

A PROCEDURE FOR THE SELECTION OF TERMS AND ESTIMATION  
OF COEFFICIENTS IN A RESPONSE SURFACE MODEL WITH  
INTEGRATION-ORTHOGONAL TERMS

by

Ronald W. Helms\*, R. J. Hader\*\*,  
and A. R. Manson\*\*

\*Department of Biostatistics  
University of North Carolina

\*\*Department of Experimental Statistics  
North Carolina State University

University of North Carolina  
Institute of Statistics Mimeo Series No. 646

✓  
August 1969

# ABSTRACT

HELMS, RONALD WILLIAM. A Procedure for the Selection of Terms and Estimation of Coefficients in a Response Surface Model with Integration-Orthogonal Terms. (Under the direction of ROBERT JOHN HADER and ALLISON RAY MANSON.)

General linear approximation theory and general linear estimation theory are combined to produce a generalization of the Minimum Bias Estimator (Karson, et al., 1969) in which a linear model is transformed to an equivalent representation in terms of integration-orthogonal functions. The integrated mean square error of an estimator  $\check{\eta}$  for a response function  $\eta$  is shown to have the form  $IMSE(\check{\eta}) = \sum E(\check{\beta}_j - \beta_j)^2$ , where  $\check{\beta}_j$  is the estimator of the coefficient  $\beta_j$ . One can thus consider the deletion or inclusion of terms in an estimation model one at a time. A term can be "deleted" from a model by setting the corresponding coefficient estimator to zero ( $\check{\beta}_j = 0$ ); various procedures are considered for using the least squares estimator  $\check{\beta}_j = \hat{\beta}_j$  for certain sets of  $\hat{\beta}_j$ -values and zero (deletion of the term;  $\check{\beta}_j = 0$ ) for other sets of  $\hat{\beta}_j$ -values. The expected value, bias, and mean square error are derived for each estimator. Techniques are given for the selection of the region on which  $\check{\beta}_j = 0$  when prior information is available about  $\beta_j$  in the form of a prior distribution.

## BIOGRAPHY

The author was born October 30, 1941, in Charlotte, North Carolina. He was reared primarily in Charlotte and in Oak Ridge, Tennessee. He graduated from Oak Ridge High School.

He received the Bachelor of Arts degree with a major in Mathematics from the University of Tennessee in 1963. After working as a laboratory assistant and programmer at the Atomic Energy Commission Agricultural Research Laboratory in Oak Ridge, he accepted a graduate assistantship at the Computation Center of the University of Tennessee. He received the Master of Arts degree in Mathematics in 1966. In September, 1966, he enrolled in the Graduate School of North Carolina State University. After two semesters of study he was awarded a National Aeronautics and Space Administration traineeship. He was inducted into Phi Kappa Phi Honor Society in 1968. He accepted a position as instructor with the Biostatistics Department of the University of North Carolina in August, 1968, and has accepted the position of Assistant Professor effective autumn 1969.

The author married Miss Mary Irene Wagner, of Knoxville, Tennessee, in 1963; they have no children.

## ACKNOWLEDGMENTS

The author wishes to express his appreciation to all who have contributed assistance in the preparation of this study. Appreciation is particularly extended to Dr. R. J. Hader and Dr. A. R. Manson for their continued guidance and to the other members of the Advisory Committee, Professors J. W. Bishir, C. P. Quesenberry, and R. G. D. Steel, for their aid. A special note of appreciation is extended to Dr. James E. Grizzle of the University of North Carolina for continual moral and intuitive support throughout the course of this study.

The computer time required for this study was granted by the Computation Center of North Carolina State University, which is supported in part by a grant from the National Science Foundation.

To Mrs. Gay Goss and Mrs. Mary Donelan, who had the unenviable task of typing final and rough drafts, respectively, of this thesis, the author expresses his special thanks for a job well done. Appreciation is also given to Mrs. Mary Flanagan who typed much of the rough draft and helped with the figures.

Finally, it is not possible for the author to express adequately his appreciation for the continual support of his wife, Mary.

## TABLE OF CONTENTS

	Page
LIST OF TABLES . . . . .	vi
LIST OF FIGURES. . . . .	vii
1. INTRODUCTION . . . . .	1
2. REVIEW OF LITERATURE . . . . .	4
3. MINIMUM BIAS ESTIMATION WITH INTEGRATION-ORTHOGONAL FUNCTIONS. . . . .	11
3.1 Least-Squares Linear Approximation with Orthogonal Functions . . . . .	12
3.1.1 Some Basic Theory from Linear Algebra. . . . .	12
3.1.2 The Approximation Theorem. . . . .	16
3.1.3 Some Results for Integration-Orthogonal Polynomials. . . . .	22
3.2 Minimum Bias Estimation: Best Linear Unbiased Estimation of the Best Approximating Function . . . . .	31
3.2.1 Definitions and Notation . . . . .	31
3.2.2 Best Linear Unbiased Estimation of the Best Approximating Function. . . . .	36
3.2.3 Equivalence of the Minimum Bias Estimator and BLUE of the Best Approximating Function . . . . .	37
4. A PROCEDURE FOR SELECTION OF TERMS IN THE MODEL WITH MINIMUM BIAS ESTIMATION AND THE INTEGRATED MEAN SQUARE ERROR CRITERION . . . . .	41
4.1 Minimization of Integrated Mean Square Error with Respect to Choice of Terms in the Model, Parameters Known. . . . .	43
4.1.1 Experimental Model. . . . .	45
4.1.2 An Example . . . . .	46
4.1.3 Results Useful for Evaluating Integrals. . . . .	50

	Page
4.2 Properties of the Estimators Based on Two-Tail Tests . . . . .	57
4.2.1 Properties with $\sigma^2$ Known . . . . .	57
4.2.2 Properties with $\sigma^2$ Unknown . . . . .	64
4.3 Estimators Based on One-Tail Tests. . . . .	81
4.3.1 Properties of the One-Tail Estimators When $\sigma^2$ is Known. . . . .	84
4.3.2 Properties of One-Tail Estimators When Variance is Unknown. . . . .	85
4.4 Generalized Estimators. . . . .	90
5. SOME TECHNIQUES FOR THE SELECTION OF CUTOFF POINTS . . . .	112
5.1 The Bayes Decision Procedure. . . . .	114
5.2 The Posterior $C=1$ Procedure . . . . .	115
5.3 Cutoff Point Selection with a Normal Prior. . . . .	116
5.4 Cutoff Point Selection with a Uniform Prior . . . . .	122
6. SUMMARY. . . . .	125
7. LIST OF REFERENCES . . . . .	130

## LIST OF TABLES

Page

5.1 A comparison of the expected values, with respect to a $N(\mu, t^2)$ prior, of the Mean Square Error function for the Bayes decision procedure and the Posterior $C=1$ Procedure . . . . .	120
--	-----

## LIST OF FIGURES

	Page
4.1 Mean Square Error of the two-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ and $C$ ; $\sigma^2$ assumed known . . . . .	61
4.2 Absolute value of the bias of the two-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ and $C$ ; $\sigma^2$ assumed known . . . . .	63
4.3 Mean Square Error of the two-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 2$ , $\sigma^2$ estimated . . . . .	71
4.4 Mean Square Error of the two-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 5$ , $\sigma^2$ estimated . . . . .	72
4.5 Mean Square Error of the two-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 10$ , $\sigma^2$ estimated . . . . .	73
4.6 Mean Square Error of the two-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 25$ , $\sigma^2$ estimated . . . . .	74
4.7 Mean Square Error of the two-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 50$ , $\sigma^2$ estimated . . . . .	75
4.8 Absolute value of the bias of the two-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 2$ , $\sigma^2$ estimated . . . . .	76
4.9 Absolute value of the bias of the two-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 5$ , $\sigma^2$ estimated . . . . .	77
4.10 Absolute value of the bias of the two-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 10$ , $\sigma^2$ estimated . . . . .	78
4.11 Absolute value of the bias of the two-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 25$ , $\sigma^2$ estimated . . . . .	79
4.12 Absolute value of the bias of the two-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 50$ , $\sigma^2$ estimated . . . . .	80



LIST OF FIGURES (continued)	Page
4.13 Mean Square Error of the one-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ and $C$ ; $\sigma^2$ assumed known . . . . .	86
4.14 Absolute value of the bias of the one-tail estimator $\theta = \beta/s.d.(\hat{\beta})$ as a function of $\theta$ and $C$ ; $\sigma^2$ assumed known . . . . .	87
4.15 Mean Square Error of the one-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 2$ , $\sigma^2$ estimated . . . . .	91
4.16 Mean Square Error of the one-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 5$ , $\sigma^2$ estimated . . . . .	92
4.17 Mean Square Error of the one-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 10$ , $\sigma^2$ estimated . . . . .	93
4.18 Mean Square Error of the one-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 25$ , $\sigma^2$ estimated . . . . .	94
4.19 Mean Square Error of the one-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 50$ , $\sigma^2$ estimated . . . . .	95
4.20 Absolute value of the bias of the one-tail estimator $\theta = \beta/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 2$ , $\sigma^2$ estimated . . . . .	96
4.21 Absolute value of the bias of the one-tail estimator $\theta = \beta/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 5$ , $\sigma^2$ estimated . . . . .	97
4.22 Absolute value of the bias of the one-tail estimator $\theta = \beta/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 10$ , $\sigma^2$ estimated . . . . .	98
4.23 Absolute value of the bias of the one-tail estimator $\theta = \beta/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 25$ , $\sigma^2$ estimated . . . . .	99
4.24 Absolute value of the bias of the one-tail estimator $\check{\theta} = \check{\beta}/s.d.(\hat{\beta})$ as a function of $\theta$ , $C$ and $\nu = 50$ , $\sigma^2$ estimated . . . . .	100
4.25 Mean Square Error of the generalized estimator $\check{\theta} = 0 \cdot 1_A(\hat{\theta}) + \hat{\theta} \cdot 1_B(\hat{\theta})$ , variance assumed known . . . . .	106

LIST OF FIGURES (continued)	Page
4.26 Mean Square Error of the generalized estimator $\hat{\theta} = 0 \ 1_A(\hat{\theta}) + \hat{\theta} \ 1_B(\hat{\theta})$ variance assumed known . . . . .	107
4.27 Mean Square Error of the generalized estimator $\hat{\theta} = 0 \ 1_A(\hat{\theta}) + \hat{\theta} \ 1_B(\hat{\theta})$ variance assumed known . . . . .	108
4.28 Absolute value of the bias of the generalized estimator $\hat{\theta} = 0 \ 1_A(\hat{\theta}) + \hat{\theta} \ 1_B(\hat{\theta})$ variance assumed known . . . . .	109
4.29 Absolute value of the bias of the generalized estimator $\hat{\theta} = 0 \ 1_A(\hat{\theta}) + \hat{\theta} \ 1_B(\hat{\theta})$ variance assumed known . . . . .	110
4.30 Absolute value of the bias of the generalized estimator $\hat{\theta} = 0 \ 1_A(\hat{\theta}) + \hat{\theta} \ 1_B(\hat{\theta})$ variance assumed known . . . . .	111
5.1 Configurations of $q(\hat{\theta})$ for various combinations of prior distribution parameters . . . . .	118

## 1. INTRODUCTION

In many experimental situations it is convenient to suppose that a functional relationship,

$$\eta = g(x_1, x_2, \dots, x_m; \theta_1, \dots, \theta_p) = g(\underline{x}, \underline{\theta}), \quad (1.1)$$

exists between a response  $\eta$  and  $m$  continuous variables  $x_1, \dots, x_m$ . When the form of the function (1.1) is known, the object of the experiment may be to estimate the parameter vector,  $\underline{\theta}$ . In many cases, however, the form of the function is unknown, or, if known, may not be linear in the parameters. In such cases the objective may be to estimate an approximation of the function  $g(\underline{x}, \underline{\theta})$  over some given region  $R$  of the  $m$ -dimensional space of the  $\underline{x}$ -variables.

