

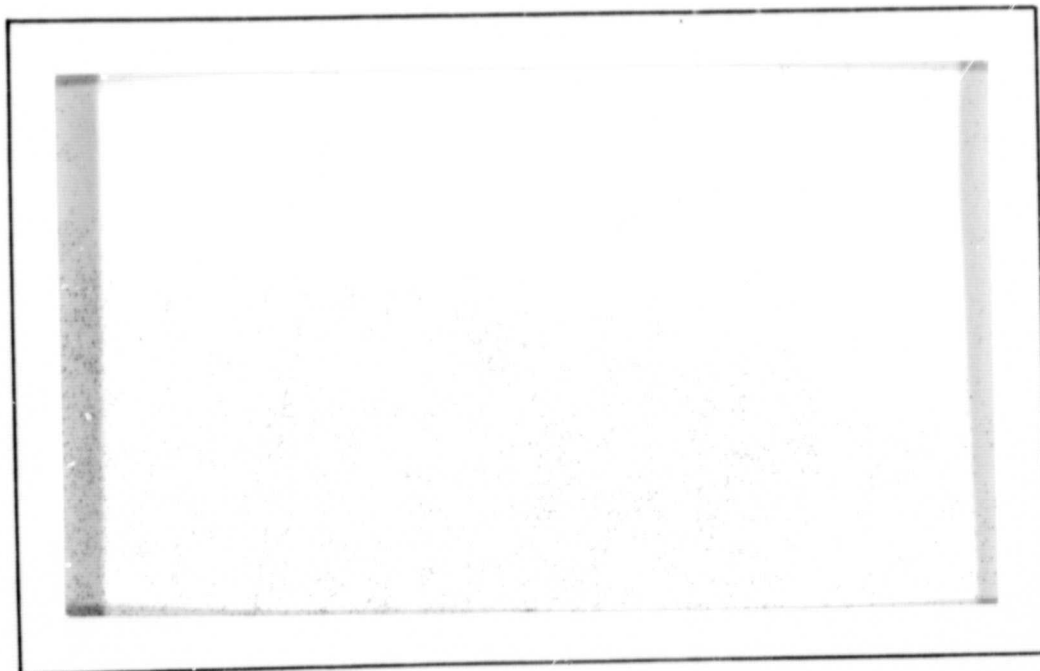
General Disclaimer

One or more of the Following Statements may affect this Document

- This document has been reproduced from the best copy furnished by the organizational source. It is being released in the interest of making available as much information as possible.
- This document may contain data, which exceeds the sheet parameters. It was furnished in this condition by the organizational source and is the best copy available.
- This document may contain tone-on-tone or color graphs, charts and/or pictures, which have been reproduced in black and white.
- This document is paginated as submitted by the original source.
- Portions of this document are not fully legible due to the historical nature of some of the material. However, it is the best reproduction available from the original submission.

N70-20906

(ACCESSION NUMBER)	(THRU)
226	1
(PAGES)	(CODE)
QF109079	10
(NASA CR OR TMX OR AD NUMBER)	(CATEGORY)



CENTER FOR
SYSTEM SCIENCE



UNIVERSITY OF ROCHESTER
ROCHESTER, NEW YORK



THE INVERSE PROBLEM
OF THE OPTIMAL REGULATOR

Ryuzo Yokoyama

December, 1969

CS 8-69-12

A Thesis Submitted
in Partial Fulfillment
of the
Requirements for the Degree
DOCTOR OF PHILOSOPHY

Supervised by Professor Edwin Kinnen
Department of Electrical Engineering
The University of Rochester
Rochester, New York

VITA

Ryuzo Yokoyama was born in [REDACTED], on [REDACTED]

[REDACTED]. During the period of 1960-1966, he received undergraduate and graduate training at Tohoku University in Japan and was awarded the Bachelor of Engineering degree and the Master of Engineering degree with a major in Electrical Engineering in 1964 and 1966 respectively.

He entered the University of Rochester in the fall of 1966 to begin a Ph.D. program in Electrical Engineering. He has been assisted in this program by the award of a graduate teaching assistantship in the Department of Electrical Engineering during the period of September 1966 to May 1967, and the award of graduate research assistantships thereafter from the National Aeronautics and Space Administration and the Office of Naval Research, Center for Naval Analyses.

ACKNOWLEDGEMENT

The author expresses his appreciation to Dr. Edwin Kinnen for his guidance and encouragement during the course of this research and to Dr. J.H.B. Kemperman for his helpful advice and assistance.

This work was performed partially with the financial support of the National Aeronautics and Space Administration under grant ^{NCR-}~~NSG-574~~/33-019-014 to the University of Rochester. It was also supported in part by The Center for Naval Analyses of the University of Rochester. This support does not imply endorsement of the content by the Navy.

Finally, I should like to thank Miss Gail Herbst for typing a very difficult thesis.

ABSTRACT

A new phase canonical form is given for a class of multi-input dynamical systems described by time invariant ordinary differential equations. This is based on a modified definition of an equivalent relation for the class of systems. It is shown that a characteristic quantity called a stage distribution defined with respect to the linear part of the system uniquely determines the structure of its canonical form.

The inverse problem of the optimal regulator is considered for this class of systems with integral type performance indices. A convenient analysis of this problem is possible, using the developed phase canonical form. A theorem is stated which asserts necessary and sufficient conditions for optimized performance indices for a specified feedback control law. Further results concern the nonnegativity of loss functions as optimized performance indices under the additional assumptions that the nonlinearities of the system are given as polynomial functions of the state variables and that the feedback control law results in a linear autonomous system. A theorem of necessary conditions for this is given. Sufficient conditions are stated for linear systems. Based on these main theorems, supplementary theorems and

corollaries are given which reveal other fundamental aspects of optimal feedback control systems.

In comparison with similar studies by other investigators, this work is directed toward more general assumptions on the inverse problem, i.e., generalizations of the system description, the specified feedback control law, and the performance indices. As a consequence, results of other investigators can be described as special cases of those resulting from this work.

TABLE OF CONTENTS

	<u>Page</u>
VITA	ii
ACKNOWLEDGEMENT	iii
ABSTRACT	iv
TABLE OF CONTENTS	vi
LIST OF FIGURES	xi
 1 INTRODUCTION	 1
1.1 Introduction to the Thesis	1
1.2 Outline of the Thesis	6
1.3 Summary of the Thesis	8
 2 FOUNDATIONS	
2.1 Fundamental Mathematical Concepts	10
2.1.1 Notation for Vectors and Matrices	 10
2.1.2 Multi-Variable Scalar Functions	11
2.1.3 Sign Definite Matrices	16
2.2 Description of Physical Systems	19
2.3 Stability and Controllability	25
2.3.1 System Stability, Stability in the Liapunov Sense	 25
2.3.2 System Controllability	28

3	CANONICAL FORMS OF LINEAR SYSTEMS	32
3.1	Canonical Forms of Systems	32
3.2	Jordan Matrix Canonical Form	37
3.3	Phase Variable Canonical Form	40
3.3.1	Phase Variable Canonical Form for Single Input Systems	40
3.3.2	Phase Variable Canonical Forms for Multi-Input Systems	43
3.3.2.1	Canonical Form by Tuel	43
3.3.2.2	Canonical Form by Luenberger	46
3.3.2.3	Canonical Form by Asseo	48
3.4	Comments	52
4	DEVELOPMENT OF A CANONICAL FORM	54
4.1	Equivalent System	54
4.2	Fundamental Theorems	57
4.3	Development of a New Canonical Form for Linear Systems	91
4.3.1	Canonical Form for Linear Systems	91
4.3.2	Controllability of the System Determined from the Canonical Form	92
4.3.3	Heuristic Explanation of the Canonical Form	94
4.4	Uniqueness of the Structure of the Canonical Form	101

4.5	Uniqueness of the Canonical Form	112
4.6	Examples	119
4.7	Application of the Canonical Form of General Systems	125
5	THE OPTIMAL FEEDBACK CONTROL LAW AND THE INVERSE PROBLEM OF THE REGULATOR	130
5.1	Formulation of the Optimal Control Problem and the Inverse Problem	130
5.2	Optimal Feedback Control Law	134
5.2.1	Statement of the Problem	134
5.2.2	Fundamental Lemma	135
5.2.3	Heuristic Approach to the Optimal Feedback Control Law	137
5.2.4	Miscellaneous Comments	145
5.3	Review of Studies on the Inverse Problem	145
5.3.1	Study of Kalman	146
5.3.2	Study of Suga	149
5.3.3	Study of Thau	155
5.4	Comments	162
6	INVERSE PROBLEM OF THE OPTIMAL REGULATOR	164
6.1	Statement of the Inverse Problem	164
6.2	An Equivalent Inverse Problem	167
6.3	Fundamental Lemmas	176
6.4	Analysis of the Inverse Problem	190

6.4.1	Hamilton-Jacobi Equation	190
6.4.2	Concerning $\nabla^0(\underline{X})$	193
6.4.3	Principal Theorem of the Inverse Problem	197
6.5	Discussion	200
6.5.1	On the General Method of Solution of the Inverse Problem	200
6.5.2	Dependency within $\underline{U}(\underline{X})$	201
6.5.3	Consideration of the Variety of $\underline{L}(\underline{X})$	201
6.5.4	Uniqueness of $\underline{L}(\underline{X})$	202
6.5.5	Linear Control Law	203
6.5.6	Nonnegative $\nabla^0(\underline{X})$ for a Linear Control Law	207
6.5.7	Necessity of Control Action	212
6.5.8	Asymptotic Stability of the Synthesized Feedback Control System	214
6.5.9	Miscellaneous Comments	216
6.6	Examples	217
7	THE INVERSE PROBLEM OF LINEARLY SYNTHESIZED FEEDBACK CONTROL SYSTEMS	225
7.1	Statements of the Modified Inverse Problem	225
7.2	Fundamental Lemma	228

7.3	Solution of the Modified Inverse Problem	230
7.4	Nonnegative Loss Function of the Optimized Performance Index	236
7.4.1	Nonnegative $L(\underline{X})$ for a Controllable System	237
7.4.2	Nonnegative $L(\underline{X})$ for an Uncontrollable System	241
8	CONCLUSIONS AND SUGGESTIONS FOR FURTHER STUDIES	243
8.1	Conclusions	243
8.1.1	The Canonical Form	243
8.1.2	The Inverse Problem of the Optimal Regulator	248
8.2	Suggestions for Further Studies	251
	REFERENCES	252

LIST OF FIGURES

	<u>Page</u>
 Chapter 2	
Figure 2-1. A physical system.	19
 Chapter 3	
Figure 3-1. Subsystem (S_i) for (3-20).	40
Figure 3-2. Single input phase variable canonical form given by (3-21) and (3-22).	41
Figure 3-3. Single input phase variable canonical form given by (3-23) and (3-24).	42
Figure 3-4. Luenberger canonical form.	49
 Chapter 4	
Figure 4-1. Canonical form of a completely controllable linear system (4-117).	97
Figure 4-2. Canonical form of an uncon- trollable linear system (4-118).	98
Figure 4-3. Subsystem of the canonical form in (4-117) and (4-118).	100
Figure 4-4. Canonical form of Example 4-1, (4-184).	126

Figure 4-5. Canonical form of Example 4-2,
(4-186). 126

Figure 4-6. Canonical form of Example 4-3,
(4-190). 126

Chapter 1

INTRODUCTION

1.1 Introduction to the Thesis

When a task in the physical world is approached, there naturally occurs the question of the best method to accomplish it. Problems of optimal control are those which attempt to find the best methods through mathematical descriptions of the task. These mathematical descriptions are composed essentially of

(i) a model of the cause-effect dynamics of the task, i.e., the system equation,

(ii) beginning and end parts of the task, i.e., initial and final conditions for the system equations,

(iii) permitted methods to accomplish the task, i.e., admissible control functions, and

(iv) a standard to measure the optimality for each admissible method, i.e., a performance index.

This search for best methods is identified with calculations of an optimal control function to transfer the conditions of the system as desired with the minimum possible value of the performance index.

An elementary problem of optimal control can be formulated from a single initial condition, a single final

condition and some control function of time to be calculated. This is generally called an open loop optimal control problem. Practically these problems are more often recognized under somewhat different circumstances, those of an optimal regulator. Thus the system function is to maintain a specified condition even though it is exposed occasionally to unforeseen disturbances. The recovery to the specified condition after each disturbance is to be in some optimal manner.

Consider, for example, a room in a building with an air temperature of 10°C . It may be desired to change the temperature to a steady 20°C as fast as possible, using a particular heating system. The open loop optimum control problem would be to design the given heating system behavior to minimize the time required for this change, recognizing the characteristics of the room and the heating system. Alternately, the regulation of the room temperature at 20°C could be considered under disturbances due to opening and shutting of doors and to outside weather conditions. To return the temperature to 20°C in some optimal sense, say minimum time, after a change due to these unpredictable causes, is a problem of an optimal regulator.

An optimal regulator problem can be expressed as a family of many individual open loop optimal control problems

with common system equations and assumptions, but with different initial conditions. If optimal control functions are required for many initial conditions, the necessity to calculate them one by one is unreasonable. An optimal control function can be developed alternately in closed form as a function of the system condition. Generally called an optimal feedback control law, this function establishes the optimal control from the present system condition regardless of its preceding career. Consequently, control action for any initial condition can be covered by a single control law.

Practical problems of optimal control may have many controlling and controlled quantities represented by complex algebraic descriptions, e.g., rocket control during space flight, utility power plant and distribution network regulation, and industrial chemical process control. Consider the regulation of a rocket flight path to a fixed trajectory in space. The control variables, representing thrust from individual rockets or combinations of rockets, could be represented in terms of three orthogonal directions of space. The output quantities are the position and velocity of the rocket, each composed of three orthogonal component variables. Consequently, the system has six controlled variables associated with the three controlling variables. The performance index might be a minimum integral

of deviations of the rocket from the fixed trajectory, or alternately, minimum fuel consumption during a specific control action.

The problem of optimal control has been intensively studied as one of the main branches of modern control theory, not only because of the interesting mathematics of the problem, but because of the practical character of the solutions. Many techniques of optimization have been developed in this field, based mostly on the calculus of variations. [1-6] Unsolved problems, however, still exist. Analytically these result from difficulties in finding sufficient conditions for optimality in a general sense and from the rapid increase of required calculations as the size of the system equations becomes realistically large. Complete answers to optimum control problems are restricted at the present time to a few specific classes of problems with relatively simple system equations. Furthermore, the optimal control function solutions to some of these problems may be impossible or inconvenient to reduce to hardware.

This thesis investigates relevant characteristics of optimal feedback control laws for a class of optimal regulators with system equations given by multi-input, time invariant ordinary differential equations. But the direction of the attack is just opposite or inverse to

the usual methods of investigation. The question is asked, what performance indices can be optimized by an assumed feedback control law? The objective is to seek all performance indices shared by a control law. This problem formulation is generally called the inverse problem of the optimal regulator,^[6] or briefly the inverse problem.

A study of the inverse problem would (1) disclose practical advantages of using specific classes of control laws in combination with specific performance indices, and (2) distinguish between control laws which are optimal in some sense and those which are not. Consequently, the results would allow future optimal design problems to start with realistic performance indices, and be helpful for understanding observed solutions to control problems that are assumed optimal in some sense.

Before an analysis of the inverse problem is given in the following chapters, however, a canonical form is developed for the class of systems of interest. A canonical form is a compact standard form for describing all systems that are mathematically similar. It is useful both to clearly expose the mathematical composition of the system structure and to allow the analysis and design of the system to conveniently proceed using a compact description. Necessarily the choice of a canonical form for a given

class of problems plays an important role in the succeeding analysis. While various canonical forms have been suggested for linear systems, the new one introduced in this study is apparently necessary for the analysis of the inverse problem considered. It is also more generally useful for demonstrating the mathematical structure of systems. Thus this thesis considers two topics, (1) the inverse problem of the optimal regulator and (2) a canonical form for a broad class of multi-input deterministic systems described by time invariant ordinary differential equations.

1.2 Outline of the Thesis

The material presented in this thesis is divided into eight chapters. Chapter 1 contains an introduction to the topics considered, an outline, and a summary of the results. Chapter 2 is an introduction to the mathematical formalism used in subsequent theoretical developments. Notations, definitions and theorems that are assumed in later chapters are given here.

Chapter 3 summarizes and reviews work that has been published in the area of canonical forms for linear time invariant systems. These canonical forms are grouped in two categories according to descriptive structures, the Jordan standard matrix canonical form and the phase variable

canonical form. The new phase canonical form developed for a class of multi-input systems is given in Chapter 4. First, the development is concerned with linear systems, defining two characteristic quantities, a stage number and a stage distribution, which determine the structure of the systems uniquely, as illustrated by three examples. Second, the results are expanded for more general non-linear systems.

Chapter 5 summarizes and reviews publications on the inverse problem of optimal regulators. Interest is focused on work for problem assumptions similar to those made during this investigation. In Chapter 6, after a precise statement of the inverse problem of the thesis, a general analysis is given. A number of interesting characteristics of optimal feedback control systems revealed by this analysis are then given.

The results of Chapter 6 are applied to linearly synthesized optimal feedback control systems in Chapter 7. Conditions for a nonnegative loss function in optimized performance indices are given, and a principle of necessity of control action is disclosed. Conclusion and suggestions for future studies are in Chapter 8.

1.3 Summary of the Thesis

There are three principal results in this thesis, applicable to a class of deterministic, dynamic systems described by time invariant ordinary differential equations either, linear or nonlinear.

(i) A new phase canonical form is developed for this class of systems. In comparison with other canonical forms, this canonical form has the advantages that (1) it is applicable to a larger class of nonlinear, uncontrollable systems, (2) its structure is uniquely determined for each system by a defined quantity, a stage distribution, and (3) it would appear to be more conveniently used for analyses of optimal control problems than other known forms.

(ii) Necessary and sufficient conditions are given for feedback control laws which are optimal for performance indices given as integral forms and loss functions as sums of penalty functions of the state variables plus positive definite quadratic forms of the control variables. From this result, additional characteristics of optimal feedback control systems are revealed.

(iii) Necessary conditions for nonnegative loss functions in optimized performance indices are given for the inverse problem resulting in linearly synthesized optimal feedback

control systems, assuming the nonlinearities of the systems and the penalty functions are given as polynomials of the state variables. These conditions are also sufficient if the system has no nonlinearity. Specifically it is shown that, for controllable linear systems with linear feedback control laws, nonnegative loss functions in optimized performance indices must be quadratic forms of the state variables and control variables. Additional relevant aspects of linearly synthesized optimal control systems are also given.

Chapter 2

FOUNDATIONS

Material basic to the theoretical developments throughout this thesis is given in this chapter. Following the introduction of mathematical notations for the abstract space descriptions of multi-variable functions, sections are given defining and explaining system modeling, solutions, stability, and controllability. A definition of approachability is introduced.

2.1 Fundamental Mathematical Concepts

2.1.1 Notation for Vectors and Matrices

Vectors and matrices are denoted by underlined capital Roman letters or Greek letters. Their dimensions are apparent either from definitions or are stated explicitly. A null matrix and a null vector are denoted by $[0]$ and $\underline{0}$ respectively. The $r \times r$ unit matrix is given by I_r and the inverse matrix by a -1 superscript, e.g., \underline{A}^{-1} . The transpose of a matrix or a vector has a T superscript, e.g., \underline{A}^T or \underline{X}^T . The scalar product of two vectors \underline{X} and \underline{Y} is

$$\underline{X}^T \underline{Y} \triangleq \sum x_i y_i . \quad (2-1)$$

The quadratic form of \underline{X} associated with a matrix \underline{A} is written as

$$\underline{x}^T \underline{A} \underline{x} \triangleq \sum_{i,j} a_{ij} x_i x_j . \quad (2-2)$$

The Euclidian norm of a vector \underline{x} , denoted by $||\underline{x}||$, is

$$||\underline{x}|| = \sqrt{\underline{x}^T \underline{x}} . \quad (2-3)$$

All subsequent discussions are assumed to be in finite dimensional Euclidian spaces or the product spaces. R^n designates an n-dimensional Euclidian space.

2.1.2 Multi-Variable Scalar Functions

With vector notation, functions of many variables are conveniently described as

$$f(\underline{x}) = f(x_1, x_2, \dots, x_n) ,$$

$$f(\underline{x}, \underline{u}) = f(x_1, x_2, \dots, x_n, u_1, u_2, \dots, u_m) ,$$

or
$$f(\underline{x}, \underline{u}, t) = f(x_1, x_2, \dots, x_n, u_1, u_2, \dots, u_m, t) \quad (2-4)$$

defined on R^n , $R^n \times R^m$ and $R^n \times R^m \times R^1$ spaces respectively.

Theorem 2-1:

If $f(\underline{x})$ is continuous on R^n , then $f(\underline{N} \underline{x})$ is also continuous on R^n , where \underline{N} is an $n \times n$ constant nonsingular matrix.

The proof follows directly from a fundamental theorem of the composition of continuous mappings. [7]

Definition 2-1: Functions of Class C_n . [7]

A scalar function $f(\underline{x})$ defined on R^n is said to be of class C_r in a region $\Gamma \subset R^n$ if it has continuous partial derivatives with respect to all x_i , ($i = 1, 2, \dots, n$), up to order r everywhere in Γ . When Γ is the entire R^n , the phrase "in Γ " is omitted.

Provided that $f(\underline{x})$ is of class C_2 , the following notations are often used

$$\frac{\partial f}{\partial \underline{x}} = \text{grad } f(\underline{x}) \triangleq \begin{bmatrix} \frac{\partial f(\underline{x})}{\partial x_1} \\ \frac{\partial f(\underline{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\underline{x})}{\partial x_n} \end{bmatrix}, \quad (2-5)$$

$$\frac{\partial^2 f}{\partial \underline{X} \partial \underline{X}} \triangleq \frac{\partial}{\partial \underline{X}} \{ \text{grad } f(\underline{X}) \} \triangleq \begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_1} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2 \partial x_2} & \cdots & \cdots \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n} & \cdots & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{bmatrix}$$

(2-6)

Theorem 2-2:

Let

$$\underline{W}(\underline{X}) = \begin{bmatrix} w_1(\underline{X}) \\ w_2(\underline{X}) \\ \vdots \\ w_n(\underline{X}) \end{bmatrix} \quad (2-7)$$

be a vector valued function of class C_1 . Then $\frac{\partial}{\partial \underline{X}} \{ \underline{W}(\underline{X}) \}$ is symmetric if and only if $\underline{W}(\underline{X})$ is a gradient of a scalar function $V(\underline{X})$, i.e.,

$$\text{grad } V(\underline{X}) = \underline{W}(\underline{X}) \quad (2-8)$$

The proof of this theorem follows directly from the theory of functions of several variables. [7, Chapter 3] Furthermore, $\nabla(\underline{X})$ can be uniquely determined by the line integral

$$\nabla(\underline{X}) = \int_0^{\underline{X}} \underline{W}^T(\underline{X}) d\underline{X} + C \quad (2-9)$$

and is independent of the path of integration. A convenient line integral is

$$\begin{aligned} \nabla(\underline{X}) = & \int_0^{x_1} w_1(\gamma_1, 0 \dots) d\gamma_1 + \int_0^{x_2} w_2(x_1, \gamma_2, 0 \dots 0) d\gamma_2 + \dots \\ & \dots + \int_0^{x_n} w_n(x_1, x_2, \dots, x_{n-1}, \gamma_n) d\gamma_n + C, \end{aligned} \quad (2-10)$$

where C is a constant. [8]

Definition 2-2: Sign Definite Functions. [9]

Consider a scalar function $\nabla(\underline{X})$ defined on R^n and an open region $\Gamma \subset R^n$ such that $\underline{0} \in \Gamma$. Then $\nabla(\underline{X})$ is said to be positive (negative) semidefinite in Γ if

$$(i) \quad \nabla(\underline{0}) = 0 \quad (2-11)$$

$$(ii) \quad \nabla(\underline{X}) \geq 0, (\nabla(\underline{X}) \leq 0), \text{ for all } \underline{X} \in \Gamma. \quad (2-12)$$

If $\nabla(\underline{X})$ satisfies

$$(ii)' \quad \nabla(\underline{x}) > 0, (\nabla(\underline{x}) < 0), \text{ for all } \underline{x} \in r \text{ of } \underline{x} \neq \underline{0} \quad (2-13)$$

instead of (ii), it is said to be positive (negative) definite in r . When r is the entire R^n , the phrase "in r " is omitted.

Theorem 2-3:

Assume a finite degree polynomial scalar function $\nabla(\underline{x})$ given by

$$\nabla(\underline{x}) = \nabla^{(2)}(\underline{x}) + \nabla^{(3)}(\underline{x}) + \dots + \nabla^{(\xi)}(\underline{x}), \quad (2-14)$$

where each $\nabla^{(i)}(\underline{x})$, ($i = 2, 3, \dots, \xi$), is an i^{th} degree homogeneous function of \underline{x} and $\nabla^{(\xi)}(\underline{x})$, $\xi \geq 2$, is not identically zero. Then, for $\nabla(\underline{x})$ to be positive semi-definite, it is necessary that $\nabla^{(\xi)}(\underline{x})$ be positive semi-definite.

Proof: Consider the contrary, that $\nabla(\underline{x})$ is positive semi-definite but there exists a certain vector $\underline{x} \in R^n$ such that

$$\nabla^{(\xi)}(\underline{x}) < 0. \quad (2-15)$$

Consider a set of vectors given $y\underline{x}$, where y is a scalar

variable. Then

$$f(y) \triangleq V(y\underline{x}) = y^2 V^{(2)}(\underline{x}) + y^3 V^{(3)}(\underline{x}) + \dots + y^\xi V^{(\xi)}(\underline{x}), \quad (2-16)$$

which is a polynomial function of y . Since $V^{(\xi)}(\underline{x})$ is negative, $f(y)$ becomes negative as $y \rightarrow +\infty$ and $V(\underline{x})$ becomes negative for $\underline{x} = y\underline{x}$. This contradicts the hypothesis. It also follows that ξ must be even for positive definiteness of $V(\underline{x})$.

2.1.3 Sign Definite Matrices

Assume \underline{Q} to be an $n \times n$ real matrix.

Definition 2-3: Sign Definite Matrix. [10]

The matrix \underline{Q} is said to be positive (positive semi) definite if the scalar function $\underline{x}^T \underline{Q} \underline{x}$ is positive (positive semi) definite. If $-\underline{x}^T \underline{Q} \underline{x}$ is positive (positive semi) definite, \underline{Q} is said to be negative (negative semi) definite.

Without loss of generality, \underline{Q} is also assumed to be a symmetric matrix in this section. [10]

Theorem 2-4: [10]

(i) Every \underline{Q} is congruent to

$$\begin{bmatrix} I_{n_1} & & \\ & -I_{n_2} & \\ & & [0] \end{bmatrix}, \text{ all other entries zero,} \quad (2-17)$$

where n_1 and n_2 are uniquely determined by \underline{Q} .

(ii) \underline{Q} is positive definite if and only if $n_1 = n$, i.e., (2-17) degenerates to the unit matrix, and \underline{Q} is positive semidefinite if and only if $n_2 = 0$.

(iii) An $n \times n$ real symmetric matrix \underline{Q}_1 is congruent to \underline{Q} if and only if \underline{Q}_1 is congruent to (2-17), i.e., sign definiteness of \underline{Q} is invariant for a congruent transformation.

Corollary 2-1:

(i) \underline{Q} is positive definite if and only if all principal minor determinants of \underline{Q} are positive. [10]

(ii) \underline{Q} is positive semidefinite if and only if there exists an $n_1 \times n$ matrix \underline{D} of rank $n_1 \leq n$ such that [11]

$$\underline{Q} = \underline{D}^T \underline{D}. \quad (2-18)$$

(iii) \underline{Q} is positive definite if and only if $n = n_1$ above.

(iv) \underline{Q} is positive definite if and only if \underline{Q}^{-1} is positive definite.

Proof: The proof of (iii) follows from (ii) of this Corollary and (ii) of the Theorem. For (iv), for \underline{Q} to be positive definite, it must be nonsingular by (i) and \underline{Q}^{-1} exists. Consider a congruent transformation of \underline{Q} by \underline{Q}^{-1} . Thus

$$(\underline{Q}^{-1})^T \underline{Q} \underline{Q}^{-1} = (\underline{Q}^{-1})^T = (\underline{Q}^T)^{-1} = \underline{Q}^{-1}, \quad (2-19)$$

which follows from the transpose of an inverse matrix^[10] and from the symmetry of \underline{Q} . From (iii) of Theorem 2-4, \underline{Q} is positive definite if and only if \underline{Q}^{-1} is positive definite.

Theorem 2-5:^[13]

Let \underline{Q} be an arbitrary $n \times n$ positive definite symmetric matrix. Then there exist two positive numbers, $\lambda_{\min}(\underline{Q})$ and $\lambda_{\max}(\underline{Q})$, such that

$$\lambda_{\min}(\underline{Q}) \|\underline{x}\|^2 \leq \underline{x}^T \underline{Q} \underline{x} \leq \lambda_{\max}(\underline{Q}) \|\underline{x}\|^2 \quad \text{for all } \underline{x} \in \mathbb{R}^n. \quad (2-20)$$

2.2 Description of Physical Systems [5,12]

Systems which behave according to the Principle of Causality in the physical world can be schematically described as in Figure 1-1. Causes are classified into

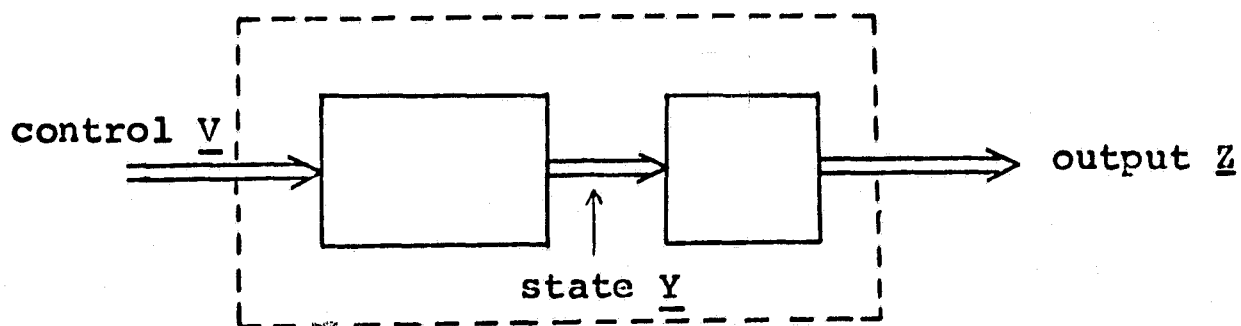


Figure 2-1. A physical system.

controls and disturbances; the former quantities can be specified and manipulated at will but the latter cannot. Effects are created by the causes on the physical system. These are outputs or directly observed quantities of effects. To some extent, outputs depend upon the preferences of the observer. To avoid any ambiguity of outputs, a third quantity, states within the system, are conveniently introduced. Always regarded as abstract quantities, the states are defined as the minimal amount of information about the past history of the system which is sufficient to predict the affects upon the future.

Systems considered in this thesis are assumed to belong to a class called deterministic, real, finite dimensional, continuous time, ordinary differential, dynamical systems described by the equations

$$\dot{\underline{Y}} = \hat{\underline{F}}(\underline{Y}, \underline{V}, t) \quad (2-21)$$

$$\underline{Z} = \underline{G}(\underline{Y}, t) . \quad (2-22)$$

\underline{Y} , \underline{V} , and \underline{Z} are called the state, control and output vectors,

$$\underline{Y} = [y_1, y_2, \dots, y_n]^T \quad (2-23)$$

$$\underline{V} = [v_1, v_2, \dots, v_m]^T \quad (2-24)$$

$$\underline{Z} = [z_1, z_2, \dots, z_{n_o}]^T . \quad (2-25)$$

By each adjective, the following is meant:

(i) Deterministic: The process described by the system is deterministic.

(ii) Real and finite dimensional: The state and control vectors are defined in real finite dimensional spaces.

(iii) Continuous time: The set of time for these equations is an open interval $(T_0, T_1) \subset \mathbb{R}^1$.

(iv) Ordinary differential: The behavior of the state of the system is given by ordinary differential equations as (2-21), where $\cdot = \frac{d}{dt}$.

(v) While a more detailed definition of a dynamical system can be given, [5,12] the following is adequate for this work. Dynamical system: From any \underline{y}_0 and t_0 and for any piecewise continuous n dimensional vector valued function $\underline{v}(t)$, each existing in defined regions, there is a unique solution $\hat{\underline{\phi}}_v(t; \underline{y}_0, t_0)$ to (2-21), i.e., an n -dimensional, vector valued function, differentiable in t , satisfying

$$(a) \quad \hat{\underline{\phi}}_v(t_0; \underline{y}_0, t_0) = \underline{y}_0, \quad (2-26)$$

$$(b) \quad \frac{d}{dt} \hat{\underline{\phi}}_v(t; \underline{y}_0, t_0) = \hat{F}(\hat{\underline{\phi}}_v(t; \underline{y}_0, t_0), \underline{v}(t), t) \\ \text{for all } t \in \mathbb{R}^1, \quad (2-27)$$

$$(c) \quad \hat{\underline{\phi}}_v(t; \underline{y}_0, t_0) = \hat{\underline{\phi}}_v(t; \hat{\underline{\phi}}_v(t_1; \underline{y}_0, t_0), t_1) \\ \text{for all } t \geq t_1 \geq t_0. \quad (2-28)$$

Generally (2-21), (2-22) and a solution $\hat{\underline{\phi}}_v(t; \underline{y}_0, t_0)$ are collectively called the system equation, the output equation, and the solution for (\underline{y}_0, t_0) and $\underline{v}(t)$. The behavior of the state variables can be identified directly

from the solution, while the output variables are calculated algebraically from (2-22). Consequently substantial analyses of the systems can proceed through system equations only, with the output equations reduced to a secondary role. Thus a dynamical system is usually represented by just the system equations, (2-21), with the output equations, (2-22), implied.

The existence of a unique solution to system equations depends both on the given system equations and a specified control function. A sufficient condition can be stated as follows.

Theorem 2-6: [5]

For the system (2-21), if

- (i) $\hat{F}(\cdot)$ is a continuous function, from $R^n \times R^m \times R^1$ into R^n ,
- (ii) the partial derivatives $\frac{\partial \hat{f}_i(\cdot)}{\partial y_j}$ for $i, j = 1, 2, \dots, n$, are continuous functions from $R^n \times R^m \times R^1$ into R^1 ,
- (iii) $\underline{V}(t)$ is a piecewise continuous function from R^1 to R^m , and
- (iv) $\underline{y}_0 \in R^n$ and $t_0 \in R^1$ are specified,

then there exists a unique solution $\hat{\underline{\phi}}_{\underline{V}}(t; \underline{y}_0, t_0)$ on a time interval containing t_0 .

Imbedding the specified $\underline{v}(t)$ as $\hat{\underline{F}}(\underline{y}, t) = \hat{\underline{F}}(\underline{y}, \underline{v}(t), t)$, another theorem can be given.

Theorem 2-7: [2,13]

Let D be a polyhedron in $R^n \times R^1$ where (\underline{y}_0, t_0) exists as an interior point. For a system

$$\dot{\underline{y}} = \hat{\underline{F}}(\underline{y}, t) , \quad (2-29)$$

if $\hat{\underline{F}}(\underline{y}, t)$ is continuous in D and there exists a number $k > 0$ which satisfies

$$||\hat{\underline{F}}(\underline{y}_1, t) - \hat{\underline{F}}(\underline{y}_2, t)|| \leq k ||\underline{y}_1 - \underline{y}_2||$$

for all $(\underline{y}_1, t), (\underline{y}_2, t) \in D$, (2-30)

then there exists a unique solution $\hat{\underline{\phi}}_{\underline{y}}(t; \underline{y}_0, t_0)$ in a domain $D^1 \subset D$.

Generally (2-30) is called a Lipschitz condition.

Depending upon characteristics, systems are described by special names.

Definition 2-4: Time Invariant, Free, Autonomous System. [14]

A system is said to be

(i) time invariant (or stationary), if (2-21) is given as

$$\dot{\underline{y}} = \hat{\underline{F}}(\underline{y}, \underline{v}) , \quad (2-31)$$

(ii) free, if (2-21) is given as

$$\dot{\underline{y}} = \hat{\underline{F}}(\underline{y}, t) , \quad (2-32)$$

and

(iii) autonomous, if (2-21) is given as

$$\dot{\underline{y}} = \hat{\underline{F}}(\underline{y}) , \quad (2-33)$$

i.e., time invariant and free.

Definition 2-5: Linear System. [12]

A system is said to be linear if (2-21) is given as

$$\dot{\underline{y}} = \hat{\underline{A}}(t) \underline{y} + \hat{\underline{B}}(t) \underline{v} \quad (2-34)$$

where $\hat{\underline{A}}(t)$ and $\hat{\underline{B}}(t)$ are $n \times n$ and $n \times m$ matrix valued functions.

For a free system, the concept of an equilibrium point is particularly convenient.

Definition 2-6: Equilibrium Point of a System. [9]

A point \underline{y}_e in R^n is said to be an equilibrium point of a free system (2-32) if

$$\underline{0} = \hat{\underline{F}}(\underline{y}_e, t) \quad \text{for all } t \in R^1. \quad (2-35)$$

Practical examples of physical systems described by this symbolism and terminology appear in the literature. [3,5,8]

2.3 Stability and Controllability

For any given system, two descriptive characteristics can be considered, system stability and system controllability. These are useful in the analyses of the system behavior and the syntheses of control functions.

2.3.1 System Stability, [9,14] Stability in the Liapunov Sense

The theory of the behavior of solutions in relation to an equilibrium point of a free system is known as stability theory. It is based largely on concepts originally proposed by Liapunov. Some of the extensive developments of this theory are particularly convenient for application to the class of systems considered during this investigation.

Assume a free system given by

$$\dot{\underline{y}} = \hat{\underline{F}}(\underline{y}, t) \quad (2-36)$$

defined on $R^n \times R^1$ and, for convenience, also assume $\underline{y} = \underline{0}$ is an equilibrium point, i.e.,

$$\hat{\underline{F}}(\underline{0}, t) = \underline{0} \quad \text{for all } t \in R^1. \quad (2-37)$$

Definition 2-7: Stability. [9]

The origin of the system given by (2-36) is said to be stable with respect to t_0 if for every $\epsilon > 0$, there exists $\delta(\epsilon, t_0) > 0$ such that

$$||\underline{y}_0|| < \delta \quad (2-38)$$

implies

$$||\hat{\phi}_f(t; \underline{y}_0, t_0)|| < \epsilon \quad \text{for all } t \geq t_0, \quad (2-39)$$

where $\hat{\phi}_f(t; \underline{y}_0, t_0)$ is a solution of (2-36) from (\underline{y}_0, t_0) .

When the system is a dynamical system, and stability exists for some t_0 , then it is stable for any other $t_1 \in R$. [14]
Accordingly, the expression "with respect to t_0 " can be eliminated from the definition.

Definition 2-8: Asymptotic Stability. [9]

The origin of the system given by (2-36) is said to

be asymptotically stable if

- (i) it is stable, and
- (ii) for every $\mu > 0$, there exists a $T(\mu, t_0, \delta) > 0$ such that

$$||\hat{\phi}_f(t; \underline{y}_0, t_0)|| < \mu \quad \text{for all } t \geq t_0 + T. \quad (2-40)$$

Heuristically every solution starting in a neighborhood of the equilibrium point at any t_0 is required to converge to $\underline{0}$ as $t \rightarrow \infty$.

Definition 2-9: Asymptotic Stability in the Large. [14]

The origin of the system of (2-36) is said to be asymptotically stable in the large if

- (i) it is stable, and
- (ii) every solution converges to $\underline{0}$ as $t \rightarrow \infty$.

These three definitions are generally considered to refer to stability in the Liapunov sense. Criteria for satisfying these definitions of stability are established in the so-called theory of the Liapunov direct method.

Theorem 2-8: [6]

Consider an autonomous system

$$\dot{\underline{Y}} = \underline{\hat{F}}(\underline{Y}) \quad (2-41)$$

of which $\underline{\hat{F}}(\underline{Y})$ is defined on Γ such that $\underline{0} \in \Gamma$. If there exists scalar function $\hat{V}(\underline{Y})$ of class C_1 in R^n such that

(i) $\hat{V}(\underline{Y})$ is positive definite in Γ ,

(ii) $\dot{\hat{V}}(\underline{Y}) = \left[\frac{\partial \hat{V}}{\partial \underline{Y}} \right]^T \underline{\hat{F}}(\underline{Y})$ is negative definite in Γ ,

then the origin is asymptotically stable. If $\dot{\hat{V}}$ is negative semidefinite in (ii), the origin is stable. If $\hat{V}(\underline{Y})$ also satisfies

(iii) $\lim_{\|\underline{Y}\| \rightarrow \infty} \hat{V}(\underline{Y}) \rightarrow \infty$ in R^n

and Γ can be selected as the entire R^n , the origin is asymptotically stable in the large.

A function $\hat{V}(\underline{Y})$ to identify stability in this sense is generally called a Liapunov function for the system. [14]

As the existence of a Liapunov function for a given system is sufficient to guarantee stability without knowledge of specific solutions, the method is particularly valuable for nonlinear systems.

2.3.2 System Controllability

When a system is to be controlled so as to transfer its initial condition to another condition, there is a

question as to the realizability of the requirement, i.e., whether a control function exists for the transition. To summarize this idea, the concept of an admissible control function is introduced first.

Definition 2-10: Admissible Control Function. [5]

A control function $\underline{v}(t)$ defined on $[t_0, t_1] \subset \mathbb{R}^1$ is said to be admissible if it is piecewise continuous on the interval.

In effect, an admissible control function provides a unique solution from an arbitrary initial condition to (2-21), as provided by Theorem 2-5.

Definition 2-11: Controllability. [5]

For a system given by

$$\dot{\underline{y}} = \hat{\underline{F}}(\underline{y}, \underline{v}, t), \quad (2-42)$$

a state \underline{y}_0 is said to be controllable at t_0 with respect to a state \underline{y}_f if there exists an admissible control function $\underline{v}(t)$ defined on $[t_0, t_f]$ such that

$$\hat{\underline{\phi}}_{\underline{v}}(t_f; \underline{y}_0, t_0) = \underline{y}_f \quad \text{for some } t_f \geq t_0. \quad (2-43)$$

If every \underline{y}_0 is controllable with respect to every \underline{y}_f at every t_0 , then the system is said to be completely controllable. When the system is linear such that

$$\dot{\underline{y}} = \underline{\hat{A}} \underline{y} + \underline{\hat{B}} \underline{v}, \quad (2-44)$$

where $\underline{\hat{A}}$ and $\underline{\hat{B}}$ are $n \times n$ and $n \times m$ constant matrices respectively, a criterion of complete controllability is known.

Theorem 2-9: [5]

The linear system given by (2-44) is completely controllable if and only if the $n \times mn$ matrix

$$[\underline{\hat{B}}, \underline{\hat{A}} \underline{\hat{B}}, \underline{\hat{A}}^2 \underline{\hat{B}}, \dots, \underline{\hat{A}}^{n-1} \underline{\hat{B}}] \quad (2-45)$$

has rank n .

The following two definitions, extending this idea, are used in later chapters.

Definition 2-12: Reachability. [5]

For a system (2-42), a state \underline{y}_1 is said to be reachable with respect to a state \underline{y}_0 at t_0 if there is an admissible control function $\underline{v}(t)$ defined on $[t_0, t_f]$ such

that it provides a solution $\hat{\phi}_{\underline{v}}(t; \underline{y}_0, t_0)$ which satisfies

$$\hat{\phi}_{\underline{v}}(t_1; \underline{y}_0, t_0) = \underline{y}_1 \quad \text{for some finite } t_1 \geq t_0. \quad (2-46)$$

For the case of an unlimited t_f , a new definition is introduced.

Definition 2-13: Approachability.

For a system given by (2-42), a state \underline{y}_1 is said to be approachable with respect to a state \underline{y}_0 at t_0 if there exist an admissible control function $\underline{v}(t)$ defined on $[t_0, \infty)$ to provide a solution $\hat{\phi}_{\underline{v}}(t; \underline{y}_0, t_0)$ such that for each $\epsilon > 0$, there exists a $T > 0$ which satisfies

$$||\hat{\phi}_{\underline{v}}(t; \underline{y}_0, t_0) - \underline{y}_1|| < \epsilon \quad \text{for all } t \geq t_0 + T. \quad (2-47)$$

If the \underline{y}_1 is approachable with respect to every state at every t_0 , \underline{y}_1 is said to be totally approachable for the system.

