10^{-7} .

Example 2

The circuit shown in Figure 4.2 has been considered by both Nakhla [11] and by Trick et al [1] who showed that the system was still far from steady-state after 75 full-period integrations. With the gradient algorithm we obtained convergence to a Φ of 10^{-12} in 16 iterations with $t_0 = 3.5$ T. The only scaling used was that of the adjoint system initial conditions, SCLD. With SCLD equal to 0.01 the convergence within the quadratic region was improved and the execution time was reduced by about 10%. In Figure 4v3 the ratio (TE / TF) for both the scaled and the unscaled cases are plotted against the iteration number. With scaling the average (TE / TF) was 463 while without scaling it was 51%.

Example 3

•The importance of applying the different types of scaling in using the gradient method is shown in the analysis of the class C amplifier of Figure 4.4. Starting the analysis after 4.9 periods, and with scaling applied as shown in Table 4.3, the steady-state

11-22-24-24-24-24-24-24-24-24-24-24-24-24-	1			
RESULTS OF	EXAMPLE	1	0	
•				
Brute Force Method	1	Gradient Met	hođ	
Period Φ	Iteration	ι Φ	×ı	*2
0 0.236	0	0.236	0.1	0.1
1 0.467 x 10^{-1}	l I o	0.992×10^{-8}	-1.1513	0.10156
2 0.925 $\times 10^{-2}$	2	0.477×10^{-9}	-1.1515	0.10156
10. 1.00×10^{-7}				•
20 0.415 x 10^{-7}				•
30 0.34 x 10^{-7}	•		a	-
40 0.28 x 10 ⁻⁷			÷.,	1
49 0.234 x 10^{-7}				
í, í		·		•
For the gradient method :-	•••	<u> </u>		
SCLX 0,5 10.	4			1
$(x_1, x_2)_{\text{final}}$: (-1.1515, 0.	10156)	-		
Execution time = 5.37 sec.	° `			
(TE / TF) average = 32 % .	h	· ·]		
NP; = 2.64.	-		75	,
	. i	-	. e	z j t
	- for a		-	
·	\		1	
	· · · ·			- ·

こうにはないないないで、このないないで、

TABLE 4.1



TABLE 4.2 EXAMPLE 2 RESULTS OF Scaling : SCLD = 0.01.Iteration Iteration 0.995×10^{-3} 4:26 9 0.490×10^{-4} 0.227 10 0.153 x 10⁻⁵ 0.226 11 2 0.173×10^{-6} 0.222 3 12 0.122 x 10⁷⁷ 0.212 13 4 0.479×10^{-9} 0.118 14 5 0.489×10^{-11} 0.111 15 0.232×10^{-12} $0.254 \times 10^{+1}$ **'16**, $0.892 \times 10^{+2}$ 8 $[x_1, x_2, x_3, x_4] = [-9.07535, 9,05648, 9.10251, 9.029 \times 10^{-3}]$ Execution time = 45.7 sec. (TB / TF) average = 46 NP

K.

()



conditions were found after 23 iterations accurate to six digits with $\Phi < 10^{-12}$. The two sets of values shown for x in Table 4.3 correspond to the final steady-state solution at t = t₀ = 4.9T and also at t = 5.T. The latter are given to permit comparison with the results of Trick et al [1] and Nakhla et al [10], [11].

When this problem was solved with the same RLS, FCTR SCLX but without SCLD, the method converged to $\Phi = 0.819 \times 10^{-11}$ and after 24 iterations with 118.4 sec execution time. This corresponds to an increase of about 15 % in the execution time due to the omission of SCLD. Figure 4.5 shows (TB / TF) as a function of iteration number without SCLD and with SCLD equal to 0.01. For the unscaled case, the increase of the execution time during the first few iterations was not balanced by the reduction during the last iterations resulting in an increase of the average (TB / TF) over that for the case where SCLD was used properly. In the absence of any kind of scaling, the method was not able to converge within the same number of iterations. After 28 iterations, and 150 secs. of exectuion time the objective function Φ was reduced to only 0.45 x 10⁻⁴. It is worth noting that when the same problem was solved, for comparison, using the conjugate gradient optimization routine, the algorithm was not able to converge and, after 28 iterations, the objective function reduced only to 0.19×10^{-2} which is very far from the steady-state condition.



يه بر ^مير ، د ^{ال}مسي



82

Scaling :

دیس

RLS = 50.

FCTR_i, i = 1, 2, ..., 10 : 1., 1., 1., 0.1, 0.1, 0.1, 1., 1., 1., 1. SCLX_i, i = 1, 2, ... 10 : 0.5, 0.5, 1., 0.1, 0.1, 0.1, 1., 1., 1., 1. SCLD : 0.01

Iteration	φ	Iteration	Φ	Iteration	х ф
0	0.138	_ B	0.307×10^{-3}	l 16	0.345×10^{-4}
1	0.329	9	0.284 x 10 ⁻¹	l. 17	0.114×10^{-4}
2	0.973×10^{-1}	10	0.272×10^{-1}	18	0.496×10^{-5}
3	0.718×10^{-1}	11	0.262×10^{-1}	19	0.546×10^{-6}
4 ~	0.532×10^{-1}	, 1 2	0.180×10^{-1}	20	0.331×10^{-7}
5	0.416×10^{-1}	13	Q.133 x 10 ⁻¹	21	0.691 x 10 ⁻⁹
6.	0.338×10^{-1}	14	0.463 x 10 ⁻²	22	0.105 x 10 ⁻¹⁰
7 '	0.315×10^{-1}	15	0. <u>6</u> 37 x 10 ⁻³	23	0.234×10^{-12}

(x_i, i = 1, 2, ..., 10) 5.0 1.4530 1.41695 -0.288972 28.3779 25.5844 , 29.4632 , -0.076572 , -0.062961 , -0.076349 , -0.149948 26.8711 28.8556 -0.079830 0.000141 -0.076101 -0.098802 Execution time -= 104.73 sec.

 $(TB / TF)_{average} = 48'$

NP = 39.5.



CHAPTER V. PROGRAMMING

PART A

DC AND TRANSIENT ANALYSIS PROGRAM SCAMPER

The transient analysis program SCAMPER, to which the periodic steady-state analysis has been added was written by Millar and Blostein [19], [20]. Its main features are accuracy, high speed and a capability of handling a wide range of branch types. Circuits of about 350 branches (which is roughly equivalent to a circuit of 20 transistors) can be analyzed.

The analysis starts by generating the dc solution. All the independent sources are set to their values at zero time and all the time derivatives are removed from the network equations, (which is equivalent to short circuiting all the inductors and open circuiting all the capacitors). Then, if required, a transient analysis is carried out over a specified time interval.

Techniques such as high order implicit integration with variable step size, the use of sparse matrices and the generation of executable machine code account for SCAMPER's high speed and accuracy. The program is written mainly in FORTRAN IV.

