

# Piecewise-linear modelling and analysis

## Citation for published version (APA):

Bokhoven, van, W. M. G. (1981). Piecewise-linear modelling and analysis. [Phd Thesis 1 (Research TU/e / Graduation TU/e), Electrical Engineering]. Technische Hogeschool Eindhoven. https://doi.org/10.6100/IR118197

DOI: 10.6100/IR118197

## Document status and date:

Published: 01/01/1981

## Document Version:

Publisher's PDF, also known as Version of Record (includes final page, issue and volume numbers)

## Please check the document version of this publication:

• A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.

• The final author version and the galley proof are versions of the publication after peer review.

• The final published version features the final layout of the paper including the volume, issue and page numbers.

Link to publication

### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- · Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
  You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

www.tue.nl/taverne

### Take down policy

If you believe that this document breaches copyright please contact us at:

openaccess@tue.nl

providing details and we will investigate your claim.

# W.M.G. van Bokhoven **Piecewise–Linear Modelling and Analysis**



# PIECEWISE-LINEAR MODELLING AND ANALYSIS

# PIECEWISE-LINEAR MODELLING AND ANALYSIS

PROEFSCHRIFT

ter verkrijging van de graad van doctor in de technische wetenschappen aan de Technische Hogeschool Eindhoven, op gezag van de rector magnificus, prof. ir. J. Erkelens, voor een commissie aangewezen door het college van dekanen in het openbaar te verdedigen op dinsdag 26 mei 1981 te 16.00 uur

door

Wilhelmus Maria Gezinus van Bokhoven geboren te Drunen

## DIT PROEFSCHRIFT IS GOEDGEKEURD DOOR DE PROMOTOREN

Prof. Dr.-Ing. J. A. G. Jess en Prof. Dr. O. Wing

**Aan Tonny** 

### CONTENTS

1.	Introduction and summary	3
2.	Piecewise-linear electrical networks	12
	<ul><li>2.1 The multiport network model</li><li>2.2 Some conditions about the network response</li><li>2.3 Piecewise-linear two-pole elements</li></ul>	12 15 18
3.	Equivalent network properties of P and P $_0$ matrices	21
	3.1 Relations between hybrid representations 3.2 The hybrid structure of M 3.3 The matrix classes P and $P_0$	24 27 31
4.	The state-model of a piecewise-linear mapping	38
	<ul><li>4.1 The structure of the state-model</li><li>4.2 Adjacent regions in a minimal state-model</li></ul>	44 48
5.	Complementary pivoting methods	53
	<ul><li>5.1 The Katzenelson algorithm</li><li>5.2 The algorithm of Cottle</li><li>5.3 The Lemke algorithm</li></ul>	54 57 63
6.	The modulus algorithm	67
	<ul> <li>6.1 The modulus transformation</li> <li>6.2 The modulus algorithm as a contraction mapping</li> <li>6.3 A polynomial algorithm</li> <li>6.4 Implementation aspects and variants of the modulus</li> </ul>	67 71 73 77
	6.5 Relations with the global model of Chua	79
7.	The inverse mapping and some homeomorphism conditions	84
	<ul> <li>7.1 The state-model of the inverse mapping</li> <li>7.2 Homeomorphic mappings</li> <li>7.3 The construction of the state-model for some simple mappings</li> </ul>	84 89 101
8.	Interconnection of PL systems and dynamic PL systems	108
	<ul><li>8.1 The state-model construction</li><li>8.2 Dynamic PL systems</li></ul>	108 114
9.	Piecewise-linear models for some electronic devices	121
Appendix I		133
Appendix II		135
References		138
Samenvatting		

Curriculum vitae

#### 1. Introduction and summary

In electrical engineering, the simulation of complex electronic networks has become one of the important and vital tasks in the design of large scale integrated circuits. Due to inevitable disturbing effects like for example stray capacitances, finite isolation or parasitic transistors it is no longer adequate or even possible to construct a physical equivalent circuit from discrete components in order to test the circuit-design for its correct operation by a sequence of measurements.

In a simulation of the circuit, these parasitic effects can in theory be incorporated easily and thus an early indication of the correctness of the design can be obtained before the circuit is actually integrated in silicon.

Then, although simulation seems to be the answer to a lot of problems, the difficulties are merely shifted from one place to another, i.e. the simulation task will become enormously complex. It has to cover logic simulation of digital circuits as well as timing simulation of these devices and in addition it must be able to predict the response of analog circuits for DC and transient excitations. Moreover, these traditionally different simulation tasks now have to be applied to the same device on account of the combination of both analog and digital circuits on a single chip. Such a complex task requires at least a hierarchical organization and subdivision in order to be completed successfully within a reasonable amount of time. Besides that, severe restrictions have to be imposed on the structure of the used database for the different levels of simulation, to enable an easy conversion of the calculated data between those levels. For example, the result of a simulation at register-transfer level must be translated to data required at the gate level for application in a timing simulation or the other way round.

Assuming all these problems to be solved, a simulation can still be unsuccessful because of instability arising by the interaction of the various simulation tasks at the different levels or types of simulation.

All these requirements indicate that an efficient simulation

can only be performed when organized in such a way that a common approach for all levels is obtained using a uniform database which avoids the conversion problems mentioned.

The main task at each simulation level is the solution of a set of nonlinear algebraic equations describing the network response after an appropriate modelling has been applied to the network components. Hence a common modelling approach for each simulation level is the key to a uniform simulation process and the associated database. Due to the diversity of the network components at the various levels with respect to the complexity of their internal structure or functional behaviour, a hierarchical and transparent organization of the simulation is necessary.

At the lowest level, a model must describe the electrical behaviour of bipolar or MOS transistors represented by a relatively simple set of nonlinear equations. At the gate level, dealing with logic signals represented by voltage or current levels, the operations on these signals are given by input-output relations in terms of tables. At a still higher level, macromodels are needed for the description of for example shift registers or full-adders in digital circuits as well as for operational amplifiers or oscillators in analog circuits. Thus, depending on the particular level or circuit type, the classical simulation of a circuit may require the solution of a set of nonlinear equations, a set of piecewise-linear equations, boolean function evaluation or even a combination of them.

Each of these topics has its own class of solution methods and related problems. The nonlinear equation set is generally solved by applying some version of the standard Newton-Raphson process [1] which is based on local linearizations valid in an infinitessimal small area of the solution space. The nonlinear equation set is then in fact transformed in a sequence of Linear equations in such a way that the sequence of solutions of these linear equations converges to a solution of the original problem. One of the main difficulties in this approach is that a sufficiently close estimate of the solution has to be supplied in advance in order to guarantee the convergence of these iterative solution algorithms.

This problem is almost as difficult as the determination of the solution itself. In particular for digital circuits where the DC state is generally non-unique, the resulting set of nonlinear equations will have multiple solutions and the standard Newton-Raphson method cannot easily cope with such a situation. In [2], Sandberg and Willson supplied a theory to test for the existence of a unique DC operating point for a certain class of networks with bipolar transistors. This test may reduce the simulation problems but is in fact of more theoretical importance as it shows whether a circuit has intrinsic possibilities to be used as a multi-stable element or can have a unique DC operating point exclusively, which is e.g. required for amplifiers and related analog circuit-blocks.

The evaluation of boolean functions in logic simulation is in principle not a difficult task. However the translation of these data in order to link up with transient simulation in some analog part of the circuit is not straightforward. An easy way out is provided by using threshold-gate equivalents which can be modelled by piecewise-linear functions.

From the foregoing observations it will be clear that a simulation program which includes different types of element-models requires a complex internal organization for dealing with the different solution methods and databases. Hence for an efficient simulation a choice has to be made for a unique type of representation for all element-models. As with increasing complexity, the DC response of electronic components shows a tendency to approach a piecewise-linear relationship, it seems most natural to use the piecewise-linear description as a basis for the modelling. Besides that it is far more easy to approximate a nonlinear function by a piecewise-linear function than the other way round.

This thesis is supposed to demonstrate that piecewise-linear modelling can be applied to solve many of the problems outlined above. Traditionally piecewise-linear modelling implies a number of problems.

As to the piecewise-linear equations, these were defined by piecewise-linear functions specified by a list of mappings and

polyhedral regions in some Euclidean space. In contrast with the Newton-Raphson method this yields a linearization valid in a collection of finite regions in the solution space. Within each region, the piecewise-linear equation set is then equivalent to a linear equation set which can easily be solved by standard methods from linear algebra. In this case the main problem is twofold. First of all the required list of polyhedral regions suffers from the disadvantage that a lot of data is required to specify the boundaries of each region and the mapping of those regions. For an approximation of sufficient accuracy the number of required regions may become quite large. In fact the description is still a collection of local linear descriptions without any strong interrelationship except for some continuity demands at the region boundaries. It can be compared to the sequence of locally linearized descriptions which arise during the run of the Newton-Raphson algorithm, although the latter linearizations are only generated along a particular path in the solution space.

The second problem is situated in the method of localizing the particular region in which the solution is to be found. This process requires a search through the list of regions, guided by some algorithm and is time-consuming as well as inefficient. One cause of inefficiency is that the boundaries of a region are at least stored twice for each adjacent region and general information about adjacency of regions is not available in a direct way. Moreover, for any linear transformation applied on the set of piecewise-linear functions, the boundaries of the regions in the transformed space have to be determined one by one in order to be able to construct the list of the mapping for the transformed equation set. In particular, an electrical interconnection of two components would require the construction of a composite list for the combined piecewise-linear equation set from the separate lists of the equation sets for both individual components. This complex operation results in an exponential growth of the storage complexity and computation.

In order for a piecewise-linear approach to be successful, the forementioned disadvantages of the piecewise-linear equations have

to be overcome. As with nonlinear equations or mappings one should look for a comparable implicit global description of PL mappings, valid in all polyhedral regions of the definition space.

One of the main achievements of the research reported in this thesis is the presentation of such a new global implicit description of piecewise-linear mappings. This description is called the statemodel of the mapping. Basically the model consists of two functionally different parts. It combines linear algebraic methods with combinatorial methods in a single model.

The first part of the state-model actually describes the piecewise linear mapping in terms of affine transformations. The second (combinatorial) part is the so called state-equation which is responsible for the subdivision of the domain of the mapping in polyhedral regions. The well known linear complementarity problem is the essential tool which is used in the state-equation to obtain the required piecewise-linear mapping.

A principal advantage of the state-model over the standard list-oriented specification of PL mappings is the fact that the boundaries of the polyhedral regions are implicitly defined by the state-equation. Therefore it is no longer required to determine explicitly the transformed boundaries when some linear transformation is applied to the PL mapping. It appears that a number of equivalent formulations of the state-model can be given with the property that each equivalent formulation behaves like a kind of natural representation for one particular polyhedral region and still contains all necessary global information about the mapping. In each polyhedral region the state-model can be considered as a system being in a particular state. The natural descriptions for the different regions are obtained by transitions of the corresponding states. The mappings in adjacent regions in space are then described by adjacent states of the state-model. The neighbour relations between the various regions play an important role in the solution algorithms of PL equations applied on the state-model. These neighbour relations are mapped on a graph which is called the "structure" of the state-model. The structure contains information about mutually reachable states or regions.

In particular, a structural degeneracy of the state-model can lead to a non-connected graph, indicating that homotopy algorithms for finding a solution of the PL equations may not function properly. Moreover in this case the mapping is in general non-continuous.

Besides the application of this state-model to the modelling of electronic devices, it also appears to be a new valuable tool in the theory of PL homeomorphisms. Due to its formulation in terms of matrices and vectors instead of regions and corresponding mappings, new theorems about homeomorphism properties can be derived, which express the conditions for a homeomorphism in certain known matrix properties which can easily be verified. The results bear a strong connection with the theory developed by Sandberg and Willson [2].

To introduce the state-model, we start in chapter 2 from the description of piecewise-linear electrical networks containing ideal diodes. The various equivalent hybrid representations of a linear lumped memoryless multiport will form the basis of the state-model and the ideal diodes can be considered as to embody a real physical state. From this approach a physical background for the state-model is supplied which may help to understand the specific properties of the model. It appears that multiport networks which yield a positive output impedance at each port under all possible passive load conditions at the remaining ports can be expected to have a unique DC response when loaded by ideal diodes. In chapter 3 this expectation is substantiated by proving the equivalence of some mathematical conditions defining the so called Class P and  ${\rm P}_{\rm O}$  matrices and certain physical properties of a multiport network regarding the positivity of a set of port immitancies. At the same time, the hybrid structure of a multiport network is introduced as a graph representing the various existing hybrid descriptions of the multiport. This graph will be used in the sequel as an adjacency structure for relating the different linear regions with respect to their topological properties such as neighbour relationship or connectivity. It is shown later on that a structural degeneracy of a multiport network indicated by singular principal submatrices of one of its

hybrid matrices implies that the represented mapping is not uniquely defined.

Based on the earlier observations of the properties of PL networks, the state-model is defined in chapter 4 as a mathematical model for PL mappings. The notions of state and equivalent representations are introduced and brought into relation with the subdivision of the domain of the mapping in polyhedral regions. In parallel with the standard state-space formulation of linear dynamic systems, the controllability and observability of the state-model is discussed. For later application in chapter 7 concerning homeomorphic mappings, the minimal state-model is introduced to restrict the class of PL mappings such that some homeomorphism theorems can be derived later on.

In chapter 5 the standard solution algorithms for the linear complementarity problem are discussed and brought into relation with the state-model formulation. In particular the application of Cottle's algorithm is extensively explained.

Chapter 6 deals with a new approach to the solution of the linear complementarity problem. By applying the so called modulus transformation, this problem is transformed into a specific type of nonlinear equations. The important advantage of this transformation is that the complementarity conditions disappear such that standard methods for the solution of nonlinear equations can be applied to solve the linear complementarity problem.

For a specific class of linear complementarity problems a contraction mapping algorithm is presented which runs in polynomial time-complexity, in contrast with the algorithms of chapter 5 which are known to be exponential in complexity. This new algorithm may save a lot of computation time in worst-case situations. Furthermore the modulus transformation is related to a global model of a PL mapping given by Chua [21], which latter model appears to be a special case of the state-model formulation.

In chapter 7 the state-model of the inverse mapping is derived by a simple linear transformation of the original model. The corresponding operation for a list-specified PL mapping would be much more complicated since the boundaries of the polyhedral regions in the range of mapping then have to be determined one by

one. The straightforward relation between the state-models of a mapping and its inverse mapping means an important advantage which is used to derive some new theorems about homeomorphic mappings. The previously defined structure of the state-model is extensively used in the proofs of the relevant theorems and appears to be a powerful tool. Furthermore explicit state-models are constructed for some relatively simple PL mappings, to be used later on in the formulation of more complex state-models.

Chapter 8 then deals with the construction of these complex state-models. To this purpose a state-model is considered as a description of a black box representing some PL system. Piecewiselinear systems of increasing complexity are then constructed by interconnection of PL subsystems just like the construction of a complex electronic circuit from standard electronic devices. The mathematical implementation of an interconnection amounts simply to the catenation of a linear relation between interconnected variables. Depending on the type of physical system that is to be represented by the state-models each interconnection may require a set of linear relations te be added to the system. For example a galvanic interconnection of electrical components requires some voltages to become equal and the sum of some currents te become zero. The proposed construction results in a hierarchically based modelling process in which each model can be used as a submodel or component in a more complex overall model.

The physical origin of the models is immaterial for the construction and analysis of the overall model, since we are always dealing with PL mappings. For example a PL segment originating from the description of a digital gate is indistinguishable from that of a bipolar transistor.

This property of the modelling is extremely important and allows for an efficient mixed-level simulation of electrical networks since a uniform database and solution algorithm can be applied.

The remaining part of chapter 8 is devoted to the implementation of dynamic PL systems, which appear to fit in a natural way to the

state-model description. By application of an integration formula, the differentiation operator can be replaced by an algebraic relation, yielding a state-model which may be considered as a PL equivalent of the transition matrix arising in the description of standard nonlinear dynamic systems.

It will be no surprise that time-cusp problems associated with the standard transient analysis of certain electronic circuits are completely absent in the state-model formulation, due to the fact that a sequence of two or more pivoting steps in one of the pivoting algorithms may yield a discontinuous response.

Finally in chapter 9 some DC models are constructed for a couple of common electronic devices to give an indication of the general line of the model construction. Two examples are included as an illustration of the capabilities of the presented theory.

#### 2. Piecewise-linear electrical networks

As discussed in the introduction, our intention is to find some kind of closed form analytical description of a piecewise-linear mapping and a specific algorithm such that the solution of piecewise-linear equations can be determined by applying this algorithm to the model of a corresponding PL mapping. This analytical description or model must be able to cope with multi-dimensional PL mappings and has to cover PL dynamic systems as well. Furthermore, working with lists of linear systems that define the PL mapping in each PL region of the domain of the mapping should be avoided since these lists may grow exponentially long when a concatenation of mappings has to be constructed. Instead of this, one global description should cover the description of all PL regions.

It is obvious that these demands are not easily met when starting from a pure mathematical point of view. As in fact the PL equations of interest originate from linearizing nonlinear electrical networks, it is easy and more natural to attempt to construct the required global model from the description of a class of piecewise-linear networks. Hopefully such a model can be generalized to satisfy the properties required for less restricted PL mappings as well.

To this purpose, we focus our attention to the description of PL electrical networks. For the time being, we will restrict ourselves to PL resistive networks and cover dynamic networks later on.

#### 2.1. The multiport network-model

In order to establish a PL behaviour, a network should contain PL elements besides all other kinds of linear static elements like resistors and independent or controlled linear sources. We will use the ideal diode, to be defined shortly, as the most primitive piecewiselinear element and consider all other PL elements as being constructed from these ideal diodes and the standard linear elements mentioned before. Let us define the ideal diode as a two-pole element with a voltage-current characteristic as given in fig.2.1. Note the orientation of the diode voltage which is opposite to the standard

reference direction. Then the voltage  $\boldsymbol{v}$  and current i of the ideal diode satisfy

$$\mathbf{v} \ge \mathbf{0}, \quad \mathbf{i} \ge \mathbf{0}, \quad \mathbf{v} \cdot \mathbf{i} = \mathbf{0} \tag{2.1}$$





Due to the previous restriction concerning the allowed element types, we can depict the PL network as given in fig.2.2, in which all ideal diodes are pulled out of the network. The remaining multiport network M is then a linear lumped memoryless network containing resistors and all types of fixed or controlled linear sources.



fig.2.2.

Let us assume that the network contains n ideal diodes. We denote the port voltages and port currents of the linear n-port network by

$$v = (v_1, v_2, \dots, v_n)^{t}, \quad i = (i_1, i_2, \dots, i_n)^{t}$$
 (2.2)

It is known that the network M can be described by a so called constraint-matrix description [3], of the following form

$$C_1 v + C_2 i = b,$$
 (2.3)  
with  $C_1 = (m \times n), \quad C_2 = (m \times n), \quad m \le n$ 

If m < n, the network response for any allowed excitation with fixed sources at the ports is apparently non-unique. We will exclude these types of degenerate networks and assume that the linear multiport M is nondegenerate such that at least one excitation with voltage or current sources exists which yield a unique response. That is m = nand at least one reordering of the columns of (2.3) exists which satisfies

$$Nw + Mz = b, \quad det(N) \neq 0$$
(2.4)

and 
$$w_k = v_k$$
,  $z_k = i_k$  or  $w_k = i_k$ ,  $z_k = v_k$ ;  $k=1,2,...,m$ 
  
(2.5)

The vectors w and z in (2.4) satisfying (2.5) are called complementary port vectors. From (2.4) we then obtain the familiar hybrid description

$$w = Hz + a,$$
with  $H = -N^{-1}M$  and  $a = N^{-1}b$ 
(2.6)

Due to (2.5), there are 2<sup>n</sup> different partitions of v and i which may yield 2<sup>n</sup> different but equivalent hybrid descriptions of the form (2.6). By virtue of the previous definition, a nondegenerate network has at least one hybrid representation.

As for the ideal diodes, we have by (2.1) and (2.5) for each ideal diode the relations  $w_k \ge 0$ ,  $z_k \ge 0$  and  $w_k z_k = 0$ , which is equivalent to

$$w \ge 0, z \ge 0, w^{t}z = 0$$
 (2.7)

the inequalities taken componentwise.

Hence the PL network of fig.2.2 is fully characterized by

$$w = Hz + a$$
  
 $w \ge 0, z \ge 0, w^{t}z = 0,$ 
(2.8)

with w and z some particular complementary partition of the diode currents and voltages.

#### 2.2. Some considerations about the network response

The network response is obtained by finding a solution pair w,z which satisfies (2.8). This general problem is known in the field of mathematical programming as the linear complementarity problem (LCP problem). Algorithms to solve this problem will be discussed in detail in chapter 5. However at this place it will be convenient to sketch the idea behind these algorithms in terms of operations on the network of fig.2.2.

It is obvious that an excitation of the network, which defines a vector a in (2.8), influences the particular states of the ideal diodes. That is, depending on the the vector a, some diodes will conduct ( $v_k = 0$ ,  $i_k \ge 0$ ) and the remaining will block ( $v_k > 0$ ,  $i_k = 0$ ). As a conducting diode forms a short-circuit and a blocking diode forms an open-circuit, the state of the diodes in turn influences the topology and thus the response of the network M. Then the network of fig.2.2 can be seen as the continuous counterpart of a finite state machine, with M corresponding to the combinatorial logic and the collection of diodes corresponding to the sequential part, memorizing the particular state (see fig.2.3).



fig.2.3.

Calculation of the response to some excitation is now equivalent to the determination of the diode state, i.e. the specific conducting state of all diodes. A direct approach to do so would be to consider all  $2^{n}$  possible diode states and solve for the response. If for some diode state a response of M will yield the same state as was initially assumed, then a solution is found. In fact, this suggests some iterative algorithm visualized in fig.2.4, which is obtained from fig.2.3 by opening the lower feedback loop.



fig.2.4.

In fig.2.4  $x_{n+1}$  is some function of  $x_n$  and the input e, say

$$\mathbf{x}_{n+1} = \mathbf{f}(\mathbf{x}_{n}, \mathbf{e}) \tag{2.9}$$

The solutions of (2.8) are then fixed-points of the mapping f(x,e) and, under certain conditions, equation (2.9) can be used as a contraction-mapping algorithm to find those solutions. The details of such an algorithm are discussed in chapter 6.

Another approach would be to move the input e gradually from a value  $e_0$  for which a solution is known, to the value  $e^*$  for which the solution is required, meanwhile keeping track of all changes in the diode states. In this way the solution is continuously embedded in a (one) parameter family of solutions which can be found by continuation methods. A particular algorithm of this class is the well known Katzenelson algorithm [4].

All algorithms have in common that they run through a certain sequence of diode states ending in the required solution sooner or later.

To facilitate further discussions we make the following definitions.

A port of a linear resistive multiport network is called a voltage (current) port if the excitation at that port is performed by a voltage (current) source. A port is an uncommitted port if the excitation is not specified. [See e.g. 5]

Let us from now on assign  $z_k = v_k$  to a voltage port and  $z_k = i_k$  to a current port, without leaving any uncommitted ports. If, for a particular port assignment, the matrix N in the relation corresponding to (2.4) is nonsingular, we say that the corresponding port assignment is admissible. Then for every excitation leading to an admissible port assignment, the network M has a unique response and there exists a hybrid representation of that network of the form (2.6).

Now, each diode state corresponds to a certain open- or shortcircuit termination of the ports of M. An open-circuit can be seen as a current source of strength zero and a short-circuit as a voltage source of strength zero. In this way, each diode state is related to a particular port assignment. Then for a diode state yielding an admissible port assignment, we have from (2.6) that w = a since z = 0.

According to (2.8), for this state to arise it is in addition required that a  $\geq 0$ . Hence we may find the solutions of (2.8) by investigating for all admissible port assignments whether they yield a nonnegative vector a in the corresponding hybrid representation (2.6).

Depending on the properties of M, multiple solutions or even no solution at all may exist as can be seen from the simple one-port example of fig.2.5. This circuit has two solution for  $e \ge 0$  and no solution for  $e \le 0$ .



It is already known that the LCP problem (2.8) has a unique solution for any vector a, if and only if  $H \in P$  [6], [10], where the properties of the matrix class P have been given by Fiedler and Ptåk [8]. If the vector a is restricted to some subspace of  $\mathbb{R}^n$ , no such general statement is known by now. However, an upper bound for the number of different diode states has been given by Belevitch [9]. The properties of class P will be discussed in equivalent network terms in chapter 3.

#### 2.3. Piecewise-linear two-pole elements

The construction of a PL equivalent of a nonlinear electronic network requires PL descriptions to be derived for various nonlinear devices such as for example bipolar transistors, diodes, tunnel diodes and MOS transistors. The nonlinear models of these devices currently employed in simulation packages have the property that the static nonlinearities can be incorporated by properly interconnected nonlinear two-pole elements. Then piecewise-linear modelling of these devices amounts to find PL models for onedimensional nonlinear functions, describing the voltage-current characteristics of the particular two-pole elements. Let us assume that a certain two-pole element has been approximated by a finite number of linear segments in a sufficiently accurate way, yielding a continuous PL voltage-current relation with positive slopes for each segment (see fig.2.6). Then a network model in the sense of section 2.1 can be obtained as follows.



Let the modelling be conducted by going bottom-up along the voltage-axis and assume that a network model  $M_k$  has been derived which describes the required v-i curve for  $v \leq v_k$ , having a fixed slope  $R_k$  for  $v \geq v_{k-1}$ . We then extend the network to  $M_{k+1}$  in such a way that the interval  $(v_{k-1}, v_k)$  is included with a slope  $R_{k+1}$  for  $v \geq v_k$ . Repetition of this step will then yield the required model after all segments have been processed.

When extending M, two situations may arise with either  $R_{k+1} > R_k$  or  $R_{k+1} < R_k$ . In the first case we define  $M_{k+1}$  to be the network of fig.2.7. Since  $R_{k+1} > R_k$ , the resistor to be added is positive.



fig.2.7.

For  $i \leq i_k$ , diode  $D_k$  conducts and the behaviour of  $M_{k+1}$  is identical to that of  $M_k$ . As soon as  $i > i_k$  the diode blocks and the slope for  $v > v_k$  becomes  $R_k + (R_{k+1} - R_k) = R_{k+1}$ , as required. When  $R_{k+1} < R_k$ , the dual situation applies, leading to the network of fig.2.8.



Since a model for the lowest (first) segment is a simple resistor in series with some fixed voltage source, the required process can be performed step by step. Note that the above synthesis also includes segments with zero or infinite slope which are just on the edge of realizability.

Hence monotonous PL voltage-current characteristics with nonnegative slope can be modelled with positive resistors, fixed sources and ideal diodes only. This important property will be used later on when the solvability of PL networks will be discussed.

For voltage-current characteristics with negative slopes such as for example arising intunnel diodes, the given synthesis leads to negative resistors. However, if the voltage-current characteristic under consideration is such that the voltage is a single-valued function of the current or the other way round, the synthesis can be performed with only one negative resistor.

As can be seen from fig.2.9, for the case v = v(i) the circuit can be realized by a negative resistor  $R^* = \min_i (R_i)$  in series with i a network having a monotonously increasing v-i curve.



## 3. Equivalent network properties of P and P matrices

As we showed in chapter 2, formulation of the network equations of a class of piecewise-linear networks results in an equation set known as the linear complementarity problem, given by

$$w = Hz + a \qquad H = (n \times n)$$

$$w \ge 0, \quad z \ge 0, \quad w^{t} = 0 \qquad w, z \in \mathbb{R}^{n}$$
(3.1)

In (3.1) the elements of the vectors w and z are the port voltages and port currents of a linear multiport network M and H represents a hybrid matrix describing the electrical behaviour of M. The independent port variables are collected in z and the complementary variables in w.

For the moment let us assume that the impedance matrix R of M does exist. With z = i and w = v, the network M is then described by

$$\mathbf{v} = \mathbf{R}\mathbf{i} + \mathbf{a} \tag{3.2}$$

Equation (3.2) leads to a Thévenin equivalent of M as depicted in fig.3.1.



fig.3.1.

fig.3.2.

The internal port resistances are given by the diagonal elements of R and the controlled voltage sources  $e_{L}$  are defined by

$$e_{k} = \sum_{\substack{j=1\\ j \neq k}}^{n} R_{kj} i_{j} + a_{k}$$
(3.3)

Next, the ports of M are loaded with ideal diodes, yielding for each port an equivalent circuit of the type given in fig.3.2. Now, for a single port the conducting state of the diode is unique for any fixed value of  $e_k$  if and only if  $R_{kk} > 0$ , because  $R_{kk} \leq 0$  would lead to a situation already depicted in fig.2.5 and  $R_{kk} > 0$ , according to fig.3.3, yields a unique solution for any  $e_k$ .



fig.3.3.

However, since  $e_k$  represents a controlled voltage source, its value cannot be chosen independently but is determined by the conducting state of the diodes at all other ports. Then a negative value of  $R_{kk}$  not necessarily yields multiple solutions since an alternate solution at port k may be prohibited by the conducting states of the diodes at the remaining ports. However for any R which is lower triangular, a necessary and sufficient condition for the existence of a unique solution for any vector a, is given by  $R_{kk} > 0$  for all k, which can be proved as follows. The first equation of (3.2) will give

$$v_1 = R_{11}i_1 + a_1$$
  
 $v_1 \ge 0, i_1 \ge 0$   $v_1i_1 = 0,$   
(3.4)

independent of the state of other diodes. The solution of (3.4) exists and is unique for any vector a if and only if  $R_{11} > 0$ .

Assume that a unique solution for the first k-1 equations exists for any vector a, independent of the state of the remaining diodes. The k-th equation now reads

$$\mathbf{v}_{\mathbf{k}} = \sum_{\mathbf{j} < \mathbf{k}} \mathbf{R}_{\mathbf{k}\mathbf{j}} \mathbf{i}_{\mathbf{j}} + \mathbf{R}_{\mathbf{k}\mathbf{k}} \mathbf{i}_{\mathbf{k}} + \mathbf{a}_{\mathbf{k}}$$

$$\mathbf{v}_{\mathbf{k}} \ge 0, \quad \mathbf{i}_{\mathbf{k}} \ge 0 \qquad \mathbf{v}_{\mathbf{k}} \mathbf{i}_{\mathbf{k}} = 0$$
(3.5)

Due to the assumption, the value of  $\sum_{j < k} R_{kj}i_j + a_k$  is fixed and does not depend on  $v_m$  and  $i_m$  for  $m \ge k$ . Hence a unique solution of (3.5), and thus of the first k equations, exists for any vector a if and only if  $R_{kk} > 0$ . Then by induction the solution of (3.2) satisfying  $v \ge 0$ ,  $i \ge 0$  and  $v^t i = 0$  exists and is unique for any vector a if and only if  $\forall_{v} [R_{kk} > 0]$ .