Chapter 3

CANONICAL FORMS OF LINEAR SYSTEMS

This chapter reviews the work that has been done to develop canonical forms for linear time-invariant systems. Some fundamental properties of canonical forms are also stated for use in subsequent chapters. A canonical form proposed by Kalman and another based on the Jordan form are described. Others that are described, called phase variable canonical forms, represent original work by Luenberger, Teul, and Asseo.

3.1 Canonical Forms of Systems

It often occurs that two different but mathematically similar systems are sufficiently related that the analyses and solutions for one can be applied to the other with relatively minor modifications. If this can be done among many systems forming a group, it is reasonable to analyze and solve the one system of the group offering the least complexity, and then translate the results to the other systems of the group.

A canonical form is a compact standard form for describing all systems that are mathematically similar. Similarity is associated with an equivalent relation on the set of systems of interest, say S . The equivalent relation used historically to develop canonical forms is as follows.

Definition 3-1: Equivalent System. [5]

Two systems are said to be equivalent if there exists an $n \times n$ nonsingular constant matrix \underline{N} such that

$$\underline{N}^{-1} \underline{Y} = \underline{X}, \quad (3-1)$$

where \underline{Y} and \underline{X} are state vectors of each system.

It can be shown that equivalent systems in S are related by topological relations called reflexive, symmetric and transitive laws, i.e., the nonsingular transformation (3-1) is a topological equivalent relation defined on S . It is known, further, that topological equivalence preserves the stability properties of dynamical systems. [12]

The value of a canonical form depends on the convenience of the specified structure in practical analyses and the extent to which it displays noteworthy characteristics. Historically, the development of canonical forms has been limited to linear systems. These forms can be described by considering the linear system

$$\dot{\underline{Y}} = \hat{\underline{A}} \underline{Y} + \hat{\underline{B}} \underline{V} \quad (3-2)$$

where $\hat{\underline{A}}$ and $\hat{\underline{B}}$ are assumed $n \times n$ and $n \times m$ matrices with r the rank of $\hat{\underline{B}}$. By (3-1), an equivalent system is

$$\dot{\underline{X}} = \underline{A} \underline{X} + \underline{\widetilde{B}} \underline{V}, \quad (3-3)$$

where

$$\underline{A} = \underline{N}^{-1} \underline{\widehat{A}} \underline{N}$$

and

$$\underline{\widetilde{B}} = \underline{N}^{-1} \underline{\widehat{B}}. \quad (3-4)$$

A canonical form is concerned with specifying a particular structure for \underline{A} and $\underline{\widetilde{B}}$.

In a general sense, Kalman suggested a canonical form for linear systems so as to conveniently display the controllability property of the system.^[12] He stated that \underline{A} and $\underline{\widetilde{B}}$ can be specified as

$$\underline{A} = \begin{bmatrix} \underline{A}_{11} & [0] \\ \underline{A}_{21} & \underline{A}_{22} \end{bmatrix} \quad (3-5)$$

and

$$\underline{\widetilde{B}} = \begin{bmatrix} [0] \\ \underline{\widetilde{B}}_2 \end{bmatrix}, \quad (3-6)$$

where the dimensions of the submatrices are

$$\begin{aligned}
 \underline{A}_{11}; \quad n_1 \times n_1 & \quad \underline{A}_{21}; \quad n_2 \times n_1 \\
 \underline{A}_{22}; \quad n_2 \times n_2 & \quad \widetilde{\underline{B}}_2; \quad n_2 \times m,
 \end{aligned}
 \tag{3-7}$$

such that $n = n_1 + n_2$ and $n_1 \geq 0$, $n_2 \geq 0$. Thus, an equivalent system to the original system is

$$\begin{aligned}
 \dot{\underline{x}}_1 &= \underline{A}_{11} \underline{x}_1 \\
 \dot{\underline{x}}_2 &= \underline{A}_{21} \underline{x}_1 + \underline{A}_{22} \underline{x}_2 + \widetilde{\underline{B}}_2 \underline{v},
 \end{aligned}
 \tag{3-8}$$

where $\underline{x}_1 = [x_1, x_2, \dots, x_{n_1}]^T$ and $\underline{x}_2 = [x_{n_1+1}, \dots, x_n]^T$.

Assuming the rank of $\hat{\underline{B}} = r > 0$, Kalman defined the controllability of systems based on this description, instead of Definition 2-11.

Definition 3-2: ^[12] Controllability (Kalman).

The system (3-2) is completely controllable if it is not equivalent to the system (3-8) with an $n_1 > 0$.

The following theorem gives the equivalency between Definition 2-11 and 3-2.

Theorem 3-1:

The system (3-2) is completely controllable in the

sense of Definition 3-2 if and only if the conditions in Theorem 2-9 are satisfied.

Proof: Assume that (3-2) is equivalent to (3-8) with an $n_1 > 0$. Then from (3-4), (3-5), and (3-6), it follows that

$$\begin{aligned} [\underline{\hat{B}}, \underline{\hat{A}} \underline{\hat{B}}, \dots, \underline{\hat{A}}^{n-1} \underline{\hat{B}}] &= [\underline{N} \underline{\tilde{B}}, \underline{N} \underline{A} \underline{\tilde{B}}, \dots, \underline{N} \underline{A}^{n-1} \underline{\tilde{B}}] \\ &= \underline{N} \left[\begin{bmatrix} [0] \\ \underline{\tilde{B}}_2 \end{bmatrix}, \begin{bmatrix} [0] \\ \underline{A}_{22} \underline{\tilde{B}}_2 \end{bmatrix}, \dots, \begin{bmatrix} [0] \\ \underline{A}_{22}^{n-1} \underline{\tilde{B}}_2 \end{bmatrix} \right], \end{aligned} \quad (3-9)$$

where each $[0]$ is $n_1 \times m$. As the rank of (3-9) is less than n , the system is not controllable by Theorem 2-9.

Conversely, assume (3-2) is completely controllable.

From Theorem 2-9, the matrix

$$[\underline{\hat{B}}, \underline{\hat{A}} \underline{\hat{B}}, \dots, \underline{\hat{A}}^{n-1} \underline{\hat{B}}] \quad (3-10)$$

must have rank n . Assume that a nonsingular \underline{N} exists such that $\underline{\hat{A}}$ and $\underline{\hat{B}}$ are transformed by (3-4) to (3-5) and (3-6).

Then

$$\underline{N}^{-1} [\underline{\hat{B}}, \underline{\hat{A}} \underline{\hat{B}}, \dots, \underline{\hat{A}}^{n-1} \underline{\hat{B}}] = [\underline{N}^{-1} \underline{\hat{B}}, \underline{N}^{-1} \underline{\hat{A}} \underline{\hat{B}}, \dots, \underline{N}^{-1} \underline{\hat{A}}^{n-1} \underline{\hat{B}}] \quad (3-11)$$

must have the rank n . But for each $0 \leq i \leq n-1$,

$$\underline{N}^{-1} \hat{\underline{A}}^i \underline{B} = \underbrace{(\underline{N}^{-1} \hat{\underline{A}} \underline{N}) (\underline{N}^{-1} \hat{\underline{A}} \underline{N}) \dots (\underline{N}^{-1} \hat{\underline{A}} \underline{N}) (\underline{N}^{-1} \underline{B})}_{i \text{ stages}} = \underline{A}^i \underline{\tilde{B}} \quad (3-12)$$

and (3-11) is reduced to

$$\left[\begin{bmatrix} [0] \\ \underline{\tilde{B}}_2 \end{bmatrix}, \begin{bmatrix} [0] \\ \underline{A}_{22} \underline{\tilde{B}}_2 \end{bmatrix}, \dots, \begin{bmatrix} [0] \\ \underline{A}_{22}^{n-1} \underline{\tilde{B}}_2 \end{bmatrix} \right], \quad (3-13)$$

where each $[0]$ is $n_1 \times m$. For (3-10) to have the rank n then, n_1 must be zero, that is, the system must not be equivalent to (3-8) with an $n_1 > 0$.

3.2 Jordan Matrix Canonical Form

A canonical form exists if \underline{A} in (3-3) is given by a Jordan canonical form of matrix. [8,15] That is, an \underline{N} can be chosen for a similar transformation of $\hat{\underline{A}}$ such that

$$\underline{A} = \underline{N}^{-1} \hat{\underline{A}} \underline{N} = \begin{bmatrix} \underline{J}_{(1)}(\lambda_1) & & & \\ & \underline{J}_{(2)}(\lambda_2) & & \\ & & \ddots & \\ & & & \underline{J}_{(v)}(\lambda_v) \end{bmatrix},$$

all other entries zero, (3-14)

where each $\underline{J}_{(i)}(\lambda_i)$, ($i = 1, 2, \dots, v$), is an $\ell_i \times \ell_i$ matrix given by

$$\underline{J}_{(i)}(\lambda_i) = \begin{bmatrix} \lambda_i & 1 & & \\ & \lambda_i & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{bmatrix}, \text{ all other entries zero,}$$

(3-15)

$\lambda_1, \lambda_2, \dots, \lambda_v$ are eigenvalues of $\hat{\underline{A}}$ and v, ℓ_1, \dots, ℓ_v are positive integers of the characteristic equation of $\hat{\underline{A}}$,

$$0 = |\hat{\underline{A}} - \lambda \underline{I}| = (\lambda - \lambda_1)^{\ell_1} \cdot (\lambda - \lambda_2)^{\ell_2} \dots (\lambda - \lambda_v)^{\ell_v}. \quad (3-16)$$

Among the numbers $\lambda_1, \lambda_2, \dots, \lambda_v$, some may be equal, however, one combination of $\underline{N}, v, \ell_1, \ell_2, \dots, \ell_v$ consistent with (3-15) and (3-16) is guaranteed to exist for the given $\hat{\underline{A}}$.^[13] It is convenient to decompose the resulting

\underline{x} and $\underline{\tilde{B}}$ as

$$\underline{x} = [\underline{x}_{(1)}^T, \underline{x}_{(2)}^T, \dots, \underline{x}_{(v)}^T]^T \quad (3-17)$$

$$\text{and } \underline{\tilde{B}} = \underline{N}^{-1} \underline{\hat{B}} = [\underline{\tilde{B}}_{(1)}, \underline{\tilde{B}}_{(2)}, \dots, \underline{\tilde{B}}_{(v)}]^T, \quad (3-18)$$

where each $\underline{\tilde{B}}_{(i)}$, ($i = 1, 2, \dots, v$), is an $\ell_i \times m$ matrix and

$$\begin{aligned} \underline{x}_{(1)} &= [x_1, x_2, \dots, x_{\ell_1}]^T \\ \underline{x}_{(2)} &= [x_{\ell_1+1}, x_{\ell_1+1}, \dots, x_{\ell_1+\ell_2}]^T \\ &\vdots \\ \underline{x}_{(v)} &= [x_{\sum_{i=1}^{v-1} \ell_i+1}, x_{\sum_{i=1}^{v-1} \ell_i+1}, \dots, x_n]^T. \end{aligned} \quad (3-19)$$

Then this canonical form can be written as

$$\begin{aligned} \dot{\underline{x}}_{(1)} &= \underline{J}_{(1)} (\lambda_1) \underline{x}_{(1)} + \underline{\tilde{B}}_{(1)} \underline{v} & (S_1) \\ \dot{\underline{x}}_{(2)} &= \underline{J}_{(2)} (\lambda_2) \underline{x}_{(1)} + \underline{\tilde{B}}_{(2)} \underline{v} & (S_2) \\ &\vdots & \vdots \\ \dot{\underline{x}}_{(v)} &= \underline{J}_{(v)} (\lambda_v) \underline{x}_{(v)} + \underline{\tilde{B}}_{(v)} \underline{v}, & (S_v) \end{aligned} \quad (3-20)$$

which can be represented by v subsystems each appearing as shown in Figure 3-1.

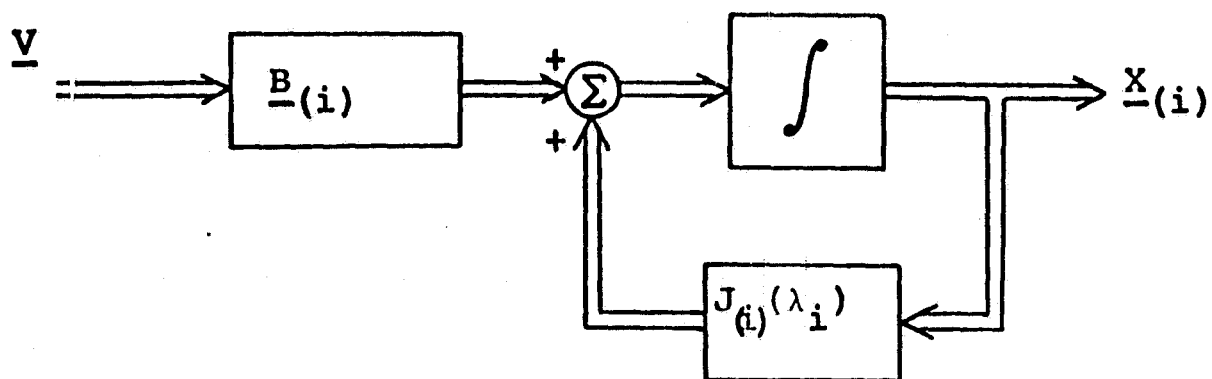


Figure 3-1. Subsystem (S_i) for (3-20).

The principal advantage of this canonical form is the convenience of calculating system solutions. This is evident as the free motion of the system is identified from the eigenvalues of $\hat{\underline{A}}$, and the system is effectively decomposed into independent subsystems with respect to the state variables. [8]

3.3 Phase Variable Canonical Form

In this section, the system (3-2) is assumed

- (i) completely controllable, and
- (ii) $\hat{\underline{B}}$ is of full rank, $r = m$.

3.3.1 Phase Variable Canonical Form For Single Input Systems [12,16]

Two canonical forms are known for single input linear systems, i.e., $m = 1$, meeting the above assumptions. One has the structures of \underline{A} and $\underline{\tilde{B}}$ given as

$$\underline{A} = \underline{N}^{-1} \hat{\underline{A}} \underline{N} = \begin{bmatrix} 0 & 1 & & & & \\ & \ddots & \ddots & & & \\ & & \ddots & \ddots & & \\ & & & \ddots & \ddots & \\ & & & & 0 & 1 \\ -a_1 & -a_2 & \cdots & \cdots & \cdots & -a_n \end{bmatrix},$$

all other entries zero, (3-21)

and $\underline{\tilde{B}} = \underline{N}^{-1} \hat{\underline{B}} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}.$ (3-22)

This form can be illustrated as shown in Figure 3-2.

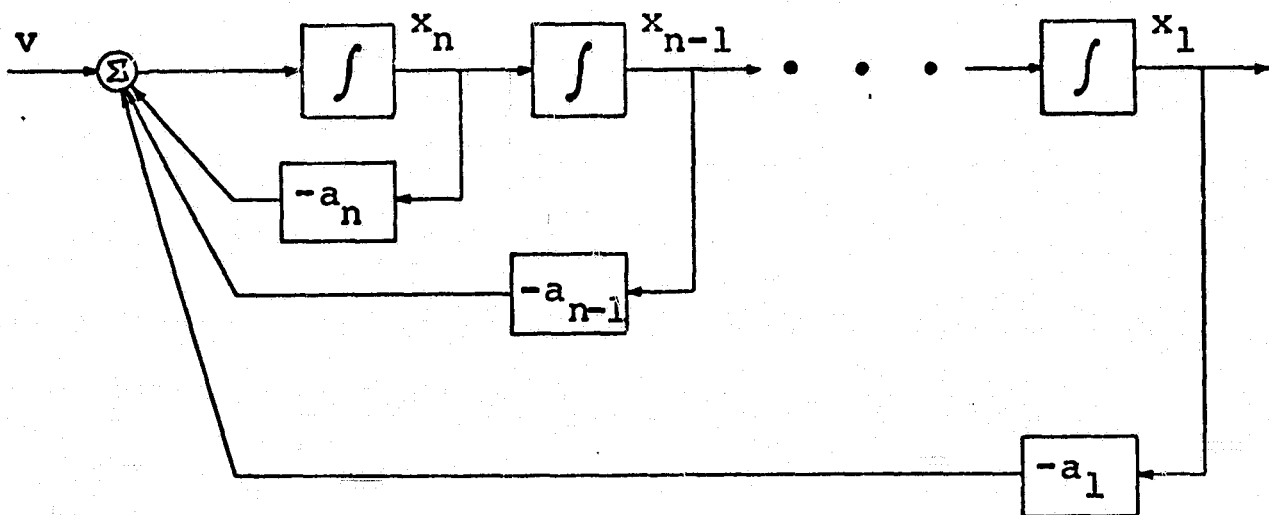


Figure 3-2. Single input phase variable canonical form given by (3-21) and (3-22).

The second canonical form has \underline{A} and $\underline{\tilde{B}}$ given as

$$\underline{A} = \begin{bmatrix} 0 & & & & -a_n \\ 1 & & & & \cdot \\ & \ddots & & & \cdot \\ & & \ddots & & \cdot \\ & & & 0 & -a_2 \\ & & & 1 & -a_1 \end{bmatrix}, \text{ all other entries zero,} \quad (3-23)$$

and

$$\underline{\tilde{B}} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ \vdots \\ 0 \end{bmatrix}, \quad (3-24)$$

which is illustrated as in Figure 3-3.

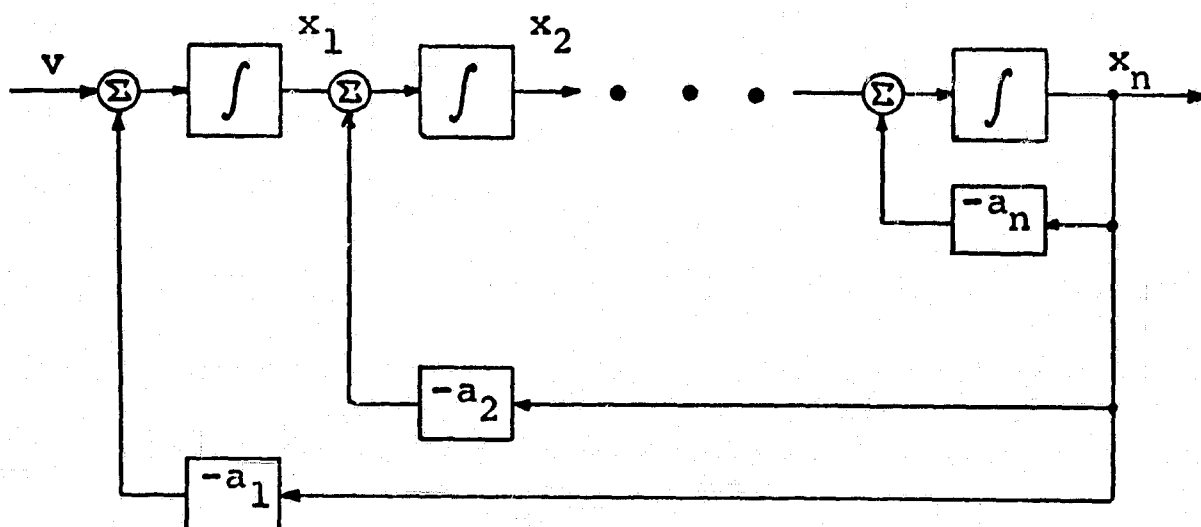


Figure 3-3. Single input phase variable canonical form given by (3-23) and (3-24).

Consequently, the system is viewed as a cascade connection of n -integrators with ordered feedback paths. Since zero and unit elements in \underline{A} and $\underline{\tilde{B}}$ are concentrated regularly, these canonical forms are conveniently used for abstract analyses of system theory. The second canonical form with (3-21) and (3-22), is particularly convenient for optimal control problems. Many papers have been published about reduction techniques for given systems into these canonical forms. [12,17-20]

3.3.2 Phase Variable Canonical Forms for Multi-Input Systems

3.3.2.1 Canonical Form by Tuel [21]

Tuel developed a canonical form, called a control canonical form, in which \underline{A} and $\underline{\tilde{B}}$ are decomposed as

$$\underline{A} = \begin{bmatrix} \underline{A}_{11} & \underline{A}_{12} \\ \underline{A}_{21} & \underline{A}_{22} \end{bmatrix} \quad (3-25)$$

and

$$\underline{\tilde{B}} = \begin{bmatrix} \underline{\tilde{B}}_{(1)} \\ \underline{\tilde{B}}_{(2)} \end{bmatrix} \quad (3-26)$$

such that

- (i) there exists a set of r positive integers ℓ_i , ($i = 1, 2, \dots, r$), (where r is the rank of $\underline{\hat{B}}$), defined

for a given system such that $\sum_{i=1}^r \ell_i = n$,

(ii) \underline{A}_{11} and \underline{A}_{12} are $(n-r) \times (n-r)$ and $(n-r) \times r$ matrices such that

$$\underline{A}_{11} = \begin{bmatrix} \underline{A}_{(1,1)} & & & \\ & \underline{A}_{(2,2)} & & \\ & & \ddots & \\ & & & \underline{A}_{(r,r)} \end{bmatrix}, \text{ all other entries zero,} \quad (3-27)$$

and

$$\underline{A}_{12} = \begin{bmatrix} \underline{E}_{(1)} & & & \\ & \underline{E}_{(2)} & & \\ & & \ddots & \\ & & & \underline{E}_{(r)} \end{bmatrix}, \text{ all other entries zero,} \quad (3-28)$$

where $\underline{A}_{(i,i)}$ is an $(\ell_i-1) \times (\ell_i-1)$ matrix

$$\underline{A}_{(i,i)} = \begin{bmatrix} 0 & 1 & & & \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ & & & \ddots & 1 \\ & & & & 0 \end{bmatrix}, \text{ all other entries zero,} \quad (3-29)$$

and $\underline{E}_{(i)}$ is an (ℓ_i-1) column vector

$$\underline{E}_{(i)} = \begin{bmatrix} 0 \\ \cdot \\ \cdot \\ \cdot \\ 0 \\ 1 \end{bmatrix}, \quad (3-30)$$

(iii) \underline{A}_{21} and \underline{A}_{22} are arbitrary $r \times (n-r)$ and $r \times r$ matrices,

(iv) $\widetilde{\underline{B}}_{(1)} = [0]$ and $\widetilde{\underline{B}}_{(2)}$ is an upper triangular matrix with unit elements on the diagonal.

When $m = r = 1$, this canonical form is reduced to the canonical form of the single input system, (3-21) and (3-22), i.e.,

$$\underline{A}_{11} = \underline{A}_{(1,1)} = \begin{bmatrix} 0 & 1 & & & \\ & \cdot & \cdot & & \\ & & \cdot & \cdot & \\ & & & \cdot & \cdot \\ & & & & 1 \\ & & & & 0 \end{bmatrix}, \text{ all other entries zero,} \quad (3-31)$$

$$\underline{A}_{12} = \underline{E}_{(1)} = \begin{bmatrix} 0 \\ \cdot \\ \cdot \\ \cdot \\ 0 \\ 1 \end{bmatrix}, \quad (3-32)$$

$$\widetilde{\underline{B}}_{(2)} = [1], \quad (3-33)$$

$\tilde{\underline{B}}_{(1)}$ is the $(n-1) \times 1$ null matrix, and \underline{A}_{21} and \underline{A}_{22} are $1 \times (n-1)$ and 1×1 matrices respectively.

This phase canonical form was developed early for solving multi-input optimum control system problems. It compactly describes the original system and illustrates its mathematical structure by an ordered array of zero and unity entries. Another comparable canonical form called an observation canonical form, can also be developed in this manner but is omitted here. [21]

3.3.2.2 Canonical Form by Luenberger [22]

Luenberger suggested a canonical form for linear multi-input systems (3-2), in which \underline{A} and $\tilde{\underline{B}}$ are decomposed into r^2 and r submatrices

$$\underline{A} = \begin{bmatrix} \underline{A}_{(1,1)} & \underline{A}_{(1,2)} & \cdot & \cdot & \cdot & \underline{A}_{(1,r)} \\ \cdot & & & & & \cdot \\ \cdot & & & & & \cdot \\ \cdot & & & & & \cdot \\ \underline{A}_{(r,1)} & \cdot & \cdot & \cdot & \cdot & \underline{A}_{(r,r)} \end{bmatrix} \quad (3-34)$$

and

$$\tilde{\underline{B}} = \begin{bmatrix} \tilde{\underline{B}}_{(1)} \\ \tilde{\underline{B}}_{(2)} \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \tilde{\underline{B}}_{(r)} \end{bmatrix} \quad (3-35)$$

such that

(i) there exists a set of r positive integers ℓ_i , ($i = 1, 2, \dots, r$), (r the rank of \hat{B}), such that

$$\sum_{i=1}^r \ell_i = n,$$

(ii) each $\underline{A}_{(i,j)}$, ($i, j = 1, 2, \dots, r$), is an $\ell_i \times \ell_j$ matrix such that

$$\underline{A}_{(i,i)} = \begin{bmatrix} 0 & 1 & & & \\ & \cdot & \cdot & & \\ & & \cdot & \cdot & \\ & & & \cdot & \cdot \\ & & & & 0 & 1 \\ -a_1 & \cdot & \cdot & \cdot & \cdot & -a_n \end{bmatrix}, \text{ all other entries zero,} \quad (3-36)$$

$$\underline{A}_{(i,j)} = \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ * & * & \cdot & \cdot & * \end{bmatrix}, \quad i \neq j, \text{ all other entries zero,} \quad (3-37)$$

where $*$'s indicate arbitrary elements,

(iii) each $\tilde{\underline{B}}_{(i)}$ ($i = 1, \dots, r$), is an $\ell_i \times r$ matrix given by

$$\tilde{\underline{B}}_{(i)} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}, \text{ all other entries zero.} \quad (3-38)$$

\uparrow
 i^{th} column

The form is illustrated in Figure 3-4. With this canonical form, a multi-input system can be viewed as interactions of r individual single-input systems, each with the canonical form given as (3-21) and (3-22). Consequently, the conveniences of the phase variable canonical form for single-input systems can be appreciated for the multi-input systems. Again a modification of this canonical form can be developed corresponding to the use of (3-23) and (3-24) instead of (3-21) and (3-22); this is omitted. [22]

3.3.2.3 Canonical Form by Asseo [23]

Asseo described a canonical form in which \underline{A} and $\tilde{\underline{B}}$ are decomposed into (3-25) and (3-26) such that the following are satisfied.

(i) \underline{A}_{11} is the $(n-r) \times r$ null matrix and \underline{A}_{12} is the $(n-r) \times (n-r)$ unit matrix,

(ii) \underline{A}_{21} and \underline{A}_{22} are $r \times r$ and $r \times (n-r)$ arbitrary matrices, and

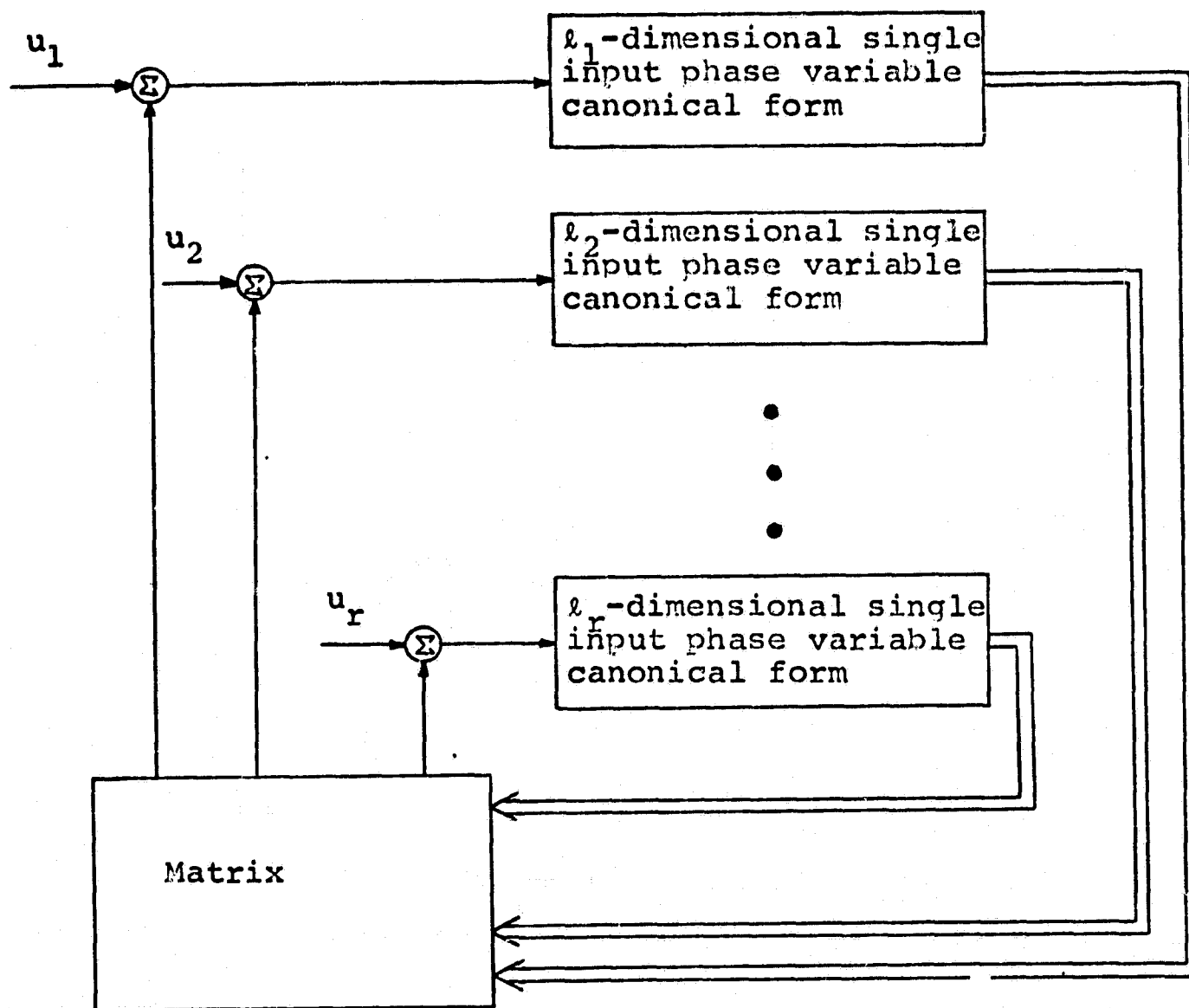


Figure 3-4. Luenberger canonical form.

(iii) $\widetilde{\underline{B}}_{(1)}$ is a $(n-r) \times r$ null matrix and $\widetilde{\underline{B}}_{(2)}$ is the $r \times r$ unit matrix.

The structure of this form appears simple and convenient for analysis in comparison with other suggested forms. Contrary to his assertion, however, the canonical form cannot be used for all systems which satisfy the two assumptions of Section 3.3.

As a counter example to illustrate this, consider a system (3-2) with $n = 4$, $r = 2$, i.e.,

$$\hat{\underline{A}} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \quad (3-39)$$

and

$$\hat{\underline{B}} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}. \quad (3-40)$$

It can be shown that this system is completely controllable by Theorem 2-9. From (3-4) let \underline{N}^{-1} for the equivalent reduction to \underline{A} and \underline{B} be

$$\underline{N}^{-1} = \begin{bmatrix} n_{11} & n_{12} & n_{13} & n_{14} \\ n_{21} & & & \cdot \\ \cdot & & & \cdot \\ n_{41} & \cdot & \cdot & n_{44} \end{bmatrix} \quad (3-41)$$

To obtain $\underline{\widetilde{B}}$ according to statement (iii) above, it follows that $n_{13} = n_{14} = n_{23} = n_{24} = n_{34} = n_{43} = 0$ and $n_{33} = n_{44} = 1$.

But in addition, from $\underline{A} = \underline{N}^{-1} \underline{\hat{A}} \underline{N}$, (3-41) and (i) above,

$$\underline{A} \underline{N}^{-1} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ * & * & * & * \\ * & * & * & * \end{bmatrix} \begin{bmatrix} n_{11} & n_{12} & 0 & 0 \\ n_{21} & n_{22} & 0 & 0 \\ n_{31} & n_{32} & 1 & 0 \\ n_{41} & n_{42} & 0 & 1 \end{bmatrix} = \begin{bmatrix} n_{31} & n_{32} & 1 & 0 \\ n_{41} & n_{42} & 0 & 1 \\ * & * & * & * \\ * & * & * & * \end{bmatrix} \quad (3-42)$$

and from (3-39)

$$\underline{N}^{-1} \underline{\hat{A}} = \begin{bmatrix} n_{11} & n_{12} & 0 & 0 \\ n_{21} & n_{22} & 0 & 0 \\ n_{31} & n_{32} & 1 & 0 \\ n_{41} & n_{42} & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & n_{11} & 0 & n_{12} \\ 0 & n_{22} & 0 & n_{22} \\ 0 & n_{31} & 0 & n_{32} \\ 0 & n_{41} & 0 & n_{42} \end{bmatrix} \quad (3-43)$$

where * indicates arbitrary elements.

The identity of (3-42) and (3-43) fails for the (1,3) element.

3.4 Comments

The canonical forms described in this chapter display the internal or elemental mathematical structure of systems from different viewpoints. The usefulness of the phase canonical forms for optimal control problems, however, could be improved if (i) a greater number of elements in \underline{A} and $\tilde{\underline{B}}$ were reduced to zero or unity, and (ii) these zero and unit elements were arranged in both a simple and unique order.

It is shown in the next chapter that the phase canonical form for single input systems is uniquely determined. It would appear, therefore, that this canonical form cannot be improved for this application. There are possibilities for improvement, however, when multi-input systems are considered. Both the Luenberger and Tuel canonical forms are valid for multi-input systems. While they have the same number of zero and unit elements in \underline{A} and $\tilde{\underline{B}}$, they both have some ambiguity or arbitrariness about the dimensions of the decomposed submatrices; these dimensions are not unique but depend upon the chosen matrix \underline{N} . While the Asseo canonical form has a particularly simplified structure and is free from submatrix dimensional

ambiguity, the application is for a more limited subclass of systems.

Applications of known phase canonical forms are also restricted to systems which are completely controllable and with full rank of \hat{B} . The possibility exists, therefore, to remove this limitation. The extension of canonical forms for use with nonlinear systems has also been avoided in past work.

The new canonical form given in the next chapter is developed to have a unique and regular distribution of zero and unit elements in A and \tilde{B} for linear systems, without the restrictive assumptions of controllability or on the rank of \hat{B} . The canonical form is also proposed for use with a class of nonlinear systems by application to the linear part of these systems.

Chapter 4

DEVELOPMENT OF A CANONICAL FORM

In this chapter, a new phase variable canonical form is developed for a class of multi-input systems. This particular canonical form is shown to be superior to those reviewed in the previous chapter. It also provides the form for the analysis of the inverse problem of the optimal regulator in subsequent chapters. First a new definition of an equivalent system is given. Section 4.2 introduces two theorems on matrix transformations. Based on these ideas, the new canonical form is then presented in Section 4.3. The uniqueness of the structure of the canonical form for linear systems is discussed, and finally the canonical form is extended for applications to a class of nonlinear systems.

Throughout this chapter only, matrices are indicated by capital Roman or Greek letters without the underline and are constant unless otherwise noted. Vectors are underlined.

4.1 Equivalent System

Systems which are considered in this chapter are given by

$$\dot{\underline{Y}} = \hat{\underline{F}}(\underline{Y}) + \hat{\underline{B}} \underline{V} \quad (4-1)$$

where $\hat{\underline{F}}(\underline{Y}) = [\hat{f}_1(\underline{Y}), \hat{f}_2(\underline{Y}), \dots, \hat{f}_n(\underline{Y})]^T$, \hat{B} is an $n \times m$ matrix with $0 < m \leq n$ and the rank \hat{B} is r , $0 < r \leq m$.

One definition of an equivalent system is given by

Definition 3-1. A more rigorous statement is possible, however, for the particular class of systems given by (4-1).

Definition 4-1: Equivalent System.

For two arbitrary systems

$$\begin{aligned} \dot{\underline{Y}} &= \hat{\underline{F}}(\underline{Y}) + \hat{B} \underline{V}, & S_1 \\ \dot{\underline{X}} &= \underline{F}(\underline{X}) + B \underline{U}, & S_2 \end{aligned} \quad (4-2)$$

S_1 is said to be equivalent to S_2 if there exist non-singular matrices N and M which are $n \times n$ and $m \times m$ respectively, satisfying

$$\begin{aligned} \underline{X} &= N^{-1} \underline{Y} \\ \underline{U} &= M^{-1} \underline{V}, \end{aligned} \quad (4-3)$$

or equivalently

$$\begin{aligned} \underline{F}(\underline{X}) &= N^{-1} \hat{\underline{F}}(N \underline{X}) \\ B &= N^{-1} \hat{B} M. \end{aligned} \quad (4-4)$$

The nonsingular transformation (4-3) can be identified as a topological equivalent relation^[24] defined for this class of systems. The reflective law is satisfied directly by choosing M and N to be unit matrices. The symmetric law is satisfied by substituting

$$\begin{aligned}\underline{Y} &= N \underline{X} \\ \underline{V} &= M \underline{U}\end{aligned}\tag{4-5}$$

into (4-4). To show that the transitive law is also satisfied, consider a system

$$\dot{\underline{Z}} = \underline{\tilde{F}}(\underline{Z}) + \underline{\tilde{B}} \underline{W}, \quad S_3,\tag{4-6}$$

to which S_2 is equivalent. Thus there exist nonsingular matrices $\underline{\tilde{N}}$ and $\underline{\tilde{M}}$ such that

$$\begin{aligned}\underline{Z} &= \underline{\tilde{N}}^{-1} \underline{X} \\ \underline{W} &= \underline{\tilde{M}}^{-1} \underline{U}.\end{aligned}\tag{4-7}$$

S_1 is made an equivalent system of S_3 by defining

$$\begin{aligned}\underline{Z} &= (N \underline{\tilde{N}})^{-1} \underline{Y} \\ \underline{W} &= (M \underline{\tilde{M}})^{-1} \underline{V},\end{aligned}\tag{4-8}$$

where both $(N \widetilde{N})$ and $(M \widetilde{M})$ are nonsingular.

If M in (4-3) is fixed as the unit matrix, Definition 4-1 is essentially reduced to Definition 3-1. Thus, Definition 4-1 is a generalization of Definition 3-1, but with application restricted to the class of systems given by (4-1). The flexibility of the additional matrix M , however, allows a more compact canonical form to be given for the class of systems (4-1).

4.2 Fundamental Theorems

Theorem 4-1:

Let A be an $n \times n$ matrix. For each positive integer $s \leq n$, there exists an $n \times n$ nonsingular matrix N which transforms A into

$$A \triangleq N^t A N \triangleq \begin{bmatrix} A_{(1,1)} & A_{(1,2)} & \cdots & A_{(1,v)} \\ \vdots & & & \vdots \\ A_{(v,1)} & \cdot & \cdots & A_{(v,v)} \end{bmatrix}$$

$$= \begin{bmatrix} A_{(1,1)} & A_{(1,2)} & [0] & \cdot & \cdot & \cdot & [0] \\ & & & \cdot & & & \cdot \\ [0] & [0] & A_{(2,3)} & \cdot & & & \cdot \\ \cdot & & \cdot & \cdot & & & \cdot \\ \cdot & & \cdot & \cdot & \cdot & & [0] \\ \cdot & & \cdot & \cdot & \cdot & \cdot & \\ [0] & \cdot & \cdot & \cdot & [0] & A_{(v-1,v)} \\ A_{(v,1)} & \cdot & \cdot & \cdot & \cdot & A_{(v,v)} \end{bmatrix}, \quad (4-9)$$

where each submatrix $A_{(i,j)}$, $(i,j = 1,2,\dots,v)$, is $\ell_i \times \ell_j$ such that

$$(i) \quad \ell_v = s; \quad (4-10)$$

(ii) $v, \ell_1, \ell_2, \dots, \ell_{v-1}$ are positive integers dependent on A and s such that

$$\sum_{i=1}^v \ell_i = n; \quad (4-11)$$

(iii) if $v \geq 2$, then either

(a) $A_{(1,1)}$ is a Jordan canonical form and

$$A_{(1,2)} = [0] , \text{ or} \quad (4-12)$$

$$(b) \quad A_{(1,1)} = [0] \quad (4-13)$$

$$A_{(1,2)} = \begin{bmatrix} [0] & I_{\ell_1} \end{bmatrix} , \quad (4-14)$$

where, if (b), then $\ell_1 \leq \ell_2$;

(iv) in addition, if $v \geq 3$, then

$$(a) \quad \ell_i \leq \ell_j \quad \text{for } 2 \leq i \leq j \leq v , \quad (4-15)$$

$$(b) \quad A_{(i,i+1)} = \begin{bmatrix} [0] & I_{\ell_i} \end{bmatrix} \quad \text{for } 2 \leq i \leq v-1 , \quad (4-16)$$

$$(c) \quad A_{(v,i)} , \quad \text{for } i = 1, 2, \dots, v, \text{ are unspecified.}$$

The proof of the theorem follows three lemmas.

Lemma 4-1:

(a) If H is a $\theta_1 \times \theta_2$ matrix with rank θ_3 , there exist nonsingular matrices H_1 and H_2 which are $\theta_1 \times \theta_1$ and $\theta_2 \times \theta_2$ respectively, satisfying

$$H_1 H H_2 = \begin{bmatrix} [0] & [0] \\ [0] & I_{\theta_3} \end{bmatrix}. \quad (4-17)$$

(b) If H is full rank, i.e., $\theta_3 = \text{Min}(\theta_1, \theta_2)$, and $\theta_1 < \theta_2$ ($\theta_1 > \theta_2$) is satisfied, $H_1(H_2)$ can be the unit matrix.

(c) If $\theta_1 = \theta_2$, then either H_1 or H_2 can be the unit matrix.

Proof: The reduction of H to (4-17) follows directly from the equivalence of matrices.^[10] To show part (b) of the lemma, assume $\theta_1 < \theta_2$. Since $\theta_3 = \theta_1$, H has θ_1 independent columns. Let $H_2^{(1)}$ be a $\theta_2 \times \theta_2$ nonsingular matrix resulting from an interchange of columns of H such that

$$H H_2^{(1)} = \begin{bmatrix} H_3 & H_4 \end{bmatrix}, \quad (4-18)$$

where H_4 is a $\theta_1 \times \theta_1$ nonsingular matrix composed of θ_1 independent columns of H . Define a $\theta_2 \times \theta_2$ nonsingular matrix such that

$$H_2^{(2)} \triangleq \begin{bmatrix} I_{\theta_2 - \theta_1} & [0] \\ -H_4^{-1} H_3 & H_4^{-1} \end{bmatrix} \quad (4-19)$$

and

$$H_2 \triangleq H_2^{(1)} H_2^{(2)} .$$

Then it follows that

$$H H_2 = \begin{bmatrix} [0] & I_{\theta_1} \end{bmatrix} \quad (4-20)$$

and H_1 is the unit matrix. If $\theta_2 < \theta_1$, a similar method can be applied. If $\theta_1 = \theta_2 = \theta_3$, then either H_1 or H_2 can be H^{-1} .