Kn (EPAKK)

85

5.1 Acceptable Element Types

There are two general types of branches which are

accepted by SCAMPER.

(i) Sources:

For current source

J = J(t)

For voltage source

E (t)

(5.1b)

(5.la)

where the time dependent function can be any of the following functions. (1) Constant (dc source). (2) Sine wave (magnitude A_0 , starting time t_0 , frequency f). (3) Pulse (general trapezoidal periodic source). (4) Tabulated source (general piece-wise linearization of non-periodic source).

(ii) Controlled Branches

These are the branches that can be described by either a relation between network variables (voltage, current, etc.), or by a relation between element values (resistance, etc.) and network variables.

Using the conventional network definitions, where X is a network variable (voltage V, current I, charge Q, flux F), the general acceptable relations are

Through this general representation, types ranging from constant, linear, self-dependent elements up to cross coupled non-linear elements are acceptable. The only restrictions on the type of dependency are that dependency is not permitted on non-representable variables such as a current through a capacitor or a voltage across an inductor, because these are not treated as equation variables. Also the nonlinearities are restricted to be of the piece-wise linear form, with the exception of a current-or charge-controlled branch which can be of the exponential type as well.

5.2 Order of Setting the Equations and Variables

The tableau representation of the network is set in the following order. First, the branch Kirchhoff's voltage law equations (KVL), in order of increasing branch number. Second, the nodal Kirchhoff's current law equations (KCL), in order of increasing node number.

Third, the branch constitutive equations (B.C), in order of increas-

C

The variables represented in the network equations, X, are not the network variables themselves, but the negated differences between their values at a new point (i) and the last accepted point (i-1), i.e.

where X is the equation variable, and Y is the corresponding network variable.

The arrangement of variables is in the following order. First, branch current or charge (I or Q), in order of increasing branch number. Second, branch voltage or flux (V or F), in order of increasing branch number. Third, node voltage (VN), in order of increasing node number.

5.3 Solution, Procedures and Numerical Techniques

 $X = Y_{i-1}^{i} - Y_{i},$

The transient analysis of the network is done by first linearizing the branch equations at a point, this being the last accepted point or a predicted one. The differentiated variables are replaced by a backward difference formula (integration formula) of a

(5.3)

user-specified order k. The linearized system of equations may take the form

where J is the Jacobian coefficient matrix, E is the forcing vector, and X is the equation variable vector. Equation (5.4) is solved using the Gaussian LU factorization,

 $J \rightarrow LU$

L (UX)

JX =

E

followed by backward and forward substitutions.

Έ

Thus, substituting forward in

LY = E

yields Y, whence, substituting back in-

UX = ,Y

gives X .

As the system of equations is set in tableau form, the Jacobian \hat{J} may be highly sparse and of large dimensions. Also, it can be shown that most of the Jacobian entries are simply ± 1 . More-

88

(5.4)

(5.5)

(5.6)

(5,7)

over, for each mode of analysis (dc or transient), the same system of equations, with a fixed sparsity of the Jacobian J, is required to be solved repeatedly with different values of J and E. In order to derive the most benefit from these properties, a sparse technique was developed which generates non-looping executable machine instructions (code) through Fortran statements. The code is generated twice, once for the dc analysis and once for the transient. When this code is executed, it performs only necessary operations, avoiding any trivial or unnecessary operations (e.g., multiplication and division by or storage of ± 1).

The linkage and accounting vectors used with the code generation are quite large. However, they are overlaid on other vectors, so that no extra storage space has to be specially reserved for these vectors.

In the following sections, the main techniques used in SCAMPER are reviewed.

5.4 Code Generation

There are two types of machine instructions used in SCAMPER [19], register-to-register (RR) and register-to-core (RX) instructions. The required RR instructions are defined in DATA statements, so that they can be added to the code stream when needed. The

RX instructions consist of an operation code, operand register number, index register number, base register number and displacement field. The operation code, operand register number and index register number occupy the first half word of the instruction, and are defined in DATA statements for every required combination. The index register is used to specify which vector is being operated on, and it contains the address of the start of this vector. The base register number and the displacement field occupy the second half of the instruction word, and define the operand in the specified vector. This is done by setting the second half word to

where I is the subscript value of the operand. The multiplication by 8 is done because vectors are of double precision. In this way vectors of up to 512 entries (=4096'/8) can be assessed. To increase the accessible length of each vector, some registers (1 through 11) are loaded by values, 4096. K, where K is the register number. This allows vectors of entries up to 6144 to be accessed correctly.

8 · (I - 1),

Considering equations (5.4) to (5.7) we can see that the code string has to consist of four basic sections, each with a specific function.

 J-code, to fill in the constant (real) entries of J.
 Also to fill in the integration formula corresponding to the linear reactive elements.

(ii)

E-code, to compute the r.h.s. E of equation (5.3).

(iii) F-code, to perform a Gaussian LU factorization of Jacobian J, with complete pivoting to ensure numerical stability and high accuracy.

(iv)

S-code to perform the solution with the factorized matrix via a forward substitution followed by a back-

In addition to these four main segments of the code, some control statements are inserted between segments (such as reloading index registers and transfer of control), as well as some constants to be initialized at the starting locations of the code.

5.5 Sparse Matrix Technique

The sparse matrix technique is used in conjunction with the code generation. It may seem quite complicated due to the large number of linkage and counting vectors used, but it should be noticed that this technique is executed only twice, once for the dc analysis and once for the transient one. Accordingly, a large portion of the time that would be consumed in generating and updating these vectors is saved due to the execution of the generated code. This is not the place for a detailed analysis of this technique, but the main idea and the description of two basic vectors will be given.

During the setting of the linearized network equations, two vectors, IC and ISGN, are constructed to link and identify the non-zero entries of Jacobian J.

(i) Linkage vector IC (IJ), IJ = 1, 2, ... etc., is a list by rows of the non-zero Jacobian entry column numbers. The entries for each row are separated by terminators, so that entries of each row can be identified. This vector is effectively a two dimensional (I, J) vector, where I represents the row number and J represents the column number.

(ii) Identity vector ISGN (IJ), $IJ = 1, 2, \dots$. This is a marker vector for the non-zero Jacobian entries and corresponds to IC. Entries which are +1 or -1 are identified by 0 and -1, respectively. Non-zero entries other than ± 1 are identified by their location in the vector VAL where they are stored. Note that ± 1 are not stored in VAL.

Using these two basic vectors IC, ISGN, some other linkage and accounting vectors are constructed and updated during the step of Gaussian LU factorization to facilitate the process of generating the machine code.

The process of LU factorization is illustrated in Figure 5.1. Assume that we have a full matrix as in Figure 5.1a. Also assume that the diagonal elements have been chosen as pivots.