Of course, the above statement also holds true for any hybrid matrix which can be reordered into lower triangular form.

For non-triangular matrices the situation is more complicated.Instead of considering the value of  $e_k$  as a function of the voltages at all other ports, it is more suitable to consider the Thévenin equivalent of port k for any allowed open- or short-circuit termination at the other ports. Then for each hybrid representation  $H^j$  of M we have an equivalent circuit at port k given by one of both alternatives of fig.3.4, depending on the particular port-type assignment.



fig.3.4.

The sources in fig.3.4 are now fixed and on account of the previous arguments one might expect a unique solution to exist if the diagonal elements of all hybrid representations  $H^{j}$  are positive. In other words, for any port k, the output immittance should be positive for any open- or short-circuit termination at all other ports. We will call such a property the "zero load positive immittance property". In the next sections we will derive some theorems concerning the solution of (3.1) in terms of such properties.

#### 3.1. Relations between hybrid representations

Consider again a partition of the constraint equations (3.1) induced by a particular admissible port-type assignment, with  $z_k = i_k$ ,  $w_k = v_k$  for a current port and  $z_k = v_k$ ,  $w_k = i_k$  for a voltage port, yielding

Nw + Mz = b,  $det(N) \neq 0$ Hence

$$w = Hz + a \tag{3.6}$$

with H being a hybrid matrix of M. We are now concerned about the existence of other hybrid representations equivalent to (3.6). We are only interested in H and therefore take a=0. Furthermore we reorder the variables in (3.6) and then partition the vectors w and z in two parts, resulting in

$$\begin{pmatrix} w^{1} \\ - \\ w^{2} \end{pmatrix} = \begin{pmatrix} H_{11} \\ H_{21} \\ H_{21} \\ H_{22} \end{pmatrix} \begin{pmatrix} z^{1} \\ - \\ z^{2} \end{pmatrix}$$
(3.7)

Let  $H_{22}$  be nonsingular. Then the second equation set of (3.7) yields

$$z^{2} = -H_{22}^{-1} H_{21} z^{1} + H_{22}^{-1} w^{2}$$
(3.8)

Substitution of (3.8) into the first equation set of (3.7) then

results in

$$w^{1} = (H_{11} - H_{12} H_{22}^{-1} H_{21}) z^{1} + H_{11} H_{22}^{-1} w^{2}$$
 (3.9)

Then (3.8) and (3.9) can be combined into

$$\begin{pmatrix} w^{1} \\ -\frac{z^{2}}{z^{2}} \end{pmatrix} = \begin{pmatrix} H_{11}^{\prime} & H_{12}^{\prime} \\ -\frac{H_{21}^{\prime}}{z^{2}} & H_{22}^{\prime} \end{pmatrix} \begin{pmatrix} z^{1} \\ -\frac{z^{2}}{w^{2}} \end{pmatrix} , \text{ with}$$
(3.10)

$$H_{11}' = H_{11} - H_{12} H_{22}' H_{21}' \quad H_{12}' = H_{11} H_{22}^{-1}$$

$$H_{21}' = -H_{22}^{-1} H_{22} \quad \text{and} \quad H_{22}' = H_{22}^{-1}$$
(3.11)

Equation (3.11) will be referred to as a "principal transformation". Since w' =  $\binom{w^1}{-\frac{z^2}{z}}$ ,  $z' = \binom{z^1}{-\frac{z^2}{w^2}}$  is again a complementary partition satisfying (1.5), equation (3.10) defines a new hybrid description with hybrid matrix H'. Hence each nonsingular principal submatrix of H yields a new hybrid description. The converse statement is also true.

## Lemma 1: Suppose both representations (3.7) and (3.10) to exist. Then $H_{22}$ and $H_{22}'$ are both nonsingular.

Proof: Subsitute  $z^1 = 0$  in (3.7) and (3.10), which yields  $w^2 = H_{22} z^2$  and  $z^2 = H_{22} w^2$ , leading to  $z^2 = H_{22} H_{22} z^2$ . Since (3.7) and (3.10) describe the same network, the last relation must hold for any  $z^2$ , i.e. we must have  $H_{22}^1 \cdot H_{22} = I$ . Therefore both  $H_{22}^1$  and  $H_{22}$  are nonsingular and  $H_{22}^1 = H_{22}^{-1}$ .

As a consequence, given a hybrid matrix H of a linear multiport network M, the set of existing hybrid matrices is fully characterized by the set of nonsingular principal minors of H. After a principal transformation, any nonsingular principal minor yields an alternate hybrid matrix and a corresponding admissible distribution of openand short-circuit terminations of the network ports. Due to lemma 1 a singular principal minor also means that the corresponding hybrid matrix does not exist. In the sequel we will adopt the following notation to circumvent the reordering of the variables in (3.6).

With the set I = {m<sub>1</sub>,m<sub>2</sub>,...,m<sub>r</sub>}  $\subset$  {1,2,...,n} we define H<sub>I</sub> as the principal submatrix of H with its elements defined by (H<sub>I</sub>)<sub>p,q</sub> = H<sub>p,mq</sub>. In other words, the elements of H<sub>I</sub> are exclusively taken from the rows and columns m<sub>1</sub>,m<sub>2</sub>,...,m<sub>r</sub> of H. Furthermore, det(H<sub>I</sub>) denotes the determinant of H<sub>I</sub>. If I =  $\phi$ , we define det(H<sub> $\phi$ </sub>) = 1. Let H be some hybrid matrix of the linear multiport M. For all H<sub>I</sub> with det(H<sub>I</sub>)  $\neq$  0 we denote the corresponding hybrid representation of M obtained after a principal transformation with H<sub>r</sub> by H|I. Then the following lemma holds.

Lemma 2: The diagonal elements  $h'_{ii}$  of H|I are given by

$$h_{ii}' = \begin{cases} \frac{\det(H_{I} \setminus \{i\})}{\det(H_{I})} & \text{for } i \in I \\ \\ \frac{\det(H_{I \cup \{i\}})}{\det(H_{I})} & \text{for } i \notin I \end{cases}$$
(3.12)

Proof: Consider the equation  $H_{I} x = y$ , with

 $y = 1_k^t = (0, 0, \dots, 0, 1, 0, \dots, 0)^t$ . Since  $H_I$  is nonsingular we have  $x = H_I^{-1}y$  and hence by (3.11)  $x_k = (H_I^{-1})_{kk} = h_{kk}^i$ . According to Cramer's rule, we also have

$$\mathbf{x}_{k} = \frac{\det(\mathbf{H}_{I} \setminus \{i_{k}\})}{\det(\mathbf{H}_{I})}, \text{ which proves (3.12) for } i \in I.$$

In particular, a principal transformation on a single diagonal element h<sub>i</sub> will yield

$$\mathbf{h}_{\mathbf{i}\mathbf{i}}' = \frac{\det(\mathbf{H}_{\phi})}{\det(\mathbf{h}_{\mathbf{i}\mathbf{i}})} = \frac{1}{\mathbf{h}_{\mathbf{i}\mathbf{i}}}$$
(3.13)

Then the value  $h'_{jj}$  for  $j \notin I$  can be obtained by a principal transformation on the indices  $I \cup \{j\}$  leading to  $\widetilde{H}$ , followed by a transformation on  $\widetilde{h}_{jj}$  yielding  $h'_{jj} = \frac{1}{\widetilde{h}_{jj}}$ . Assume that  $H_{I\cup\{j\}}$  is nonsingular. From the first part of the proof we then have

$$\widetilde{\mathbf{h}}_{jj} = \frac{\det(\mathbf{H}_{\mathbf{I}\cup\{j\}}\setminus\{j\})}{\det(\mathbf{H}_{\mathbf{I}\cup\{j\}})} = \frac{\det(\mathbf{H}_{\mathbf{I}})}{\det(\mathbf{H}_{\mathbf{I}\cup\{j\}})} \text{, yielding}$$

$$\mathbf{h}_{jj}^{*} = \frac{1}{\widetilde{\mathbf{h}}_{jj}} = \frac{\det(\mathbf{H}_{\mathbf{I}\cup\{j\}})}{\det(\mathbf{H}_{\mathbf{I}})} \text{, j \notin \mathbf{I}}$$
(3.14)

When  $\det(H_{I\cup\{j\}}) = 0$  then  $h'_{jj}$  has to be zero, otherwise it would be possible to transform H|I into a new  $\widetilde{H}$  with the portvariables with index j also interchanged. By lemma 1 the existence of  $\widetilde{H}$  would imply that  $H_{I\cup\{j\}}$  is nonsingular, which is a contradiction. Thus (3.14) also holds for  $\det(H_{I\cup\{j\}}) = 0$ . End of proof.

A principal transformation using a single diagonal element will be called a one-step principal transformation. Using  $h_{ii}$  to this purpose, we find from lemma 2 that

$$h'_{ii} = \frac{1}{h_{ij}}, \quad h'_{jj} = \frac{\det(H_{\{i,j\}})}{h_{ij}}, \quad j \neq i$$
 (3.15)

From (3.15) it can be seen that such a transformation is identical with a single Gauss-Jordan elimination step. One such step can be seen as changing the type of port from current port to voltage port or vice versa.

#### 3.2. The hybrid structure of M

Given two existing hybrid matrices H and H', it will be useful to investigate whether or not H' can be derived from H by a sequence of one-step principal transformations. In fact the question is raised whether H|I, obtained by partial inversion of H on the matrix  $H_I$ , can be found by a sequence of Gaussian elimination steps on the diagonal pivots of  $H_I$ . To this purpose we define a partial ordering on the hybrid matrices of the network M. By lemma 1 this amounts to defining a partial ordering on all nonsingular principal submatrices of some hybrid matrix H. Let the determinants of these submatrices, considered as functions of their coefficients, be the elements of a set V and let the relation  $\langle$  be defined on V in the following way:  $det(H_I) \langle det(H_J)$  if  $I \subset J$ . Then the relation  $\langle$  satisfies

a) det(H<sub>1</sub>) < det(H<sub>1</sub>)

b) 
$$\det(H_{I}) \prec \det(H_{J})$$
 and  $\det(H_{J}) \prec \det(H_{K}) \rightarrow \det(H_{I}) \prec \det(H_{K})$   
c)  $\det(H_{I}) \prec \det(H_{I})$  and  $\det(H_{I}) \prec \det(H_{I}) \rightarrow \det(H_{I}) = \det(H_{I})$ 

Hence  $\prec$  indeed defines a partial ordering on the determinants of all principal minors of H, which can be visualized in a graph called the Hasse diagram. In this diagram, there is an edge between vertices corresponding to det(H<sub>T</sub>) and det(H<sub>T</sub>) iff

$$det(H_{I}) \prec x \prec det(H_{J}) \Rightarrow (x = det(H_{I})) \lor (x = det(H_{J}))$$

The principal minors corresponding to a pair of vertices connected by a single edge will be called adjacent minors. As an example fig.3.5 depicts the Hasse diagram for a  $2\times 2$  matrix H.



Now, by lemma 1, for each vertex with  $\det(H_{I}) \neq 0$  there exists a corresponding hybrid representation H|I of M. If in addition an adjacent minor  $H_{J}$  exists with  $\det(H_{J}) \neq 0$ , then the corresponding hybrid representation H|J can be derived from H|I by a one-step principal transformation on the diagonal element with index  $k = (I\cup J) \setminus (I\cap J)$  of H|I. By lemma 2, this diagonal element then has a value  $\theta(I,J)$  given by

$$\theta(\mathbf{I},\mathbf{J}) = (\mathbf{H} \mid \mathbf{I})_{\mathbf{k}\mathbf{k}} = \frac{\det(\mathbf{H}_{\mathbf{J}})}{\det(\mathbf{H}_{\mathbf{I}})}$$
(3.16)

with  $k = (I\cup J) \setminus (I\cap J)$ 

On the other hand H|I can also be obtained from H|J by pivoting on

$$(\mathbf{H} \mid \mathbf{J})_{\mathbf{kk}} = \theta(\mathbf{J}, \mathbf{I}) = \frac{\det(\mathbf{H}_{\mathbf{I}})}{\det(\mathbf{H}_{\mathbf{J}})} = \frac{1}{\theta(\mathbf{I}, \mathbf{J})}$$

We can use this property to assign an arbitrary direction and a corresponding real number  $\theta$  to each edge in the Hasse diagram for which the determinants at the corresponding vertices are nonzero, as indicated in fig.3.6.



For those vertices with  $det(H_I) = 0$  we know that the hybrid matrix  $H \mid I$  does not exist. From now on, we delete these vertices and the incident edges .

Whenever appropriate the above modified Hasse diagram will be called *the hybrid structure of H or the hybrid structure of the network M*. Note that the hybrid structure of M is not unique, however all hybrid structures are isomorphic. For a proof see appendix I ,where also the following more general relation is derived.

$$det(H|I)_{J} = \frac{det(H_{I-J})}{det(H_{I})}$$
(3.17)

As an example, the hybrid structure for  $H = \begin{pmatrix} 1 & 2 \\ 1 & 2 \end{pmatrix}$  is given in fig.3.7.



fig.3.7.

From the previous considerations we have found that the weights  $\theta$  on the edges of the structure of H are given by all nonzero diagonal elements  $h_{ii}^{j}$  of all existing hybrid representations  $H^{j}$  of M. Due to its construction, the structure of H may be non-connected, that is in general its structure is a collection of disjunct but connected substructures.

Let us denote a directed edge between two adjacent minors  $H_I$ and  $H_J$  by e(I,J) (see fig.3.6). Furthermore, the number of incident edges at any vertex  $H_I$  in a subgraph of the hybrid structure is denoted by d(I). Within each connected substructure S we define a directed path P(I,J) between two vertices  $det(H_I)$  and  $det(H_J)$  as a subgraph of S satisfying

a)  $P(I,I) = \Phi$ 

b)  $P(I,K) = P(I,L) \cup e(L,K)$ .  $H_{I}, H_{L}, H_{K} \in S$ 

c)  $P(I,K) \rightarrow d(I) = d(K) = 1$ .

Then we have the following theorems.

- Theorem 1: The nonzero diagonal elements of any hybrid representation H<sup>j</sup> of M are positive if and only if, within each connected substructure of H, the determinants of all nonsingular principal submatrices of a particular hybrid representation H of M have the same sign.

The theorem then follows from the fact that, due to lemma 2, the set of all  $\theta(K,L)$  is identical with the set of all nonzero diagonal elements of all existing hybrid representations of M.

- <u>Corollary</u>: The product of all weights on the edges of a directed path between two vertices is the same for each directed path between those vertices.
- Theorem 2: All2<sup>n</sup>hybrid descriptions H<sup>J</sup> of Mexist and have positive diagonal elements if and only if the determinants of all principal minors of some H<sup>k</sup> are all positive.
- Proof: Under the conditions of the theorem, the structure of M is connected and contains  $det(H_{\{i\}}^k) = h_{ii}^k$ . Then the theorem follows from theorem 1.
- <u>Corollary</u>: If the determinants of all principal minors of some H are positive, then the determinants of all principal minors of any other hybrid representation are also positive.

Due to the above theorms we have the important result that two hybrid representations H' = H|I and H'' = H|J can be obtained from each other by a sequence of one-step principal transformations if and only if there exists a path between the vertices corresponding to det( $H_{\tau}$ ) and det( $H_{\tau}$ ) in a hybrid structure of M.

### 3.3. The matrix classes P and Po

As a generalization of the concept of positive definite and positive semi-definite matrices, Fiedler and Ptåk introduced two matrix classes called class-P and class-P<sub>0</sub>, with the property that for a square matrix H

$$H \in P \quad \Leftrightarrow \forall_{x \neq 0} \quad \exists_{k} x_{k} \cdot (Hx)_{k} > 0 \tag{3.18}$$

$$H \in P_{0} \quad \leftrightarrow \forall \qquad \exists \qquad k \cdot (Hx)_{k} \geq 0 \land x_{k} \neq 0$$
 (3.19)

From these definitions it is obvious that  $P \subseteq P_0$ . In [8] they derived a number of equivalent properties for a matrix to belong to class-P or class-P<sub>0</sub>. The most important of them are listed below. Let D be a diagonal matrix. We denote  $\forall_i d_{ii} > 0$  by D > 0 and analogously  $\forall_i [d_{ii} \ge 0 \land D \neq 0]$  by  $D \ge 0$ .
a) H  $\epsilon$  P  $\leftrightarrow$  the determinant of every principal submatrix of H is positive.

b) H  $\in$  P  $\leftrightarrow$  every real eigenvalue of H is positive.

c)  $H \in P_0 \iff$  the determinant of every principal submatrix of H (3.20) is nonnegative.

d) H  $\epsilon$  P<sub>0</sub>  $\leftrightarrow$  every real eigenvalue of H is nonnegative.

e)  $H \in P_0 \leftrightarrow \det(H+D) \neq 0$  for any D > 0.

From the above properties it can be easily shown that

$$H \in P_{0} \rightarrow \begin{cases} H + D \in P & \text{if } D > 0 \\ \\ H + D \in P_{0} & \text{if } D \ge 0 \end{cases}$$
(3.21)

Furthermore, if  $H \in P(P_0)$  then any principal submatrix of H is in  $P(P_0)$ . The matrix classes P and  $P_0$  play an important role in the solvability and uniqueness of solutions of network equations arising from electronic circuits with bipolar transistors modelled by their Ebers-Moll model, as well as in the linear complementarity problem. In particular we have the following theorem.

Theorem 3: The linear complementarity problem (3.1) has a unique solution for any vector a, if and only if  $H \in P$ .

For a proof see [6] and [10].

In this section we will derive some relations equivalent to (3.20) in terms of network properties. To this purpose we consider the matrix H to represent a  $(n\times n)$  hybrid matrix of some sourceless multiport network M, described by

$$w = Hz$$
, (3.22)

with w and z some complementary partition of the port currents and voltages. Next we load the ports of M in the following way. For all ports except port k we connect an immittance  $d_i$  such that  $w_i = -d_i z_i$  and connect a fixed source to port k such that  $z_k = b_k$ . See fig.3.8 for the case k = 1.



Fig.3.8.

Let the matrix  $D_{\nu}$  be defined by

 $D_{k} = diag(d_{1}, d_{2}, \dots, d_{k-1}, d_{k+1}, \dots, d_{n}).$ 

We say that  $D_k$  forms a positive load of M with respect to port k when  $D_k > 0$ . Next we define  $u_k$  as the immittance seen into port k under the previous load conditions, that is

$$u_{k} = \frac{w_{k}}{z_{k}} = \frac{w_{k}}{b_{k}}$$
(3.23)

With the above  $u_k$  we now define the following properties. M has the "positive load nonnegative immittance property" (M  $\epsilon$  PLNNI) if and only if M is nondegenerate and for any existing hybrid matrix H and for any port k, every positive load with respect to port k yields a nonnegative immittance  $u_k$ . In other words

$$M \in PLNNI \leftrightarrow H_{H} \forall_{k} \forall_{D_{k}>0} [u_{k} \ge 0].$$

M has the "zero load positive immittance property" (M  $\epsilon$  ZLPI) if and only if M is nondegenerate and for some port k, every open- or short-circuit termination at all other ports is admissable and yields  $u_k > 0$ . With the above definitions we have the following theorems. <u>Theorem 4</u>: Given M, then any existing hybrid matrix H  $\epsilon$  P<sub>0</sub> if and only if M  $\epsilon$  PLNNI.

Proof: Necessity: Assume H  $\epsilon$  P<sub>0</sub> and arbitrarily take k = 1 and D<sub>k</sub> = D<sub>1</sub> > 0. The network equations describing the loaded multiport network can now be written as

$$\begin{pmatrix} w_1 \\ \cdots \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = Hz + Dz ,$$

$$(3.24)$$

$$\begin{pmatrix} 0 \mid 0 \end{pmatrix}$$

with 
$$z_1 = b_1$$
 and  $D = \begin{pmatrix} - & - & - \\ 0 & D_1 \end{pmatrix}$ 

Let us denote the matrix obtained from H after deleting the first row and first column by  $H^{(1)}$ . Then Cramer's rule applied to (3.24) gives

$$w_1 \det(H^{(1)} + D_1) = z_1 \det(H + D)$$
 (3.25)

Since  $H \in P_0 \rightarrow H^{(1)} \in P_0$ , we have from (3.21) and  $D_k > 0$ ,  $D \ge 0$  $H^{(1)} + D_1 \in P$  and  $H + D \in P_0$ .

Hence det(H<sup>(1)</sup> + D<sub>1</sub>) > 0 and det(H + D)  $\ge$  0. Then (3.25) yields

$$u_1 = \frac{w_1}{z_1} = \frac{w_1}{b_1} = \frac{\det(H + D)}{\det(H^{(1)} + D_1)} \ge 0.$$

Since k was arbitrarily set to one, we have

$$H \in P_{0} \rightarrow \forall_{k} \forall_{D_{k} \geq 0} [u_{k} \geq 0] \rightarrow M \in PLNNI$$

Sufficiency:

Let  $H \notin P_0$ . Then by (3.20e) there exists a diagonal matrix  $\overline{D}$  such that

$$\det(\mathbf{H} + \mathbf{\bar{D}}) = 0, \ \mathbf{\bar{D}} = \operatorname{diag}(\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_n) > 0$$
(3.26)

Next, for  $i \neq 1$  load all ports i with immittances  $d_i > 0$  and connect a source  $z_1 = b_1$  to port 1. In the same way as for the necessary part we then find corresponding to (3.24)

$$\begin{pmatrix} w_1 \\ \vdots \\ 0 \end{pmatrix} = \begin{pmatrix} H_{11} \vdots H_{12} \\ \vdots \\ H_{21} \vdots H_{22} \end{pmatrix} \begin{pmatrix} z_1 \\ \vdots \\ z_2 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ \vdots \\ 0 & D_1 \end{pmatrix} \begin{pmatrix} z_1 \\ \vdots \\ z_2 \end{pmatrix}$$
(3.27)

with  $z_1 = b_1$  and  $D_1 = diag(d_2, d_3, \dots, d_n) > 0$ . As in (3.25) Cramer's rule will yield

$$w_1 \det(H_{22} + D_1) = b_1 \det(H + \widetilde{D})$$
, (3.28)

with  $\widetilde{D} = \begin{pmatrix} 0 & | & 0 \\ -- & -- \\ 0 & | & D_1 \end{pmatrix}$ .

Now, from (3.26) we find

$$det(H + \vec{D}) = d_1 det(H_{22} + D_1) + det(H + \vec{D}) = 0, \quad d_1 > 0$$
(3.29)

Then from (3.29) either

 a) both determinants in (3.29) are nonzero and have opposite signs;

b) both determinants are zero.

When case a) applies, we obtain from (3.28) that  $u_1 = \frac{w_1}{b_1} < 0$ . In case b) we consider the value of  $h_{11}$ . If  $h_{11} < 0$ , we can find a sufficient high load  $D_1^*$  for the remaining ports such that the immittance seen port 1 is negative since

 $u_1 = h_{11} - H_{12}(H_{22} + D_1^*)^{-1}H_{21} < 0$  for  $D_1^*$  sufficiently large. Thus we assume that  $h_{11} \ge 0$  and have  $h_{11} + d_1 > 0$ . Next we load port 1 with the immittance  $d_1$  and consider the (n-1)-port network formed by the other ports of M. By eliminating  $z_1$  from the n-port equations it is easily seen that the new hybrid matrix describing the (n-1)-port becomes

w = Hz

$$\tilde{H} = H_{22} - \frac{H_{21}H_{12}}{h_{11} + d_1}$$
, (3.30)

and  $w = (w_2, w_3, \dots, w_n)^t$ ,  $z = (z_2, z_3, \dots, z_n)^t$ . Now from (3.26) we have

$$det(H + D) = (h_{11} + d_1) \cdot det(H_{22} + D_1 - \frac{h_{21}h_{12}}{h_{11}+d_1}) = 0,$$
$$h_{11}+d_1 \neq 0.$$

Then with (3.30) we have

$$\det(\widetilde{H} + D_1) = 0 \rightarrow \widetilde{H} \notin P_0 \text{ since } D_1 > 0 \qquad (3.31)$$

Then the above procedure can again be applied on the (n-1)-port network, ultimately yielding some  $u_k < 0$  or ending in a one port network with hybrid matrix  $h \notin P_0$ , i.e. h < 0 in which case the immittance of this port is obviously negative. Hence we have

$$H \notin P_0 \rightarrow \exists_k \exists_{D_k > 0} [u_k < 0] \rightarrow M \notin PLNNI.$$

- Theorem 5: Given M, then for any hybrid matrix H, H  $\epsilon$  P if and only if M  $\epsilon$  ZLPI.
- Proof: Consider some port k of M. The values of u<sub>k</sub> for all open- or short-circuit terminations at all other ports are now identical with the set of all diagonal elements h<sup>j</sup><sub>kk</sub> of all hybrid matrices H<sup>j</sup> of M with port k considered as voltage port and as current port. The theorem then follows immediately from theorem 2.

The conditions for theorem 4 can be relaxed in the following way. From (3.17) we observe that  $\exists H \in P_0 \Rightarrow$  all existing  $H^j \in P_0$ . Since the sufficiency proof for theorem 4 holds for every existing  $H^j$ the PLNNI property can be relaxed to M  $\in$  PLNNI  $\Leftrightarrow \exists_H \forall_k \forall_{D_k} > 0 [u_k \ge 0].$ 

In [5] some additional theorems have been derived along a different line. From theorem 4 we conclude that a multiport network M with a hybrid matrix  $H \in P_0$  can never become "active" at a certain

port when loaded by positive resistors in the sense that no "small signal gain" can be obtained at the input of the port where the excitation is connected to. Thus the network shows a certain kind of passivity which we will call  $P_0$ -passivity. Furtermore , a purely resistive network constructed from PL two-pole elements with a monotone increasing voltage-current characteristic always has a unique solution for any excitation. This is so because according to section 2.3, these elements can be constructed with fixed sources, positive resistors and ideal diodes only. In the model of fig.2.2, the network M is then a multiport network containing only positive resistors and fixed sources and thus its hybrid matrix H is positive definite (PD). Since PD  $\subset$  P, we have H  $\epsilon$  P and by theorem 3 the network solution is unique.

## 4. The state-model of a piecewise-linear mapping

In chapter 2 and 3 we discussed some properties of PL networks formed by linear resistive n-port networks loaded by ideal diodes. It appeared that a hybrid description of the resistive n-port network M played an essential role in the formulation of the network equations and resulted in the LCP-problem (2.8) describing the conducting state of the ideal diodes. Once this state has been determined, the voltages and currents at the ports of M are known and then it is also possible to determine any other voltage or current in the inside of the network M.

Obviously these internal voltages and currents are also piecewiselinear functions of the excitation of the network M and the intention of this chapter is to use these variables to model piecewise-linear mappings. Without loss of generality we may consider the internal excitations and responses to be applied or measured on M at some additional network ports. Hence for a PL mapping  $f : x \in \mathbb{R}^k \rightarrow y \in \mathbb{R}^m$ , let us represent the variables x and y as voltages or currents at some ports of M. In particular let x represent a set of current sources applied to a set of ports #1, let y represent the opencircuit voltages at a set of ports #2 and assume that the ideal diodes are connected to a set of ports #3, as indicated in fig.4.1.



fig.4.1.

Furthermore, let us assume that the impedance matrix Z of M exists and is given by

$$\begin{pmatrix} \mathbf{v}_{1} \\ \mathbf{v}_{2} \\ \mathbf{v}_{3} \end{pmatrix} = \begin{pmatrix} \mathbf{z}_{11} & \mathbf{z}_{12} & \mathbf{z}_{13} \\ \mathbf{z}_{21} & \mathbf{z}_{22} & \mathbf{z}_{23} \\ \mathbf{z}_{31} & \mathbf{z}_{32} & \mathbf{z}_{33} \end{pmatrix} \begin{pmatrix} \mathbf{i}_{1} \\ \mathbf{i}_{2} \\ \mathbf{i}_{3} \end{pmatrix} + \begin{pmatrix} \mathbf{e}_{1} \\ \mathbf{e}_{2} \\ \mathbf{e}_{3} \end{pmatrix}$$
(4.1)

With the load conditions of fig.4.1 we then have

$$i_1 = x, \quad i_2 = 0, \quad v_2 = y$$
 (4.2)

As we are not interested in  ${\bf v}_1^{},$  the remaining equations of (4.1) and (4.2) yield

$$y = z_{21}x + z_{23}i_3 + e_2$$

$$v = z_{31}x + z_{33}i_3 + e_3$$
(4.3)

Of course, some voltages of y may be measured at some ports of #1, in which case  $z_{21}$  in (4.3) would become some submatrix of  $\begin{pmatrix} z_{11} & z_{12} \\ z_{21} & z_{22} \end{pmatrix}$ Let us rename the matrices and variables in (4.3) such that this equation system is rewritten as

$$y = Ax + Bi + g$$

$$(4.4)$$

$$v = Cx + Di + h$$

Due to the ideal diodes connected to the set of ports #3, the network equations of interest finally become

 $y = Ax + Bi + g \tag{4.5a}$ 

 $\mathbf{v} = \mathbf{C}\mathbf{x} + \mathbf{D}\mathbf{i} + \mathbf{h} \tag{4.5b}$ 

$$i \ge 0, v \ge 0, i^{T}v = 0,$$
 (4.5c)

with  $A = (m \times k)$ ,  $B = (m \times n)$ ,  $C = (n \times k)$  and  $D = (n \times n)$ .

In analogy with the state-space formulation of linear dynamic systems given by the equation set

$$y = Ax + Bu$$
  
 $\dot{u} = Cx + Du$ 

we call the system (4.5) the state-model of the mapping  $f : x \in \mathbb{R}^k \rightarrow y \in \mathbb{R}^m$ , and denote the system by S(A,B,C,D,g,h) or shortly S.

Given x, the image vector y is obtained as follows. First for a given x, the pair v, i in the state-model is determined by solving the LCP problem

$$v = Cx + Di + h, v \ge 0, i \ge 0, v^{T}i = 0$$
 (4.6)

This pair v,i then determines the conducting state of the connected diodes, i.e. the state of the system (4.5). Once this state is known, the vector y is obtained from (4.5a). Thus the LCP problem (4.6) plays the same role in the state-model as the state-equation  $\dot{u} = Cx + Du$  in the description of linear dynamic systems. Therefore we will call the equation set (4.5b) and (4.5c) the state-equations of the state-model S. As an example of such a formulation consider the network of fig.4.2, with the nonlinear resistive elements  $R_1$  and  $R_2$  as given in fig.4.3.