Lemma 4-2:

Let H_1 and H_2 be $\theta_2 \times \theta_1$ and $\theta_2 \times \theta_3$ matrices respectively with

$$0 \leq \theta_2 \leq \theta_3 , \quad (4-21)$$

and let the rank of H_2 be θ_2 . Then for any $\theta_1 > 0$, there exists a $\theta_3 \times \theta_1$ matrix K satisfying

$$H_1 = H_2 K . \quad (4-22)$$

Proof: The proof is complete if a construction of K can be demonstrated. Since H_2 is of full rank, from (4-21) and Lemma 4-1 there exists a $\theta_3 \times \theta_3$ nonsingular matrix K_1 such that

$$H_2 K_1 = \begin{bmatrix} [0] & I_{\theta_2} \end{bmatrix} . \quad (4-23)$$

Define K_2 to be a $\theta_3 \times \theta_1$ matrix such that

$$K_2 \triangleq \begin{bmatrix} [0] \\ H_1 \end{bmatrix} . \quad (4-24)$$

and define

$$K \triangleq K_1 K_2 . \quad (4-25)$$

Then K is $\theta_2 \times \theta_1$ and satisfies

$$H_2 K = H_1 . \quad (4-26)$$

Lemma 4-3:

Let A be an $n \times n$ matrix. For each positive integer $s \leq n$, there exists an $n \times n$ nonsingular matrix \hat{N} to transform A into

$$\tilde{A} \triangleq \hat{N}^{-1} A \hat{N} \triangleq \begin{bmatrix} \tilde{A}_{(1,1)} & \tilde{A}_{(1,2)} & \cdot & \cdot & \cdot & \tilde{A}_{(1,v)} \\ \tilde{A}_{(2,1)} & & & & & \cdot \\ \cdot & & & & & \cdot \\ \cdot & & & & & \cdot \\ \tilde{A}_{(v,1)} & \cdot & \cdot & \cdot & \cdot & \tilde{A}_{(v,v)} \end{bmatrix}$$

$$= \begin{bmatrix} \tilde{A}_{(1,1)} & \tilde{A}_{(1,2)} & [0] & \cdot & \cdot & \cdot & [0] \\ \cdot & & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & & \cdot & \cdot & \cdot & \cdot & [0] \\ \cdot & & & & & & \tilde{A}_{(v-1,v)} \\ \cdot & & & & & & \tilde{A}_{(v,v)} \\ \tilde{A}_{(v,1)} & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{bmatrix}, \quad (4-27)$$

where each submatrix $\tilde{A}_{(i,j)}$, $(i,j = 1,2,\dots,v)$, is $\ell_i \times \ell_j$ such that

$$(i) \quad \ell_v = s; \quad (4-28)$$

(ii) $v, \ell_1, \ell_2, \dots, \ell_{v-1}$ are positive integers dependent on A and s such that $\sum_{i=1}^v \ell_i = n$;

(iii) if $v \geq 2$, then either

$$(a) \quad \tilde{A}_{(1,2)} = [0] , \text{ or} \quad (4-29)$$

$$(b) \quad \tilde{A}_{(1,2)} = \begin{bmatrix} [0] & I_{\theta_1} \end{bmatrix} , \quad (4-30)$$

where, for (b) it follows that $\ell_1 \leq \ell_2$; (4-31)

(iv) in addition, if $v \geq 3$,

$$(a) \quad \ell_i \leq \ell_j , \text{ for } 2 \leq i \leq j \leq v , \quad (4-32)$$

$$(b) \quad \tilde{A}_{(i,i+1)} = \begin{bmatrix} [0] & \hat{N}_{ii}^{(i)-1} \end{bmatrix} \text{ for } 2 \leq i \leq v-1 , \quad (4-33)$$

where $\hat{N}_{ii}^{(i)-1}$ is a $\ell_i \times \ell_i$ nonsingular matrix,

(c) all other submatrices are unspecified.

Proof: (1) First, if $s = n$, define

$$v = 1$$

(4-34)

$$\ell_1 = s .$$

The statement is satisfied by $\hat{N} = I$.

(2) If $s < n$, A can be decomposed as

$$A \triangleq \left[\begin{array}{c|c} \overbrace{A_{22}}^{n-s} & \overbrace{A_{21}}^s \\ \hline \overbrace{A_{12}}^{n-s} & \overbrace{A_{11}}^s \end{array} \right] \begin{array}{l} \}^{n-s} \\ \}^s \end{array}, \quad (4-35)$$

where each A_{ij} , $(i, j = 1, 2)$, has the dimensions indicated.

According to Lemma 4-1 $\hat{N}_{11}^{(1)}$ and $\hat{N}_{22}^{(1)}$ are $s \times s$

and $(n-s) \times (n-s)$ nonsingular matrices satisfying

$$\hat{N}_{22}^{(1)} A_{21} \hat{N}_{11}^{(1)} = \begin{bmatrix} [0] & [0] \\ [0] & I_{n_1} \end{bmatrix}, \quad (4-36)$$

where n_1 is the rank of A_{21} . If $n_1 = s \leq n-s$, $\hat{N}_{11}^{(1)}$

can be the unit matrix, according to Lemma 4-1. Define

an $n \times n$ nonsingular matrix such that

$$\hat{N}^{(1)} \triangleq \begin{bmatrix} \hat{N}_{22}^{(1)-1} & [0] \\ [0] & \hat{N}_{11}^{(1)} \end{bmatrix}. \quad (4-37)$$

Then

$$A^{(1)} \triangleq \hat{N}^{(1)-1} A \hat{N}^{(1)} = \begin{bmatrix} \hat{N}_{22}^{(1)} A_{22} \hat{N}_{22}^{(1)-1} & \hat{N}_{22}^{(1)} A_{21} \hat{N}_{11}^{(1)} \\ \hat{N}_{11}^{(1)-1} A_{12} \hat{N}_{22}^{(1)-1} & \hat{N}_{11}^{(1)-1} A_{11} \hat{N}_{11}^{(1)} \end{bmatrix}.$$

(4-38)

(3) As a subclass of (2), if $n_1 = 0$ or $n_1 = n - s \leq s$ in (4-36),

$$\hat{N}_{22}^{(1)} A_{21} \hat{N}_{11}^{(1)} = \begin{cases} [0] ; & \text{if } n_1 = 0 \\ \begin{bmatrix} [0] & I_{n_1} \end{bmatrix} ; & \text{if } n_1 = n - s \leq s . \end{cases} \quad (4-39)$$

(4-40)

Then define

$$v = 2$$

$$l_1 = n - s$$

$$l_2 = s$$

(4-41)

$$\hat{N} = \hat{N}^{(1)}.$$

It follows that

$$\tilde{A} = \begin{bmatrix} \tilde{A}_{(1,1)} & \tilde{A}_{(1,2)} \\ \tilde{A}_{(2,1)} & \tilde{A}_{(2,2)} \end{bmatrix}, \quad (4-42)$$

where $\tilde{A}_{(1,2)}$ is either (4-39) or (4-40). For (4-40), it follows that

$$\ell_1 = \text{rank of } \tilde{A}_{(1,2)} = \text{Min}(\ell_1, \ell_2) \leq \ell_2. \quad (4-43)$$

Therefore the assertion is proved for the case of $s < n$ and $v = 2$.

(4) The remaining subclass of (2) is $0 < n_1 < n-s$.

Again $A^{(1)}$ can be decomposed, using (4-36), as

$$A^{(1)} = \begin{bmatrix} A_{33}^{(1)} & A_{32}^{(1)} & A_{31}^{(1)} \\ A_{23}^{(1)} & A_{22}^{(1)} & A_{21}^{(1)} \\ A_{13}^{(1)} & A_{12}^{(1)} & A_{11}^{(1)} \end{bmatrix} = \begin{bmatrix} A_{33}^{(1)} & A_{32}^{(1)} & [0] \\ A_{23}^{(1)} & A_{22}^{(1)} & A_{21}^{(1)} \\ A_{13}^{(1)} & A_{12}^{(1)} & A_{11}^{(1)} \end{bmatrix} \left\{ \begin{array}{l} n-s-n_1 \\ n_1 \\ s \end{array} \right\}, \quad (4-44)$$

$\underbrace{\hspace{1.5cm}}_{(n-s-n_1)} \quad \underbrace{\hspace{1.5cm}}_{n_1} \quad \underbrace{\hspace{1.5cm}}_s$

where each $A_{ij}^{(1)}$ has the dimension indicated. It also

follows from (4-36) that

$$A_{21}^{(1)} = \begin{bmatrix} [0] & I_{n_1} \end{bmatrix}. \quad (4-45)$$

Using the previous process, let $\hat{N}_{22}^{(2)}$ and $\hat{N}_{33}^{(2)}$ be $n_1 \times n_1$ and $(n-s-n_1) \times (n-s-n_1)$ nonsingular matrices satisfying

$$\hat{N}_{33}^{(2)} A_{32}^{(1)} \hat{N}_{22}^{(2)} = \begin{bmatrix} [0] & [0] \\ [0] & I_{n_2} \end{bmatrix}, \quad (4-46)$$

where n_2 is the rank of $A_{32}^{(1)}$. If $n_2 = n_1 \leq n - s - n_1$, $\hat{N}_{22}^{(2)}$ can be the unit matrix according to Lemma 4-1.

Define an $n \times n$ nonsingular matrix

$$N^{(2)} \triangleq \begin{bmatrix} \hat{N}_{33}^{(2)-1} & [0] & [0] \\ [0] & \hat{N}_{22}^{(2)} & [0] \\ [0] & [0] & I_s \end{bmatrix} \quad (4-47)$$

and calculate

$$A^{(2)} \triangleq \hat{N}^{(2)-1} A^{(1)} \hat{N}^{(2)}$$

$$= \begin{bmatrix} \hat{N}_{33}^{(2)} A_{33}^{(1)} \hat{N}_{33}^{(2)-1} & \hat{N}_{33}^{(2)} A_{32}^{(1)} \hat{N}_{22}^{(2)} & [0] \\ \hat{N}_{22}^{(2)-1} A_{23}^{(1)} \hat{N}_{33}^{(2)-1} & \hat{N}_{22}^{(2)-1} A_{22}^{(1)} \hat{N}_{22}^{(2)} & \begin{bmatrix} [0] & \hat{N}_{22}^{(2)-1} \end{bmatrix} \\ A_{13}^{(1)} \hat{N}_{33}^{(2)-1} & A_{12}^{(1)} \hat{N}_{22}^{(2)} & A_{11}^{(1)} \end{bmatrix},$$

(4-48)

where (4-45) is used.

(5) A subclass of (4) is $n_2 = 0$ or $n_2 = n - s - n_1 \leq n_1$ in (4-48). From Lemma 4-1,

$$\hat{N}_{33}^{(2)} A_{32}^{(1)} \hat{N}_{22}^{(2)} = \begin{cases} [0] ; & \text{if } n_2 = 0 \\ \begin{bmatrix} [0] & I_{n_2} \end{bmatrix} ; & \text{if } n_2 = n - s - n_1 \leq n_1 . \end{cases}$$

(4-49)

(4-50)

Then define

$$v = 3$$

$$l_1 = n_2 = n - n_1 - s$$

$$l_2 = n_1$$

(4-51)

$$l_3 = s$$

$$\hat{N} = \hat{N}^{(1)} \hat{N}^{(2)}$$

and

$$\tilde{A} \triangleq A^{(2)} \triangleq \begin{bmatrix} \tilde{A}_{(1,1)} & \tilde{A}_{(1,2)} & [0] \\ \tilde{A}_{(2,1)} & \tilde{A}_{(2,2)} & \tilde{A}_{(2,3)} \\ \tilde{A}_{(3,1)} & \tilde{A}_{(3,2)} & \tilde{A}_{(3,3)} \end{bmatrix}, \quad (4-52)$$

where $\tilde{A}_{(1,2)}$ is either (4-49) or (4-50) and

$$\tilde{A}_{(2,3)} = \begin{bmatrix} [0] & \hat{N}_{22}^{(2)-1} \end{bmatrix} \quad (4-53)$$

with $\hat{N}_{22}^{(2)-1}$ nonsingular. With respect to the rank of $\tilde{A}_{(2,3)}$

$$l_2 = \text{rank of } \tilde{A}_{(2,3)} = \text{Min}(l_2, l_3) \leq l_3. \quad (4-54)$$

If, in addition, (4-50) is satisfied, it is necessary from

the rank of $\tilde{A}_{(1,2)}$ that

$$\ell_1 \leq \ell_2. \quad (4-55)$$

Therefore the assertion is proved for the case of $s < n$ and $v = 3$.

(6) This process can be continued for the remaining subclass of (4). Consider, instead, a general description of $A^{(k)}$ as a transformed matrix from $A^{(k-1)}$ by a nonsingular matrix $\hat{N}^{(k)}$ such that

$$A^{(k)} = \hat{N}^{(k)-1} A^{(k-1)} \hat{N}^{(k)} \quad (4-56)$$

with

$$A^{(k-1)} = \begin{bmatrix} A_{k+1 \ k+1}^{(k-1)} & \cdot & \cdot & \cdot & A_{k-1 \ 1}^{(k-1)} \\ \cdot & & & & \cdot \\ \cdot & & & & \cdot \\ \cdot & & & & \cdot \\ A_{1 \ k+1}^{(k-1)} & \cdot & \cdot & \cdot & A_{11}^{(k-1)} \end{bmatrix} \quad (4-57)$$

and

$$\hat{N}^{(k)} = \begin{bmatrix} \hat{N}_{k+1 \ k+1}^{(k)-1} & & \\ & \hat{N}_{kk}^{(k)} & \\ & & I_{n-h_1^{(k)}-h_2^{(k)}} \end{bmatrix},$$

all other entries zero,

(4-58)

where

(a) each $A_{ij}^{(k-1)}$, $(i, j = 1, 2, \dots, k+1)$, is an $h_i^{(k-1)} \times h_j^{(k-1)}$,

$h_1^{(k-1)} = s$, $h_2^{(k-1)} = n_1$, \dots , $h_h^{(k-1)} = n_{k-1}$, and

$h_{h+1}^{(k-1)} = n - s - \sum_{i=1}^{k-1} n_i$ such that

$$\text{rank } A_{i \ i-1}^{(k-1)} = n_{i-1}, \quad \text{for } 2 \leq i \leq k. \quad (4-59)$$

(b) $\hat{N}_{kk}^{(k)}$ and $\hat{N}_{k+1 \ k+1}^{(k)}$ are nonsingular matrices which are $n_{k-1} \times n_{k-1}$ and $(n - s - \sum_{i=1}^{k-1} n_i) \times (n - s - \sum_{i=1}^{k-1} n_i)$

such that

$$\hat{N}_{k+1, k+1}^{(k)} A_{k+1, k}^{(k-1)} \hat{N}_{kk}^{(k)} = \begin{bmatrix} [0] & [0] \\ [0] & I_{n_k} \end{bmatrix}, \quad (4-60)$$

where the rank of $A_{k+1, k}^{(k-1)}$ is n_k . If $n_k = n_{k-1} \leq n - s - \sum_{i=1}^{k-1} n_i$,

$\hat{N}_{kk}^{(k)}$ can be the unit matrix, from Lemma 4-1. According to the preceding transformations and the decomposition of $A^{(k)}$, it must follow that

$$\begin{aligned} A_{ij}^{(k-1)} &= [0] \quad \text{for } 1 \leq j \leq i+2 \leq k+2 \\ A_{k, k-1}^{(k-1)} &= \begin{bmatrix} [0] & I_{n_{k-1}} \end{bmatrix} \quad (4-61) \\ A_{i+1, i}^{(k-1)} &= \begin{bmatrix} [0] & \hat{N}_{ii}^{(i)-1} \end{bmatrix} \quad \text{for } i = 2, 3, \dots, k-2, \end{aligned}$$

that is, (4-57) becomes

$$A^{(k-1)} = \begin{bmatrix} A_{k+1, k+1}^{(k-1)} & A_{k+1, k}^{(k-1)} & & & \\ \cdot & \cdot & & & \\ & [0] \ I_{n_k} & & & \\ \cdot & \cdot & \cdot & & \\ & \cdot & [0] \ \hat{N}_{k+1, k-1}^{(k-1)-1} & & \\ \cdot & \cdot & \cdot & \cdot & \\ & \cdot & \cdot & \cdot & \\ & & & [0] \ \hat{N}_{22}^{(2)-1} & \\ A_{1, k+1}^{(k-1)} & \cdot & \cdot & \cdot & A_{11}^{(k-1)} \end{bmatrix}$$

all other entries zero,

(4-62)

A(k)

75

all other entries zero.

Repeating the transformation n times, either the rank of

$A_{n+1\ n}^{(n-1)}$ becomes

(a) zero, when it follows that

$$\hat{N}_{n+1\ n+1}^{(n)} A_{n+1\ n}^{(n-1)} \hat{N}_{nn}^{(n)} = [0] , \quad (4-64)$$

or (b) full, i.e., the rank $A_{n+1\ n}^{(n-1)} = n - s - \sum_{i=1}^{n-1} n_i$,

when it follows that

$$\hat{N}_{n+1\ n+1}^{(n)} A_{n+1\ n}^{(n-1)} \hat{N}_{nn}^{(n)} = \begin{bmatrix} [0] & I_{n_n} \end{bmatrix} . \quad (4-65)$$

Define

$$v = n + 1$$

$$l_1 = n - s - \sum_{i=1}^{n-1} n_i$$

$$l_2 = n_{v-2}$$

$$\vdots$$

$$l_{v-1} = n_1$$

$$l_v = s$$

(4-66)

$$\hat{N} = \hat{N}^{(1)} \hat{N}^{(2)} \dots \hat{N}^{(v-1)} ,$$

where $\ell_1, \ell_2, \dots, \ell_v$ are positive integers and \hat{N} is non-singular. With $k = n$ in (4-63), and

$$A^{(n)} = A^{(v-1)} = \tilde{A}, \quad (4-67)$$

\tilde{A} given by (4-27), it follows that

$$(a) \quad \tilde{A}_{(1,2)} = \begin{cases} [0], & \text{from (4-64), or} \\ \begin{bmatrix} [0] & I_{\ell_1} \end{bmatrix}, & \text{from (4-65),} \end{cases} \quad (4-68)$$

where, for the second case, from the rank of $\tilde{A}_{(1,2)}$,

$$\ell_1 \leq \ell_2 ;$$

$$(b) \quad \tilde{A}_{(i,i+1)} = \begin{bmatrix} [0] & \hat{N}_{v+1-i, v+1-i}^{(v+1-i)^{-1}} \end{bmatrix}, \quad \text{for } i = 2, 3, \dots, v-1, \quad (4-69)$$

where each $\hat{N}_{v+1-i, v+1-i}^{(v+1-i)^{-1}}$ is nonsingular;

(c) according to the rank of (4-69),

$$\ell_i \leq \ell_{i+1}, \quad \text{for } i = 2, 3, \dots, v-1 ; \quad (4-70)$$

$$(d) \quad \tilde{A}_{(i,j)} = [0], \quad \text{for } 2 \leq i \leq j+2 \leq v+2. \quad (4-71)$$

Thus the lemma is proved.

From the structure of each $\hat{N}^{(i)}$ in (4-58), \hat{N} in (4-66) is

$$\hat{N} = \hat{N}^{(1)} \hat{N}^{(2)} \dots \hat{N}^{(v-1)} = \begin{bmatrix} \hat{N}_{22}^{(1)-1} & & & \\ & \hat{N}_{33}^{(2)-1} & & \\ & & \hat{N}_{22}^{(2)} & \\ & & & \hat{N}_{11}^{(1)} \\ & & & & \ddots & & \hat{N}_{33}^{(3)} & \\ & & & & & & & \hat{N}_{44}^{(3)-1} & \\ & & & & & & & & I_{\ell_v} & \\ & & & & & & & & & I_{\ell_v + \ell_{v-1}} & \end{bmatrix}$$

$$\begin{bmatrix} \hat{N}_{vv}^{(v-1)} & & & \\ \dots & \hat{N}_{v-1, v-1}^{(v-1)} & & \\ & & \hat{N}_{11} & \\ & & \hat{N}_{21} & \\ & & & \hat{N}_{22} \end{bmatrix} \triangleq \begin{bmatrix} [0] \\ \hat{N}_{22} \end{bmatrix}, \quad (4-72)$$

where \hat{N}_{11} and \hat{N}_{22} are $(n-s) \times (n-s)$ and $s \times s$ respectively.

Proof of Theorem 4-1:

It is sufficient to show the existence of an $n \times n$ nonsingular matrix \tilde{N} to similarly transform \tilde{A} as given in Lemma 4-3 such that

$$A = \tilde{N}^{-1} \tilde{A} \tilde{N}, \quad (4-73)$$

where A is given by the statement of the theorem. Then N of the theorem is given by

$$N = \hat{N} \tilde{N}. \quad (4-74)$$

(1) If $v = 1$, let $\tilde{N} = I_n$.

(2) Assume $v \geq 2$. For a null $\tilde{A}_{(1,2)}$, (4-29), define

$$\tilde{N}^{(1)} = \begin{bmatrix} \tilde{N}_{11}^{(1)} & [0] \\ [0] & I_{n-\ell_1} \end{bmatrix}, \quad (4-75)$$

where $\tilde{N}_{11}^{(1)}$ is $\ell_1 \times \ell_1$ nonsingular to provide a Jordan

canonical form for $\tilde{N}_{11}^{(1)-1} \tilde{A}_{(1,1)} \tilde{N}_{11}^{(1)}$. For $\tilde{A}_{(1,2)}$

given by (4-30), define

$$\tilde{N}^{(1)} = \begin{bmatrix} I_{\ell_1} & [0] & & \\ & \tilde{N}_{21}^{(1)} & \tilde{N}_{22}^{(1)} & [0] \\ & & & \\ & [0] & & I_{n-\ell_1-\ell_2} \end{bmatrix}, \quad (4-76)$$

where $\tilde{N}_{22}^{(1)} = I_{\ell_2}$ and $\tilde{N}_{21}^{(1)}$ is an $\ell_2 \times \ell_1$ matrix satisfying

$$\tilde{A}_{(1,1)} + \tilde{A}_{(1,2)} \tilde{N}_{21}^{(1)} = [0]. \quad (4-77)$$

This $\tilde{N}_{21}^{(1)}$ exists from Lemma 4-2, where $\tilde{A}_{(1,2)}$ is of full rank and $\ell_1 \leq \ell_2$. Define

$$\tilde{A}^{(1)} = \tilde{N}^{(1)-1} \tilde{A} \tilde{N}^{(1)} \triangleq \begin{bmatrix} \tilde{A}_{(1,1)}^{(1)} & \tilde{A}_{(1,2)}^{(1)} & & \\ \cdot & \cdot & \cdot & \\ \cdot & \cdot & \cdot & \\ \cdot & \cdot & \cdot & \\ \tilde{A}_{(v,1)}^{(1)} & \cdot & \cdot & \cdot \\ & & & \tilde{A}_{(v-1,v)}^{(1)} \\ & & & \tilde{A}_{(v,v)}^{(1)} \end{bmatrix},$$

all other entries zero,

(4-78)

where (i) each $\tilde{A}_{(i,j)}^{(1)}$ is an $\ell_i \times \ell_j$ matrix;

(ii) either

(a) for (4-29),

$$\tilde{A}_{(1,1)}^{(1)} = \tilde{N}_{11}^{(1)-1} \tilde{A}_{(1,1)} \tilde{N}_{11}^{(1)}, \text{ a Jordan canonical form}$$

$$\tilde{A}_{(1,2)}^{(1)} = [0] \quad (4-79)$$

or

(b) for (4-30)

$$\tilde{A}_{(1,1)}^{(1)} = \tilde{A}_{(1,1)} + \tilde{A}_{(1,2)} \tilde{N}_{21}^{(1)} = [0],$$

(4-80)

$$\tilde{A}_{(1,2)}^{(1)} = \begin{bmatrix} [0] & I_{\ell_1} \end{bmatrix};$$

$$(iii) \quad \tilde{A}_{(i,i+1)}^{(1)} = \tilde{A}_{(i,i+1)} = \begin{bmatrix} [0] & \hat{N}_{v+1-i, v+1-i}^{(v+1-i)-1} \end{bmatrix},$$

$$\text{for } i = 2, 3, \dots, v-1, \quad (4-81)$$

as given by (4-27) and (4-69). Therefore the theorem is proved for $v = 2$ by defining $\tilde{N} = \tilde{N}^{(1)}$.

(3) If $v \geq 3$, consider

$$\tilde{N}^{(2)} \triangleq \begin{bmatrix} I_{\ell_1} & & & \\ & I_{\ell_2} & & \\ \tilde{N}_{31}^{(2)} & \tilde{N}_{32}^{(2)} & \tilde{N}_{33}^{(2)} & \\ & & & I_{n-\ell_1-\ell_2-\ell_3} \end{bmatrix},$$

all other entries zero,

(4-82)

where $\tilde{N}_{3i}^{(2)}$ are chosen to satisfy

$$\tilde{A}_{(2,1)}^{(1)} + \tilde{A}_{(2,3)}^{(1)} \tilde{N}_{31}^{(2)} = [0],$$

$$\tilde{A}_{(2,2)}^{(1)} + \tilde{A}_{(2,3)}^{(1)} \tilde{N}_{32}^{(2)} = [0],$$

(4-83)

$$\tilde{N}_{33}^{(2)} = \begin{bmatrix} I_{\ell_3-\ell_2} & [0] \\ [0] & \hat{N}_{v-1, v-1}^{(v-1)} \end{bmatrix}$$

and $\hat{N}_{v-1, v-1}^{(v-1)}$ is given by (4-81). $\tilde{N}_{21}^{(2)}$ and $\tilde{N}_{32}^{(2)}$ exist

from Lemma 4-2. Since $\tilde{N}_{33}^{(2)}$ is nonsingular, $\tilde{N}^{(2)}$ is

nonsingular and its inverse is

$$\tilde{N}^{(2)-1} = \begin{bmatrix} I_{\ell_1} & & & \\ & I_{\ell_2} & & \\ & & -\tilde{N}_{33}^{(2)-1} \tilde{N}_{31}^{(2)} & -\tilde{N}_{33}^{(2)-1} \tilde{N}_{32}^{(2)} & \tilde{N}_{33}^{(2)-1} \\ & & & & I_{n-\ell_1-\ell_2-\ell_3} \end{bmatrix},$$

all other entries zero.

(4-84)

Define

$$\tilde{N}^{(2)-1} \tilde{A}^{(1)} \tilde{N}^{(2)} \triangleq \tilde{A}^{(2)} \triangleq \begin{bmatrix} \tilde{A}_{(1,1)}^{(2)} & \tilde{A}_{(1,2)}^{(2)} & & & \\ \tilde{A}_{(2,1)}^{(2)} & \tilde{A}_{(2,2)}^{(2)} & \tilde{A}_{(2,3)}^{(2)} & & \\ \cdot & & \cdot & \cdot & \\ \cdot & & & \cdot & \tilde{A}_{(v-1,v)}^{(2)} \\ \cdot & & & & \\ \tilde{A}_{(v,1)}^{(2)} & \cdot & \cdot & \cdot & \tilde{A}_{(v,v)}^{(2)} \end{bmatrix},$$

all other entries zero,

(4-85)

where

(i) each $\tilde{A}_{(i,j)}^{(2)}$ is an $\ell_i \times \ell_j$ matrix,

(ii)

$$\tilde{A}_{(1,1)}^{(2)} = \tilde{A}_{(1,1)}^{(1)}$$

$$\tilde{A}_{(1,2)}^{(2)} = \tilde{A}_{(1,2)}^{(1)}$$

(4-86)

$$\tilde{A}_{(i,i+1)}^{(2)} = \tilde{A}_{(i,i+1)}^{(1)} = \begin{bmatrix} [0] & \hat{N}_{v+1-i, v+1-i}^{(v+1-i)^{-1}} \end{bmatrix}$$

for $i = 4, 5, \dots, v-1$, if $v \geq 5$

by (4-79~81),

(iii)

$$\tilde{A}_{(2,1)}^{(2)} = \tilde{A}_{(2,1)}^{(1)} + \tilde{A}_{(2,3)}^{(1)} \tilde{N}_{31}^{(2)} = [0]$$

$$\tilde{A}_{(2,2)}^{(2)} = \tilde{A}_{(2,2)}^{(1)} + \tilde{A}_{(2,3)}^{(1)} \tilde{N}_{32}^{(2)} = [0] \quad (4-87)$$

$$\tilde{A}_{(2,3)}^{(2)} = \tilde{A}_{(2,3)}^{(1)} \tilde{N}_{33}^{(2)} = \begin{bmatrix} [0] & I_{\ell_2} \end{bmatrix}$$

from (4-83),

(iv) if $v \geq 4$

$$\tilde{A}_{(3,4)}^{(2)} = \tilde{N}_{33}^{(2)^{-1}} \tilde{A}_{(3,4)}^{(1)} = \begin{bmatrix} [0] & \tilde{N}_{33}^{(2)^{-1}} \hat{N}_{v-2, v-2}^{(v-2)^{-1}} \end{bmatrix},$$

(4-88)

where $\hat{N}_{v-2, v-2}^{(v-2)}$ is given by (4-69). The theorem is proved

for the case of $v = 3$ by defining $\tilde{N} = \tilde{N}^{(1)} \tilde{N}^{(2)}$.

(4) For $v \geq 4$, this process can be repeated $(v - 1)$ times. Generally, defining

$$\tilde{N}^{(i+1)-1} \tilde{A}^{(i)} \tilde{N}^{(i+1)} \triangleq \tilde{A}^{(i+1)} \triangleq \begin{bmatrix} \tilde{A}_{(1,1)}^{(i+1)} & \cdot & \cdot & \cdot & \tilde{A}_{(1,i+1)}^{(i+1)} \\ \cdot & & & & \cdot \\ \cdot & & & & \cdot \\ \cdot & & & & \cdot \\ \tilde{A}_{(i+1,1)}^{(i+1)} & \cdot & \cdot & \cdot & \tilde{A}_{(i+1,i+1)}^{(i+1)} \end{bmatrix},$$

(4-89)

where

(i) $\tilde{A}^{(i)}$ has a structure after $(i-1)$ reductions, such that

$$\tilde{A}^{(i)} = \begin{bmatrix} \tilde{A}_{(1,1)}^{(1)} & \tilde{A}_{(1,2)}^{(1)} & [0] & \cdot & \cdot & \cdot & \cdot & \cdot & [0] \\ [0] & [0] & \begin{bmatrix} [0] & I_{\ell_2} \end{bmatrix} & \cdot & & & & & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ [0] & \cdot & \cdot & [0] & \begin{bmatrix} [0] & I_{\ell_i} \end{bmatrix} & \cdot & \cdot & \cdot & \cdot \\ \tilde{A}_{(i+1,1)}^{(i)} & \cdot & \cdot & \cdot & \cdot & \cdot & \tilde{A}_{(i+1,i+2)}^{(i)} & \cdot & \cdot \\ \cdot & & & & & & \cdot & \cdot & [0] \\ \cdot & & & & & & \cdot & \cdot & \tilde{A}_{(v-1,v)}^{(i)} \\ \cdot & & & & & & \cdot & \cdot & \tilde{A}_{(v,v)}^{(i)} \\ \tilde{A}_{(v,1)}^{(i)} & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{bmatrix}$$

(4-90)

where

$$\tilde{A}_{(i+1,i+2)}^{(i)} = \begin{bmatrix} [0] & \tilde{N}_{i+1 \ i+1}^{(i)-1} \hat{N}_{v-i \ v-i}^{(v-i)-1} \end{bmatrix},$$

$$(\tilde{N}_{i+1 \ i+1}^{(i)-1} \hat{N}_{v-i \ v-i}^{(v-i)-1} \text{ is } \ell_{i+1} \times \ell_{i+1} \text{ nonsingular}),$$

(4-91)

$$\tilde{A}_{(i,i+1)}^{(i)} = \begin{bmatrix} [0] & \hat{N}_{v-j \ v-j}^{(v-j)-1} \end{bmatrix} \text{ for } j = i+1, i+2, \dots, v-1,$$

(ii) $\tilde{N}^{(i+1)}$ is such that

$$\tilde{N}^{(i+1)} = \begin{bmatrix} I_{\ell_1} & & & & & & \\ & I_{\ell_2} & & & & & \\ & & \ddots & & & & \\ & & & I_{\ell_i} & & & \\ & \tilde{N}_{i+2 \ 1}^{(i+1)} & & & \tilde{N}_{i+2 \ i+2}^{(i+1)} & & \\ & & & & & I_{\ell_{i+2}} & \\ & & & & & & I_{\ell_v} \end{bmatrix},$$

all other entries zero,

(4-92)

where each $\tilde{N}_{i+2 \ j}^{(i+1)}$ satisfies

$$(a) \quad \tilde{A}_{(i+1,j)}^{(i)} + \tilde{A}_{(i+1,i+2)}^{(i)} \tilde{N}_{i+2 \ j}^{(i+1)} = [0] \quad \text{for } j = 1, 2, \dots, i+1, \quad (4-93)$$

where existence is by Lemmas 4-2, and by using (4-91),

$$(b) \quad \tilde{N}_{i+2 \ i+2}^{(i+1)} = \begin{bmatrix} I_{l_{i+2}-l_{i+1}} & [0] \\ [0] & \hat{N}_{v-i \ v-i}^{(v-i)} \tilde{N}_{i+1 \ i+1}^{(i)} \end{bmatrix}, \quad (4-94)$$

Since (4-94) is nonsingular, $\tilde{N}^{(i+1)}$ is nonsingular. For the last transformation of the sequence, $i = v - 1$ and defining

$$\tilde{N} \triangleq \tilde{N}^{(1)} \tilde{N}^{(2)} \dots \tilde{N}^{(v-1)}, \quad (4-95)$$

which is nonsingular, Theorem 4-1 is proved.

From the structure of each $\tilde{N}^{(i)}$, ($i = 1, 2, \dots, v-1$), (4-92), it follows that

$$\tilde{N} = \begin{bmatrix} \tilde{N}_{11} & [0] \\ \tilde{N}_{21} & \tilde{N}_{22} \end{bmatrix}, \quad (4-96)$$

where \tilde{N}_{11} and \tilde{N}_{22} are $(n-s) \times (n-s)$ and $s \times s$ respectively. From (4-42) and (4-89),

$$N = \hat{N} \tilde{N} = \begin{bmatrix} N_{11} & [0] \\ N_{21} & N_{22} \end{bmatrix}, \quad (4-97)$$

where N_{11} and N_{22} are $(n-s) \times (n-s)$ and $s \times s$ respectively and N_{22} is nonsingular because both \hat{N}_{22} and \tilde{N}_{22} are nonsingular.

Theorem 4-2:

Let \hat{A} and \hat{B} be $n \times n$ and $n \times m$ matrices with $\text{rank } \hat{B} = r$ and $0 < r \leq m \leq n$. There exist nonsingular matrices N and M which are $n \times n$ and $m \times m$ respectively such that

$$N^{-1} \hat{A} N = A \quad (4-98)$$

and

$$N^{-1} \hat{B} M = \begin{bmatrix} [0] & [0] \\ [0] & I_r \end{bmatrix} \triangleq B, \quad (4-99)$$

where A is given by Theorem 4-1, (4-9), with $s = r$.

Proof: First consider nonsingular matrices $N^{(1)}$ and $M^{(1)}$, $n \times n$ and $m \times m$, such that

$$N^{(1)-1} \hat{B} M^{(1)} = \begin{bmatrix} [0] & [0] \\ [0] & I_r \end{bmatrix}, \quad (4-100)$$

by Lemma 4-1. Define

$$A \triangleq N^{(1)-1} \hat{A} N^{(1)}. \quad (4-101)$$

Then from Theorem 4-1, there exists a nonsingular matrix N such that

$$A = N^{-1} \hat{A} N, \quad (4-102)$$

where the structure of N is given in (4-97). Also define

$$M^{(2)} \triangleq \begin{bmatrix} I_{m-r} & [0] \\ [0] & N_{22} \end{bmatrix}, \quad (4-103)$$

where N_{22} is defined in (4-88). If

$$\begin{aligned} N &= N^{(1)} N \\ M &= M^{(1)} M^{(2)}, \end{aligned} \quad (4-104)$$

then (4-98) and (4-99) follow.

4.3 Development of a New Canonical Form for Linear Systems

4.3.1 Canonical Form for Linear Systems

For the class of linear systems given by

$$\dot{\underline{Y}} = \hat{A} \underline{Y} + \hat{B} \underline{V}, \quad (4-105)$$

such that

- (a) \hat{A} and \hat{B} are $n \times n$ and $n \times m$ respectively with $0 < m \leq n$,
- (b) $\text{rank } \hat{B} = r$ with $0 < r \leq m$, a new canonical form is suggested by Theorems 4-1 and 4-2.

Precisely, the new canonical form of (4-105) is given by

$$\dot{\underline{X}} = \underline{A} \underline{X} + \underline{B} \underline{U}$$

$$= \begin{bmatrix} A_{(1,1)} & A_{(1,2)} & & & & \\ & & A_{(2,3)} & & & \\ & & & \ddots & & \\ & & & & A_{(v-1,v)} & \\ A_{(1,v)} & \cdot & \cdot & \cdot & \cdot & A_{(v,v)} \end{bmatrix} \underline{X} + \begin{bmatrix} \\ \\ \\ I_r \end{bmatrix} \underline{U},$$

all other entries zero, (4-106)

where

$$\begin{aligned} \underline{Y} &= \underline{N} \underline{X} & \underline{A} &= \underline{N}^{-1} \hat{\underline{A}} \underline{N} \\ \underline{V} &= \underline{M} \underline{U}, \text{ or equivalently} & \underline{B} &= \underline{N}^{-1} \hat{\underline{B}} \underline{M}. \end{aligned}$$

(4-107)

The existence of \underline{N} and \underline{M} are guaranteed by Theorem 4-2.

4.3.2 Controllability of the System Determined from the Canonical Form

From Theorem 2-9, the controllability of (4-105) can be determined by the rank of

$$\hat{\underline{G}} = [\hat{\underline{B}}, \hat{\underline{A}} \hat{\underline{B}}, \dots, \hat{\underline{A}}^{n-1} \hat{\underline{B}}]. \quad (4-108)$$

Using the nonsingular M in (4-107), define an $mn \times mn$ matrix such that

$$\hat{M} = \begin{bmatrix} M & & & \\ & M & & \\ & & \ddots & \\ & & & M \end{bmatrix}, \text{ all other entries zero,} \quad (4-109)$$

n stages

which is also nonsingular. Define

$$\begin{aligned} G &\triangleq N^{-1} \hat{G} \hat{M} = [N^{-1} \hat{B} M, N^{-1} \hat{A} \hat{B} M, \dots, N^{-1} \hat{A}^{n-1} \hat{B} M] \\ &= [B, A B, \dots, A^{n-1} B] \end{aligned} \quad (4-110)$$

from (4-107), by applying the technique used in (3-12). The rank of \hat{G} is equal to that of G because N^{-1} and M are nonsingular. Thus the following theorem is proved.

Theorem 4-3:

The given system (4-105) is completely controllable if and only if its canonical form (4-106) is completely controllable.

The alternate definition of controllability can also be considered, using Theorem 3-1. If $v = 1$, for (4-106), $B = I_r$ and the system is completely controllable by this

theorem. If $v \geq 2$, the controllability of the system depends upon the structures of the $A_{(1,1)}$ and $A_{(1,2)}$ submatrices in the canonical form given in Theorem 4-3.

Corollary 4-1:

Provided that the canonical form has $v \geq 2$, the canonical form is controllable if and only if

$$\begin{aligned} A_{(1,1)} &= [0] \\ A_{(1,2)} &= \begin{bmatrix} [0] & I_{\ell_1} \end{bmatrix} . \end{aligned} \quad (4-111)$$

Equivalently the canonical form is uncontrollable if and only if

$$A_{(1,2)} = [0] . \quad (4-112)$$

The proof follows directly from Theorems 3-1 and 4-3 and the structure of the canonical form.

4.3.3 Heuristic Explanation of the Canonical Form

Conveniently decompose \underline{U} of (4-106) to

$$\begin{aligned} \underline{U}_d &= [u_1, u_2, \dots, u_{m-r}]^T \\ \underline{U}_e &= [u_{m-r+1}, u_{m-r+2}, \dots, u_m]^T \end{aligned} \quad (4-113)$$

and define an $n \times r$ matrix

$$B_e \triangleq \begin{bmatrix} [0] \\ I_r \end{bmatrix}; \quad (4-114)$$

then the canonical form (4-106) is reduced to

$$\dot{\underline{X}} = A \underline{X} + B_e \underline{U}_e. \quad (4-115)$$

Thus only r control variables out of m are effective in multi-input systems, where r is the rank of \hat{B} .

Conveniently call each ℓ_i of Theorem 4-1 the i^{th} stage number of the canonical form and decompose the state vector

$$\underline{X} = \begin{bmatrix} \underline{X}_{(1)} \\ \underline{X}_{(2)} \\ \vdots \\ \underline{X}_{(v)} \end{bmatrix}, \quad (4-116)$$

where each $\underline{X}_i = [x_{p_i+1}, x_{p_i+2}, \dots, x_{p_i+\ell_i}]^T$ with $p_i = \sum_{k=1}^{i-1} \ell_k$.

Then the canonical form (4-106) becomes

(i) if the original system is controllable,

$$\begin{aligned}
\dot{\underline{x}}_{(1)} &= A_{(1,2)} \underline{x}_{(2)} \\
\dot{\underline{x}}_{(2)} &= A_{(2,3)} \underline{x}_{(3)} \\
&\vdots \\
\dot{\underline{x}}_{(v-1)} &= A_{(v-1,v)} \underline{x}_{(v)} \\
\dot{\underline{x}}_{(v)} &= A_{(v,1)} \underline{x}_{(1)} + \dots + A_{(v,v)} \underline{x}_{(v)} + \underline{u}_e, \text{ or}
\end{aligned}
\tag{4-117}$$

(ii) if it is uncontrollable,

$$\begin{aligned}
\dot{\underline{x}}_{(1)} &= A_{(1,1)} \underline{x}_{(1)} \\
\dot{\underline{x}}_{(2)} &= A_{(2,3)} \underline{x}_{(3)} \\
&\vdots \\
\dot{\underline{x}}_{(v-1)} &= A_{(v-1,v)} \underline{x}_{(v)} \\
\dot{\underline{x}}_{(v)} &= A_{(v,1)} \underline{x}_{(1)} + \dots + A_{(v,v)} \underline{x}_{(v)} + \underline{u}_e.
\end{aligned}
\tag{4-118}$$

Provided $v \geq 2$, these equations appear schematically as shown in Figures 4-1 and 4-2.

In (4-116), the stable variables in $\underline{x}_{(1)}$ are uncontrollable because they behave as $e^{A_{(1,1)}(t-t_0)} \underline{x}_{(1)0}$ from an initial condition $(\underline{x}_{(1)0}, t_0)$. However, the other state variables, $\underline{x}_{(2)}, \underline{x}_{(3)}, \dots, \underline{x}_{(v)}$, can be controlled everywhere by \underline{u}_e . To show this, consider an initial states \underline{x}_0 , and define a control function

$$\underline{u}_e(t) = \underline{u}(t) - A_{(v,1)} e^{A_{(1,1)}(t-t_0)} \underline{x}_{(1)0} \tag{4-119}$$

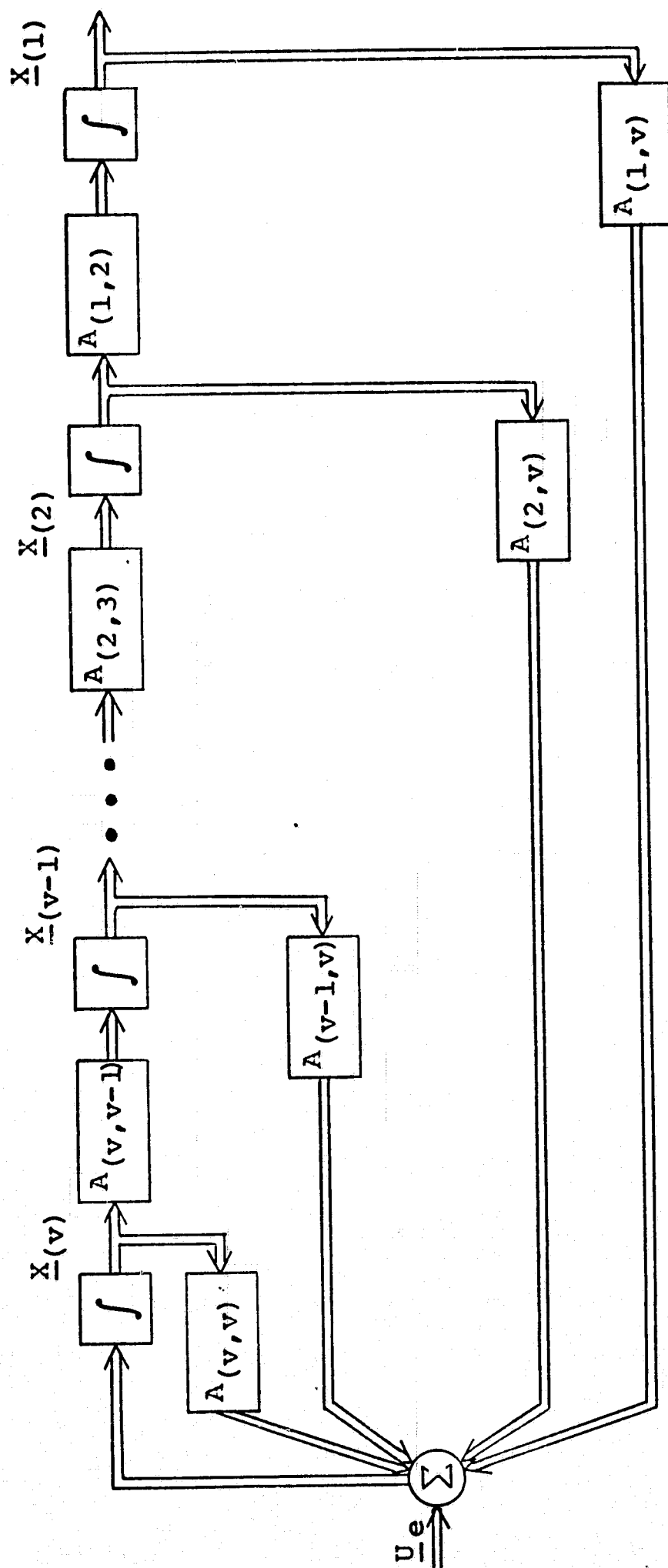


Figure 4-1. Canonical form of a completely controllable linear system (4-117).