92

At

the beginning, the whole matrix is linked as in Figure 5.1b. As the process of factorization proceeds the linkage is split into three different paths: (i) the upper trapezoidal, (ii) the lower triangular, and, (iii) the unreduced matrix (Figure 5.1c). The start and the end of these paths are marked by pointers. At the end of the LU factori-, zation both the upper U and the lower L triangular matrices are linked and identified. In general, the matrix is sparse and the pivot elements may not be the diagonals, so that the linkage paths will be overlapped, but the same principle remains.

5.6 Pivoting Procedure

The pivoting process is almost a complete pivoting, except for some minor constraints which, from experience [21], are found to give better numerical stability. The algorithm used starts by eliminating the rows corresponding to the B.C. equations, then the KVL equations and the KCL (if possible). Also the choice of the pivot column (if possible) is forced not to correspond first to the node voltage variables. During each step of the factorization, the next

Taking into consideration the above remarks the following is a review of the piyoting algorithm.



ſ

(1) Set an upper bound KREF on the number of entries in a row to be accepted as a pivot row. (At the start, this limit is set during the construction of the linkage vectors). This acceptance limit might be modified at any step.

(2) During each subsequent step of factorization, compare the number of entries of each row of the unreduced matrix with KREF. If a row of entries equal to KREF or less is found, terminate the search. If this condition is not satisfied after a complete scan choose the row with the smallest number of entries as the pivot row and reset KREF.

From the selected row, choose the pivot column as follows:-

If the first entry is ± 1 (corresponding to ISGN (•) = 0, -1), pick it.

(ii)

(i)

(3)

If not, pick the first unit magnitude (+1)before the entries corresponding to the node voltage variables.

If none, choose the largest entry which is greater than a certain lower limit (10^{-60}) .

(iv)

(iii)

If none, choose any entry of unit magnitude (+, 1) corresponding to a node voltage

variable.

As soon as a point is chosen, a step of factorization is carried out: (1) processing the entries of the coefficient Jacobian vector VAL, (2) generating the necessary machine instructions for this step of factorization, (3) updating the linkage and accounting vectors. In order to ensure numerical stability and accuracy the representative Jacobian entries are chosen to correspond to the worst case (e.g., the coefficient of an exponential branch is set to zero).

5.7 Numerical Considerations

À variable x = x (t) can in general, be represented by a polynomial having an infinite number of terms. Practical considerations limit the number of terms k to $1 \le k \le 6$ [28]. The differentiated variable x is replaced by the slope of this polynomial at the new time point, which is computed using a "corrector" formula, while the error estimation is done using a "predictor" formula ' [19]. ' The error introduced by this "truncation" of the polynomial is called the truncation error. Moreover, when the nonlinear equations are replaced at a given time point by a linear model, the error introduced by this linearization is called the nonlinear Accordingly, there are two types of error control, the truncaerror () tion error control and the nonlinear error control. The step size is chosen to limit the largest error to be less than or equal to a certain user-specified limit.

5.7.1 Trancation Error

The predictor-corrector algorithm replaces the differentiated variable \dot{x} by the slope of the truncated polynomial of order k at the new time point t_n . The associated local truncation error E_{tr} is defined by [18].

$$(x_{n \text{ true}} - x_{n \text{ approx}}) = E_{tr} + 0 (h^{k+2}) b,$$
 (5.8)

where x_n approx denotes the value yielded by the corrector formula. Neglecting the term 0 (h +2) it can be proved [18] that

$$E_{tr} = \frac{h}{h_t} \left| x_n - x_n^p \right|$$
 (5.9)

where $h_t = t_n - t_{n-k-1}$, $h = t_n - t_{n-1}$, and x_n^p is a predicted value of x_n .

5.7.2 <u>Nonlinear Error</u>

of nonlinearity. If the equation for, say, an exponential capacitor

$$q' = a (e^{bv} - 1)$$

is linearized at a predicted value v^p , it can be proved [19], that the nonlinear error E_n can be approximated by

98 /

$$E_{n} = \frac{b}{2} (v - v^{p})^{2}$$
 (5.10)

5.7.3 Step Size Control

Step size h, is determined according to two criteria: the truncation error criterion (with a lower limit h = H0.), and the nonlinear error criterion (with a lower limit h = H00 where H00 is related to H0 and k).

Let us consider

X, the normalized maximum truncation error,

Y, the normalized maximum nonlinear error,

h , the last (accepted) step size, \ddot{h}_{n+1} , is the new step size.

In SCAMPER, the step size is chosen such that the maximum error is less than or equal to half a user-specified error.

Accordingly, the step size determined by the truncation error criterion is

$$h_{n+1} = h_n \left(\frac{0.5}{X}\right)^{\frac{1}{k+1}}$$
 (5.11)

and, that determined by the nonlinear error criterion is

 $h'_{n+1} = h_{n'} \left(\frac{0.5}{y}\right)^{\frac{1}{k+1}}$

(5,12)

The smaller of h_{n+1} and h_{n+1} is taken as the new step size. There are some other constraints on the step size such as: (1) avoiding wide change of h_{new} / h_{old} (say, < 10, 0.1 >) to achieve better numerical stability ; (2) choosing h_{new} to yield a point just after the cusp change of independent sources ; (3) choosing h_{new} to yield points at specified times to output results at these points.

STEADY-STATE ANALYSIS ROUTINE

5.8 Code Generation for the Adjoint System

One of the main reasons for the superiority of the and transient analysis algorithms used in the program SCAMPER is dc the successful marriage between a highly sophisticated sparse matrix technique and a compiler-like routine for generating machine instruction (code) strings. Once the code is generated it is used repeatedly in constructing and solving the network equations. When executed, this code does only the necessary operations, avoiding any trivial work such as multiplication or division by unity or the storage thereof. Since the steady-state analysis method considered here involves the proposed gradient, algorithm 3.1, the adjoint system has to be integrated as well as the original one. To achieve a performance consistent with that of it was necessary to modify the compiler-like routine to SCAMPER generate another code for integrating the adjoint system in an efficient way, exploiting the common properties of both systems.

Let us first consider the actual representation of the original system. The linearized form of the network relations is

Substituting the definitions given in Part A of this chapter, this becomes

 $X = E_{c}$

100

(5.13)


where D (\cdot) symbolizes the entries of a block of the Jacobian J which are differentiated, while 1 and n represent, respectively, linear and nonlinear terms.

This system of linearized equations is solved by using Gaussian LU factorization followed by forward and backward substitution. The pivoting procedure [19] forces the pivot row to be taken in the followorder: first, rows corresponding to the branch, constitutive relations, second KVL, third the KCL (if possible). Also, for the prot columns, those corresponding to the node voltage variables VN are forced to be the last (if possible).