## fig.4.2.



fig.4.3.

The network  $\boldsymbol{M}$  is for this case given in fig.4.4 and yields the impedance matrix

$$\mathbf{v} = \begin{pmatrix} 3/2 & 1/2 & 0 & -1/2 \\ 1/2 & 5/2 & 1 & -1/2 \\ 0 & 1 & 1 & 0 \\ -1/2 & -1/2 & 0 & 3/2 \end{pmatrix} \mathbf{i} + \begin{pmatrix} 0 \\ -1 \\ -1 \\ 2 \end{pmatrix}$$
(4.7)



fig.4.4.

From (4.7) the state-model then becomes after renaming the vectors  $v_3^{},\ i_3^{}$  and  $v_4^{}$  ,  $i_4^{}$  in  $v_1^{},i_1^{}$  and  $v_2^{},i_2^{}$ 

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 3/2 & 1/2 \\ 1/2 & 5/2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} 0 & -1/2 \\ 1 & -1/2 \end{pmatrix} \begin{pmatrix} i_1 \\ i_2 \end{pmatrix} + \begin{pmatrix} 0 \\ -1 \end{pmatrix}$$

$$\begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ -1/2 & -1/2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} 1 & 0 \\ 0 & 3/2 \end{pmatrix} \begin{pmatrix} i_1 \\ i_2 \end{pmatrix} + \begin{pmatrix} -1 \\ 2 \end{pmatrix} (4.8)$$

$$v \ge 0, \quad i \ge 0, \quad v^t i = 0.$$

For example, with  $x_1 = 1$  and  $x_2 = 1$ , a solution is given by

$$v_1 = 0$$
,  $i_1 = 0$ ,  $v_2 = 1$ ,  $i_2 = 0$ ;  $y_1 = 2$ ,  $y_2 = 2$ .

Observe that (4.8) itself describes a PL two-port network by its impedance matrix. It can again be used as a PL element for constructing more complex PL networks.

Up to now, we did not show that the mapping f as modelled by (4.5) is indeed PL. This will be postponed until section 4.2. At this place it is however convenient to discuss some properties of (4.5) in network terms.

As in chapter 3, we observe that the (n×n) matrix D considered as a hybrid matrix may yield maximally 2<sup>n</sup> different descriptions equivalent to (4.6). Of course for each description the term Cx + h representing the combined effect of all sources at the ports #3 is in general different too. Let us assume that  $Cx + h = a \ge 0$ . Then a solution of (4.6) is obviously given by v = a, i = 0, yielding y = Ax + g from (4.5a). Hence in a region  $R_0 \subset \mathbb{R}^k$  with  $R_0 = \{x \mid Cx + h \ge 0\}$  we have

 $\mathbf{x} \in \mathbf{R}_0 \rightarrow \mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{g}$ .

Let us denote the other hybrid representations equivalent to (4.6) by

$$w^{j} = C^{j}x + h^{j} + D^{j}z^{j}.$$

Then for these hybrid descriptions we have a corresponding equation for y given by

 $y = A^{j}x + B^{j}z^{j} + g^{j}.$ 

Then

$$\mathbf{x} \in \mathbf{R}^{\mathbf{j}} = {\mathbf{x} | \mathbf{C}^{\mathbf{j}}\mathbf{x} + \mathbf{h}^{\mathbf{j}} \ge 0} \rightarrow \mathbf{y} = \mathbf{A}^{\mathbf{j}}\mathbf{x} + \mathbf{g}^{\mathbf{j}},$$

since  $w^j = C^j x + h^j$ ,  $z^j = 0$  is a solution of the state-equations. Thus, we have maximally  $2^n$  regions in  $\mathbb{R}^k$  where (4.5) defines maximally  $2^n$  affine mappings  $f^j : x \to y$ . Then, in principle the system (4.5) using n diodes is equivalent to a list of  $2^n$  mappings in  $2^n$  regions.

Should two systems S and  $\tilde{S}$  each with  $n_1$  respectively  $n_2$  diodes be interconnected, i.e. some concatenation of two PL mappings be defined, then the overall system contains  $n_1 + n_2$  diodes. This system may describe  $2^{n_1+n_2} = 2^{n_1} \cdot 2^{n_2}$  affine mappings in as many

regions. Hence, the complexity of the state-model description grows linearly while a list of regions and corresponding mappings would grow exponentially. This is a very important advantage of the state-model description over all other descriptions that work with a list of all mappings.

As to the various regions, it is possible that these regions intersect, that is

$$\mathbf{x} \in \mathbf{R}_{\mathbf{i}} \cap \mathbf{R}_{\mathbf{j}} \neq \mathbf{y} = \mathbf{A}^{\mathbf{i}}\mathbf{x} + \mathbf{g}^{\mathbf{i}} \vee \mathbf{y} = \mathbf{A}^{\mathbf{j}}\mathbf{x} + \mathbf{g}^{\mathbf{j}},$$

and the mapping is one to many. The particular value of y which will be found, will then in general depend on the history of x, i.e. on the particular way in which we arrived at x. Such a system has an inherent memory and systems of this type will be used later to model static electronic memory devices like flip-flops or Schmitt-triggers. It may also happen that two different diode states yield the same vector y, which means that some diode states may become unobservable. On the other hand, if  $\bigcup_{i=1}^{k} R_{i} \neq \mathbb{R}^{k}$ , there are some regions in  $\mathbb{R}^k$  for which the mapping  $x \rightarrow y$  is not defined. In network terms this means that no hybrid matrix exists for that particular excitation. Furthermore, if  $\exists_i R_i = \phi$ , this means that the corresponding diode state is uncontrollable by the input vector x, for example due to the fact that a particular part of the network becomes disconnected for that x or two or more diodes are in some way connected in series such that they all have to conduct at the same time.

## 4.1. The structure of the state-model

From now on, we consider the equation set (4.5) as a pure mathematical model describing a piecewise-linear mapping  $f : x \rightarrow y$ , without any specific electrical equivalent in mind for the related variables. However all other equivalent descriptions that originated from our previous network approach, will still be called equivalent hybrid representations of the state-model (4.5). These equivalent representations are based on the different partitions

of the vector pairs v,i in complementary pairs w,z. Due to (4.5c) it is obvious that there exist in total  $2^n$  different classes of complementary pairs v,i if we only consider the distribution of the zero components over i and v and pay no attention to the specific numerical value of their nonzero components. We now say that the system (4.5) can be in one of  $2^n$  different states, depending on the particular complementary partition upon which a certain hybrid description or representation is based.

To facilitate the notation we we will map these partitions on the set of  $2^n$  boolean n-vectors  $s_p$ ,  $p = 0, 1, \ldots, 2^n - 1$  of the boolean n-cube, i.e.

$$s_{p} = (s_{p_{1}}, s_{p_{2}}, \dots, s_{p_{n}}),$$

with  $s_{p_2} \in \{0,1\}$  and  $p = 0, 1, 2, \dots, 2^n - 1$ .

The vectors  $s_p$  are called *state-vectors* and are numbered such that the index p equals the decimal equivalent of  $s_p$  considered as a p binary number. For example for n = 3 we have  $s_6 = (1,1,0)$ .

The particular mapping is defined as follows

$$i_k > 0, v_k = 0 \rightarrow s_p = 1$$
  
 $i_k = 0, v_k \ge 0 \rightarrow s_p = 0$ 

$$(4.9)$$

For example for n = 4 and a complementary solution of (4.5b) and (4.5c) given by i = (2,0,1.5,0), v = (0,0.7,0,1.2) we have from (4.9) s = (1,0,1,0) =  $s_{10}$ . For each  $s_p$  there may exist a corresponding so called *natural equivalent hybrid representation* for (4.5), which describes the same PL mapping and is obtained from (4.5) by a principal transformation using D<sub>I</sub>, with I = {m | s<sub>pm</sub> = 1}. This means that all v<sub>m</sub> and i<sub>m</sub> for which s<sub>pm</sub> = 1 have to be interchanged. In state s<sub>p</sub> the natural equivalent hybrid representation of (4.5) is then defined by

$$y = A^{p}_{x} + B^{p}z^{p} + g^{p}$$

$$w = C^{p}_{x} + D^{p}z^{p} + h^{p}$$

$$w^{p} \ge 0, \quad z^{p} \ge 0 \quad \text{and} \quad w^{p}_{z}z^{p} = 0,$$

$$with \begin{vmatrix} w_{k}^{p} = i_{k} \\ z_{k}^{p} = v_{k} \end{vmatrix} \text{ for } s_{p} = 1 \quad \text{and} \quad z_{k}^{p} = i_{k} \end{vmatrix} \text{ for } s_{p} = 0$$

$$(4.11)$$

The superscript p in (4.10) refers to state s and the representation (4.10) is denoted by S. From (4.11) we observe that (4.5) is the representation in state S. In the above way, each solution of (4.5b) and (4.5c) is coupled to a particular hybrid representation (4.10) by means of its corresponding state-vector. We now investigate for which set of vectors x a particular class of vector-pairs v,i is produced as a solution of the state equations. To this purpose let us consider the class v,i which maps on state s. Due to (4.9) and (4.11) this class satisfies  $z^{p} = 0$ ,  $w^{p} \ge 0$  and hence from representation S. this implies

$$\mathbf{C}^{\mathbf{P}}\mathbf{x} + \mathbf{h}^{\mathbf{P}} \ge 0 \qquad \mathbf{C} = (\mathbf{n} \times \mathbf{k}) \tag{4.12}$$

Hence, when x satisfies (4.12), the particular class v,i which maps on state s contains a solution of the state equations. The n inequalities in (4.12) represent n halfspaces bounded by k-dimensional hyperplanes in  $\mathbb{R}^k$ . The rows of  $\mathbb{C}^p$  represent the normal vectors on these hyperplanes. Then (4.12) is satisfied in the intersection of these halfspaces which defines a convex polyhedron in  $\mathbb{R}^k$ . Let us denote these polyhedral regions by  $\mathbb{R}_p$ , with

$$R_{p} = \{x \mid C^{p}x + h^{p} \ge 0\}$$
(4.13)

In a polyhedron R  $_{\rm p}$  we have w  $^{\rm p} \ge 0$  and  $z^{\rm p}$  = 0. Using the first equation of (4.10), this yields

$$\mathbf{x} \in \mathbf{R}_{\mathbf{p}} \to \mathbf{y} = \mathbf{A}^{\mathbf{p}}\mathbf{x} + \mathbf{g}^{\mathbf{p}}$$
 (4.14)

Therefore in each polyhedron  $R_p$ , the domain x is mapped on the range y by an affine mapping, yielding the following theorem.

Theorem 6: The state-model represents a piecewise-linear mapping  $f : x \in \mathbb{R}^k \rightarrow y \in \mathbb{R}^m$  over a union of polytopes  $\bigcup_{i=1}^{k} R_i \subset \mathbb{R}^k$ .

Of course, all  $S_p$  are equivalent, that is each representation (4.10) can be used to describe the PL mapping in all regions, which means that the state-model is a global model.

As to the regions (4.13), the same considerations as in the previous section apply. Hence there may be empty regions, unobservable or uncontrollable states or states  $s_p$  for which the representation  $S_p$  does not exist. In the latter case it means that a certain principal transformation on D is not possible. In terms of the definitions of chapter 3 we then have

$$S_{p} \xrightarrow{p} det(D_{I}) \neq 0$$

4.15)

with  $I_p = \{m \mid s_p = 1\}$ 

We can use (4.15) to relate each existing  $S_p$  to nonsingular principal submatrices of D and visualize this relation by mapping these  $S_p$  or their corresponding regions  $R_p$  on the vertices of the hybrid structure of D as defined in chapter 3. The structure obtained in this way with the vertices det( $D_I$ ) now labelled by  $S_p$  or  $R_p$  according to (4.15), will be called the structure of the state-model S. A path between two vertices  $S_p$  and  $S_q$  then implies that representation  $S_q$  can be obtained from  $S_p$  by a sequence of one-step principal transformations using only diagonal elements of intermediate matrices D. As a simple example, the state-model of a one-dimensional mapping  $f : x \in \mathbb{R} \rightarrow y \in [-1,+1]$  is given in fig.4.5.



fig.4.5.

The structure of the state-model of fig.4.5 is given in fig.4.6 below.



fig.4.6.

The various "polytopes" are obtained from the given  $S_i$ , yielding

 $R_0 = \{x \mid x + 1 \ge 0, -x + 1 \ge 0\}$   $R_1 = \{x \mid 0.x + 2 \ge 0, -x + 1 \ge 0\}$  $R_2 = \{x \mid x + 1 \ge 0, 0.x + 2 \ge 0\}$ 

Observe that  $S_3$  does not exist since  $det(D_{\{1,2\}}) = 0$ . Furthermore,  $R_0$  overlaps with  $R_1$  and  $R_2$ . As we will see later, this is due to the negative weights on the edges of the structure of S.

## 4.2. Adjacent regions in a minimal state-model

In this section, we will discuss some properties of the statemodel S for a piecewise-linear mapping  $f : x \in \mathbb{R}^m \to y \in \mathbb{R}^m$ , given by

y = Ax + Bi + g v = Cx + Di + h  $A = (m \times m), B = (m \times n)$  (4.16)  $v \ge 0, i \ge 0, v^{t}i = 0$   $C = (n \times m), D = (n \times n)$ 

To facilitate further notations let us define  $\textbf{H}^{p}$  and  $\textbf{I}_{p}$  in representation  $\textbf{S}_{p}$  by

$$H^{P} = \begin{pmatrix} A^{P} & B^{P} \\ - & - & - \\ C^{P} & D^{P} \\ & & \end{pmatrix}$$
(4.17)

$$I_{p} = \{k \mid s_{p_{k}} = 1\}$$
 (4.18)

By the definition of S , any other existing  $H^{\rm q}$  can be derived from  $H^{\rm p}$  by a principal transformation on  $D^{\rm p}_{\rm M}$ , with

$$M = I_{p} - I_{q} \triangleq (I_{p} \cup I_{q}) \setminus (I_{p} \cap I_{q})$$

$$(4.19)$$

From (3.10) and (3.11) this will yield

$$A^{q} = A^{p} - B^{p}_{M} D^{-1}_{M} C^{p}_{M} , \qquad (4.20)$$

with  $B_M^p$  consisting of all columns of  $B^p$  with their indices in M considered as an ordered set and  $C_M^p$  consisting of all rows for  $C^p$  under the same conditions. For a particular example see fig.4.7.



Next consider region R given by

$$R_{p} = \{x \mid C^{p}x + h^{p} \ge 0\}$$
(4.21)

Let the hyperplane  $K_k^p = \{x \mid c_k^p, x + h_k^p = 0\}$ , with  $c_k^p$ , representing the k-th row of  $c^p$ , be a boundary of  $R_p$ . If for some  $1 \neq k$ ,  $K_1^p = K_k^p$ , the boundary hyperplane  $\kappa_{k}^{p}$  would be redundantly determined, yielding an in some sense unnecessarily complicated state-model. Therefore we define a *minimal state-model* as a state-model for which

 $\forall_p \forall_{l \neq k} [\kappa_l^p \neq \kappa_k^p]$  and in addition for each boundary hyperplane the corresponding diagonal element of D is nonzero. In the remaining of this section we will only consider such minimal models. Then we have

$$\exists \sum_{\substack{\mathbf{x} \in \mathbb{R}_{p} \cap \mathbb{K}_{k}^{p}} [c_{k}^{p}.\bar{\mathbf{x}} + h_{k}^{p} = 0 \land c_{1}^{p}.\bar{\mathbf{x}} + h_{1}^{p} > 0, \ 1 \neq k]$$
(4.22)

Let us now apply a one-step principal transformation using  $D_{kk}^p$ , which yields a new representation  $S_1$  with  $H = H^q$  and  $k = I_p - I_q$ . In this representation we have by (3.10) and (3.11)

$$A^{q} = A^{p} - \frac{B^{p}_{\cdot k} C^{p}_{k}}{D^{p}_{kk}}$$

$$C^{q}_{r.} = \begin{cases} C^{p}_{r.} - \frac{D^{p}_{rk}}{D^{p}_{kk}} C^{p}_{k.}, & r \neq k \\ - \frac{1}{D^{p}_{kk}} C^{p}_{k.}, & r = k \end{cases}$$

$$(4.23a)$$

$$(4.23b)$$

$$(4.23b)$$

$$(4.23b)$$

$$(4.23c)$$

$$h_{r}^{q} = \begin{cases} n_{r} & D_{kk}^{p} & r \neq k \\ -\frac{h_{k}^{p}}{D_{kk}^{p}} & r = k \end{cases}$$
(4.23)

and 
$$g^{q} = g^{p} - \frac{h_{k}^{p}}{b_{kk}^{p}} B^{p}_{\cdot k}$$
 (4.23d)

From (4.23b) and (4.23c) this yields

$$C_{1}^{q} x + h_{1}^{q} = C_{1}^{p} x + h_{1}^{p} - \frac{D_{1k}^{p}}{D_{kk}^{p}} (C_{k}^{p} x + h_{k}^{p}) , 1 \neq k$$
 (4.24)

Next consider again the points  $\bar{x}$  from (4.22). For these  $\bar{x}$  we have by (4.24)

$$C_{1}^{q} \cdot \vec{x} + h_{1}^{q} = C_{1}^{p} \cdot \vec{x} + h_{1}^{p} > 0 , \quad 1 \neq k$$

$$C_{k}^{q} \cdot \vec{x} + h_{k}^{q} = C_{k}^{p} \cdot \vec{x} + h_{k}^{p} = 0$$

$$(4.25)$$

Hence  $K_k^p = \{x \mid C_{k}^p x + h_k^p = 0\}$  is also a boundary hyperplane of  $R_q$ . Then  $R_p$  and  $R_q$  are adjacent regions which have at least a subset of the hyperplane  $K_k^p = K_k^q$  in common (see fig.4.8).



As to the affine mappings in both regions we have from (4.23a) and (4.23d)

$$A^{q}x + g^{q} = A^{p}x + g^{p} - \frac{1}{D^{p}_{kk}} B^{p}_{\cdot k} (C^{p}_{k} x + h^{p}_{k})$$
(4.26)

Then for any  $\overline{x}$  on the boundary between  $\underset{p}{R}$  and  $\underset{q}{R}$ , we have from (4.22) and (4.26)

$$A^{q}\bar{x} + g^{q} = A^{p}\bar{x} + g^{p}$$
(4.27)

Hence a PL mapping f modelled by a minimal state-model is continuous since it is continuous within each region as well as on the boundaries between the regions. Furthermore the representation in adjacent regions can be found by a one-step principal transformation using the diagonal element of D with the same index as the common boundary hyperplane which separates them.

In network terms we could think of each boundary hyperplane as being coupled to a single ideal diode which has to reverse its state when the exitation traverses from one region to an adjacent region. In a nonminimal model at least one other ideal diode could change its state at the same time, which in general implies that the mapping f is not continuous.

The minimal state-models will play an important role in describing homeomorphisms as we will see in chapter 7.

# 5. Complementary pivoting methods

In the previous sections the state-model was presented as a global model for a PL mapping f :  $x \in \mathbb{R}^m \to y \in \mathbb{R}^k$ . The calculation of the image vector y for a given vector x by this model is reduced to the determination of a solution of the linear complementarity problem given by the state equations

$$v = Di + Cx + h$$
 (5.1)  
 $v \ge 0, i \ge 0, v^{t}i = 0$   $D = (n \times n), C = (n \times m)$ 

The existance of efficient solution algorithms for the LCP problem is then of vital importance for a practical application of the statemodel. Fortunately various such solution algorithms exist and there is a vast amount of literature on this topic (see e.g. [11] and [12]).

The algorithms to solve the LCP problem can in general be divided into three different classes, depending on the particular equivalent formulation of the LCP problem.

The first class of algorithms are so called homotopy algorithms. In these algorithms, the problem (5.1) is embedded in a one-parameter family of problems  $H(v,i,\lambda)$  defined on  $\mathbb{R}^n \times [0,1]$ , in such a way that the solution H(v,i,0) is known and the solution of H(v,i,1) is a solution of (5.1). The algorithms in this class generate a family of solutions  $v(\lambda), i(\lambda)$  along some path in  $\mathbb{R}^n \times [0,1]$ .

In the most familiar homotopy methods, the solution path is generated by a sequence of one-step principal transformations (Gauss-Jordan pivoting steps), in each step extending the already obtained path by some segment, meanwhile keeping the solution complementary or almost complementary. Such methods will be called complementary pivoting methods and will be the subject of this chapter. In particular the Katzenelson algorithm, the Cottle algorithm and the Lemke algorithm will be discussed.

A second class of algorithms is formed by variational algorithms which minimize some constrained cost-function in  $\mathbb{R}^n$  which is formulated in such a way that its solution is equivalent to a solution of (5.1). As an example, Cryer [13] reformulated (5.1)

for symmetric D into

minimize  $f(i) = \frac{1}{2}i^{t}Di + a^{t}i$ , subject to  $i \ge 0$ 

In this way a lot of constraint minimization-algorithms can be used to solve the LCP problem. However in most situations these algorithms only work properly when f(i) is convex, which limits the applicability of these methods. Therefore we will not discuss these methods any further.

The third class of methods are the iterative methods. In this case, the LCP problem is transformed into an equivalent set of nonlinear equations in  $\mathbb{R}^n$ . By using adequate fixed-point algorithms, the LCP problem can be solved iteratively. A special representative of this class is the modulus algorithm which will be treated in chapter 6. Other algorithms of this type can be found in the simplicial fixed-point methods [14].

### 5.1. The Katzenelson algorithm

Katzenelson [4] originally devised his algorithm for solving a set of PL equations f(x) = y defined by

$$x \in R_{i} \rightarrow f(x) = f^{i} = A^{i}x + b^{i}, A^{i} = (n \times n), i = 0, 1, \dots, k,$$
  
(5.2)

with  $R_i$  representing convex polyhedral regions in  $\mathbb{R}^n$ . The particular set of PL equations evolved from a description of a PL resistive network containing uncoupled resistors with strictly monotonously increasing v-i characteristics. For these equations it could easily be shown that

$$\forall_{i,j} \operatorname{sign}(\det A^{i}) = \operatorname{sign}(\det A^{j})$$
(5.3)

From our discussions in section 3.3 we already know that in such a case always a unique solution exists.

In [15] and [16] Fujisawa and Kuh respectively Kuh and Chien extended this algorithm for less restricted conditions than (5.3).

In this section we will describe a version of the algorithm which can be used to solve the LCP problem (5.1) for matrices C and D such that the pivots that are taken in the algorithm are always positive. A sufficient condition to this purpose is given by  $D \in P$ .

Let us assume that we are interested in the solution of (5.1) for  $\mathbf{x} = \mathbf{x}^*$ . We start from a point  $\mathbf{x}_0$  for which the solution of (5.1) is known. Let this point be such that  $C\mathbf{x}_0 + \mathbf{h} > 0$ , yielding  $\mathbf{v} = C\mathbf{x}_0 + \mathbf{h}$ ,  $\mathbf{i} = 0$ . Next we construct a path  $\mathbf{x}(\lambda) = \mathbf{x}_0 + \lambda(\mathbf{x}^* - \mathbf{x}_0)$ ,  $0 \le \lambda \le 1$  in  $\mathbb{R}^n$  and try to find the solution of (5.1) for any  $\lambda$  within the given interval. From (5.1) and the given  $\mathbf{x}(\lambda)$  we then have

$$v(\lambda) = Di(\lambda) + Cx_0 + h + \lambda C(x^* - x_0) ,$$

$$Cx_0 + h > 0, \quad i \ge 0, \quad v \ge 0, \quad i^{\dagger}v = 0$$
(5.4)

If  $Cx_0 + h + C(x^*-x_0) \ge 0$ , we take  $\lambda = 1$  and immediately find a solution

$$v(1) = Cx^* + h$$
,  $i(1) = 0$ 

Hence we assume that some components of  $C\left(\mathbf{x^{*}-x}_{0}\right)$  are negative such that

$$\exists : \begin{bmatrix} \lambda &= \min_{\mathbf{k}} \lambda \\ \mathbf{s} & \mathbf{s} \end{bmatrix} \begin{pmatrix} \mathbf{C}_{\mathbf{k}}, \mathbf{x}_{0} + \mathbf{h}_{\mathbf{k}} + \lambda_{\mathbf{k}} \mathbf{C}_{\mathbf{k}}, (\mathbf{x}^{*} - \mathbf{x}_{0}) = 0 \end{pmatrix} \land 0 < \lambda_{\mathbf{s}} < 1 \end{bmatrix} (5.5)$$

The situation that the index s in (5.5) is nonunique will be discussed later. Then from (5.5) we have

$$C_{k} \cdot x(\lambda_{s}) + h_{k} > 0, \quad k \neq s$$
  
 $C_{s} \cdot x(\lambda_{s}) + h_{s} = 0, \quad C_{s} \cdot (x^{*} - x_{0}) < 0$ 
(5.6)

Hence  $x = x(\lambda_s)$  is a point on the boundary hyperplane  $K_s^0$  of  $R_0$  given by

$$\kappa_{s}^{0} = \{ \mathbf{x} \mid C_{s}, \mathbf{x} + \mathbf{h}_{s} = 0 \}$$

Increasing  $\lambda$  above the value  $\lambda_{s}$  means that a new region  $R_{p}$  (p = 2<sup>n-S</sup>) is entered for which the natural hybrid representation can be found from (5.4) by pivoting on the diagonal element D<sub>SS</sub>. The representation in question then yields

$$w^{p} = D^{p}z^{p} + C^{p}x(\lambda) + h^{p},$$
with  $p = 2^{n-s}$ ,  $w^{p} \ge 0$ ,  $z^{p} \ge 0$  and  $w^{p}z^{p} = 0$ 

$$(5.7)$$

From the relation (4.23) we easily find

$$C_{k}^{p} \cdot x(\lambda) + h_{k}^{p} = C_{k} \cdot x(\lambda_{s}) + h_{k} + (\lambda - \lambda_{s}) (C_{k} - \frac{D_{ks}}{D_{ss}} C_{s}) (x^{*} - x_{0})$$
(5.8a)

for  $k \neq s$ 

and

$$C_{s}^{p} \mathbf{x}(\lambda) + h_{s}^{p} = -(\lambda - \lambda_{s}) \frac{C_{s} (\mathbf{x}^{*} - \mathbf{x}_{0})}{D_{ss}}$$
(5.8b)

As stated before, we assume  $\forall_{p,k} [D_{kk}^p > 0]$ . Then by (5.6) and (5.8) we have

$$c^{p}x(\lambda) + h^{p} > 0$$
 for  $\lambda > \lambda_{s}$  and  $(\lambda - \lambda_{s}) << 1$ .

Hence we may continue to increase  $\lambda$  until either some boundary hyperplane  $K_t^p$  of  $R_p$  is reached for  $\lambda = \lambda_t$ , indicated by

$$\lambda_{t} = \min_{k} (\lambda_{k} \mid C_{k}^{p} \cdot \mathbf{x}_{0} + h_{k}^{p} + \lambda_{k} C_{k}^{p} \cdot (\mathbf{x}^{*} - \mathbf{x}_{0}) = 0 \land \lambda_{k} > \lambda_{s}),$$

or  $\lambda_t$  does not exist and  $\lambda$  can be increased indefinitely. In the last case as well as for  $\lambda_t > 1$ , we take  $\lambda = 1$  and from (5.7) have found a solution  $w^P = C^P x(1) + h^P$ ,  $z^P = 0$ . If  $\lambda_t < 1$ we repeat the previous process by pivoting on  $D_{tt}^P$ . After each pivoting, we enter a new region in which the value of  $\lambda$  can again be increased. This guarantees that the algorithm never re-enters the same region (see the corresponding proof in [16]). Since the number of regions is finite, sooner or later a region R is entered for which  $x(1) = x^*$  is an internal point, yielding  $w^{q} = c^{q}x^* + h^{q}$ ,  $z^{q} = 0$  as a solution of (5.1).

Should the index s in (5.5) be nonunique, then at least two boundary hyperplanes are crossed at the same time, in which case the path  $x(\lambda)$  hits some corner of a region. In that case the initial point  $x_0$  can either be perturbated such that this situation no longer happens or other actions equivalent to those described in [16] can be taken, which will not be repeated here.

#### 5.2. The algorithm of Cottle

In [17] R.W. Cottle described a complementary pivoting algorithm for the solution of (5.1) when D  $\epsilon$  P. In our terminology of hybrid matrices and ideal diodes, this algorithm can be seen as a method to determine the solution of a multiport network loaded with ideal diodes by gradually connecting these diodes one by one to the multiport network. The algorithm to solve (5.1) a given x = x<sup>\*</sup> goes as follows.

Without loss of generality we assume that

$$v = Di + a, a = Cx^* + h$$
  
with  $a_j \ge 0$  for  $j \le k$  (5.9)  
 $a_j < 0$  for  $k < j \le n$ 

From the above representation with k nonnegative components in the vector a, we try to find a new representation by a sequence of one-step principal transformations (simplex pivoting steps) such that k + 1 components of a are nonnegative. When we are able to do so, the solution is obtained by repeated application of the same step until no negative components are left in the vector a.