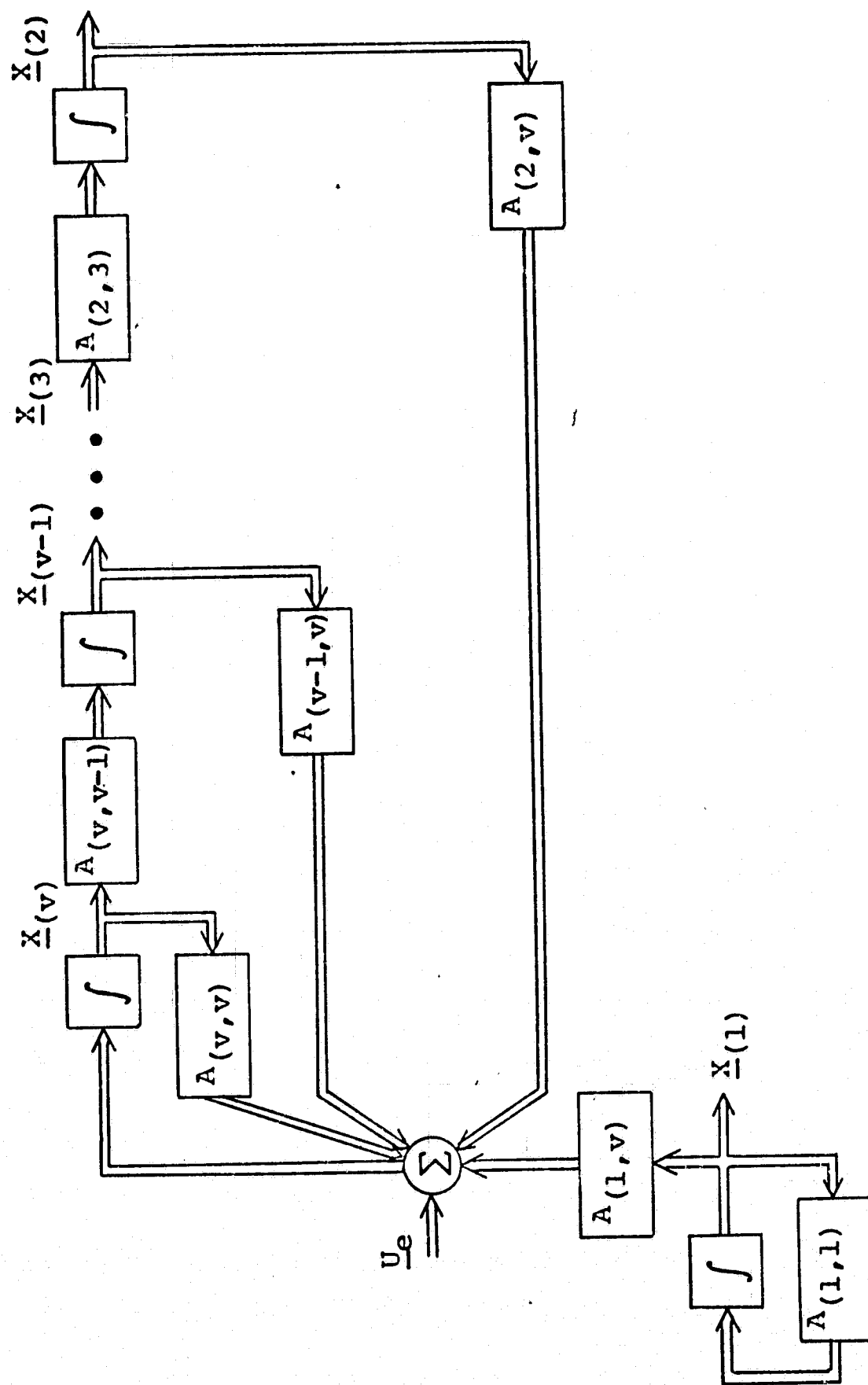


Figure 4-2. Canonical form of an uncontrollable linear system (4-118).

Then, from (4-116), it follows that

$$\begin{aligned}
 \dot{\underline{x}}_{(2)} &= A_{(2,3)} \underline{x}_{(3)} \\
 \dot{\underline{x}}_{(3)} &= A_{(3,4)} \underline{x}_{(4)} \\
 &\vdots \\
 \dot{\underline{x}}_{(v-1)} &= A_{(v-1,v)} \underline{x}_{(v)} \\
 \dot{\underline{x}}_{(v)} &= A_{(v,2)} \underline{x}_{(2)} + A_{(v,3)} \underline{x}_{(3)} + \dots + A_{(v,v)} \underline{x}_{(v)} + \underline{u}(t).
 \end{aligned}
 \tag{4-120}$$

The state variables in $\underline{x}_{(2)}, \underline{x}_{(3)}, \dots, \underline{x}_{(v)}$ are completely controllable, by (4-117). Conveniently call the state variables in $\underline{x}_{(1)}$ of (4-118) the uncontrollable state variables.

For subsystems of (4-117) or (4-118) given by

$$\dot{\underline{x}}_{(i)} = A_{(i,i+1)} \underline{x}_{(i+1)}, \tag{4-121}$$

with $A_{(i,i+1)} = \begin{bmatrix} [0] & I_{\ell_i} \end{bmatrix}$, a more detailed representation

can be made than is shown in Figure 4-4 and Figure 4-2, i.e., Figure 4-3.

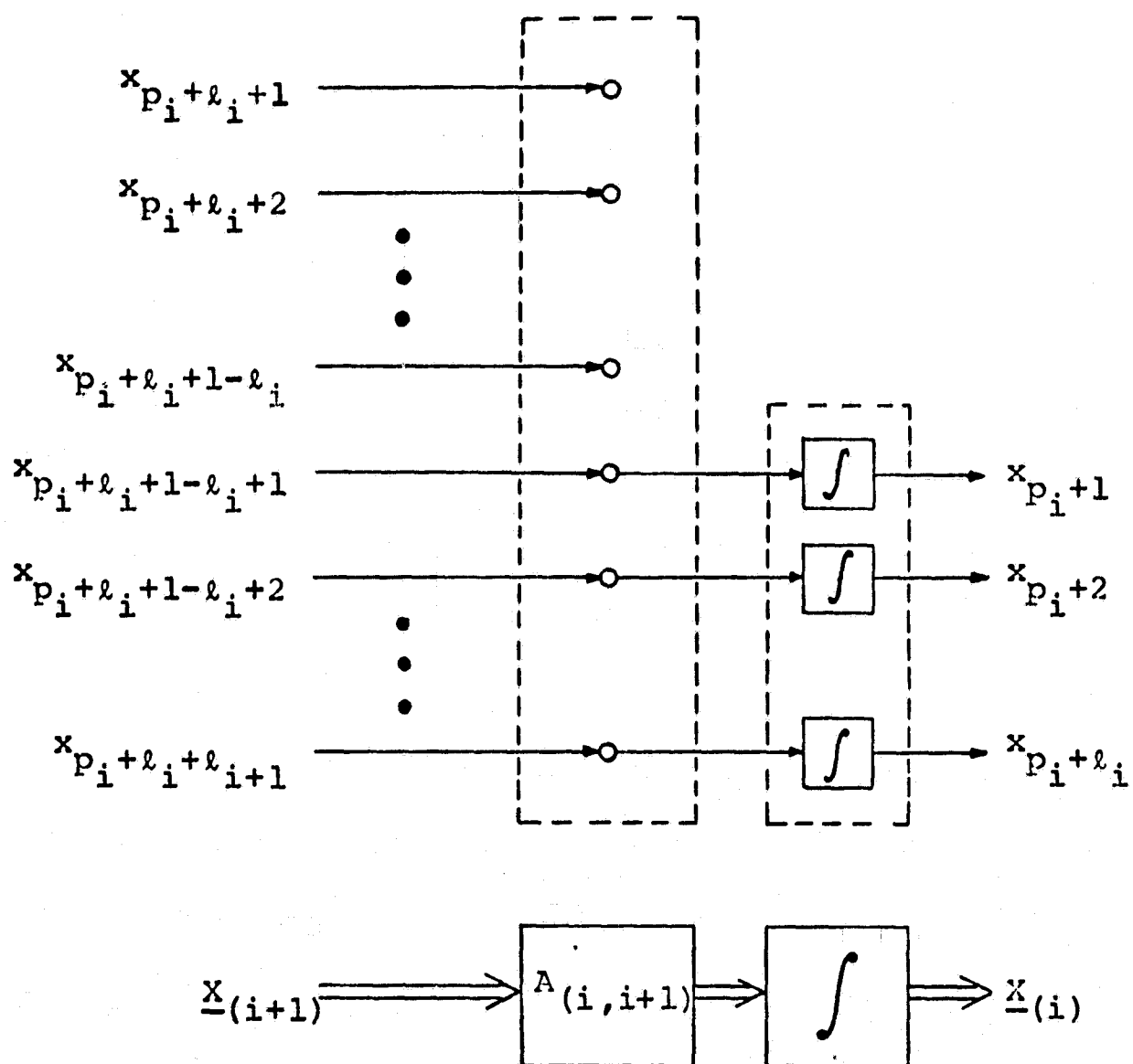


Figure 4-3. Subsystem of the canonical form in (4-117) and (4-118).

That is, the last ℓ_i state variables of $\underline{X}_{(i+1)}$ are integrated to become $\underline{X}_{(i)}$. Let the ordered set $\{\ell_1, \ell_2, \dots, \ell_v\}$ be the state distribution of the canonical form.

Physically, this refers to the numbers of integrators at each stage of the partitioned $\underline{X}_{(i)}$, in sequence.

4.4 Uniqueness of the Structure of the Canonical Form

A given system does not have a unique canonical form, (4-106), but one that is dependent on the choice of different combinations of the M and N matrices in (4-107). The stage distributions for these different forms, however, are unique.

Theorem 4-4:

The stage distribution $\{\ell_1, \ell_2, \dots, \ell_v\}$ for all canonical forms of a given system is uniquely determined.

Proof: If $\ell_1 = n = r$, $v = 1$ is determined from the construction and the stage distribution is unique. Consider $v \geq 2$, i.e., $0 < r < n$ and two transformations

$$\underline{X} = N^{-1} \underline{Y}$$

$$\underline{U} = M^{-1} \underline{V}$$

(4-122)

and

$$\underline{x} = n^{-1} \underline{y}$$

(4-123)

$$\underline{u} = m^{-1} \underline{v}$$

such that they provide the canonical form

$$\dot{\underline{x}} = A \underline{x} + B \underline{u}$$

(4-124)

with a stage distribution

$$\{\ell_1, \ell_2, \dots, \ell_v\}$$

(4-125)

and

$$\dot{\underline{x}} = A \underline{x} + B \underline{u}$$

(4-126)

with a stage distribution

$$\{\ell'_1, \ell'_2, \dots, \ell'_\mu\},$$

(4-127)

where $\ell_v = \ell'_\mu = r$. Define

$$\begin{aligned}\Theta &= \mathcal{N}^{-1} N \\ \pi &= \mathcal{M}^{-1} M.\end{aligned}\tag{4-128}$$

Then from the transitive law of the equivalent relation for the canonical form, (4-126) must be transformed into (4-124) by

$$\begin{aligned}\underline{\chi} &= \Theta \underline{x} \\ \underline{u} &= \pi \underline{u},\end{aligned}\tag{4-129}$$

which is equivalent to

$$\begin{aligned}A &= \Theta^{-1} A \Theta \\ B &= \Theta^{-1} B \pi.\end{aligned}\tag{4-130}$$

Assume the contrary of the theorem statement, that is, in (4-125) and (4-127), there exists a positive integer $\rho \leq \text{Min}(v, \mu)$ such that

$$\begin{aligned}l_{v-i} &= l'_{\mu-i}, \quad \text{for } i = 0, 1, \dots, \rho, \\ l_{v-\rho-1} &\neq l'_{\mu-\rho-1}.\end{aligned}\tag{4-131}$$

Identify

$$A = \begin{bmatrix} A_{(0,0)} & A_{(0,\alpha)} & & & & \\ & & A_{(\alpha,\alpha+1)} & & & \\ & & & \cdot & & \\ & & & & \cdot & \\ & & & & & \cdot \\ & & & & & A_{(\mu-1,\mu)} \\ A_{(\mu,0)} & \cdot & \cdot & \cdot & \cdot & A_{(\mu,\mu)} \end{bmatrix},$$

all other entries zero, (4-132)

and

$$A = \begin{bmatrix} A_{(0,0)} & A_{(0,\alpha)} & & & & \\ & & A_{(\beta,\beta+1)} & & & \\ & & & \cdot & & \\ & & & & \cdot & \\ & & & & & \cdot \\ & & & & & A_{(v-1,v)} \\ A_{(v,0)} & \cdot & \cdot & \cdot & \cdot & A_{(v,v)} \end{bmatrix},$$

all other entries zero, (4-133)

where

$$\alpha = \mu - \rho + 1 \quad (4-134)$$

$$\beta = \nu - \rho + 1, \quad (4-135)$$

each $A_{(i,j)}$, $(i,j = 0, \alpha, \alpha+1, \dots, \mu)$, is $\gamma_i \times \gamma_j$ such that

$$\gamma_0 = n - \sum_{i=\alpha}^{\mu} \ell'_i = n - \sum_{i=\beta}^{\mu} \ell_i \quad (4-136)$$

$$\gamma_i = \ell'_i = \ell_{\nu-\mu+i}, \quad \mu-\rho+1 \leq i \leq \mu,$$

and A is decomposed in the same partitioning as A .

From (4-121), it must follow that

$$A_{(\mu-i, \mu-i+1)} = A_{(\nu-i, \nu-i+1)} = \begin{bmatrix} [0] & I_{\ell_{\nu-i}} \end{bmatrix} \quad (4-137)$$

for $i = 1, 2, \dots, \rho-1,$

but the rank of $A_{(0,\alpha)}$ differs from the rank of $A_{(0,\alpha)}$.

Let

$$\Theta \triangleq \begin{bmatrix} \Theta_{(0,0)} & \Theta_{(0,\alpha)} & \cdot & \cdot & \cdot & \Theta_{(0,\mu)} \\ \cdot & & & & & \cdot \\ \cdot & & & & & \cdot \\ \cdot & & & & & \cdot \\ \cdot & & & & & \cdot \\ \Theta_{(\mu,0)} & \Theta_{(\mu,\alpha)} & \cdot & \cdot & \cdot & \Theta_{(\mu,\mu)} \end{bmatrix} \quad (4-138)$$

and

$$\Theta^{-1} \triangleq \begin{bmatrix} \tilde{\Theta}_{(0,0)} & \tilde{\Theta}_{(0,\alpha)} & \cdot & \cdot & \cdot & \tilde{\Theta}_{(0,\mu)} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \tilde{\Theta}_{(\mu,0)} & \tilde{\Theta}_{(\mu,\alpha)} & \cdot & \cdot & \cdot & \tilde{\Theta}_{(\mu,\mu)} \end{bmatrix}, \quad (4-139)$$

where the decomposition corresponds to that of A (4-132).

From (4-130),

$$B = \begin{bmatrix} \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ I_r \end{bmatrix} = \Theta^{-1} B \mathcal{P} = \begin{bmatrix} [0] \tilde{\Theta}_{(0,\mu)} \mathcal{P} \\ [0] \tilde{\Theta}_{(\alpha,\mu)} \mathcal{P} \\ \cdot \\ \cdot \\ [0] \tilde{\Theta}_{(\mu,\mu)} \mathcal{P} \end{bmatrix}. \quad (4-140)$$

Comparing the submatrices in (4-140), it follows that

$$[0] = \begin{bmatrix} [0] \tilde{\Theta}_{(i,\mu)} \mathcal{P} \end{bmatrix}, \quad \text{for } i = 0, \alpha, \alpha+1, \dots, \mu-1, \quad (4-141)$$

$$\begin{bmatrix} [0] & I_r \end{bmatrix} = \begin{bmatrix} [0] \tilde{\Theta}_{(\mu,\mu)} \mathcal{P} \end{bmatrix}.$$

As \mathcal{P} is nonsingular,

$$[0] = \tilde{\Theta}_{(i,\mu)} , \text{ for } i = 0, \alpha, \alpha+1, \dots, \mu-1 , \quad (4-142)$$

and $\tilde{\Theta}_{(\mu,\mu)}$ must be nonsingular. Subsequently, from the inverse relation between (4-138) and (4-139),

$$\begin{aligned} \Theta_{(i,\mu)} &= [0] , \text{ for } i = 0, \alpha, \alpha+1, \dots, \mu-1, \\ \Theta_{(\mu,\mu)}^{-1} &= \tilde{\Theta}_{(\mu,\mu)} . \end{aligned} \quad (4-143)$$

From the first equation of (4-130), it follows that

$$A = \Theta^{-1} A \Theta = \begin{bmatrix} \tilde{\Theta}_{(0,0)} & \cdot & \cdot & \cdot & \tilde{\Theta}_{(0,\mu-1)} & [0] \\ \cdot & & & & \cdot & \cdot \\ \cdot & & & & \cdot & \cdot \\ \cdot & & & & \cdot & [0] \\ \tilde{\Theta}_{(\mu,0)} & \cdot & \cdot & \cdot & \tilde{\Theta}_{(\mu,\mu-1)} & \tilde{\Theta}_{(\mu,\mu)} \end{bmatrix}$$

$$\begin{bmatrix} A_{(0,0)} & A_{(0,\mu)} & [0] & \cdot & \cdot & \cdot & [0] \\ [0] & [0] & A_{(\mu,\mu+1)} & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & [0] \\ [0] & \cdot & \cdot & \cdot & \cdot & [0] & A_{(\mu-1,\mu)} \\ A_{(\mu,0)} & \cdot & \cdot & \cdot & \cdot & \cdot & A_{(\mu,\mu)} \end{bmatrix}$$

$$\begin{bmatrix} \Theta_{(0,0)} & \cdot & \cdot & \cdot & \Theta_{(0,\mu-1)} & [0] \\ \cdot & & & & \cdot & \cdot \\ \cdot & & & & \cdot & \cdot \\ \cdot & & & & \cdot & [0] \\ \Theta_{(\mu,0)} & \cdot & \cdot & \cdot & \Theta_{(\mu,\mu-1)} & \Theta_{(\mu,\mu)} \end{bmatrix} \cdot$$

(4-144)

Comparing the submatrices in the second last columns of (4-144) and (4-134),

$$[0] = \tilde{\Theta}_{(i, \mu-1)} A_{(\mu-1, \mu)} \Theta_{(\mu, \mu)} = \begin{bmatrix} [0] & \tilde{\Theta}_{(i, \mu-1)} \end{bmatrix} \Theta_{(\mu, \mu)},$$

$$\text{for } i = 0, \alpha, \dots, \nu-2 \quad (4-145)$$

$$\begin{aligned} A_{(\nu-1, \nu)} &= \begin{bmatrix} [0] & I_{\ell_{\nu-1}} \end{bmatrix} = \tilde{\Theta}_{(\mu-1, \mu-1)} A_{(\mu-1, \mu)} \Theta_{(\mu, \mu)} \\ &= \begin{bmatrix} [0] & \tilde{\Theta}_{(\mu-1, \mu-1)} \end{bmatrix} \Theta_{(\mu, \mu)}. \end{aligned}$$

As $\Theta_{(\mu, \mu)}$ is nonsingular, it must be concluded that

$$[0] = \tilde{\Theta}_{(i, \mu-1)}, \quad \text{for } i = 0, \alpha, \alpha+1, \dots, \nu-2, \quad (4-146)$$

and from the inverse relation between (4-138) and (4-139),

$$\begin{aligned} \Theta_{(i, \mu-1)} &= [0], \quad \text{for } i = 0, \alpha, \alpha+1, \dots, \mu-1, \\ \Theta_{(\mu-1, \mu-1)} &= \tilde{\Theta}_{(\mu-1, \mu-1)}^{-1}. \end{aligned} \quad (4-147)$$

Repeating the decomposition of Θ ρ times, it follows that for (4-126) to be satisfied,

$$\begin{aligned} \Theta_{(i, j)} &= \tilde{\Theta}_{(i, j)} = [0], \quad \text{for } i < j, \\ \Theta_{(i, i)}^{-1} &= \tilde{\Theta}_{(i, i)}, \quad \text{for } i = 0, \alpha, \alpha+1, \dots, \mu, \end{aligned} \quad (4-148)$$

where $\Theta_{(i,i)}$ is nonsingular. From (4-134) and (4-147)

$$A_{(0,\alpha)} = \tilde{\Theta}_{(0,0)} A_{(0,\alpha)} \Theta_{(\alpha,\alpha)} \quad (4-149)$$

and the rank of $A_{(0,\alpha)}$ must be equal to that of $A_{(0,\alpha)}$,
i.e.,

$$l_{\mu-\rho-1} = l_{v-\rho-1} \quad (4-150)$$

as both $\tilde{\Theta}_{(0,0)}$ and $\tilde{\Theta}_{(\alpha,\alpha)}$ are nonsingular. But this contradicts the hypothesis (4-131).

Consequently, the canonical forms (4-124) and (4-126) differ only for elements in submatrices of A and A with indices $(v,1), (v,2), \dots, (v,v)$, and also $(1,1)$ if the system is uncontrollable, i.e.,

$$A = \begin{bmatrix} A_{(1,1)} & A_{(1,2)} & & & & \\ & & A_{(2,3)} & & & \\ & & & \ddots & & \\ & & & & \ddots & \\ & & & & & A_{(v-1,v)} \\ A_{(v,1)} & \cdot & \cdot & \cdot & \cdot & A_{(v,v)} \end{bmatrix},$$

all other entries zero.

(4-151)

The uniqueness of the structure of canonical forms guaranteed for each given system is due to the uniqueness of the stage distribution. The ambiguity in sizes of the submatrices of A and B is thereby avoided in contrast to other suggested canonical forms.

From Theorem 4-4 and (4-146), a corollary follows directly from (4-148).

Corollary 4-2:

The matrix \oplus in (4-138) must be such that

$$\oplus = \begin{bmatrix} \oplus_{(1,1)} & & & & \\ \cdot & \oplus_{(2,2)} & & & \\ \cdot & & \cdot & & \\ \cdot & & & \cdot & \\ \oplus_{(v,1)} & \cdot & \cdot & \cdot & \oplus_{(v,v)} \end{bmatrix},$$

all other entries zero.

(4-152)

where each $\oplus_{(i,j)}$, $(i,j = 1,2,\dots,v)$, is $\ell_i \times \ell_j$ and

$\oplus_{(i,i)}$ is nonsingular.

4.5 Uniqueness of the Canonical Form

Furthermore, consider a set of canonical forms for a given system with a fixed M (4-107), and observe the variety of matrices A corresponding to various N . This can be done effectively by observing the differences in A and \hat{A} of (4-124) and (4-126) under the restriction $M = \mathcal{M}$ in (4-122) and (4-123), or equivalently $\mathcal{L} = I_m$.

Theorem 4-5:

Consider a completely controllable system (4-105) with the stage distribution

$$\{r, r, \dots, r\}, \quad (4-153)$$

where r is the rank of \hat{B} and assume that a transformation (4-107) is made. Then N is unique.

Proof: Assume two transformations, (4-124) and (4-126), with $\mathcal{M} = M$. Then the proof is based on (4-130), showing that Θ in (4-138) is the unit matrix.

For the case of $v = 1$, since $n = 4$, N^{-1} must be $(\hat{B} M)^{-1}$ to satisfy (4-107) with $B = I_r = I_n$, and the uniqueness follows. For $v \geq 2$, consider A and \hat{A} as given by (4-124) and (4-126). Then from the controllability

of the system and Corollary 4-1, (4-153) and (4-151), it follows that

$$A = \begin{bmatrix} [0] & I_r & & & \\ & \cdot & \cdot & & \\ & & \cdot & \cdot & \\ & & & \cdot & \cdot \\ & & & & [0] & I_r \\ A_{(v,1)} & \cdot & \cdot & \cdot & & A_{(v,v)} \end{bmatrix}, \text{ other entries zero,} \quad (4-154)$$

and

$$A = \begin{bmatrix} [0] & I_r & & & \\ & \cdot & \cdot & & \\ & & \cdot & \cdot & \\ & & & \cdot & \cdot \\ & & & & [0] & I_r \\ A_{(v,1)} & \cdot & \cdot & \cdot & & A_{(v,v)} \end{bmatrix}, \text{ other entries zero,} \quad (4-155)$$

where each submatrix is $r \times r$. Then from (4-130) and

$$\pi = I_m,$$

$$A \oplus = \oplus A \quad (4-156)$$

and

$$B = \oplus B, \quad (4-157)$$

where Θ is given as (4-152). Then from (4-157)

$$\Theta_{(v,v)} = I_r. \quad (4-158)$$

Alternately, from (4-154~156),

$$A \Theta = \begin{bmatrix} \Theta_{(2,1)} & \Theta_{(2,2)} & [0] & \cdot & \cdot & \cdot & [0] \\ & & & \cdot & & & \cdot \\ \Theta_{(3,1)} & & \Theta_{(3,3)} & \cdot & & & \cdot \\ \cdot & & & \cdot & & \cdot & \cdot \\ \cdot & & & & \cdot & & [0] \\ \cdot & & & & & \cdot & \cdot \\ \Theta_{(v,1)} & \cdot & \cdot & \cdot & \cdot & \cdot & \Theta_{(v,v)} \\ \Delta & \Delta & \cdot & \cdot & \cdot & \cdot & \Delta \end{bmatrix}, \quad (4-159)$$

$$\Theta A = \begin{bmatrix} [0] & \Theta_{(1,1)} & [0] & \cdot & \cdot & \cdot & [0] \\ \cdot & \Theta_{(2,1)} & \Theta_{(2,2)} & \cdot & & & \cdot \\ \cdot & \cdot & & \cdot & & \cdot & \cdot \\ \cdot & \cdot & & & \cdot & & [0] \\ \cdot & \cdot & & & & \cdot & \cdot \\ [0] & \Theta_{(v-1,1)} & \cdot & \cdot & \cdot & \cdot & \Theta_{(v-1,v-1)} \\ \Delta & \Delta & \Delta & \cdot & \cdot & \cdot & \Delta \end{bmatrix}, \quad (4-160)$$

where Δ refers to entries of no importance to the following.

Comparing submatrices in the first $v - 1$ rows,

$$\Theta_{(i,j)} = [0] , \quad \text{if } i > j , \quad (4-161)$$

and $\Theta_{(i,i)} = \Theta_{(i+1,i+1)} , \text{ for } i = 1, 2, \dots, (v-1) .$
(4-162)

But from (4-153), (4-162) reduces to

$$\Theta_{(i,i)} = I_r , \text{ for } i = 1, 2, \dots, v , \quad (4-163)$$

and Θ must be the unit matrix.

Theorem 4-5 establishes the uniqueness of the canonical form for a given completely controllable system, with a stage distribution given by (4-153). If it is possible to select M as a unit matrix, then this canonical form is reduced to the canonical forms discussed in Chapter 3. Sufficient conditions which allow this choice of M are given in the next theorem.

Theorem 4-6:

Consider the system (4-105) with a stage distribution

$$\{\ell_1, r, r, \dots, r\} , \quad 0 < \ell_1 \leq r , \quad (4-164)$$

$$r = \text{rank } \hat{B} = m , \quad (4-165)$$

where it is assumed that $\ell_1 = r$ if $v = 1$. Then the canonical transformation (4-107) is possible with

$$M = I_r = I_m.$$

Proof: The resulting structure of A is

$$A = \begin{bmatrix} [0] & \begin{bmatrix} [0] & I_{\ell_1} \end{bmatrix} & & & \\ & & I_r & & \\ & & & \ddots & \\ & & & & I_r \\ A_{(v,1)} & A_{(v,2)} & \cdot & \cdot & \cdot & A_{(v,v)} \end{bmatrix}, \text{ all other entries zero.} \quad (4-166)$$

If $v = 1$, then $n = r$, \hat{B} is an $n \times n$ nonsingular matrix by assumption, and the proof follows by specifying $N = \hat{B}^{-1}$. Consider the case of $v \geq 2$. Since $\ell_v = r$, $M^{(1)}$ in (4-100) can be chosen as I_r according to Lemma 4-1. For the transformation of A (4-101) into \tilde{A} (4-27), \hat{N} of (4-72) can be chosen as

$$\hat{N} = \begin{bmatrix} \hat{N}_{11} & [0] \\ \hat{N}_{21} & I_r \end{bmatrix}. \quad (4-167)$$

This follows as each $\hat{N}_{ii}^{(i)}$ of (4-58), $i = 1, 2, \dots, v$, can

be chosen as I_r according to Lemma 4-1 due to the null rank of each $A_{(i+1,i)}^{(i-1)}$, $i = 2, 3, \dots, v$, by assumption (4-164). Then from (4-63) and (4-67),

$$\tilde{A} = A^{(v-1)} = \begin{bmatrix} \tilde{A}_{(1,1)} & \tilde{A}_{(1,2)} & & & \\ \tilde{A}_{(2,1)} & & I_r & & \\ \vdots & & & \ddots & \\ \vdots & & & & I_r \\ \tilde{A}_{(v,1)} & & & & \tilde{A}_{(v,v)} \end{bmatrix},$$

all other entries zero.

(4-168)

Then from (4-92) and (4-05), \tilde{N} can be chosen as

$$\tilde{N} = \begin{bmatrix} \tilde{N}_{11} & [0] \\ \tilde{N}_{21} & I_r \end{bmatrix} \quad (4-169)$$

because each $\tilde{N}_{i+1, i+1}^{(i)}$, $i = 1, 2, \dots, v-1$, can be chosen as I_r , according to (4-94) with all $\hat{N}_{ii}^{(i)} = I_r$. Subsequently, it follows that

$$N = \hat{N} \tilde{N} = \begin{bmatrix} N_{11} & [0] \\ N_{21} & I_r \end{bmatrix} \quad (4-170)$$

and $M^{(2)}$ of (4-103) is also the unit matrix. Therefore

$$M = M^{(1)} M^{(2)} = I_r. \quad (4-171)$$

Combining Theorems 4-5 and 4-6, the following is evident.

Corollary 4-3:

Consider a completely controllable system (4-105) with the stage distribution given by (4-153) and $\text{rank } \hat{B} = r = m$. Then there exists a unique canonical form for the system with $M = I_r$.

If the stage distribution of the system is given by (4-164), the canonical form corresponds to Asseo's canonical form, A given (4-166). However, in the reduction to the Asseo's canonical form in the sense of Definition 4-1, it is necessary that:

- (i) the stage distribution be given as (4-164), and
- (ii) the canonical transformation (4-107) is possible with the unit matrix M .

Consequently, the application of Asseo's unique compact decomposition to his canonical form is limited to subclass of systems (4-105) for which this decomposition is

applicable and equally unique and compact.

The transformation in Theorem 4-6 corresponds to the canonical transformations used in Definition 3-1. Thus the stage distribution of a single input completely controllable system must be $\{1,1,\dots,1\}$, and it follows from Theorem 4-6 that the canonical form of the system uniquely exists in the sense of Definition 3-1, i.e.,

$$\dot{\underline{x}} = \begin{bmatrix} 0 & 1 & 0 & 0 & \cdot & \cdot & 0 \\ 0 & 0 & 1 & 0 & \cdot & \cdot & 0 \\ \cdot & & & \cdot & \cdot & \cdot & \cdot \\ \cdot & & & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & \cdot & \cdot & 0 & 1 & 0 \\ a_1 & \cdot & \cdot & \cdot & a_{n-1} & a_n & \cdot \end{bmatrix} \underline{x} + \begin{bmatrix} 0 \\ \cdot \\ \cdot \\ \cdot \\ 0 \\ 1 \end{bmatrix} u. \quad (4-172)$$

Consequently, the statement of the uniqueness of this canonical form in Section 3.4 is verified.

4.6 Examples

Example 4-1. Consider the system

$$\dot{y}_1 = y_1 + 3y_3 - y_4 + v_2$$

$$\dot{y}_2 = 2y_1 + y_2 - y_3 + 4y_4 + v_1$$

$$\dot{y}_3 = -y_1 + y_2 - y_3 + 2y_4 + v_1$$

$$\dot{y}_4 = y_1 + 2y_3 - y_4 + v_2$$

(4-173)

and from (4-100)

$$\hat{A} = \begin{bmatrix} 1 & 0 & 3 & -1 \\ -2 & 1 & -1 & 4 \\ -1 & 1 & -1 & 2 \\ 1 & 0 & 2 & -1 \end{bmatrix} \quad \text{and} \quad \hat{B} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (4-174)$$

For

$$N^{(1)-1} \hat{B} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (4-175)$$

it is sufficient to assume

$$N^{(1)} = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (4-176)$$

Then

$$A = \tilde{A} = N^{(1)-1} \hat{A} N^{(1)} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ -1 & 0 & 0 & 1 \\ -1 & 1 & 0 & 1 \\ 1 & 0 & 2 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} \tilde{A}_{(1,1)} & \tilde{A}_{(1,2)} \\ \tilde{A}_{(2,1)} & \tilde{A}_{(2,2)} \end{bmatrix}. \quad (4-177)$$

Comparing this to (4-78) and (4-80), $\tilde{A}_{(1,1)}$ can be made [0] if

$$N^{(1)} = N = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix}. \quad (4-178)$$

Then

$$N^{-1} N^{(1)-1} \hat{A} N^{(1)} N = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{bmatrix}, \text{ and} \quad (4-179)$$

$$N^{-1} N^{(1)-1} \hat{B} = B, \quad (4-180)$$

which is the desired canonical form. Defining

$$N \triangleq N^{(1)} N = \begin{bmatrix} 2 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix} \quad (4-181)$$

and

$$M = I_2, \quad (4-182)$$

then, by the transformation,

$$\underline{Y} = N \underline{X} \quad (4-183)$$

$$\underline{V} = M \underline{U},$$

the canonical form of the system becomes

$$\begin{aligned} \dot{x}_1 &= x_3 \\ \dot{x}_2 &= x_4 \\ \dot{x}_3 &= x_2 + x_4 + u_1 \\ \dot{x}_4 &= x_1 + x_3 + u_2. \end{aligned} \quad (4-184)$$

The system is controllable by Theorem 4-3. As $M = I_2$, (4-184) is the unique canonical form of (4-173), from Corollary 4-3.

Example 4-2. Consider

$$\begin{aligned}\dot{y}_1 &= y_1 + 3y_2 - 3y_3 - y_4 + v_2 \\ \dot{y}_2 &= -2y_1 + 4y_4 + v_1 \\ \dot{y}_3 &= -y_1 + 2y_4 + v_1 \\ \dot{y}_4 &= y_1 + 2y_2 - 2y_3 - y_4 + v_2.\end{aligned}\tag{4-185}$$

With the same transformation of Example 4-1,

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= x_4 \\ \dot{x}_3 &= x_4 + u_1 \\ \dot{x}_4 &= x_1 + x_2 + u_2,\end{aligned}\tag{4-186}$$

which is controllable. This is not a unique canonical form, however. Selecting an alternate N for (4-178) such that

$$N = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix},\tag{4-187}$$

(4-185) becomes

$$\dot{x}_1 = x_2$$

$$\dot{x}_2 = x_4$$

$$\dot{x}_3 = -x_2 + x_4 + u_1$$

$$\dot{x}_4 = x_1 + x_2 + u_2 ,$$

(4-188)

where the first two equations of (4-188) are identical to those of (4-186) as determined by the structure of the canonical form.

Example 4-3. Consider

$$\dot{y}_1 = -y_1 + y_2 + y_4 + v_2$$

$$\dot{y}_2 = -y_1 - y_2 + y_3 + 2y_4 + v_1$$

$$\dot{y}_3 = -y_2 + y_3 + v_1$$

$$\dot{y}_4 = y_2 + v_2 .$$

(4-189)

With the same transformation as Example 4-1, i.e., (4-181) and (4-182),

$$\begin{aligned}
 \dot{x}_1 &= -x_1 \\
 \dot{x}_2 &= x_4 \\
 \dot{x}_3 &= -x_2 + u_1 \\
 \dot{x}_4 &= x_1 + x_2 + x_3 + u_2 ,
 \end{aligned}
 \tag{4-190}$$

which is uncontrollable as x_1 is isolated.

The diagrams of these example canonical systems are shown in Figures 4-4~6.

4.7 Application of the Canonical Form to General Systems

Consider a class of systems given by (4-1) which can be expressed by

$$\hat{\underline{F}}(\underline{y}) \triangleq \hat{\underline{A}} \underline{y} + \hat{\underline{F}}(\underline{y}) , \tag{4-191}$$

where $\hat{\underline{A}}$ is an $n \times n$ matrix such that $\hat{\underline{A}} \underline{y}$ describes the first degree homogeneous function of \underline{y} in $\hat{\underline{F}}(\underline{y})$, with $\hat{\underline{F}}(\underline{y})$ the remainder. A canonical form for this class of systems, using the development in this chapter, is

$$\dot{\underline{x}} = \underline{A} \underline{x} + \underline{F}(\underline{x}) + \underline{B} \underline{u} \tag{4-192}$$

$$\triangleq \underline{F}(\underline{x}) + \underline{B} \underline{u} ,$$

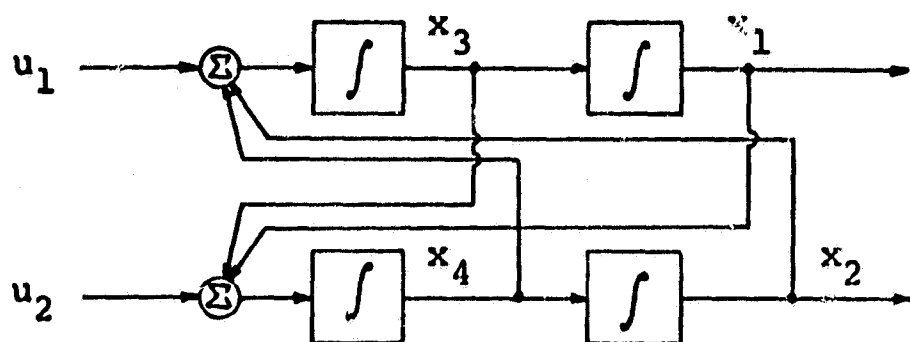


Figure 4-4. Canonical form of Example 4-1, (4-184).

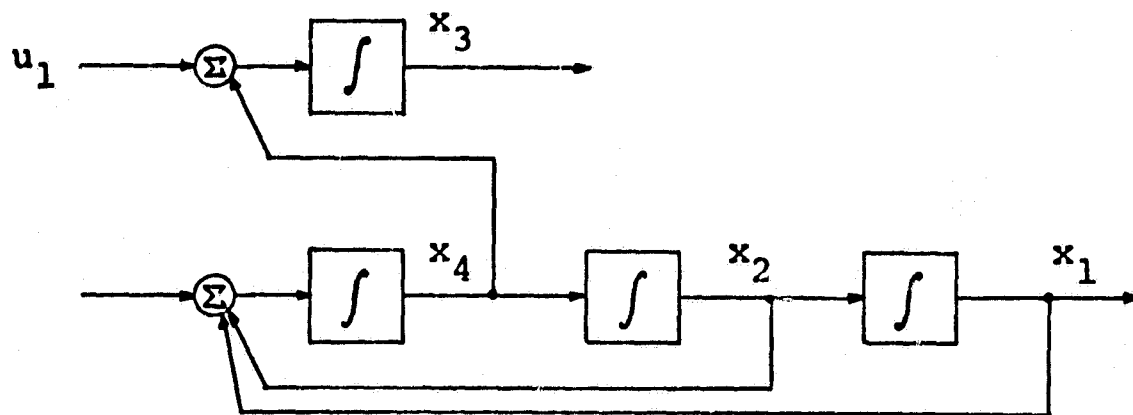


Figure 4-5. Canonical form of Example 4-2, (4-186).

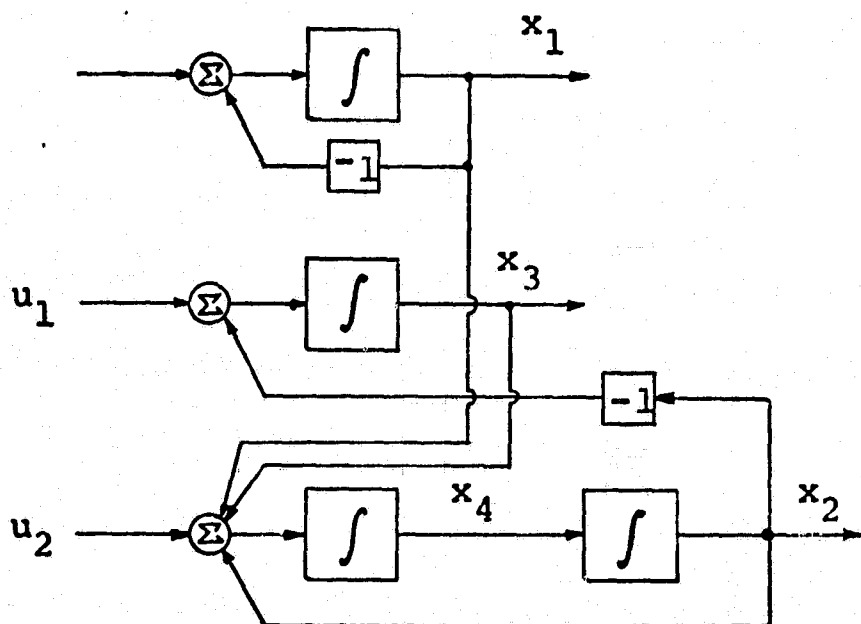


Figure 4-6. Canonical form of Example 4-3, (4-190).

where A and B are given by (4-106). For the canonical transformation (4-107),

$$A = N^{-1} \hat{A} N$$

$$B = N^{-1} \hat{B} M \quad (4-193)$$

$$\underline{F}(\underline{X}) = N^{-1} \hat{F}(N\underline{X}) .$$

The characteristics of the canonical form discussed in Sections 4.4 and 4.5 then exist for the linear part of (4-192). Furthermore, if $\hat{F}(\underline{Y})$ is of class C_2 , then $\underline{F}(\underline{X})$ in (4-193) is also of class C_2 (according to Theorem 2-1), and the uniqueness of solutions for the canonical form (4-192) is preserved according to Theorem 2-6.

Finally, a solution for the canonical form can be characterized as follows.

Lemma 4-4:

For the system (4-1), arbitrarily assume an initial condition (\underline{Y}_0, t_0) , an admissible control function $\underline{V}(t)$, and a solution $\hat{\phi}_{\underline{V}}(t; \underline{Y}_0, t_0)$. Then the solution of the canonical form (4-185) with the initial condition $(N^{-1}\underline{Y}_0, t_0)$ and the control function $\underline{U}(t) \triangleq M^{-1} \underline{V}(t)$ is

$$\underline{\phi}_U(t; N^{-1}\underline{y}_0, t_0) = N^{-1}\hat{\underline{\phi}}_V(t; \underline{y}_0, t_0) . \quad (4-194)$$

Proof: By the hypothesis, $\hat{\underline{\phi}}_V(t; \underline{y}_0, t_0)$ satisfies the characteristics of the solution given in Chapter 2.

From (2-26),

$$\underline{\phi}_U(t_0; N^{-1}\underline{y}_0, t_0) = N^{-1}\hat{\underline{\phi}}_V(t; \underline{y}_0, t_0) = N^{-1}\underline{y}_0 = \underline{x}_0 . \quad (4-195)$$

From (2-27),

$$\begin{aligned} \frac{d}{dt} \{ NN^{-1} \hat{\underline{\phi}}_V(t; \underline{y}_0, t_0) \} &= \underline{F}(NN^{-1} \hat{\underline{\phi}}_V(t; \underline{y}_0, t_0) \\ &+ B MM^{-1} \underline{V}(t) , \end{aligned} \quad (4-196)$$

which reduces to

$$\begin{aligned} \frac{d}{dt} \{ N^{-1} \hat{\underline{\phi}}_V(t; \underline{y}_0, t_0) \} &= N^{-1} \underline{F}(NN^{-1} \hat{\underline{\phi}}_V(t; \underline{y}_0, t_0) \\ &+ N^{-1} B MM^{-1} \underline{V}(t) = N^{-1} \underline{F}(N\hat{\underline{\phi}}_V(t; \underline{y}_0, t_0) + B \underline{U}(t) . \end{aligned} \quad (4-197)$$

Therefore,

$$\frac{d}{dt} \{ \underline{\phi}_U(t; N^{-1} \underline{y}_0, t_0) \} = \underline{F}(\underline{\phi}_U(t; N^{-1} \underline{y}_0, t_0)) + \underline{B} \underline{U}(t) .$$

(4-199)

From (2-28),

$$\begin{aligned} \underline{\phi}_U(t; N^{-1} \underline{y}_0, t_0) &= N^{-1} \hat{\underline{\phi}}_V(t; \underline{y}_0, t_0) \\ &= N^{-1} \hat{\underline{\phi}}_V(t, N^{-1} \hat{N} \hat{\underline{\phi}}_V(t_1; \underline{y}_0, t_0), t_1) \\ &= N^{-1} \hat{\underline{\phi}}_V(t, N^{-1} \underline{\phi}_U(t_1; \underline{x}_0, t_0), t_1) \\ &= \underline{\phi}_U(t; \underline{\phi}_U(t_1; \underline{x}_0, t_0), t_1), \text{ for all } t \geq t_1 \geq t_0 . \end{aligned}$$

(4-200)

From (4-195), (4-199), and (4-200), the function

$\underline{\phi}_U(t; N^{-1} \underline{y}_0, t_0)$ satisfies the characteristics of the solution of (4-192); therefore the assertion is justified.