Consider now the adjoint system,

or

.τ^t λ

(5.14)



This adjoint system may be solved using procedures similar to those for the original one, e.g., by using Gaussian LU factorization, and forward and backward substitution. The additional requirements for generating a solution code are:

(i) to manipulate the transposition of the matrix in such a way as
 to use the same J-code (which fills in the linear entries of the Jacobian)
 for both the original and the adjoint system;

(ii) to make use of the relation between the two systems to simplify the procedure for generating the E-code (which computes the r.h.s. vector E_b);

(iii) to modify the pivot selection algorithm in order to achieve requirements similar to those imposed on the pivoting order of the original system of equations, and also to use the same routine which generates the F- and S- codes for both systems.

• The way in which these requirements were satisfied is outlined in the following sections.

5.8.1 Transposition of the Jacobian

As discussed in Part A, the Jacobian is mainly described by three vectors:

(i) VAL, the coefficient Jacobian vector which holds the real values of J (not including + 1);

(iii) ISGN, the identity vector, which describes the type of entries located by IC (real, ± 1).

A typical sketch of these vectors is shown in Figure 5.2 for a 9×9 matrix.



Figure 5.2. Typical representation of some entries of vectors IC, ISGN and VAL.

Since the execution of the J-code fills in the linear entries of VAL (with the nonlinear entries to be filled via Fortran statements), hence, if IC and ISGN are rearranged in such a way as to keep the order of VAL unchanged, the J-code can serve both the original and the adjoint system. This is achieved by means of the following steps :

(i) Count the number of entries in each row of the adjoint system using IC ;

(ii) use the above counting vector to set the terminators of each row ;

(iii) scan the linkage vector IC to fill the entries of each row of the corresponding linkage vector

ICJ in order of increasing column number and, at the same time, for each entry of IC fill in the corresponding identity vector ISGNJ, i.e., if IC (NL) \rightarrow ICJ (NLJ), hence take ISGN (NL) \rightarrow ISGNJ (NLJ).

In the actual program, not only ICJ and ISGNJ are constructed, but the other linkage and counting vectors required for the LU factorization are generated simultaneously.

5.8.2 E-Code Generation

For the original system, the equation variables are taken as the (negated) differences between the current network variables x_n and the last accepted solution x_{n-1} , i.e., $\Delta x = x_{n-1} - x_n$. Furthermore, the differentiated variable x_n is replaced by an integration formula (polynomial) of order k, $1 \le k \le 6$,

$$\mathbf{x}_{n} = \mathbf{C}_{0} \quad \Delta \mathbf{x} + BDF \quad , \qquad (5.15)$$

where $\Delta x = x_{n-1} - x_n$ and BDF is the explicit part of x_n (backward difference formula) which depends on the last k differences as well as on the step sizes. The term, BDF will be added to the r.h.s of the equation.

106

Now, for the adjoint system a general form of the equa-

tion is 🐳

$$\sum_{j=1}^{n} a_{j} \lambda_{j} + \sum_{j=1}^{n} b_{j} \lambda_{j} = 0$$
(5.16)

If we use the difference $\Delta \lambda$ as variable, $\Delta \lambda = \lambda_{old} - \lambda_{new}$, equation (5.16) takes the form

$$\sum_{i} a_{i} \Delta \lambda_{i} + \sum_{j} b_{j} C_{0} \Delta \lambda_{j} = \sum_{i} a_{i} \lambda_{i} + \sum_{j} b_{j} BDF_{j} \qquad (5.17)$$

Equation (5.17) is used to compute the r.h.s. entries of equation (5.14) corresponding to the (I, Q) and (V, F) rows, while for the equations corresponding to VN the corresponding entries are set to zero. The special functional structure of the Jacobian is exploited to significantly reduce the complexity of the E-code generator.

5.8.3 F- and S-Code Generation and Pivoting

To formulate the Gaussian LU factorization, it was decided to use a pivoting procedure similar to that used in the forward integration. Thus, for the adjoint system, the rows corresponding to node voltages VN and the columns corresponding to the KCL variables are forced to be the last (if possible). This has been achieved by modifying the pivoting algorithm in such a way that it can handle both systems, the original and its adjoint. By doing this, one routine SGEN is used for generating the F- and S- codes of the original system and its adjoint as well.

In summary, the original Jacobian J is transposed by changing the linkage and the identity vectors IC, ISGN in such a way as to keep the coefficient Jacobian vector VAL unchanged. Thus, the J-code generated for the original system serves the adjoint system as well. The functional structure of the Jacobian is utilized to simplify the E-code generating routine. The pivoting procedure has been modified so that it can handle both the original and the adjoint system in a similar manner and so that the same routine for generating the Fand S-codes can be used for both systems.

5.9 Data Management

Dealing with the integration of two different systems requires some kind of data management (storing, retrieving and processing). Some of these data are :

(i) the necessary values of the nonlinear entries of the Jacobian J as well as the corresponding step size (or the time) at different time points over one full period ;

(ii) the two machine instruction strings (codes) generated for solving both the original and the adjoint systems ;

(iii) vector of $(1/_{C}, 1/_{L}, C(0), L(0))_{J}$, in the order of reactive branch number, where C(0) and L(0) are the values of the nonlinear capacitive and inductive elements respectively at the starting time point t_0 . Even though these values can be retrieved from the data in item (i), it is found more efficient to have this vector separately and arranged in a particular way;

(iv) the initial condition vector of the reactive variables q_0 as well as the solution of the network variables at the beginning of the first steady-state iteration.

(ii) to (iv) are each treated as one unit Items (sequential) and their manipulation is a trivial matter. This is not the case for data item (i). The nonlinear entries and the step sizes are evaluated at a series of time points over one full period during. the integration of the original system. For the integration of the adjoint system, this data has to be retrieved in the opposite order and, hence, care must be taken to minimize both storage and the time for retrieving data from disk. We found it most efficient to défine a temporary storage vector 'AREA with suitable pointers and linkage vectors to be used with a particular scheme for handling this problem. The size of AREA was chosen large enough to ensure that in the worst case it can hold data of at least three time points.

* By worst case we mean the largest problem that can be solved by the program (380 branches) with most of the branches nonlinear.

Suppose that the temporary storage vector AREA has NT entries and the number of entries-per time point is NETP (the number of nonlinear entries plus one for step size). Hence, the maximum number of time points that can be held in AREA is , NTPA = NT/NETP (rounded to the least integer number). Since interpolation between time points may be required for computing the Jacobian J, it is necessary to keep the data of the last time point in AREA before reloading it from disk with a new record. Accordingly, the "dynamic" size of AREA, is (NTPA-1) time points. That means the transfer of data from and into disk takes place every (NTPA-1) points. If the data is retrieved from disk more than once, and because we have to keep the last point of the previous record, the starting (and the end) location of the newly retrieved record will move in a loop. sketch of this situation is given in Figure 5.3. This process is controlled by defining pointers that change in a closed loop.

5.10 Some Aspects of the Backward Integration of the Adjoint System

Figure 5.4 is a sketch of the adjoint system integration loop. It represents the integration of equation (3.19) over a period in the opposite order to that of integrating the original system given by equation (3.5). Requirements for the adjoint system integration are quite different from those for the original system.