It would be convenient if an exchange between  $v_{k+1}$  and  $i_{k+1}$  could do the job. Taking  $D_{k+1,k+1}$  as a pivot in any case guarantees that that the new  $a_{k+1}$  becomes positive since, after the pivoting has

taken place, its new value becomes  $-a_{k+1}/D_{k+1,k+1}$  and  $D_{k+1,k+1} > 0$ because  $D \in P$ . However, performing that pivoting step may produce negative values among the first k components of the new vector a, such that nothing is gained by such a pivoting step. In order to overcome this problem let us view upon the one-step principal transformation by the pivot  $D_{k+1,k+1}$  as a continuous proces instead of an immediate exchange between  $v_{k+1}$  and  $i_{k+1}$ . To this purpose we allow the variable  $i_{k+1}$  to grow continuously from zero to positive values, keeping all other components of i equal to zero. Furthermore, for the time being we focus our attention to the first k + 1 equations of (5.9) only. With  $i_j = 0$ for  $j \le k$  we may write these equations as

$$\begin{pmatrix} v_{1} \\ v_{2} \\ \vdots \\ \vdots \\ v_{k} \\ \cdots \\ v_{k+1} \end{pmatrix} = \begin{pmatrix} a_{1} \\ a_{2} \\ \vdots \\ a_{k} \\ \cdots \\ a_{k+1} \end{pmatrix} + i_{k+1} \begin{pmatrix} D_{1,k+1} \\ D_{2,k+1} \\ \vdots \\ D_{k,k+1} \\ \cdots \\ D_{k+1,k+1} \end{pmatrix}$$
(5.10)

In (5.10) every component  $v_j$  is given as a linear function of  $i_{k+1}$ . These functions  $v_j(i_{k+1})$  can be plotted as depicted in fig.5.1 for some values of j. For  $i_{k+1} = -a_{k+1}/D_{k+1,k+1} = d_{k+1}$  we have  $v_{k+1}(d_{k+1}) = 0$ . The corresponding values  $v_j(d_{k+1})$  for  $j \le k$  appear to become the components  $\tilde{a}_j$  of the new vector  $\tilde{a}$  which is obtained from (5.9) after exchanging  $v_{k+1}$  and  $i_{k+1}$ . In the case that  $v_j(d_{k+1}) \ge 0$  for  $j \le k$ , this pivoting step immediately leads to the desired k + 1 nonnegative components of  $\tilde{a}$ . It may however happen that one or more components  $v_j$  become negative before  $v_{k+1}$  reaches zero.



fig.5.1.

with 
$$d_j = -a_j/D_{j,k+1}$$
  
and  $d_m = \min_j (d_j \mid d_j \ge 0)$ , (5.11)

this happens when  $d_m < d_{k+1}$ . In this case we keep  $i_{k+1} = d_m$  and take  $D_{m,m}$  as a pivot, exchanging  $i_m$  and  $v_m$ . The new values  $v_j = v_j(d_m)$ ,  $j \neq m$  can be found from the intersections with the dotted axis in fig.5.1. It will be clear that for  $i_{k+1} = d_m$  we have

 $v_m = 0$ ,  $i_m = 0$  and  $v_j > 0$  for  $k \ge j \ne m$ Furthermore, due to  $D_{k+1,k+1} > 0$ ,  $v_{k+1}$  has become less negative at this point. The performed principal transformation results in the new equation system

 $w = \widetilde{D}z + \widetilde{a}, w \ge 0, z \ge 0, w^{t}z = 0$  (5.12) with

For  $z_j = 0$ ,  $j \le k$  the equations equivalent to (5.10) are now given by

$$\begin{pmatrix} \mathbf{w}_{1} \\ \mathbf{w}_{2} \\ \vdots \\ \mathbf{w}_{k} \\ \cdots \\ \mathbf{v}_{k+1} \end{pmatrix} = \begin{pmatrix} \widetilde{a}_{1} \\ \widetilde{a}_{1} \\ \vdots \\ \widetilde{a}_{k} \\ \cdots \\ \widetilde{a}_{k+1} \end{pmatrix} + \mathbf{i}_{k+1} \begin{pmatrix} \widetilde{D}_{1,k+1} \\ \widetilde{D}_{2,k+1} \\ \vdots \\ \widetilde{D}_{k,k+1} \\ \cdots \\ \widetilde{D}_{k+1,k+1} \end{pmatrix}$$
(5.13)

The various coefficients in (5.13) are easily obtained from (4.23) yielding

$$\widetilde{a}_{m}^{a} = -a_{m}^{\prime}/D_{m,m} \qquad \widetilde{D}_{m,k+1}^{a} = -D_{m,k+1}^{\prime}/D_{m,m}$$

$$\widetilde{a}_{j}^{a} = a_{j}^{a} - a_{m}^{a}D_{m,j}^{\prime}/D_{m,m}$$

$$\widetilde{D}_{j,1}^{a} = D_{j,1}^{a} - D_{j,m}^{a}D_{m,1}^{\prime}/D_{m,m}$$

$$(5.14)$$

From (5.14) and  $d_m \doteq -a_m / D_{m,m}^{\dagger}$  we find

$$v_{j}(d_{m}) = a_{j} + d_{m}D_{j,k+1} = \tilde{a}_{j} + d_{m}\tilde{D}_{j,k+1} = w_{j}(d_{m}), j = 1, 2, \dots, k+1$$

Hence the values of  $v_j$  and  $i_j$  as a function of  $i_{k+1}$  join continuously at  $i_{k+1} = d_m$  when a change in the basis between  $v_m$  and  $i_m$  is performed at this point. Therefore the picture of fig.5.1 may be extended to that given in fig.5.2.



fig.5.2.

From fig.5.2 we may see that the pivoting step with  $D_{m,m}$  may result into negative components of the vector  $\tilde{a}$  as determined by the dotted extensions of the new linear relations. However, this is no problem since the relevant values of  $v_j$  or  $i_j$ ,  $j \neq k$  are all positive immediately to the right of  $i_{k+1} = d_m$ , until a new intersection of one of these lines with the  $i_{k+1}$ -axis occurs. As before we should look for this possible intersection among the lines with negative slope, i.e. among the equations with negative entries in the k + 1-th column of  $\tilde{D}$ . We thus look for a  $\tilde{d}_1$  satisfying

$$\widetilde{d}_{1} = \min_{j} (\widetilde{d}_{j} = -\widetilde{a}_{j}/\widetilde{D}_{j,k+1} \mid \widetilde{d}_{j} > 0)$$
(5.15)

If such a  $\tilde{d}_1$  does not exist or if  $\tilde{d}_1 > \tilde{d}_{k+1}$ , we take  $\tilde{D}_{k+1,k+1}$  as a pivot and obtain a vector  $\bar{a}$  with at least k + 1 positive elements. However if  $\tilde{d}_1 < \tilde{d}_{k+1}$  we are again forced to perform a one-step principal transformation with  $\tilde{D}_{1,1}$  as pivot and repeat the above procedure on the new representation. This process ultimately leads to a basis in which  $v_{k+1}$  can be exchanged for  $i_{k+1}$ , because the piecewise-linear curve  $v_{k+1}(i_{k+1})$  arising in fig.5.2, from all pivoting steps (base transformations), always has a positive slope since in every representation (base)  $\widetilde{D}_{k+1,k+1} > 0$  when  $D \in P$ . Then this monotonically increasing function guarantees that no set of base variables occurs twice since the maximum value of  $v_{k+1}$  is uniquely determined by the particular base. Hence in the worst case the algorithm runs through all possible bases before a new representation is produced with at least k + 1 positive components in the corresponding vector a.

Should the vector a at the beginning of the algorithm satisfy a < 0 then any pivoting step produces at least one positive entry in the vector a. Hence the algorithm can always start with k = 1.

During the run of the algorithm it may happen that the index l in (5.15) is not unique, i.e. two or more variables w<sub>j</sub> become zero at the same time, which is in fact equivalent to the corner problem in the Katzenelson algorithm. The question is then which pivots should be taken in order to get all slopes of these variables positive for  $i_{k+1} > \tilde{d}_1$ . To this purpose let us assume without loss of generality that the m variables  $w_1, w_2, \ldots, w_m$  all become zero if  $i_{k+1} = \tilde{d}_1$  and consider the subsystem

$$\begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{pmatrix} = \begin{pmatrix} \widetilde{a}_1 \\ \widetilde{a}_2 \\ \vdots \\ \widetilde{a}_m \end{pmatrix} + i_{k+1} \begin{pmatrix} \widetilde{D}_{1,k+1} \\ \widetilde{D}_{2,k+1} \\ \vdots \\ \widetilde{D}_{m,k+1} \end{pmatrix} + \begin{pmatrix} \widetilde{D}_{11} & \widetilde{D}_{12} & \cdots & \widetilde{D}_{1m} \\ \widetilde{D}_{21} & \widetilde{D}_{22} & \cdots & \widetilde{D}_{2m} \\ \vdots & \vdots & \vdots \\ \widetilde{D}_{m1} & \widetilde{D}_{m2} & \cdots & \widetilde{D}_{mm} \end{pmatrix} \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{pmatrix},$$

(5.16)

With  $\tilde{a}_j + \tilde{d}_1 \tilde{p}_{j,k+1} = 0$  j = 1, 2, ..., m m < kFor  $i_{k+1} > \tilde{d}_1$  we then have

$$\begin{pmatrix} w_{1} \\ w_{2} \\ \vdots \\ w_{m} \end{pmatrix} = (i_{k+1} - \tilde{d}_{1}) \begin{pmatrix} \tilde{D}_{1,k+1} \\ \tilde{D}_{2,k+1} \\ \vdots \\ \tilde{D}_{m,k+1} \end{pmatrix} + \begin{pmatrix} \tilde{D}_{11} & \tilde{D}_{12} & \cdots & \tilde{D}_{1m} \\ \tilde{D}_{21} & \tilde{D}_{22} & \cdots & \tilde{D}_{2m} \\ \vdots & \vdots & & \vdots \\ \tilde{D}_{m1} & \tilde{D}_{m2} & \cdots & \tilde{D}_{mm} \end{pmatrix} \begin{pmatrix} z_{1} \\ z_{2} \\ \vdots \\ z_{m} \end{pmatrix}$$
(5.17)

In order for all slopes of  $w_j(i_{k+1})$  to become positive we have to pivot the system (5.17) such a way that the vector  $(\tilde{D}_{1,k+1},\ldots\tilde{D}_{m,k+1})^{t}$  becomes positive. However this problem is exactly of the same type as the problem we are currently solving but of smaller dimension. Hence the same Cottle algorithm can be applied to make all slopes nonnegative. Then for  $i_{k+1} = \tilde{d}_1$  all other variables  $w_j$  with j > m,  $j \neq k + 1$  have remained unchanged after the pivoting performed on (5.17) since  $w_j(\tilde{d}_1) = 0$  for  $j \le m$ .

In a practical application of the above algorithm it is of advantage to check after each pivoting step whether by accident the current number q of nonnegative components of the vector a already exceeds k. In such a case the algorithm becomes more efficient when restarted on the last obtained representation with q nonnegative components. Furthermore in selecting the negative component of the vector a which has to be made positive it is advisable to choose that one for which the corresponding  $d_{k+1}$  is as small as possible. This reduces the chance that other pivots have to be taken before  $v_{k+1}$  and  $i_{k+1}$  can be exchanged and benefits the numerical stability of the algorithm.

#### 5.3. The Lemke algorithm

The most familiar algorithm to solve the LCP problem is the well known Lemke algorithm [7]. A review of the matrix classes on which the algorithm can be applied successfully is given by Karamardian [18] and includes class P matrices. For a given x and

$$w = Dz + a, \quad a = Cx + h$$
  
 $w \ge 0, \quad z \ge 0, \quad w^{t}z = 0, \quad D = (n \times n), \quad C = (n \times m),$ 
(5.18)

the algorithm runs as follows. First of all, the vector a in (5.18) is extended by a multiple  $z_0$  of some positive vector e, yielding

$$w = Dz + a + z_0 e, e > 0$$
 (5.19)

Next by defining

$$\begin{split} \mathbf{w} &= \begin{pmatrix} \mathbf{w}_0 \\ \cdots \\ \mathbf{w} \end{pmatrix} , \quad \mathbf{z} &= \begin{pmatrix} \mathbf{z}_0 \\ \cdots \\ \mathbf{z} \end{pmatrix} , \quad \mathbf{M} &= \begin{pmatrix} 0 : 0 \cdots 0 \\ \vdots \\ \mathbf{e} : \mathbf{D} \end{pmatrix} \quad \text{and} \quad \mathbf{a} &= \begin{pmatrix} 0 \\ \cdots \\ \mathbf{a} \end{pmatrix} \end{split}$$

We consider the problem

w = a + Mz (5.20a)

$$\mathbf{w} \ge \mathbf{0} \qquad \mathbf{z} \ge \mathbf{0} \tag{5.20b}$$

$$\mathbf{w}^{\mathsf{L}}\mathbf{z} = 0 \tag{5.20c}$$

Obviously, if after some pivoting steps a solution of (5.20) is obtained for which  $z_j = z_0 = 0$ , then the components w and z of that solution also satisfy (5.18).

Before continuing our description we make the following definitions. In any pivoted equivalent relation (5.20a), the variables in  $\underline{w}$  will be called basic variables and those in  $\underline{z}$  are nonbasic variables. Furthermore a pair  $\underline{w}, \underline{z}$  which satisfies (5.20a) and (5.20b) will be called *feasible* and a pair  $\underline{w}, \underline{z}$  for which  $\underline{w}^{t}\underline{z} = 0$  and  $\underline{z}$  contains exactly one pair of components  $\underline{w}_{j}, \underline{z}_{j}$  will be called *almost complementary*.

Let us now consider a value  $\boldsymbol{z}_{\cap}$  such that

$$z_0^* = \min_j \left(-\frac{a_j}{e_j} \mid a_j < 0\right)$$
 (5.21)

and let this minimum be obtained for j = k. Then from (5.19) and (5.20) we obtain for z = 0

$$w_{j} = a_{j} + z_{0}^{*}e_{j} > 0 \quad j \neq k$$
  
 $w_{k} = 0$ 
(5.22)

Next the system (5.20) is pivoted on the element  $M_{k,1}$  yielding

$$w' = a' + M'z'$$
, (5.23)

with the basic variables  $w^\prime$  and the nonbasic variables  $z^\prime$  given by

$$\underline{w}' = (w_0, w_1, \dots, w_{k-1}, z_0, w_{k+1}, \dots, w_n)^{t}$$

$$\underline{z}' = (w_k, z_1, \dots, z_{k-1}, z_k, z_{k+1}, \dots, z_n)^{t}$$
(5.24)

From (5.22) to (5.24) and  $z_0 = z_0^*$  we then have that an almost complementary feasible solution of (5.23) is given by

$$w' = a', z' = 0, a' > 0$$

since

$$a'_{j} = a_{j} + z^{*}_{0}e_{j} > 0 \qquad j \neq k \text{ and } a'_{k} = z^{*}_{0} > 0 \qquad (5.25)$$

From (5.24) we observe that the variable  $w_k$  has become nonbasic. We then start to make its complement  $z_k$  positive until some component  $w'_r = a'_r + M'_{rk}z_k = 0$ . Then the system (5.23) is pivoted on  $M'_{rk}$  yielding

$$\underline{w}^{"} = \underline{a}^{"} + \underline{M}^{"}\underline{z}^{"}$$

with 
$$\underline{w}^{"} = (w_{0}, w_{1}, \dots, w_{r-1}, z_{k}, w_{r+1}, \dots, z_{0}, \dots)^{t}$$
  
 $\underline{z}^{"} = (w_{k}, z_{1}, \dots, z_{r-1}, w_{r}, z_{r+1}, \dots, z_{r}, \dots)^{t}$ 
(5.26)

As before it is easy to show that  $\underline{a}^{"} > 0$  and  $\underline{w}^{"} = \underline{a}^{"}$ ,  $\underline{z}^{"} = 0$  is then an almost complementary solution of (5.20) since  $\underline{z}^{"}$  contains  $w_{r}$  as well as  $z_{r}$ . Again the complement  $z_{r}$  of the variable  $w_{r}$  which became nonbasic is increased until some other component of  $\underline{w}^{"}$ becomes zero, which is then exchanged with  $z_r$ . This process is repeated until either  $z_0$  becomes again nonbasic or no component of  $\underline{w}$  can be made zero since the column of M corresponding the the variable to be increased, only contains nonnegative elements.

In the first case we have found a complementary feasible solution of (5.20) for which the elements  $w_{i}, z_{i}$  satisfy (5.18) and thus a solution is found.

In the second case the algorithm ends in a so called almost complementary ray. Lemke has shown that the latter case never arises when  $D \in P$  and also showed that no cycling can occur in this situation.

A disadvantage of the algorithm is that the pivots to be taken are not principal pivots which makes the implementation of the algorithm more difficult, especially when sparse matrix techniques are to be used.

## 6. The modulus algorithm

The complementary pivoting algorithms as described in chapter 5 all work reasonable efficient in practice when applied to those classes of matrices for which they are known to produce a solution. As to the complexity of these methods it is known that the Katzenelson-Cottle- and the Lemke-algorithm are nonpolynomial in time-complexity [19], [20]. Yet as for the simplex method in linear programming, the average complexity is polynomial.

In this section we will describe an algorithm which is based on contraction methods and has the property that for a specific class of matrices with a specified bounded condition number the time-complexity is polynomial in the dimension of the problem. At the same time, the formulation of the LCP problem in this particular algorithm can be used to present an alternate state-model description for which the global model as presented by Chua [21] is included as a special case.

## 6.1. The modulus transformation

Consider the LCP problem

$$\mathbf{v} = \mathbf{M}\mathbf{i} + \mathbf{a}$$
  $\mathbf{M} = (\mathbf{m} \times \mathbf{m})$   
 $\mathbf{v} \ge 0$ ,  $\mathbf{i} \ge 0$ ,  $\mathbf{v}^{\dagger}\mathbf{i} = 0$  (6.1)

We will transform (6.1) into another equivalent problem by the so called modulus transformation. We will define this modulus transformation on the general nonlinear complementary problem:

Given v,i  $\epsilon \mathbb{R}^{m}$ , with v = v(i). Determine the points v,i satisfying  $v \ge 0$ ,  $i \ge 0$ ,  $v^{t}i = 0$ . (6.2)

Then we apply this transformation on (6.1) as a special case. For this purpose, with  $z \in \mathbb{R}^m$  we define

$$|\mathbf{z}| = (|\mathbf{z}_1|, |\mathbf{z}_2|, \dots, |\mathbf{z}_m|)^{\mathsf{T}}$$
 (6.3)
Furthermore, for a scalar function  $\varphi : \mathbb{R} \to \mathbb{R}$  and a vector  $\mathbf{x} \in \mathbb{R}^{m}$  we define  $\varphi[\mathbf{x}]$  by

 $\varphi[\mathbf{x}] = (\varphi(\mathbf{x}_1), \varphi(\mathbf{x}_2), \dots, \varphi(\mathbf{x}_m))^{\mathsf{t}}$ 

With the above definitions we have the following theorem.

- <u>Theorem 7</u>: For any strictly monotone  $\theta(t)$  :  $\mathbb{R}_{+} \rightarrow \mathbb{R}_{+}$ , satisfying  $\theta(0) = 0$ , the nonlinear complementarity problem is completely equivalent to the solution of a set of m nonlinear equations in  $\mathbb{R}^{m}$  given by  $|\theta[\mathbf{v}] - \theta[\mathbf{i}]| = \theta[\mathbf{v}] + \theta[\mathbf{i}]$
- Proof: Let us denote the property "i,v is a solution of (6.2)"
  by CP. We will first show that

$$CP \leftrightarrow \exists [i = h[|z|+z] \land v = h[|z|-z]]$$
(6.4)

with a strictly monotone increasing  $h(t) : \mathbb{R}_+ \to \mathbb{R}_+$ and h(0) = 0. For this purpose we assume that we have a vector i satisfying (6.2). For this vector i we define the vectors z and w according to

$$z = \frac{1}{2} \left( \theta[i] - \theta[v] \right), \quad w = \frac{1}{2} \left( \theta[i] + \theta[v] \right)$$
(6.5)

and  $\theta$  as given in the theorem. Due to the definition of  $\theta(t)$ , its inverse on  $\mathbb{R}_+$  exists and is denoted by h(t). Then from (6.5) we obtain for a pair i,v satisfying (6.2)

$$\begin{aligned} \mathbf{z}_{\mathbf{k}} &= \frac{1}{2} \, \theta(\mathbf{i}_{\mathbf{k}}), \quad \mathbf{w}_{\mathbf{k}} &= \frac{1}{2} \, \theta(\mathbf{i}_{\mathbf{k}}) \quad \text{if} \quad \mathbf{v}_{\mathbf{k}} = 0, \quad \mathbf{i}_{\mathbf{k}} \geq 0 \\ \mathbf{z}_{\mathbf{k}} &= -\frac{1}{2} \, \theta(\mathbf{v}_{\mathbf{k}}), \quad \mathbf{w}_{\mathbf{k}} = \frac{1}{2} \, \theta(\mathbf{v}_{\mathbf{k}}) \quad \text{if} \quad \mathbf{i}_{\mathbf{k}} = 0, \quad \mathbf{v}_{\mathbf{k}} \geq 0 \end{aligned}$$
(6.6)

Hence from (6.6) we obtain

$$CP \rightarrow w = |z|$$
(6.7)

Subsitution of (6.7) in (6.5) and eliminating  $\theta$ [i] and

$$\theta[\mathbf{i}] = |\mathbf{z}| + \mathbf{z}$$

$$\theta[\mathbf{v}] = |\mathbf{z}| - \mathbf{z}$$
(6.8)

With the previous h(t), (6.7) and (6.8) result into

$$CP \rightarrow \exists_{z}[i = h[|z|+z] \land v = h[|z|-z]],$$

with the particular z given by (6.5). On the other hand, if the r.h.s of (6.4) holds, it is immediately obvious that  $i \ge 0$ ,  $v \ge 0$  and  $i^{t}v = 0$ , due to the definition of h(t). This concludes the proof of (6.4).

Next, from the r.h.s of (6.4) we eliminate |z| and z yielding

$$\begin{array}{l}
\exists \quad [\mathbf{i} = \mathbf{h}[|\mathbf{z}|+\mathbf{z}] \land \mathbf{v} = \mathbf{h}[|\mathbf{z}|-\mathbf{z}] \leftrightarrow \\
\mathbf{z} \in \mathbb{R}^{m} \\
|\mathbf{z}| = \frac{1}{2} (\theta[\mathbf{i}] + \theta[\mathbf{v}]), \\
\mathbf{z} = \frac{1}{2} (\theta[\mathbf{i}] - \theta[\mathbf{v}]) \\
\end{array}$$
(6.9)

Eliminating z from the r.h.s of (6.9) then yields with (6.4)

$$CP \leftrightarrow |\theta[i] - \theta[v]| = \theta[i] + \theta[v],$$

which concludes the proof of theorem 7.

The result of theorem 7 is closely related to a theorem of Mangasarian [22] who proved that for the same conditions on  $\theta(t)$  :  $\mathbb{R} \to \mathbb{R}$ 

 $CP \leftrightarrow \theta[|v-i|] = \theta[v] + \theta[i]$ 

From the proof of theorem 7 we also derive the following theorem. <u>Theorem 8</u>: If  $h : \mathbf{R}_{+} \neq \mathbf{R}_{+}$  is a strictly increasing function with h(0) = 0, then i := h[|z| + z] v := h[|z| - z]is a solution of the nonlinear complementarity problem iff z is a solution of the system of m nonlinear equations For the special case that h(t): = t, the equations in theorem 8 are a different notation for the equivalent formulation of the nonlinear complementarity problem given by Megiddo in terms of the so called extension [23, p.113].

Definition: The transformation i,v defined by i: = 
$$h[|z| + z]$$
,  
v: =  $h[|z| - z]$  with a strictly increasing  $h : \mathbb{R}_+ \to \mathbb{R}_+$   
and  $h(0) = 0$ , is called *the modulus transformation*.

By theorem 8, application of the modulus transformation to a function v = v(i) automatically guarantees that  $i \ge 0$ ,  $v \ge 0$  and  $i^{t}v = 0$ .

<u>Theorem 9</u>: The linear complementarity problem (6.1) with det(I+M)  $\neq$  0 is equivalent to the determination of a vector  $z \in \mathbb{R}^m$ which satisfies

 $z = (I+M)^{-1} (I-M) |z| - (I+M)^{-1} a.$ 

The solutions  $z^{j}$  of this equation determine the solutions  $v^{j}$ ,  $i^{j}$  of the LCP problem by  $v^{j} = |z^{j}| - z^{j}$ ,  $i^{j} = |z^{j}| + z^{j}$ .

Proof: With h(t): = t we apply the modulus transformation
on v = Mi + a, yielding

|z| - z = M(|z|+z) + a, with v = |z| - z, i = |z| + z.

Rearranging terms then gives

$$(I+M)z = (I-M)|z| - a$$
 (6.10)

By a suitable scaling of i in (6.1) we can always arrange that I + M is nonsingular. Hence from (6.10) we find

$$z = (I+M)^{-1}(I-M)|z| - (I+M)^{-1}a$$
, (6.11)  
with  $v = |z| - z$ ,  $i = |z| + z$ .  
The theorem then follows from theorem 8.

In the next section we will present an algorithm to solve the nonlinear equation (6.11) by some contraction-mapping algorithm. Due to theorem 9, we do not have to bother about the complementarity conditions for v and i and in fact we now have the large field of algorithms to solve nonlinear equations at our disposal for solving the LCP problems as well.

# 6.2. The modulus algorithm as a contraction mapping

In this section we will show that (6.11) can be solved by contraction mapping for certain conditions on the matrix M.

Let us rewrite (6.11) in the form

$$z = b + D|z|, \qquad (6.12)$$

(6.13)

with  $b = -(I+M)^{-1}a$  and  $D = (I+M)^{-1}(I-M)$ . In the sequel we will use an absolute norm over  $\mathbb{R}^{M}$  satisfying

a)  $||\mathbf{x}|| \ge 0$ ,  $||\mathbf{x}|| = 0$  iff  $\mathbf{x} = 0$ b)  $||\alpha\mathbf{x}|| = |\alpha| \cdot ||\mathbf{x}||$ c)  $||\mathbf{a} + \mathbf{b}|| \le ||\mathbf{a}|| + ||\mathbf{b}||$ d)  $|||\mathbf{a}||| = ||\mathbf{a}||$ e)  $0 \le \mathbf{a} \le \mathbf{b} + ||\mathbf{a}|| \le ||\mathbf{b}||$ 

As an example the standard  $\|\cdot\|_{\infty}$  and  $\|\cdot\|_{2}$  norm are absolute norms. Furthermore, the norm  $\|A\|$  of a matrix A is defined to be the induced norm

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

Lemma 3:  $||a| - |b|| \le |a-b|$ ,  $a, b \in \mathbb{R}^{m}$ , the inequality holding componentwise.

<u>Proof</u>: We partition a and b in four parts such that any possible sign combination of the components of a and b is represented. With  $a_i, b_i \ge 0$  this is accomplished by

$$a = (a_1, a_2, -a_3, -a_4)^{t}$$
 and  $b = (b_1, -b_2, b_3, -b_4)^{t}$ 

Then we have

$$p = ||a| - |b|| = (|a_1 - b_1|, |a_2 - b_2|, |a_3 - b_3|, |a_4 - b_4|)^{t}$$
(6.14)

anđ

$$q = |a-b| = (|a_1-b_1|, |a_2+b_2|, |a_3+b_3|, |a_4-b_4|)^{t}$$
  
(6.15)

From (6.14) and (6.15) we see that every component of p is less or equal to the corresponding component of q and thus  $0 \le p \le q$ , which proves the lemma.

Theorem 10: If ||D|| < 1, the iterative process  $z^{n+1} = b + D|z^n|$ converges and the limit  $z^*$  of the sequence  $\{z^n\}$  for  $n \to \infty$  is the unique solution of z = b + D|z|.

Proof: With f(z) = b + D|z|, equation (6.12) can be written as

$$z = f(z) \tag{6.16}$$

Now consider the function f(z) for which we derive

$$\|\mathbf{f}(z^{1}) - \mathbf{f}(z^{2})\| = \|D(|z^{1}| - |z^{2}|)\| \le \|D\| \cdot \||z^{1}| - |z^{2}|\|$$
(6.17)

By lemma 3 and the norm properties d) and e) we have  $|||z^{1}| - |z^{2}||| = ||||z^{1}| - |z^{2}|||| \le |||z^{1} - z^{2}|| = ||z^{1} - z^{2}||$ 

Thus (6.17) yields

 $||f(z^{1})-f(z^{2})|| \leq ||D|| \cdot ||z^{1}-z^{2}||$ ,

which means that f is Lipschitz bounded for all  $z^1, z^2 \in \mathbb{R}^m$  since ||D|| does not depend on  $z^1$  or  $z^2$ . Then if ||D|| < 1, the iteration  $z^{n+1} = f(z^n)$  converges by the Banach contraction mapping theorem [24] and  $z^* = \lim_{n \to \infty} z^n$  is the unique solution of z = f(z),

which proves theorem 10.

- Definition: The iteration  $z^{n+1} := b + D|z^n|$  will be called the modulus algorithm.
- <u>Theorem 11</u>: If  $||(I+M)^{-1}(I-M)|| < 1$ , the modulus algorithm  $z^{n+1} = (I+M)^{-1}(I-M)|z^n| - (I+M)^{-1}a$  yields the unique solution of the LCP problem
  - $\begin{aligned} \mathbf{v} &= \mathbf{Mi} + \mathbf{a}, \quad \mathbf{v} \geq \mathbf{0}, \quad \mathbf{i} \geq \mathbf{0}, \quad \mathbf{v}^{\mathsf{t}} \mathbf{i} = \mathbf{0} \quad \text{via} \\ \mathbf{z}^{\mathsf{t}} &= \lim_{n \to \infty} \mathbf{z}^{n} \text{ and } \mathbf{v} = \left| \mathbf{z}^{\mathsf{t}} \right| \mathbf{z}^{\mathsf{t}}, \quad \mathbf{i} = \left| \mathbf{z}^{\mathsf{t}} \right| + \mathbf{z}^{\mathsf{t}}. \end{aligned}$

Proof: The proof is obvious from theorem 9 and theorem 10.

#### 6.3. A polynomial algorithm

In this section we will bound the number of operations of the modulus algorithm to obtain an exact solution of the LCP problem in the case of symmetric positive definite matrices M.

In the sequel  $||\cdot||$  will denote the standard  $||\cdot||_2$  norm unless stated otherwise.

#### Definition: The radius $\theta(A)$ of a matrix A is defined by

$$\theta(A) := || (I+A)^{-1} (I-A) ||_{2}$$

Next we present the main theorem of this section.