Chapter 5

THE OPTIMAL FEEDBACK CONTROL LAW AND THE INVERSE PROBLEM OF THE OPTIMAL REGULATOR

The purpose of this chapter is to review work that has been done on the inverse problem of the optimal regulator. Initially, however, the problem of optimal control as explained in Chapter 1 is mathematically restated. A theoretical background for the problem of the optimal regulator is given in Section 5.2.2 based on the principle of optimality and Caratheodory's lemma. In Section 5.3 studies of the inverse problem by Kalman, Suga and Thau are reviewed. Finally comments about these studies are given.

5.1 Formulation of the Optimal Control Problem and the Inverse Problem

Initial and final conditions for a system of objects are generally defined as manifolds in $R^n \times R^1$. It is convenient to call them starting manifolds M_s and terminating manifolds M_t , and the space $R^n \times R^1$ a motion space. The magnitudes of the control variables and the state variables may be restricted during control action to subdomains of $R^m \times R^1$ and $R^n \times R^1$ for practical reasons, e.g., structural design limitations. Call these admissible domains

the available control region, A.C.R., and the available state region, A.S.R.

Definition 5-1: Suitable Control.

An admissible control function $\underline{v}(t)$ defined on $[t_0, t_1] \in \mathbb{R}^1$ is said to be suitable for the problem specifications, or simply a suitable control, if it remains in the A.C.R. and provides a solution to (2-21) such that

$$(\hat{\phi}_{\underline{v}}(t_1; \underline{y}_0, t_0), t_1) \in M_t \quad (5-1)$$

which remains in A.S.R. When the time interval is given by $[t_0, \infty)$, (5-1) can be restated: for each $\epsilon > 0$, there exists a $T > 0$ satisfying

$$\inf_{\underline{y}_1 \in M_t} \{ ||\hat{\phi}_{\underline{v}}(\tau; \underline{y}_0, t_0) - \underline{y}_1|| \} < \epsilon, \text{ for all } \tau \geq t_0 + T. \quad (5-2)$$

A performance index for control action is usually given as

$$\tilde{J}[\underline{y}_0, t_0, \underline{v}(t)] = \tilde{K}[\underline{y}_0, t_0, \underline{y}_1, t_1] + \int_{t_0}^{t_1} \tilde{L}(\hat{\phi}_{\underline{v}}(\tau; \underline{y}_0, t_0), \underline{v}(\tau), \tau) d\tau, \quad (5-3)$$

where t_1 is defined as

$$t_1 \triangleq \inf_t [\{t \in \mathbb{R}^1 / t \geq t_0, \inf_{\underline{y}_1 \in M_t} \|\phi_{\underline{v}}(t; \underline{y}_0, t_0) - \underline{y}_1\| = 0\}] . \quad (5-4)$$

The function $\tilde{K}[\]$, called a terminal cost, is a penalty for the choice of the starting and terminating points on M_s and M_t and is usually assumed to be nonnegative valued on $\mathbb{R}^n \times \mathbb{R}^1 \times \mathbb{R}^n \times \mathbb{R}^1$. If \tilde{K} is constant on M_s for each fixed $(\underline{y}_1, t_1) \in M_t$ or on M_t for each fixed $(\underline{y}_0, t_0) \in M_s$, or is constant on both M_s and M_t , the penalty function is constant and generally omitted from (5-3). The loss function, $\tilde{L}(\underline{y}, \underline{v}, t)$, can be considered as a penalty for each point in motion space and is generally assumed to be nonnegative valued on $\mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^1$. The problem of optimal control for an open loop control function is stated as follows. For a given set of problem specifications, i.e., a system equation, $M_s, M_t, A.C.R., A.S.R.$ and a performance index, find a $(\underline{y}_0, t_0)_0 \in M_s$ and a suitable control function $\underline{v}^0(t)$ to minimize the value of the performance index. The function $\underline{v}^0(t)$ is called the open loop optimal control function.

For the problem of the optimal regulator, there exists a collection of M_s , say M_s , but a unique M_t and optimal control functions are required for each element of M_s . For this problem, an optimal feedback

control law, $\underline{v}^0(\underline{y}, t)$, generally provides the optimal control function. Thus for each $(\underline{y}_0, t_0) \in M_s$,

$$\underline{v}^0(t; \underline{y}_0, t_0) = \underline{v}^0(\hat{\phi}_f(t; \underline{y}_0, t_0), t), \quad (5-5)$$

where $\hat{\phi}_f(t, \underline{y}_0, t_0)$ is a free solution of

$$\dot{\underline{y}} = \hat{\underline{F}}(\underline{y}, \underline{v}^0(\underline{y}, t), t). \quad (5-6)$$

Assuming suitable feedback control law instead of a performance index (5-3), it is possible to attempt to find all performance indices for which the assumed control law is optimal. This is called the inverse problem of the optimal regulator. Specifically, consider \mathcal{L} and \mathcal{V} to be spaces of all performance indices and all suitable feedback control laws for the given optimal regulator problem. The usual or forward problem of the optimal regulator can be given as a mapping in an optimum sense from \mathcal{L} to \mathcal{V} , while the inverse problem is from \mathcal{V} to \mathcal{L} . As the space \mathcal{L} is too large for analytical treatment, some additional assumptions usually restrict the objects to a subset of \mathcal{L} and \mathcal{V} , e.g., the restriction of \mathcal{L} to a sum of quadratic forms in \underline{y} and \underline{v} . [25]

5.2 Optimal Feedback Control Law

5.2.1 Statement of the Problem

The fundamental characteristics of an optimal feedback control law are considered under the assumptions that

(i) a system is given by

$$\dot{\underline{Y}} = \hat{\underline{F}}(\underline{Y}, \underline{V}, t), \quad (5-7)$$

where $\hat{\underline{F}}(\underline{Y}, \underline{V}, t)$ is defined on $R^n \times R^m \times R^1$ and is of class C_2 with respect to all arguments;

(ii) M_t , the final condition of the system, is a smooth manifold in $R^n \times R^1$ and M_s , the set of initial conditions composed of all reasonable points in $R^n \times R^1$;

(iii) A.S.R. and A.C.R. are the entire $R^n \times R^1$ and $R^m \times R^1$ spaces respectively;

(iv) a performance index is given by

$$\tilde{J}[\underline{Y}_0, t_0, \underline{V}(t)] = \tilde{K}[\underline{Y}_1, t_1] + \int_{t_0}^{t_1} \tilde{L}(\underline{Y}, \underline{V}, t) dt, \quad (5-8)$$

where

(iva) the terminal cost $\tilde{K}[\underline{Y}, t]$, considered only for final conditions, is of class C_2 with respect to all arguments, and

(ivb) the loss function $\tilde{L}(\underline{y}, \underline{v}, t)$ is of class C_1 with respect to all arguments.

Since M_s is a set of starting points, there is no terminal cost with respect to initial condition. It is possible to imbed the terminal cost of (5-8) into the integral. It has been shown^[27] that an optimal control function from each $(\underline{y}_0, t_0) \in M_s$ for the performance index (5-8) must be equal to that for

$$\hat{J}[\underline{y}_0, t_0, \underline{v}(t)] = \int_{t_0}^{t_1} \hat{L}(\underline{y}, \underline{v}, t) dt, \quad (5-9)$$

where

$$\hat{L}(\underline{y}, \underline{v}, t) = \tilde{L}(\underline{y}, \underline{v}, t) + \{\text{grad } \tilde{K}[\underline{y}, t]\}^T \hat{F}(\underline{y}, \underline{v}, t) + \frac{\partial \tilde{K}[\underline{y}, t]}{\partial t}. \quad (5-10)$$

$\hat{L}(\underline{y}, \underline{v}, t)$ is of class C_1 with respect to all arguments by the assumptions (iva) and (ivb). For convenience, therefore, the following analysis proceeds with (5-9) instead of (5-8).

5.2.2 Fundamental Lemma

The principle of optimality states that any portion of an optimal solution is also an optimal solution.^[2] Mathematically, let $\underline{v}^0(t)$, $t_0 \leq t \leq t_f$, be an optimal

control function from (\underline{y}_0, t_0) to (\underline{y}_f, t_f) . Then, if $t_0 \leq t_1 \leq t_2 \leq t_f$, the control $\underline{v}^0(t)$ considered on the interval $t_1 \leq t \leq t_2$ is an optimal control function from $(\hat{\phi}_{\underline{v}^0}(t_1; \underline{y}_0, t_0), t_1)$ to $(\hat{\phi}_{\underline{v}^0}(t_2; \underline{y}_0, t_0), t_2)$ with $\hat{\phi}_{\underline{v}^0}(t; \underline{y}_0, t_0)$, $t_1 \leq t \leq t_2$, the corresponding trajectory.

Caratheodory's lemma is a sufficiency statement of optimality.

Lemma 5-1: [5,26] Caratheodory's Lemma.

If there exists a suitable feedback control law $\underline{v}^*(\underline{y}, t)$ for this problem of the optimal regulator such that for all $(\underline{y}, t) \in \mathbb{R}^n \times \mathbb{R}^1$

$$\begin{aligned} \hat{L}(\underline{y}, \underline{v}^*(\underline{y}, t), t) &= 0 \\ \hat{L}(\underline{y}, \underline{v}, t) &\geq 0 \quad \text{if } \underline{v} \neq \underline{v}^*(\underline{y}, t), \end{aligned} \tag{5-11}$$

then the function $\underline{v}^*(\underline{y}, t)$ is an optimal feedback control law.

Necessarily the corresponding optimal performance index from every initial condition is identically zero.

5.2.3 Heuristic Approach to the Optimal Feedback Control Law [4,5,26,27]

Define a function $\hat{V}^o(\underline{y}, t)$, called an optimal performance index function, such that, for each $(\underline{y}_o, t_o) \in M_s$,

$$\begin{aligned}\hat{V}^o(\underline{y}_o, t_o) &= \inf_{\underline{v}(t; \underline{y}_o, t_o)} \{ \hat{J}[\underline{y}_o, t_o, \underline{v}(t; \underline{y}_o, t_o)] \} \\ &= \hat{J}[\underline{y}_o, t_o, \underline{v}^o(t; \underline{y}_o, t_o)] ,\end{aligned}\quad (5-12)$$

where $\underline{v}(t; \underline{y}_o, t_o)$ is any suitable open loop control function from (\underline{y}_o, t_o) . In the following, it is assumed that $\hat{V}^o(\underline{y}, t)$ is of class C_2 with respect to all arguments. From the definition of t_1 in (5-4), if $(\underline{y}_o, t_o) \in M_t$, then

$$\hat{V}^o(\underline{y}_o, t_o) = 0 . \quad (5-13)$$

Consider an arbitrary $(\underline{y}_o, t_o) \in M_s$ and assume a corresponding optimal control function $\underline{v}^o(t; \underline{y}_o, t_o)$.

Consider also a perturbed control function $\underline{v}_d(t; \underline{y}_o, t_o)$ from (\underline{y}_o, t_o) such that for an incremental Δt_o ,

$$\underline{v}_d(\tau; \underline{y}_o, t_o) \triangleq \underline{v}_a \in R^m \quad \text{for } \tau \in [t_o; t_o + \Delta t_o] \quad (5-14)$$

transforms the system condition to $(\underline{y}_0 + \Delta \underline{y}_0, t_0 + \Delta t_0)$

and

$$\underline{v}_d(\tau; \underline{y}_0, t_0) = \underline{v}^0(\tau; \underline{y}_0 + \Delta \underline{y}_0, t_0 + \Delta t_0) \quad \text{for } \tau > t_0 + \Delta t_0. \quad (5-15)$$

The performance index becomes

$$\hat{J}[\underline{y}_0, t_0, \underline{v}_d(t; \underline{y}_0, t_0)] = \int_{t_0}^{t_0 + \Delta t_0} \hat{L}(\hat{\underline{v}}_{\underline{v}_d}(\tau; \underline{y}_0, t_0), \underline{v}_\alpha, \tau) d\tau + \hat{V}^0(\underline{y}_0 + \Delta \underline{y}_0, t_0 + \Delta t_0). \quad (5-16)$$

Since Δt_0 is small and $\hat{L}(\underline{y}, \underline{v}, t)$ and $\hat{F}(\underline{y}, \underline{v}, t)$ are of class C_1 , this integral can be approximated as

$$\hat{L}(\underline{y}_0, \underline{v}_\alpha, t_0) \Delta t_0 + o(\Delta t_0), \quad (5-17)$$

and $\hat{V}^0(\underline{y}_0 + \Delta \underline{y}_0, t_0 + \Delta t_0)$ as

$$\begin{aligned} \hat{V}^0(\underline{y}_0, t_0) + \{ [\text{grad } \hat{V}^0(\underline{y}_0, t_0)]^T \hat{F}(\underline{y}_0, \underline{v}_\alpha, t_0) + \hat{V}_t^0(\underline{y}_0, t_0) \} \Delta t_0 \\ + o(\Delta t_0), \end{aligned} \quad (5-18)$$

where

$$\Delta \underline{y}_0 = \hat{F}(\underline{y}_0, \underline{v}_d, t_0) \Delta t_0, \quad (5-19)$$

$o(\Delta t_0)$ and $\sigma(\Delta t_0)$ are higher orders of Δt_0 , i.e.,

$$\lim_{\Delta t_0 \rightarrow 0} \frac{o(\Delta t_0)}{\Delta t_0} = 0 \quad (5-20)$$

$$\lim_{\Delta t_0 \rightarrow 0} \frac{\sigma(\Delta t_0)}{\Delta t_0} = 0 ,$$

and

$$\hat{v}_t^o(\underline{y}, t) \triangleq \frac{\partial \hat{v}^o(\underline{y}, t)}{\partial t} . \quad (5-21)$$

From the assumption that $\underline{v}^o(t; \underline{y}_0, t_0)$ is an optimum control function

$$\hat{v}^o(\underline{y}_0, t_0) = \hat{J}[\underline{y}_0, t_0, \underline{v}^o(t; \underline{y}_0, t_0)] \leq \hat{J}[\underline{y}_0, t_0, \underline{v}_d(t; \underline{y}_0, t_0)]$$

for $\underline{v}_d \in R^m$. (5-22)

But from (5-18~20),

$$\begin{aligned} & \{ \hat{L}(\underline{y}_0, \underline{v}_\alpha, t_0) + [\text{grad } \hat{v}^o(\underline{y}_0, t_0)]^T \hat{F}(\underline{y}_0, \underline{v}_\alpha, t_0) + \hat{v}_t^o(\underline{y}_0, t_0) \} \Delta t_0 \\ & + o(\Delta t_0) + \sigma(\Delta t_0) \geq 0 . \end{aligned} \quad (5-23)$$

Dividing by Δt_0 and taking $\lim_{\Delta t_0 \rightarrow 0}$,

$$\hat{L}(\underline{y}_0, \underline{v}_\alpha, t_0) + [\text{grad } \hat{V}^0(\underline{y}_0, t_0)]^T \hat{F}(\underline{y}_0, \underline{v}_\alpha, t_0) + \hat{V}_t^0(\underline{y}_0, t_0) \geq 0, \quad (5-24)$$

where the equality is satisfied if

$$\underline{v}_\alpha = \underline{v}^0(t_0; \underline{y}_0, t_0). \quad (5-25)$$

As any point of any motion can be regarded as an initial condition for an optimal control problem, (5-24) is valid for all points in $R^n \times R^1$, according to the principle of optimality. The important point is that an optimal feedback control law $\underline{v}^0(\underline{y}, t)$ is a suitable control law which satisfies the equality of (5-24) at every point in $R^n \times R^1$.

The function $\hat{V}^0(\underline{y}, t)$ describes the value of the performance index for the optimal control function for each (\underline{y}, t) . Its time derivative, governed by the system equation, is

$$\dot{\hat{V}}^0(\underline{y}, t) = [\text{grad } \hat{V}^0(\underline{y}, t)]^T \hat{F}(\underline{y}, \underline{v}, t) + \hat{V}_t^0(\underline{y}, t), \quad (5-26)$$

which depends upon \underline{v} at each (\underline{y}, t) . Thus (5-24) for any (\underline{y}_0, t_0) , is generally written as

$$\hat{L}(\underline{y}, \underline{v}, t) + \dot{\hat{V}}^0(\underline{y}, t) \geq 0 \quad \text{for } (\underline{y}, t) \in R^n \times R^1. \quad (5-27)$$

The equality is satisfied when an optimal control is used. The optimal feedback control law $\underline{v}^0(\underline{y}, t)$ can be described as the suitable feedback control law which satisfies the equality in (5-27) at every $(\underline{y}, t) \in \mathbb{R}^n \times \mathbb{R}^1$.

For a calculation of the function $\hat{V}^0(\underline{y}, t)$, it is convenient to introduce a function called the Hamiltonian of the problem such that

$$\hat{H}(\underline{y}, \underline{\psi}, \underline{v}, t) \triangleq \underline{\psi}^T \hat{F}(\underline{y}, \underline{v}, t) + \hat{L}(\underline{y}, \underline{v}, t), \quad (5-28)$$

where $\underline{\psi}$ is an arbitrary n -dimensional vector of variables ψ_i . Then from (5-13) and (5-24), $\hat{V}^0(\underline{y}, t)$ is a solution of a specific partial differential equation such that

$$\begin{aligned} \min_{\underline{v} \in \mathbb{R}^m} \left\{ \hat{H}(\underline{y}, \frac{\partial \hat{V}(\underline{y}, t)}{\partial \underline{y}}, \underline{v}, t) \right\} + \hat{V}_t(\underline{y}, t) &= 0 \\ \text{for all } (\underline{y}, t) &\in \mathbb{R}^n \times \mathbb{R}^1 \end{aligned} \quad (5-29)$$

with the boundary condition

$$\hat{V}(\underline{y}, t) = 0 \quad \text{for all } (\underline{y}, t) \in M_t. \quad (5-30)$$

Identifying $\hat{H}(\underline{y}, \underline{\psi}, \underline{v}, t)$ as a function of \underline{v} at each $(\underline{y}, \underline{\psi}, t) \in \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^1$, denote $\widetilde{V}(\underline{y}, \underline{\psi}, t)$ as a function to provide the absolute minimum value to $\hat{H}(\underline{y}, \underline{\psi}, \underline{v}, t)$ everywhere.

Then (5-29) can be written

$$\hat{H}(\underline{y}, \frac{\partial \hat{V}(\underline{y}, t)}{\partial \underline{y}}, \underline{V}(\underline{y}, \frac{\partial \hat{V}(\underline{y}, t)}{\partial \underline{y}}, t)) + \hat{V}_t(\underline{y}, t) = 0, \quad (5-31)$$

which is generally called the Hamilton-Jacobi equation of the problem. As the fundamental condition for optimality in dynamic programming, Bellman called (5-29) Bellman's equation and its solution a Bellman function.^[1]

Alternately, it is possible to recognize $\text{grad}[\hat{V}(\underline{y}, t)]$ as an independent variable in the Hamilton-Jacobi equation. Pontryagin^[2] developed a different technique of calculating open loop optimal control functions, in effect, by doing so. This is called his maximum principle (or sometimes the minimum principle). Systematically calculating open loop optimal control functions from various initial conditions and observing their common characteristics, a synthesis of an optimal feedback control law is possible.^[2] However, at this point two major difficulties exist for these calculations. From a practical aspect, no general method for solving the Hamilton-Jacobi equation is known and solutions can only be calculated for a few classes of problems with restrictive assumptions. Secondly, from a theoretical aspect, an optimal performance index function $\hat{V}^o(\underline{y}, t)$ must be a solution to the Hamilton-Jacobi equation but this is not a sufficient condition. Thus a technique

to identify the $\hat{V}^0(\underline{Y}, t)$ among solutions must be developed when more than one solution exists. If the number of solutions to the Hamilton-Jacobi equation is finite and small, $\hat{V}^0(\underline{Y}, t)$ may be identified by comparing each of the solutions. These difficulties are not avoided if the maximum principle is used.

This second difficulty can be avoided, however, if a unique solution can be shown to exist for the Hamilton-Jacobi equation. But this depends upon the specified $\hat{L}(\underline{Y}, \underline{V}, t)$ and $\hat{F}(\underline{Y}, \underline{V}, t)$ in the problem.

Definition 5-2: A Normal Hamiltonian. [5]

If the Hamiltonian of a problem is minimized by a unique value of $\underline{V} \in R^m$ at each $(\underline{Y}, \underline{\psi}, t) \in R^n \times R^n \times R^1$, then the Hamiltonian is said to be normal. In this case, the function $\underline{V}^0(\underline{Y}, \underline{\psi}, t)$ which provides the absolute minimum to the Hamiltonian is called the H-minimal control law.

In order to be normal, the Hamiltonian must be a strict convex function of \underline{V} , i.e., the matrix

$$\frac{\partial^2 \hat{H}}{\partial \underline{V} \partial \underline{V}} = \begin{bmatrix} \frac{\partial^2 \hat{H}}{\partial v_1 \partial v_1} & \cdot & \cdot & \cdot & \frac{\partial^2 \hat{H}}{\partial v_m \partial v_1} \\ \cdot & & & & \cdot \\ \cdot & & & & \cdot \\ \cdot & & & & \cdot \\ \frac{\partial^2 \hat{H}}{\partial v_1 \partial v_m} & \cdot & \cdot & \cdot & \frac{\partial^2 \hat{H}}{\partial v_m \partial v_m} \end{bmatrix} \quad (5-32)$$

must be positive definite for every argument. [5] Then the H-minimal control law is calculated from

$$\begin{aligned} \underline{0} &= \left. \frac{\partial \hat{H}(\underline{Y}, \underline{\psi}, \underline{V}, t)}{\partial \underline{V}} \right|_{\underline{V} = \underline{V}^0(\underline{Y}, \underline{\psi}, t)} \\ &= \{ \underline{\psi}^T \left[\frac{\partial \hat{F}(\underline{Y}, \underline{V}, t)}{\partial \underline{V}} \right] + \frac{\partial \hat{L}(\underline{Y}, \underline{V}, t)}{\partial \underline{V}} \} \quad (5-33) \\ &\quad \underline{V} = \underline{V}_0(\underline{Y}, \underline{\psi}, t) \end{aligned}$$

Theorem 5-1: [5]

If the Hamiltonian of a problem is normal, then a solution of the Hamilton-Jacobi equation is the optimal performance index function $\hat{V}^0(\underline{Y}, t)$.

This theorem is proved by the uniqueness of the solution to the Hamilton-Jacobi equation and Lemma 5-1. The

corresponding optimal feedback control law is given by

$$\underline{v}^o(\underline{y}, \frac{\partial \hat{V}^o(\underline{y}, t)}{\partial \underline{y}}, t) . \quad (5-34)$$

Examples of this theory are given in the literature. [5,6,26]

5.2.4 Miscellaneous Comments

The preceding considerations are given under restrictive assumptions for simplicity of discussion. Studies have been made for fewer restrictions. Based on measure theory, Bridgland^[28] generalized the theory with relaxed assumptions on $\hat{F}(\underline{y}, \underline{v}, t)$ and $\hat{L}(\underline{y}, \underline{v}, t)$ and generalized the integral interval of the performance index to $[t_0, \infty)$. Boltyanskii^[29] extended the theory for a continuous $\hat{V}^o(\underline{y}, t)$, (not necessarily of class C_2), but under other conditions. Thus the existence of continuous $\hat{V}^o(\underline{y}, t)$ (not necessarily of C_2) for the time optimal control problem,^[2] as mentioned by Pontryagin, is justified.

5.3 Review of Studies on the Inverse Problem

Results reported for the inverse problem are reviewed but with interest directed to those with similar problem assumptions as those made Chapters 6 and 7.

5.3.1 Study of Kalman [25]

Kalman considered an inverse problem with the assumptions that

- (i) system is completely controllable and given by

$$\dot{\underline{y}} = \underline{\hat{A}} \underline{y} + \underline{\hat{B}} v, \quad (5-35)$$

where $\underline{\hat{A}}$ and $\underline{\hat{B}}$ is constant and the control vector v is one dimensional,

- (ii) the control law is time-invariant such that

$$v(\underline{y}) = - \underline{\hat{K}}^T \underline{y}, \quad (5-36)$$

where $\underline{\hat{K}}$ is an $n \times 1$ constant matrix and all real parts of the eigenvalues of $(\underline{\hat{A}} - \underline{\hat{B}} \underline{\hat{K}}^T)$ are negative, thus predicting the final desired condition of system to be $\underline{y} = \underline{0}$,

- (iii) the form of the performance index is restricted to

$$\int_0^{\infty} (\underline{y}^T \underline{\hat{H}}^T \underline{\hat{H}} \underline{y} + v^2) dt, \quad (5-37)$$

where $\underline{\hat{H}}$ is a $n_1 \times n$ constant matrix with the rank $n_1 \leq n$ and

$$\text{rank } [\underline{\hat{H}}, \underline{\hat{A}} \underline{\hat{H}}, \dots, \underline{\hat{A}}^{n-1} \underline{\hat{H}}] = n, \quad (5-38)$$

$\underline{\hat{H}}^T \underline{\hat{H}}$ is positive semidefinite by (ii) of Corollary 2-1. In another study, [26] Kalman proved that the optimal performance index function for this problem is given by

$$\hat{V}^0(\underline{y}) = \underline{y}^T \underline{\hat{P}} \underline{y}, \quad (5-39)$$

where $\underline{\hat{P}}$ is a symmetric positive semidefinite matrix. Using Theorem 5-1, he showed the corresponding Hamilton-Jacobi equation to be

$$\underline{y}^T \{ \underline{\hat{P}} \underline{\hat{A}} + \underline{\hat{A}}^T \underline{\hat{P}} + \underline{\hat{H}}^T \underline{\hat{H}} \} \underline{y} + 2 \underline{v}^T \underline{\hat{B}}^T \underline{\hat{P}} \underline{y} + \underline{v}^2 = 0 \quad (5-40)$$

and the optimal feedback control law

$$\underline{v}^0(\underline{y}) = - \underline{\hat{B}}^T \underline{\hat{P}} \underline{y}. \quad (5-41)$$

Theorem 5-2: [25]

For the completely controllable system (5-35) and the performance index (5-37), the necessary and sufficient condition for the control law (5-36) to be a stable optimal control law is that there exists a positive definite, symmetric matrix $\underline{\hat{P}}$ which uniquely satisfies the algebraic relations

$$\hat{\underline{P}} \hat{\underline{B}} = \hat{\underline{K}} \quad (5-42)$$

and

$$-\hat{\underline{P}} (\hat{\underline{A}} - \hat{\underline{B}} \hat{\underline{K}}^T) - (\hat{\underline{A}} - \hat{\underline{B}} \hat{\underline{K}}^T)^T \hat{\underline{P}} = \hat{\underline{H}}^T \hat{\underline{H}} + \hat{\underline{K}} \hat{\underline{K}}^T. \quad (5-43)$$

Restated, if an arbitrary symmetric positive definite matrix $\hat{\underline{P}}$ determined by

$$\hat{\underline{P}} \hat{\underline{B}} = \hat{\underline{K}} \quad (5-44)$$

also allows a solution for $\hat{\underline{H}}$ from (5-43) which also satisfies (5-38), then the resulting performance index (5-37) is optimized by the given control law (5-36).

This study of Kalman was the first published on the inverse problem. The results revealed the positive definiteness of $\hat{\underline{V}}^0(\underline{y})$ for the loss function of (5-37) with the condition of (5-38). Although the problem assumptions are relatively simple, others were subsequently encouraged to attempt to generalize them. In Kalman's original paper, [25] the solution to the inverse problem was also discussed from the viewpoint of the frequency domain techniques of synthesizing optimal feedback control systems.

5.3.2 Study of Suga^[30]

Suga considered an inverse problem under the assumptions that:

(i) the system is given by

$$\dot{\underline{y}} = \underline{\hat{A}}(t) \underline{y} + \underline{\hat{B}}(t) \underline{v} , \quad (5-45)$$

where $\underline{\hat{A}}(t)$ and $\underline{\hat{B}}(t)$ have continuous first derivatives and

$$\text{rank } \underline{\hat{B}}(t) = r = m \leq n \text{ (full rank)} , \quad (5-46)$$

(ii) the control law is given by

$$\underline{v}(\underline{y}, t) = - \underline{\hat{K}}^T(t) \underline{y} , \quad (5-47)$$

where $\underline{\hat{K}}(t)$ has a continuous first derivative,

(iii) the form of the performance index is restricted to

$$\int_{t_0}^T \{ \underline{\hat{L}}(\underline{y}, t) + \underline{v}^T \underline{\hat{R}}(t) \underline{v} \} dt , \quad (5-48)$$

where T is fixed and $\underline{\hat{R}}(t)$ is a given $r \times r$ positive definite symmetric matrix with a continuous time derivative at every $t \in R^1$. The Hamiltonian is normal from these

assumptions. Thus, the author stated the following.

Theorem 5-3:

Suppose that $\hat{\underline{K}}(t)$ is specified so that $\hat{\underline{R}}(t) \hat{\underline{K}}(t) \hat{\underline{B}}(t)$ is a symmetric matrix. Then the performance index is optimized if and only if it is given by

$$\begin{aligned} \hat{\underline{L}}(\underline{Y}, t) = & \underline{Y}^T [\hat{\underline{K}}(t) \hat{\underline{R}}(t) \hat{\underline{K}}^T(t) - \hat{\underline{A}}^T(t) \hat{\underline{P}}(t) - \hat{\underline{P}}(t) \hat{\underline{A}}(t) - \hat{\underline{P}}(t)] \underline{Y} \\ & + \left[\frac{\partial \hat{\underline{r}}(\underline{Y}, t)}{\partial \underline{Y}} \right]^T \hat{\underline{A}}(t) \underline{Y} + \frac{\partial \hat{\underline{r}}(\underline{Y}, t)}{\partial t}, \end{aligned} \quad (5-49)$$

where $\hat{\underline{P}}(t)$ is an $n \times n$ symmetric matrix of class C_2 satisfying

$$\hat{\underline{P}}(t) \hat{\underline{B}}(t) = \hat{\underline{K}}(t) \hat{\underline{R}}(t) \quad (5-50)$$

and

$$\hat{\underline{P}}(T) = [0], \quad (5-51)$$

and $\hat{\underline{r}}(\underline{Y}, t)$ is an arbitrary scalar function of class C_2 in all arguments satisfying

$$\hat{\underline{B}}^T(t) \frac{\partial \hat{\underline{r}}(\underline{Y}, t)}{\partial \underline{Y}} = \underline{0} \quad (5-52)$$

and

$$\left. \frac{\partial \hat{r}(\underline{y}, t)}{\partial \underline{y}} \right|_{t=T} = \underline{0} . \quad (5-53)$$

Then the resulting optimal performance index function is

$$\hat{V}^0(\underline{y}, t) = \underline{y}^T \hat{P}(t) \underline{y} - \hat{r}(\underline{y}, t) + \hat{r}(\underline{y}_1, T) , \quad (5-54)$$

where \underline{y}_1 is a final state of the system and the last term is constant, say $\gamma(T)$, because of (5-53) and the fixed T .

If a stable control law is a suitable control law which provides asymptotic stability in the large for the synthesized feedback control system relative to the origin, then the following corollary exists.

Corollary 5-1:

Assume that a stable control law is given and $T \rightarrow \infty$. Then the theorem is still valid with a change of (5-51) to

$$\lim_{t \rightarrow \infty} \hat{P}(t) \hat{\Phi}_{\underline{y}}(t; \underline{y}_0, t_0) = \underline{0} , \quad \text{for every } (\underline{y}_0, t_0) \in \mathbb{R}^n \times \mathbb{R}^1 . \quad (5-55)$$

Based on this work, Suga observed the following points for his problem.

(1). The symmetry of $\hat{R}(t) \hat{K}^T(t) \hat{B}(t)$ is necessary for the inverse problem to be meaningful. In the study of Kalman, this condition is trivial because of the single input.

(2). The additive terms of $\hat{L}(\underline{y}, t)$ which are associated with the function $\hat{r}(\underline{y}, t)$ don't affect the optimal feedback control law. To show this, the value of the optimal performance index for each initial condition $(\underline{y}_0, t_0) \in R^n \times R^1$ becomes, from (5-54),

$$\hat{V}^0(\underline{y}_0, t_0) = \underline{y}_0^T \hat{P}(t_0) \underline{y}_0 - \hat{r}(\underline{y}_0, t_0) - \gamma(T) . \quad (5-56)$$

Only the first term is sensitive to the control function; the last two are independent to $\underline{v}(t)$. If $r = n$, $\hat{B}(t)$ is nonsingular and there can be no $\hat{r}(\underline{y}, t)$ because of the restrictions of (5-52). The larger the value of $n - r$, the more flexibility of $\hat{L}(\underline{y}, t)$ exists through this \hat{r} . Suga expressed this idea as a flexibility of loss functions.

As an example, consider a system given by

$$\begin{bmatrix} \dot{y}_1 \\ \dot{y}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} v , \quad (5-57)$$

a feedback control law given by

$$U(\underline{y}) = - [1, \sqrt{3}] \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}, \quad (5-58)$$

and a performance index with the form

$$\int_0^{\infty} \{ \hat{L}(\underline{y}, t) + U^2 \} dt. \quad (5-59)$$

(5-58) is a stable control law and the integral interval of the performance index is $[0, \infty)$; thus Corollary 5-1 can be applied. From (5-51), it follows that

$$\hat{P}(t) = \begin{bmatrix} g(t) & 1 \\ 1 & \sqrt{3} \end{bmatrix}, \quad (5-60)$$

where $g(t)$ is undetermined. Subsequently (5-49) becomes

$$\hat{L}(\underline{y}, t) = \underline{y}^T \begin{bmatrix} 1-\dot{g}(t) & \sqrt{3}-g(t) \\ \sqrt{3}-g(t) & 1 \end{bmatrix} \underline{y} + \frac{\partial \hat{r}(\underline{y}, t)}{\partial y_1} y_2 + \frac{\partial \hat{r}(\underline{y}, t)}{\partial t} \quad (5-61)$$

and (5-54) is

$$\hat{V}^0(\underline{y}, t) = \underline{y}^T \begin{bmatrix} g(0) & 1 \\ 1 & \sqrt{3} \end{bmatrix} \underline{y} - \hat{r}(\underline{y}, 0) - \hat{r}(0, \infty) . \quad (5-62)$$

From Corollary 5-1,

$$\lim_{t \rightarrow \infty} \begin{bmatrix} g(t) & 1 \\ 1 & \sqrt{3} \end{bmatrix} \begin{bmatrix} \hat{\phi}_1(t; \underline{y}_0, t_0) \\ \hat{\phi}_2(t; \underline{y}_0, t_0) \end{bmatrix} = \underline{0} , \text{ for every } (\underline{y}_0, t_0) \in \mathbb{R}^n \times \mathbb{R}^1 . \quad (5-63)$$

From (5-52) and (5-53),

$$\hat{B}^T(t) \frac{\partial \hat{r}(\underline{y}, t)}{\partial \underline{y}_1} = \frac{\partial \hat{r}(\underline{y}, t)}{\partial \underline{y}_2} = 0 , \quad \text{for all } t \in \mathbb{R}^1 , \quad (5-64)$$

and

$$\left. \frac{\partial \hat{r}(\underline{y}, t)}{\partial \underline{y}} \right|_{t \rightarrow \infty} = \underline{0} . \quad (5-65)$$

By various choices of $g(t)$ and $\hat{r}(\underline{y}, t)$, the following optimized $\hat{L}(\underline{y}, t)$ were found.

	$g(t)$	$\hat{r}(t)$	$\hat{L}(\underline{y}, t)$	
(a)	$\sqrt{3} - \frac{\alpha}{2}$	0	$y_1^2 + \alpha y_1 y_2 + y_2^2$	α : arbitrary real
(b)	$\sqrt{3} - e^{-t}$	0	$(1 - e^{-t}) y_1^2 + 2e^{-t} y_1 y_2 + y_2^2$	
(c)	$\sqrt{3} + t$	0	$y_2^2 - 2t y_1 y_2$	
(d)	$\sqrt{3} - 1$	$\frac{y_1^{m+1}}{m+1}$	$(y_1 + y_2)^2 + y_1^m y_2$	m : positive integer

By generalizing the assumptions of Kalman, Suga discovered the fundamental composition of $\hat{L}(\underline{y}, t)$ in relation to the dummy function $\hat{r}(\underline{y}, t)$ and the symmetry condition of $\hat{R}(t)$ $\hat{K}(t)$ $\hat{B}(t)$. However, as in the demonstrated examples, the calculation does not increase a nonnegative $\hat{L}(\underline{y}, t)$ on $R^n \times R^1$. From a practical viewpoint for the forward optimal control problem, $\hat{L}(\underline{y}, t)$ is is usually assumed to be nonnegative as a penalty function with respect to regulating errors. Thus Suga's examples tend to be unrealistic.

5.3.3 Study of Thau^[31]

Thau considered an inverse problem such that

- (i) a system is time invariant and given by

$$\dot{\underline{y}} = \hat{F}(\underline{y}) + \hat{B} \underline{v}, \quad (5-66)$$

where $\hat{F}(\underline{y})$ is of class C_2 ,

- (ii) a control law is given by $\underline{v}(\underline{y})$, a stable control law of class C_2 ,

- (iii) the form of performance indices is restricted to

$$\int_0^{\infty} \{\hat{L}(\underline{y}) + \hat{R}(\underline{v})\} dt, \quad (5-67)$$

where $\hat{L}(\underline{Y})$ and $\hat{R}(\underline{V})$ are of class C_2 and

$$(a) \quad \hat{L}(0) = 0 \quad (5-68)$$

$$(b) \quad \hat{R}(0) = 0 \quad (5-69)$$

$$(c) \quad \frac{\partial \hat{R}(\underline{V})}{\partial \underline{V}} \text{ has a one-to-one correspondence from } R^m \text{ to } R^m.$$

$$(d) \quad \left[\frac{\partial^2 \hat{R}(\underline{V})}{\partial \underline{V} \partial \underline{V}} \right] \text{ is positive definite for a normal Hamiltonian.}$$

It was then asserted, based on Theorem 5-1, that $\hat{V}^o(\underline{Y})$ of class C_2 is the optimal performance index function and $\underline{V}(\underline{Y})$ is the optimal feedback control law for

$$\hat{L}(\underline{Y}) = - \hat{R}(\underline{V}(\underline{Y})) - \left[\frac{\partial \hat{V}^o}{\partial \underline{Y}} \right]^T \{ \underline{F}(\underline{Y}) + \underline{B} \underline{V}(\underline{Y}) \} \quad (5-70)$$

if and only if $\hat{V}^o(\underline{Y})$ satisfies the Hamilton-Jacobi equation of

$$\hat{H}(\underline{Y}, \frac{\partial \hat{V}^o}{\partial \underline{Y}}, \underline{V}(\underline{Y}), t) = 0, \quad \text{for all } (\underline{Y}, t) \in R^n \times R^1 \quad (5-71)$$

and $\underline{V}(\underline{Y})$ is given as

$$\underline{V}(\underline{Y}) = \hat{\mathcal{R}}(-B^T \frac{\partial \hat{V}^o}{\partial \underline{Y}}), \quad (5-72)$$

where $\hat{\mathcal{R}}(\cdot)$ is the H-minimal of the Hamiltonian.

Further investigations were directed to restrict the resulting $\hat{V}^0(\underline{Y})$ to a quadratic form $\underline{Y}^T \hat{\underline{P}} \underline{Y}$. Then (5-70) and (5-71) become

$$\hat{L}(\underline{Y}) = -\hat{R}(\underline{V}(\underline{Y})) - \underline{Y}^T \hat{\underline{P}} \{ \hat{\underline{F}}(\underline{Y}) + \hat{\underline{B}} \underline{V}(\underline{Y}) \} \quad (5-73)$$

and

$$\underline{V}(\underline{Y}) = \hat{\mathcal{R}}(-\underline{B}^T \underline{P} \underline{Y}) . \quad (5-74)$$

The following three cases were investigated.

(1) If (a) the system (5-67) is linear,

$$\dot{\underline{Y}} = \hat{\underline{A}} \underline{Y} + \hat{\underline{B}} \underline{V} , \quad (5-75)$$

(b) the form of the loss function is restricted to

$$\begin{aligned} \hat{L}(\underline{Y}) &= \underline{Y}^T \hat{\underline{H}}^T \hat{\underline{H}} \underline{Y} \\ \hat{R}(\underline{V}) &= \underline{V}^T \underline{V} , \end{aligned} \quad (5-76)$$

where $\hat{\underline{H}}$ is any $n_1 \times n$ matrix with the rank $n_1 \leq n$,

(c) a linear control law is assumed

$$\underline{V}(\underline{Y}) = -\hat{\underline{K}}^T \underline{Y} , \quad (5-77)$$

then (5-73) and (5-74) reduce to

$$\underline{\hat{K}} = \underline{\hat{P}} \underline{\hat{B}} \quad (5-78)$$

and

$$-\underline{\hat{P}} (\underline{\hat{A}} - \underline{\hat{B}} \underline{\hat{K}}^T) - (\underline{\hat{A}} - \underline{\hat{B}} \underline{\hat{K}}^T)^T \underline{\hat{P}} = \underline{\hat{H}}^T \underline{\hat{H}} + \underline{\hat{K}} \underline{\hat{K}}^T, \quad (5-79)$$

which are generalizations of the equations of Kalman, (5-42) and (5-43). Theoretical justifications are not given, however.

(2) Assume in addition to (1) that (5-75) has a single input with

$$\underline{\hat{B}} = [0, 0, \dots, 0, 1]^T. \quad (5-80)$$

Then for a positive definite $\underline{\hat{P}}$, it is necessary for \hat{k}_n
 $\underline{\hat{K}} = [\hat{k}_1, \hat{k}_2, \dots, \hat{k}_n]^T$ to be positive. In fact, from
 (5-78)

$$\hat{k}_n = \hat{p}_{nn} \quad (5-81)$$

of the matrix $\underline{\hat{P}}$. If $\hat{k}_n = 0$, then the sign-definiteness of $\underline{\hat{P}}$ fails. More generally, if the elements of $\underline{\hat{B}}$ are

either 0 or 1, for a positive definite $\hat{\underline{P}}$ it is necessary to have positive elements $\hat{\underline{K}}$ which correspond to the unit elements of $\hat{\underline{B}}$.