(i) For the original system the required data are computed at each step using the currently available values of the variables, while for the adjoint system these data already exist but may need some manipulation (retrieval, interpolation).

(ii) The number of differentiated variables of the adjoint system λ_2 may be different from that of the original system. Accordingly, the truncation error computation as well as setting the initial conditions for both systems will be different.

(iii) In the original system, there are two error criteria to be applied : the truncation and nonlinear error criteria. For the adjoint system the only criterion is the truncation error control since the system is linear (time varying if the original system is nonlinear).

(iv) For the original system the step size is determined by the error control as well as the "cusp" changes (if any) of the independent sources, while for the adjoint the step size is determined by the error control as well as by the interpolation process.

5.10.1 Setting the Initial Conditions of the Adjoint System

As the initial conditions of the algebraic variables λ_1 can be chosen arbitrarily we simply put λ_1 (T) = 0. The initial conditions of the differentiated variables λ_2 (T), are defined by equation (3.17), with the detailed equations for the different types of



reactive elements given by equations (3.22a, b, c, d). Corresponding to each inductive branch there is only one differentiated variable. In the case of capacitive branches we are faced with two problems: (1) corresponding to each branch there may be two differentiated variables (one for each node); (2) there may be loops of only capacitors and/or independent voltage sources and/or short circuit branches. In order to overcome the second problem an algorithm was formed to identify such loops and to ignore one of the capacitors in such a loop in computing the initial conditions. To get around the first problem an algorithm was designed to construct linkage vectors which are used to set the initial conditions in the correct and proper way. The outline of this algorithm is as follows.

Algorithm 5.1

- 1. Count the number of trees that consist of only capacitors and/or independent voltage sources and/or short circuits.
- 2. Count the number of capacitors connected to each node in these trees.
- 3. If there is a tree which contains the ground node (if the ground node is specified) consider it to be the first tree and consider the ground node the first node in this tree. For any other tree consider the first node to be the one which connects the largest number of capacitors.

Put the variable corresponding to this node equal to zero and proceed.

Following from the previous node NI choose the next node NJ separated from NI by one branch, C, and j set the initial condition of this node, λ_{NI} according to the relation

$$\lambda_{NJ} = \pm \left(\frac{\omega_j}{C_j} \left(v_j \left(0 \right) - v_j \left(T \right) \pm \lambda_{NI} \right) \right)$$
 (5.18)

5.10.2 • Truncation Error Computation

 $z = \lambda_{N1} - \lambda_{N2}$

The computation of the truncation error of variables related to inductive branches is done in a similar manner to that for the original system. Corresponding to each capacitive branch, the <u>difference</u> between the two node variables is defined and the initial conditions are set such that this difference is equal to $(\Delta v_0 / C)$. Since the capacitors in the circuit may differ by orders of magnitude, the control of the absolute values of the variables is meaningless. Hence, the control is carried out on the differences corresponding to each capacitive branch. Consider a capacitor, connected between two nodes NI and N2. Define a variable

(5.19)

Hence

$$= \lambda_{N1} - \lambda_{N2}$$
 (5.20)

If the differentiated variable z_{k} is replaced by an integration formula (polynomial) of order k, the predicted value will be

$$z^{P} = \lambda_{N1}^{P} - \lambda_{N2}^{P}$$
 (5.21)

Accordingly, the truncation error control of the differences corresponding to each capacitive element is done according to equations (5.19) and (5.21) i.e., on the difference $|z - z^p|$, and not on $|\lambda_{N1} - \lambda_{N1}^p|$ or $|\lambda_{N2} - \lambda_{N2}^p|$.

5.10.3 Interpolation and Step Size Control

Suppose the largest normalized truncation error computed at a time point is X. (If the last step size is h_n , the new step size h_{n+1} is chosen to produce half the normalized truncation error (the, same as for the original system integration) i.e.,

$$h_{n+1} = h_{pn} \left(\frac{0.5}{X}\right)^{\frac{1}{k+1}}$$

The step size is subject to other restrictions such as the ratio h_{n+1} / h_n having to be within certain limits $\langle \epsilon_1, \epsilon_2 \rangle$ (say °<10, 0.1>). Also, if the last point was rejected the step size is not allowed to increase.

Another restriction is set by the interpolation between points and by the end-of-record condition as will be shown next.

For each time point the step size h_n is stored along with the nonlinear entries. Considering the backward integration, suppose the precomputed step size is H, and IP defines the location of the last accepted point, where the time difference between IP and the preceding point is DHl as shown below. A pointer J



is moved starting from value I and scanning step sizes h_{I+1} , h_{I+2} , \cdots etc., accumulating a step h

 $h = (h_{I} + h_{I+1} + \dots + h_{J-1}) - DHI$

The accumulation of h- terminates when

 $h + h_{J} > H$.

A new point JP will be taken between (or coincident with) points J and J+1. If JP is within a certain limit (say | H - h | / H < 0.01) of points J or J+1, the nearest point will be taken as the new point. If the new point is accepted, pointers will be reset to I = J, IP = JP and DH1 = DH2. For the special case when the last accepted point is `equal to or just before the last point I of the temporary storage vector AREA, and

 $h = h_{I} - DHI < H$,

there are two possibilities:

(1) if h > HO, where HO is a lowest limit on the step size,
H will be set to H = h and IP to I;

(2) if h < H0, this means that the last point I in AREA precedes any new acceptable point. In this case new data will be retrieved from disk in place of the previous data, except for the last point I. Pointer IP will be set to I and DH1 = DH1 - h_{I-1} . The process of interpolation will then continue normally.

It should be mentioned that we use a linear interpolation process. For example, if a point JP is chosen between two points J and J+1 with corresponding values X_J and X_{J+1} respectively, the value corresponding to JP, X_{JP} is given by

 $x_{JP} = x_{J} + \frac{DH2}{h_{J}} (x_{J+1} - x_{J})$

5.11 Resetting the Initial Conditions of the Original System

The values of the differentiated variables of the original system corresponding to the reactive elements, before and after the call of the optimization routine, may differ considerably. Due to the nonlinearity of the network, the integration procedure may fail completely to converge at the first step of integration if there are such big changes. To avoid this problem, the solution may proceed as follows.