If the response function is linear in the parameters,

$$\eta = g(\underline{x}; \underline{\theta}) = \sum_{i=1}^m \theta_i f_i(\underline{x}) \quad (1.2)$$

where the  $f_i$  are known functions (polynomials, for example), then, under certain assumptions about the nature of the experimental errors, it is possible to apply traditional general linear model (least squares) estimation theory to obtain Best (minimum variance) Linear Unbiased Estimators for the parameters  $\underline{\theta}$  and the function  $\eta$ . However, if the assumed model (1.2) is inadequate, in the sense that the true model is of the form

$$\eta = g(\underline{x}; \underline{\theta}) = \sum_{i=1}^{m+q} \theta_i f_i(\underline{x}) \quad (1.3)$$

then the estimators for the parameters  $\theta$  and for the function are biased. In the presence of bias the use of the minimum variance criterion for the selection of an estimation procedure is subject to question.

Since bias attributable to an inadequate model is usually not constant over the region of interest,  $R$ , the mean square error of an estimator,  $\hat{\eta}$ , integrated over the region  $R$ , viz.,

$$\text{IMSE}(\hat{\eta}) = \int_R E \{ [\eta(\underline{x}) - \hat{\eta}(\underline{x})]^2 \} d\underline{x} \quad (1.4)$$

would appear to be a good criterion for use in the selection of an estimator. In this thesis the IMSE is used as a criterion for the comparison of estimators in inadequate model situations.

It is useful to notice that the IMSE can be decomposed into bias and variance components:

$$\text{IMSE}(\hat{\eta}) = \int_R \{ E[\hat{\eta}(\underline{x})] - \eta(\underline{x}) \}^2 d\underline{x} + \int_R \text{Var}[\hat{\eta}(\underline{x})] d\underline{x} \quad (1.5)$$

Box and Draper (1959) showed that in considering the selection of a design, the bias contribution to integrated mean square error generally far outweighed the variance contribution. Motivated by this result, Karson et al. (1969) derived a Minimum Bias Estimator for  $\eta$  by the technique of first considering a class of estimators which minimize the bias contribution in (1.5) and then selecting the estimator from that class which has minimum variance. From (1.5) one can see that  $E\{\hat{\eta}(\underline{x})\}$  is an approximating function for  $\eta$ . Thus, the minimum bias estimator is the minimum variance estimator of the minimum bias

approximating function. Linear approximation with respect to the integrated squared error criterion is considered in Section 3 of this thesis. Some underlying connections are established between least squares estimation based on integration-orthogonal functions and the minimum bias estimator.

An important problem in minimum bias estimation and in many other regression problems is the selection of the approximating function, i.e., selection of the independent variables or terms to be included in the analysis. For the case in which one is using minimum bias estimation based on integration-orthogonal functions ("independent variables"), a particularly simple term-by-term procedure can be used. This procedure and its properties are derived and discussed in Section 4. Techniques for selection of the parameters of the estimation procedure are discussed in Section 5.

## 2. REVIEW OF LITERATURE

Box and Draper (1959) motivated the study of minimum bias estimation. They investigated the problem of finding optimal designs for least squares estimation in the following setting. It is desired to fit a response function  $\eta(\underline{x})$ , assumed to be a polynomial of degree  $d_2$  over the region of interest  $R$ , by a polynomial  $\hat{y}(\underline{x})$  of degree  $d_1 < d_2$ . The polynomial  $\hat{y}(\underline{x})$  is a least squares estimator. Since  $\hat{y}(\underline{x})$  is of lower degree than  $\eta(\underline{x})$ ,  $E(\hat{y}(\underline{x})) \neq \eta(\underline{x})$  and both variance error (due to sampling error) and bias error must be considered in the selection of an estimation procedure and the construction of designs. Define the Integrated Mean Square Error of  $\hat{y}$  as

$$\begin{aligned} \text{IMSE}(\hat{y}) &= \int_R E\{\hat{y}(\underline{x}) - \eta(\underline{x})\}^2 d\underline{x} \\ &= \int_R \text{Var}[\hat{y}(\underline{x})] d\underline{x} + \int_R \{E[\hat{y}(\underline{x}) - \eta(\underline{x})]\}^2 d\underline{x} \\ &= V + B. \end{aligned}$$

(The  $V$  and  $B$  above differ from Box and Draper's  $V$ ,  $B$  by a normalizing scale factor.) In their paper Box and Draper considered designs for the situation in which  $\eta(\underline{x})$  is a second degree polynomial and  $\hat{y}(\underline{x})$  is a first degree polynomial. In a later paper, Box and Draper (1963), they extended their results to fitting a quadratic polynomial,  $\hat{y}(\underline{x})$ , to a cubic response,  $\eta(\underline{x})$ . In both studies they found that the

contribution to IMSE due to bias (B, above) far outweighed the contribution due to variance (V, above), in the sense that (Box and Draper, 1959, p. 622) "the optimal design in typical situations in which both variance and bias occur is very nearly the same as would be obtained if variance were ignored completely and the experiment designed so as to minimize bias alone." (The emphasis is theirs.)

Karson, et al. (1967), questioned the use of minimum variance as a criterion for choice of a biased estimator. In the Box and Draper setting the least squares estimators are certainly biased, although one can construct designs to minimize the bias contribution to IMSE. Karson, et al. (1969), took the approach of considering design and estimation as two phases of the problem of producing estimators with small IMSE. Recognizing from the work of Box and Draper (1959, 1963) that the bias contribution to IMSE exerts much more influence on the choice of design than does the variance contribution, Karson et al. (1969) derived a class of estimators which minimize the bias contribution to IMSE, and from that class selected the estimator with smallest variance. The resulting estimator is called the Minimum Bias Estimator (MBE) for  $\eta$ . Since the MBE minimizes the bias contribution to IMSE for any design, one is free to choose designs which minimize the variance contribution to IMSE. Karson, et al. (1969) illustrated this technique by demonstrating optimum designs for fitting a linear polynomial (in one variable) to a quadratic over a three point design, to a quadratic over a four point design, and to a cubic model over a four point design.

Karson, et al. (1969), assumed throughout their study the point of view of choosing in advance the terms to be included in the fitted model,  $\hat{y}(x)$ . Specifically, they assumed that all terms of degree less than or equal to  $d_1$  (which is specified in advance) will be included in the fitted model. No consideration was given to the problem of choosing the terms in the fitted model as a result of tests based on the data.

However, considerable work has been performed on the problem of selection of terms in the model in the general linear model setup (least squares estimation) and with problems involving an ordering of the terms in the model (as in a polynomial, for example).

Draper and Smith (1967, Chapter 6) discuss the following six techniques "in current use" for the selection of the "best fit": (1) comparison of the  $R^2$  (coefficient of determination) for all possible regressions, (2) backward elimination, (3) forward selection, (4) stepwise regression, (5) two variations on the four methods above, and (6) stagewise regression. The expected values and mean square errors of the coefficient estimators and function estimators ( $\hat{y}$ ) produced by these techniques apparently have not been investigated except in special cases.

Bancroft (1944) produced an expression for the bias of the estimator  $b^*$  for  $\beta_1$  under the assumption that the full (true) model is

$$y = \beta_1 x_1 + \beta_2 x_2 + e,$$

where  $b^*$  is obtained by the following procedure:

- (1) Obtain the usual least squares estimators  $(\hat{\beta}_1, \hat{\beta}_2)$  for  $(\beta_1, \beta_2)$ ;



- (2) Use an F-test to test the hypothesis  $H: \beta_2 = 0$ ; vs. the alternative  $K: \beta_2 \neq 0$ .
- (3) If the F-test is significant at the predetermined  $\alpha$ -level,  $b^* = \beta_1$ ; otherwise fit the model  $Y = \beta_1 x_1 + e$  to the data and  $b^*$  is the resulting least squares estimator for  $\beta_1$ .

Bancroft did not consider estimators based on one-tail tests, bias of the resulting estimator for  $\beta_2$ , mean square errors or variances, or integrated mean square errors.

Larson and Bancroft (1936b) considered an extension of the above model, viz.,

$$y = x_1 \beta_1 + x_2 \beta_2 + e$$

with essentially the same procedure as above, with  $\beta_1$  replaced by the vector  $\beta_i$ ;  $i = 1, 2$ . In this situation one performs an F-test of the compound hypothesis  $H: \beta_2 = 0$ . They assumed that the elements of the least squares vector estimator,  $\hat{\beta}$ , were uncorrelated. Thus,  $\beta_1^*$  is  $\hat{\beta}_1$ ,  $\beta_2^*$  is either  $\hat{\beta}_2$  or  $0$ , depending on the outcome of the test. The particular elements of  $x_1$  and  $x_2$  are specified in advance. They derived expressions for the bias and mean square error of

$$y^* = x_1 \beta_1^* + x_2 \beta_2^*$$

under the assumption that the errors,  $e$ , are normally and independently distributed with zero mean and variance  $\sigma^2$ . The derivation includes a derivation of the bias and mean square error of the elements of  $\beta_2^*$ . In a companion paper (Larson and Bancroft, 1963a) they considered successive tests of the hypotheses  $H_j: \beta_j = 0$ ,  $j = 1, 2, \dots$ , stopping at the first hypothesis rejected.

Gorman and Toman (1966), Hocking and Leslie (1967), and Schatzoff, et al., (1968) have presented several efficient methods for calculating and comparing subsets of all possible regressions. All these techniques are based on choosing a small subset of the independent variables which produces a great reduction in the residual sum of squares. Since the actual selection procedure is somewhat imprecise in each case, the properties of the resultant estimators are difficult to investigate.

Toro and Wallace (1968) considered the comparison of the usual least squares estimator,  $\hat{\beta}$ , for  $\beta$  in the general linear model (Graybill, 1961) with  $\beta^*$ , the least squares estimator subject to the linear constraints  $H\beta^* = h$ . They suggested as a final estimator whichever one of  $\hat{\beta}$ ,  $\beta^*$  produces the smallest mean square error, i.e., if  $E\{(\hat{\beta} - \beta)(\hat{\beta} - \beta)'\} - E\{(\beta^* - \beta)(\beta^* - \beta)'\}$  is positive definite, use  $\beta^*$  as the estimator. They also derived the UMP test of the hypothesis that the above condition holds. However, they did not derive the properties of the resulting estimator.

A number of procedures have been proposed for problems in which there is a natural ordering of the coefficients, as in polynomials of a single variable, for example. Graybill (1961) discussed a procedure for finding the degree of a polynomial which "describes" a set of data, and proposed the use of summation-orthogonal polynomials (which yield uncorrelated coefficients) as a labor saving technique. He did not discuss the properties of the resultant estimators. Hoel (1968) proposed a sequential procedure for determining the degree of a polynomial, when the model is written in terms of summation-orthogonal

polynomials. A feature of Hoel's procedure (and other sequential decision procedures) is that it allows control of the probability of both types of error in hypothesis testing.

Sclove (1968) discussed estimators which have uniformly smaller mean square error than the usual least squares estimator when there are more than three terms in the model. He also discussed a stepwise procedure for successively testing the hypotheses  $H_j: \beta_j = 0$ ,  $j = p, p-1, \dots$ , (until one hypothesis is rejected), for the case in which the terms are ordered, as in polynomial models in one variable. Although he referred to proofs that the proposed estimators have uniformly smaller mean square error than least squares estimators, he did not evaluate the mean square errors, noting that (Sclove, 1968, p. 599), "Computation of [the mean square error] seems difficult."