(3) For a single input linear system, consider a control law such that

$$\mathcal{V}(\hat{\underline{K}}^T \underline{y}), \quad (5-82)$$

where the $\hat{\underline{K}}$ is a $n \times 1$ matrix and the $\mathcal{V}(\sigma)$ is given by

$$\mathcal{V}(\sigma) = \sum_{\substack{i=1 \\ i \text{ odd}}}^{\infty} a_i \sigma^i, \quad a_i > 0 \text{ for all } i. \quad (5-83)$$

It was shown that if the inverse function of $\mathcal{V}(\sigma)$ can be expressed as a power series

$$\tilde{\mathcal{V}}(\sigma) = \sum_{\substack{i=1 \\ i \text{ odd}}}^{\infty} C_i \sigma^i, \quad (5-84)$$

and if $\hat{\underline{P}}$ is positive definite, then each coefficient C_i can be determined explicitly in terms of the coefficients a_i and the components of $\hat{\underline{P}}$, and $\hat{\mathcal{R}}(\mathcal{V}(\sigma))$ can be calculated as

$$\hat{R}(V(\sigma)) = \sum_{\substack{i=2 \\ \text{even}}}^{\infty} d_i \sigma^i . \quad (5-85)$$

Substituting,

$$\begin{aligned} \hat{L}(\underline{y}) = & -\frac{1}{2} \underline{y}^T (\hat{\underline{P}} \hat{\underline{A}} + \hat{\underline{A}}^T \hat{\underline{P}}) \underline{y} - \sum_{\substack{i=1 \\ i \text{ odd}}}^{\infty} \{ \underline{y}^T \hat{\underline{P}} \hat{\underline{B}} a_i (\hat{\underline{K}}^T \underline{y})^i \\ & + d_{i+1} (\hat{\underline{K}}^T \underline{y})^{i+1} \} . \end{aligned} \quad (5-86)$$

Practically, two examples were demonstrated. One considered the system

$$\begin{bmatrix} \dot{y}_1 \\ \dot{y}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -1 & -a \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} v \quad (5-87)$$

with $a > 0$ and

$$V(\sigma) = \sigma^3 , \quad (5-88)$$

with

$$\sigma = \hat{\underline{K}}^T \underline{y} = -(y_1 + \frac{a}{2} y_2) . \quad (5-89)$$

Then it can be shown that

$$\hat{R}(\sigma) = \frac{3a\hat{P}_{22}}{8} \mathcal{U}^{4/3} \quad (5-90)$$

and

$$\begin{aligned} \hat{L}(\underline{y}) = \frac{1}{2} \underline{y}^T & \begin{bmatrix} a\hat{P}_{22} & -\hat{P}_{11} + \hat{P}_{22} + \frac{a^2}{2} \hat{P}_{22} \\ -\hat{P}_{11} + \hat{P}_{22} + \frac{a^2}{2} \hat{P}_{22} & a\hat{P}_{22} \end{bmatrix} \underline{y} \\ & + \frac{a}{2} \hat{P}_{22} \left(\frac{y_1^4}{4} + \frac{2}{a} y_1^3 y_2 + \frac{6}{a} y_1^2 y_2^2 + \frac{8}{a^3} y_1 y_2^3 \right) + \frac{2\hat{P}_{22}}{a^3} y_2^4, \end{aligned} \quad (5-91)$$

with

$$\hat{\underline{P}} = \begin{bmatrix} \hat{P}_{11} & \frac{a}{2} \hat{P}_{22} \\ \frac{a}{2} \hat{P}_{22} & \hat{P}_{22} \end{bmatrix}. \quad (5-92)$$

Thau pointed out that the result, case (3), could be applied to the construction of a Liapunov function for a Lur'e system

$$\dot{\underline{y}} = \hat{\underline{A}} \underline{y} + \hat{\underline{B}} \mathcal{U}(\hat{\underline{K}}^T \underline{y}). \quad (5-93)$$

Considering (5-93) as a synthesized feedback control system for a control law $\mathcal{U}(\hat{\underline{K}}^T \underline{y})$ in the inverse problem,

$\hat{\underline{P}}$, $\hat{\underline{L}}(\underline{y})$ and $\hat{\underline{R}}(\underline{u})$ can be calculated according to the method of the inverse problem. If there exists a $\hat{\underline{P}}$, positive definite, such that

$$\hat{\underline{L}}(\underline{y}) + \hat{\underline{R}}(\underline{u}(\underline{K}^T \underline{y})) \quad (5-94)$$

is positive definite, then the origin of (5-93) is asymptotically stable in the large from Theorem 2-8, and a Liapunov function is

$$\hat{\underline{V}}^0(\underline{y}) = \underline{y}^T \hat{\underline{P}} \underline{y} . \quad (5-95)$$

In fact, from (5-70), $\dot{\hat{\underline{V}}}^0(\underline{y})$ is the negative of (5-94).

Summarizing Thau's work, the form of his loss function was less restrictive than that considered by Kalman and Suga. However, the nonnegative character of $\hat{\underline{L}}(\underline{y})$ was not discussed. The application of the results to construct Liapunov functions is a unique contribution, although the necessary assumptions are quite restrictive.

5.4 Comments

Based on Kalman's study of the inverse problem, Suga generalized the assumptions to include time varying multi-input linear systems (5-45), a performance index with any time interval, a time-varying loss function and a

time-varying linear control law. Thau further generalized the assumptions to include multi-input nonlinear systems (5-66), a broader class of loss functions and control laws. Although relevant characteristics of the optimal feedback control system are revealed by this work, the nonnegative property of either the loss function or the optimal performance index function was not adequately considered. Accordingly the results sometimes seem unrealistic from the viewpoint of optimal control theory, as illustrated by examples.

For the inverse problem considered in the next two chapters, generalizations of Kalman's assumptions are made with respect to the nonlinear, multi-input systems, the form of loss function and the nonlinear control law. Furthermore, nonnegativity of an optimized loss function and an optimized performance index function are considered. The canonical form developed in Chapter 4 contributes to the efficient analyses and the compact descriptions of the development.

Chapter 6

INVERSE PROBLEM OF THE OPTIMAL REGULATOR

In this chapter, the inverse problem of the optimal regulator is considered in a general context, i.e., for a class of multi-input systems with an unspecified non-linearity and feedback control law. Following a precise description of the problem in Section 6.1, an equivalent problem is defined in Section 6.2, using the canonical form developed in Chapter 4. Fundamental lemmas for the analysis are given in Section 6.3. Based on the Hamilton-Jacobi theory and Caratheodory's lemma, a principal theorem for the inverse problem is stated in Section 6.4. Section 6.5 has a discussion of the relevant aspects of this theorem to optimal feedback control systems and to work by other authors as special cases. Two practical examples of the application of the theorem are demonstrated in Section 6.6.

6.1 Statement of the Inverse Problem

The inverse problem of the optimal regulator is considered in this chapter under the assumptions such that

- (i) the system equation is given by

$$\dot{\underline{y}} = \hat{\underline{F}}(\underline{y}) + \hat{\underline{B}} \underline{v}, \quad (6-1)$$

where $\hat{\underline{F}}(\underline{Y})$ is an n -dimensional vector valued function of class C_2 satisfying $\hat{\underline{F}}(\underline{0}) = \underline{0}$ and $\hat{\underline{B}}$ is an $n \times m$ matrix such that

$$0 < \text{rank } \hat{\underline{B}} = r \leq m \leq n, \quad (6-2)$$

(ii) the desired final condition of the system is $\underline{Y} = \underline{0}$ in $R^n \times R^1$, with a stable feedback control law given by an m -dimensional vector valued function $\underline{V}(\underline{Y})$ of class C_2 with $\underline{V}(\underline{0}) = \underline{0}$ (thus the origin of the synthesized system

$$\dot{\underline{Y}} = \hat{\underline{F}}(\underline{Y}) + \hat{\underline{B}} \underline{V}(\underline{Y}) \quad (6-3)$$

is asymptotically stable in the large),

(iii) the form of performance index is restricted to

$$\int_{t_0}^{\infty} \{ \hat{\underline{L}}(\underline{Y}) + \underline{V}^T \hat{\underline{R}} \underline{V} \} dt, \quad (6-4)$$

where $\hat{\underline{R}}$ is an $m \times m$ symmetric positive definite matrix and $\hat{\underline{L}}(\underline{Y})$ is of class C_2 satisfying

$$\hat{\underline{L}}(\underline{0}) = 0. \quad (6-5)$$

Equivalently an initial condition (\underline{y}_0, t_0) and a suitable control function $\underline{v}(t)$ defined on $[t_0, \infty)$ are assumed such that

$$\int_{t_0}^{\infty} \{ \hat{L}(\hat{\phi}_{\underline{v}}(t; \underline{y}_0, t_0)) + \underline{v}^T(t) \hat{R} \underline{v}(t) \} dt$$

$$= \lim_{\tau \rightarrow \infty} \int_{t_0}^{\tau} \{ \hat{L}(\hat{\phi}_{\underline{v}}(t; \underline{y}_0, t_0)) + \underline{v}^T(t) \hat{R} \underline{v}(t) \} dt, \quad (6-6)$$

(iv) M_s and A.S.R. are the whole $R^n \times R^1$ and A.C.R. is the whole $R^m \times R^1$.

The inverse problem is to find performance indices (6-4) or, equivalently, loss functions that are optimized by the assumed control law, under assumptions (i)-(iv). This inverse problem is a generalization of the inverse problems considered by other authors, as reviewed in Chapter 5, i.e., the assumptions of the problem in Section 6.1 are less restrictive than those previously made. Specifically:

- (a) (6-1) is a nonlinear multi-input system with \hat{B} of a general rank, in comparison with (5-35) used by Kalman and (5-45) by Suga;
- (b) the control law (6-3) is unspecified, in comparison with (5-36) and (5-47) assumed by Kalman and Suga;
- (c) the performance index (6-4) has a general penalty function $\hat{L}(\underline{y})$, in comparison with the quadratic form

used by Kalman, (5-37).

Consequently, Kalman's problem is completely generalized in this chapter. However, all assumptions are limited to time invariancy in comparison to those of Suga. Also the quadratic form in \underline{V} in the performance index (6-4) is more restrictive than the $\hat{R}(\underline{V})$ used by Thau (5-67).

6.2 An Equivalent Inverse Problem

The analysis of the inverse problem can be facilitated by using the canonical form given in Chapter 4. Consider the canonical form of (6-1) as

$$\dot{\underline{X}} = \underline{F}(\underline{X}) + \underline{B} \underline{U} \quad (6-7)$$

for the transformation

$$\begin{aligned} \underline{X} &= \underline{N}^{-1} \underline{Y} \\ \underline{U} &= \underline{M}^{-1} \underline{V} , \end{aligned} \quad (6-8)$$

and specifically define

$$\underline{U}(\underline{X}) \triangleq \underline{M}^{-1} \underline{V}(\underline{N} \underline{X}) \quad (6-9)$$

as the feedback control law, equivalent to that given in (6-3).

The problem assumptions given in Section 6.1 are invariant with this transformation (6-8) and can be expressed in terms of the transformed variables. This is shown in the remainder of the section.

The rank of \underline{B} is r from (4-193) and

$$\underline{F}(\underline{0}) = \underline{N}^{-1} \hat{\underline{F}}(\underline{N} \underline{0}) = \underline{0}. \quad (6-10)$$

Also $\underline{F}(\underline{X})$ is of class C_2 , as each function $f_i(\underline{X})$,

$\frac{\partial f_i(\underline{X})}{\partial x_j}$ or $\frac{\partial^2 f_i(\underline{X})}{\partial x_k \partial x_j}$, $i, j, k = 1, 2, \dots, n$, is a linear combination of

$\hat{f}_i(\underline{N} \underline{X})$, $\frac{\partial \hat{f}_i(\underline{N} \underline{X})}{\partial x_j}$ and $\frac{\partial^2 \hat{f}_i(\underline{N} \underline{X})}{\partial x_k \partial x_j}$

respectively, and Theorem 2-1 can be applied. As $\underline{V}(\underline{0}) = \underline{0}$, from (6-8)

$$\underline{U}(\underline{0}) = \underline{M}^{-1} \underline{V}(\underline{N} \underline{0}) = \underline{0}, \quad (6-11)$$

and $\underline{U}(\underline{X})$ can be similarly shown to be of class C_2 .

To establish $\underline{U}(\underline{X})$ as a stable control law, a lemma is introduced.

Lemma 6-1:

If $\underline{V}(\underline{Y})$ is a stable control law for (6-1), then $\underline{U}(\underline{X})$ is a stable control law for the transformed system (6-7).

Proof: From Lemma 4-4, a solution for (6-7) with an initial condition (\underline{x}_0, t_0) is

$$\underline{\phi}_f(t; \underline{x}_0, t_0) = \underline{N}^{-1} \hat{\underline{\phi}}_f(t; \underline{N} \underline{x}_0, t_0), \quad (6-12)$$

where $\hat{\underline{\phi}}_f(t; \underline{N} \underline{x}_0, t_0)$ is a solution of (6-3) with an initial condition $(\underline{N} \underline{x}_0, t_0)$. Since the origin of (6-3) is asymptotically stable in the large by assumption, for arbitrary $\epsilon > 0$, $\mu > 0$ and t_0 , there exist a $\hat{\delta}(\epsilon, t_0) > 0$ and a $\hat{T}(\delta, \mu, t_0) > 0$ such that if

$$||\underline{y}_0|| \leq \hat{\delta}(\epsilon, t_0), \quad (6-13)$$

then

$$(a) \quad ||\hat{\underline{\phi}}_f(t; \underline{y}_0, t_0)|| \leq \epsilon, \text{ for all } t \geq t_0, \quad (6-14)$$

and

$$(b) \quad ||\hat{\underline{\phi}}_f(t; \underline{y}_0, t_0)|| \leq \mu, \text{ for all } t \geq t_0 + \hat{T}. \quad (6-15)$$

Define

$$\delta(\epsilon, t_0) \triangleq \frac{\hat{\delta}\left(\frac{\epsilon}{\lambda_{\max}((\underline{N}^{-1})^T \underline{N}^{-1})}, t_0\right)}{\lambda_{\max}(\underline{N}^T \underline{N})} \quad (6-16)$$

and

$$T(\delta, \mu, t_0) = \hat{T} \left(\delta, \frac{\mu}{\lambda_{\max}((\underline{N}^{-1})^T \underline{N}^{-1})}, t_0 \right), \quad (6-17)$$

where $\lambda_{\max}(\)$ is defined in Theorem 2-5 with $\underline{Q} = \underline{N}^T \underline{N}$.

Then for the $\epsilon > 0$, $\mu > 0$ and t_0 chosen, consider

$$||\underline{x}_0|| \leq \delta(\epsilon, t_0). \quad (6-18)$$

From (6-16), (6-18) and Theorem 2-5,

$$\begin{aligned} \hat{\delta} \left(\frac{\epsilon}{\lambda_{\max}((\underline{N}^{-1})^T \underline{N}^{-1})}, t_0 \right) &\geq \lambda_{\max}(\underline{N}^T \underline{N}) ||\underline{x}_0|| \\ &\geq ||\underline{N} \underline{x}_0|| = ||\underline{y}_0||. \end{aligned} \quad (6-19)$$

Therefore, from (6-13) and (6-19), (6-14) becomes

$$||\hat{\phi}_f(t; \underline{y}_0, t_0)|| \leq \frac{\epsilon}{\lambda_{\max}((\underline{N}^{-1})^T \underline{N}^{-1})}, \text{ for all } t \geq t_0, \quad (6-20)$$

or

$$\begin{aligned} \epsilon &\geq \lambda_{\max}((\underline{N}^{-1})^T \underline{N}^{-1}) ||\hat{\phi}_f(t; \underline{y}_0, t_0)|| \geq ||\underline{N}^{-1} \hat{\phi}_f(t; \underline{N} \underline{x}_0, t_0)|| \\ &= ||\phi_f(t; \underline{x}_0, t_0)||, \text{ for all } t \geq t_0, \end{aligned} \quad (6-21)$$

using Theorem 2-5 and Lemma 4-4. Also, from (6-15) and (6-17), it follows that

$$||\hat{\phi}_f(t; \underline{y}_0, t_0)|| \leq \frac{\mu}{\lambda_{\max}((\underline{N}^{-1})^T \underline{N}^{-1})}, \text{ for all } t \geq t_0 + T, \quad (6-22)$$

or

$$\mu \geq \lambda_{\max}((\underline{N}^{-1})^T \underline{N}^{-1}) ||\hat{\phi}_f(t; \underline{y}_0, t_0)|| \geq ||\phi_f(t; \underline{x}_0, t_0)||, \quad (6-23)$$

for all $t \geq t_0 + T$,

again using Theorem 2-5 and Lemma 4-4. By (6-18), (6-21) and (6-23), the lemma is proved.

The value of the performance index given by (6-5) for a suitable control $\underline{v}(t)$ and from an initial condition (\underline{y}_0, t_0) is

$$\hat{J}[\underline{y}_0, t_0, \underline{v}(t)] = \int_{t_0}^{\infty} \{ \hat{L}(\hat{\phi}_{\underline{v}}(t; \underline{y}_0, t_0)) + \underline{v}^T(t) \hat{R} \underline{v}(t) \} dt. \quad (6-24)$$

This reduces by Lemma 4-4 for (6-8), to

$$\hat{J}[\underline{N} \underline{x}_0, t_0, \underline{M} \underline{u}(t)] \triangleq \int_{t_0}^{\infty} \{ \hat{L}(\underline{N} \phi_{\underline{u}}(t; \underline{x}_0, t_0)) + \underline{u}^T(t) \underline{M}^T \hat{R} \underline{M} \underline{u}(t) \} dt. \quad (6-25)$$

Define

$$J[\underline{X}_0, t_0, \underline{U}(t)] \triangleq \hat{J}[\underline{N} \underline{X}_0, t_0, \underline{M} \underline{U}(t)]$$

$$L(\underline{X}) \triangleq \hat{L}(\underline{N} \underline{X}) \quad (6-26)$$

$$\underline{R} \triangleq \underline{M}^T \hat{\underline{R}} \underline{M}.$$

$L(\underline{X})$ is of class C_2 because $\hat{L}(\underline{Y})$ is of class C_2 and Theorem 2-1. Also $L(0) = \hat{L}(\underline{N} 0) = 0$, and \underline{R} is symmetric and positive definite, from part (iii) of Theorem 2-4. Then the performance index (6-4) is equivalently written as

$$J[\underline{X}_0, t_0, \underline{U}(t)] = \int_{t_0}^{\infty} \{L(\underline{X}) + \underline{U}^T \underline{R} \underline{U}\} dt \quad (6-27)$$

for the canonical form (6-7), and assumption (iii) for $\hat{L}(\underline{Y})$ and $\hat{\underline{R}}$ is completely preserved for $L(\underline{X})$ and \underline{R} .

Since M_s , A.S.R. and A.C.R. in the \underline{X} and \underline{U} coordinates are whole $R^n \times R^1$ and $R^m \times R^1$ because of the bijective mappings of (6-8), assumption (iv) is also preserved. Thus, the original inverse problem stated in Section 6.1 can be considered in the canonical form (6-7) under the same mathematical assumptions without loss of generality. The recovery of the solution for the original system follows from the inverse transformations of (6-8) and (6-26).

For convenience of analysis, the following notation is used for (6-7), (6-9) and (6-27).

(1) \underline{x} is decomposed into

$$\begin{aligned}\underline{x}_1 &\triangleq [x_1, x_2, \dots, x_{n-r}]^T \\ \underline{x}_2 &\triangleq [x_{n-r+1}, \dots, x_n]^T,\end{aligned}\tag{6-28}$$

and, if $n = r$, $\underline{x} = \underline{x}_2$.

(2) \underline{u} is decomposed into

$$\begin{aligned}\underline{u}_d &\triangleq [u_1, u_2, \dots, u_{m-r}]^T \\ \underline{u}_e &\triangleq [u_{m-r+1}, \dots, u_m]^T,\end{aligned}\tag{6-29}$$

and, if $m = r$, $\underline{u} = \underline{u}_e$.

(3) $\underline{F}(\underline{x})$ is decomposed into

$$\begin{aligned}\underline{F}_1(\underline{x}) &\triangleq [f_1(\underline{x}), f_2(\underline{x}), \dots, f_{n-r}(\underline{x})]^T \\ \underline{F}_2(\underline{x}) &\triangleq [f_{n-r+1}(\underline{x}), \dots, f_n(\underline{x})]^T,\end{aligned}\tag{6-30}$$

with

$$\underline{F}(\underline{x}) = \underline{A} \underline{x} + \underline{F}'(\underline{x}),\tag{6-31}$$

where $\underline{A} \underline{X}$ defines the first degree homogeneous function in $\underline{F}(\underline{X})$, and $\underline{F}(\underline{X})$ the remainder. \underline{A} and $\underline{F}(\underline{X})$ are further defined as

$$\underline{A} \triangleq \begin{bmatrix} \underline{A}_1 \\ \underline{A}_2 \end{bmatrix} \triangleq \left\{ \begin{bmatrix} \underline{A}_{11} & \underline{A}_{12} \\ \underline{A}_{21} & \underline{A}_{22} \end{bmatrix} \right\} \begin{matrix} n-r \\ r \end{matrix} \quad (6-32)$$

$\underbrace{\hspace{1.5cm}}_{n-r} \quad \underbrace{\hspace{1.5cm}}_r$

and

$$\underline{F}(\underline{X}) = \left\{ \begin{bmatrix} \underline{F}_1(\underline{X}) \\ \underline{F}_2(\underline{X}) \end{bmatrix} \right\} \begin{matrix} n-r \\ r \end{matrix} \quad (6-33)$$

Subsequently (6-7) can be written as

$$\dot{\underline{X}}_1 = \underline{F}_1(\underline{X}) \quad (6-34)$$

$$\dot{\underline{X}}_2 = \underline{F}_2(\underline{X}) + \underline{U}_e$$

from (6-30), or

$$\dot{\underline{X}}_1 = \underline{A}_1 \underline{X} + \underline{F}_1(\underline{X}) \quad (6-35)$$

$$\dot{\underline{X}}_2 = \underline{A}_2 \underline{X} + \underline{F}_2(\underline{X}) + \underline{U}_e$$

from (6-32) and (6-33).

(4) Let

$$\underline{R} \triangleq \left[\begin{array}{cc} \underbrace{\underline{R}_{11}}_{m-r} & \underbrace{\underline{R}_{12}}_r \\ \underbrace{\underline{R}_{12}^T}_{m-r} & \underbrace{\underline{R}_{22}}_r \end{array} \right] \begin{array}{l} \} m-r \\ \} r \end{array} \quad (6-36)$$

(5) If $\nabla^0(\underline{x})$ is a scalar function of class C_2 , then define

$$\begin{aligned} \frac{\partial \nabla^0}{\partial \underline{x}_1} &\triangleq [w_1(\underline{x}), w_2(\underline{x}), \dots, w_{n-r}(\underline{x})]^T \\ \frac{\partial \nabla^0}{\partial \underline{x}_2} &\triangleq [w_{n-r+1}(\underline{x}), \dots, w_n(\underline{x})]^T \end{aligned} \quad (6-37)$$

and

$$\frac{\partial \nabla^0}{\partial \underline{x} \partial \underline{x}} = \left[\begin{array}{cc} \frac{\partial^2 \nabla^0}{\partial \underline{x}_1 \partial \underline{x}_1} & \frac{\partial^2 \nabla^0}{\partial \underline{x}_2 \partial \underline{x}_1} \\ \frac{\partial^2 \nabla^0}{\partial \underline{x}_1 \partial \underline{x}_2} & \frac{\partial^2 \nabla^0}{\partial \underline{x}_2 \partial \underline{x}_2} \end{array} \right],$$

$$= \begin{bmatrix} \frac{\partial w_1}{\partial x_1} & \frac{\partial w_2}{\partial x_2} & \cdot & \cdot & \cdot & \frac{\partial w_1}{\partial x_n} \\ \frac{\partial w_2}{\partial x_1} & & & & & \cdot \\ \cdot & & & & & \cdot \\ \cdot & & & & & \cdot \\ \cdot & & & & & \cdot \\ \frac{\partial w_n}{\partial x_1} & \cdot & \cdot & \cdot & \cdot & \frac{\partial w_n}{\partial x_n} \end{bmatrix} \cdot \quad (6-38)$$

6.3 Fundamental Lemmas

Lemma 6-2:

Let m and r be integers such that $m > r > 0$ and \underline{R} $m \times m$ symmetric matrix. Defining \underline{R} as (6-36), \underline{R} is positive definite if and only if

(i) \underline{R}_{11} and $\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}$ are positive definite,

or

(ii) \underline{R}_{22} and $\underline{R}_{11} - \underline{R}_{12} \underline{R}_{22}^{-1} \underline{R}_{12}^T$ are positive definite.

Proof: Only (i) is proved; (ii) is proved similarly.

For \underline{R} to be positive definite, \underline{R}_{11} must be positive definite by (i) of Corollary 2-1. Thus \underline{R}_{11}^{-1} exists. Define a matrix

$$\underline{S} \triangleq \begin{bmatrix} \underline{I}_{m-r} & -\underline{R}_{11}^{-1} \underline{R}_{12} \\ [0] & \underline{I}_r \end{bmatrix} \quad (6-39)$$

which is nonsingular. According to (iii) of Theorem 2-4, \underline{R} is positive definite if and only if

$$\underline{S}^T \underline{R} \underline{S} = \begin{bmatrix} \underline{R}_{11} & [0] \\ [0] & \underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12} \end{bmatrix} \quad (6-40)$$

is positive definite. Since the characteristic equation of (6-40) is

$$0 = |\lambda \underline{I} - \underline{S}^T \underline{R} \underline{S}| = |\lambda \underline{I} - \underline{R}_{11}| \cdot |\lambda \underline{I} - (\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12})|, \quad (6-41)$$

eigenvalues of (6-40) are those of \underline{R}_{11} and of $(\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12})$. From Theorem 2-4, the assertion is immediate.

Lemma 6-3:

For \underline{R} given as in (6-36) and for the rank of \underline{R}_{22} as $r_1 \leq r$, \underline{R} is positive semidefinite if and only if

- (i) \underline{R}_{22} is positive semidefinite,
- (ii) $\underline{\Omega} \triangleq \underline{R}_{11} - \underline{R}_{12} \underline{D}_{22}^T \underline{D}_{22} \underline{R}_{22} \underline{D}_{22}^T \underline{D}_{22} \underline{R}_{12}^T$ is positive semidefinite,

where \underline{D}_{22} is an $r \times r$ nonsingular congruent transformation matrix such that

$$\underline{D}_{22}^T \underline{R}_{22} \underline{D}_{22} = \begin{bmatrix} \underline{I}_{r_1} & [0] \\ [0] & [0] \end{bmatrix}, \quad (6-42)$$

(the existence of \underline{D}_{22} follows from Theorem 2-4 and (i) above),

and

(iii) the last $r - r_1$ columns of $\underline{R}_{12} \underline{D}_{22}$ are null, i.e.,

$$\underline{R}_{12} \underline{D}_{22} \triangleq \left\{ \begin{bmatrix} \underline{\Sigma} & \underline{\mathcal{R}} \end{bmatrix} \right\}_{m-r}, \text{ with } \underline{\mathcal{R}} = [0]. \quad (6-43)$$

$\underbrace{\quad}_{r_1} \quad \underbrace{\quad}_{r-r_1}$

Proof: Assume an m -dimensional vector $\underline{z} \triangleq [z_1, z_2, \dots, z_m]^T$ and

$$\begin{aligned} \underline{z}_1 &= [z_1, z_2, \dots, z_{m-r}]^T \\ \underline{z}_2 &= [z_{m-r+1}, z_{m-r+2}, \dots, z_{m-r+r_1}]^T \\ \underline{z}_3 &= [z_{m-r+r_1+1}, z_{m-r+r_1+2}, \dots, z_m]^T. \end{aligned} \quad (6-44)$$

Then \underline{R} is positive semidefinite if and only if the quadratic form $\underline{z}^T \underline{R} \underline{z}$ is positive semidefinite. Fix $\underline{z}_1 = 0$ and then

$$\underline{z}^T \underline{R} \underline{z} = [\underline{z}_2^T \quad \underline{z}_3^T] \underline{R}_{22} \begin{bmatrix} \underline{z}_2 \\ \underline{z}_3 \end{bmatrix} \quad (6-45)$$

and (i) is necessary. For (iii), define

$$\underline{D} \triangleq \begin{bmatrix} \underline{I}_{m-r} & [0] \\ -\underline{\Delta} & \underline{D}_{22} \end{bmatrix} \quad (6-46)$$

with

$$\underline{\Delta} \triangleq \underline{D}_{22} \underline{D}_{22}^T \underline{R}_{22} \underline{D}_{22} \underline{D}_{22}^T \underline{R}_{12}^T. \quad (6-47)$$

\underline{D} is nonsingular because of the nonsingular \underline{D}_{22} given by (6-42). From (iii) of Theorem 2-4, \underline{R} is positive semidefinite if and only if

$$\underline{D}^T \underline{R} \underline{D} = \begin{bmatrix} \underline{R}_{11} - \underline{\Delta}^T \underline{R}_{12} \underline{\Delta} & \underline{R}_{12} \underline{D}_{22} - \underline{\Delta}^T \underline{R}_{22} \underline{D}_{22} \\ + \underline{\Delta}^T \underline{R}_{22} \underline{\Delta} & \\ \underline{D}_{22}^T \underline{R}_{12} - \underline{D}_{22}^T \underline{R}_{22} \underline{\Delta} & \underline{D}_{22}^T \underline{R}_{22} \underline{D}_{22} \end{bmatrix} \quad (6-48)$$

is positive semidefinite. But

$$\begin{aligned} \underline{R}_{11} - \underline{\Delta}^T \underline{R}_{12}^T - \underline{R}_{12} \underline{\Delta} + \underline{\Delta}^T \underline{R}_{22} \underline{\Delta} &= \underline{R}_{11} - \underline{R}_{12} \underline{D}_{22} \underline{D}_{22}^T \underline{R}_{22} \underline{D}_{22} \underline{D}_{22}^T \underline{R}_{12}^T \\ &= \underline{\Omega} \end{aligned} \quad (6-49)$$

and, from (6-43),

$$\underline{R}_{12} \underline{D}_{22} - \underline{\Delta}^T \underline{R}_{22} \underline{D}_{22} = \underline{R}_{12} \underline{D}_{22} \begin{bmatrix} [0] & [0] \\ [0] & \underline{I}_{r-r_1} \end{bmatrix} = \begin{bmatrix} [0] & \underline{R} \end{bmatrix}, \quad (6-50)$$

where \underline{R} is the last $(r - r_1)$ columns of $\underline{R}_{12} \underline{D}_{22}$. Substituting (6-47), (6-49) and (6-50) into (6-48),

$$\underline{D}^T \underline{R} \underline{D} = \left[\begin{array}{ccc} \underline{\Omega} & [0] & \underline{R} \\ [0] & \underline{I}_{r_1} & [0] \\ \underbrace{\underline{R}^T}_{m-r} & \underbrace{[0]}_{r_1} & \underbrace{[0]}_{r-r_1} \end{array} \right] \begin{array}{l} \} \quad m-r \\ \} \quad r_1 \\ \} \quad r-r_1 \end{array} \quad (6-51)$$

In order for $\underline{D}^T \underline{R} \underline{D}$ to be positive semidefinite, \underline{R} must be the null matrix. If it is not, consider $\underline{z}_2 = \underline{0}$ and \underline{z} satisfying $\underline{z}_1 \underline{R} \neq \underline{0}$. Then

$$\underline{z}^T \underline{D}^T \underline{R} \underline{D} \underline{z} = \underline{z}_1^T \underline{\Omega} \underline{z}_1^T + 2 \underline{z}_1^T \underline{R} \underline{z}_3. \quad (6-52)$$

As \underline{z}_3 can be chosen to provide a negative value for (6-52), the positive semidefiniteness of $\underline{D}^T \underline{R} \underline{D}$ fails. Therefore $\underline{R} = [0]$ and (iii) is necessary. Alternately, fixing $\underline{z}_2 = \underline{0}$ and $\underline{z}_3 = \underline{0}$,

$$\underline{z}^T \underline{D}^T \underline{R} \underline{D} \underline{z} = \underline{z}_1^T \underline{\Omega} \underline{z}_1, \quad (6-53)$$

and (ii) is required.

Conversely if (i)-(iii) are satisfied,

$$\underline{z}^T \underline{D}^T \underline{R} \underline{D} \underline{z} = \underline{z}_1^T \underline{\Omega} \underline{z}_1 + \underline{z}_2^T \underline{z}_2, \quad (6-54)$$

and $\underline{D}^T \underline{R} \underline{D}$ are positive semidefinite as $\underline{\Omega}$ is positive semidefinite; by Theorem 2-4, \underline{R} is then positive semidefinite.

For the loss function in (6-27), define

$$V(\tau; \underline{x}_0, t_0) = \int_{t_0}^{\tau} \{ L(\underline{\phi}_f(t; \underline{x}_0, t_0)) + \underline{U}^T(\underline{\phi}_f(t; \underline{x}_0, t_0)) \underline{R} \underline{U}(\underline{\phi}_f(t; \underline{x}_0, t_0)) \} dt \quad (6-55)$$

and

$$V(\underline{x}_0, t_0) = \lim_{\tau \rightarrow \infty} V(\tau; \underline{x}_0, t_0), \quad (6-56)$$

where $\underline{\phi}_f(t; \underline{x}_0, t_0)$ is the solution of

$$\dot{\underline{X}} = \underline{F}(\underline{X}) + \underline{B} \underline{U}(\underline{X}) \quad (6-57)$$

from (\underline{X}_0, t_0) .

Lemma 6-4:

Provided that $\nabla(\underline{X}, t)$, (6-56), is well defined on $R^n \times R^1$, for an arbitrary $(\underline{X}_0, t_0) \in R^n \times R^1$ and $\alpha > 0$,

$$\nabla(\underline{X}_0, t_0) = \nabla(\underline{X}_0, t_0 + \alpha) , \quad (6-58)$$

that is, $\nabla(\underline{X}, t)$ is independent of t .

Proof: From (6-56), the statement is justified by proving that for each (\underline{X}_0, t_0) and $\alpha > 0$,

$$\lim_{\tau \rightarrow \infty} \{\nabla(\tau; \underline{X}_0, t_0) - \nabla(\tau; \underline{X}_0, t_0 + \alpha)\} = 0 , \quad (6-59)$$

or, for an arbitrary $\epsilon > 0$, that there exists a $T > 0$ such that

$$|\nabla(\tau; \underline{X}_0, t_0) - \nabla(\tau; \underline{X}_0, t_0 + \alpha)| \leq \epsilon, \quad \text{for all } \tau > T . \quad (6-60)$$

As (6-57) is autonomous, it is known that^[13]

$$\underline{\phi}_f(t + \alpha, \underline{X}_0, t_0 + \alpha) = \underline{\phi}_f(t, \underline{X}_0, t_0) . \quad (6-61)$$

Then, for arbitrary $\tau > t_0 + \alpha$, it follows from (6-55) that

$$\begin{aligned} V(\tau; \underline{X}_0, t_0) &= \int_{t_0 + \alpha}^{\tau + \alpha} \{ L(\underline{\phi}_f(t; \underline{X}_0, t_0 + \alpha) \\ &\quad + \underline{U}^T(\underline{\phi}_f(t; \underline{X}_0, t_0 + \alpha)) \underline{R} \underline{U}(\underline{\phi}_f(t; \underline{X}_0, t_0 + \alpha)) \} dt \\ &= V(\tau; \underline{X}_0, t_0 + \alpha) + \int_{\tau}^{\tau + \alpha} \{ L(\underline{\phi}_f(t; \underline{X}_0, t_0 + \alpha) \\ &\quad + \underline{U}^T(\underline{\phi}_f(t; \underline{X}_0, t_0 + \alpha)) \underline{R} \underline{U}(\underline{\phi}_f(t; \underline{X}_0, t_0 + \alpha)) \} dt . \end{aligned} \quad (6-62)$$

As $L(\underline{X})$ and $\underline{U}(\underline{X})$ are continuous and zero at $\underline{X} = \underline{0}$ by the assumptions, i.e., (6-5) and (6-11), there exist $\mu_1 > 0$ and $\mu_2 > 0$ satisfying

$$|L(\underline{X})| < \frac{\epsilon}{2\alpha} , \text{ if } ||\underline{X}|| < \mu_1 , \quad (6-63)$$

and

$$||\underline{U}(\underline{X})|| < \sqrt{\frac{\epsilon}{2\alpha\lambda_{\max}(\underline{R})}} , \text{ if } ||\underline{X}|| < \mu_2 , \quad (6-64)$$

where $\lambda_{\max}(\underline{R})$ is given in Theorem 2-5. Since the origin of (6-59) is assumed asymptotically stable in the large, for the given (\underline{X}_0, t_0) , there exists a $T_1 > 0$ such that

$$||\underline{\phi}_f(t; \underline{X}_0, t_0 + \alpha)|| \leq \text{Min}(\mu_1, \mu_2), \text{ for all } t \geq t_0 + T_1. \quad (6-65)$$

From (6-62~65), it follows that

$$\begin{aligned} & |V(\tau; \underline{X}_0, t_0) - V(\tau; \underline{X}_0, t_0 + \alpha)| = \\ & \leq \left| \int_{\tau}^{\tau + \alpha} L(\underline{\phi}_f(t; \underline{X}_0, t_0 + \alpha)) dt \right| \\ & + \left| \int_{\tau}^{\tau + \alpha} \underline{U}^T(\underline{\phi}_f(t; \underline{X}_0, t_0 + \alpha)) \underline{R} \underline{U}(\underline{\phi}_f(t; \underline{X}_0, t_0 + \alpha)) dt \right| \\ & \leq \frac{\epsilon}{2\alpha} \cdot \alpha + \frac{\epsilon}{2\alpha} \cdot \alpha = \epsilon, \text{ for all } \tau \geq t_0 + \alpha + T_1. \end{aligned} \quad (6-66)$$

Thus (6-60) follows if

$$T \triangleq t_0 + \alpha + T_1. \quad (6-67)$$

Accordingly, $V(\underline{X}, t)$ in (6-56) can be simply described as $V(\underline{X})$.

Consider a class of nonlinear systems given by

$$\dot{\underline{X}} = \underline{A} \underline{X} + \underline{B} \underline{U} + \underline{B}_e \underline{F}_2(\underline{X}) \quad , \quad (6-68)$$

that is, $\underline{F}_1(\underline{X})$ in (6-35) is identically zero and \underline{B}_e given by (4-114). The controllability of this particular class of systems can be established by the application of Lemma 6-5, based on Theorem 2-9.

Lemma 6-5:

A system (6-68) is completely controllable if and only if the system

$$\dot{\underline{X}} = \underline{A} \underline{X} + \underline{B} \underline{U} \quad (6-69)$$

is completely controllable.

Proof: Assume that (6-68) is completely controllable.

Then, from Definition 2-11, for arbitrary $\underline{X}_0, \underline{X}_1 \in \mathbb{R}^n$ and $t_0 \in \mathbb{R}^1$, there exists a control function $\underline{U}(t)$ to provide a solution

$$\underline{\phi}_{\underline{U}}(t; \underline{X}_0, t_0) \quad (6-70)$$

satisfying

$$\begin{aligned} \dot{\underline{\phi}}_{\underline{u}}(t; \underline{x}_0, t_0) &= \underline{A} \underline{\phi}_{\underline{u}}(t; \underline{x}_0, t_0) \\ &+ \underline{B} \underline{u}(t) + \underline{B}_e \underline{F}_2(\underline{\phi}_{\underline{u}}(t; \underline{x}_0, t_0)) \end{aligned} \quad (6-71)$$

and at some $t_1 \geq t_0$

$$\underline{\phi}_{\underline{u}}(t_1; \underline{x}_0, t_0) = \underline{x}_1. \quad (6-72)$$

Consider a control function

$$\underline{u}(t) + \begin{bmatrix} \underline{0} \\ \underline{F}_2(\underline{\phi}_{\underline{u}}(t; \underline{x}_0, t_0)) \end{bmatrix} \quad (6-73)$$

for (6-69), i.e.,

$$\dot{\underline{x}} = \underline{A} \underline{x} + \underline{B} \underline{u}(t) + \underline{B}_e \underline{F}_2(\underline{\phi}_{\underline{u}}(t; \underline{x}_0, t_0)), \quad (6-74)$$

where $\underline{0}$ is $(m - r)$ dimensional. Then the solution from (\underline{x}_0, t_0) is (6-70) and (6-69) is completely controllable if (6-68) is completely controllable.

The same arguments show that (6-68) is completely controllable if (6-69) is completely controllable.

Thus, according to Theorems 2-9 and 4-3, the controllability of (6-68) can be simply identified by the structure of $\underline{A}_{(1,1)}$ and $\underline{A}_{(1,2)}$ through the application of Corollary 4-1.

Lemma 6-6:

Consider \underline{A}_1 as defined in (6-32) with \underline{A} given by (4-9) for $v \geq 2$, and assume an arbitrary scalar function $\varphi(\underline{x}_1)$ of class C_2 and $\varphi(\underline{0}) = 0$ with \underline{x}_1 defined by (6-28). Then for

$$- \left[\frac{\partial \varphi(\underline{x}_1)}{\partial \underline{x}_1} \right]^T \underline{A}_1 \underline{x} \quad (6-75)$$

to be positive semidefinite in \mathbb{R}^n , $\varphi(\underline{x}_1)$ must be

(i) identically zero if $\underline{A}_1 \underline{x}$ is from a completely controllable system,

(ii) a function of only $\underline{x}_{(1)}$ defined by (4-116) if $\underline{A}_1 \underline{x}$ is from an uncontrollable system.

Proof: Define

$$\hat{\underline{w}}(\underline{x}_1) \triangleq \left[\frac{\partial \varphi(\underline{x}_1)}{\partial \underline{x}_1} \right]^T. \quad (6-76)$$

From Theorem 2-2,

$$\frac{\partial}{\partial \underline{x}_1} \{ \hat{\underline{w}}(\underline{x}_1) \} \quad (6-77)$$

must be symmetric. Then, using (4-9),

$$-\hat{\underline{w}}^T(\underline{x}_1)\underline{A}_1\underline{x} = -[\hat{\underline{w}}_{(1)}^T(\underline{x}_1), \hat{\underline{w}}_{(2)}^T(\underline{x}_1), \dots, \hat{\underline{w}}_{(v-1)}^T(\underline{x}_1)]$$

$$\begin{bmatrix} \underline{A}_{(1,1)} & \underline{A}_{(1,2)} & & & \\ & & \underline{A}_{(2,3)} & & \\ & & & \ddots & \\ & & & & \underline{A}_{(v-1,v)} \end{bmatrix} \begin{bmatrix} \underline{x}_{(1)} \\ \underline{x}_{(2)} \\ \vdots \\ \underline{x}_{(v)} \end{bmatrix}$$

$$\begin{aligned} &= -\{\hat{\underline{w}}_{(1)}^T(\underline{x}_1)\underline{A}_{(1,1)}\underline{x}_{(1)} + \hat{\underline{w}}_{(1)}^T(\underline{x}_1)\underline{A}_{(1,2)}\underline{x}_{(2)} \\ &+ \hat{\underline{w}}_{(3)}^T(\underline{x}_1)\underline{A}_{(2,3)}\underline{x}_{(3)} + \dots + \hat{\underline{w}}_{(v-1)}^T(\underline{x}_1)\underline{A}_{(v-1,v)}\underline{x}_{(v)}\} , \end{aligned}$$

(6-78)

where

$$\hat{\underline{w}}_{(i)}(\underline{x}_1) = [\hat{w}_{p_i+1}(\underline{x}_1), \hat{w}_{p_i+2}(\underline{x}_1), \dots, \hat{w}_{p_i+\ell_i}(\underline{x}_1)]^T \quad (6-79)$$

with

$$p_i = \sum_{j=1}^{i-1} \ell_j . \quad (6-80)$$

(1) Assume $v = 2$. Then it follows, from (6-28) and (4-116), that

$$\underline{x}_1 = \underline{x}_{(1)}$$

(6-81)

$$\underline{x}_2 = \underline{x}_{(2)},$$

If the system is completely controllable, $A_{(1,1)} = [0]$ from Corollary 4-1, and (6-78) becomes

$$-\hat{\underline{w}}_{(1)}^T(\underline{x}_{(1)}) A_{(1,2)} \underline{x}_{(2)} = -\hat{\underline{w}}_{(1)}^T(\underline{x}_{(1)}) \begin{bmatrix} [0] & I_{\ell_1} \end{bmatrix} \underline{x}_{(2)} .$$

(6-82)

Therefore, if $\hat{\underline{w}}_{(1)}(\underline{x}_{(1)})$ is not identically zero, $\underline{x}_{(1)}$ can be selected to provide a nonzero value to $\hat{\underline{w}}_{(1)}(\underline{x}_{(1)})$. Then (6-82) becomes a linear function of $\underline{x}_{(2)}$ with non-zero coefficients and can have negative values for a proper choice of $\underline{x}_{(2)}$. Therefore $\hat{\underline{w}}_{(1)}(\underline{x}_{(1)})$ must be identically zero. As $\varphi(0) = 0$, then $\varphi(\underline{x}_{(1)}) = 0$.

(2) Assume a general case of $v > 2$. Then for (6-78) to be positive semidefinite it similarly follows that $\hat{\underline{w}}_{(v-1)}(\underline{x}_{(1)})$ must be identically zero. According to the symmetry in (6-76), it must follow that

$$\left[\frac{\partial \hat{\underline{w}}_{(i)}(\underline{x}_{(1)})}{\partial \underline{x}_{(v-1)}} \right] = \left[\frac{\partial \hat{\underline{w}}_{(v-1)}(\underline{x}_{(1)})}{\partial \underline{x}_{(i)}} \right]^T = [0], \text{ for } i = 1, 2, \dots, v-1 .$$

(6-83)

Therefore, $\hat{\underline{w}}(\underline{x}_1)$ cannot be a function of $\underline{x}_{(v-1)}$ for the

positive semidefiniteness of (6-78).