- <u>Theorem 12</u>: Given a symmetric positive definite matrix  $M = (m \times m)$ with a bounded condition number  $K = ||M|| \cdot ||M^{-1}||$ . Then the LCP problem v = Mi + a,  $v \ge 0$ ,  $i \ge 0$ ,  $v^{t}i = 0$  is solvable in polynomial time by the modulus algorithm and the number N of multiplications or additions is asymptotically bounded by  $N \sim \frac{1}{7} K^{1/2} m^{7/2}$ .
- <u>Proof</u>: Since M is symmetric and positive definite, M has positive real eigenvalues, i.e.  $\lambda_i(M) > 0$ . The eigenvalues  $\mu_i$  of D = (I+M)<sup>-1</sup>(I-M) are given by

$$\begin{split} \mu_{i}(D) &= \frac{1-\lambda_{i}}{1+\lambda_{i}} \text{ and hence the } \mu_{i} \text{ are real and satisfy} \\ |\mu_{i}| < 1 \text{ since } \lambda_{i} > 0. \text{ Let } \rho(A) \text{ denote the maximum} \\ \text{modulus of the eigenvalues of A. Then since D is also} \\ \text{symmetric, we have } \theta &= \theta(M) = ||D||_{2} = \sqrt{\rho(D^{t}D)} = \\ &= \sqrt{\rho(D)^{2}} = \rho(D) = \max |\mu_{i}| < 1 \text{ and by theorem 10} \\ \text{ the iteration } z^{n+1} = b + D|z^{n}| \text{ converges to a unique} \\ \text{point } z^{*}. \text{ The contraction mapping theorem then yields} \end{split}$$

$$||z^{*}-z^{n}|| \leq \frac{\theta^{n}}{1-\theta} ||z^{1}-z^{0}||,$$
 (6.18)

with z\* satisfying

$$z^* = b + D|z^*|$$
 (6.19)

From theorem 10 we know that the solution  $z^*$  is unique. From (6.19) we find

$$||z^{*}|| = ||b + D|z^{*}|| \ge ||b|| - ||D|z^{*}|| \ge ||b|| - ||D|| \cdot ||z^{*}||$$

Thus

$$||z^{\star}|| \geq \frac{||\mathbf{b}||}{1+||\mathbf{D}||} = \frac{||\mathbf{b}||}{1+\theta} = \gamma ||\mathbf{b}||$$
 (6.20)

If we start the iteration with  $z^0 = 0$  we obtain  $z^1 = b$  and thus by (6.18)

$$||z^{*}-z^{n}|| \leq \frac{\theta^{n}}{1-\theta} ||b|| = \delta(n) ||b||$$
 (6.21)

Then  $||z^{\star}|| - ||z^{n}|| \le ||z^{\star}-z^{n}|| \le \delta(n) ||b||$ , or using (6.20)

$$||\mathbf{z}^{\mathbf{n}}|| \geq ||\mathbf{z}^{\star}|| - \delta ||\mathbf{b}|| \geq (\gamma - \delta) ||\mathbf{b}|| = \varepsilon(\mathbf{n}) ||\mathbf{b}||$$
(6.22)

From  $||z^{n}||_{2} \ge \varepsilon(n) ||b||_{2}$  , we have

$$\mathbf{J}_{i}\left[|\mathbf{z}_{i}^{n}| \geq \frac{\varepsilon(n) ||\mathbf{b}||}{\sqrt{m}}\right]$$
(6.23)

On the other hand, from (6.21) we find

$$|z_{i}^{\star} - z_{i}^{n}| \leq \delta(n) ||b||$$
 (6.24)

As soon as  $\frac{\varepsilon(n)}{\sqrt{m}} > \delta(n)$ , we see from (6.23) and (6.24) that the component  $z_i^*$  of  $z^*$  must have the same sign as the component  $z_i^n$  of  $z^n$ . This happens as soon as

$$\frac{1}{1+\theta} > (1+\sqrt{m}) \frac{\theta^{n}}{1-\theta} , \text{ or}$$

$$n > N_{0} = \left[ \frac{\ln \frac{1-\theta_{0}}{1+\theta_{0}} - \ln (1+\sqrt{m})}{\ln \theta_{0}} \right] , \theta_{0} = \theta(M) < 1 \quad (6.25)$$

Hence after N<sub>0</sub> iterations, the sign of at least one component of  $z^*$  is known. Then from (6.11) we also know whether  $v_k$  or  $i_k$  in the complementary solution has to be zero.

In the case that  $i_k = 0$ , we delete the k-th equation from (6.1) as well as the k-th column of M and are left with an identical problem of dimension m - 1 with a matrix  $\overline{M} = (m-1) \times (m-1)$ . Should  $v_k = 0$  we exchange  $v_k$  and  $i_k$  in (6.1) and again reduce the problem to dimension m - 1with a matrix  $\widetilde{M}$ .

This phase concludes cycle number zero of the algorithm. Now we apply the same process in each following cycle. In appendix II it is shown that  $\theta(M_{j}) < 1$  for all matrices  $M_{j}$ that arise this way in every subsequent cycle of the algorithm. Hence the process results in the signs of all components of  $z^*$  after running through m such cycles. In the case that  $v_k + i_k = 0$  for some k, we get a reduced system sooner or later for which the b column is zero,

75

leading to z = 0 and hence  $v_k = i_k = 0$ . From (6.25) the number of iterations in cycle k is limited by

$$N_{\mathbf{k}} \leq \left[ \frac{\ln \frac{1-\theta_{\mathbf{k}}}{1+\theta_{\mathbf{k}}} - \ln (1+\sqrt{m-\mathbf{k}})}{\frac{\ln \theta_{\mathbf{k}}}{2}} \right] , \quad \mathbf{k}=0,1,2,\ldots,m-1$$

With 
$$K_1 = \frac{\ln \frac{1-\overline{\theta}}{1+\overline{\theta}}}{\ln \overline{\theta}} > 0$$
,  $K_2 = \frac{-1}{\ln \overline{\theta}}$  and  $\overline{\theta} = \max_i \theta_i$ 

this yields

$$N_{k} \leq K_{1} + K_{2}\sqrt{m-k}$$
(6.26)

Of course, the reduced problem in the k+1-th cycle is most efficiently obtained by eliminating the component of  $z^*$ , whose sign was determined in the k-th cycle, from the equations in the k-th cycle. This amounts in total for at most  $O(m^3)$  multiplications (long operations). The number of long operations in the required contractions equals

$$N_{C} = \sum_{i=0}^{m-1} N_{i}(m-i)^{2} \le \sum_{i=0}^{m-1} (K_{1} + K_{2}\sqrt{m-i}) (m-i)^{2}$$

By using the Euler-Maclaurin summation formula we find

$$N_{\rm C} \le \frac{2}{7} \kappa_2^{-m} + 1 \text{ lower order terms}$$
 (6.27)

In appendix II it is derived that  $K_2 \leq \frac{1}{2} (||M|| \cdot ||M^{-1}||)^{1/2} = \frac{1}{2} \kappa^{1/2} \text{ can be realized by}$ a suitable scaling of M, yielding  $N \sim \frac{1}{2} \kappa^{1/2} m^{7/2}$ 

This concludes the proof of theorem 12.

From our previous considerations about networks containing uncontrolled resistive elements with strictly monotone increasing characteristics, it will be clear that the above algorithm will yields a polynomial-time solution process since the hybrid matrix of a purely resistive network (e.g. the impedance matrix) is symmetric and positive definite.

## 6.4. Implementation aspects and variants of the modulus algorithm

In a lot of practical situations, the matrix M in (6.1) will be sparse and the iteration  $z^{n+1} = (I+M)^{-1}(I-M)|z^n| - (I+M)^{-1}a$ is not efficient since  $(I+M)^{-1}(I-M)$  may become non-sparse. To be able to use sparse matrix techniques it is better to use the iteration in the equivalent formulation of

$$(I+M)(z^{n+1}+|z^n|) = 2|z^n| + a$$
 (6.28)

If M is sparse, then I + M is also sparse and an L\U factorization of (I+M) enables us to solve (6.28) by sparse matrix techniques with a minimum amount of work. Furthermore it will be convenient during the iteration to monitor the residue vector  $e^n$  given by

$$e^{n} = Mi^{n} + a - v^{n} \tag{6.29}$$

with  $\mathbf{i}^{n} = |\mathbf{z}^{n}| + \mathbf{z}^{n}$ ,  $\mathbf{v}^{n} = |\mathbf{z}^{n}| - \mathbf{z}^{n}$ . From (6.29) we then find  $\mathbf{e}^{n} = \mathbf{M}(|\mathbf{z}^{n}|+\mathbf{z}^{n}) + \mathbf{a} - |\mathbf{z}^{n}| + \mathbf{z}^{n} = (\mathbf{M}+\mathbf{I})\mathbf{z}^{n} - (\mathbf{I}-\mathbf{M})|\mathbf{z}^{n}| + \mathbf{a}$ , or  $(\mathbf{I}+\mathbf{M})^{-1}\mathbf{e}^{n} = \mathbf{z}^{n} - (\mathbf{I}+\mathbf{M})^{-1}(\mathbf{I}-\mathbf{M})|\mathbf{z}^{n}| + (\mathbf{I}+\mathbf{M})^{-1}\mathbf{a} = \mathbf{z}^{n} - \mathbf{z}^{n+1}$ 

Hence

$$e^{n} = (I+M)(z^{n}-z^{n+1})$$
 (6.30)

From (6.30) we obtain

$$||e^{n}|| \leq ||I+M|| ||z^{n}-z+z-z^{n+1}|| \leq ||I+M|| \left(\frac{\theta^{n}}{1-\theta} + \frac{\theta^{n+1}}{1-\theta}\right) ||b||, \quad (6.31)$$
  
with  $b = (I+M)^{-1}a$ .

Then, with  $\theta < 1$ , (6.31) results into a relative error

$$\frac{||e^{n}||}{||a||} \leq 2||I+M|| \cdot ||(I+M)^{-1}|| \frac{\theta^{n}}{1-\theta}$$
(6.32)

For a prescribed accuracy, the required number of iterations can now be found from (6.32).

As to the contraction mapping iteration  $z^{n+1} = D|z^n| + b$  to solve z = D|z| + b, it will be clear that refinements to this brute force attack should be possible which can increase the convergence speed of this algorithm, depending on the properties of D. From the derivations in appendix II, we may obtain that the eigenvalues of each principal submatrix of D are less than one in modulus if M is symmetric and positive definite. Hence all diagonal elements of D then satisfy  $|D_{kk}| < 1$ . For that case it may be interesting to consider the iteration

$$z^{n+1} - L|z^{n+1}| = U|z^{n}| + b$$
, (6.33)

with L + U = D and U<sub>ij</sub> = 0 for  $i \le j$ . Then the left hand side of (6.33) can be uniquely solved for  $z^{n+1}$  by going from top to bottom through this "triangular" system since each equation is of the form

$$z_{k}^{n+1} \sim L_{kk}|z_{k}^{n+1}| = \sum_{j=k+1}^{m} U_{kj}|z_{j}^{n}| + \sum_{j=1}^{k-1} L_{kj}|z_{j}^{n+1}| = k=1,2,\ldots,m$$

with  $|L_{kk}| < 1$ .

Practical experiments indeed showed a speed-up of this Gauss-Seidellike iteration. A different approach to the solution of (6.1) is found by applying the modulus transformation with  $h(t) = t^2$ . We then have

$$i = h[|z|+z] = (|z|+z)^{2} = 2z^{2} + 2z|z| ,$$

$$v = h[|z|-z] = (|z|-z)^{2} = 2z^{2} - 2z|z| ,$$
(6.34)

with the multiplications in  $z^2$  and z|z| understood to be taken componentwise. Substitution of (6.34) in (6.1) then yields

$$z|z| = (I+M)^{-1} (I-M) z^{2} - \frac{1}{2} (I+M)^{-1} a$$
 (6.35)

From (6.35) we immediately observe that this nonlinear set of equations, written in the form F(z) = 0 is continuously differentiable for all  $z \in \mathbb{R}^m$  since the scalar functions t|t| and  $t^2$  both have this property. Hence we can apply Newton-Raphson iteration on (6.35) to find the solution  $z^*$  and benefit from the quadratic convergence. A disadvantage is of course that some initial guess, sufficiently close to the final solution, must be given.

## 6.5. Relations with the global model of Chua

The modulus transformation as defined in section 6.1 can of course also be used in the state-model to describe some PL mapping  $f : x \rightarrow y$ . To this purpose we apply the transformation i = |z| + z, v = |z| - zto the equations

y = Ax + Bi + g  $v = Cx + Di + h , v \ge 0, i \ge 0, v^{t}i = 0$ (6.36)

yielding

$$y = \overline{A}x + \overline{B}|z| + \overline{g}$$
(6.37a)  
$$z = \overline{C}x + \overline{D}|z| + \overline{h}$$
(6.37b)

with

$$\vec{A} = A - B(I+D)^{-1}C, \quad \vec{B} = 2B(I+D)^{-1}$$

$$\vec{C} = -(I+D)^{-1}C, \quad \vec{D} = (I+D)^{-1}(I-D) \quad (6.38)$$

$$\vec{h} = -(I+D)^{-1}h, \quad \vec{g} = g - B(I+D)^{-1}h$$

Due to our previous derivations, the equation system (6.37) is completely equivalent to the state-model (6.36). As an additional advantage no extra complementarity conditions are necessary for the state-equations (6.37b). A specific state in (6.36) is now equivalent to a specific sign distribution over the components  $z_i$  of z in (6.37b). As with the state-model (6.36), there may exist  $2^n$  different equivalent descriptions of the form (6.37) which also may be classified by a state vector  $s_p$  corresponding to

$$z_{k} > 0 \rightarrow s_{p_{k}} = 1$$
$$z_{k} \le 0 \rightarrow s_{p_{k}} = 0$$

Now consider the special situation that the matrix D in (6.36) is equal to the unit matrix I. The state equations in (6.36) then become

$$v = Cx + h + i$$
,  $C = (n \times m)$ 

from which we obtain for each component

$$v_k = (Cx+h)_k + i_k = C_k \cdot x + h_k + i_k, k=1,2,...,n$$
 (6.39)

Equation (6.39) then yields

$$C_{k} \cdot x + h_{k} \ge 0 \Rightarrow v_{k} = C_{k} \cdot x + h_{k}, \quad i_{k} = 0$$

$$C_{k} \cdot x + h_{k} < 0 \Rightarrow v_{k} = 0, \quad , \quad i_{k} = -(C_{k} \cdot x + h_{k})$$
(6.40)

From (6.40) we see that the domain  $\mathbb{R}^m$  is divided into a number of halfspaces by hyperplanes  $C_k \cdot x + h_k = 0$ , which extend through whole  $\mathbb{R}^m$  (see fig.6.1).



fig.6.1.

In each region, we only have to determine on which side of each hyperplane the region is in order to find a solution of the state equations. In our previous terminology we have with

$$I_{p} = \{m \mid s_{p_{m}} = 1\}$$

$$R_{p} = \{x \mid C_{i}, x + h_{i} \ge 0, i \notin I_{p} \land C_{i}, x + h_{i} < 0, i \in I_{p}\} (6.41)$$

which is precisely the definition of the regions in Chua's global model [21]. Then for this specific case with D = I we obtain an alternate state model from (6.37) and (6.38) given by

$$y = \bar{A}x + \bar{B}[\bar{C}x + \bar{h}] + \bar{g} , \qquad (6.42)$$

since  $\overline{D} = 0$ . The expression for a single component  $y_j(x)$  from eq. (6.42) is exactly identical to eq. (3) in [21].

It is easy to see that the global description (6.42) became possible because the state equations (6.39) or (6.37b) could explicitly be solved in this special case. However the model (6.37) is more general than Chua's global model, due to the implicit equation (6.37b) which defines z. As we already showed in section 6.4, the state equation can also be solved explicitly if  $\overline{D}$  is upper or lower triangular, yielding explicit expressions for the corresponding PL mappings. Moreover when dealing with a mapping f :  $x \in \mathbb{R}^{m} \rightarrow y \in \mathbb{R}^{k}$ and k > 1, the global model of Chua cannot be used in general although the regions  $R_{i}$  satisfy (6.41), as can be seen from the example in fig.6.2, for which the state-model is given by

81





The corresponding affine mapping  $y = A^{p}x$  in each region  $R_{p}$  is indicated in the figure. The maps of the component functions  $y_{1}(x)$  and  $y_{2}(x)$  can be easily obtained from fig.6.2 and are depicted in fig.6.3. From this figure, it is immediately obvious that (6.41) does not hold for the regions of the components  $y_{1}(x)$ and  $y_{2}(x)$ , hence modelling by Chua's global model is not possible although the regions of the vector function y(x) do satisfy (6.41). In this case the state equation in (6.43) cannot explicitly be solved to obtain eq. (6.42).





As we will show in the next chapter, the state-model formulation has additional advantages which makes it more versatile than Chua's model.

### 7. The inverse mapping and some homeomorphism conditions

Up to now, the state-model has been used to define a PL mapping f : x  $\in \mathbb{R}^m \to y \in \mathbb{R}^k$  and some algorithms have been described to obtain the image vector v of a given vector x. These state-model descriptions can be used to model static nonlinear electronic devices such as bipolar transistors, digital gates etc. With these models available, it will be possible to describe a static nonlinear electrical network constructed with these devices by an overall state-model. For general nonlinear networks, the calculation of the DC network response due to some applied excitation amounts to solving a set of nonlinear equations. In the same way, for piecewise-linear networks with a given state-model description, the DC network response is obtained by solving a set of piecewiselinear equations. With the state-model defining y = f(x), it is then required to determine the vector  $\mathbf{x}$  for a prescribed vector  $\mathbf{y}$ , which is however equivalent to the determination of the inverse mapping  $x = f^{-1}(y)$ . If we are able to determine the state-model of this inverse mapping, the same algorithms as before can be used to find the image x for a given y. In the next sections we will show how to construct such a state-model and discuss some homeomorphism conditions for a PL function y = f(x).

#### 7.1. The state model of the inverse mapping

When discussing an inverse mapping in an Euclidean space, it is obvious that the domain and the range of the particular mapping have to be of the same dimension, otherwise the inverse mapping cannot be defined properly. Hence let us consider a PL mapping f from  $\mathbb{R}^{m}$  into  $\mathbb{R}^{m}$ , defined by the state-model S(A,B,C,D,g,h).

y = Ax + Bi + g v = Cx + Di + h  $v \ge 0, \quad i \ge 0, \quad v^{\dagger}i = 0$   $A = (m \times m), \quad B = (m \times n), \quad C = (n \times m), \quad D = (n \times n)$   $x, y \in \mathbb{R}^{m}, \quad y, i \in \mathbb{R}^{n}$ (7.1)

81

The domain X of f consists of a number of polytopes  $R_i$  in which f is a linear mapping, i.e.

$$\mathbf{x} \in \mathbf{R}_{i} \rightarrow \mathbf{y} = \mathbf{A}^{i}\mathbf{x} + \mathbf{g}^{i}$$
,  $\mathbf{A}^{i} = (\mathbf{m} \times \mathbf{m})$ , (7.2)

with  $A^{i}$  and  $g^{i}$  defined by the equivalent representation  $S_{i}$ . For a proper inverse mapping to exist, it is required that at least one of the matrices  $A^{i}$  in (7.2) is nonsingular, otherwise the whole domain  $X = \bigcup_{i=1}^{M} R_{i}$  will be mapped into a subspace of  $\mathbb{R}^{m}$  and the inverse mapping is nowhere defined. Without loss of generality we assume that in representation  $S_{0}$ : det(A) = det(A^{0}) \neq 0. From the first equation of (7.1) we then eliminate x, yielding

$$x = A^{-1}y - A^{-1}Bi - A^{-1}g$$
 (7.3)

Substitution of (7.3) in the remaining equations then results into

$$x = \widetilde{A}y + \widetilde{B}i + \widetilde{g}$$

$$v = \widetilde{C}y + \widetilde{D}i + \widetilde{h}$$

$$v \ge 0, \quad i \ge 0, \quad v^{t}i = 0,$$

$$(7.4)$$

with

$$\widetilde{A} = A^{-1}, \quad \widetilde{B} = -A^{-1}B, \quad \widetilde{C} = CA^{-1}, \quad \widetilde{D} = D - CA^{-1}B$$

$$\widetilde{g} = -A^{-1}g, \quad \widetilde{h} = h - CA^{-1}g$$
(7.5)

Since (7.4) is in fact nothing else than the equation set (7.1), written in a different way, it is equivalent to (7.1) and describes the same mapping. However (7.4) is also a state-model describing a PL mapping  $h : y \in \mathbb{R}^m \to x \in \mathbb{R}^m$ , hence (7.4) is the state-model of the required inverse mapping  $f^{-1}$  and is denoted by  $s^{-1}(A,B,C,D,g,h)$  or  $s(\widetilde{A},\widetilde{B},\widetilde{C},\widetilde{D},\widetilde{g},\widetilde{h})$ .

Observe that the existance of only one nonsingular matrix  $A^{i}$  was sufficient to obtain  $s^{-1}$ . All other  $A^{i}$  may have been singular, which means that  $f^{-1} : y \to x$  may be a one to many mapping. If some

 $\widetilde{A}^{i}$  is singular, the same holds for  $f : x \rightarrow y$ , yielding the somewhat surprising result that a mapping  $f : x \rightarrow y$  which is one to many, has an inverse mapping  $f^{-1} : y \rightarrow x$  which can still be described by a single state-model.

As an example consider the mapping  $f : \mathbb{R} \rightarrow \mathbb{R}$  defined by

$$y = x - 2i$$
  
 $v = -x + i + 1$ ,  $v \ge 0$ ,  $i \ge 0$ ,  $v \cdot i = 0$ 
(7.6)

From (7.5) the state-model  $s^{-1}$  becomes

$$x = y + 2i$$
  
 $v = -y - i + 1$ ,  $v \ge 0$ ,  $i \ge 0$ ,  $v \cdot i = 0$ 
(7.7)

By considering all representations  $S_i$  and  $S_i^{-1}$ , the particular mappings f and  $f^{-1}$  can be depicted as indicated in fig.7.1.





From this figure we see that  $f^{-1} : y \to x$  is one to many because  $\widetilde{R}_0$  and  $\widetilde{R}_1$  overlap. Furthermore from (7.7) is is easily found that the state-equation v = -y + 1 - i has no solution for y > 1, indicating that  $f^{-1}$  is not defined for y > 1, which is obvious from fig.7.1.

As for the state-model S, the state-model  $S^{-1}$  also has a number of equivalent representations  $S_i^{-1}$ , yielding a natural description for the mapping  $f^{-1}$  in polyhedral regions  $\widetilde{R}_i$ . It is easy to see that a pair  $\overline{v}, \overline{i}$  which satisfies the state equation of S for a given  $x_0$  also satisfies the state-equations of  $S^{-1}$ , for  $y = Ax_0 + B\overline{i} + g$ . Hence the regions  $R_i$  and  $\widetilde{R}_i$  both correspond to the same state in S and S<sup>-1</sup> respectively. However when some representation  $S_k$  does not exist, this does not mean that the corresponding  $S_k^{-1}$  does not exist, as can be seen from the following example with the various regions depicted in fig.7.2.

$$\begin{pmatrix} y \\ y \\ -y \\ v \end{pmatrix} = \begin{pmatrix} 1 & 0 & | & 3 & 2 \\ 0 & 2 & | & -2 & -1 \\ -1 & 0 & | & 1 & 1 \\ 0 & -1 & | & 1 & 1 \end{pmatrix} \begin{pmatrix} x \\ -1 \\ y \end{pmatrix} = \begin{pmatrix} 1 & 0 & | & -3 & -2 \\ 0 & 1/2 & | & 1 & 1/2 \\ -1 & 0 & | & 4 & 3 \\ 0 & -1/2 & | & 0 & 1/2 \end{pmatrix} \begin{pmatrix} y \\ -1 \\ y \end{pmatrix} (7.8)$$

$$v \ge 0, \quad i \ge 0, \quad v^{t}i = 0$$

$$\begin{pmatrix} R_{1} \\ (1 & 2) \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ -2 & 2 \end{pmatrix} \begin{pmatrix} x \\ 1 & -2 \\ -2 & 2 \end{pmatrix} \begin{pmatrix} y \\ 2 \\ 1 & -2 \\ -2 & -1 \end{pmatrix} \begin{pmatrix} y \\ -2 \\ -1 & 0 \\ -2 & 2 \end{pmatrix} \begin{pmatrix} y \\ 2 \\ R_{1} \end{pmatrix} \begin{pmatrix} 1 & -2 \\ -2 & 2 \end{pmatrix} \begin{pmatrix} y \\ 2 \\ -2 & -2 \end{pmatrix} \begin{pmatrix} y \\ -2 & -2 \\ -1 & 0 \\ -2 & 2 \end{pmatrix} \begin{pmatrix} y \\ 2 \\ R_{1} \end{pmatrix} \begin{pmatrix} y \\ -2 \\ -2 & -2 \end{pmatrix} \begin{pmatrix} y \\ -2 & -2 \\ -1 & 0 \\ -2 & 2 \end{pmatrix} \begin{pmatrix} y \\ R_{1} \end{pmatrix} \begin{pmatrix} y \\ -2 \\ -2 & -2 \\ -1 & 0 \\ -2 & 2 \end{pmatrix} \begin{pmatrix} y \\ R_{1} \end{pmatrix} \begin{pmatrix} y \\ -2 \\ -2 & -2 \\ -1 & 0 \\ -2 & 2 \end{pmatrix} \begin{pmatrix} y \\ R_{2} \end{pmatrix} \begin{pmatrix} y \\ R_{1} \end{pmatrix} \begin{pmatrix} y \\ R_{2} \end{pmatrix}$$

fig.7.2.

From (7.8) we find  $\det(D_{\{1,2\}}) = 0$ ,  $\det(\widetilde{D}_{\{1,2\}}) \neq 0$ . Hence  $S_3$ does not exist while  $S_3^{-1}$  does. Using the relations (3.17) it can easily be shown that for every k for which  $S_k^{-1}$  exists and  $S_k$  does not exist, the determinant of  $\widetilde{A}^k$  will be zero. Hence such regions  $\widetilde{R}_k$  will be mapped on some subspace (hyperplane or corner) of the X-space. In the particular example  $\widetilde{R}_3$  is mapped into the hyperplane  $x_1 = x_2 = 1/4(y_1+y_2)$ . Thus although  $S_3$  does not exist, there is a region  $R_3$  corresponding to  $\widetilde{R}_3$  given by  $x_1 - x_2 = 0$ . Observe that the mapping  $f : x \neq y$  is not continuous on the boundary between  $R_1$  and  $R_2$ , due to the fact that two pivoting steps are necessary to enter region  $R_2$  from region  $R_1$ . From the above considerations it will be clear that the statemodel description of a PL mapping has important advantages over descriptions based on a list of mappings. In the latter case, the description of the inverse mapping requires the boundaries of all . regions  $\widetilde{R}_i$  corresponding to  $R_i$  to be determined before such an explicit list of regions for  $f^{-1}$  can be produced. Moreover, due to the linear-algebraic nature of the state-model descriptions, it is a simple matter to model all kinds of linear transformations of a PL mapping by performing these transformations on its statemodel description. In particular the model for a concatenation  $t \circ f$  of two PL mappings  $f : x \to z$  and  $t : z \to y$  is easily derived. To this purpose considere f and t to be defined by

 $z = Ax + Bi_{1} + g \qquad y = Ez + Fi_{2} + k$   $v_{1} = Cx + Di_{1} + h \qquad v_{2} = Gz + Hi_{2} + 1 \qquad (7.9)$   $v_{1} \ge 0, \quad i_{1} \ge 0, \quad v_{1}^{t}i_{1} = 0 \qquad v_{2} \ge 0, \quad i_{2} \ge 0, \quad v_{2}^{t}i_{2} = 0$ 

Then a simple substitution of the first model into the second one will yield the state-model  $\bar{S}$  for tof:  $x \to y$ , given by

$$y = \overline{A}x + \overline{B}i + \overline{g}$$

$$v = \overline{C}x + \overline{D}i + \overline{h}$$

$$v = \left(-\frac{v_1}{v_2}\right) \ge 0 , \quad i = \left(-\frac{i_1}{i_2}\right) \ge 0 , \quad v^{t_i} = 0$$
(7.10)

with

$$\vec{A} = EA, \quad \vec{B} = (EB \mid F)$$

$$\vec{C} = \begin{pmatrix} C \\ --- \\ GA \end{pmatrix}, \quad \vec{D} = \begin{pmatrix} D \mid 0 \\ --- \\ --- \\ GB \mid H \end{pmatrix}$$

$$\vec{g} = k + Eg, \quad \vec{h} = \begin{pmatrix} h \\ --- \\ 1+Gg \end{pmatrix}$$
(7.11)

In the next section we will use the state-model  $S^{-1}$  for presenting some homeomorphism conditions for a mapping  $f : \mathbb{R}^m \to \mathbb{R}^m$ .

## 7.2. Homeomorphic mappings

Consider a PL mapping  $f : \mathbb{R}^m \to \mathbb{R}^m$  which can be modelled by a state-model S such that  $S^{-1}$  exists. The mapping f is said to be a homeomorphism if f is continuous, one to one and onto.

Homeomorphic mappings  $f : x \rightarrow y$  play an important role in the analysis of nonlinear or piecewise-linear electrical networks. If, for an excitation y and a response x, the particular network equations have the form y = f(x), then this response exists and is unique for any excitation y if f is a homeomorphism. Thus, circuits with multi-stable DC operating points like flip-flops or circuits with hysteresis effects cannot be described by a homeomorphic mapping. With the availability of the state-models S and S<sup>-1</sup> it is then of interest to derive conditions on the matrices of these models, which will guarantee that the modelled PL mapping is a homeomorphism.

Without any constraints on the state-model in connection with the class of PL mappings to be modelled, it is a difficult task to find general necessary and sufficient conditons for a homeomorphism. However a general sufficient condition is easily obtained and is given by the following theorem.

- <u>Theorem 13</u>: Let the PL mapping  $f : x \in \mathbb{R}^m \to y \in \mathbb{R}^m$  be modelled by a state-model S for which  $S^{-1}$  exists. Then f is a homeomorphism if the matrices D and  $\widetilde{D}$  from S and  $S^{-1}$  satisfy D, $\widetilde{D} \in P$ .
- Proof: From theorem 3 in chapter 3.3 and (7.4), we know that the solution of the state-equations for each  $x, y \in \mathbb{R}^{m}$ exists and is unique when  $D, D \in P$ . Hence for each  $x \in \mathbb{R}^{m}$  the model S produces only one y = y(x) and for each  $y \in \mathbb{R}^{m}$  the model S<sup>-1</sup> only yields one x = x(y). Then the mapping f is one to one and onto. The continuity of f is implied by the unique solutions of the state equations.