Each of the above mentioned procedures for selection of terms to be included in the model uses a point-wise criterion to guide the selection. The criteria are "point-wise" or "local" in the sense that they are based solely on information obtained at the experimental points; no consideration is given to how well the resulting function estimator fits the true function over the region of interest. In contrast, the IMSE criterion may be referred to as a "global" criterion. In addition, the properties of some of the resultant estimators above are difficult to evaluate; some of the estimators are difficult to evaluate in comparison to the least squares estimator.

Michaels (1969) considered the problems of estimation and design for linear versus quadratic estimation (for an assumed quadratic model) for univariate polynomials over the scaled region of interest,

$R = [-1, 1]$ , using the integrated mean square error criterion. Considering the expectation of the IMSE function with respect to the prior density of  $\beta_2$ , the coefficient of  $x^2$ , he found that for normal or uniform prior densities one would optimally choose between least-squares linear estimation or least-squares quadratic (full model) estimation by comparing  $E[\beta_2^2]$  (with respect to the prior) and  $\sigma_2^2$ , the variance of  $\hat{\beta}_2$ . Minimum Bias Estimation of the linear model was considered, but for all  $\beta_2$  values the MBE was dominated by either the linear or quadratic (full model) least squares procedure. Michaels (1969) also investigated the selection of an optimum design-estimation procedure with respect to the IMSE criterion and various prior distributions for  $\beta_2$ .

In Chapter 4 of this paper a procedure is presented which allows one to select terms in the model with respect to the IMSE criterion. All possible "regressions" are considered; the calculations are straightforward (only the least squares estimator and estimated variance need be calculated), and the bias, mean square error and IMSE of the estimates are derived and computed.

Moreover, the resulting estimator for the response function,  $\eta$ , is, conditional upon the model chosen, the Minimum Bias Estimator for  $\eta$ .

Chapter 3 of this paper presents a generalization of the Minimum Bias Estimation scheme for quite general weighting functions, over a general region of interest, and for any finite set of continuous, linearly independent functions in several variables.

### 3. MINIMUM BIAS ESTIMATION WITH INTEGRATION-ORTHOGONAL FUNCTIONS

The material in this chapter is mostly expository in nature; the theory of linear approximation and the theory of linear estimation are standard material in the fields of numerical analysis and statistics, respectively. A summary of some pertinent sections of the two bodies of theory is presented here because the combination apparently does not appear elsewhere and because the combined results lead to an elegant and useful description of minimum bias estimation.

In the first subsection some results are presented for vector spaces with inner products. Examples are given for finite dimensional Euclidean vector spaces,  $E^n$ , and for infinite dimensional vector spaces of real, continuous, square-integrable functions. Orthogonal functions are defined for each of these spaces, and are illustrated with orthogonal polynomials in several variables. Since the primary application of this theory involves the use of integration-orthogonal polynomials, some properties are derived for these functions, including algorithms for their construction and for conversions from standard polynomial models to and from orthogonal polynomial models.

In the second subsection, a model is defined for a response function as a linear combination of orthogonal functions (polynomials, for example). Standard least squares estimation theory is applied to find estimators for the response function and for "best" approximations to the response function. The equivalence of minimum bias estimation and least squares estimation of the best approximating function is

established, and some of the properties of minimum bias estimators which follow from this equivalence are given.

### 3.1. Least-Squares Linear Approximation with Orthogonal Functions

In this section approximation of functions which are linear in the parameters is considered. The sense in which the approximation is "least-squares" will be explained in the development of the theory.

Orthogonal functions, particularly orthogonal polynomials in one variable, are basic tools in numerical analysis and in linear approximation theory in particular. Treatments of least-squares linear approximation of a function of one variable are given in several introductory numerical analysis texts, such as Todd (1962), Hildebrand (1956) and Handscomb (1966). The present treatment will assume the function being approximated may be a function of several variables.

In introductory texts the region,  $R$ , over which the function is to be approximated is usually assumed to be the interval  $[-1,1]$  for univariate approximations and the "unit cube" (the same interval in each variable) for multivariate approximations. Our assumptions about  $R$  are somewhat more general.

#### 3.1.1. Some Basic Theory from Linear Algebra

Throughout this section general vectors (unspecified vector space) will be denoted by lower case Latin letters, with or without subscripts; most will be from the end of the alphabet ( $x, y, z$ ). Scalars will also be lower case Latin letters, with or without subscripts, mostly from the first of the alphabet ( $a, b, c$ ). Vectors from finite dimensional vector spaces, such as  $E^n$ , will be

underscored (x, z); the letters f, g, h, p and q with or without subscripts, will denote vectors from infinite dimensional vector spaces (i.e., functions), and will not be underscored. The letters i, j, k, m and n will be used for sub- and super-scripts and for limits of summation.

Definition: Inner product. Let  $V$  denote a finite or infinite dimensional real vector space. An inner product is a function, denoted  $(\ , \ )$ , from  $V \times V$  to  $E^1$ , the real line, such that for all  $x, y, z$  in  $V$  and for all real constants  $a$ , the following four conditions hold:

- (i)  $(x, y) = (y, x)$ ;
- (ii)  $(x, x) \geq 0$ , and  $= 0$  only if  $x = 0$  (the zero vector in  $V$ );
- (iii)  $(ax, y) = a(x, y)$ ;
- (iv)  $(x+y, z) = (x, z) + (y, z)$ .

Corollary: Let  $x_1, \dots, x_m, y_1, \dots, y_n$  all be vectors in  $V$ , and let  $a_1, \dots, a_m, b_1, \dots, b_n$  be real constants. Then

$$\left( \sum_{i=1}^m a_i x_i, \sum_{j=1}^n b_j y_j \right) = \sum_{i=1}^m \sum_{j=1}^n a_i b_j (x_i, y_j). \quad (3.1)$$

Definition: Norm based on an inner product. With the setting above, define the norm of the vector  $x$ , denoted  $||x||$ , by:

$$||x|| = \sqrt{(x, x)}, \text{ or } ||x||^2 = (x, x). \quad (3.2)$$

It should be noted that there are norms which are not based on inner products.

Definition: Orthogonal, orthonormal vectors. Two vectors  $x, y$  in  $V$  are said to be orthogonal if  $(x, y) = 0$ , and are said to be orthonormal if they are orthogonal and  $(x, x) = \|x\|^2 = (y, y) = \|y\|^2 = 1$ . A set of vectors is said to be an orthogonal set if all pairs of vectors in the set are orthogonal; it is said to be an orthonormal set of vectors if it is an orthogonal set and the norm of each vector in the set is one.

Consider two rather different vector spaces and inner products. Since both will be used in the following sections, the  $(x, y)$  notation for these special inner products will be altered slightly.

Let  $V = E^n = n$ -dimensional Euclidean space over the reals. Let  $H$  be a real, symmetric,  $n \times n$  positive definite matrix. For all  $\underline{x}, \underline{y}$  in  $V$ , define and denote the inner product:

$$S_H(\underline{x}, \underline{y}) = (\underline{x}, \underline{y}) = \underline{x}' H \underline{y} = \sum_{i=1}^n \sum_{j=1}^n x_i h_{ij} y_j. \quad (3.3)$$

The  $S$  refers to summation and the  $H$  denotes the matrix of the inner product. A special case is for  $H = I$ , the  $n \times n$  identity matrix. In this case,

$$S_I(\underline{x}, \underline{y}) = \underline{x}' I \underline{y} = \underline{x}' \underline{y}$$

This is the "usual" inner product (also referred to as the "dot product") over  $E^n$ . Returning to the general case,  $\underline{x}$  and  $\underline{y}$  in  $E^n$  are orthogonal with respect to the inner product  $S_H$  if

$$S_H(\underline{x}, \underline{y}) = \underline{x}' H \underline{y} = 0.$$



The norm of a vector  $\underline{x}$  in  $E^n$  with respect to the inner product  $S_H$  is

$$\sqrt{S_H(\underline{x}, \underline{x})} = \sqrt{\underline{x}' H \underline{x}} = ||\underline{x}||_{S(H)}.$$

Notice the special notation for the norm with respect to  $S_H$ .

For the second example, let  $R$  be some non empty subset of  $E^n$ .

$R$  does not denote the real line except as a special case. Let  $\underline{x}$  denote a vector in  $E^n$ , and let  $W(\underline{x})$  denote either a finite measure over the Borel subsets ( $B$ ) of  $R$  or a distribution function of finite variation defined over  $R$ . (See Halmos (1950) and Loeve (1955) for discussions of finite measures, distribution functions of finite variation, and integration with respect to such measures and distribution functions.) Let  $V$  be the space of continuous functions which are square-integrable with respect to  $W$  over  $R$ , i.e.,  $f$  is in  $V$  if  $f$  is continuous and

$$0 \leq \int_R f^2(\underline{x}) dW(\underline{x}) < + \infty. \quad (3.4)$$

Denote  $V$  by  $V = L_2(R, B, W) = L_2$ . Functions which are equal almost everywhere with respect to  $W$  are considered to be identical, i.e., if  $f(\underline{x}) = g(\underline{x})$  except on a set  $B$ , such that  $W(B) = 0$ , then  $f$  and  $g$  are considered to be the same function. Under these conditions, the function  $I_W$  defined by:

$$I_W(f, g) = \int_R f(\underline{x}) g(\underline{x}) dW(\underline{x}) \quad (3.5)$$

for all  $f, g$  in  $L_2$  is an inner product over  $L_2$  (Loeve, 1955). Here, the  $I$  denotes integration (to distinguish from the summation inner product) and  $W$  denotes the measure.

There are two special cases of this inner product which are very useful. First, let  $w(\underline{x})$  be the Radon-Nikodym derivative of  $W(\underline{x})$  with respect to Lebesgue measure;  $w(\underline{x})$  is a finite, non-negative, integrable weight function over  $R$ . Then (3.5) becomes

$$I_W(f, g) = \int_R f(\underline{x}) g(\underline{x}) w(\underline{x}) d\underline{x}, \quad (3.6)$$

which is a multidimensional integral over the set  $R$ . This is the most common case and the one which will be applied most often.

The second special case of interest arises when  $W$  is a discrete measure which assigns the weight  $w_i$  to the point  $\underline{x}_i$ ,  $i = 1, 2, \dots, m$ . In such a case,

$$I_W(f, g) = \sum_{i=1}^m f(\underline{x}_i) g(\underline{x}_i) w_i.$$

This is, of course, just a special case of (3.3).

In the general case of  $V = L_2(R, B, W)$ , the norm will be denoted:

$$||f||_{I(W)} = \sqrt{I_W(f, f)}.$$

### 3.1.2. The Approximation Theorem

Theorem 3.1. Let  $V$  be a real vector space with inner product  $(, )$  and associated norm  $|| \cdot ||$ . Let  $\{x_1, x_2, \dots, x_n\}$  be a fixed, finite set of orthonormal vectors from  $V$ . Then for any  $y$  in  $V$ , the minimum of

$$||y - \sum_{i=1}^n b_i x_i|| \quad (3.7)$$

is attained if and only if

$$b_i = a_i \equiv (y, x_i). \quad (3.8)$$

The number  $a_i$  is called the Fourier coefficient of  $x_i$  and  $y$  with respect to the inner product  $(\cdot, \cdot)$ .

Proof. The quantity (3.7) is minimized if and only if the following is minimized:

$$\|y - \sum_{i=1}^n b_i x_i\|^2 = (y - \sum b_i x_i, y - \sum b_i x_i) \quad (3.9)$$

$$= (y, y) - 2\sum b_i (y, x_i) + \sum b_i^2 (x_i, x_i)$$

$$= (y, y) - 2\sum a_i b_i + \sum b_i^2$$

$$= \|y\|^2 + \sum (a_i - b_i)^2 - \sum a_i^2 \quad (3.9a)$$

All summations are over the range  $i = 1, 2, \dots, n$ . In (3.9a) only the term  $\sum (a_i - b_i)^2$  depends upon the  $b_i$ . Being a sum of real squares, it takes on its minimum value of zero if and only if all  $b_i = a_i$ ,  $i = 1, 2, \dots, n$ .