A similar process can be repeated for each subvector $\underline{x}_{(v-1)}, \underline{x}_{(v-2)}, \dots, \underline{x}_{(1)}$ succeedingly. Then it is concluded that, for (6-78) to be positive semidefinite, $\hat{\underline{w}}_{(i)}(\underline{x}_1)$, $i = v-2, v-3, \dots, 2$, must be identically zero and $\hat{\underline{w}}_{(1)}(\underline{x}_1)$ must be a function of only $\underline{x}_{(1)}$, say $\hat{\underline{w}}_{(1)}(\underline{x}_{(1)})$. If the system is completely controllable, (6-78) reduces to

$$- \hat{\underline{w}}_{(1)}^T(\underline{x}_{(1)}) \underline{A}_{(1,2)} \underline{x}_{(2)} \quad (6-84)$$

Then applying the results of (1), $\hat{\underline{w}}_{(1)}(\underline{x}_{(1)})$ must be identically zero for the positive semidefiniteness of (6-75).

If the system is uncontrollable, the same argument follows, except, from Corollary 4-1, (6-78) reduces to

$$- \hat{\underline{w}}_{(1)}^T(\underline{x}_{(1)}) \underline{A}_{(1,1)} \underline{x}_{(1)} \quad (6-85)$$

6.4 Analysis of the Inverse Problem

6.4.1 Hamilton-Jacobi Equation

Assume that (a) a specific loss function is given as

$$L(\underline{x}) + \underline{u}^T \underline{R} \underline{u} \quad (6-86)$$

for the inverse problem and (b) the resulting optimal performance index function $\nabla^0(\underline{x})$ is of class C_2 . The Hamilton-Jacobi equation becomes from (5-29) and (6-7),

$$0 = \min_{\underline{u}} \left\{ \left[\frac{\partial \nabla^0}{\partial \underline{x}} \right]^T \underline{B} \underline{u} + \underline{u}^T \underline{R} \underline{u} \right\} + \left[\frac{\partial \nabla^0}{\partial \underline{x}} \right]^T \underline{F}(\underline{x}) + L(\underline{x}) . \quad (6-87)$$

As \underline{R} is positive definite by the assumption, the Hamiltonian is normal from Definition 5-1. Thus, its minimum at each $\underline{x} \in \mathbb{R}^n$ is uniquely provided by \underline{u} satisfying

$$\begin{aligned} 0 &= \frac{\partial}{\partial \underline{u}} \left\{ \left[\frac{\partial \nabla^0}{\partial \underline{x}} \right]^T \underline{B} \underline{u} + \underline{u}^T \underline{R} \underline{u} \right\} \\ &= \left[\frac{\partial \nabla^0}{\partial \underline{x}} \right]^T \underline{B} + 2 \underline{u}^T \underline{R} , \end{aligned} \quad (6-88)$$

from (5-33). Identifying this \underline{u} in a closed form as a function of \underline{x} , the Hamilton-Jacobi equation is realized by a feedback control law such that

$$\underline{u}(\underline{x}) = - \frac{1}{2} \underline{R}^{-1} \underline{B}^T \left[\frac{\partial \nabla^0}{\partial \underline{x}} \right] . \quad (6-89)$$

Define

$$\underline{R}^{-1} \triangleq \underline{\tilde{R}} \triangleq \begin{bmatrix} \underline{\tilde{R}}_{11} & \underline{\tilde{R}}_{12} \\ \underline{\tilde{R}}_{12}^T & \underline{\tilde{R}}_{22} \end{bmatrix}, \quad (6-90)$$

where the dimensions of each $\underline{\tilde{R}}_{ij}$ are the same to those of \underline{R}_{ij} in (6-36). Then it follows, [32] that

$$\begin{aligned} \underline{\tilde{R}}_{11} &= [\underline{R}_{11} - \underline{R}_{12}^T \underline{R}_{22}^{-1} \underline{R}_{12}]^{-1}, \\ \underline{\tilde{R}}_{12} &= -\underline{R}_{11}^{-1} \underline{R}_{12} [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}]^{-1}, \\ \underline{\tilde{R}}_{22} &= [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}]^{-1}. \end{aligned} \quad (6-91)$$

For \underline{B} given by (4-99), (6-89) can be reduced to

$$\begin{aligned} \underline{U}_d(\underline{x}) &= -\frac{1}{2} \underline{\tilde{R}}_{12} \left[\frac{\partial \nabla^0}{\partial \underline{x}_2} \right] \\ \underline{U}_e(\underline{x}) &= -\frac{1}{2} \underline{\tilde{R}}_{22} \left[\frac{\partial \nabla^0}{\partial \underline{x}_2} \right]. \end{aligned} \quad (6-92)$$

Since $\underline{\tilde{R}}_{22}$ must be nonsingular to insure the positive definiteness of \underline{R} from Lemma 6-2, it follows from (6-91) that

$$\left[\frac{\partial \nabla^0}{\partial \underline{x}_2} \right] = -2 \underline{\tilde{R}}_{22}^{-1} \underline{U}_e(\underline{x}) = -2 [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{U}_e(\underline{x}) \quad (6-93)$$

and

$$\underline{U}_d(\underline{X}) = \underline{\tilde{R}}_{12} \underline{\tilde{R}}_{22}^{-1} \underline{U}_e(\underline{X}) = -\underline{R}_{11}^{-1} \underline{R}_{12} \underline{U}_e(\underline{X}) . \quad (6-94)$$

Substituting (4-91), (6-30), (6-92), (6-93) and (6-94) into (6-89),

$$\begin{aligned} 0 = & \left[0, 0, \dots, 0, \left[\frac{\partial \nabla^0}{\partial \underline{X}_2} \right] \right] \underline{U}(\underline{X}) + \underline{U}^T(\underline{X}) \begin{bmatrix} \underline{R}_{11} & \underline{R}_{12} \\ \underline{R}_{12}^T & \underline{R}_{22} \end{bmatrix} \underline{U}(\underline{X}) \\ & + \begin{bmatrix} \left[\frac{\partial \nabla^0}{\partial \underline{X}_1} \right]^T & \left[\frac{\partial \nabla^0}{\partial \underline{X}_2} \right]^T \end{bmatrix} \begin{bmatrix} \underline{F}_1(\underline{X}) \\ \underline{F}_2(\underline{X}) \end{bmatrix} + L(\underline{X}) \end{aligned} \quad (6-95)$$

or

$$\begin{aligned} L(\underline{X}) = & \underline{U}_e^T(\underline{X}) [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{U}_e(\underline{X}) + 2 \underline{U}_e^T(\underline{X}) [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{F}_2(\underline{X}) \\ & - \left[\frac{\partial \nabla^0}{\partial \underline{X}_1} \right]^T \underline{F}_1(\underline{X}) . \end{aligned} \quad (6-96)$$

6.4.2 Concerning $\nabla^0(\underline{X})$

Since $\nabla^0(\underline{X})$ is a scalar of class C_2 , the functional matrix $\frac{\partial^2 \nabla^0(\underline{X})}{\partial \underline{X} \partial \underline{X}}$ is symmetric (by Theorem 2-2). From (6-38),

$$w_{ij}(\underline{x}) \triangleq \frac{\partial w_i(\underline{x})}{\partial x_j} = \frac{\partial w_j(\underline{x})}{\partial x_i} \triangleq w_{ji}(\underline{x}), \text{ for all } i, j = 1, 2, \dots, n, \quad (6-97)$$

and, with (6-92), it follows that

$$\left\{ \frac{\partial^2 \nabla^0}{\partial \underline{x}_1 \partial \underline{x}_1} \right\}^T = \left\{ \frac{\partial^2 \nabla^0}{\partial \underline{x}_1 \partial \underline{x}_1} \right\}, \quad (6-98)$$

$$\left\{ \frac{\partial^2 \nabla^0}{\partial \underline{x}_2 \partial \underline{x}_1} \right\}^T = \left\{ \frac{\partial^2 \nabla^0}{\partial \underline{x}_1 \partial \underline{x}_2} \right\} = -2 [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \left[\frac{\partial \underline{U}_e(\underline{x})}{\partial \underline{x}_1} \right] \quad (6-99)$$

and

$$\left\{ \frac{\partial^2 \nabla^0}{\partial \underline{x}_2 \partial \underline{x}_2} \right\}^T = \left\{ \frac{\partial^2 \nabla^0}{\partial \underline{x}_2 \partial \underline{x}_2} \right\} = -2 [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \left[\frac{\partial \underline{U}_e(\underline{x})}{\partial \underline{x}_2} \right]. \quad (6-100)$$

From (2-10), $w_i(\underline{x})$ can be calculated from its gradient as

$$\begin{aligned} w_i(\underline{x}) &= \int_0^{\underline{x}} \left[\frac{\partial w_i(\underline{x})}{\partial \underline{x}} \right]^T d\underline{x} \\ &= \int_0^{x_1} w_{1i}(\gamma_1, 0, \dots, 0) d\gamma_1 + \int_0^{x_2} w_{2i}(x_1, \gamma_2, 0, \dots, 0) d\gamma_2 + \dots \\ &\quad + \int_0^{x_{n-r}} w_{n-ri}(x_1, x_2, \dots, x_{n-r-1}, \gamma_{n-r}, 0, \dots, 0) d\gamma_{n-r} \end{aligned}$$

$$\begin{aligned}
& + \int_0^{x_{n-r+1}} w_{n-r+1i}(x_1, x_2, \dots, x_{n-r}, \gamma_{n-r+1}, 0, \dots, 0) d\gamma_{n-r+1} + \dots \\
& + \int_0^{x_n} \bar{w}_{ni}(x_1, x_2, \dots, x_{n-1}, \gamma_n) d\gamma_n. \quad (6-101)
\end{aligned}$$

Let conveniently describe the sum of the first $(n-r)$ integrals (6-101) as $w_i(\underline{x}_1)$, a function of \underline{x}_1 only. From (6-97), the last r integrals can be calculated with $w_{ij}(\underline{x})$, instead of $w_{ji}(\underline{x})$. Thus, from (6-98),

$$\begin{aligned}
\frac{\partial \nabla^0}{\partial \underline{x}_1} &= \begin{bmatrix} w_1(\underline{x}) \\ w_2(\underline{x}) \\ \vdots \\ w_{n-r}(\underline{x}) \end{bmatrix} = \begin{bmatrix} \tilde{w}_1(\underline{x}_1) \\ \tilde{w}_2(\underline{x}_1) \\ \vdots \\ \tilde{w}_{n-r}(\underline{x}_1) \end{bmatrix} \\
&+ \int_0^{\underline{x}_2} \begin{bmatrix} w_{n-r+1 \ 1}(\underline{x}), w_{n-r+2 \ 1}(\underline{x}), \dots, w_{n1}(\underline{x}) \\ w_{n-r+1 \ 2}(\underline{x}), w_{n-r+2 \ 2}(\underline{x}), \dots, w_{n2}(\underline{x}) \\ \vdots \\ w_{n-r+1 \ n-r}(\underline{x}), \dots, w_{nn-r}(\underline{x}) \end{bmatrix} d\underline{x}_2 \\
\Delta &= \tilde{w}(\underline{x}_1) + \int_0^{\underline{x}_2} \left\{ \frac{\partial^2 \nabla^0}{\partial \underline{x}_1 \partial \underline{x}_2} \right\}^T d\underline{x}_2
\end{aligned}$$

$$= -2 \int_0^{\underline{x}_2} \left[\frac{\partial \underline{U}_e(\underline{x})}{\partial \underline{x}_1} \right]^T [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] d\underline{x}_2 + \widetilde{\underline{w}}(\underline{x}_1) . \quad (6-102)$$

Then the gradient of $\nabla^0(\underline{x})$ becomes, from (6-93) and (6-102),

$$\left[\frac{\partial \nabla^0}{\partial \underline{x}} \right] = \begin{bmatrix} \frac{\partial \nabla^0}{\partial \underline{x}_1} \\ \vdots \\ \frac{\partial \nabla^0}{\partial \underline{x}_n} \end{bmatrix}$$

$$= \begin{bmatrix} -2 \int_0^{\underline{x}_2} \left[\frac{\partial \underline{U}_e(\underline{x})}{\partial \underline{x}_1} \right]^T [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] d\underline{x}_2 + \widetilde{\underline{w}}(\underline{x}_1) \\ -2 [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{U}_e(\underline{x}) \end{bmatrix} \quad (6-103)$$

and the Hamilton-Jacobi equation (6-96) becomes

$$\begin{aligned} L(\underline{x}) &= \underline{U}_e^T(\underline{x}) [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{U}_e(\underline{x}) + 2 \underline{U}_e^T [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{F}_2(\underline{x}) \\ &+ 2 \left\{ \int_0^{\underline{x}_2} \left[\frac{\partial \underline{U}_e(\underline{x})}{\partial \underline{x}_1} \right]^T [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] d\underline{x}_2 \right\}^T \underline{F}_1(\underline{x}) - \widetilde{\underline{w}}^T(\underline{x}_1) \underline{F}_1(\underline{x}) . \end{aligned}$$

(6-104)

6.4.3 Principal Theorem of the Inverse Problem

With the preceding results, the inverse problem can be investigated to determine combinations of $L(\underline{X})$ and \underline{R} to satisfy the Hamilton-Jacobi equation (6-104). The fundamental result can be stated as the following.

Theorem 6-1: The Inverse Problem.

For the inverse problem as described in Section 6.2 such that the optimal performance index function is of class C_2 , a performance index can be optimized by the given $\underline{U}(\underline{X})$ if and only if

$$(i) \quad \underline{U}_d(\underline{X}) = -\underline{R}_{11}^{-1} \underline{R}_{12} \underline{U}_e(\underline{X}) , \quad (6-105)$$

$$(ii) \quad [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \left[\frac{\partial \underline{U}_e(\underline{X})}{\partial \underline{X}_2} \right] \quad (6-106)$$

is symmetric, and

(iii) there exists an $(n - r)$ dimensional vector valued function $\underline{w}(\underline{X}_1)$ of class C_1 and insuring symmetry in

$$-2 \frac{\partial}{\partial \underline{X}_1} \left\{ \int_0^{\underline{X}_2} \left[\frac{\partial \underline{U}_e(\underline{X})}{\partial \underline{X}_1} \right]^T [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] d\underline{X}_2 \right\} + \left[\frac{\partial \underline{\tilde{w}}(\underline{X}_1)}{\partial \underline{X}_1} \right] . \quad (6-107)$$

The corresponding $\nabla^0(\underline{X})$ is given from (6-103) as

$$\nabla^0(\underline{x}) = \int_0^{\underline{x}} \left[\frac{\partial \nabla^0}{\partial \underline{x}} \right]^T d\underline{x} . \quad (6-108)$$

Proof: The necessity of the conditions has been shown by the previous work of the chapter. That is, (i) comes from the absolute minimum condition of the Hamiltonian (6-94). The symmetry in the functional matrix of (6-100) corresponds to (ii), and (iii) follows from the Hamilton-Jacobi equation and the symmetry in (6-98), using (6-102).

According to Theorem 5-1, the sufficiency of the conditions can be proved by showing that the existence of a unique function $\tilde{w}(\underline{x}_1)$ satisfying the conditions of the theorem can exist for each combination of $L(\underline{x})$ and \underline{R} as the Hamiltonian is normal. Assuming the contrary, that there exist two different functions, say $\tilde{w}_a(\underline{x}_1)$ and $\tilde{w}_b(\underline{x}_1)$, satisfying (iii) for a combination of $L(\underline{x})$ and \underline{R} . Necessarily, from (6-103) and (6-96),

$$\tilde{w}_a^T(\underline{x}_1) \underline{F}_1(\underline{x}) = \tilde{w}_b^T(\underline{x}_1) \underline{F}_1(\underline{x}) . \quad (6-109)$$

Describe the resulting optimal performance index functions as $\nabla_a^0(\underline{x})$ and $\nabla_b^0(\underline{x})$. Then from (6-103) and (6-108), it follows that

$$\nabla_b^0(\underline{x}) = \nabla_a^0(\underline{x}) + \varphi(\underline{x}_1), \quad (6-110)$$

where $\varphi(\underline{x}_1)$ is such that

$$\varphi(\underline{x}_1) = \int_0^{\underline{x}_1} \{-\tilde{\underline{w}}_a(\underline{x}_1) + \tilde{\underline{w}}_b(\underline{x}_1)\}^T d\underline{x}_1. \quad (6-111)$$

This is not identically zero by the contrary hypothesis.

However $\varphi(0) = 0$ as $\nabla_a^0(0) = \nabla_b^0(0) = 0$. Thus there exists a specific $(n - r)$ dimensional vector

$\underline{x} = [\underline{x}_1, \underline{x}_2, \dots, \underline{x}_{n-r}]^T$ satisfying

$$\varphi(\underline{x}_1) = k, \text{ a nonzero constant.} \quad (6-112)$$

Consider the hypersurface $\varphi(\underline{x}_1) = k$ in $R^n \times R^1$ which does not include $\underline{x} = 0$. However, consider the time derivative of $\varphi(\underline{x}_1)$ governed by the synthesized system equation. This follows from (6-111), (6-34), and (6-109),

$$\begin{aligned} \dot{\varphi}(\underline{x}_1) &= \left[\frac{\partial \varphi(\underline{x}_1)}{\partial \underline{x}} \right]^T \{ \underline{F}(\underline{x}) + \underline{B} \underline{U}(\underline{x}) \} \\ &= \left[[-\tilde{\underline{w}}_a(\underline{x}_1) + \tilde{\underline{w}}_b(\underline{x}_1)]^T \quad 0 \right] \begin{bmatrix} \underline{F}_1(\underline{x}) \\ \underline{F}_2(\underline{x}) + \underline{U}_e(\underline{x}) \end{bmatrix} \\ &= [-\tilde{\underline{w}}_a(\underline{x}_1) + \tilde{\underline{w}}_b(\underline{x}_1)]^T \underline{F}_1(\underline{x}) = 0 \end{aligned} \quad (6-113)$$

from (6-109). This implies that every motion of the synthesized system from points on $\varphi(\underline{x}_1) = k$ stays on this hypersurface and can never approach to origin. This contradicts the asymptotic stability in the large of the origin of the synthesized system, Lemma 6-1 and assumption (ii) of Section 6.1.

6.5 Discussion

6.5.1 On the General Method of Solution of the Inverse Problem

A solution to the inverse problem is obtained by determining all combinations of $L(\underline{x})$ and \underline{R} satisfying the conditions of Theorem 6-1. \underline{R} must be determined to meet conditions (i) and (ii) with respect to the given $\underline{U}(\underline{x})$. The corresponding $L(\underline{x})$ are then determined from (6-104) by choosing various $\tilde{w}(\underline{x}_1)$ satisfying condition (iii). If no positive definite symmetric \underline{R} exists for the given $\underline{U}(\underline{x})$, then the $\underline{U}(\underline{x})$ cannot be an optimal control law.

6.5.2 Dependency within $\underline{U}(\underline{X})$

As $\underline{U}_d(\underline{X})$ is dependent upon $\underline{U}_e(\underline{X})$, (6-105), there can exist at most r independent elements in an optimal feedback control law, where r is the rank of $\hat{\underline{B}}$. If $m = r$, i.e., $\hat{\underline{B}}$ is full rank, condition (i) of the theorem is nonexistent because $\underline{U}_d(\underline{X})$ is reduced to dimension zero.

In Thau's problem statement, Section 5.3.3, the rank of $\hat{\underline{B}}$ is not mentioned. If it is assumed to be either full rank or less than full rank, an additional condition corresponding to (6-105) must be given.

6.5.3 Consideration of the Variety of $\underline{L}(\underline{X})$

For an \underline{R} satisfying Theorem 6-1, a variety of $\underline{L}(\underline{X})$ may exist for which the given $\underline{U}(\underline{X})$ is an optimal feedback control law. These are associated with various $\hat{\underline{w}}(\underline{x}_1)$.

Corollary 6-1:

For an \underline{R} satisfying Theorem 6-1, assume an $\underline{L}_a(\underline{X})$ can be found from (6-104) as an optimized performance index for the given $\underline{U}(\underline{X})$. Let the resulting optimal performance index function be $\nabla_a^0(\underline{X})$. Then an $\underline{L}(\underline{X})$ can also be found from (6-104) for \underline{R} if and only if

$$L(\underline{X}) = L_a(\underline{X}) - \left[\frac{\partial \varphi(\underline{X}_1)}{\partial \underline{X}_1} \right]^T \underline{F}_1(\underline{X}) , \quad (6-114)$$

where $\varphi(\underline{X}_1)$ is an arbitrary scalar function of class C_2 and $\varphi(0) = 0$. Necessarily, the resulting optimal performance index function is

$$v^0(\underline{X}) = v_a^0(\underline{X}) + \varphi(\underline{X}_1) . \quad (6-115)$$

Proof: $L_a(\underline{X})$ corresponds to $\tilde{w}_a(\underline{X}_1)$ in (6-104). Then from (6-96), it follows that

$$\tilde{w}(\underline{X}_1) = \tilde{w}_a(\underline{X}_1) + \left[\frac{\partial \varphi(\underline{X}_1)}{\partial \underline{X}_1} \right] , \quad (6-116)$$

and the assertion follows directly from Theorem 6-1.

Possibilities for a nonnegative $L(\underline{X})$ can be considered by trying various $\varphi(\underline{X}_1)$ in (6-114). In the next chapter, the nonnegative property of $L(\underline{X})$ is examined in some detail under additional problem assumptions.

6.5.4 Uniqueness of $L(\underline{X})$

Consider the case of $n = m = r$ in (6-2), (6-28) and (6-29). Then the dimension of \underline{X}_1 is zero and the general method of solution stated in Section 6.5.1 can be simplified.

Condition (iii) of Theorem 6-1 is nonexistent because $\widetilde{W}(\underline{X}_1)$ is reduced to zero dimension. Thus, a corollary follows directly from Theorem 6-1.

Corollary 6-2:

For the inverse problem, if $n = m = r$, then a unique $L(\underline{X})$ corresponds to each \underline{R} and, from (6-104),

$$L(\underline{X}) = \underline{U}^T(\underline{X}) \underline{R} \underline{U}(\underline{X}) + 2 \underline{U}^T(\underline{X}) \underline{R} \underline{F}(\underline{X}) . \quad (6-117)$$

6.5.5 Linear Control Law

Let $\underline{U}(\underline{X})$ be specified as a linear feedback control law,

$$\underline{U}(\underline{X}) = -\underline{K}^T \underline{X} , \quad (6-118)$$

where \underline{K} as an $n \times m$ matrix decomposed as

$$\underline{K} \triangleq \left[\begin{array}{cc} \underline{K}_{11} & \underline{K}_{12} \\ \underline{K}_{21} & \underline{K}_{22} \end{array} \right] \begin{array}{l} \left. \vphantom{\begin{array}{c} \underline{K}_{11} \\ \underline{K}_{21} \end{array}} \right\} n-r \\ \left. \vphantom{\begin{array}{c} \underline{K}_{12} \\ \underline{K}_{22} \end{array}} \right\} r \end{array} . \quad (6-119)$$

$\underbrace{\hspace{1.5cm}}_{m-r} \qquad \underbrace{\hspace{1.5cm}}_r$

Thus by (6-29), (6-118) is

$$\underline{U}_d(\underline{X}) = -\underline{K}_{11}^T \underline{X}_1 - \underline{K}_{21}^T \underline{X}_2 \quad (6-120)$$

$$\underline{U}_e(\underline{X}) = -\underline{K}_{12}^T \underline{X}_1 - \underline{K}_{22}^T \underline{X}_2$$

and

$$\frac{\partial \underline{U}_e(\underline{X})}{\partial \underline{X}_1} = -\underline{K}_{12}^T \quad (6-121)$$

$$\frac{\partial \underline{U}_e(\underline{X})}{\partial \underline{X}_2} = -\underline{K}_{22}^T \quad .$$

The conditions in Theorem 6-1 can be identified directly with these submatrices.

Assuming a linear feedback control law (6-118),
(6-104) becomes

$$\begin{aligned} L(\underline{X}) &= \underline{X}^T \begin{bmatrix} \underline{K}_{12} \\ \underline{K}_{22} \end{bmatrix} [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] [\underline{K}_{12}^T \quad \underline{K}_{22}^T] \underline{X} \\ &\quad - 2\underline{X}^T \begin{bmatrix} \underline{K}_{12} \\ \underline{K}_{22} \end{bmatrix} [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{F}_2(\underline{X}) \\ &\quad - 2\underline{X}_2^T [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{K}_{12}^T \underline{F}_1(\underline{X}) - \underline{\tilde{W}}^T(\underline{X}_1) \underline{F}_1(\underline{X}) \quad . \end{aligned} \quad (6-122)$$

For the equivalent statement of condition (i), by (6-105), it can be shown from (6-120) that

(i)'

$$\begin{aligned} -\underline{K}_{11}^T &= \underline{R}_{11}^{-1} \underline{R}_{12} \underline{K}_{12}^T \\ -\underline{K}_{21}^T &= \underline{R}_{11}^{-1} \underline{R}_{12} \underline{K}_{22}^T . \end{aligned} \quad (6-123)$$

Conditions (ii) and (iii) follow directly, i.e.,

(ii)'

$$-[\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{K}_{22}^T \quad (6-124)$$

is symmetric, and

(iii)'

$$\left[\frac{\partial \tilde{w}(\underline{x}_1)}{\partial \underline{x}_1} \right] \quad (6-125)$$

is symmetric.

From (6-118) and (6-93), it follows that

$$\frac{\partial \nabla^0}{\partial \underline{x}_2} = 2 [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] [\underline{K}_{12}^T \quad \underline{K}_{22}^T] \underline{x} . \quad (6-126)$$

Therefore

$$\frac{\partial^2 \nabla^0}{\partial \underline{x}_1 \partial \underline{x}_2} = \left\{ \frac{\partial^2 \nabla^0}{\partial \underline{x}_2 \partial \underline{x}_1} \right\}^T = 2 [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{K}_{12}^T$$

$$\frac{\partial^2 \nabla^0}{\partial \underline{x}_2 \partial \underline{x}_2} = 2 [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{K}_{22}^T \quad (6-127)$$

$$\frac{\partial^2 \nabla^0}{\partial \underline{x}_1 \partial \underline{x}_1} = \frac{\partial^2 \widetilde{w}(\underline{x}_1)}{\partial \underline{x}_1}$$

Consequently, it follows that

$$\nabla^0(\underline{x}) = \underline{x}^T \begin{bmatrix} [0] & \underline{K}_{12} [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \\ [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{K}_{12}^T & [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{K}_{22}^T \end{bmatrix} \underline{x} + \varphi(\underline{x}_1) , \quad (6-128)$$

where

$$\varphi(\underline{x}_1) = \int_0^{\underline{x}_1} \widetilde{w}^T(\underline{x}_1) d\underline{x}_1 . \quad (6-129)$$

From (6-128), the structure of $\nabla^0(\underline{x})$ is a sum of an arbitrary scalar function $\varphi(\underline{x}_1)$ and a quadratic form in \underline{x} determined by \underline{K} and \underline{R} . If $n = m = r$, then only one $\nabla^0(\underline{x})$, a quadratic form, can exist for each \underline{R} , according to Corollary 6-2.

Considering the results of Suga for a linear control law, Section 5.3.2, from the viewpoint of the above results, the symmetry in (6-124) corresponds to that of $\hat{R}(t) \hat{K}(t) \hat{R}(t)$ in Theorem 5-3. It is interesting to note that the structure of $\hat{V}^0(\underline{y}, t)$ in (5-56), i.e., a quadratic in \underline{y} determined by the given $\underline{V}(\underline{y})$ plus an arbitrary function, is invariant for the general non-linear system, (6-128). Thus the flexibility of $L(\underline{x})$ due to the function $\hat{W}(\underline{x}_1)$ corresponds to that of $\hat{L}(\underline{y}, t)$ due to $\hat{r}(\underline{y}, t)$ in (5-49).

6.5.6 Nonnegative $\nabla^0(\underline{x})$ for a Linear Control Law

For a linear control law, a definitive statement is possible for the sign definiteness of $\nabla^0(\underline{x})$.

Theorem 6-2:

For a linear feedback control law (6-118) in the inverse problem, the resulting $\nabla^0(\underline{x})$, (6-128), is positive semi-definite in R^n if and only if

- (i) the last $r - r_1$ columns of $\underline{R}_{12} \underline{D}_{22}$ are null, r_1 the rank of $[\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}]$, and
- (ii) a function

$$\underline{x}_1^T \tilde{\underline{\Omega}} \underline{x}_1 + \varphi(\underline{x}_1) \quad (6-130)$$

is positive semidefinite in R^{n-r} , where

$$\tilde{\underline{\Omega}} \triangleq -\underline{K}_{12}^T \underline{R}_0 \underline{D}_{22} \underline{D}_{22}^T \underline{R}_0 \underline{K}_{22}^T \underline{D}_{22} \underline{D}_{22}^T \underline{R}_0 \underline{K}_{12}^T, \quad (6-131)$$

$$\underline{R}_0 \triangleq [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] , \quad (6-132)$$

and \underline{D}_{22} is an $r \times r$ nonsingular matrix for the congruent transformation

$$\underline{D}_{22}^T [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}] \underline{K}_{22}^T \underline{D}_{22} = \begin{bmatrix} \underline{I}_{r_1} & [0] \\ [0] & [0] \end{bmatrix} . \quad (6-133)$$

Proof: For (6-128) to be semidefinite, \underline{R}_0 must be positive semidefinite, from

$$\nabla^0([0]^T, \underline{x}_2^T) = \underline{x}_2^T \underline{R}_0 \underline{K}_{22}^T \underline{x}_2 . \quad (6-134)$$

Then the matrix \underline{D}_{22} satisfying (6-133) exists, from Theorem 2-4. Define a nonsingular matrix

$$\underline{D} \triangleq \begin{bmatrix} \underline{I}_{n-r} & [0] \\ -\underline{A} & \underline{D}_{22} \end{bmatrix} \quad (6-135)$$

and a function

$$\tilde{V}^0(\underline{X}) = V^0(\underline{D} \underline{X}) , \quad (6-136)$$

where

$$\underline{A} = \underline{D}_{22} \underline{D}_{12}^T \underline{R}_0 \underline{K}_{22}^T \underline{D}_{22} \underline{D}_{22}^T \underline{R}_0 \underline{K}_{12}^T . \quad (6-137)$$

Then

$$(a) \quad \tilde{V}^0(\underline{0}) = V^0(\underline{0}) = 0 ,$$

(b) if $\tilde{V}^0(\underline{X})$ is positive semidefinite, then

$$V^0(\underline{X}) = \tilde{V}^0(\underline{D}^{-1} \underline{X}) \geq 0 , \text{ for all } \underline{X} \in \mathbb{R}^n ,$$

and (c) if $V^0(\underline{X})$ is positive semidefinite, then

$$\tilde{V}^0(\underline{X}) = V^0(\underline{D} \underline{X}) \geq 0 , \text{ for all } \underline{X} \in \mathbb{R}^n .$$

Thus $V^0(\underline{X})$ is positive semidefinite if and only if $\tilde{V}^0(\underline{X})$ is positive semidefinite. Therefore, from (6-128), it follows that

$$\begin{aligned}
V^0(D \underline{x}) &= \underline{x}^T \begin{bmatrix} \underline{\Omega} & [0] & \underline{R} \\ [0] & I_{r_1} & [0] \\ \underline{R}^T & [0] & [0] \end{bmatrix} \underline{x} + \varphi(\underline{x}_1) \\
&= \underline{x}_1^T \underline{\Omega} \underline{x}_1 + 2\underline{x}_1^T \begin{bmatrix} [0] & \underline{R} \end{bmatrix} \underline{x}_2 \\
&\quad + \underline{x}_2^T \begin{bmatrix} I_{r_1} & [0] \\ [0] & [0] \end{bmatrix} \underline{x}_2 + \varphi(\underline{x}_1), \quad (6-138)
\end{aligned}$$

where \underline{R} is defined by the last $r - r_1$ column of $K_{12} \underline{R}_0 D_{22}$. Applying a similar argument as in Lemmas 6-3 for (6-138), the necessity of the conditions can be justified. The sufficiency is apparent from (6-138) if $\underline{R} = [0]$.

Theorem 6-3:

For a linear feedback control law (6-118) in the inverse problem, the resulting $V^0(\underline{x})$, (6-128), is positive definite if and only if

- (i) $\underline{R}_0 K_{22}^T$ is positive definite, and

$$(ii) \quad -\underline{x}_1^T \underline{K}_{12} \underline{K}_{22}^{-1} \underline{R}_{012} \underline{K}_{12}^T \underline{x}_1 + \varphi(\underline{x}_1) \quad (6-139)$$

is positive definite in \mathbb{R}^{n-r} and $\underline{R}_0 = [\underline{R}_{22} - \underline{R}_{12}^T \underline{R}_{11}^{-1} \underline{R}_{12}]$.

Proof: The necessity of (i) follows directly from the positive definiteness of $\nabla^0(\underline{x})$ for $\underline{x}_1 = \underline{0}$ or

$$\nabla^0([0^T, \underline{x}_2^T]) = \underline{x}_2^T \underline{R}_0 \underline{K}_{22}^T \underline{x}_2. \quad (6-140)$$

For the necessity of (ii), define

$$\underline{D}_1 \triangleq \begin{bmatrix} \underline{I}_{n-r} & [0] \\ -(\underline{K}_{22}^{-1})^T \underline{K}_{12}^T & \underline{I}_r \end{bmatrix}. \quad (6-141)$$

The inverse of $\underline{R}_0 \underline{K}_{22}^T$ exists and is symmetric by (6-124).

Then with the same argument as used in the proof of

Theorem 6-2, $\nabla^0(\underline{x})$ is positive definite if and only if

$\nabla^0(\underline{D}_1 \underline{x})$ is positive definite. It follows that

$$\nabla^0(\underline{D}_1 \underline{x}) = \underline{x}^T \begin{bmatrix} -\underline{K}_{12} \underline{K}_{22}^{-1} \underline{R}_{012} \underline{K}_{12}^T & [0] \\ [0] & \underline{R}_0 \underline{K}_{22}^T \end{bmatrix} \underline{x} + \varphi(\underline{x}_1)$$

$$= -\underline{x}_1^T \underline{R}_{12} \underline{K}_{22}^{-1} \underline{R}_{012} \underline{K}_{12}^T \underline{x}_1 + \varphi(\underline{x}_1) + \underline{x}_2^T \underline{R}_0 \underline{K}_{22}^T \underline{x}_2,$$

(6-142)

and (ii) is necessary. The sufficiency of the conditions follows from (6-142).

The nonnegative characteristic of $\nabla^0(\underline{x})$ is completely established by the above analysis. This is in contrast to Thau's partial consideration of the topic, as discussed in Section 5.3.3.

6.5.7 Necessity of Control Action

The problem of whether optimal control can exist for "No control action" for a nonnegative $L(\underline{x})$ in the performance index can be considered through the corresponding inverse problem. For the class of systems given by (6-68), let $\underline{U}(\underline{x})$ be identically zero, that is, the equation of the synthesized feedback control system is

$$\dot{\underline{x}} = \underline{A} \underline{x} + \underline{B}_e \underline{F}_2(\underline{x}) . \quad (6-143)$$

Also assume that the origin of this system is asymptotically stable in the large. For $v \geq 2$, substituting $\underline{U}(\underline{x}) = \underline{0}$ into (6-104), $L(\underline{x})$ for the optimized performance index is

$$L(\underline{x}) = - \tilde{\underline{w}}^T(\underline{x}_1) \underline{A}_1 \underline{x} . \quad (6-144)$$

Then, from (6-136) and Lemma 6-6, for (6-144) to be positive semidefinite,

(a) $\tilde{w}(\underline{x}_1)$ must be identically zero, if the given system is completely controllable, or

(b) $\tilde{w}(\underline{x}_1)$ must be a function of $\underline{x}_{(1)}$ only, if the given system is uncontrollable, say $\tilde{w}(\underline{x}_{(1)})$.

Subsequently, (6-144) reduces to

$$L(\underline{x}) = \begin{cases} 0, & \text{if the given system is completely controllable,} \\ L(\underline{x}_{(1)}), & \text{if the given system is uncontrollable,} \\ & \text{where the state variables in } \underline{x}_{(1)} \text{ are the} \\ & \text{uncontrollable state variables, as discussed} \\ & \text{in Section 4.3.3.} \end{cases}$$

If $v = 1$, then $\left[\frac{\partial V^0}{\partial \underline{x}}\right] = \left[\frac{\partial V^0}{\partial \underline{x}_2}\right] = \underline{0}$ from (6-103), and

$L(\underline{x})$ is identically zero.

Observing these results from the viewpoint of the forward problem, an important characteristic of an optimal feedback control system is evident.

Principle of Necessary Control Action:

Consider the optimal regulator problem such that

- (i) the system is given as (6-68),
- (ii) the desired final condition of the system is $\underline{x} = \underline{0}$,
- (iii) the performance index is

$$\int_0^{\infty} \{L(\underline{x}) + \underline{u}^T \underline{R} \underline{u}\} dt ,$$

where \underline{R} is positive definite, and $L(\underline{x})$ is a function of the controllable state variables and is positive semi-definite, then some control action is necessary for optimality, i.e., the optimal feedback control law cannot be identically zero.

6.5.8 Asymptotic Stability of the Synthesized Feedback Control System

The feedback control law $\underline{u}(\underline{x})$ by assumption (ii) of Section 6.1 requires that the origin of the synthesized system be asymptotically stable in the large. No definitive criterion exists to verify this, except for linearly synthesized systems. Practically, if \underline{R} , $L(\underline{x})$ and $\nabla^0(\underline{x})$ are

determined according to the procedure of Section 6.5.1 for some given $\underline{U}(\underline{X})$, it is then necessary to check the synthesized system for asymptotic stability. It may be possible to do this by applying the Liapunov direct method.

Theorem 6-4:

Let \underline{R} , $L_a(\underline{X})$ and $\nabla_a^0(\underline{X})$ be calculated for some $\underline{U}(\underline{X})$, following the procedure of Section 6.5.1. If there is a scalar function $\tilde{\Phi}(\underline{x}_1)$ such that it is of class C_2 satisfying

$$\nabla^0(\underline{X}) \triangleq \nabla_a^0(\underline{X}) + \tilde{\Phi}(\underline{x}_1), \quad (6-145)$$

and $\nabla^0(\underline{X})$ and

$$L_a(\underline{X}) + \underline{U}^T(\underline{X}) \underline{R} \underline{U}(\underline{X}) - \left[\frac{\partial \tilde{\Phi}(\underline{x}_1)}{\partial \underline{x}_1} \right]^T \underline{F}_1(\underline{X}) \quad (6-146)$$

are positive definite, then the origin of the synthesized system is asymptotically stable in the large.

Proof: From Corollary 6-1, $L(\underline{X})$ and $\nabla^0(\underline{X})$ in (6-114) and (6-115) can be determined as a function of \underline{R} and $\underline{U}(\underline{X})$. Since the time derivative of (6-145) as governed by the synthesized system is

$$\begin{aligned} \dot{V}^0(\underline{x}) = \frac{d}{dt} \{ V_a^0(\underline{x}) + \tilde{\Phi}(\underline{x}_1) \} = & -L_a(\underline{x}) - \underline{U}^T(\underline{x}) R \underline{U}(\underline{x}) \\ & + \left[\frac{\partial \tilde{\Phi}(\underline{x}_1)}{\partial \underline{x}_1} \right]^T \underline{F}_1(\underline{x}) , \end{aligned} \quad (6-147)$$

using (6-103), the proof follows directly from Theorem 2-8.

This theorem provides a sufficient condition only; accordingly a failure of the conditions does not necessarily mean that the origin of the synthesized system is not asymptotically stable in the large.

6.5.9 Miscellaneous Comments

As a generalization of the inverse problem first considered by Kalman (Chapter 5), Theorem 6-1 and the succeeding developments of this chapter are shown to include results of other authors. In addition, this work reveals new important characteristics of optimal feedback control systems, i.e., Sections 6.5.2 and 6.5.7. Moreover these results are presented very compactly as a result of the developed canonical form of Chapter 4.

6.6 Examples

Example 6-1: Consider a system given by

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ -x_3^3 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad (6-148)$$

and assume a feedback control law

$$\begin{bmatrix} u_1(\underline{x}) \\ u_2(\underline{x}) \end{bmatrix} = \begin{bmatrix} -1 & -2 & 0 \\ 0 & 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 0 \\ -x_3^3 \end{bmatrix} . \quad (6-149)$$

As $n = 3$ and $m = r = 2$, it follows for the canonical form, from (6-28~30), that

$$\underline{x}_1 = x_1 , \quad (6-150)$$

$$\underline{x}_2 = [x_2, x_3]^T ,$$

$$\underline{u}_d(\underline{x}) \quad \text{nonexisting} , \quad (6-151)$$

$$\underline{u}_e(\underline{x}) = \underline{u}(\underline{x}) ,$$

and

$$\begin{aligned} \underline{F}_1(\underline{x}) &= x_2, \\ \underline{F}_2(\underline{x}) &= \begin{bmatrix} 0 \\ -2x_3 - x_3^3 \end{bmatrix}. \end{aligned} \quad (6-152)$$

Thus,

$$\left[\frac{\partial \underline{U}_e}{\partial \underline{x}_1} \right] = \left[\frac{\partial u_1(\underline{x})}{\partial x_1} \quad \frac{\partial u_2(\underline{x})}{\partial x_1} \right]^T = [-1 \quad 0]^T \quad (6-153)$$

and

$$\left[\frac{\partial \underline{U}_e}{\partial \underline{x}_2} \right] = \begin{bmatrix} \frac{\partial u_1(\underline{x})}{\partial x_2} & \frac{\partial u_1(\underline{x})}{\partial x_3} \\ \frac{\partial u_2(\underline{x})}{\partial x_2} & \frac{\partial u_2(\underline{x})}{\partial x_3} \end{bmatrix} = \begin{bmatrix} -2 & 0 \\ 0 & -1 - 3x_3^2 \end{bmatrix}. \quad (6-154)$$

Since \underline{x}_1 is one dimensional, define

$$\underline{\tilde{w}}(\underline{x}_1) \triangleq \tilde{w}(x_1). \quad (6-155)$$

Substituting (6-149~155) into (6-103) and (6-104),

$$\frac{\partial \nabla^0}{\partial \underline{x}} = \begin{bmatrix} \frac{\partial \nabla^0}{\partial x_1} \\ \frac{\partial \nabla^0}{\partial x_2} \\ \frac{\partial \nabla^0}{\partial x_3} \end{bmatrix} = \begin{bmatrix} 2 \int_0^{\underline{x}_2} [1, 0] \underline{R} \, d\underline{x}_2 + \underline{\tilde{w}}(x_1) \\ \underline{2R} \begin{bmatrix} x_1 + 2x_2 \\ x_3 + x_3^3 \end{bmatrix} \end{bmatrix} \quad (6-156)$$

and

$$L(\underline{X}) = [x_1 + 2x_2, x_3 + x_3^3] \underline{R} \begin{bmatrix} x_1 + 2x_2 \\ 5x_3 + 3x_3^3 \end{bmatrix} \\ - 2 \left\{ \int_0^{\underline{X}_2} [1, 0] \underline{R} d\underline{X}_2 \right\} x_2 - \{\tilde{w}(x_1)\} x_2, \quad (6-157)$$

where $\underline{R} = \underline{R}_{22} - \underline{R}_{12} \underline{R}_{11}^{-1} \underline{R}_{12}^T = \underline{R}_{22}$ because \underline{E} is of full rank in (6-148). Referring to the statement of Theorem 6-1, (i) is nonexistent and the symmetry for (6-108) is satisfied since \underline{X}_1 is one dimensional. Consequently, for the system (6-148), a performance index

$$\int_{t_0}^{\infty} \{L(\underline{X}) + \underline{U}^T \underline{R} \underline{U}\} dt \quad (6-158)$$

can be optimized by the feedback control law (6-149) if and only if symmetry exists for

$$\underline{R} = \begin{bmatrix} -2 & 0 \\ 0 & -1-3x_3^2 \end{bmatrix}, \quad (6-159)$$

and $\tilde{w}(x_1)$ is of class C_1 .

Arbitrarily choose $\tilde{w}(x_1) = 2x_1$ and

$$\underline{R} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad (6-160)$$

which satisfy the above conditions. Then (6-156 and 157) become

$$\left[\frac{\partial \nabla^0}{\partial \underline{X}} \right] = \begin{bmatrix} 2x_1 + 2x_2 \\ 2x_1 + 4x_2 \\ 2x_3 + 2x_3^3 \end{bmatrix} \quad (6-161)$$

and

$$\begin{aligned} L(\underline{X}) &= x_1^2 + 2x_1x_2 + 2x_2^2 + 5x_3^2 + 8x_3^4 + 3x_3^6 \\ &= (x_1+x_2)^2 + x_2^2 + x_3^2(5+8x_3^2+3x_3^4) \quad (6-162) \end{aligned}$$

From (6-161) and (6-108), the optimal performance index function is

$$\begin{aligned} \nabla^0(\underline{X}) &= x_1^2 + 2x_1x_2 + 2x_2^2 + x_3^2 + \frac{x_3^4}{2} \\ &= (x_1+x_2)^2 + x_2^2 + x_3^2 + \frac{x_3^4}{2} \quad (6-163) \end{aligned}$$

From Corollary 6-1, the performance indices from (6-114), that is

$$\int_{t_0}^{\infty} \{ (x_1+x_2)^2 + x_2^2 + x_3^2(5+8x_3^2+3x_3^4) - \left\{ \frac{\partial \tilde{\Phi}(x_1)}{\partial x_1} \right\} x_2 + u_1^2 + u_2^2 \} dt, \quad (6-164)$$

can be optimized by the feedback control law (6-149).