89

The sufficient conditions of theorem 13 differ from those given by Fujisawa and Kuh [15] and the conditions given by Chien [25], as can be seen from the following two examples, borrowed from Chien. Example I (fig.7.3) satisfies the sufficient conditions of Chien and does not satisfy those of Fujisawa and Kuh, with the reverse situation for example II from fig.7.4. However both examples satisfy the condition of theorem 13 as can be found from the given state-models.



$$\begin{pmatrix} y_1 \\ y_2 \\ \hline v_1 \end{pmatrix} = \begin{pmatrix} 1/2 & 1 & | & 1 \\ 0 & 1/2 & | & 1 \\ \hline 1 & 1 & - & - & - \\ 1 & 1 & | & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \hline \vdots \\ 1 \end{pmatrix}$$

$$\begin{pmatrix} (-1 & 0) \\ (1 & 0) \\ (0 & 1) \end{pmatrix} \begin{pmatrix} (1 & 0) \\ (2 & 1) \end{pmatrix} \begin{pmatrix} (0 & 1) \\ (0 & 1) \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ v_1 \end{pmatrix} = \begin{pmatrix} (1 & 2 & 0) & -2 \\ (0 & 1 & 2 & 0) \\ -1 & -2 & 1 & 2 \\ 0 & -1 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

# fig.7.4.

One might expect that the conditions of theorem 13 are necessary as well because D  $\notin$  P implies that for at least one value of the vector a in v = a + Di, two different solutions of the stateequations exist [6], which could possibly yield two different vectors y for a given vector x. However this need not be the case since the particular vector a may not belong to the subspace spanned by Cx + h for  $x \in \mathbb{R}^{m}$ . Furthermore, if it does belong to that subspace, the second solution of the state-equations not necessarily yields a different value for y because the difference of both solutions may ly in the null-space of the matrix B.

To illustrate the problems we are faced with, we will discuss some examples.

As a first example we again consider the mapping given by (7.8) and depicted in fig.7.2, for which we find D  $\epsilon$  P<sub>0</sub> and  $\tilde{D} \epsilon$  P. The mapping is not a homeomorphism although the conditions are close to the sufficient conditions of theorem 13.

The fact that the conditions of theorem 13 are not necessary can be seen from the next example given in (7.12) and depicted in fig.7.5, for which both matrices D and  $\widetilde{D}$  satisfy  $D, \widetilde{D} \in P_0$  and yet the mapping is a homeomorphism.

$$y = \begin{pmatrix} 2 & -1 \\ 1 & 0 \end{pmatrix} x + \begin{pmatrix} 2 & 1 \\ 2 & 1 \end{pmatrix} i \qquad x = \begin{pmatrix} 0 & 1 \\ -1 & 2 \end{pmatrix} y + \begin{pmatrix} -2 & -1 \\ -2 & -1 \end{pmatrix} i$$
$$\Rightarrow$$
$$v = \begin{pmatrix} -1/2 & -1/2 \\ 1 & -1 \end{pmatrix} x + \begin{pmatrix} 1 & 1/2 \\ 2 & 1 \end{pmatrix} i \qquad v = \begin{pmatrix} 1/2 & -3/2 \\ 1 & -1 \end{pmatrix} y + \begin{pmatrix} 3 & 3/2 \\ 2 & 1 \end{pmatrix} i$$

$$v \ge 0, i \ge 0, v^{t}i = 0$$
 (7.12)





91

In contrast with example I, the homeomorphism is due to the fact that in this case the matrix  $\widetilde{D}$  has the same singular principal submatrices as the matrix D.

Finally, the third example is given by

$$y = x + (2 - 7)i \qquad x = y + (-2 - 7)i$$
$$v = \begin{pmatrix} 1 \\ 1 \end{pmatrix} x + \begin{pmatrix} 1 & -2 \\ 3/4 & -1 \end{pmatrix} i \qquad v = \begin{pmatrix} 1 \\ 1 \end{pmatrix} y + \begin{pmatrix} -1 & 5 \\ -5/4 & 6 \end{pmatrix} i$$

The corresponding mapping is depicted in fig.7.6.



Although D  $\notin$  P,P<sub>0</sub> and  $\widetilde{D} \notin$  P,P<sub>0</sub> the mapping is a homeomorphism. However, the model is redundant since a one-dimensional stateequation would already be sufficient to divide the domain in two half-spaces.

The given three examples all have the common property that there is at least one boundary hyperplane between two adjacent regions for which two pivoting steps have to be performed in order to obtain the natural equivalent representations of both adjacent regions from each other. In other words, although two regions are adjacent in space, their corresponding vertices in the structure of D are not adjacent, which indicates some redundancy in the state-model.

In order to eliminate this situation, we again use the minimal state-model of section 4.2 for a class of piecewise-linear mappings, which is now more accurately defined as follows. <u>Definition</u>: A state-model is called a *minimal state-model* if and only if the natural equivalent representations of two adjacent regions in space yield different mappings and can be obtained from each other by a single pivoting step (one-step principal transformation).

Now consider a state-model S with a state-equation v = Cx + h + Di, C = (n×m), for which the k-th component of Cx + h is denoted by  $C_k \cdot x + h_k$ . The region R<sub>j</sub> corresponding to representation S<sub>j</sub> is then defined by

$$R_{j} = \{x \mid C_{k}^{j}, x + h_{k}^{j} \ge 0, k=1,2,...,n\}$$
(7.13)

Let some hyperplanes  $K_k^{j}$ , defined by

$$\kappa_{k}^{j} = \{ \mathbf{x} \mid c_{k}^{j}, \mathbf{x} + h_{k}^{j} = 0 \}$$
 (7.14)

be depicted in fig.7.7.



From this figure it is obvious that the hyperplane  $K_1^{j}$  is not an actual boundary of the region R. In terms of linear algebra this means that the inequality  $C_{1}^{j} \cdot x + h_{1}^{j} \ge 0$  depends on the remaining inequalities of (7.13) in such a way that it is satisfied if the remaining inequalities are satisfied. This implies that the row  $C_{1}^{j}$  is a positive linear combination of the remaining rows of  $C^{j}$  and hence rank  $(C^{j}) < n$ . We will call such hyperplanes virtual boundary hyperplanes.

Definition: The hyperplane  $K_k^j$  is called a proper boundary hyperplane of R<sub>j</sub>, if and only if R<sub>j</sub>  $\cap K_k^j$  is non-empty.

The proper boundary hyperplanes of a region  $R_j$  can be determined by an algorithm given by Tschernikov [26], which determines whether a system of linear inequalities has a solution and in addition detects the redundant inequalities. Next consider a minimal statemodel S for which  $S^{-1}$  exists. Then for each proper boundary hyperplane  $K_k^j$ , the corresponding diagonal element  $D_{kk}^j$  is non-zero otherwise the corresponding representation in the adjacent region cannot be obtained by a single pivoting step. For the same reason,  $K_k^j \neq K_1^j$  for  $k \neq 1$ . As each proper boundary hyperplane  $K_k^j$  is coupled to a diagonal element  $D_{kk}^j$  which in turn is related to an edge in the structure of S, we are able to modify this structure by deleting all edges that do not correspond to proper boundary hyperplanes. Furthermore, we also delete all vertices and incident edges corresponding to empty regions caused by uncontrollable states. The remaining structure will be called *the proper* 

- Theorem 14: Let the mapping  $f : x \rightarrow y$  be modelled by a minimal state-model S. Then f is a homeomorphism from  $\mathbb{R}^{m}$  onto  $\mathbb{R}^{m}$  only if the proper structure of S is connected and each edge in the proper structure has a positive weight.

94

$$C_{k} \cdot x_{0} + h_{k} = \varepsilon > 0, \varepsilon \text{ arbitrarily small}$$

$$C_{1} \cdot x_{0} + h_{1} > 0 \quad , l \neq k$$

$$(7.15)$$

Next, this system is pivoted on  ${\rm D}_{kk}$  yielding a new representation S for the adjacent region R , with a state equation

$$v = C^{q}x_{0} + h^{q} + D^{q}i$$
 (7.16)

and

$$C_{1}^{q} = C_{1} - C_{k} \frac{D_{1k}}{D_{kk}}, \quad h_{1}^{q} = h_{1} - h_{k} \frac{D_{1k}}{D_{kk}}, \quad l \neq k$$

$$C_{k}^{q} = -\frac{C_{k}}{D_{kk}}, \quad h_{k}^{q} = -\frac{h_{k}}{D_{kk}}$$
(7.17)

From (7.15) and (7.17) we obtain

$$C_{k}^{q} \cdot x_{0} + h_{k}^{q} = \frac{-\varepsilon}{D_{kk}}$$

$$C_{1}^{q} \cdot x_{0} + h_{1}^{q} = (C_{1} \cdot x_{0} + h_{1}) - \frac{\varepsilon D_{1k}}{D_{kk}} \cdot 1 \neq k$$
(7.18)

Now assume  $D_{kk} < 0$ . Then from (7.18) we can always have  $C^{q}x_{0} + h^{q} > 0$  by taking  $\varepsilon$  sufficiently small, which means that  $x_{0} \in R_{q}$ . Hence the adjacent regions  $R_{0}$  and  $R_{q}$  overlap each other yielding two different images for a single point  $x_{0}$ . Then f cannot be a homeomorphism and thus  $D_{kk}$  has to be positive. The same arguments hold for any representation  $S_{j}$  with  $R_{j}$  non-empty, which proves the theorem.

Next we will derive a relation between the determinants of two matrices  $A^{j}$  and  $A^{k}$  from the affine mappings in  $R_{j}$  and  $R_{k}$ . To this purpose, consider a state-model S, for which  $S^{-1}$  exists, and S and  $H^{p}$  defined by the relations (4.16) and (4.17) in section 4.2. With (4.18) and (4.20) we obtain with p = 0

$$A^{q} = A - B_{M} D_{M}^{-1} C_{M}$$
,  
 $M = \{k \mid s_{qk} = 1\}$ 
(7.19)

From (7.19) we have

$$\det(A^{q}) = \det(A - B_{M} D_{M}^{-1} C_{M}) = \det(A) \cdot \det(1 - A^{-1} B_{M} D_{M}^{-1} C_{M})$$
(7.20)

Using the relation det(1 + PQ) = det(1 + QP) this yields

$$det(A^{q}) = det(A) \cdot det(1 - D_{M}^{-1}C_{M}A^{-1}B_{M}) =$$
$$= det(A) \cdot det(D_{M} - C_{M}A^{-1}B_{M})/det(D_{M})$$
(7.21)

However from (7.5) we obtain  $det(D_M - D_M A^{-1}B_M) = det(\widetilde{D}_M)$ , which results with (7.21) into

$$\frac{\det(A^{q})}{\det(A)} = \frac{\det(\widetilde{D}_{M})}{\det(D_{M})}$$
(7.22)

With  $J = \{k \mid s_k = 1\}$ , eq(7.22) then yields

$$\frac{\det(\mathbf{A}^{\mathbf{q}})}{\det(\mathbf{A}^{\mathbf{r}})} = \frac{\det(\widetilde{\mathbf{D}}_{\mathbf{M}})}{\det(\widetilde{\mathbf{D}}_{\mathbf{J}})} \cdot \frac{\det(\mathbf{D}_{\mathbf{J}})}{\det(\mathbf{D}_{\mathbf{M}})}$$
(7.23)

Now assume that  $R_r$  and  $R_q$  are adjacent with a common boundary hyperplane  $K_k^q = K_k^r$ , i.e. k = M - J. Then (7.23) yields with (3.12)

$$\frac{\det(\mathbf{A}^{\mathbf{q}})}{\det(\mathbf{A}^{\mathbf{r}})} = \frac{\det(\widetilde{\mathbf{D}}_{\mathbf{J}} - \{\mathbf{k}\})}{\det(\widetilde{\mathbf{D}}_{\mathbf{J}})} \cdot \frac{\det(\mathbf{D}_{\mathbf{J}})}{\det(\mathbf{D}_{\mathbf{J}} - \{\mathbf{k}\})} = \frac{\widetilde{\mathbf{D}}_{\mathbf{kk}}^{\mathbf{r}}}{\mathbf{D}_{\mathbf{kk}}}$$
(7.24)

We are now able to prove the following theorem.

Theorem 15: Let  $f : \mathbb{R}^m \to \mathbb{R}^m$  be modelled by a minimal state-model S for which S<sup>-1</sup> exists. Then f is a homeomorphism only if the proper structures of S and S<sup>-1</sup> are isomorphic and the weights on all edges are positive, Proof:

Let f be a homeomorphism, which requires that det( $A^{j}$ )  $\neq 0$  for each non-empty region R<sub>j</sub>. Then each non-empty region R<sub>j</sub> is mapped by  $y = A^{j}x + g^{j}$  onto a region  $\widetilde{R}_{j}$  in Y-space, corresponding to an existing vertex in the proper structure of  $S^{-1}$ . By theorem 14 and (7.24), an edge in  $S^{-1}$  which corresponds to a proper boundary hyperplane between two regions  $\widetilde{R}_{r}$ and  $\widetilde{R}_{q}$  then has a nonzero weight  $\widetilde{D}_{kk}^{r}$ . Furthermore, in Y-space there cannot exist a region  $\widetilde{R}_{p}$  adjacent to  $\widetilde{R}_{q}$  which has no corresponding region  $R_{p}$  in S, because this could only be the case when det( $\widetilde{A}^{p}$ ) = 0.

By (7.24), using  $S = (S^{-1})^{-1}$ , this would yield a zero weight on the corresponding edge in S which contradicts the minimality of S. Hence the proper structures of S and  $S^{-1}$  are isomorphic. In addition for a homeomorphism, adjacent regions in Y-space may not overlap, which results in the requirement that all weights on the edges of the proper structure of  $S^{-1}$  have to be positive as well.

<u>Corollary</u>: Let f be modelled by a minimal statemodel for which  $s^{-1}$  exists. Then f is a homeomorphism only if the determinants of the Jacobians A<sup>j</sup> in the regions R<sub>j</sub> have the same sign.

This corollary easily follows from theorem 15 and relation (7.24) and has already been proved by Fujiwawa and Kuh [15]. However the conditions of theorem 15 are stronger since two corresponding edges in S and S<sup>-1</sup> with negative weight would also yield the same sign for the determinants of the adjacent regions. By theorem 15 this is not allowed for a homeomorphism.

Next we consider a PL mapping y = f(x), mapping  $\mathbb{R}^{m}$  into  $\mathbb{R}^{m}$ . For a given  $\bar{y}$  there are at most finitely many solutions of the equation  $\bar{y} = f(x)$ . Let these solutions be given by  $\bar{x}^{1}$ ,  $\bar{x}^{2}$ ,..., $\bar{x}^{k}$ . With  $\bar{x}^{j} \in \mathbb{R}_{n_{j}}$ , j=1,2,...,k, the Jacobian of the mapping in each region  $\mathbb{R}_{n_{j}}$  is given by  $\mathbb{A}^{n_{j}}$ . Following Ohtsuki et al [27], we define the degree of the above mapping f at  $\bar{y}$ , as the integer

$$deg(f, \overline{y}) = \sum_{j=1}^{k} sign(det(A^{j}))$$
(7.25)

Now, the following invariance property holds

Let  $f : \mathbb{R}^m \to Q \subset \mathbb{R}^m$  and  $\overline{y}^1, \overline{y}^2 \in Q$ . Then  $\deg(f, \overline{y}^1) = \deg(f, \overline{y}^2)$ , if  $\overline{y}^1$  and  $\overline{y}^2$  can be connected by a path in Q. (7.26)

We next prove the following theorem.

Theorem 16: Let 
$$f : \mathbb{R}^m \to \mathbb{R}^m$$
 be modelled by a minimal state-model   
S for which S<sup>-1</sup> exists. Then f is a homeomorphism from   
 $\mathbb{R}^m$  onto  $\mathbb{R}^m$  if and only if

 a) the proper structures of S and S<sup>-1</sup> are isomorphic, connected and have only positive weights on the edges,

b) there exists a pair 
$$\bar{x}, \bar{y}$$
 with  $\bar{y} = f(\bar{x})$  for which  
 $|\deg(f, \bar{y})| = |\deg(f^{-1}, \bar{x})| = 1.$ 

Proof: Necessity: the necessity is obvious from theorem 15 and the definition (7.25).

Sufficiency: When condition a) holds, the system S produces a vector y for any vector  $x \in \mathbb{R}^m$  because Katzenelson's algorithm always works since only positive pivots will show up during the run of the algorithm which always allows for an increase of the parameter  $\lambda$ . As the same holds for S<sup>-1</sup> and any  $y \in \mathbb{R}^m$ , the mapping is onto. Conditon b), together with (7.26) then guarantees that the mappings is also one to one because on account of a) the determinants of the Jacobian have the same sign. Thus the mapping is a homeomorphism.

To demonstrate the application of the foregoing theorems we consider an example given by Kuh [15] which is depicted in fig.7.8 and can be described by:



fig.7.8.

fig.7.9.

The proper structure corresponding to (7.27) is given in fig.7.9. The inverse system can be found from (7.27) and yields

$$\mathbf{x} = \begin{pmatrix} 0 & 1/\sqrt{3} \\ -\sqrt{3} & 0 \end{pmatrix} \mathbf{y} + \begin{pmatrix} 0 & 1/\sqrt{3} & -1/\sqrt{3} \\ -4 & 1 & 1 \end{pmatrix} \mathbf{i}$$

$$\mathbf{v} = \begin{pmatrix} -\sqrt{3} & 0 \\ -\sqrt{3} & -1 \\ -\sqrt{3} & 1 \end{pmatrix} \mathbf{y} + \begin{pmatrix} -3 & 1 & 1 \\ -4 & 1 & 2 \\ -4 & 2 & 1 \end{pmatrix} \mathbf{i}$$

$$\mathbf{v} \ge 0, \quad \mathbf{i} \ge 0, \quad \mathbf{v}^{\mathsf{t}} \mathbf{i} = 0$$

$$(7.28)$$

The corresponding regions in Y-space and the proper structure of  $s^{-1}$  are given in fig.7.10 and 7.11 respectively.





fig.7.10

fig.7.11

By this mapping the whole X-space is mapped twice on Y-space as can be seen from fig.7.10. Hence  $\forall_{\overline{y}} \in \mathbb{R}^2 \deg(f,\overline{y}) = 2$ . Furthermore both proper structures of S and S<sup>-1</sup> satisfy condition a) of theorem 16. Due to the violation of condition b) the mapping is not a homeomorphism.

The conditions in theorem 16 can be relaxed for  $C = (n \times m)$  when rank(C) = n. For such a case we have the following theorem.

- <u>Theorem 17</u>: Let  $f : \mathbb{R}^{m} \to \mathbb{R}^{m}$  be modelled by a minimal state model S for which S<sup>-1</sup> exists and let rank(C) = n, for C = (n×m). Then f is a homeomorphism from  $\mathbb{R}^{m}$  onto  $\mathbb{R}^{m}$  if and only if  $D, \widetilde{D} \in P$ .
- **Proof:** Necessity: Assume that  $D \notin P$ . Since rank(C) = n, we can construct any vector  $a \in \mathbb{R}^n$  by some particular x from  $\mathbb{R}^m$ . Next we pivot the system S in any possible state. Then each state  $s_p$  corresponds to a region  $\mathbb{R}_p \subset \mathbb{R}^m$  which is non-empty since we can always find a vector  $\bar{x}$  such that  $C^p \bar{x} + h^p = \bar{a} > 0$ . As a result each row of  $C^p x + h^p = 0$  has to be a proper boundary

hyperplane of  $R_p$ . Furthermore all representations  $s^p$  have to exist, otherwise we would be able to find a region  $R_q$  and a particular proper boundary hyperplane  $K_j^q$  which corresponds to a zero diagonal element of  $D^q$ , contradicting the minimality of S. Hence the proper structure of S is a full Hasse diagram with  $2^n$  vertices in which each edge represents a proper boundary hyperplane. When D  $\notin$  P, at least one edge will have a negative weight and then by theorem 15 the mapping is not a homeomorphism. The same arguments hold for  $s^{-1}$ .

Sufficiency: Let  $D, \widetilde{D} \in P$ . Then as before the proper structures of S and S<sup>-1</sup> are both full Hasse-diagrams representing 2<sup>n</sup> states with a positive weight on all edges. Thus condition a) of theorem 16 is satisfied. When  $D \in P$ , the solution of the state equations of S is unique for any x and thus each vector x produces only one vector y, i.e.  $\deg(f^{-1},x) = \pm 1$ . In the same way  $\widetilde{D} \in P$  will yield  $\deg(f,y) = \pm 1$ . Then by theorem 16 the mapping is a homeomorphism.

As an example, the mappings depicted before in fig.7.3 and 7.4 are homeomorphism on account of theorem 17.

7.3. The construction of the state-model for some simple mappings

Up to now a general construction of the state-model from a listspecified PL mapping is not known. For a relatively small number of segments or regions the state-model can readily be found by trial and error. Only for those mappings for which the regions are defined by parallel hyperplanes a general method does exist and will be discussed below. A specific example of this class of mappings is given by one-dimensional PL functions, which class will be treated first. The models of more complex mappings have to be constructed by concatenation and interrelation of these onedimensional mappings, which is treated in chapter 8. Let a one-dimensional continuous PL function y = f(x) be defined as depicted in fig.7.12. The linear segments are defined in the intervals  $[-\infty, a_1]$ ,  $[a_1, a_2]$ ,...,  $[a_{n-1}, a_n]$ ,  $[a_n, \infty]$  and have a slope  $e_k$  in interval  $[a_{k-1}, a_k]$ .



fig.7.12.

In order to find the state-model for f(x) we could use the method of chapter 2.3 to construct a PL two-pole element having a v-i characteristic corresponding to v = f(i), which can then be modelled by some hybrid matrix description of the obtained resistive network loaded with ideal diodes. However a faster method is obtained by considering the following state-model.

$$y = e_{1}x + (b_{1} \ b_{2} \ \dots \ b_{n})i + g_{1}$$

$$v = \begin{pmatrix} -1 \\ -1 \\ \cdot \\ \cdot \\ -1 \end{pmatrix} x + \begin{pmatrix} 1 & 0 & \cdot & 0 \\ 0 & 1 & \cdot \\ \cdot & 1 & 0 \\ 0 & \cdot & 0 & 1 \end{pmatrix} i + \begin{pmatrix} a_{1} \\ a_{2} \\ \cdot \\ \cdot \\ a_{n} \end{pmatrix}$$

$$v \ge 0, \quad i \ge 0, \quad v^{t}i = 0$$
with  $a_{1} < a_{2} < \dots < a_{n}$ 

$$(7.29)$$

From (7.29) we easily see that the natural hybrid representation of (7.29) for  $x \in [a_{k-1}, a_k]$  is found by pivoting on all pivots  $D_{11}$ 

with j < k since each pivot D<sub>jj</sub> only affects the relation between  $v_j$  and  $i_j$ . Hence when x runs from the left to right along the x-axis, a pivot D<sub>kk</sub> must be used when x crosses the boundary  $a_k$ . From (7.29) the slope  $e_{k+1}$  of y = f(x) in the new interval  $[a_k, a_{k+1}]$  then becomes  $e_{k+1} = e_k + b_k$ . Then, in order for (7.29) to be the state-model of f(x) in fig.7.12, we have to take  $b_k = e_{k+1} - e_k$ , k=1,2,...,n. The value of  $g_1$  is obtained from the representation in state  $s_0$  when  $x \le a_1$ , yielding  $y = e_1 x + g_1$ .

Observe that the pivoting of D proceeds from top-left to bottom-right along the diagonal when x goes from left to right along the x-axis. Hence the states corresponding to existing regions are given by the state-vectors

$$(0,0,\ldots,0), (1,0,\ldots,0), (1,1,0,\ldots,0), \ldots, (1,1,1,\ldots,1)$$

All other states will yield empty regions. Obviously the mapping y = f(x) will be a homeomorphism if all slopes in fig.7.12 have the same sign, with no slope being zero. Using  $e_{k+1} = e_k + b_k$ , this means that (7.29) represents a homeomorphism if and only if

$$sign(e_1 + \sum_{k=1}^{m} b_k) = sign(e_1), m=1,2,...,n$$
 (7.30)

Next for a PL mapping  $f : x \in \mathbb{R}^m \to y \in \mathbb{R}^m$  such that each region in y-space is bounded by parallel hyperplanes, the state-model takes the following form

$$y = Ax + Bi + g$$

$$v = \begin{pmatrix} -n^{t} \\ -n^{t} \\ \cdot \\ \cdot \\ -n^{t} \end{pmatrix} x + 1i + \begin{pmatrix} a_{1} \\ a_{2} \\ \cdot \\ \cdot \\ a_{k} \end{pmatrix}$$

$$v \ge 0, \quad i \ge 0, \quad v^{t}i = 0$$
and  $a_{1} \le a_{2} \le \dots \le a_{k}$ ,  $A = (m \times m), \quad B = (m \times k)$ 

$$(7.31)$$
In (7.31) the parallel hyperplanes are given by  $n^{t}x = a_{j}$ ,  $j=1,2,\ldots,k$ , with n representing the common normal vector on these hyperplanes.

Let us consider a vector  $\bar{x} \in R_0$  satisfying  $n^{t}\bar{x} < a_1$ . Then  $R_0$ is a halfspace bounded by the hyperplane  $n^{t}x = a_1$ . From  $\bar{x}$  we are going to trace a line  $x = \bar{x} + \lambda n$ , orthogonal to all parallel hyperplanes and traversing all strip-shaped regions in x-space. Each time a boundary hyperplane is crossed, the corresponding diagonal element of the unit matrix (matrix D) is used as a pivot to obtain the new natural representation in the adjacent region, resulting in an update of the affine mapping. Let us denote the regions and hyperplanes encountered along the previously given line as given in fig.7.13.



fig.7.13.

With the mapping in  $R_k$  defined by  $y = A_k x + g_k$  we then find from (7.31) and the discussed process

$$A_{1+1} = A_1 + B_{1+1}n^{t}$$
 (7.32)

with  $B_1$  representing the 1-th column of B. Thus the matrix A and the columns of B are easily found from the various mappings in the different regions. For mappings of the above type, the following theorem holds.

- -

- <u>Theorem 18</u>: A continuous PL mapping  $f : x \in \mathbb{R}^m \to y \in \mathbb{R}^m$  such that the regions in x-space are bounded by parallel hyperplanes is a homeomorphism if and only if the determinants of the Jacobian in all regions have the same sign.
- **Proof:** From the foregoing arguments we know that the mapping can be modelled by the state-model of eq.(7.31). Obviously y = f(x) is a homeomorphism if and only if  $\tilde{y} = Mf(Nx)$  is a homeomorphism for some nonsingular pair of matrices M.N because M and N only perform linear transformations on the whole x- and y-spaces and do not affect the one to one and onto properties of the mapping. In particular we take  $M = NA^{-1}$  and define  $\tilde{x}$  by  $\tilde{x} = Nx$ , with

$$\mathbf{N} = \begin{pmatrix} -n_1 & -n_2 & \cdot & \cdot & \cdot & \cdot & n_m \\ 0 & 1 & 0 & \cdot & \cdot & 0 \\ 0 & 0 & 1 & \cdot & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & 1 & 0 \\ 0 & \cdot & \cdot & 0 & 1 \end{pmatrix} ,$$
(7.33)

such that  $n^{t}N^{-1} = (-1 \ 0 \ 0 \ ... \ 0)$ . With (7.31) this results in a state-model for  $\tilde{y} = Mf(\tilde{x})$  given by  $\tilde{y} = i_{m}\tilde{x} + \tilde{B}i + g$   $v = \begin{pmatrix} -1 \ 0 \ 0 \ ... \ 0 \\ -1 \ 0 \ 0 \ ... \ 0 \\ ... \ ... \\ ... \\ -1 \ 0 \ 0 \ ... \ 0 \end{pmatrix} \qquad \tilde{x} + 1_{k}i + \begin{pmatrix} a_{1} \\ a_{2} \\ ... \\ a_{k} \end{pmatrix}$  (7.34) with  $v \ge 0$ ,  $i \ge 0$ ,  $v^{t}i = 0$ 

In stead of the original problem we now prove that (7.34) is a homeomorphism if and only if the determinants

of the Jacobian in all regions have the same sign. To this purpose we partition the vectors  $\tilde{x}$  and  $\tilde{y}$  as follows

$$\widetilde{\mathbf{y}} = \begin{pmatrix} \widetilde{\mathbf{y}}_1 \\ - & - \\ & \widetilde{\mathbf{y}}_r \end{pmatrix} \} \quad 1 \text{ component} \qquad \qquad \widetilde{\mathbf{x}} = \begin{pmatrix} \widetilde{\mathbf{x}}_1 \\ - & - \\ & & \widetilde{\mathbf{x}}_r \end{pmatrix}$$

This partition then yields with (7.34)

$$\widetilde{\mathbf{y}}_{1} = \widetilde{\mathbf{x}}_{1} + \widetilde{\mathbf{B}}_{1}\mathbf{i} + \mathbf{g}_{1}$$
(7.35a)

$$\mathbf{v} = \begin{pmatrix} -1 \\ -1 \\ \vdots \\ -1 \end{pmatrix} \stackrel{\sim}{\underset{\mathbf{x}_{1}}{\times}} + 1_{\mathbf{k}}\mathbf{i} + \begin{pmatrix} \mathbf{a}_{1} \\ \mathbf{a}_{2} \\ \vdots \\ \mathbf{a}_{\mathbf{k}} \end{pmatrix}$$
(7.35b)

$$\widetilde{Y}_{\mathbf{r}} = \widetilde{\mathbf{x}}_{\mathbf{r}} + \widetilde{B}_{\mathbf{r}}\mathbf{i} + g_{\mathbf{r}}$$
(7.35c)
$$\mathbf{v} \ge 0, \quad \mathbf{i} \ge 0, \quad \mathbf{v}^{\mathsf{t}}\mathbf{i} = 0,$$
with 
$$\widetilde{B} = \begin{pmatrix} \widetilde{B}_{1} \\ -\widetilde{B}_{\mathbf{r}} \end{pmatrix} \quad \text{and} \quad g = \begin{pmatrix} g_{1} \\ -\widetilde{g}_{\mathbf{r}} \end{pmatrix}$$

A comparison of (7.35a) and (7.35b) with (7.29) shows that  $\widetilde{Y}_1$  is a one-dimensional PL function of  $\widetilde{x}_1$ , which will be denoted by  $\widetilde{Y}_1 = h(\widetilde{x}_1)$ . We next claim that (7.35a) to (7.35c) is a homeomorphism if and only if  $\widetilde{Y}_1 = h(\widetilde{x}_1)$  is a homeomorphism. Obviously we only have to prove the sufficiency of the above claim. Hence let us assume that  $\widetilde{Y}_1 = h(\widetilde{x}_1)$  is a homeomorphism. Then for each  $\widetilde{x}_1 \in \mathbb{R}$  we have a unique  $\widetilde{Y}_1 \in \mathbb{R}$  and because  $1_k \in P$  we also have a unique solution of the state equation (7.35b). For any  $\widetilde{x}_r \in \mathbb{R}^{m-1}$  we then produce by (7.35c) a unique  $\widetilde{Y}_r \in \mathbb{R}^{m-1}$ . Thus any  $\widetilde{x} \in \mathbb{R}^m$  will yield a unique  $\widetilde{Y} \in \mathbb{R}^m$ .