Most of the theory of linear approximation and of linear estimation can be derived from this theorem. The result is now extended to include the case of orthogonal (not necessarily orthonormal) vectors. These results are given in many linear algebra and numerical analysis texts.

The approximation

$$\tilde{y} = \sum_{i=1}^n a_i x_i$$

is called a "least-squares" approximation because it minimizes the square of the norm (3.9). In most applications, since the  $x_i$  are orthogonal the square of the norm is either a "sum of squares" or the integral of the squared error. Throughout the remainder of this thesis, the least squares approximator of a vector (or function) will be denoted by the use of a tilde ( $\sim$ ) e.g.  $\tilde{y}$  is the least squares approximator of  $y$ .

Corollary 1. The Approximation Theorem for Orthogonal Vectors.

Let  $\{z_1, z_2, \dots, z_n\}$  be a finite set of non-zero orthogonal vectors from the vector space  $V$ , with inner product  $(\cdot, \cdot)$  and associated norm,  $|| \cdot ||$ . Then for any  $y$  in  $V$ , the minimum of

$$||y - \sum_{j=1}^n c_j z_j|| \quad (3.10)$$

is attained if and only if

$$c_j = \frac{(y, z_j)}{(z_j, z_j)} = \frac{(y, z_j)}{||z_j||^2}, \quad j = 1, 2, \dots, n. \quad (3.11)$$

Proof. The set of vectors  $x_i$ ,  $i = 1, 2, \dots, n$ , defined by

$$x_i = \frac{1}{||z_i||} z_i, \quad i = 1, 2, \dots, n.$$

is clearly an orthonormal set. Set  $c_i = a_i / ||z_i||$ , where  $a_i = (y, x_i)$  is the Fourier coefficient of  $x_i$  and  $y$ ; then  $a_i x_i = c_i z_i$  and the minimum of

$$||y - \sum a_i x_i|| \equiv ||y - \sum c_i z_i||$$

is attained. But

$$c_i = \frac{a_i}{||z_i||} = (y, x_i) / ||z_i||$$

$$= \left( \frac{y, z_i}{||z_i||} \right) \frac{1}{||z_i||} = \frac{(y, z_i)}{||z_i||^2},$$

as was to be shown.

The following corollary establishes a result which will be useful in consideration of minimum bias estimation.

Corollary 2. The Truncation Approximation Theorem. Let

$Q = \{x_j, j = 1, 2, \dots, m\}$  be a fixed orthogonal set of nonzero vectors from  $V$  and define

$$y = \sum_{j=1}^m c_j x_j \quad (3.12)$$

where the  $c_j$  are also fixed. Of course,  $y \in V$ . Now approximate  $y$  by a subset of the vectors in  $Q$ . Let  $n$  be such that  $0 < n < m$ . The quantity

$$||y - \sum_{i=1}^n b_i x_i|| = ||\sum_{i=1}^m c_i x_i - \sum_{i=1}^n b_i x_i|| \quad (3.13)$$

is minimized if and only if

$$b_j = c_j, j = 1, 2, \dots, n < m. \quad (3.14)$$

Proof: By Corollary 1, the  $b_j$  which minimize (3.13) satisfy, for  $j = 1, 2, \dots, n$ :

$$\begin{aligned}
 b_j ||x_j||^2 &= (y, x_j) \\
 &= \left( \sum_{i=1}^m c_i x_i, x_j \right) = \sum_{i=1}^m c_i (x_i, x_j) = c_j ||x_j||^2,
 \end{aligned}$$

which implies  $b_j = c_j$ ,  $j = 1, 2, \dots, n$ .

Example. As an example let  $V = L_2(R, B, W)$  as described above. Let  $Q = \{f_i, i = 1, 2, \dots, m\}$  be an orthogonal set of nonzero functions from  $V$ , and let  $y = f(\underline{x})$  be in  $V$ . From the orthogonality,

$$\begin{aligned}
 I_W(f_i, f_j) &= \int_R f_i(\underline{x}) f_j(\underline{x}) dW(\underline{x}) \\
 &= \begin{cases} 0 & \text{if } i \neq j \\ ||f_i||^2 > 0 & \text{if } i = j \end{cases}.
 \end{aligned} \tag{3.15}$$

Also

$$I_W(y, f_j) \equiv I_W(f, f_j) = \int_R f(\underline{x}) f_j(\underline{x}) dW(\underline{x}).$$

Suppose one decides to approximate  $y = f(\underline{x})$  with a linear combination of a subset  $Q^*$  of the vectors (functions) in  $Q$ . Consider the following notation which will be useful in a later example. Let  $K = \{1, 2, \dots, m\}$ .  $K$  is the set of subscripts of functions in  $Q$ . Let  $K^*$  be the subset of  $K$  which contains the subscripts of the functions to be used in approximating  $y$ , i.e., the subscripts of the functions in  $Q^*$ . Then the best linear approximation of  $y$  with respect to the set  $Q^*$ , i.e., the set of coefficients  $b_i$  which minimizes

$$||y - \sum_{i \in K^*} b_i f_i||$$

or, equivalently, which minimizes

$$||y - \sum_{i \in K} b_i f_i||^2 = \int_R [f(\underline{x}) - \sum_{i \in K} b_i f_i(\underline{x})]^2 dW(\underline{x}), \quad (3.16)$$

is defined by

$$b_i = \frac{I_W(y, f_i)}{I_W(f_i, f_i)} = \frac{1}{||f_i||^2} \int_R f(\underline{x}) f_i(\underline{x}) dW(\underline{x}).$$

for each  $i \in K^*$ . Since these  $b_i$  values minimize the integral of the square of the "error" (3.16), the approximation is called a least squares approximation.

Now suppose further that  $y$  has the following structure:

$$y = f(\underline{x}) = \sum_{i \in K} a_i f_i(\underline{x}) \quad (3.17)$$

That is,  $y$  is a linear combination of the functions in  $Q$ . By Corollary 2, the best ("least-squares") approximation of  $y$  with respect to the set  $Q^*$  is given by

$$\tilde{y} = \tilde{f}(\underline{x}) = \sum_{i \in K} a_i f_i(\underline{x}). \quad (3.18)$$

Consider the special case where  $\underline{x}$  is a real variable ( $\underline{x} \in E^1$ ) and  $R = [-1, 1]$ . Also suppose the weight function is  $w(x) = 1$  so that the inner product of functions  $f, g$  over  $V$  is given by

$$I_W(f, g) = \int_{-1}^1 f(x)g(x)dx \quad (3.19)$$

If the set of orthogonal functions is the set of polynomials of degree 0, 1, 2, ...,  $m$ , which are orthogonal with respect to the inner product (3.19), the functions we obtain are the Legendre Polynomials (Abramowitz and Stegun, 1964, Chapter 25). The same region,  $R$ , with the

weighting function  $dW(x) = (1 - x^2)^{-1/2}$ , generates the family of Tchebychev Polynomials. Note also that if  $R = [-\pi, \pi]$  and  $dW(x) = dx$ , the family of functions  $\{\sin nx, \cos nx, n = 0, 1, 2, \dots\}$  forms a set of orthogonal functions.

### 3.1.3. Some Results for Integration-Orthogonal Polynomials

Although the approximation theorem can be extended to cover non-orthogonal sets of vectors, approximation with orthogonal sets of vectors is easier and, in many cases, less subject to round-off error in computational processes (Davis, 1962). Before proceeding to the presentation of some useful results for integration orthogonal polynomials, two techniques are presented for producing orthogonal (or orthonormal) vectors from sets of linearly independent vectors. While the two methods are equivalent, one, the Gram-Schmidt orthonormalization process, is standard material in linear algebra texts, and the other, based on the square root decomposition of positive definite symmetric matrices, is a more useful algorithm in the present study.

Theorem 3.2. Gram-Schmidt Orthogonalization. Let  $V$  be a real vector space with inner product  $(\cdot, \cdot)$  and associated norm,  $|| \cdot ||$ . If  $Z = \{z_1, z_2, \dots, z_m\}$  is any linearly independent set in  $V$ , i.e.,

$$\sum_{i=1}^m a_i z_i = 0 \in V \quad (3.20)$$

is impossible for any set of real scalars  $\{a_i\}$  unless all  $a_i = 0$ , then there exists an orthonormal set  $X = \{x_1, x_2, \dots, x_m\}$  such that

$$x_k = \sum_{i=1}^k a_{ik} z_i. \quad (3.21)$$



The proof consists of the Gram-Schmidt process, which is reproduced here for reference purposes. The procedure is recursive.

Let  $x_1 = \frac{1}{||z_1||} z_1$ . Clearly  $||x_1|| = 1$ ,  $a_{11} = \frac{1}{||z_1||}$ . (Notice that  $||z_1|| = 0$  is impossible, for  $||z_1|| = 0 \rightarrow z_1 = 0$ , and the zero vector cannot belong to any linearly independent set.) Suppose that  $\{x_1, \dots, x_r\}$  has been found so that it is an orthonormal set satisfying (3.21). Let

$$u_{r+1} = z_{r+1} - \sum_{i=1}^r x_i (x_i, z_{r+1}) \quad (3.22)$$

Clearly for  $j = 1, 2, \dots, r$ ,

$$\begin{aligned} (x_j, u_{r+1}) &= (x_j, z_{r+1}) - \sum_{i=1}^r (x_j, x_i) (x_i, z_{r+1}) \\ &= (x_j, z_{r+1}) - (x_j, z_{r+1}) = 0. \end{aligned}$$

Set 
$$x_{r+1} = \frac{u_{r+1}}{||u_{r+1}||}; \quad (3.23)$$

the set  $\{x_1, \dots, x_{r+1}\}$  is orthonormal and  $x_{r+1}$  satisfies (3.21), as was to be shown.

The primary application of the Gram-Schmidt process in this study will involve transformations of sets of "standard" polynomials to orthogonal polynomials. The actual calculations will be performed in a slightly different format. In order to motivate the calculation procedure which follows, consider a set of standard univariate polynomials,

$$F = \{f_0, f_1, \dots, f_m\}. \quad (3.23)$$

where

$$f_j(z) = z^j, \quad j = 0, 1, \dots, m.$$

Consider the vector space  $V = L_2(R, B, W)$  where  $R = [-1, 1]$ ,  $dW(z) = dz$ . As noted earlier, the application of the Gram-Schmidt procedure to the set  $F$  (3.23) produces the Legendre Polynomials. Denote the set of Legendre polynomials through degree  $m$  by  $P = \{p_0, p_1, \dots, p_m\}$ . From (3.21),

$$p_k(x) = \sum_{i=0}^k a_{ik} f_i(z) = \sum_{i=0}^k a_{ik} z^i. \quad (3.24)$$

Form a  $1 \times (m+1)$  row vector (with values in  $E^{(m+1)}$ ) of the functions  $\{p_k\}$  and  $\{f_i\}$ :

$$\begin{aligned} \underline{f} &= (f_0, f_1, \dots, f_m) \\ \underline{p} &= (p_0, p_1, \dots, p_m). \end{aligned} \quad (3.25)$$

Then, equation (3.24) can be written in matrix notation using the matrix  $A = [a_{ij}]$ :

$$\underline{p} = \underline{f}A; \quad \underline{p}(x) = \underline{f}(x)A. \quad (3.26)$$

Note that the  $(m+1) \times (m+1)$  matrix  $A$  is upper triangular;  $a_{ij} = 0$  for  $j < i$ . As noted above, the  $A$  matrix is non-singular; letting  $B = A^{-1}$ ,

$$\underline{f} = \underline{p}A^{-1} = \underline{p}B; \quad \underline{f}(x) = \underline{p}(x)B. \quad (3.27)$$

The  $B$  matrix is also upper triangular.