The optimal performance index function from (6-115) then becomes

$$V^0(\underline{x}) = (x_1+x_2)^2 + x_2^2 + x_3^2 + \frac{x_3^4}{2} + \tilde{\Phi}(x_1), \quad (6-165)$$

where $\tilde{\Phi}(x_1)$ is any function of class C_2 , e.g., $x_1 + x_1^3$.

To check the asymptotic stability of the synthesized system, combine (6-148) and (6-149) for

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} x_2 \\ -x_1 - 2x_2 \\ -3x_3 - 2x_3^3 \end{bmatrix}. \quad (6-166)$$

Then $V^0(\underline{x})$ for $\tilde{\Phi}(x_1) = 0$ is positive definite in R^n and it follows, from (6-164), that

$$\dot{V}^0(\underline{x}) = \left[\frac{\partial V^0}{\partial \underline{x}} \right]^T \{ \underline{F}(\underline{x}) + \underline{B} \underline{U}(\underline{x}) \} = -2 \{ (x_1 + \frac{3}{2}x_2)^2 + \frac{3}{4}x_2^2 + x_3^2(3+5x_2^2+2x_3^4) \}, \quad (6-167)$$

which is negative definite in R^n . From Theorem 6-4, the origin of (6-166) is asymptotically stable in the large.

Example 6-2: Consider a system given by

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_1 + x_2 - x_1^3 - x_1 x_2^2 \\ -x_1 + x_2 - x_1^2 x_2 - x_2^3 \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad (6-168)$$

and assume a linear feedback control law given by

$$\underline{U}(\underline{X}) = \begin{bmatrix} u_1(\underline{X}) \\ u_2(\underline{X}) \end{bmatrix} = \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad (6-169)$$

As $n = m = r$, it follows for the canonical form that

\underline{x}_1 is nonexistent and $\underline{x} = \underline{x}_2$,

$$\left[\frac{\partial \underline{U}}{\partial \underline{X}} \right] = \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix}, \quad (6-170)$$

$$\left[\frac{\partial \nabla^0}{\partial \underline{X}} \right] = \begin{bmatrix} \frac{\partial \nabla^0}{\partial x_1} \\ \frac{\partial \nabla^0}{\partial x_2} \end{bmatrix} = \underline{R} \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \underline{x}, \quad (6-171)$$

and

$$L(\underline{X}) = [x_1, x_2] \underline{R} \begin{bmatrix} -x_1 - 2x_2 + 2x_1^3 + 2x_1x_2^2 \\ 2x_1 - x_2 + 2x_1^2x_2 + 2x_2^3 \end{bmatrix}, \quad (6-172)$$

from (6-106). Arbitrarily assume

$$\underline{R} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}. \quad (6-173)$$

Then the symmetry of $\underline{R} \left[-\frac{\partial U(\underline{X})}{\partial \underline{X}} \right]$ is satisfied. From

Corollary 6-2, a performance index given by

$$\int_{t_0}^{\infty} \{L(\underline{X}) + u_1^2 + u_2^2\} dt \quad (6-174)$$

can be optimized by the linear feedback control law

(6-169) only if

$$L(\underline{X}) = -(x_1^2 + x_2^2) + 2(x_1^2 + x_2^2)^2. \quad (6-175)$$

Then the optimal performance index function becomes

$$V^0(\underline{X}) = x_1^2 + x_2^2. \quad (6-176)$$

To check the asymptotic stability of the synthesized feedback control system, calculate

$$\dot{V}^0(\underline{x}) = \left[\frac{\partial V^0}{\partial \underline{x}} \right]^T \{ \underline{F}(\underline{x}) + \underline{B} \underline{u} \} = -(x_1^2 + x_2^2)^2 \quad (6-177)$$

from (6-168), (6-169) and (6-171). Since $V^0(\underline{x})$ and $-\dot{V}^0(\underline{x})$ are positive definite, the origin of the synthesized system are asymptotically stable in the large.

Chapter 7

THE INVERSE PROBLEM OF LINEARLY SYNTHESIZED FEEDBACK CONTROL SYSTEMS

The inverse problem considered in this chapter is a subclass of the inverse problem of the previous chapter. This subclass is identified basically by the additional assumption that the synthesized feedback control system is linear. The precise statement of this problem, called the modified inverse problem, is given in Section 7.1. After a lemma is presented in Section 7.2, the results of Chapter 6 are restated for the specific subclass considered in this chapter. Finally the nonnegative property of the loss functions in an optimized performance index is discussed. The purpose for considering this modified inverse problem is to establish more general conclusions about optimal feedback control systems synthesized as linear systems, than have been presented in the literature.

7.1 Statement of the Modified Inverse Problem

The modified inverse problem considered in this chapter has additional assumptions to those stated for the inverse problem. In addition to assumptions (i)-(iv) in Section 6.1, it is assumed for the modified inverse problem that:

(i) For the system equation given in the canonical form (6-7),

(a) $\hat{\underline{B}}$ is of full rank, i.e., $m = r$,

(b) $\underline{F}_1(\underline{x})$ is identically zero, and

(c) $\underline{F}_2(\underline{x})$ is an r -dimensional, vector valued, finite degree polynomial function given as

$$\underline{F}_2(\underline{x}) = \underline{F}_2^{(2)}(\underline{x}) + \underline{F}_2^{(3)}(\underline{x}) + \dots + \underline{F}_2^{(\psi)}(\underline{x}) , \quad (7-1)$$

where each $\underline{F}_2^{(i)}(\underline{x})$ is i^{th} degree homogeneous and $\underline{F}_2^{(\psi)}(\underline{x})$ is not identically zero. (If a linear system is given, it is convenient to set $\psi = 1$ in (7-1) by identifying $\underline{A}_2 \underline{x} = \underline{F}_2^{(1)}(\underline{x})$ in (6-35).) Thus the system (6-7) can be written as

$$\dot{\underline{x}} = \underline{A} \underline{x} + \underline{B} \{ \underline{F}_2(\underline{x}) + \underline{u} \} ,$$

or, from (4-192) and (4-106),

$$\begin{bmatrix} \dot{\underline{x}}_{(1)} \\ \dot{\underline{x}}_{(2)} \\ \cdot \\ \cdot \\ \cdot \\ \dot{\underline{x}}_{(v)} \end{bmatrix} = \begin{bmatrix} \underline{A}_{(1,1)} & \underline{A}_{(1,2)} & & & & \\ & & \underline{A}_{(2,3)} & & & \\ & & & \cdot & & \\ & & & & \cdot & \\ & & & & & \cdot \\ & & & & & \underline{A}_{(v-1,v)} \\ \underline{A}_{(v,1)} & \cdot & \cdot & \cdot & \cdot & \underline{A}_{(v,v)} \end{bmatrix} \begin{bmatrix} \underline{x}_{(1)} \\ \underline{x}_{(2)} \\ \cdot \\ \cdot \\ \cdot \\ \underline{x}_{(v)} \end{bmatrix} \\
 + \begin{bmatrix} [0] \\ \underline{F}_2(\underline{x}) \end{bmatrix} + \begin{bmatrix} [0] \\ \underline{I}_r \end{bmatrix} \underline{u}, \text{ all other entries zero,} \quad (7-2)$$

which has the nonlinear functions in the last r equations corresponding to those equations which also have independent control variables,

(ii) For $\underline{u}(\underline{x})$ of (6-9)

$$\underline{u}(\underline{x}) = -\underline{K}^T \underline{x} - \underline{F}_2(\underline{x}), \quad (7-3)$$

where \underline{K} is an $n \times r$ matrix such that

$$\underline{K} = \left[\begin{array}{c} \underline{K}_{12} \\ \underline{K}_{22} \end{array} \right] \left\{ \begin{array}{l} n - r \\ r \end{array} \right\} \quad (7-4)$$

to provide negative real parts for all eigenvalues of $[\underline{A} - \underline{B} \underline{K}^T]$. The synthesized feedback control system, therefore, is the linear autonomous system

$$\dot{\underline{X}} = [\underline{A} - \underline{B} \underline{K}^T] \underline{X} \triangleq \underline{A}_K \underline{X} \quad (7-5)$$

and the origin is asymptotically stable in the large.

If the system (7-2) is uncontrollable, then all diagonal elements of $\underline{A}_{(1,1)}$ must have negative real parts for asymptotic stability, i.e., state variables in $\underline{X}_{(1)}$ behave as $\dot{\underline{X}} = \underline{A}_{(1,1)} \underline{X}_{(1)}$.

(iii)' $L(\underline{X})$, given by (6-27), is restricted to the form

$$L(\underline{X}) = L^{(2)}(\underline{X}) + L^{(3)}(\underline{X}) + \dots + L^{(\sigma)}(\underline{X}), \quad (7-6)$$

where each $L^{(i)}(\underline{X})$ is i^{th} degree homogeneous and $\sigma \geq 2$ is a given integer.

7.2 Fundamental Lemma

Lemma 7-1:

The form of the optimal performance index function

$\nabla^0(\underline{x})$ for the modified inverse problem is

$$\nabla^0(\underline{x}) = \nabla^{(2)}(\underline{x}) + \nabla^{(3)}(\underline{x}) + \dots + \nabla^{(\xi)}(\underline{x}) , \quad (7-7)$$

where

$$\xi \triangleq \max(\sigma, 2\psi) , \quad (7-8)$$

and σ and ψ are given by (7-1) and (7-6).

Proof: From (7-3) and (7-6), the loss function with a feedback control law is

$$\begin{aligned} L(\underline{x}) + \underline{U}^T(\underline{x}) \underline{R} \underline{U}(\underline{x}) &= L^{(2)}(\underline{x}) + L^{(3)}(\underline{x}) + \dots \\ &+ L^{(\sigma)}(\underline{x}) + \{-\underline{K}^T \underline{x} - \underline{F}_2(\underline{x})\}^T \underline{R} \{-\underline{K}^T \underline{x} - \underline{F}_2(\underline{x})\} \\ &\triangleq \mathcal{L}^{(2)}(\underline{x}) + \mathcal{L}^{(3)}(\underline{x}) + \dots + \mathcal{L}^{(\xi)}(\underline{x}) \\ &\triangleq \mathcal{L}(\underline{x}) , \end{aligned} \quad (7-9)$$

where each $\mathcal{L}^{(i)}(\underline{x})$ is i^{th} degree homogeneous. Assume a solution for the synthesized system (7-5) for an arbitrary (\underline{x}_0, t_0) ,

$$\underline{\phi}_f(t; \underline{x}_0, t_0) = e^{\underline{A}_K(t-t_0)} \underline{x}_0 . \quad (7-10)$$

The value of the performance index from (\underline{x}_0, t_0) is then

$$V^0(\underline{x}_0) = \int_{t_0}^{\infty} \mathcal{L}(e^{\underline{A}_K(t-t_0)} \underline{x}_0) dt. \quad (7-11)$$

But the integral

$$\int_{t_0}^{\infty} \mathcal{L}^{(i)}(e^{\underline{A}_K(t-t_0)} \underline{x}_0) dt, \text{ for } i = 2, 3, \dots, \xi, \quad (7-12)$$

is i^{th} degree homogeneous in \underline{x}_0 because the integrand is i^{th} degree in \underline{x}_0 . Defining (7-12) to be $V^{0(i)}(\underline{x}_0)$, the lemma is proved.

7.3 Solution of the Modified Inverse Problem

According to the assumptions of the modified inverse problem, $\frac{\partial V^0}{\partial \underline{x}}$, (6-103), and the Hamilton-Jacobi equation, (6-104), are reduced to

$$\left[\frac{\partial V^0}{\partial \underline{x}} \right] = \begin{bmatrix} \frac{\partial V^0}{\partial \underline{x}_1} \\ \frac{\partial V^0}{\partial \underline{x}_2} \end{bmatrix} = \begin{bmatrix} 2\underline{K}_{12}^R \underline{x}_2 + 2 \int_0^{\underline{x}_2} \left[\frac{\partial \underline{F}_2(\underline{x})}{\partial \underline{x}_1} \right]^T \underline{R} d\underline{x}_2 + \widetilde{w}(\underline{x}_1) \\ 2\underline{R} \underline{K}^T \underline{x} + 2\underline{R} \underline{F}_2(\underline{x}) \end{bmatrix}, \quad (7-13)$$

and

$$\begin{aligned}
 L(\underline{X}) &= \underline{X}^T \underline{K} \underline{R} \underline{K}^T \underline{X} - 2 \underline{X}^T \underline{K} \underline{R} \underline{A}_2 \underline{X} - 2 \underline{X}_2^T \underline{R} \underline{K}_{12}^T \underline{A}_1 \underline{X} \\
 &- 2 \underline{F}_2^T(\underline{X}) \underline{R} \underline{A}_2 \underline{X} - \underline{F}_2^T(\underline{X}) \underline{R} \underline{F}_2(\underline{X}) \\
 &- 2 \left\{ \int_0^{\underline{X}_2} \left[\frac{\partial \underline{F}_2(\underline{X})}{\partial \underline{X}_1} \right]^T \underline{R} d\underline{X}_2 \right\}^T \underline{A}_1 \underline{X} - \underline{\tilde{W}}^T(\underline{X}_1) \underline{A}_1 \underline{X} . \quad (7-14)
 \end{aligned}$$

Then Theorem 6-1 can be restated as follows.

Corollary 7-1:

For the modified inverse problem, a performance index $\int_{t_0}^{\infty} (L(\underline{X}) + \underline{U}^T \underline{R} \underline{U}) dt$ is optimized by a given $\underline{U}(\underline{X})$ if and only if

$$(i) \quad \underline{R} \underline{K}_{22}^T + \underline{R} \left[\frac{\partial \underline{F}_2(\underline{X})}{\partial \underline{X}_2} \right] \quad (7-15)$$

is symmetric,

(ii) there exists an $(n - r)$ dimensional vector valued function $\underline{\tilde{W}}(\underline{X}_1)$ such that

$$(a) \quad \underline{\tilde{W}}(\underline{X}) = \underline{\tilde{W}}^{(1)}(\underline{X}_1) + \underline{\tilde{W}}^{(2)}(\underline{X}_1) + \dots + \underline{\tilde{W}}^{(\xi-1)}(\underline{X}_1) , \quad (7-16)$$

where each $\tilde{w}^{(i)}(\underline{x}_1)$ is i^{th} degree homogeneous and ξ is given by (7-19),

$$(b) \quad 2 \frac{\partial}{\partial \underline{x}_1} \left\{ \int_0^{\underline{x}_2} \left[\frac{\partial F_2(\underline{x})}{\partial \underline{x}_1} \right] \underline{R} d\underline{x}_2 \right\} + \left[\frac{\partial \tilde{w}(\underline{x}_1)}{\partial \underline{x}_1} \right] \quad (7-17)$$

is symmetric.

Proof: As $m = r$, (i) of Theorem 6-1 is nonexistent.

The required symmetry of (7-15) and (7-17) follow from the requirements on (6-106) and (6-107). Since the highest and lowest degrees in $\nabla^0(\underline{x})$ are 2 and ξ respectively from (7-7), $\left[\frac{\partial \nabla^0}{\partial \underline{x}} \right]$ must be the sum of homogeneous polynomial functions with degree from 1 to $(\xi - 1)$. As $F_2(\underline{x})$ is also a polynomial, from (7-13) condition (7-16) must be satisfied.

Corollary 7-2:

Let \underline{R} and $L_a(\underline{x})$ be calculated for an optimized performance index for a given $\underline{U}(\underline{x})$ and a $\tilde{w}_a(\underline{x}_1)$ as (7-16), and let $\nabla_a^0(\underline{x})$ be the resulting optimal performance index function. Then, a performance index with \underline{R} and

$$L(\underline{x}) = L_a(\underline{x}) - \left[\frac{\partial \nabla_a^0(\underline{x}_1)}{\partial \underline{x}_1} \right]^T \underline{A}_1 \underline{x} \quad (7-18)$$

can be optimized by the same $\underline{U}(\underline{X})$, where $\varphi(\underline{X}_1)$ is an arbitrary ξ^{th} degree polynomial function from degree 2 to $(\xi - 1)$. Then the resulting optimal performance index function $\nabla^0(\underline{X})$ is

$$\nabla^0(\underline{X}) = \nabla_a^0(\underline{X}) + \varphi(\underline{X}_1) . \quad (7-19)$$

The proof follows directly from Corollary 6-1 with the additional assumptions given.

Consider $\sigma = 2$ in (7-6) and the given system is linear, i.e., $\underline{F}_2(\underline{X})$ in (7-2) is identically zero, and $\psi = 1$. Then the control law (7-3) is a linear control law and the optimal performance index function $\nabla^0(\underline{X})$ is a quadratic form, from Lemma 7-1. Necessarily $\widetilde{w}(\underline{X}_1)$ of (7-16) is a linear function of \underline{X}_1 , say

$$\widetilde{w}(\underline{X}_1) \triangleq 2\underline{P}_{11} \underline{X}_1 , \quad (7-20)$$

where \underline{P}_{11} is an $(n-r) \times (n-r)$ symmetric matrix according to (iib) of Corollary 7-1. Therefore, (7-13) is

$$\left[\frac{\partial \nabla^0}{\partial \underline{X}} \right] = 2 \begin{bmatrix} \underline{P}_{11} & \underline{K}_{12}^R \\ \underline{R} \underline{K}_{12}^T & \underline{R} \underline{K}_{22}^T \end{bmatrix} \begin{bmatrix} \underline{X}_1 \\ \underline{X}_2 \end{bmatrix} , \quad (7-21)$$

which results in

$$\nabla^0(\underline{x}) = \underline{x}^T \begin{bmatrix} \underline{P}_{11} & \underline{K}_{12}\underline{R} \\ \underline{R} \underline{K}_{12}^T & \underline{R} \underline{K}_{22}^T \end{bmatrix} \underline{x} \triangleq \underline{x}^T \underline{P} \underline{x} . \quad (7-22)$$

Substituting (7-20) into (7-14), $L(\underline{x})$ can be written as a quadratic form,

$$L(\underline{x}) \triangleq \underline{x}^T \underline{Q} \underline{x} \triangleq [\underline{x}_1^T \quad \underline{x}_2^T]^T \begin{bmatrix} \underline{Q}_{11} & \underline{Q}_{12} \\ \underline{Q}_{12}^T & \underline{Q}_{22} \end{bmatrix} \begin{bmatrix} \underline{x}_1 \\ \underline{x}_2 \end{bmatrix} , \quad (7-23)$$

where

$$\begin{aligned} \underline{Q}_{11} &= \underline{K}_{12}\underline{R} \underline{K}_{12}^T - \underline{K}_{12}\underline{R} \underline{A}_{21} - \underline{A}_{21}^T \underline{R} \underline{K}_{12}^T - \underline{P}_{11}\underline{A}_{11} - \underline{A}_{11}^T \underline{P}_{11} \\ \underline{Q}_{12} &= \underline{K}_{12}\underline{R} \underline{K}_{22}^T - \underline{A}_{21}^T \underline{R} \underline{K}_{22}^T - \underline{K}_{12}\underline{R} \underline{A}_{22} - \underline{A}_{11}^T \underline{K}_{12}\underline{R} - \underline{P}_{11}\underline{A}_{12} \\ \underline{Q}_{22} &= \underline{K}_{22}\underline{R} \underline{K}_{22}^T - \underline{K}_{22}\underline{R} \underline{A}_{22} - \underline{A}_{22}^T \underline{R} \underline{K}_{22}^T - \underline{R} \underline{K}_{12}^T \underline{A}_{12} - \underline{A}_{12}^T \underline{K}_{12}\underline{R} , \end{aligned} \quad (7-24)$$

and \underline{A}_{ij} is given in (6-32). Then Corollary 7-1 can be restated for this case.

Corollary 7-3:

For the modified inverse problem with the given system

linear and $L(\underline{X})$ restricted to a quadratic form, a performance index is optimized if and only if

- (i) $L(\underline{X})$ is given as (7-23 and 24),
- (ii) $\underline{R} \underline{K}_{22}^T$ is symmetric, and
- (iii) \underline{P}_{11} is symmetric.

Then $\nabla^0(\underline{X})$ is reduced to (7-22).

The sign definiteness of $L(\underline{X})$ and $\nabla^0(\underline{X})$ can be identified from \underline{Q} and \underline{P} , according to Lemmas 6-2 and 3. Thus

Corollary 7-4: The $\nabla^0(\underline{X})$ of Corollary 7-3, (7-22), is positive definite if and only if $\underline{R} \underline{K}_{22}^T$ and $\underline{P}_{11} - \underline{K}_{12} \underline{R} \underline{K}_{22}^{-1} \underline{K}_{12}^T$ are positive definite.

Corollary 7-5:

The $L(\underline{X})$ in (7-23) is positive semidefinite if and only if in (7-24)

- (i) \underline{Q}_{22} is positive semidefinite, say of rank r_1 ,
- (ii) the last $r - r_1$ columns of $\underline{Q}_{12} \underline{C}_{22}$ are null, and
- (iii) $\underline{Q}_{11} - \underline{Q}_{12} \underline{C}_{22}^T \underline{C}_{22} \underline{Q}_{22} \underline{C}_{22}^T \underline{C}_{22} \underline{Q}_{12}^T$ is positive semidefinite,

where \underline{C}_{22} is an $(n-r) \times (n-r)$ nonsingular for the congruent transformation

$$\underline{C}_{22}^T \underline{Q}_{22} \underline{C}_{22} = \begin{bmatrix} \underline{I}_{r_1} & [0] \\ [0] & [0] \end{bmatrix}. \quad (7-25)$$

Corollary 7-6:

The $L(\underline{X})$ in (7-23) is positive definite if and only if in (7-23)

- (i) \underline{Q}_{22} is positive definite, and
- (ii) $\underline{Q}_{11} - \underline{Q}_{12} \underline{Q}_{22}^{-1} \underline{Q}_{12}^T$ is positive definite.

Thus the inverse problem that Kalman originally presented has been generalized. While Kalman considered a controllable, single input, linear system with a linear control law, the results of this section are also applicable to uncontrollable, multi-input and not necessarily linear systems with a more general control law.

7.4 Nonnegative Loss Function of the Optimized Performance Index

Generally the loss function in (6-28) is assumed nonnegative. As \underline{R} in this equation is restricted to be positive definite, the nonnegative property of the loss function depends on that of $L(\underline{X})$. The nonnegative property

of $L(\underline{X})$ is considered in this section.

7.4.1 Nonnegative $L(\underline{X})$ for a Controllable System

Theorem 7-1:

For the modified inverse problem, assume that the system (7-2) is completely controllable. Then $L(\underline{X})$ in an optimized performance index can be positive semi-definite only as a quadratic form of \underline{X} .

Proof: Assume the contrary, that an $L(\underline{X})$ in an optimized performance index is positive semidefinite polynomial

$$L(\underline{X}) = L^{(2)}(\underline{X}) + L^{(3)}(\underline{X}) + \dots + L^{(\beta)}(\underline{X}), \quad (7-26)$$

where $L^{(\beta)}(\underline{X})$ is not identically zero for an arbitrary β , $2 < \beta \leq \sigma$ and σ is given by (7-6). For (7-26) to be positive semidefinite, β must be even, from Theorem 2-3.

- (1) Assume for the canonical form that $v = 1$. Then \underline{A}_1 , $\underline{\tilde{W}}(\underline{X}_1)$ and \underline{X}_1 are nonexistent. It follows from (7-2) and (7-14) that

$$L^{(\beta)}(\underline{X}) = - \underline{F}_2^{(\psi)T}(\underline{X}) \underline{R} \underline{F}_2^{(\psi)}(\underline{X})$$

$$= (\underline{D} \underline{F}_2^{(\psi)}(\underline{x}))^T [-\underline{I}_r] (\underline{D} \underline{F}_2^{(\psi)}(\underline{x})) , \quad (7-27)$$

where \underline{D} is an $r \times r$ nonsingular matrix satisfying

$$\underline{R} = \underline{D}^T \underline{D} . \quad (7-28)$$

The existence of \underline{D} is from (ii) of Corollary 2-1, as \underline{R} is positive definite. However (7-27) can have negative values. Thus (7-26) can also have negative values (Theorem 2-3), contradicting the hypothesis.

(2) Assume $v \geq 2$ for the canonical form. Let $\widetilde{w}^{(k)}(\underline{x}_1)$ be the highest degree, nonidentically zero, homogeneous function in (7-16).

(2a) For the case of

$$2 \leq k + 1 < 2\psi , \quad (7-29)$$

it follows from (7-14) that

$$L^{(\beta)}(\underline{x}) = -\underline{F}_2^{(\psi)}(\underline{x})^T \underline{R} \underline{F}_2^{(\psi)}(\underline{x}) . \quad (7-30)$$

Arguing as in (1), (7-26) cannot be positive semidefinite and the hypothesis fails.

(2b) For the case of

$$2 \leq 2\psi < k + 1, \quad (7-31)$$

it follows from (7-13) and (7-14) that

$$L^{(\beta)}(\underline{x}) = -\underline{\tilde{w}}^{(k)T} \underline{A}_1 \underline{x}, \quad (\beta = k + 1), \quad (7-32)$$

and

$$V^{o(k+1)}(\underline{x}) = \int_0^{\underline{x}_1} \underline{\tilde{w}}^{(k)T}(\underline{x}_1) d\underline{x}_1. \quad (7-33)$$

According to (7-33) and Lemma 6-6, (7-32) can have negative values since it is not identically zero by the assumption and the system is controllable. Thus (7-26) can have negative values (Theorem 2-3) and the hypothesis fails.

(2c) For the case of

$$2 \leq 2\psi = k + 1, \quad (7-34)$$

it follows that

$$L^{(\beta)}(\underline{x}) = -\underline{F}_2^{(\psi)T}(\underline{x}) \underline{R} \underline{F}_2^{(\psi)}(\underline{x}) - \underline{\tilde{w}}^{(k)T}(\underline{x}_1) \underline{A}_1 \underline{x} \quad (7-35)$$

and

$$\nabla^{o(k+1)}(\underline{x}_1) = \int_0^{\underline{x}_1} \underline{\tilde{w}}^{(k)T}(\underline{x}_1) d\underline{x}_1 . \quad (7-36)$$

From the results of (2a) and (2b), (7-35) can have negative values if it is not identically zero and the hypothesis fails.

To provide a quadratic form for $L(\underline{X})$ consistent with Theorem 7-1, $\underline{F}_2(\underline{X})$ must be identically zero, from (7-30). Thus for a linearly synthesized feedback control system with a completely controllable system, a nonnegative $L(\underline{X})$ is possible only if $L(\underline{X})$ is a quadratic form.

Optimal controls may be designed to minimize measures both of errors and energy during the control action. A performance index often used for these designs has the form

$$\int_{t_0}^{\infty} (\underline{x}^T \underline{Q} \underline{x} + \underline{u}^T \underline{R} \underline{u}) dt , \quad (7-37)$$

where \underline{Q} and \underline{R} are at least positive semidefinite. This choice of a performance index is due primarily to practical aspects of the problem, e.g., for mathematical convenience. From Theorem 7-1, however, this choice is seen to be

particularly appropriate if a completely controllable system is to be synthesized as a linear feedback control system.

For Kalman's inverse problem, Section 5.3.1, $L(\underline{X})$ is restricted to a quadratic form. From Theorem 7-1, it follows that no other nonnegative polynomial $L(\underline{X})$ can exist if the given system is completely controllable. In contrast, the example in Section 5.3.3, illustrating Suga's work, showed optimized nonquadratic polynomial $L(\underline{X})$. Theorem 7-1 explains why these nonquadratic polynomials are not positive semidefinite.

7.4.2 Nonnegative $L(\underline{X})$ for an Uncontrollable System

Assume that the given system (7-2) is uncontrollable. Then it follows that $v \geq 2$ and $\underline{X}_{(1)}$ represents the uncontrollable state variables governed by

$$\dot{\underline{X}}_{(1)} = \underline{A}_{(1,1)} \underline{X}_{(1)} \quad (7-38)$$

Corollary 7-7:

Provided that the given system is uncontrollable in the modified inverse problem, then it is necessary for $L^{(\beta)}(\underline{X})$ in (7-26) to be a function of only $\underline{X}_{(1)}$ if $L(\underline{X})$ is positive semidefinite for some $\beta > 0$.

Proof: For $L(\underline{X})$ to be positive semidefinite with some $\beta > 2$, $L^{(\beta)}(\underline{X})$ must be nonnegative from Theorem 2-3.

For $2\psi > k+1 \geq 2$, $L^{(\beta)}(\underline{X})$ can have negative values, as shown in (2a) of the proof of Theorem 7-1. For $k+1 \geq 2\psi \geq 2$, it follows that

$$V^{(k+1)}(\underline{X}) = \int_0^{\underline{X}_1} \underline{\tilde{w}}^{(k)T}(\underline{X}_1) d\underline{X}_1 \quad (7-39)$$

and either

$$(a) \quad L^{(\beta)}(\underline{X}) = -\underline{\tilde{w}}^{(k)T}(\underline{X}_1) \underline{A}_1 \underline{X}, \text{ if } k+1 > 2\psi, \quad (7-40)$$

or

$$(b) \quad L^{(\beta)}(\underline{X}) = -\underline{F}_2^{(\psi)T}(\underline{X}) \underline{R} \underline{F}_2^{(\psi)}(\underline{X}) - \underline{\tilde{w}}^{(k)T}(\underline{X}_1) \underline{A}_1 \underline{X},$$

$$\text{if } k+1 = 2\psi. \quad (7-41)$$

The last terms of either (7-40) or (7-41) may be positive semidefinite if it is a function of $\underline{X}_{(1)}$, from (7-39) and Lemma 6-6. Therefore for (7-41) to be positive semidefinite, $\underline{F}_2^{(\psi)}(\underline{X})$ must be a function only of $\underline{X}_{(1)}$.

For the case of the uncontrollable system, the consequences of the nonnegative property of $L(\underline{X})$ are not as definitive as those for the completely controllable system.

Chapter 8

CONCLUSIONS AND SUGGESTIONS FOR FURTHER STUDIES

8.1 Conclusions

Systems considered throughout in this study belong to the class given by

$$\dot{\underline{Y}} = \hat{\underline{F}}(\underline{Y}) + \hat{\underline{B}} \underline{V},$$

as described in Section 4.1.

8.1.1 The Canonical Form

A new canonical form is presented for this class of systems through a nonsingular transformation such that

$$\underline{Y} = \underline{N} \underline{X}$$

$$\underline{V} = \underline{M} \underline{U}$$

following the definition of an equivalent system,

Definition 4-1. This canonical form for a linear system,

i.e., $\hat{\underline{F}}(\underline{Y}) = \hat{\underline{A}} \underline{Y}$, is

$$\dot{\underline{X}} = \underline{A} \underline{X} + \underline{B} \underline{U}$$

$$= \begin{bmatrix} \underline{A}_{(1,1)} & \underline{A}_{(1,2)} & & & & \\ & & \underline{A}_{(2,3)} & & & \\ & & & \cdot & & \\ & & & & \cdot & \\ & & & & & \underline{A}_{(v-1,v)} \\ \underline{A}_{(v,1)} & \cdot & \cdot & \cdot & \cdot & \underline{A}_{(v,v)} \end{bmatrix} \begin{bmatrix} \underline{X}_{(1)} \\ \underline{X}_{(2)} \\ \cdot \\ \cdot \\ \cdot \\ \underline{X}_{(v)} \end{bmatrix}$$

$$+ \begin{bmatrix} & \\ & \underline{I}_r \end{bmatrix} \underline{U}, \text{ all other entries zero,}$$

as characterized by the following statements.

(i) Each $\underline{A}_{(i,j)}$, $(i,j = 1,2,\dots,v)$, is an $\ell_i \times \ell_j$ submatrix (Theorem 4-2).

(ii) v and ℓ_i , $(i = 1,2,\dots,v)$, for the decomposition are positive integers uniquely determined by $\hat{\underline{A}}$ and $\hat{\underline{B}}$ of the given system such that

$$\sum_{i=1}^v \ell_i = n,$$

where each ℓ_i and the ordered set $\{\ell_1, \ell_2, \dots, \ell_v\}$ are called the i^{th} stage number and the stage distribution (Theorem 4-2).

(iii) It is possible to let $v = 1$ if and only if $n = m = r$ when the system is completely controllable (Theorem 4-2).

(iv) If $v \geq 2$ and the system is completely controllable, then

$$\underline{A}_{(1,1)} = [0]$$

and

$$\underline{A}_{(1,2)} = \begin{bmatrix} [0] & I_{\ell_1} \end{bmatrix},$$

where

$$\ell_i \leq \ell_j, \quad 1 \leq i < j \leq v.$$

In addition, if $v \geq 3$, then

$$\underline{A}_{(i,i+1)} = \begin{bmatrix} [0] & I_{\ell_i} \end{bmatrix}, \quad \text{for } i = 2, 3, \dots, v-1 \text{ (Theorem 4-2).}$$

(v) If the system is uncontrollable, then $v \geq 2$,

$\underline{A}_{(1,1)}$ is given in Jordan canonical form, and $\underline{A}_{(1,2)} = [0]$.

In addition, if $v \geq 3$,

$$\underline{A}_{(i,i+1)} = \begin{bmatrix} [0] & I_{\ell_i} \end{bmatrix}, \text{ for } i = 2, 3, \dots, v-1, \text{ and}$$

$$\ell_i \leq \ell_j, \quad 2 \leq i < j \leq v \text{ (Theorem 4-2).}$$

In this case, only state variables $x_1, x_2, \dots, x_{\ell_i}$ are uncontrollable, as discussed in Section 4.3.3.

(vi) The property of controllability is invariant for the transformation to the canonical form (Theorem 4-3).

(vii) The stage distribution of the system, i.e., the ordered set $\{\ell_1, \ell_2, \dots, \ell_v\}$, is a unique characteristic (Theorem 4-4).

(viii) Provided $m = r$, the canonical transformation is possible with $\underline{M} = I_r$ if the given system has the stage distributions $\{r\}$ or $\{\ell_1, r, \dots, r\}$ (Theorem 4-6).

(ix) If the given system has the stage distribution $\{r, r, \dots, r\}$, then only one \underline{N} can exist with each possible \underline{M} for the canonical transformation (Theorem 4-5).

From these characteristics, a number of observations can be made. The structure of \underline{B} in the canonical form discloses the fundamental fact that only r independent control variables out of the m control variables contained in \underline{v} can be effective in the control action. The number

v and the structures of $\underline{A}_{(1,1)}$ and $\underline{A}_{(1,2)}$ simply identify the controllability property of the given system from (iii)-(vi) above. If a given system is completely controllable and has $\hat{\underline{B}}$ of full rank and a stage distribution $\{r, r, \dots, r\}$, then a unique canonical transformation is possible with $\underline{M} = \underline{I}_r$, from (viii) and (ix). This particular canonical transformation coincides with the more familiar canonical transformation of Definition 4-1. Furthermore a single input completely controllable system has a unique phase variable canonical form and this unique form is the familiar form proposed by earlier investigators.

For a nonlinear system, the canonical form can be applied to the linear part and the above 8 characteristics are preserved for the linear part of the transformed system.

In comparison with the other phase canonical forms described in Chapter 3, the new canonical form has the following advantages.

- (1) It can be applied to the entire class of systems given by (6-1) while other suggested canonical forms are applicable essentially to subclasses.
- (2) The many elements of the matrices \underline{A} and \underline{B} describing the linear part of the new canonical form can be reduced to units and zeros arranged in a simple and unique order. These decouple the matrices uniquely

due to the uniqueness of the stage distribution.

Thus the mathematical structure of a general class of systems is compactly presented.

This canonical form has value for simplifying studies of optimal control problems for multi-input systems. Other known canonical forms are essentially subclasses of this new one.

8.1.2 The Inverse Problem of the Optimal Regulator

The inverse problem of the optimal regulator is considered for the class of systems given by (6-1). It is shown that the problem can be equivalently considered through the new canonical form under similar mathematical assumptions and without loss of generality. The recovery of the results for the originally given system is possible by the inverse of the canonical transformation. The analysis of the problem is efficiently performed with the compact structure of the canonical form.

Restricting the form of performance indices to

$$\int_{t_0}^{\infty} \{L(\underline{X}) + \underline{U}^T \underline{R} \underline{U}\} dt$$

with positive definite \underline{R} , the necessary and sufficient

conditions for optimized performance indices corresponding to a given $\underline{U}(\underline{X})$ are presented in a theorem (Theorem 6-1). From this theorem, new aspects of optimal feedback control systems are disclosed.

(i) At most, r functions out of m in the feedback control law are independent in optimal feedback control, i.e., there are r effective control functions for the optimal control action.

(ii) Various $L(\underline{X})$ can be paired with an \underline{R} for optimized performance indices based on given $\underline{U}(\underline{X})$ (Corollary 6-1). However, if $n = m = r$, the $L(\underline{X})$ is unique for an \underline{R} and given $\underline{U}(\underline{X})$ (Corollary 6-2).

(iii) When a linear feedback control law is given, the structure of the resulting optimal performance index function $\nabla^0(\underline{X})$ is the sum of a quadratic form determined by the given feedback control law $\underline{U}(\underline{X})$ and an arbitrary function of \underline{X}_1 . State variables in \underline{X}_1 are not exposed to \underline{U} directly. The nonnegative property of these $\nabla^0(\underline{X})$ is detailed (Theorems 6-2 and 6-3).

(iv) The controllability of nonlinear systems given by (6-68) is determined by the linear part of the nonlinear systems as (6-128). For this class of systems, the Principle of Necessity of Control Action is introduced.

Specially for the problem of the optimal regulator such that: the desired final condition of the system as $\underline{X} = \underline{0}$, the performance index is given as

$$\int_{t_0}^{\infty} \{L(\underline{X}) + \underline{U}^T \underline{R} \underline{U}\} dt$$

with a positive definite \underline{R} , $L(\underline{X})$ is nonnegative definite and a nonzero function of the controllable variables, then some control action is necessary for the optimality, i.e., a feedback control law not identically zero is necessary.

(v) When a feedback control law is given, the Liapunov direct method can be applied to identify the property of asymptotic stability in the large (Theorem 6-4).

With additional assumptions to the inverse problem, the modified inverse problem allows observation of general characteristics of linearly synthesized feedback control systems with a polynomial $L(\underline{X})$.

(vi) If the given system is controllable, a nonnegative $L(\underline{X})$ can exist only as a quadratic form (Theorem 7-1).

(vii) If the given system is uncontrollable, it is necessary for a nonnegative $L(\underline{X})$ that the lowest and highest degree homogeneous functions in $L(\underline{X})$ be functions

only of uncontrollable state variables (Corollary 7-7).

Thus this statement of the inverse problem and its solution is a generalization of the problem proposed by Kalman. It is also a generalization of Suga's work excepting his time varying assumption, and of Thau's results excepting his more general assumption of $\hat{R}(V)$ in (5-67).

8.2 Suggestions for Further Studies

Throughout this work, problems are considered only under time invariant assumptions. The concepts and techniques appear to be extendable for time varying systems with some modification.

For the inverse problem, the matrix R is restricted to be positive definite to insure a normal Hamiltonian. Studies can be directed to attempt to relax this assumption, e.g., consider only positive semidefiniteness for R .

More generally, the inverse problem can be considered for basically different problem assumptions, e.g., more general systems descriptions and a different form of performance indices. Exhaustive studies of the inverse problem of the optimal regulator will reveal new characteristics of optimal feedback control systems.

REFERENCES

- [1] Bellman, R., Dynamic Programming, Princeton University Press, Princeton, N. J., 1957.
- [2] Pontryagin, L. S., et al., The Mathematical Theory of Optimal Processes, Interscience Pub. Inc., N. Y., 1962.
- [3] Leitman, G., Optimization Technique, Academic Press, New York, N. Y., 1962.
- [4] Drefus, S. E., Dynamic Programming and the Calculus of Variations, Academic Press, New York, N. Y., 1965.
- [5] Athans, M. and P. L. Falb, Optimal Control, McGraw-Hill Book Co., Inc., New York, N. Y., 1966.
- [6] Lee, E. B. and L. Markus, Foundations of Optimal Control Theory, John Wiley and Sons, Inc., New York, N. Y., 1967.
- [7] Fleming, W. H., Functions of Several Variables, Addison-Wesley Pub. Co., Inc., Reading, Mass., 1965.
- [8] Gibson, J. E., Nonlinear Automatic Control, McGraw-Hill Book Co., Inc., New York, N. Y., 1963.
- [9] Hahn, W., Theory and Applications of Liapunov's Direct Method, (English Translation), Prentice-Hall, Inc., Englewood Cliffs, N. J., 1963.
- [10] Thrill, R. M. and L. Tornheim, Vector Spaces and Matrices, John Wiley and Sons, Inc., New York, N. Y., 1957.
- [11] Ayres, F., Jr., Theory and Problems of Matrices, Schaum Pub. Co., New York, N. Y., 1962.
- [12] Kalman, R. E., "Mathematical Description of Linear Dynamical System," J. SIAM Control, Series A, Vol. 1, pp. 152-192, 1963.
- [13] Pontryagin, L. S., Ordinary Differential Equations, (English Translation), Addison-Wesley Pub. Co., Inc., Reading, Mass., 1962.

- [14] Kalman, R. E. and J. E. Bertram, "Control System Analysis and Design via the Second Method of Liapunov," Trans. ASME, Series D, J. Basic Eng., Vol. 82, pp. 371-393, 1960.
- [15] Kalman, R. E., Y. C. Ho, and K. S. Narendra, Controllability of Linear Dynamical Systems, in Contribution to Differential Equations, Vol. 1, John Wiley and Sons, Inc., New York, N. Y., 1962.
- [16] Wonham, W. M. and C. D. Johnson, "Optimal Bang-Bang Control with Quadratic Index of Performance," J. Basic Eng., Vol. 86, pp. 107-115, 1964.
- [17] Johnson, C. D. and W. M. Wonham, "A Note on the Transformation to Canonical Form," Trans. IEEE, Vol. AC-9, pp. 312-313, 1964.
- [18] Chidambara, M. R., "The Transformation to Canonical Form," Trans. IEEE, Vol. AC-10, pp. 206-207, 1965.
- [19] Tuel, W. G., "On the Transformation to Canonical Form," Trans. IEEE, Vol. AC-11, pp. 607-608, 1966.
- [20] Rane, D. S., "A Simplified Transformation to Canonical Form," Trans. IEEE, Vol. AC-11, pp. 608, 1966.
- [21] Tuel, W. G., "Canonical Forms for Linear Systems - I," IBM Research Paper RJ-375, 1966.
- [22] Luenberger, D. G., "Canonical Forms for Linear Multi-Variable Systems," Trans. IEEE, Vol. AC-12, pp. 290-293, 1967.
- [23] Asseo, S. J., "Phase-Variable Canonical Transformation of Multicontroller Systems," Trans. IEEE, Vol. AC-13, pp. 129-131, 1968.
- [24] Dugundji, J., Topology, Allyn and Bacon, Inc., Boston, Mass., 1966.
- [25] Kalman, R. E., "When is a System Optimal?", J. Basic Eng., Vol. 86, pp. 51-60, 1964.
- [26] Kalman, R. E., "Contributions to the Theory of Optimal Control," Bol. Soc. Mat. Mex., Vol. 5, pp. 102-119, 1960.

- [27] Snow, D. R., "Determining Reachable Regions and Optimal Controls," *Advances in Control Systems*, pp. 133-196, Academic Press, New York, N. Y., 1967.
- [28] Bridgman, T. F., Jr., "On the Existence of Optimal Feedback Controls," *J. SIAM, Series A, Control* Vol. 1, pp. 261-274, 1963.
- [29] Boltyanskii, V. G., "Sufficient Conditions for Optimality and the Justification of the Dynamic Programming Method," *J. SIAM, Series A, Control*, Vol. 4, No. 2, pp. 326-361, 1966.
- [30] Suga, I., "On Inverse Problem of Optimal Control," *J. Soc. Inst. and Control Engrs. (Japan)*, Vol. 6, No. 8, pp. 549-557, 1967.
- [31] Thau, F. E., "On Inverse Optimal Control Problem for a Class of Non-linear Autonomous Systems," *Trans. IEEE*, Vol. AC-12, No. 6, pp. 674-681, 1967.
- [32] Hsu, J. C. and A. U. Meyer, Modern Principles and Applications, pp. 696-697, McGraw-Hill Book Co., New York, N. Y., 1968.