Corresponding to each reactive element an auxiliary source is added: an independent voltage source E_c in series with a capacitive branch and an independent current source J_l in parallel with an inductive branch. Suppose the new value of an initial condition is x_{new} and an old value is x_{old} (in our case the old value is that corresponding to the starting point of the first steady-state iteration). Assume that this previous value satisfies the system equations, and that the solution vector V_{old} at that point is available (or at least the variables corresponding to the nonlinear elements). We set the value of the corresponding source x_s at the starting point t_0 to be

 $x_s (t_0) = x_{old} - x_{new}$

and at $t_0 + \delta t$ to be

118 .



is chosen small, relative to the period of integration. However, the initial conditions of the corresponding reactive element is set equal "to x_{new} . As a result, the terminal variable, x of each reactive element (including the auxiliary source) start with $x = x_{old}$ at $t = t_0$ and approaches x_{new} near $t = t_0 + \delta t$. In this way, the integration of the network equations starts with initial conditions which satisfy the nonlinear error criteria and which can then change to the new initial conditions x_{new} in several steps within the interval from t_0 to $t_0 + \delta t$.

CHAPTER VI

CONCLUSIONS AND RECOMMENDATIONS '

The gradient method for computing the periodic steadystate response of nonlinear networks has been successfully implemented with the dc and transient analysis program SCAMPER which can be used to analyze large networks by virtue of its being based on the sparse tableau formulation of network equations. This periodic response solver has been tested with several circuits including those which were used by others working on this problem, in particular, Nakhla and Branin [10], [11] who first proposed the gradient method but used the state variable approach to solve the network equations, which is only practical for small systems.

'A by product of the derivation of the equations for comcomputing the gradient vector was the proof that they are valid for circuits having nonlinear elements involving not only single [6] but also double dependencies. This has been verified for several examples [Appendix II].

As predicted, the time per iteration never exceeded twice the time for one forward integration over one period. The rate of convergence depends very much on the characteristics of the optimization routine, the scaling of the variables and the associated gradients, and on the choice of the starting time t_0 .

The high standards set by the authors of SCAMPER were adhered to in coding the gradient algorithm. Thus, the setting up and solution of the adjoint system of equations is also achieved by a routine which generates executable machine instructions, and which was designed to exploit features common to both the original and the adjoint systems. Furthermore a large number of the subroutines used in the integration are common to both systems.

Several important and potentially fruitful areas of investigation present themselves as natural extensions of the present work. Various improvements that can be made in the existing program are also apparent. These are discussed below in an arbitrary order.

(1) To reduce the computation time, the error control should be modified. Assume a user supplies an error limit E_{max} which, correspondingly, causes the results to be accurate to, say, n digits. During the first few iterations lower accuracy may be satisfactory, hence the error limit can be larger which, in turn, causes a reduction of the integration time. Then, as the solution converges to the steady-state condition, the error limit can be reduced in steps to the user-specified limit. If this would be done, there would be need for the extra control imposed on the adjoint system initial conditions for changing the accuracy of computation.

(2) The initial conditions from which the steady-state analysis starts are determined by integrating the system equations up to a

specified time point t_0 . In some cases it may be more advantageous to start the analysis with a set of user-defined initial conditions. Work on this feature is now in progress.

(3) A linear interpolation process is used for computing the Jacobian of the adjoint system. According to the tested problems, this process works quite well, but cases may be encountered where a higher order of interpolation (2nd, 3rd, ... etc.) may be necessary.

(4) In Chapter IV some comments are made concerning the scaling problem in the light of our admittedly limited experience with the program. Further work is needed to achieve a more systematic scheme for using the scaling facilities that are made available to the user.

(5) Considerable work remains to be done on the problem of the design of an efficient function minimization routine. Although the variable metric routine was more successful in solving the test problems than the conjugate gradients one, its initial convergence was, nevertheless, quite slow in some cases. Moreover, the variable metric algorithm will present storage problems when large networks are tackled.

(6) The present investigation was concerned with the periodic steadystate analysis with periodic input (nonautonomous). The starting time t_0 of the analysis was regarded as a constant, in fact, the gradient depends on the choice of the starting point and, accordingly, the convergence criteria of the gradient method depends on that choice. In

other words, by letting the starting point t_0 change, a point can be found at which the computed gradient yield better convergence [11]. To do so, we redefine the objective function and the variable space as well as the gradient to include t_0 as a variable, i.e.,

 $\Phi = \Phi (X_0, t_0)$

Hence the gradient with respect to (x_0, t_0) , will be

$$G = \left(\frac{\partial \Phi}{\partial X_0} , \frac{\partial \Phi}{\partial t_0} \right) '.$$

By passing the objective function Φ , gradient G and the variables (X_0, t_0) to the optimization routine not only the variables X_0 are updated but also t_0 is updated towards an "optimal" point. With this modification the gradient method is expected to give better results. Also, an extension of the method to the oscillatory (autonomous) case can be achieved with little effort, especially after including optimization of the starting point.

APPENDIX I

ACCEPTABILITY OF DOUBLE DEPENDENCY REACTIVE ELEMENT

To show that the term in equation (3.14)

 $\frac{\partial \mathbf{u}}{\partial \mathbf{t}}$, $\frac{\partial \mathbf{x}}{\partial \mathbf{p}}$, $\frac{\partial \mathbf{u}}{\partial \mathbf{p}}$, $\frac{\partial \mathbf{x}}{\partial \mathbf{t}}$

is identically zero for a class of systems that have a unique solution over (t) for a given range of parameter {p}. Let the system variables (x, u) can be represented as functions of the parameters p, t, i.e.,

The parametric relation between x and u may take any of two implicit forms

$$\psi(\mathbf{x}, \mathbf{u}, \mathbf{t}) = 0$$

 $\mathbf{x} = \mathbf{x} (\mathbf{p}_i, \mathbf{t})$

u = u(p, t).

where t is taken as the running parameter, or

 ε (x, u, p) = 0,

where p is the running parameter.

We are interested in practical problems where both the

variables x and u, and the p are within a finite range. Let us

124

(I.1)

(I.2)

(I.7)

consider a point in the (x, u) plane where $p = \hat{p}$ and $t = \hat{t}$. The two functions

$$\psi(x, u, t) = 0.$$
 (I.4)

$$\epsilon(x, u, p) = 0$$
 (1.5)

will, in general, intersect at one or more points a_1 , a_2 , ..., as shown in Figure I.1. The interpretation of this figure is that, if the two curves intersect at more than one point the system has more than one solution defined by points a_1 , a_2 , ... Hence, for our proof, we consider that class of systems which can be represented by Figure I.2 where both curves are tangent at a single point. In other words, for a given value of p, $p = \hat{p}$ a group of curves (I), corresponding to equation (I.2) can be found such that curve (II), representing equation (I.3), for a given value of $p = \hat{p}$, is the envelope of this group.