Now consider a linear combination of the  $f_j$ , a polynomial:

$$g(z) = \sum_{j=0}^m c_j f_j(z) = \sum_{j=0}^m c_j z^j. \quad (3.28a)$$

$$= \underline{f}(z)\underline{c} \quad (3.28b)$$

where the  $(m+1) \times 1$  vector  $\underline{c} = (c_0, c_1, \dots, c_m)^T$ .

Evaluate  $g$  in terms of the orthogonal functions:

$$g(z) = \underline{f}(z)\underline{c} = \underline{p}(z)B\underline{c} \quad (3.29)$$

$$= \sum_{j=0}^m d_j p_j(z)$$

where

$$\underline{d} = B\underline{c} \quad (3.30)$$

which implies

$$\underline{c} = (B)^{-1} \underline{d} = A\underline{d}. \quad (3.31)$$

Using the relations (3.30) and (3.31) a polynomial function  $g(z)$  is easily converted from a standard polynomial representation (3.28a) to an orthogonal polynomial representation (3.29) and vice versa.

It should be evident that the results above hold for conversions from any linearly independent set to an orthogonal set. In particular, the functions may be polynomials in any number of variables. The following theorem demonstrates a computational algorithm for finding the  $A$  and  $B$  matrices in the general case.

Theorem 3.3. Assume the general setting of Theorem 3.2; let  $A$  be the  $m \times m$  nonsingular, upper triangular matrix defined by (3.21) and let  $B = A^{-1}$ . Define the  $m \times m$  symmetric matrix  $M$  with typical element  $m_{ij}$ , where

$$m_{ij} = (z_i, z_j), \quad i = 1, 2, \dots, m; \quad j = 1, 2, \dots, m. \quad (3.32)$$

Then  $M = B'B = (A^{-1})'(A^{-1})$ . Therefore  $M$  is positive definite.

Proof. As before, form the  $1 \times m$  row vectors

$$\underline{z} = (z_1, z_2, \dots, z_m)$$

$$\underline{x} = (x_1, x_2, \dots, x_m).$$

Of course, each element of each vector is a member of the vector space

$V$ . From (3.21),

$$\underline{x} = \underline{z}A$$

or, equivalently,

$$\underline{z} = \underline{x}A^{-1} = \underline{x}B. \quad (3.33)$$

That is, since  $B$  is also triangular,

$$z_i = \sum_{k=1}^m b_{ki} x_k = \sum_{k=1}^i b_{ki} x_k. \quad (3.34)$$

Consider the inner product (the  $(i, j)$  - element of  $M$ ),

$$\begin{aligned} m_{ij} &= (z_i, z_j) = \left( \sum_{k=1}^m b_{ki} x_k, \sum_{\ell=1}^m b_{\ell j} x_\ell \right) \\ &= \sum_{k=1}^m \sum_{\ell=1}^m b_{ki} b_{\ell j} (x_j, x_\ell). \end{aligned} \quad (3.35)$$

But  $\{x_k, k = 1, 2, \dots, m\}$  is orthonormal:

$$(x_k, x_\ell) = \begin{cases} 0 & \text{if } k \neq \ell \\ 1 & \text{if } k = \ell \end{cases}.$$

Therefore (3.35) becomes

$$m_{ij} = (z_i, z_j) = \sum_{k=1}^m b_{ki} b_{kj}. \quad (3.36)$$

But the summation in (3.36) defines the  $(i,j)$ -element of the matrix  $B'B$ , as was to be shown.

The theorem provides a computational method for finding the  $A$  and  $B$  matrices by straightforward matrix computations. Given the positive definite symmetrix  $M$ , one can apply the "square root algorithm" for matrices (see Faddeeva, 1959, for example) to compute an upper triangular matrix  $B$  such that  $B'B = M$ . One can readily compute  $A = B^{-1}$ , since triangular matrices are easily inverted.

The following result will be useful in the next subsection.

Theorem 3.4. Assume the general setting of Theorem 3.2. Suppose it is desired to approximate

$$y = \sum_{i=1}^m \gamma_i z_i$$

by

$$y = \sum_{i=1}^n c_i z_i \quad (3.37)$$

where  $0 < n < m$ . Partition the vectors of  $c_i$ ,  $\gamma_i$ ,  $z_i$ , and  $x_i$  as follows:

$$\underline{c} = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}, \quad \underline{\gamma} = \begin{bmatrix} \gamma_1 \\ \gamma_2 \end{bmatrix}, \quad \underline{z} = (\underline{z}_1, \underline{z}_2); \quad \underline{x} = (\underline{x}_1, \underline{x}_2)$$

where each of  $\underline{c}$ , and  $\underline{\gamma}$  is  $m \times 1$ . Also partition the  $A$  and  $B = A^{-1}$  matrices:

$$A = \begin{bmatrix} A_{11} & \vdots & A_{12} \\ \dots & \dots & \dots \\ 0 & \vdots & A_{22} \end{bmatrix} \begin{matrix} n \\ (m-n) \end{matrix} \quad B = \begin{bmatrix} B_{11} & \vdots & B_{12} \\ \dots & \dots & \dots \\ 0 & \vdots & B_{22} \end{bmatrix} \begin{matrix} n \\ (m-n) \end{matrix}$$

$\begin{matrix} n & m-n \end{matrix} \qquad \begin{matrix} n & m-n \end{matrix}$

If  $\underline{\delta} = \begin{pmatrix} \delta_1 \\ \delta_2 \end{pmatrix}$  is the vector of coefficients of the  $\underline{x}$ -vector,

$$\text{i.e., } \underline{\delta} = B\underline{Y} = \begin{bmatrix} B_{11} & B_{12} \\ 0 & B_{22} \end{bmatrix} \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} = \begin{bmatrix} B_{11} Y_1 + B_{12} Y_2 \\ B_{22} Y_2 \end{bmatrix} \quad (3.38)$$

Then the vector of coefficients  $\underline{c}_1$  which gives the best approximation,  $\tilde{y}$ , is given by

$$\begin{aligned} \underline{c}_1 &= Y_1 + A_{11} B_{12} Y_2 \\ &= \begin{bmatrix} I & : & A_{11} B_{12} \\ n \times n & & \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} . \end{aligned}$$

Proof. Since  $y = \sum_{i=1}^m \delta_i x_i = \underline{x} \underline{\delta}$ , where  $\underline{\delta} = (\delta_1, \dots, \delta_m)'$  is given by

(3.38), by the Approximation Theorem,

$$\tilde{y} = \sum_{i=1}^n \delta_i x_i = \underline{x} \underline{d}$$

$$\text{where } \underline{d} = \begin{pmatrix} d_1 \\ d_2 \end{pmatrix} = \begin{pmatrix} \delta_1 \\ \underline{0} \end{pmatrix} .$$

Then by (3.31),

$$\underline{c} = \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = A \underline{d} = \begin{bmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{bmatrix} \begin{bmatrix} \delta_1 \\ \underline{0} \end{bmatrix}$$

$$\begin{aligned} c_1 &= A_{11} \delta_1 = A_{11} [B_{11} Y_1 + B_{12} Y_2] \\ &= Y_1 + A_{11} B_{12} Y_2, \end{aligned}$$

since  $A_{11} B_{11} = I$ .

Although the above result is useful, it is not nearly as general as it might appear. The approximation it produces depends upon the order in which the  $z$ -vectors are arranged in the orthogonalization process; those  $z$ -vectors (functions) which are to receive zero coefficients must be the last  $(m-n)$  vectors in the Gram-Schmidt procedure for the result above to hold.

The point is that the order in which the  $z$ -vectors are entered into the orthogonalization procedure determines the resulting orthogonal vectors; different orderings produce different sets of orthogonal vectors. The orthogonal vectors are linear combinations of the original vectors; different orderings produce different sets of linear combinations (different  $A$  matrices).

In many problems, especially those involving sets of polynomials in one variable, there is a natural ordering. However, for polynomials in several variables the ordering is somewhat arbitrary. For example, consider, for two independent variables, the following ordered set of functions:

$$(1, x_1, x_2, x_1^2, x_2^2, x_1x_2).$$

Another "reasonable" ordering is:

$$(1, x_1, x_2, x_1x_2, x_1^2, x_2^2).$$

The difference is real, but not altogether easy to grasp. If  $y$  is a linear combination of the functions above and one decides to use an approximation  $\tilde{y}$  which is a linear combination of the functions

$$\{1, x_1, x_2, x_1x_2\}$$

then the approximations produced using the two different orderings above would be identical. However, Theorem 4 applies only to the second ordering.

Now consider how different approximations might be produced. Suppose the orthogonal functions produced from the first ordering are:

$$(p_1, p_2, p_3, p_4, p_5, p_6)$$

and the orthogonal functions produced from the second ordering are denoted:

$$(q_1, q_2, q_3, q_4, q_5, q_6).$$

The results will be identical,  $p_i = q_i$ , for  $i = 1, 2$ , and  $3$ . For the other functions we have the "correspondences" (which are not exact, of course):

$$p_4 \leftrightarrow x_1^2 \leftrightarrow q_5$$

$$p_5 \leftrightarrow x_2^2 \leftrightarrow q_6$$

$$p_6 \leftrightarrow x_1 x_2 \leftrightarrow q_4.$$

Now consider approximation of  $y$  by  $\tilde{y}_1$ , a linear combination of  $p_1, p_2, p_3$ , and  $p_6$ ; and by  $\tilde{y}_2$ , a linear combination of  $q_1, q_2, q_3$ , and  $q_4$ ; the approximations would, in general, be different:

$$||y - \tilde{y}_1|| \neq ||y - \tilde{y}_2||.$$

The approximations are different even though both are linear combinations of terms which "correspond" to the same subset of original functions.



However, once the ordering has been chosen and the orthogonal vectors (functions) have been produced, the approximation theorem is more general than Theorem 4, and order is unimportant in the following sense. If  $X = \{x_1, x_2, \dots, x_m\}$  is an orthogonal set of vectors, and if

$$y = \sum_{i=1}^m d_i x_i,$$

then for any non-empty subset  $X^* \subset X$ ,

$$y = \sum_{x_i \in X^*} d_i x_i$$

is the "best" approximation of  $y$  over the subset  $X^*$ , regardless of which subset was chosen.

### 3.2. Minimum Bias Estimation: Best Linear Unbiased Estimation of the Best Approximating Function

#### 3.2.1. Definitions and Notation

Some consideration is now given to the estimation of approximating functions. The following assumptions, definitions and notation will be useful through the remainder of the study.