On the other hand, for each  $\tilde{y}_1 \in \mathbb{R}$  we find a unique  $\tilde{x} \in \mathbb{R}$ and again a unique solution of the state equations, which produces by (7.35c) a unique  $\tilde{x}_r \in \mathbb{R}^{m-1}$  for any  $\tilde{y}_r \in \mathbb{R}^{m-1}$ . Thus any  $\tilde{y} \in \mathbb{R}^m$  produces a unique  $\tilde{x} \in \mathbb{R}^m$ . Then the mapping is one to one and onto, i.e. it is homeomorphism.

From (7.35) with  $\tilde{A}^0 = 1_m$  and (7.30) we know that  $\tilde{Y}_1 = h(\tilde{x}_1)$ 

is a homeomorphism if and only if  $\operatorname{sign}(\widetilde{A}_{11}^j) = \operatorname{sign}(\widetilde{A}_{11}^0)$ for all regions R<sub>j</sub> with representation S<sub>j</sub>. However from (7.34) we also observe that pivoting on D = 1<sub>k</sub> only affects the first column of  $\widetilde{A}_1$  which remains of the form

$$\widetilde{A}^{j} = \begin{pmatrix} \widetilde{A}_{11}^{j} & 0 & 0 & . & . & 0 \\ - & - & - & - & - & - & - & - & - \\ \widetilde{A}_{21}^{j} & | & 1 & 0 & . & . & 0 \\ & | & & & & & \\ \cdot & | & 0 & 1 & & . & \\ \cdot & | & \cdot & & \cdot & \cdot & \cdot \\ \cdot & | & \cdot & & \cdot & \cdot & \cdot \\ \widetilde{A}_{m1}^{j} & | & 0 & . & . & . & 1 \end{pmatrix}$$

Hence  $det(\widetilde{A}^{j}) = \widetilde{A}_{11}^{j}$ . Then  $\widetilde{Y}_{1} = h(\widetilde{x}_{1})$  is a homeomorphism if and only if sign  $det(\widetilde{A}^{j}) = 1$  for all regions  $R_{j}$ . By our previous considerations this implies that y = f(x) is a homeomorphism if and only if the determinants of the Jacobians of all regions have the same sign. End of proof.

The above theorem was also given by Chien [25]. However the proof given there relies on another theorem for which the proof seems to be doubtful; at least the proof is not as obvious as was suggested.

With the above theorem we conclude our discussions on homeomorphism properties of mappings modelled by the state-model.

## 8. Interconnection of PL systems and dynamic PL systems

Standard nonlinear electrical networks are in general constructed by interconnection of basic network elements or building blocks such as for example transistors, operational amplifiers, gates etc. In order to construct a piecewise-linear description of such a network in the form of an overall state-model, a method has to be supplied to assemble the required state-model from the individual state-models for the employed building blocks. Since the state-model formulation has a standard structure independent of the kind of mapping which has been modelled, such an approach then allows for a hierarchical modelling and description of electrical networks, as each state-model on its own can be considered as a macro-model. Furthermore, a combination of analog and digital circuits can easily be modelled since the state-model itself has no inherently different nature for both types of circuits and a single algorithm is used to solve the PL equations which arise in the analysis of such circuits. Thus the state-model description allows for a uniform description for different types of building blocks of diverse complexity, resulting in a uniform database and solutionalgorithm as well.

In the next section we will present a method to obtain such an overall state-model for a collection of interconnected subsystems. For reasons of simplicity each in- or output variable is considered as a "signal" passed from one block to the other, where as in practical network simulations the interconnection of two nodes will require two voltages to be set equal and the sum of two currents to be set to zero.

### 8.1. The state-model construction

Consider a collection of components  $c^1, c^2, \ldots, c^k$  each of which has been modelled by an appropriate state-model  $s^1, s^2, \ldots, s^k$  such that the corresponding PL mapping gives a sufficiently accurate description of the nonlinear static input-output relation of each component. Let the in- and output variables (signals) of component  $\ensuremath{\textbf{C}}^{\ensuremath{\textbf{m}}}$ 

be denoted by the vectors  $x^m$  and  $y^m$  respectively. Interconnection of these components then amounts to equate the input and output variables which correspond to the particular interconnection.

For example the interconnections of the "circuit" of fig.8.1 are determined by  $x_2^1 = y_2^2$  and  $x_1^2 = y_2^1$ . In addition we may define the overall in- and output variables by  $I_1 = x_1^1$ ,  $I_2 = x_2^2$ ,  $O_1 = y_1^1$  and  $O_2 = y_1^2$ 



fig.8.1.

Together with the PL description of  $c^1$  and  $c^2$  these relations fully define the behaviour of the overall circuit with inputs  $I_1$ ,  $I_2$  and output  $O_1$ ,  $O_2$ . Our task is now to obtain the state-model description of such a compound circuit.

To this purpose we shall not limit ourselves to the above type of interconnections but allow for more general relations of the form

$$\sum_{i,j}^{n} a_{ij} x_{j}^{i} + \sum_{i,j}^{n} b_{ij} y_{j}^{i} = 0$$
(8.1)

Using (8.1) we also have the possibility to interconnect two inputs or two outputs. In standard electronic networks such a connection in general makes no sense. However in our modelling by the state-model the difference between input variables and output variables is more or less artificial since a pivoting step on the element  $A_{k1}$  of some state-model S exchanges  $y_k$  for  $x_1$  and hence it then appears as if  $y_k$  is an input- and  $x_1$  is an output variable.

Therefore we will consider the construction of compound models or macro-models as symbolized in fig.8.2, with the overall in- and output vectors denoted by  $\mathbf{x}^0$  and  $\mathbf{y}^0$  and the interconnection realized within the block labelled #INT and described by linear relations of the form (8.1).



fig.8.3.

First of all we rewrite the state-model equations of the individual blocks  $\textbf{C}^{\text{m}}$  in the form

$$\bar{s}^{m} \cdot \begin{pmatrix} y^{m} \\ v^{m} \\ x^{m} \\ i^{m} \end{pmatrix} + g^{m} = 0$$

ł

Then the coefficients of the matrices  $\overline{S}^{m}$  and the vectors  $g^{m}$  are entered in a matrix  $\overline{S}$  as indicated in fig.8.3 for k = 3. At the same time a pointer system is set up to memorize the exact locations of all variables in  $\overline{S}$ . In the next step all interconnections are processed by appending  $\overline{S}$  with a number of rows representing equations of the type (8.1) for each interconnection. In fig.8.3 these rows are depicted in the lower border-block.

The overall model is now easily obtained from  $\overline{s}$  by eliminating all variables of the vectors  $x^{m}$  and  $y^{m}$  which do not belong to  $x^{0}$  or  $y^{0}$ , by using pivots from the lower border block. The remaining equations are then reordered and eventually pivoted to obtain the standard form

$$\bar{\mathbf{s}}^{0} \cdot \begin{pmatrix} \mathbf{y}^{0} \\ \mathbf{v}^{0} \\ \mathbf{x}^{0} \\ \mathbf{i}^{0} \end{pmatrix} + \mathbf{g}^{0} = \mathbf{0}$$

As a simple example consider the circuit of fig.8.4 which contains a saturating amplifier in  $C^1$  and a PL feedback network in  $C^2$ , with the diode in  $C^2$  considered to be an ideal diode.



fig.8.4.

fig.8.5.

The variables in x and y represent voltages with respect to ground. Let the amplifier in  $C^1$  have a characteristic as given in fig.8.5, for which the state-model is easily obtained as

$$y_{1}^{1} = (5 -5) \begin{pmatrix} x_{1}^{1} \\ x_{2}^{1} \end{pmatrix} + (-5 -5) \begin{pmatrix} i_{1}^{1} \\ i_{2}^{1} \end{pmatrix}$$

$$\begin{pmatrix} v_{1}^{1} \\ v_{2}^{1} \end{pmatrix} = \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} x_{1}^{1} \\ x_{2}^{1} \end{pmatrix} + \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} i_{1}^{1} \\ i_{2}^{1} \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$
(8.2)

Furthermore the state-model of  $c^2$  is found to become

$$y_1^2 = x_1^2 + i_1^2$$

$$y_1^2 = -x_1^2$$
(8.3)

y <sup>1</sup>	v	1	2	۶ <sup>1</sup> .	i	1	у²	2 v <sup>2</sup>	x <sup>2</sup>	i <sup>2</sup>	g
-1	0	0	5	-5	-5	5	0	0	0	0	0
0	-1	0	-1	1	1	0	0	0	0	0	1
0	0	-1	1	-1	0	1	0	0	0	0	1
0	0	0	0	0	0	0	-1	0	1	1	0
0	0	0	0	0	0	0	0	-1	-1	0	0
0	0	0	0	-1	0	0	0	0	1	0	0
-1	0	0	0	0	0	0	1	0	0	0	0

With the given interconnections the matrix  $\bar{S}$  now becomes

The first five rows of  $\vec{s}$  contain the descriptions of  $c^1$  and  $c^2$ and the last two rows contain the interconnections. Then using row 1 and row 4, we elimintate  $y_1^1$  and  $y_1^2$  from the last interconnection equations yielding the following interconnection block.

Next we use the two encircled columns in (8.4) to eliminate  $x_2^1$  and  $x_1^2$  from the upper five rows of  $\overline{s}$ . After deleting the fourth row since we are not interested in  $y_1^2$  and reordering the remaining columns we obtain the overall state-model

$$y_{1}^{0} = 5/6 x_{1}^{0} + (-5/6 - 5/6 - 5/6) i^{0}$$

$$v^{0} = \begin{pmatrix} -1/6 \\ 1/6 \\ -5/6 \end{pmatrix} x_{1}^{0} + \begin{pmatrix} 1/6 - 5/6 - 1/6 \\ 5/6 - 1/6 - 1/6 \\ 5/6 - 5/6 - 1/6 \end{pmatrix} i^{0} + \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$$
with  $v^{0} = (v_{1}^{1}, v_{2}^{1}, v_{1}^{2})^{t}, i^{0} = (i_{1}^{1}, i_{2}^{1}, i_{1}^{2})^{t}$ 
and  $v^{0} \ge 0, i^{0} \ge 0, v^{0^{t}} i^{0} = 0.$ 

$$(8.5)$$





#### fig.8.6.

As can be seen for  $x_1^0 > 0$  the input-output relation is the same as for the amplifier since in this case the diode in C<sup>2</sup> is conducting which prevents feedback from output to input. For  $x_1^0 < 0$  the diode does not conduct and hence there is a full feedback around the amplifier yielding a gain of 5/6 until the input voltage becomes less than -6 Volt in which case the amplifier goes into saturation.

It will be clear that the number of overall state equations equals the sum of all state equations of the individual state-models of the components  $C^1, C^2, \ldots, C^k$ . Thus the storage complexity of the description grows linearly with the storage complexity of the components. Furthermore it will be obvious that linear components can be added as well since their algebraic input-output relations can simply be considered as a state-model description without state equations.

In the above way we are able to construct an overall state-model or macro-model for a system of interconnected PL subsystems which allows for a hierarchically structured modelling process for which the nature and the complexity of the components is immaterial.

#### 8.2. Dynamic PL systems

Up to now the state-model was used to describe static PL mappings and could serve as a macro-model for nonlinear static electronic components. Since most electronic networks include energy storage elements like capacitors and inductors, the state-model has to be extended in such a way that these elements can be included in order to become a feasible network analysis tool. As these energy storage elements have voltage-current relationships which are expressed by means of time-derivatives of some physical quantities like charge or flux, the extended state-model must contain some variables with associated time derivatives. Rather than trying to model separate capacitors or inductors we take a more general approach by including standard linear dynamic networks that are described by a common state-space formulation of the type

$$y = Ex + Fu$$
(8.6)  
$$\dot{u} = Gx + Hu$$

Systems with the above description can be considered as linear algebraic systems with an input vector  $\begin{pmatrix} x \\ -- \\ u \end{pmatrix}$  and an output vector  $\begin{pmatrix} y \\ -- \\ \dot{u} \end{pmatrix}$ . From this point of view, these linear systems can easily be interconnected with a static PL system using the methods presented in the previous chapter. After performing this interconnection, the final state-model description then takes the following form

$$\begin{pmatrix} Y \\ \dot{u} \end{pmatrix} = A \begin{pmatrix} x \\ -- \\ u \end{pmatrix} + Bi + g$$

$$v = C \begin{pmatrix} x \\ -- \\ u \end{pmatrix} + Di + f$$

$$v \ge 0, \quad i \ge 0, \quad v^{t}i = 0$$

$$(8.7)$$

Relation (8.7) is symbolized as depicted in fig.8.7.



fig.8.7.

These dynamic PL state-models can now be interconnected just like the standard static PL systems.

As an example we will derive the dynamic state-model for the multivibrator circuit of fig.8.8, which consists of a saturating amplifier with a non-inverting integrator in the feedback loop.



fig.8.8.

fig.8.9.

The characteristic of the saturating amplifier is depicted in fig.8.9 and yields in the shorthand notation  $\begin{array}{c|c} A & B & g \\ \hline C & D & f \end{array}$  the following coefficients for the state-model S(A,B,C,D,g,f).

	$x_1^1$	$x_2^1$	ļi		g/f
y1 1	2	-2	-2	2	0
v	-1	1	1	0	1
	1	-1	0	1	1

The integrator in  $C^2$  is defined by  $y_1^2 = u$ ,  $\dot{u} = x_1^2$ , resulting in the following array of coefficients.

Interconnection of (8.8) and (8.9) using  $x_2^1 = y_1^2$ ,  $x_1^1 = y_1^1$  and  $y_1^1 = x_1^2$  then yields after elimination of  $x_1^1$ ,  $x_2^1$  and  $x_1^2$  the following coefficient array for the overall dynamic state-model.

	u		i	g/f
y10	2	2	-2	0
у <sup>0</sup> У2	1	0	0	0
ů	2	2	~2	0
v	-1	-1	2	1
	1	2	-1	1

The above model describes an autonomous system since no input vector x is left in the model (8.10).

For a final transient analysis of PL dynamic systems described by (8.7), a relation between  $\dot{u}$  and u has to be added to the model equations. Like in the case of transient analysis of standard nonlinear networks, an integration formula is applied to eliminate the derivatives. This results in a pure algebraic PL system describing the behaviour of the dynamic system at discrete time instants.

Because the differential equations originating from descriptions of electronic networks are in general stiff equations, the applied integration formula should at least be stiffly stable. Furthermore it is of advantage to use linear multistep formulas since methods are available to implement these formulas with automatic order and stepsize control for a good compromise between accuracy, stability and efficiency [28]. Because we are dealing with piecewise-linear systems for which the first derivatives will be non-continuous, it is obvious that integration methods of order higher than two do not improve, the integration process since these formulas rely on the property that locally a sufficient accurate Taylor-series expansion of the derivatives must exist [29]. Thus maximum accuracy is obtained with second order linear multistep methods like for example the trapezoidal rule. For simplicity, without loss of generality we will use the first order backward Euler integration formula with fixed stepsize h. To this purpose let us consider the system (8.7) at time instant  $t_{n+1} = t_n + h$ , which is denoted by adding the subscript n+1 to each variable. Partitioning A, B, C and g as

$$\begin{pmatrix} A_1 & A_2 \\ -A_3 & A_4 \end{pmatrix}, \begin{pmatrix} B_1 \\ -B_2 \end{pmatrix}, (C_1 & C_2), \text{ and } \begin{pmatrix} g_1 \\ -g_2 \end{pmatrix},$$

we have from (8.7)

$$y_{n+1} = A_1 x_{n+1} + A_2 u_{n+1} + B_1 i + g_1$$
 (8.11a)

$$\dot{u}_{n+1} = A_3 x_{n+1} + A_4 u_{n+1} + B_2 i + g_2$$
 (8.11b)

$$v = C_1 x_{n+1} + C_2 u_{n+1} + Di + f$$
 (8.11c)  
 $v \ge 0, \quad i \ge 0, \quad v^{t} i = 0$ 

The backward Euler formula is given by

$$u_{n+1} = u_n + h\dot{u}_{n+1}$$
, (8.12)

which introduces a discretization error of  $O(h^2)$ . Using (8.11b) and (8.12) we find

$$u_n = -hA_3x_{n+1} + (1 - hA_4)u_{n+1} - hB_2i - hg_2$$
,

which equation is used to replace (8.11b) in the above state-model, leading to

$$y_{n+1} = A_1 x_{n+1} + A_2 u_{n+1} + B_1 i + g_1$$

$$u_n = -hA_3 x_{n+1} + (1 - hA_4) u_{n+1} - hB_2 i - hg_1$$

$$v = C_1 x_{n+1} + C_2 u_{n+1} + Di + f$$

$$v \ge 0, \quad i \ge 0, \quad v^{t}i = 0$$
(8.13)

Equation (8.13) now describes a PL mapping  $\psi$  :  $(x_{n+1}, u_{n+1}) \rightarrow (y_{n+1}, u_n)$ . Pivoting (8.13) on the full matrix  $(1 - hA_4)$ , exchanges the variables  $u_{n+1}$  and  $u_n$  yielding the final mapping  $\overline{\psi}$  :  $(x_{n+1}, u_n) \rightarrow (y_{n+1}, u_{n+1})$ , denoted by

$$\begin{pmatrix} \mathbf{y}_{n+1} \\ -\mathbf{u}_{n+1} \end{pmatrix} = \overline{A} \begin{pmatrix} \mathbf{x}_{n+1} \\ -\mathbf{u}_{n} \end{pmatrix} + \overline{B}\mathbf{i} + \overline{g}$$

$$\mathbf{v} = \overline{C} \begin{pmatrix} \mathbf{x}_{n+1} \\ -\mathbf{u}_{n} \end{pmatrix} + \overline{D}\mathbf{i} + \overline{f}$$

$$\mathbf{v} \ge 0, \quad \mathbf{i} \ge 0, \quad \mathbf{v}^{\mathsf{T}}\mathbf{i} = 0$$

$$(8.14)$$

Equation (8.14) can be seen as the overall state-model of the compound system given in fig.8.10, with the lower box representing the backward Euler formula.



fig.8.10.

From (8.14), for prescribed  $x_{n+1}$  and  $u_n$ , the response  $y_{n+1}$  and  $u_{n+1}$  can be found by using one of the previously discussed solution methods for the LCP problem. In fact eq. (8.14) represents the PL analogon of the transition matrix description of nonlinear dynamic networks, which allows us to proceed by one time-step. Systems modelled by (8.14) can again be interconnected by the methods of section 8.1, allowing for the construction of models for larger systems. Observe however that the time-step h has been fixed by the transformation from (8.7) to (8.14). For maximum flexibility it is advisable to perform the latter transformation only on the final model.

Let us now continue our example of the multivibrator given by (8.10). With h = 0.1 we obtain in analogy with (8.13)

	un+1		i	g/f
y <sub>1</sub> 0	2	2	-2	0
y20	1	0	0	0
u n	0.8	-0.2	0.2	0
v	-1	-1	2	1
-	1	2	-1	1

Pivoting on the indicated pivot then yields the final model

	u <sub>n</sub>	Ŀ	L	g/f
у <mark>0</mark>	2.5	2.5	-2.5	0
у <sup>0</sup> 2	1.25	0.25	-0.25	0
<sup>u</sup> n+1	1.25	0.25	-0.25	0
	-1.25	-1.25	2.25	1
v	1.25	2.25	-1.25	1

(8.15)

The response calculated by (8.15) using  $u_0 = 0.9$  is depicted in fig.8.11.



#### fig.8.11.

Nonlinear energy storage elements like e.g. nonlinear capacitors or inductors can easily be included in the state-model by constructing a PL approximation of their charge-voltage or fluxcurrent characteristic.

With the above tools, common electronic devices can be modelled by the state-model. In the next chapter some examples will be given to aid the implementation of these models in an analysis package and to demonstrate the general line in the modelling proces.

## 9. Piecewise-linear models for some electronic devices

In this section we present PL models for a voltage controlled switch, a MOS transistor and some digital gates. In addition we will give the results of the PL analysis of some small circuits as an indication of the applicability of the PL approach.

A voltage controlled switch can be realized by the network given in fig.9.1, with the diodes considered to be ideal.



fig.9.1.

For E > 0, both diodes will conduct as long as  $|j| \le E$  and the circuit is then equivalent to a shortcircuit. For E < 0, both diodes do not conduct as long as  $|u| \le -E$  and the circuit is an open circuit. From fig.9.1 the state-model is readily found to become

	u	Е	i <sub>1</sub>	i <sub>2</sub>	
j	0	0	1	-1	0
v <sub>1</sub>	-1	-1	1	.1	0
v2	1	-1	1	1	0

The mapping corresponding to (9.1) is given in fig.9.2 and indeed represents a voltage controlled switch as long as the electrical parameters u and j remain within the above given bounds.





The PL model for the voltage controlled switch can for example be used in modelling switched capacitor filters.

A model for a MOS transistor can be obtained as follows. We start from a symmetric model for a standard MOS transistor [30], which is depicted in fig.9.3.  $i_{p}$ 



fig.9.3.

The nonlinear resistive elements in the above model are defined by

 $i_1 = f(V_{GD}, V_T)$  and  $i_2 = f(V_{GS}, V_T)$ ,

with  $\mathbf{V}_{\mathrm{T}}$  representing the threshold voltage. In the simplest approximation f is given by

$$f(\mathbf{v}_{1},\mathbf{v}_{2}) = \kappa [\mathbf{v}_{1} - \mathbf{v}_{2}]_{+}^{2},$$
with  $[\mathbf{x}]_{+} = \begin{cases} \mathbf{x} & \text{if } \mathbf{x} \ge 0 \\ 0 & \text{if } \mathbf{x} < 0 \end{cases}$ 

$$(9.2)$$

For some  $v_2 = V_T$  this characteristic is sketched in fig.9.4 and in analogy with the Ebers-Moll model for bipolar transistors it can be seen as the voltage-current characteristic of a "MOS diode".



fig.9.4.

Then the overall characteristic of the MOS transistor can apparently be constructed from a PL approximation of the one-diomensional mapping of fig.9.4. To this purpose we consider a PL Tchebychev-like approximation  $\overline{y}(x)$  of the mapping  $y = [x]_{+}^{2}$  as indicated in fig.9.5 and fig.9.6.



fig.9.5.



fig.9.6.

For a given  $x_n$  and  $\delta$  we easily derive that the approximation in the next interval is given by

$$\bar{\mathbf{y}}(\mathbf{x}) = a_n \mathbf{x} + b_n, \quad \mathbf{x}_n \le \mathbf{x} \le \mathbf{x}_{n+1},$$

with

$$a_{n} = 2x_{n} + \varepsilon$$

$$x_{n+1} = x_{n} + \varepsilon$$

$$b_{n+1} = \delta - a_{n}^{2}/4 \text{ and } \varepsilon = \sqrt{8\delta}$$

For the first interval we find  $a_0 = 2\sqrt{\delta}$ ,  $x_1 = 1/2$   $(a_0 + \varepsilon)$ . For  $\delta = 0.16$  this yields in two digits accuracy

$$\varepsilon = 1,13$$
,  $a_0 = 0.8$ ,  $x_1 = 0.97$ ,  $a_{n+1} - a_n = 2.26$ ,  $x_{n+1} - x_n = 1.13$ 

and results in a state-model given by

	x	1 <sup>1</sup> 1	<sup>i</sup> 2	i <sub>3</sub>	1 <sub>4</sub>	<sup>i</sup> 5	Ĺб	[	
У	0.8	0.8	2.26	2.26	2.26	2.26	2.26	0	
v <sub>1</sub>	1	1						0	
v <sub>2</sub>	-1					$\bigcirc$	)	0.97	
3	-1			``、				2.10	(9.
4	-1		$\sim$	`	``.			3.23	
, 5	- 1		C	)		```		4.36	
'6	-1						1	5.49	

Eq.(9.3) is a PL approximation of  $y = [x]_{+}^{2}$  in 7 segments with an absolute error  $\delta \le 0.16$  for  $x \le 6.62$ .

Next the full MOS transistor model can be constructed by interconnection of the blocks given in fig.9.7.



fig.9.7.

For  $K = 1A/V^2$  and  $V_T = 2V$  the interconnection method of chapter 8 yields with (9.3) the following state-model (9.4)

	V <sub>GS</sub>	V <sub>DS</sub>	<sup>i</sup> 1	<sup>i</sup> 2	<sup>i</sup> 3	1. 4	i <sub>5</sub>	<sup>1</sup> 6	i <sub>7</sub>	i <sub>8</sub>	<sup>1</sup> 9	<sup>i</sup> 10	í <sub>11</sub>	i 12	
1 D	0	0.8	0.8	2.26	2.26	2.26	2.26	2.26	-0.8	-2.26	-2.26	-2.26	-2.26	-2.26	0
v <sub>1</sub>	1	0	1												-2
v <sub>2</sub>	-1	0		1											2.97
۳ <sub>3</sub>	-1	0			```.										4.10
<sup>v</sup> 4	-1	0				````									5.23
<sup>v</sup> 5	-1	o					``								6.36
¥6	-1	0						``							7.49
v <sub>7</sub>	1	~1							`						-2
<b>v</b> 8	-1	1								`.					2.97
v <sub>9</sub>	-1	1									``.				4.10
<sup>v</sup> 10	-1	1										``			5.23
۷11	-1	1											<b>`</b>		6.36
v <sub>12</sub>	-1	1											*	` 1	7.49

An indication of the accuracy of the model (9.4) is given by the table of fig.9.8 which compares the current  $i_{DPL}$  obtained from (9.4) with the value  $i_{DN}$  obtained from the MOS model using (9.2)

V <sub>GS</sub>	V <sub>DS</sub>	i <sub>DPL</sub>	i DN
2	*	0	0
3	0.5	0.48	0.75
3	1	0.88	1
3	2	0.88	1
4	0.5	1.53	1.75
4	1	3.06	3
4	1.5	3.54	3.75
4	2	3.94	4
4	3	3.94	4
5	1	5.11	5
5	2	8.17	8
5	3	9.05	9
5	4	9.05	9
6	1	7.02	7
6	2	12.2	12
6	3	15.2	15
6	4	16.1	16

#### fig.9.8.

In practice  $f(v_1, v_2)$  differs from (9.2) but can always be modelled as long as it remains a one-dimensional function of  $v_1 - v_2$ . Furthermore, the dependence of the threshold voltage  $v_2$  of the source to substrate voltage can also be modelled by a PL model and included in the above model (9.4).

Next we will derive a model for a digital nand-gate. The model is derived from the circuit given in fig.9.9, using ideal diodes and a 3-segment ideal invertor.



From this circuit we observe that  $\boldsymbol{x}_3$  is defined by

$$\mathbf{x}_3 = \min(\mathbf{x}_1, \mathbf{x}_2)$$

which can be modelled by the mapping depicted in fig.9.10.



fig.9.10.

The state-model of this mapping becomes

	×1	×2	i <sub>1</sub>
×3	0	1	-1
v <sub>1</sub>	1	-1	1

Furthermore, the invertor is modelled by

	× <sub>3</sub>	i <sub>2</sub>	i <sub>3</sub>	
У	-50	-50	50	30
v <sub>2</sub>	1	1	0	-0.5
v <sub>3</sub>	- 1	0	1	0.6

The overall model for the circuit of fig.9.9 using (9.5) and (9.6) then becomes

	×1	×2	1 1	<sup>1</sup> 2	i <sub>3</sub>	
У	0	-50	50	-50	50	30
v <sub>1</sub>	1	- 1	1	0	0	0
v <sub>2</sub>	0	1	- 1	1	0	-0.5
v <sub>3</sub>	0	- 1	1	0	1	0.6

The corresponding mapping is depicted in fig.9.11. Note the nonconvex region in which y = 5, which is due to an unobservable state coupled to the "boundary"  $x_1 = x_2$ .



fig.9.11.

Next let us consider a model for the so-called threshold gates which can be used to realize various types of digital circuit blocks like e.g. flip-flops or ring-counters and lend themselves particularly well for a PL modelling. These gates are symbolized as given in fig.9.12 and their output is defined by

$$y = \begin{cases} 1 & \text{if } a_1 x_1 + a_2 x_2 + \dots + a_n x_n \ge t_1 \\ 0 & \text{if } a_1 x_1 + a_2 x_2 + \dots + a_n x_n \le t_2 \end{cases}$$
(9.8)



fig.9.12.

From (9.8) the general state-model is easily derived to become

	×1	<b>x</b> <sub>2</sub> · · ·	*n	i <sub>1</sub>	i <sub>2</sub>		
v	a_1	a <sub>2</sub>	n	1	-1	2	
*	t <sub>1</sub> -t <sub>2</sub>	t <sub>1</sub> -t <sub>2</sub>	$t_1 - t_2$	t <sub>1</sub> -t <sub>2</sub>	$t_1 - t_2$	t <sub>1</sub> -t <sub>2</sub>	(9,9)
v <sub>1</sub>	a 1	a <sub>2</sub>	an	1	0	-t <sub>2</sub>	(,
v <sub>2</sub>	-a <sub>1</sub>	-a <sub>2</sub>	-a n	0	1	t <sub>1</sub>	

Using (9.9) a set-reset flip-flop could for example be modelled from the circuit given in fig.9.13.



fig.9.13.

We will conclude by presenting the results obtained by calculating the response of two small circuits which have been modelled by the state-model approach. The first circuit is depicted in fig.9.14 and forms a Wien-bridge oscillator, stabilized by an amplitude controlloop using a MOS transistor.



fig.9.14.

Transient analysis of the above circuit produced the response given in fig.9.15. As can be seen the amplitude becomes stationary in about 4 cycles.



The second example is given by a PL van der Pol oscillator as depicted in fig.9.16.



fig.9.16.

The nonlinear resistive element shows a negative resistance in the range  $|v| \leq 1$ . The response is given in fig.9.17 and shows a highly nonlinear current flow through the PL resistor, producing a large third harmonic current component.



fig.9.17.