If the group of curves $\psi = 0$ is expressed parametrically in the form of equation (I.1), then equation (I.2) can be replaced by a function ϕ (p, t) where

 ϕ (p, t) = ψ (x (p, t), u (p, t), t) = 0 (I.6)

Hence, at any point on any curve of the group $\psi = 0$

 $\sigma = \frac{\sigma}{q6} + \frac{\psi}{r6} + \frac{x}{q6} + \frac{\psi}{x6} = \frac{\phi}{q6} = \frac{\phi}{q6}$

and



127 $\frac{\partial \phi}{\partial t} = \frac{\partial \psi}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial \psi}{\partial u} \frac{\partial u}{\partial t} + \frac{\partial \psi}{\partial t} = 0$ (I.8) For a given value $p = \hat{p}$, the solution of the system (x, u) for t as a parameter will be given by curve (II), i.e., by the envelope of group At any point on the envelope [24]^h (I). $\frac{\partial \Psi}{\partial t} = 0$ (1.9) hence $\frac{\partial \phi}{\partial t} = \frac{\partial \psi}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial \psi}{\partial u} \frac{\partial u}{\partial p}$ (I.10) If we multiply equation (1.7) by $\frac{\partial x}{\partial t}$ and equation (1.10) by $\frac{\partial x}{\partial p}$ and subtract, we obtain $\frac{\partial \psi}{\partial u} \left[\begin{array}{c} \frac{\partial u}{\partial t} & \frac{\partial x}{\partial p} - \begin{array}{c} \frac{\partial u}{\partial p} & \frac{\partial x}{\partial t} \end{array} \right] = 0$ may not equal zero But Hence $\frac{\partial \mathbf{u}}{\partial \mathbf{t}} \frac{\partial \mathbf{x}}{\partial \mathbf{p}} - \frac{\partial \mathbf{u}}{\partial \mathbf{p}} \frac{\partial \mathbf{x}}{\partial \mathbf{t}}$

APPENDIX II

VERIFICATION OF ACCEPTABILITY OF DOUBLE DEPENDENCY

Several nonlinear and time varying differential equations which require the use of nonlinear circuit elements with single or double dependencies when simulated on SCAMPER, and for which solutions can be obtained in closed form, were used to test the validity of the gradient derivation in Chapter III. Results of direct hand calculations and of the use of the program described in this thesis are compared below. It should be noted that, due to the inherent nature of the SCAMPER simulation, such as the restrictions imposed by the set up routine of this program, exact agreement can not be expected. Moreover, lower accuracy should be expected in computing the gradient G because it is obtained after two steps of integration, the forward integration and the backward one.

The results will be compared on the basis of the following definitions.

Let the differential equation be

the starting time point t and the period T . 0

 $x_0 = x(t_0)$,

 $= x (t_0 + T)$

= f(x, t),

We define

(II.2)

(II.1)

- |

From the definition of the objective function ϕ

$$= \frac{1}{2} (x_0 - x_T)^2$$
 (11.3)

the gradient is

$$G = \frac{\partial \phi}{\partial x_0} = (x_0 - x_T) \left(1 - \frac{\partial x_T}{\partial x_0}\right)$$
(II.4)

Example 1

Consider the differential equation

$$x = -t (x - 1)$$
 (II.5)

solving we get

$$-\frac{1}{2}(t^{2}-t_{0}^{2})$$

$$x = 1 + (x_{0} - 1) e$$
(II.6)

The problem was simulated by simulating

 $\mathbf{x} + \mathbf{t}\mathbf{x} - \mathbf{t} = 0$

twice, once as a node equation (Figure II.1a) and once by using its dual loop equation. In the simulation, t_0 and T were taken to be $t_0 = 1$, T = 2. Step size control, in fact, forced the starting point to be $t_0 = 1.0062$. Substituting these values and choosing x_0 such that x (0) = 0 we obtain

0.3974 , x х_т 0.9891 , 0.5805 . G

= 0.3972 ,

The results from the simulation were

x₀

 $x_{rr} = 0.9891$, t^{i} = 0.5811 G (II.8) Thus, the gradients are in agreement to three digits (< 0.1 %). In Figure II.1 the representation of the nonlinear resistive element

Gl was

$$I_{G1} = G1 (I_{J1}) V_{G1}$$

 $F_{L1} = L1 (I_{R2}) \cdot I_{L1}$

In another simulation the equation was rewritten as

$$(\frac{1}{t})$$
 x' + x - 1 = 0 , (II.9)

where the derivative term was represented as in Figure II.2 by a nonlinear inductor L (or C in the dual network) of value equal to 1/t. More specifically, the reactive element is modelled through a double dependency

(II.7)

In this case agreement is only to two significant figures, primarily due to the modelling compromises that must be made with 1 / t at

Solution of the differential equation

$$x = -x^2 t$$
 (II.10)

is

t = 0 .

Example 2

$$x = \left[\frac{1}{2}t^{2} + \left(\frac{1}{x_{0}} - \frac{1}{2}t_{0}^{2}\right)\right]^{-1}$$
 (II.11)

The simulation of this problem was done on the basis of the following form of (II.10)

 $x + (x \cdot t) x = 0$

where the nonlinear coefficient was represented by a nonlinear conductance Gl (or its corresponding element in the dual network) where, $Gl = x \cdot t$. Taking $t_0 \simeq 0$ ($\approx 0.2 \times 10^{-7}$), $T = \sqrt{2}$ and $x_0 = 1$, the hand computations yielded

$$x_0 = 1.$$
 ,
 $x_T = 0.5$,
 $G = 0.375$

(II.12)

The results of simulation as in Figure II.3 gave

$$x_{T} = 1.$$
,
 $x_{T} = 0.5$,
 $G_{*} = 0.3734$.

Both results of the gwadient are in agreement to within < 0.4 %. The representation of the nonlinear element Gl and R2 is done according to the relations

$$I_{G1} = G1 (V_{R2}) V_{G1}$$

 $V_{R2} = R_2 (V_{01}) I_{R2}$

where Ol is an open circuit branch across the capacitor C and the independent source (initial condition) El.

Example 3

Consider the differential equation

x

3

Solving we get

$$x = [t^{2} + (x_{0}^{2} - t_{0}^{2})]$$

The simulation was done on the basis of

= t/x

132

(II.13)

(II.14)

(II.15)

(II.16)

(II.17)

(x) $\cdot x - t = 0$

1.

2.

-0.5 .

1.002

2.003

-0.5009

where the derivative term was simulated by a nonlinear capacitor C (or inductor L in the dual network) where c = x.

Taking $t_0 = 1$, T = 1, x(0) = 0.

=

The hand computation gives 🔤

×0

x_T

G

×∩

G

while the results of simulation as in Figure II.4 yielded

The results for the gradient show an agreement between the simulation and the hand computations within three digits.



(a)

Jl	GND	Nl	TABLE	, 0. 0.	10.	10	, - ,	*
Gl _	Nl	GND	TABLE	V (G1)	I (J1)	0.	I.E - 8 10.	- 10
cl	Nl	GND	1		-		7	

134

Figure II.1.

Simulation of equation x = -t(x-1)using nonlinear resistance.

(a) The simulating circuit. , \sim

(b)

(b) The input circuit description.