It is assumed that  $\eta = \eta(\underline{x})$  is a response function which belongs to the function vector space  $L_2(R, B, W)$  described in section 3.1.1; i.e.,  $\eta$  is a continuous function of the  $r$  variables  $\underline{x} = (x_1, x_2, \dots, x_r)$ , and is square integrable with respect to the finite measure (or distribution function)  $W$  over  $R \subset E^r$ . It is assumed that  $\eta$  has the structure

$$\eta = \sum_{j=1}^m \alpha_j f_j, \quad (3.50)$$

where the  $\alpha_j$  are unknown parameters and the  $f_j$  are known functions (usually polynomials) of  $\underline{x}$ . It is also assumed that each

$$f_j \in L_2(R, B, W).$$

Arrange the functions  $f_j$  in a  $1 \times m$  row vector

$$\underline{f} = (f_1, f_2, \dots, f_m).$$

The Gram-Schmidt process is applied to produce the  $I_w$ -orthogonal (not necessarily orthonormal) vector of functions

$$\underline{p} = (p_1, p_2, \dots, p_m)$$

with the conversion matrices  $A, B = A^{-1}$ , such that

$$\underline{p} = \underline{f}A, \quad \underline{f} = \underline{p}B \quad (3.51)$$

as explained in the previous section. Thus, on setting

$$\underline{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)'$$

$$\underline{\beta} = (\beta_1, \beta_2, \dots, \beta_m)' = B \underline{\alpha}, \quad (3.52)$$

from (3.50) - (3.52):

$$\eta = \underline{f} \underline{\alpha} = \underline{p} \underline{\beta} = \sum_{j=1}^m \beta_j p_j(\underline{x}) \quad (3.53)$$

The function  $\eta$  is not directly observable. Experiments are performed at  $N > m$  points  $\underline{x}_i, i = 1, 2, \dots, N$ . Compute the  $N \times m$  "design matrix," traditionally denoted by  $X$ :

$$X = \begin{bmatrix} f_1(\underline{x}_1) & f_2(\underline{x}_1) & \dots & f_m(\underline{x}_1) \\ f_1(\underline{x}_2) & f_2(\underline{x}_2) & \dots & f_m(\underline{x}_2) \\ \dots & \dots & \dots & \dots \\ f_1(\underline{x}_N) & f_2(\underline{x}_N) & \dots & f_m(\underline{x}_N) \end{bmatrix} = \begin{bmatrix} \underline{f}(\underline{x}_1) \\ \underline{f}(\underline{x}_2) \\ \dots \\ \underline{f}(\underline{x}_N) \end{bmatrix}. \quad (3.54)$$

In terms of the orthogonal functions the design matrix is

$$P = XA = \begin{bmatrix} p(\underline{x}_1) \\ p(\underline{x}_2) \\ \dots \\ p(\underline{x}_N) \end{bmatrix}; \quad (3.55)$$

$$p_{ij} = p_j(\underline{x}_i); 1 \leq j \leq m; 1 \leq i \leq N.$$

It is assumed that  $X$  and  $P$  are of rank  $m$ , which implies that  $X'X$  and  $P'P$  are nonsingular.

At the point  $\underline{x}_i$  the experimenter observes

$$y_i = \eta(\underline{x}_i) + \varepsilon_i, i = 1, 2, \dots, N.$$

These observations are arranged into an  $N \times 1$  vector with the following structure

$$\underline{y} = (y_1, y_2, \dots, y_N)' = X\underline{\alpha} + \underline{\varepsilon} = P\underline{\beta} + \underline{\varepsilon} \quad (3.56)$$

where the vector  $\underline{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N)'$  is a continuous  $N$ -variate random vector with  $E(\underline{\varepsilon}) = \underline{0}$ ,  $E(\underline{\varepsilon} \underline{\varepsilon}') = \sigma_\varepsilon^2 \underline{I}$ . One may assume  $\sigma_\varepsilon^2$  is known or unknown. When  $\underline{\varepsilon}$  is assumed to have the multivariate normal distribution, this assumption will be explicitly stated.

Under these assumptions the application of classical least squares estimation theory (Graybill, 1961) yields the Best Linear Unbiased Estimators (BLUE) for  $\underline{\alpha}$  and  $\underline{\beta}$ :

$$\begin{aligned}\hat{\underline{\alpha}} &= (\underline{X}'\underline{X})^{-1}\underline{X}'\underline{y} \\ \hat{\underline{\beta}} &= (\underline{P}'\underline{P})^{-1}\underline{P}'\underline{y} = (\underline{A}'\underline{X}'\underline{X}\underline{A})^{-1}(\underline{A}'\underline{X}')\underline{y} \\ &= (\underline{A}^{-1})'(\underline{X}'\underline{X})^{-1}(\underline{A}^{-1})'\underline{A}'\underline{X}'\underline{y} = \underline{B}\hat{\underline{\alpha}}.\end{aligned}\quad (3.57)$$

Also:

$$\begin{aligned}E(\hat{\underline{\alpha}}) &= \underline{\alpha}; E[(\hat{\underline{\alpha}} - \underline{\alpha})(\hat{\underline{\alpha}} - \underline{\alpha})'] = \sigma^2 (\underline{X}'\underline{X})^{-1}; \\ E(\hat{\underline{\beta}}) &= \underline{\beta}; E[(\hat{\underline{\beta}} - \underline{\beta})(\hat{\underline{\beta}} - \underline{\beta})'] = \sigma^2 (\underline{P}'\underline{P})^{-1};\end{aligned}\quad (3.58)$$

and, moreover,  $\hat{\alpha}_i$  is the BLUE of  $\alpha_i$ ,  $\hat{\beta}_i$  is the BLUE of  $\beta_i$ ,  $i = 1, 2, \dots, m$ . (If one assumes  $\underline{\varepsilon}$  has the multivariate normal distribution, the estimators above are all Minimum Variance Unbiased Estimators.)

Also, the BLUE for  $\eta$  at the point  $\underline{x}$  is

$$\eta(\underline{x}) = \underline{f}(\underline{x})\hat{\underline{\alpha}} = \underline{p}(\underline{x})\hat{\underline{\beta}} \quad (3.59)$$

where

$$\underline{f}(\underline{x}) = (f_1(\underline{x}), f_2(\underline{x}), \dots, f_m(\underline{x}))$$

and

$$\underline{p}(\underline{x}) = (p_1(\underline{x}), p_2(\underline{x}), \dots, p_m(\underline{x})).$$

Now consider the effect of an inadequate model. Suppose that only  $n < m$  of the  $f_j$  functions are used to "fit" the model (3.56). For

convenience, assume the functions used are  $f_1, f_2, \dots, f_n$  and partition the design matrix and coefficient vectors so the true model is

$$\underline{y} = \underset{N \times n}{X_1} \underset{n \times 1}{\underline{\alpha}_1} + \underset{N \times (m-n)}{X_2} \underset{(m-n) \times 1}{\underline{\alpha}_2} + \underline{\varepsilon}. \quad (3.60)$$

If only the first  $n$  functions (those represented in  $X_1$ ) are used in the estimation process, the estimator  $\underline{a}$  for  $\underline{\alpha}$  is:

$$\underline{a}_1 = (X_1' X_1)^{-1} X_1' \underline{y} \quad (3.61)$$

and

$$\begin{aligned} E(\underline{a}_1) &= (X_1' X_1)^{-1} X_1' [X_1 \underline{\alpha}_1 + X_2 \underline{\alpha}_2 + E(\underline{\varepsilon})] \\ &= \underline{\alpha}_1 + (X_1' X_1)^{-1} X_1' X_2 \underline{\alpha}_2 \end{aligned} \quad (3.62)$$

$$\underline{a}_2 = \underline{0}.$$

That is,  $\underline{a}_1$  is biased by an amount depending on the design matrix and  $\underline{\alpha}_2$ . Denote by  $\hat{\eta}_L$  the estimator for  $\eta$  based on  $\underline{a}_1$ , viz.,

$$\hat{\eta}_L(\underline{x}) = \underline{f}_1(\underline{x}) \underline{a}_1, \quad (3.63)$$

where  $\underline{f}_1(\underline{x}) = (f_1(\underline{x}), f_2(\underline{x}), \dots, f_n(\underline{x}))$  the expected value of  $\hat{\eta}_L(\underline{x})$  is:

$$E[\hat{\eta}_L(\underline{x})] = \underline{f}_1(\underline{x}) E(\underline{a}_1) = \underline{f}_1(\underline{x}) \underline{\alpha}_1 + \underline{f}_1(\underline{x}) (X_1' X_1)^{-1} X_1' X_2 \underline{\alpha}_2 \quad (3.64)$$

Unless  $\underline{\alpha}_2 = 0$ ,  $\hat{\eta}_L$  is biased. In the presence of bias, the variance criterion would seem to be inappropriate as a criterion for comparing estimators of  $\eta$ . The following criterion is more appropriate. Let

$\eta^*$  be any estimator for  $\eta$ ; the Integrated Mean Square Error of  $\eta^*$  is:

$$\text{IMSE}(\eta^*) = \int_R E\{[\eta(\underline{x}) - \eta^*(\underline{x})]^2\} dW(\underline{x}) \quad (3.65)$$

$$\begin{aligned} &= \int_R E\{[\eta(\underline{x}) - E(\eta^*(\underline{x})) + E(\eta^*(\underline{x})) - \eta^*(\underline{x})]^2\} dW(\underline{x}) \\ &= \int_R E\{[\eta(\underline{x}) - E(\eta^*(\underline{x}))]^2\} dW(\underline{x}) + \int_R E\{[\eta^*(\underline{x}) - E(\eta^*(\underline{x}))]^2\} dW(\underline{x}) \end{aligned} \quad (3.66)$$

### 3.2.2. Best Linear Unbiased Estimation of the Best Approximating Function

The function actually being estimated by  $\hat{\eta}_L$  above (3.63), i.e.,  $E(\hat{\eta}_L)$ , is an approximating function for  $\eta$ . Since  $E(\hat{\eta}_L)$  depends upon the design matrix, it may or may not be a "good" approximating function. Consideration is now given to estimation of the best approximating function in a somewhat more general setting.

Suppose one chooses to approximate  $\eta$  as a linear combination of a subset of the functions  $\underline{p} = (p_1, p_2, \dots, p_m)$ . Let  $\mathcal{J}^*$  be the (non-empty) set of subscripts of the selected functions. By the Approximation Theorem the best linear approximating function for  $\eta$  with respect to the chosen subset of functions is obtained by simply ignoring terms in the orthogonal-function representation of  $\eta$ :

$$\eta = \sum_{i \in \mathcal{J}^*} \beta_i p_i. \quad (3.67)$$

Notice that  $\tilde{\eta}$  is not equivalent to

$$\hat{\eta}_L = \sum_{i \in \mathcal{J}^*} \alpha_i f_i,$$

which would be obtained by ignoring terms in the representation of  $\eta$  as a linear combination of the original nonorthogonal functions. By classical least squares estimation theory (Graybill, 1961, Theorem 6.3), the best linear estimator for  $\tilde{\eta}$  at the point  $\underline{x}$  is:

$$\hat{\tilde{\eta}}(\underline{x}) = \sum_{i \in \mathcal{I}^*} \hat{\beta}_i p_i(\underline{x}). \quad (3.68)$$

This representation can be converted back to a representation in terms of the original  $\underline{f}$  functions as follows. Define the  $m \times 1$  vector  $\tilde{\underline{\beta}}$  by:

$$\tilde{\beta}_i = \begin{cases} \hat{\beta}_i & \text{if } i \in \mathcal{I}^* \\ 0 & \text{if } i \notin \mathcal{I}^* \end{cases} \quad (3.69)$$

Then,

$$\tilde{\underline{\alpha}} = A \tilde{\underline{\beta}} \quad (3.70)$$

and

$$\hat{\tilde{\eta}}(\underline{x}) = \sum_{i=1}^m \tilde{\alpha}_i f_i(\underline{x}). \quad (3.71)$$

### 3.2.3. Equivalence of the Minimum Bias Estimator and BLUE of the Best Approximating Function

In general, in the work above,  $\beta_i = 0$  does not imply  $\tilde{\alpha}_i = 0$ ; that is, in general all  $m$  terms will be used in the representation of  $\hat{\tilde{\eta}}$  in terms of the original nonorthogonal functions.