More complex circuits like dynamic 4-bit shift registers and phase-locked loops have also been analysed and gave a response in agreement with the theory. Due to the non-sparse implementation, the complexity of the circuits which could be modelled was limited such that the overall model had to fit into a 42×42 matrix. A sparse implementation will eliminate this constraint and an adequate package based on the state-model approach can deliver the designer of electronic circuits a powerful tool for a mixed level analysis of complex combined analog-digital networks. Appendix I

Consider a hybrid structure for the matrix  $\widetilde{H} = H|I$ , derived from H. We are going to map this structure on the structure of H by mapping the vertices of  $\widetilde{H}$  onto those of H according to

$$\psi: \det((H|I)_{J}) \rightarrow \det(H_{I-J}), \tag{A1}$$
with I-J = (IUJ) \ (IOJ).

Then we have:

For any edge in the structure of  $\widetilde{H}$  there is a corresponding edge in the structure of H with the same weight.

<u>Proof</u>: Consider the vertices  $\det(\widetilde{H}_{J})$  and  $\det(\widetilde{H}_{J-\{k\}})$  in the structure of  $\widetilde{H}$ . The mapping  $\psi$  from (A1) maps these vertices on  $\det(H_{I-J})$  and  $\det(H_{I-J-\{k\}})$  of the structure of H (See fig.A1).



fig.A1.

Hence adjacent vertices remain adjacent under this mapping. Next consider the weights  $\tilde{\theta}$  and  $\theta$  in fig.Al. Using (3.12) we have

$$\widetilde{\theta} = \frac{\det(\widetilde{H}_{J-\{k\}})}{\det(\widetilde{H}_{J})} = (\widetilde{H}|J)_{kk}$$
(A2)

However  $(\widetilde{H}|J)_{kk} = ((H|I)|J)_{kk} = (H|I-J)_{kk}$ , due to the definition of (H|I). Furthermore, we have

$$\theta = \frac{\det(H_{I-J-\{k\}})}{\det(H_{I-J})} = (H|I-J)_{kk}$$

and thus  $\tilde{\theta} = \theta$ . End of proof.

Hence the structures of  $\widetilde{H}$  and H are isomorphic and the weights on the edges are invariant under the mapping  $\psi$ . Next we consider a path  $\widetilde{P}$  in the structure of  $\widetilde{H}$  from  $\widetilde{H}_{\phi}$  to  $\widetilde{H}_{J}$ . Then holds

$$\frac{\det(\widetilde{H}_{J})}{\det(\widetilde{H}_{\phi})} = \frac{\det((H|I)_{J})}{\det((H|I)_{\phi})} = \Pi \quad (all weights on the edges of the path \widetilde{P})$$
(A3)

From the corresponding path P in the structure of H and the invariancy of the weights we then obtain using (A1) and (A3)

$$\frac{\det(\widetilde{H}_{J})}{\det(\widetilde{H})} = \Pi \quad (all weights on the edges of the path P) = = \frac{\det(H_{I-J})}{\det(H_{I})}$$
(A4)

Because by definition  $det(\widetilde{H}_{\phi}) = 1$ , equ.(A4) yields

$$\det((\mathbf{H}|\mathbf{I})_{\mathbf{J}}) = \det(\widetilde{\mathbf{H}}_{\mathbf{J}}) = \frac{\det(\mathbf{H}_{\mathbf{I}-\mathbf{J}})}{\det(\mathbf{H}_{\mathbf{T}})}$$

# Appendix II

Let us assume that the reduced problem in the j<sup>th</sup> cycle is obtained from M by exchanging the variables  $i_{l_1}, i_{l_2}, \dots, i_{l_k}$  with the corresponding  $v_{l_1}, v_{l_2}, \dots, v_{l_k}$  (k≤j). By a reordering we may arrange that  $i_{l_1}$  to  $i_{l_k}$  form the first k components of i. Next partition the reordered equations as

$$\begin{pmatrix} \mathbf{v}_1 \\ \cdots \\ \mathbf{v}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{M}_1 & \cdots & \mathbf{M}_2 \\ \cdots & \cdots & \cdots \\ \mathbf{M}_3 & \cdots & \mathbf{M}_4 \end{pmatrix} = \begin{pmatrix} \mathbf{i}_1 \\ \cdots \\ \mathbf{i}_2 \end{pmatrix} + \begin{pmatrix} \mathbf{a}_1 \\ \cdots \\ \mathbf{a}_2 \end{pmatrix}, \quad \mathbf{M}_1 = (\mathbf{k} \times \mathbf{k}), \quad \mathbf{M}_2 = \mathbf{M}_3^{\mathsf{t}} \quad (B1)$$

From (B1) we derive the equivalent formulation

$$\begin{pmatrix} i_1 \\ ... \\ v_2 \end{pmatrix} = \begin{pmatrix} M_1^{-1} & ... & ... \\ ... & ... \\ M_3 M_1^{-1} & ... & M_4^{-M_3} M_1^{-1} M_2 \end{pmatrix} \begin{pmatrix} v_1 \\ ... \\ i_2 \end{pmatrix} + \begin{pmatrix} -M_1^{-1} a_1 \\ ... \\ a_2^{-M_3} M_1^{-1} a_1 \end{pmatrix}$$
(B2)

In the sequel we use the notation  $A \subset B$  to denote that A is a principal submatrix of B.

The reduced matrix  $M_{j}$  in the j<sup>th</sup> cycle is now a principal submatrix of  $\overline{M} = M_{4} - M_{3}M_{1}^{-1}M_{2}$  since at this state  $v_{1} = 0$  and (j-k) components of  $i_{2}$  are zero, hence

$$M_{j} \subset \overline{M} = M_{4} - M_{3}M_{1}^{-1}M_{2} = \overline{M}^{t}$$
 (B3)

Next, from (B1) consider the matrix  $M^{-1}$ , which is easily found to become

Note that  $M^{-1}$ ,  $M_1^{-1}$  and  $M_4^{-1}$  exist since M is positive definte.

Then from (B3) and (B4) we have

$$M_j \subset \overline{M} = \widetilde{M}_4^{-1}$$
, with  $\widetilde{M}_4 \subset M^{-1}$  (B5)

Now, due to (B4),  $\widetilde{M}_4$ ,  $\widetilde{M}_4^{-1}$  and  $M_j$  are symmetric and positive definite since  $M^{-1}$  is symmetric and positive definite. Therefore all eigenvalues of  $M_j$ ,  $\widetilde{M}_4^{-1}$  and  $\widetilde{M}_4$  are positive real. For a symmetric positive definite matrix A let us denote the samlles eigenvalue by  $\sigma(A)$  and the largest eigenvalue by  $\rho(A)$ . Furthermore the interval  $[\sigma(A), \rho(A)]$  is denoted by R(A). Then since  $\widetilde{M}_4 \subset M^{-1}$ , by the Sturmian separation theorem [31], we have

$$\mathbb{R}(\widetilde{M}_4) \subset \mathbb{R}(M^{-1}), \text{ i.e. } \sigma(\widetilde{M}_4) \geq \sigma(M^{-1}), \rho(\widetilde{M}_4) \leq \rho(M^{-1})$$
 (B6)

Then for  $\overline{M} = \widetilde{M}_{4}^{-1}$  we have, due to (B6)

$$R(\overline{M}) = \left[\frac{1}{\rho(\widetilde{M}_{4})}, \frac{1}{\sigma(\widetilde{M}_{4})}\right] \subset \left[\frac{1}{\rho(M^{-1})}, \frac{1}{\sigma(M^{-1})}\right] \equiv \left[\sigma(M), \rho(M)\right]$$
(B7)

or  $R(\overline{M}) \subset R(M)$ 

By the same Sutrm theorem, we have from (Bt) and (B7)

$$R(\underline{M}_{i}) \subset R(\overline{\underline{M}}) \subset R(\underline{M})$$
(B8)

Hence the eigenvalues  $\lambda_{\substack{i}}(M_{\substack{j}})$  of the matrix  $M_{\substack{j}}$  from the  $j^{\textstyle \text{th}}$  cycle satisfy

$$\sigma(\mathbf{M}) \leq \lambda_{\mathbf{i}}(\mathbf{M}_{\mathbf{j}}) \leq \rho(\mathbf{M})$$
(B9)

Then the corresponding eigenvalues of  $\mu_i$  of  $D_j = (1+M_j)^{-1}(1-M_j)$  satisfy

$$\mu_{i}(D_{j}) = \frac{1 - \lambda_{i}(M_{j})}{1 + \lambda_{i}(M_{j})}$$
(B10)

hence

$$|\mu_{i}(D_{j})| \leq \max\left(\left|\frac{\sigma(M)-1}{\sigma(M)+1}\right|, \left|\frac{\rho(M)-1}{\rho(M)+1}\right|\right) = \theta_{M} < 1$$
 (B11)

Then since  $D_{j}$  is symmetric, we have

$$\begin{split} & ||D_{j}||_{2} = \max_{i} |\mu(D_{j})| \leq \theta_{M} < 1 \\ & \text{or } \theta(M_{j}) < 1 \end{split}$$

When the matrix M is scaled in advance by a factor  $\alpha$ , we have from (B11)

$$\sigma(M_j) \leq \max\left( \left| \frac{\alpha \sigma(M) - 1}{\alpha \sigma(M) + 1} \right| , \left| \frac{\alpha \rho(M) - 1}{\alpha \rho(M) + 1} \right| \right)$$

For  $\alpha = (\sigma(M)\rho(M))^{-1/2}$  this yields

$$\theta(M_{j}) \leq \theta_{M} = \frac{\rho(M)}{\rho(M)} \frac{1/2}{1/2} - \sigma(M) \frac{1/2}{1/2} = \frac{1-d}{1+d}$$

with d =  $\sqrt{\frac{\sigma(M)}{\rho(M)}}$  = K<sup>-1/2</sup> and K being the condition number of M. Then  $1/\ln(1/\theta_M) = \frac{1}{\ln\frac{1+d}{1-d}} \le \frac{1}{2d} = \frac{1}{2}K^{1/2}$ .

The values of  $\sigma\left(M\right)$  and  $\rho\left(M\right)$  can also be bound by

$$\rho(\mathbf{M}) \leq ||\mathbf{M}||_{\infty} \text{ and } \sigma(\mathbf{M}) \geq \frac{1}{||\mathbf{M}^{-1}||_{\infty}}$$
(B12)

From (B9), (B10) and (B12) it is then easily found that

$$\left\| D_{\mathbf{j}} \right\|_{2} \leq \overline{\theta} = \max \left( \left\| \frac{\left\| \mathbf{M}^{-1} \right\|_{\infty} - 1}{\left\| \mathbf{M}^{-1} \right\|_{\infty} + 1} \right|, \left\| \frac{\left\| \mathbf{M} \right\|_{\infty} - 1}{\left\| \mathbf{M} \right\|_{\infty} + 1} \right) \right) < 1$$

which yields an easy computable bound for all  $\left\|\mathbf{D}_{j}\right\|_{2}$  .

References

- [ 1] J.M. Orthega and W.C. Rheinboldt, Iterative solution of nonlinear equations in several variables, New York, Academic Press, 1970.
- [2] I.W Sandberg and A.N. Willson, "Some theorems on properties of DC equations of nonlinear networks", Bell Syst. Tech. J., Vol. 48, January 1969, pp. 1-34.
- [ 3] P.M. Lin, "Formulation of hybrid matrices for linear multiports containing controlled sources", IEEE Trans. on Circ. and Syst., Vol. CAS-21, March 1974, pp. 169-175.
- [4] J. Katzenelson, "An algorithm for solving nonlinear resistive networks", Bell Syst. Tech. J., Vol. 44, October 1965, pp. 1605-1620.
- [5] W.M.G. van Bokhoven and J.A.G. Jess, "Some new aspects of P and P<sub>0</sub> matrices and their application to networks with ideal diodes", Proc. IEEE ISCAS, 1978, pp. 806-810.
- [ 6] G. Verkroost, "On the existence and uniqueness of solutions for networks containing linear memoryless elements and ideal diodes", Proc. of IEE Europ. Conf. on Circ. Th., IEE 116, 1974.
- [7] C.E. Lemke, "On complementary pivot theory", in: Mathematics of the decision sciences, Part I, Eds. G.B Dantzig and A.F. Veinott, Jr., New York, Academic Press, 1970.
- [8] M. Fiedler and V. Pták, "Some generalisations of positive definiteness and monotonicity", Num. Math., Vol. 9, December 1966, pp. 163-172.
- [ 9] V. Belevitch, "The number of states of rectifier networks", IRE Trans. on Circ. Th., vol. CT-9, 1962, pp. 93-94.
- [10] H. Samuelson et al., "A partition theorem for Euclidean n-space", Proc. Amer. Math. Soc., Vol. 9, 1958, pp. 805-807.
- [11] Complementarity and fixed point problems. Mathematical Programming Study 7, Eds. M.L. Balinski and R.W. Cottle, North-Holland Publ. Comp., Amsterdam, 1978.
- [12] L.T. Watson, "Solving the nonlinear complementarity problem by a homotopy method", SIAM J. of Contr. and Opt., Vol. 17, No. 1, 1979, pp. 36-46.

- [13] C.W. Cryer, "The solution of a quadratic programming problem using systematic overrelaxation", SIAM J. of Contr., Vol. 9, No. 3, 1971, p. 385.
- [14] M.L. Fisher and F.J. Gould, "A simplicial algorithm for the nonlinear complementarity problem", Math. Progr., Vol. 6, 1974, pp. 281-300.
- [15] T. Fujisawa and E.S. Kuh, "Piecewise-linear theory of nonlinear networks", SIAM J. Appl. Math., Vol. 22, No. 2, March 1972.
- [16] M.J. Chien and E.S. Kuh, "Solving piecewise-linear equations for resistive networks", Int. Jrn. Circ. Th. and Appl., Vol. 4, 1976, pp. 3-24.
- [17] R.W. Cottle and G.B. Dantzig, "Complementary pivot theory of mathematical programming", Lin. Algebra and its Appl., Vol. 1, 1968, pp. 103-125.
- [18] S. Karamardian, "The complementarity problem", Math. Progr., Vol. 2, 1972, pp. 107-129.
- [19] Y. Fathi, "Computational complexity of LCPs associated with positive definite symmetric matrices", Math. Progr. Vol. 17, 1979, pp. 335-344.
- [20] K.G. Murty, "Computational complexity of complementary pivot methods", in: Mathematical Programming Study 7, Eds. M.L. Balinski and R.W. Cottle, North-Holland Publ. Comp., Amsterdam, 1978.
- [21] S.M. Kang and L.O. Chua, "A global representation of multidimensional piecewise-linear functions with linear partitions", IEEE Trans. on Circ. and Syst., Vol. CAS-25, November 1978, pp. 938-940.
- [22] O.L. Mangasarian, "Equivalence of the complementarity problem to a system of nonlinear equations", SIAM J. Appl. Math., Vol. 31, No. 1, July 1976, pp. 89-91.
- [23] N. Megiddo and M. Kojima, "On the existence and uniqueness of solutions in nonlinear complementarity theory", Math. Progr., Vol. 12, 1977, pp. 110-130.
- [24] L.B. Rall, Computational solution of nonlinear operator equations,J. Wiley & Sons Inc., New York, 1969, p. 65.
- [25] M.J. Chien, "Piecewise-linear theory and computation of solutions of homeomorphic resistive networks", IEEE Trans. on Circ. and Syst., Vol. CAS-24, No. 3, March 1977, pp. 118-127.
- [26] S.N. Tschernikow, Lineare Ungleichungen, Hochschulbücher für Mathematik, VEB Verlag der Wissenschaften, 1971.
- [27] T. Ohtsuki, T, Fujisawa and S. Kumagai, "Existence theorems and a solution algorithm for piecewise-linear resistor networks", SIAM J. Math. Anal., Vol. 8, No. 1, February 1977, pp. 69-99.
- [28] W.M.G. van Bokhoven, "Linear implicit differentiation formulas of variable step and order", IEEE Trans. on Circ. and Syst., Vol. CAS-22, No. 2, February 1975, pp. 109-115.
- [29] P. Henrici, Discrete variable methods in ordinary differential equations, J. Wiley & Sons Inc., New York, 1962.
- [30] J.E. Meyer, "MOS models and circuit simulation", RCA Rev., Vol. 32, March 1971, pp. 42-63.
- [31] R. Bellman, Introduction to matrix analysis, McGraw-Hill Book Comp. Inc., New York, 1960, pp. 115-116.

### Samenvatting

In de elektrotechniek moet de simulatie van complexe elektronische schakelingen als een vitaal onderdeel van het huidige ontwerpproces van LSI schakelingen beschouwd worden. Vanwege onvermijdelijke effecten,zoals strooi capaciteiten of parasitaire transistoren,die in een geïntegreerde schakeling optreden is het niet langer mogelijk of aanbevelenswaardig om de schakeling zodanig uit discrete componenten op te bouwen dat het ontwerp door middel van metingen volledig op zijn correctheid getoetst kan worden.

Deze parasitaire effecten kunnen echter in theorie vrij eenvoudig in een simulatie worden meegenomen waardoor een vroegtijdige indicatie betreffende de correctheid van het ontwerp verkregen kan worden voordat de schakeling daadwerkelijk in silicium geintegreerd wordt.

Alhoewel simulatie dus het antwoord lijkt op een groot aantal problemen,worden de moeilijkheden in feite slechts doorgeschoven naar deze simulatie,die dan zeer complex zal worden. Naast logische simulatie van digitale circuits zal de simulatie in staat moeten zijn ten behoeve van een correcte timing het tijdsverloop van de diverse logische signalen te berekenen. Tevens zal de responsie van analoge circuits ten aanzien van inschakel verschijnselen en gelijkspanningsgedrag moeten kunnen worden bepaald.

Omdat in de huidige toepassingen veelal analoge circuits in combinatie met digitale circuits op één chip geIntegreerd worden, moeten traditioneel verschillende simulatie taken nu op een enkel circuit uitgevoerd worden.Dit probleem stelt stringente eisen aan de onderdelen en de interne organisatie van een simulatie programma. Een hierarchische structuur van zo'n programma is een eerste vereiste. Verder dienen de diverse taken binnen een simulatie programma dusdanig op elkaar afgestemd te worden dat een zo uniform mogelijke datastructuur toegepast kan worden teneinde conversie problemen bij de vertaling van gegevens benodigd voor de verschillend gerichte taken te voorkomen.

Een van de belangrijkste onderdelen van een simulatie programma is te vinden in het oplossen van stelsels niet-lineaire vergelijkingen welke,door toepassing van geschikte modellen voor de aanwezige componenten,ontstaan bij het formuleren van de netwerkvergelijkingen.

Voor schakelingen met bipolaire of MOS transistoren kan het gelijkspanningsgedrag van deze componenten door middel van een relatief klein stelsel niet-lineaire vergelijkingen vastgelegd worden. Het stelsel van niet-lineaire vergelijkingen dat het gedrag van de totale schakeling beschrijft wordt dan meestal opgelost door toepassing van een zekere variant van de Newton-Raphson methode.

Het gedrag van digitale schakelingen wordt daarentegen vaak vastgelegd door middel van tabellen of logische functies en vereist in het algemeen een heel andere aanpak,zowel ten aanzien van de modellering als van de analyse methodes. Dit bemoeilijkt de onderlinge uitwisselbaarheid van de berekende responsies. Door echter gebruik te maken van modellen gebaseerd op de zogenaamde drempel-logica ontstaat een beschrijving met stuksgewijslineaire vergelijkingen welke beter bij de voorafgaande beschrijving aansluit.

In complexe analoge schakelingen ten slotte wordt het gedrag van bijvoorbeeld operationele versterkers, hysterese schakelingen e.d. vaak weergegeven door een zogenaamd macro-model dat te beschrijven is met behulp van stuksgewijs-lineaire karakteristieken.

Ieder van bovenstaande modelleringstechnieken leidt tot een specifieke analyse methode met de daaraan verbonden eigen problematiek zoals convergentie van de oplossings algoritme en complexiteit van de datastructuur.

Het zal duidelijk zijn dat de eerder gewenste uniforme aanpak van de simulatie alleen bereikt kan worden indien de modellering op ieder analyse-niveau van dezelfde basis structuur uitgaat.

In dit proefschrift is gekozen voor een aanpak op basis van stuksgewijs-lineaire afbeeldingen. Enerzijds omdat de meer complexe circuits zich gemakkelijker lenen voor een stuksgewijs-

lineaire beschrijving dan voor een beschrijving door middel van algemene niet-lineaire functies en anderzijds omdat het zeker eenvoudiger is om niet-lineaire functies te benaderen door stuksgewijs-lineaire functies dan omgekeerd. Vooral de relatief eenvoudige wijze van macro-modellering van complexe bouwstenen opent aantrekkelijke mogelijkheden voor de simulatie van grote electronische schakelingen. Te meer daar deze macro-modellen gecombineerd kunnen worden met modellen voor afzonderlijke transistoren in termen van stuksgewijs-lineaire afbeeldingen,zodat een analyse van deze schakelingen op gemengd niveau kan worden uitgevoerd op basis van een gemeenschappelijke datastructuur.

De traditionele beschrijving van stuksgewijs-lineaire afbeeldingen is echter behept met een aantal nadelen. Allereerst worden deze afbeeldingen meestal vastgelegd door middel van een lijst waarin zowel de diverse gebieden (polveders) in een Euclidische ruimte als de bijbehorende lineaire afbeeldingen aangegeven zijn. Dit leidt tot een redundante formulering van de afbeelding omdat begrenzende hypervlakken minstens twee maal in de lijst voorkomen. Daarnaast kunnen burenrelaties tussen de verschillende gebieden niet op een impliciete wijze worden vastgelegd, waardoor de gangbare methodes voor het oplossen van stuksgewijs-lineaire vergelijkingen tot een inefficiënt zoekproces in de betreffende datastructuur leiden. Tevens is een lineaire transformatie op de variabelen van een stuksgewijs-lineaire afbeelding moeilijk te verwerken in de desbetreffende lijst en leidt een sequentie van stuksgewijslineaire afbeeldingen zelfs tot een exponentiële groei van deze lijst.

De basis van dit proefschrift wordt gevormd door een nieuwe beschrijvingsmethode voor stuksgewijs-lineaire afbeeldingen die de meeste van de voornoemde bezwaren wegneemt. Daartoe wordt een toestands-model voor een stuksgewijs-lineaire afbeelding geintroduceerd. Dit toestands-model levert een globale impliciete beschrijving van de afbeelding die voor alle gebieden geldig is.Het desbetreffende model is opgebouwd te denken

uit twee verschillende onderdelen. Het ene deel beschrijft de lineaire afbeelding voor elk gebied terwijl het andere deel (de zogenaamde toestands-vergelijking) verantwoordelijk is voor de opdeling van het definitie gebied van de afbeelding in afzonderlijke polyeders. Het lineaire complementariteits-probleem vormt de essentiële grondslag van het model,zodat bekende algoritmen voor de oplossing van dit probleem nu ook gebruikt kunnen worden om stuksgewijs-lineaire vergelijkingen op te lossen.

Naast de behandelde toepassing in de elektrotechniek kan het toestands-model ook met vrucht gebruikt worden in de theorie van homeomorfismen van stuksgewijs-lineaire afbeeldingen.

In de diverse hoofdstukken van dit proefschrift komen de volgende zaken aan de orde.

Hoofdstuk 2 behandelt de beschrijving van stuksgewijslineaire netwerken met ideale diodes. Deze netwerkbeschrijving wordt later gebruikt om het toestands-model te definieren.

In hoofdstuk 3 worden, aan de hand van de mogelijke hybriede beschrijvingen van een lineair geheugenloos multipoort netwerk, enkele theorema's geformuleerd die de equivalentie aantonen van de eigenschappen van klasse P en P<sub>0</sub> matrices en zekere passiviteits eigenschappen van voornoemde netwerken. Verder wordt een graaf gedefinieerd die de "structuur" van het netwerk genoemd wordt, waarmee burenrelaties tussen diverse representaties vastgelegd worden.

Hoofdstuk 4 behandelt het toestands-model van een stuksgewijslineaire afbeelding. De begrippen "toestand" en "equivalente representatie" worden uiteen gezet en gerelateerd met de in hoofdstuk 3 behandelde hybriede representaties. Tevens wordt een minimaal toestands-model gedefinieerd waarvan later in hoofdstuk 7 gebruik gemaakt wordt bij de afleiding van homeomorfisme eigenschappen van stuksgewijs-lineaire afbeeldingen.

In hoofstuk 5 worden een drietal bekende algoritmen voor de oplossing van het lineaire complementariteits-probleem beschreven. Deze algoritmen kunnen worden toegepast om stuksgewijslineaire vergelijkingen op te lossen.

Hoofstuk 6 beschrijft een nieuwe algoritme voor de oplos-

sing van het lineaire complementariteits-probleem. Door de zogenaamde "modulus transformatie" wordt dit probleem eerst omgezet in een stelsel niet-lineaire vergelijkingen welke vervolgens door contractie afbeelding kunnen worden opgelost. Voor een bepaalde klasse van problemen leidt dit, in tegenstelling met de algoritmen uit hoofdstuk 5, tot een polynome tijdcomplexiteit van de algoritme. Tevens wordt deze methode in verband gebracht met het "globale" model van Chua.

Hoofdstuk 7 behandelt het toestands-model van de inverse afbeelding. Er worden een aantal theorema's over homeomorfismen van stuksgewijs-lineaire afbeeldingen afgeleid,die tot een uitbreiding van de bestaande theorie leiden. Daarnaast worden enkele eenvoudige expliciete toestands-modellen bepaald.

In hoofdstuk 8 komt de constructie van het toestands-model voor stuksgewijs-lineaire systemen aan de orde. De modellering geschiedt door de modellen voor de diverse componenten te combineren aan de hand van een specificatie van de onderlinge connecties. Tevens wordt het toestands-model uitgebreid tot een model voor stuksgewijs-lineaire dynamische systemen door gebruik te maken van een lineaire meerstaps integratie formule.

Tenslotte worden in hoofdstuk 9 enige voorbeelden van toestands-modellen gegeven voor enkele gangbare elektronische componenten. Aan de hand van een Wien-brug oscillator worden de perspectieven van een stuksgewijs-lineaire modellering gedemonstreerd.

# Curriculum Vitae

De schrijver van dit proefschrift werd op 9 juli 1942 geboren te Drunen. In 1959 behaalde hij het HBS-b diploma aan het Dr. Moller College te Waalwijk. Van 1962 tot 1966 was hij als kort-verband officier verbonden aan de Koninklijke Luchtmacht in het dienstvak Bewapening en Elektronica. Van november 1966 tot oktober 1970 was hij als werkstudent en als Technisch Ambtenaar werkzaam op de Technische Hogeschool te Eindhoven. In september 1970 behaalde hij aldaar met lof het diploma Elektrotechnisch Ingenieur. Van oktober 1970 tot februari 1972 was hij als medewerker verbonden aan het Natuurkundig Laboratorium van de N.V. Philips, waar hij onderzoek verrichtte aan niet-lineaire oscillatoren. Sinds februari 1972 is hij als Wetenschappelijk Medewerker verbonden aan de vakgroep Automatisch Systeem Ontwerpen in de afdeling der Elektrotechniek van de Technische Hogeschool te Eindhoven.

## STELLINGEN BIJ HET PROEFSCHRIFT VAN

W.M.G. van BOKHOVEN.

26-5-1981

De uitspraak van D.Meyer-Ebrecht "Jedoch wurde bisher keine allgemeine Theorie der amplituden-geregelten harmonischen Oszillator aufgestellt" is onwaar.

Meyer-Ebrecht, D.

Schnelle Amplitudenregelung Harmonischer Oszillatoren. Diss. T.U. Braunschweig, 1974, p.2

II

Bij de behandeling van de gesynchroniseerde Van der Pol oscillator heeft Nayfeh de invloed van de amplituderegellus niet correct verwerkt. De door hem aangegeven voorwaarde a > i voor stabiele synchronisatie bij snelle regeling is daardoor onjuist en dient vervangen te worden door a >  $\sqrt{2}$ .

Nayfeh, A.H.

Forced Oscillations of the Van der Pol Oscillator with delayed Amplitude Limiting.

IEEE Trans. on Circ.Th., CT-15, 1968, pp. 192-200

### III

Volgens de gevoeligheidsanalyse van 2e orde aktieve RC-filters met één positieve versterker zoals gepubliceerd door Weyten,zou de belangrijke door Fleischer behandelde klasse van aktieve filters niet kunnen functioneren. Dit is onjuist en in tegenspraak met de praktijk. Weyten,L.

Q- and  $\omega_0$ -Sensitivities in Positive Gain Second-order RC Active Filters. Proc.IEEE,vol.60,1972,pp.1462-1463

## IV

Een minimum-fase spanningsoverdrachtsfunctie waarvan de Bode amplitudeplot slechts uit segmenten bestaat met hellingen van 0, +6 of -6 dB/oct., kan gerealiseerd worden met behulp van klasse-M immittantie functies. Bokhoven,W.M.G. van

A Canonical Synthesis for Immittance Functions that can be decomposed into the Product of an RL and an RC Impedance. IEEE Trans. on Circ. and Syst., CAS-23, 1976, pp. 36-39. De wijdverbreide opvatting dat de spanningsoverdrachtsfunctie van 2e orde aktieve RC-filters met één versterker bij negatieve terugkoppeling ongevoeliger zou zijn voor variaties in versterking dan bij positieve terugkoppeling, is onjuist.

Bokhoven, W.M.G. van

Sensitivity Bounds and Explicit Polynomial Decomposition for Biquadratic Single Amplifier Active Filters.

Proc. 1976 IEEE Int. Symp. on Circ. and Syst., pp.118-121

VI

De directe toepassing van een aantal Newton-Cotes integratieformules bij het oplossen van gewone differentiaalvergelijkingen leidt tot een instabiele oplossingsmethode. De op deze formules gebaseerde impliciete Runge-Kutta procedures leiden daarentegen tot A-stabiele integratieformules.

Bokhoven, W.M.G. van

Efficient higher order implicit one-step methods for integration of stiff differential equations.

BIT 20,1980,pp.34-43.

#### VII

Door toepassing van nieuwe op systeemtheorie gebaseerde CAD hulpmiddelen zal de rol van de soldeerbout in het ontwerpproces van elektronische schakelingen voor een groot deel door een lichtpen worden overgenomen.

## VIII

De wijziging van het verenigingsrecht en de recente plannen van de regering om bestuursleden van verenigingen hoofdelijk aansprakelijk te stellen voor financiële tekorten in de exploitatie, brengt het voortbestaan van vele amateur-sportverenigingen in gevaar doordat op de bereidheid van de leden om tot dan toe zonder vergoeding diensten te verrichten een te grote aanslag gepleegd wordt.