(a)

ECI	N4	GND	TABLE	0.	0.	10.	10			
E2	N3	GND	l°.	-			·			,
R2.	N3	GND	TABLE	I (R2)	v	(EÇ1)	0.]	E-8,	. 10.	`10 .
El	Nl	GND ·	TABLE	0. 0.		1.E-7	1.			
Rl	Nl	N2 _.	1.		rana n m fr			-	_	
ГІ,	N2	GND	TABLE	I(L1)	I	(R2)	0.	ο.	1.E8.	1.E8.

(b)

Figure II.2. Simulation of equation x = -t (x-1)using nonlinear reactive (L) circuit.

(a) Simulation circuit.

(b) Input description.



Figure II.3. Simulation of equation $\dot{x} = -x^2 t$. (a) Simulation circuit.

(b) Circuit description.


GND Nl TABLE 0. 10. **J**1 0. 10. 1 Rl TABLE 0. 1.E5 1.E-6. 1.E10. Nl GND сl Nl GND TABLE V (C1) 1.E-5. 10. Ò. 10.

(b)

Figure II.4.

(a) Simulation circuit.

Simulation of equation $\dot{x} = t / x$.

(b) Circuit description.

137

REFERENCES

- [1] T.J. Aprille, Jr., and T.N. Trick, "Steady-State Analysis
 of Nonlinear Circuits with Periodic Inputs," Proceedings of the IEEE, vol. 60, pp. 108-114, January 1972.
- [2] , "A Computer Algorithm to Determine the Steady-State Response of Nonlinear Oscillators," IEEE Trans., Circuit Theory, vol. CT-19, pp. 354-360, July 1972.

[3]

[4]

[5]

-[6]

- F.R. Colon, and T.N. Trick, "Fast Periodic Steady-State Analysis for Earge-Signal Electronic Circuits," IEEE J. Solid-State Circuits, vol. SC-8, pp. 260-69, August 1973.
 - T,N. Trick, F.R. Colon, and S.P. Fan, "Computation of Capacitor Voltage and Inductor Current Sensitivities with Respect to Initial Conditions for the Steady-State Analysis of Nonlinear Periodic Circuits," IEEE Trans. Circuits and Systems, vol. CAS-22, pp. 391 -396, May 1975.

S.W. Director, A.J. Brodersen, and b.A. Wayne, "A Method for Quick Determination of the Periodic Steady-State in Nonlinear Networks," Proceedings of the Ninth Allerton Conference on Circuit and System Theory, pp. 131-139, October 1971.

S.W. Director, and R.A. Rohrer, "The Generalized Adjoint Network and Network Sensitivities," IEEE Trans. Circuit . Theory, vol. CT-16, pp. 318-323, August 1969.

138

[7]

R.K. Brayton, and S.W. Director, "Computation of Delay Time Sensitivities for Use in Time Domain Optimization," IEEE Trans. Circuits and Systems, vol. CAS-22, pp. 910-920, December 1975.

[8] S.W. Director, and K.W. Current, "Optimization of Forced Nonlinear Periodic Circuits," IEEE Trans. Circuits and Systems, vol. CAS-23, pp. 329-334, June 1976.

[9] K.W. Current, "Optimization of Nonlinear Periodic Circuits," Ph.D. Dissertation, University of Florida, December 1974.
[10] M.S. Nakhla and F.H. Branin, Jr., "Determining the Periodic Response of Nonlinear Systems by a Gradient Method," Int. J. Cir. Theor. Appl., vol. 5, no. 3, pp. 255-273, July 1977.
[11] M.S. Nakhla, "Steady-State Analysis of Nonlinear Periodic

Systems," Ph.D. Dissertation, University of Waterloo,

[12] C. Brezinski, and A.C. Rieu, "The Solution of Systems of Equations Using the ε-Algorithm, and an Application to Boundary-Value Problems," Math. Comp., vol. 28, pp. 731-741, 1974.

[13]

E. Gekeler, "On the Solution of Systems of Equations by the Epsilon Algorithm of Wynn," Math. Comp., vol. 26, pp. 427-436, 1972.

[14] J.B. McLeod, "A Note on the e-Algorithm," Computing, vol. 7, pp. 17-24, 1971. [15]

P. Wynn, "Acceleration Techniques for Iterated Vector and Matrix Problems," Math. Comp., vol. 16, pp. 301-322, 1962.

[16] S. Skelboe, "Extrapolation Methods for Computation of the Periodic Steady-State Response of Nonlinear Circuits," 1977 IEEE Int. Symp. on Circuits and Systems, Pheonix, Arizona, pp. 64-67.

[17] G.D. Hachtel, R.K. Brayton, and F.G. Gustavson, "The Sparse Tableau Approach to Network Analysis and Design," IEEE Trans. Circuit Theory, vol. CT-18, pp. 101-113, January 1971.

[18] R.K. Brayton, F.G. Gustavson, and G.D. Hachtel, "A New Efficient Algorithm for Solving Differential-Algebraic Systems Using Implicit Backward Differentiation Formulas," Proc. IEEE, vol. 60, no.1, pp. 98, 108, January 1972,
[19] M.L. Blostein, and D.J. Millar, "SCAMPER: Information Manual," Electrical Engineering Report, McGill University, Montreal, Canada, 1976.

[20] " ———, "SCAMPER: User Manual," Electrical Engineering Report, McGill University, Montreal, Canada, 1976.

[21] M.L. Blostein, Private Communication.

[22] R. Fletcher, "New Approach to Variable Metric Algorithms," The Computer Journal, vol. 13, pp. 317-322, August 1970.

140

[23]	M.J.D. Powell. "Restart Procedures for the Conjugate Gradient
4	
	Method, " Report C.S.S. 24 (A.E.R.A., Hartwell,
,	November 1975).
	· · · · · · · · · · · · · · · · · · ·
[24]	Carl E. Pearson, "Handbook of Applied Mathematics," Van Nos-
ττ3 •	trand Reinhold, pp. 387-391, 1974.
ŧ	
[25]	M.S. Nakhla, and J. Vlach, "A Piecewise Harmonic Balance
-	Technique for Determination of Periodic Response of
,	Nonlinear Systems," IEEE Trans. Circuits and Systems,
<i>,</i>	vol. CAS-23, no. 2, pp. 85-91.
I.	
[26]	R.S. Norin, and C. Pattle, "Effective Ordering of Sparse Matrices
	Arising from Nonlinear Networks," IEEE Trans. Circuit
	Theory, vol. CT-18, pp. 139-145, January 1971.
. 4	
[27]	C.A. Desoer and E.S. Kuh, "Basic Circuit Theory," vol. 2,
	New York : McGraw-Hill, pp. 292-300, 1967.
; [28]	C.W. Gear, "The Control of Parameters in the Automatic, Integra-
Q '	tion of Ordinary Differential Equations," Department of
	Computer Science, University of Illinois, Int. Rep.,

May 1968.

Ľ

· ~ · +

「前頭」

Read the state of the state of

1.1.1. B

Ē

.141