Consider, however, the following special case which occurs, for example, if the  $f_i$  are polynomials in one variable, arranged in order of increasing degree, and if the approximator  $\tilde{\eta}$  is of degree  $n-1 < m-1$ . (The special case being discussed arises in any situation in which  $\tilde{\eta}$

is a linear function of the first  $n$  elements of the ordered set

$\underline{f} = (f_1, f_2, \dots, f_m)$ . Here,  $\mathcal{A}^* = \{1, 2, \dots, n\}$  and

$$\tilde{\eta} = \sum_{i=1}^n \beta_i p_i.$$

Partition the coefficient vectors

$$\underline{\beta}_1 = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}, \quad \underline{\alpha}_1 = \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix},$$

where  $\underline{\alpha}_1$  and  $\underline{\beta}_1$  are each  $n \times 1$ . Partition the estimator vectors  $\hat{\underline{\alpha}}, \hat{\underline{\beta}}$ , the function vectors  $\underline{f}, \underline{p}$  and the  $A$  and  $B$  matrices to correspond. By Theorem 4 of the previous section,

$$\tilde{\eta}_1 = \underline{p}_1 \underline{\beta}_1 = \underline{f}_1 [\underline{\alpha}_1 + A_{11} B_{12} \underline{\alpha}_2].$$

Thus the BLUE of  $\tilde{\eta}_1$  is

$$\begin{aligned} \hat{\tilde{\eta}}_1 &= \underline{p}_1 \hat{\underline{\beta}}_1 = \underline{f}_1 (\hat{\underline{\alpha}}_1 + A_{11} B_{12} \hat{\underline{\alpha}}_2) = \underline{f}_1 \hat{\underline{\alpha}} \\ &= \underline{f}_1 [I : A_{11} B_{12}] \begin{bmatrix} \hat{\underline{\alpha}}_1 \\ \hat{\underline{\alpha}}_2 \end{bmatrix} = \underline{f}_1 [I : A_{11} B_{12}] \hat{\underline{\alpha}}. \end{aligned} \quad (3.72)$$

where

$$\hat{\underline{\alpha}} = [I : A_{11} B_{12}] \underline{\alpha}. \quad (3.73)$$

The setting described above is equivalent to the setting Karson et al. (1969) used in deriving the Minimum Bias Estimator under the following restrictions.  $W(\underline{x})$  represents the uniform weighting function over  $R$ , i.e.,  $dW(\underline{x}) = dx_1 dx_2 \dots dx_r$ . The integration used in defining the



inner product is Riemann integration. The functions  $f_i(\underline{x})$  are standard polynomials in the  $r$ -variables, arranged in order of increasing degree in the following sense. If  $i=n$  ( $f_i$  to be used in the approximation) and  $n < j \leq m$  ( $f_j$  not to be used in the approximation), then it was assumed that the degree of  $f_i$  in  $x_k$  is less than or equal to the degree of  $f_j$  in  $x_k$ ,  $k = 1, 2, \dots, r$ . For example one might have:

$$\underline{f}_1 = (1, x_1, x_2), \quad \underline{f}_2 = (x_1 x_2, x_1^2, x_2^2).$$

With such functions the  $M$  matrix of inner products becomes the matrix of "moments" of the region  $R$ .

$$m_{ij} = \int_R f_i(\underline{x}) f_j(\underline{x}) d\underline{x}.$$

In the above example, one element of  $M$  is

$$m_{24} = \int \int_R x_1 (x_1 x_2) dx_1 dx_2 = \int \int_R x_1^2 x_2 dx_1 dx_2.$$

The MBE for  $\underline{\alpha}$ , under the above assumptions is

$$\hat{\underline{\alpha}}_1 = [I: M_{11}^{-1} M_{12}] \begin{pmatrix} \hat{\underline{\alpha}}_1 \\ \hat{\underline{\alpha}}_2 \end{pmatrix} = [I: M_{11}^{-1} M_{12}] \hat{\underline{\alpha}} \quad (3.74)$$

Where  $\hat{\underline{\alpha}}$  is the BLUE of  $\underline{\alpha}$  (using the full model) and

$$M = \begin{bmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{bmatrix}$$

is partitioned to match the partition of  $\underline{\alpha}$ . ( $M_{11}$  is  $n \times n$ .) To show the equivalence of the estimators (3.73) and (3.74) it is only necessary to show

$$A_{11} B_{12} = M_{11}^{-1} M_{12} \quad (3.75)$$

Recall that  $A_{11} = B_{11}^{-1}$ ; (3.75) becomes

$$B_{11}^{-1} B_{12} = M_{11}^{-1} M_{12} \quad (3.76)$$

Since  $B'B = M$ , we have

$$\begin{bmatrix} B'_{11} & 0 \\ B'_{12} & B'_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ 0 & B_{22} \end{bmatrix} = \begin{bmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{bmatrix}.$$

Hence

$$B'_{11} B_{11} = M_{11}$$

$$B'_{11} B_{12} = M_{12}$$

and

$$\begin{aligned} M_{11}^{-1} M_{12} &= (B'_{11} B_{11})^{-1} B'_{11} B_{12} = B_{11}^{-1} (B'_{11})^{-1} B'_{11} B_{12} \\ &= B_{11}^{-1} B_{12} = A_{11} B_{12}. \end{aligned}$$

This demonstrates the equivalence of the estimators in the special case considered by Karson et al. (1969). Although they did not explicitly derive the Minimum Bias Estimator as the BLUE of the Best Approximating Function, this is clearly the principle which guided their work. Thus, in the present exposition an estimation procedure has been developed which represents a generalization of the Minimum Bias Estimator.

#### 4. A PROCEDURE FOR SELECTION OF TERMS IN THE MODEL WITH MINIMUM BIAS ESTIMATION AND THE INTEGRATED MEAN SQUARE ERROR CRITERION

Many procedures have been proposed for attack on the problem of selecting independent variables in a regression problem. References to papers describing many of these techniques have been given in the Review of Literature section of this paper. All of the procedures mentioned there suffer one common affliction, which derives from the fact that the procedures are based on "local" criteria, such as the  $R^2$  criterion. Let the true response (or "model", or "regression equation") be denoted by

$$\eta(\underline{x}) = \sum_{j=1}^m \alpha_j f_j(\underline{x}),$$

where  $\underline{x}$  is an  $n \times 1$  column vector,  $\underline{x} = (x_1, x_2, \dots, x_n)'$  and  $f_j(\underline{x})$  is a polynomial in the  $n$  variables  $x_1, \dots, x_n$ . For example, if  $n = 4$ , one might have

$$f_1(\underline{x}) = x_1 x_2^2 x_3 x_4^3.$$

Let  $R$ , the "region of interest", be a subset of Euclidean  $n$ -space,  $E^n$ . In essence, procedures for selecting terms in a model produce an estimator for  $\eta$ , say  $\check{\eta}$ , of the form

$$\check{\eta}(\underline{x}) = \sum_{j=1}^m \check{\alpha}_j f_j(\underline{x}).$$

If the  $j$ -th term is excluded from the model (by whatever procedure is being applied), the effect is to set  $\check{\alpha}_j = 0$ . Suppose, for example, that application of the procedure being used results in setting  $\check{\alpha}_j = 0$ ,  $j = m_1 + 1, m_1 + 2, \dots, m$ , so that the estimator is

$$\check{\eta}_L(\underline{x}) = \sum_{j=1}^{m_1} \check{\alpha}_j f_j(\underline{x}) + \sum_{j=m_1+1}^m 0 \cdot f_j(\underline{x}).$$

The point is that the function really being estimated by  $\check{\eta}_L$  (that is,  $E[\check{\eta}_L]$ ) is not (except in special cases) the best approximating function for  $\eta$ . That is, define

$$\eta_L = \sum_{j=1}^{m_1} \alpha_j f_j(\underline{x}) \quad (*)$$

the function being estimated by  $\check{\eta}_L$ .  $\check{\eta}_L$  is not (usually) the best approximating function of  $\eta$  of the form (\*), as was shown in the previous chapter. That is, in general the estimator  $\check{\eta}_L$ , produced by the procedures mentioned above, is not a minimum bias estimating function.

In this chapter a procedure is developed which is "global" in the sense that the criterion takes into account the properties of the estimation-approximation procedure over the whole region of interest,  $R$ . Moreover, as will be seen, the procedure is applicable to full general linear models including polynomials in many variables. The procedure allows one to test any term in the model (low-order polynomials as well as high order polynomials) for "significance", and the one "best" regression equation of all possible regression equations (as measured by the criterion defined below) is simply obtained. The properties of the coefficient estimators are also derived.

#### 4.1. Minimization of Integrated Mean Square Error with Respect to Choice of Terms in the Model, Parameters Known

Assume that the true response function  $\eta$  is a polynomial function of  $n$  variables,  $\underline{x} = (x_1, \dots, x_n)'$ , and assume a "region of interest",  $R \subset E^n$ . Assume also a weight function or finite measure  $W$  is defined over  $R$  and that  $\eta \in L_2(R, B, W)$  the vector space discussed in the previous chapter. That is

$$0 < \int_R \eta^2(\underline{x}) dW(\underline{x}) < +\infty \quad (4.1)$$

Let  $P_k(\underline{x})$ ,  $k = 1, \dots, m$  be integration-orthogonal polynomials in  $L_2$ , i.e.,

$$\int_R P_i(\underline{x}) P_j(\underline{x}) dW(\underline{x}) = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \quad (4.2)$$

Assume that the degree of each  $P_k(\underline{x})$  in each of the variables  $x_1, \dots, x_n$  is specified. Assume also that the true response function,  $\eta$ , can be written as

$$\eta(\underline{x}) = \sum_{j=1}^m \beta_j P_j(\underline{x}) \quad (4.3)$$

Note that  $\beta_j = 0$  for any  $j \in \{1, 2, \dots, m\}$  is explicitly allowed.

Let  $\check{\beta}_j$ ,  $j = 1, 2, \dots, m$  be estimators of the  $\beta_j$ , and define

$$\check{\eta}(\underline{x}) = \sum_{j=1}^m \check{\beta}_j P_j(\underline{x}). \quad (4.4)$$

Certain of the estimators  $\check{\beta}_j$  are explicitly allowed to be zero with non-zero probability, i.e.,

$$P[\check{\beta}_j = 0] > 0.$$

Note that setting  $\beta_j = 0$  is equivalent to deleting the  $j$ -th term from the model (or the "regression equation"). It is assumed that  $E[(\check{\beta}_j)^2] < +\infty$  for each  $j = 1, 2, \dots, m$ , and that each  $\beta_j$  is estimable.

Define the integrated mean square error (IMSE) for  $\check{\eta}$  to be

$$\text{IMSE}(\check{\eta}) = \int_R E\{[\eta(\underline{x}) - \check{\eta}(\underline{x})]^2\} dW(\underline{x}). \quad (4.5)$$

Theorem 4.1. With the setting described above

$$\text{IMSE}(\check{\eta}) = \sum_{j=1}^m E[(\beta_j - \check{\beta}_j)^2]. \quad (4.6)$$

Proof. Since all terms are finite, the integral and the expected value operator can be interchanged in (4.5):

$$\begin{aligned} \text{IMSE} &= E \int_R [\eta(\underline{x}) - \check{\eta}(\underline{x})]^2 dW(\underline{x}) \\ &= E \int_R \left[ \sum_{j=1}^m (\beta_j - \check{\beta}_j) P_j(\underline{x}) \right]^2 dW(\underline{x}) \\ &= E \int_R \left[ \sum_{j=1}^m \sum_{i=1}^m (\beta_i - \check{\beta}_i)(\beta_j - \check{\beta}_j) P_i(\underline{x}) P_j(\underline{x}) \right] dW(\underline{x}) \\ &= E \left[ \sum_{j=1}^m \sum_{i=1}^m (\beta_i - \check{\beta}_i)(\beta_j - \check{\beta}_j) \int_R P_i(\underline{x}) P_j(\underline{x}) dW(\underline{x}) \right] \\ &= E \left[ \sum_{j=1}^m (\beta_j - \check{\beta}_j)^2 \right] = \sum_{j=1}^m E[(\beta_j - \check{\beta}_j)^2] \end{aligned} \quad (4.7)$$

as was to be shown.

Some remarks on equation (4.7) are in order. First notice that even though the estimators  $\hat{\beta}_j$  may be correlated, no covariance terms appear in this expression of the IMSE. Also, notice that the summation form of equation (4.7) allows one to consider the coefficient estimation one term at a time.

#### 4.1.1. Experimental Model

The following experimental model will be used throughout the paper.

It is assumed that a "true" response function

$$\eta(\underline{x}) = \sum_{j=1}^m \beta_j P_j(\underline{x}) \quad (4.8)$$

is defined over the region of interest  $R \subset E^n$  and that the  $\{P_j(\underline{x})\}$  are integration-orthonormal functions (usually polynomials) of the  $n$  variables  $\underline{x} = (x_1, \dots, x_n)'$ , i.e.,

$$\int_R P_i(\underline{x}) P_j(\underline{x}) dW(\underline{x}) = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \quad (4.9)$$

where  $W(\underline{x})$  is either a finite measure defined over  $R$  or a distribution function of finite variation over  $R$  as discussed in the previous chapter.

It is further assumed that  $N$  observations are taken at the points  $\underline{x}_i$ :

$$Y_i = \eta(\underline{x}_i) + \varepsilon_i, \quad i = 1, 2, \dots, N. \quad (4.10)$$

Also,  $\underline{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_N)'$  has the multivariate normal distribution,

$N(\underline{0}, \sigma_{\varepsilon}^2 \underline{I})$ , where  $\sigma_{\varepsilon}^2$  may or may not be known. Thus the vector of observations may be expressed as

$$\begin{matrix} \underline{Y} \\ N \times 1 \end{matrix} = \begin{matrix} P \underline{\beta} \\ N \times m \quad m \times 1 \end{matrix} + \begin{matrix} \underline{\varepsilon} \\ N \times 1 \end{matrix}, \quad (4.11)$$

where the matrix  $P = [p_{ij}]$  is defined by

$$p_{ij} = P_j(\underline{x}_i), \quad j = 1, 2, \dots, m; \quad i = 1, 2, \dots, N. \quad (4.12)$$

The matrix  $P$  is assumed to have full rank so that the minimum variance unbiased estimator of  $\underline{\beta}$  is

$$\hat{\underline{\beta}} = (P'P)^{-1}P'\underline{Y} \quad (4.13)$$

which has the multivariate normal distribution  $N(\underline{\beta}, \sigma_{\varepsilon}^2 (P'P)^{-1})$ . The minimum variance unbiased estimator of  $\sigma_{\varepsilon}^2$  is

$$\hat{\sigma}_{\varepsilon}^2 = \frac{1}{N-m} (\underline{Y}'\underline{Y} - \hat{\underline{\beta}}'P'\underline{Y}) \quad (4.14)$$

and  $\frac{(N-m)\hat{\sigma}_{\varepsilon}^2}{\sigma_{\varepsilon}^2}$  has the central  $\chi^2$  distribution with  $N-m$  degrees of freedom, and is distributed independently of  $\hat{\underline{\beta}}$ .

#### 4.1.2. An Example

Suppose one decides in advance of an experiment that certain ones of the  $\beta_j^v = 0$  (i.e., certain terms are to be deleted from the model), and that others are to be set equal to the corresponding least squares estimators (terms to be included in the model).

Under these assumptions, for coefficients corresponding to terms "in the model",



$\check{\beta}_j = \hat{\beta}_j$ , the least squares estimator, and

$$E[(\check{\beta}_j - \beta_j)^2] = E[(\hat{\beta}_j - \beta_j)^2] = \text{Var}(\hat{\beta}_j). \quad (4.15)$$

For terms not "in the model", i.e., those terms for which  $\check{\beta}_j = 0$  (set in advance of the experiment),

$$E[(\check{\beta}_j - \beta_j)^2] = E[(0 - \beta_j)^2] = \beta_j^2. \quad (4.16)$$

Then:

$$\begin{aligned} \text{IMSE}(\check{\eta}) &= \int_R E[(\eta(\underline{x}) - \check{\eta}(\underline{x}))^2] dW(\underline{x}) \\ &= \sum_{\text{terms in the model}} \text{Var}(\hat{\beta}_j) + \sum_{\text{terms not in the model}} \beta_j^2. \end{aligned} \quad (4.17)$$

As explained in Chapter 3, the procedure described above is just the Minimum Bias Estimation procedure.

The result above suggests an intriguing possibility for obtaining a smaller IMSE. If it were known in advance that  $\beta_j^2 > \text{Var}(\hat{\beta}_j)$ , one would set  $\check{\beta}_j = \hat{\beta}_j$  (keep the  $j$ -th term in the model) and obtain a smaller IMSE than would be obtained if the  $j$ -th term were deleted from the model. The procedure would be applied individually to each term in the model.

An additional bonus is obtained: no matter which terms are deleted from the model (by setting  $\check{\beta}_j = 0$ ) the resulting estimator,  $\check{\eta}(\underline{x})$  is the minimum bias estimator for  $\eta(\underline{x})$  of the form chosen, and  $\tilde{\eta}(\underline{x}) = E[\check{\eta}(\underline{x})]$  is the best approximating function of  $\eta$  of the form chosen.

The procedure above can be summarized as follows:

$$\check{\beta}_j = \begin{cases} \hat{\beta}_j & \text{if } \beta_j^2 > \text{var}(\hat{\beta}_j) \Leftrightarrow |\beta_j| > \sqrt{\text{var}(\hat{\beta}_j)} \\ 0 & \text{otherwise, i.e., } |\beta_j| < \sqrt{\text{var}(\hat{\beta}_j)} \end{cases}$$

for  $j = 1, 2, \dots, m$ .

Of course, if one knew the  $\beta_j$  values in advance there would be no need for experimentation and estimation. However, the above procedure motivates the following estimation procedure for each term ( $\beta_j$ ) in the model:

(1) Test the hypothesis:

$$H_0: |\beta| < \sqrt{\text{Var}(\hat{\beta}_j)}$$

(2) Define the estimator:

$$\check{\beta}_j = \begin{cases} 0 & \text{if } H_0 \text{ is accepted} \\ \hat{\beta}_j & \text{if } H_0 \text{ is rejected } (\hat{\beta}_j \text{ is the least squares estimator.}) \end{cases}$$

The test in step (1) will depend on whether  $\sigma_j^2 = \text{var}(\hat{\beta}_j)$  is known and on the level ( $\alpha$ ) of the test. The properties of the estimator above for the two cases ( $\sigma_E^2$  known,  $\sigma_E^2$  unknown) are derived in the next two sections. When these properties are developed (based on UMP tests), it is seen that the estimator is of the form:

$$\check{\beta}_j = \begin{cases} \hat{\beta}_j & \text{if } |\hat{\beta}_j| > c \sqrt{\hat{\sigma}_j^2} \\ 0 & \text{if } |\hat{\beta}_j| < c \sqrt{\hat{\sigma}_j^2} \end{cases} \quad (4.18)$$

(if  $\sigma_j^2$  is known, replace  $\hat{\sigma}_j^2$  by  $\sigma_j^2$ ), and one can discuss the "cutoff point,"  $C$ , rather than the level,  $\alpha$ , of the test. A discussion of the choice of the cutoff points follows development of the properties of the estimators.

The procedure outlined above is based on a two-tail test. There are occasions when a one-tail procedure is appropriate. In such a case, if it is known that  $\beta_j \geq 0$ , the estimator is:

$$\check{\beta}_j = \begin{cases} \hat{\beta}_j & \text{if } \hat{\beta}_j > c \sqrt{\hat{\sigma}_j^2} \\ 0 & \text{if } \hat{\beta}_j < c \sqrt{\hat{\sigma}_j^2} \end{cases}$$

(where  $\hat{\sigma}_j^2$  is replaced by  $\sigma_j^2$ , if known). If  $\beta_j \leq 0$  one considers  $-\beta_j$ .

The properties of "one tail estimators" are developed in section 4.3.

Consider the following generalized procedure. Let  $A$  be a subset of the real line, and define:

$$\check{\beta}_j = \begin{cases} 0 & \text{if } \frac{\hat{\beta}_j}{\sqrt{\hat{\sigma}_j^2}} \in A \\ \hat{\beta}_j & \text{otherwise} \end{cases}$$

This type of procedure is suggested by Bayesian techniques for selection of regions on which  $\check{\beta}_j = 0$ . The properties of this generalized procedure for certain simple types of  $A$ -sets (intervals, complements of intervals, and half lines) are developed in section 4.4.

#### 4.1.3. Results Useful for Evaluating Integrals

The evaluation of expectations and mean square errors in succeeding sections will require the numerical evaluation of the cumulative distribution function and partial first and second moments of the standard normal probability density function. The algorithms given here are intended for use on a computer; the computer used for computing charts in this study was the IBM 360 Model 75. Since IBM and most other computer manufacturers supply very accurate algorithms for evaluating the error function, defined by:

$$\text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt,$$

this function is used as a basis for evaluating the functions below.

The cumulative normal distribution function, denoted by  $\Phi$ , can be evaluated in terms of the error function (Abramowitz and Stegun, 1964, equation 7.1.22):

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-t^2/2} dt = \frac{1}{2} [1 + \text{erf}(x/\sqrt{2})] \quad (4.19)$$

The partial first moment (denoted PFM) of the standard normal probability density, defined by

$$\text{PFM}(a,b) = \frac{1}{\sqrt{2\pi}} \int_a^b t e^{-t^2/2} dt, \quad (4.20a)$$

can be evaluated directly by observing that

$$\frac{d}{dt} (-\exp(-t^2/2)) = t \exp(-t^2/2).$$

Thus the PFM function may be evaluated by:

$$\text{PFM}(a,b) = \frac{1}{\sqrt{2\pi}} [\exp(-a^2/2) - \exp(-b^2/2)]. \quad (4.20b)$$

The partial second moment (PSM) of the standard normal probability density, defined by:

$$\text{PSM}(a,b) = \frac{1}{\sqrt{2\pi}} \int_a^b t^2 \exp(-t^2/2) dt \quad (4.21a)$$

can be evaluated by integration-by-parts. The result is:

$$\frac{1}{\sqrt{2\pi}} \int t^2 \exp(-t^2/2) dt = \frac{-t \exp(-t^2/2)}{\sqrt{2\pi}} + \int \frac{\exp(-t^2/2)}{\sqrt{2\pi}} dt. \quad (4.21b)$$

Therefore the PSM function may be evaluated as:

$$\text{PSM}(a,b) = \frac{1}{\sqrt{2\pi}} \left[ a \exp\left(-\frac{a^2}{2}\right) - b \exp\left(-\frac{b^2}{2}\right) \right] + \Phi(b) - \Phi(a) \quad (4.21c)$$

The  $\Phi$ , PSM, and PFM functions were programmed as double precision FORTRAN subprograms for use in evaluating results discussed in the following sections.

The following lemmas will also be useful in the following sections.

Lemma 4.1. Let  $(x,u)$  be jointly independently distributed random variables such that  $x$  has the (marginal) normal distribution,  $N(\theta,1)$ , and  $vu$  has the (marginal) chi-square distribution with  $v$  degrees of freedom, i.e., the joint density function of  $(x,vu)$  is:

$$\begin{aligned} f(x,vu; \theta, v) &= \frac{\exp [-(x-\theta)^2/2]}{\sqrt{2\pi}} \cdot \frac{(vu)^{v/2-1} \exp (-u/2)}{2^{v/2} \Gamma(v/2)} \\ &= f_1(x; \theta, 1) f_2(u; v) = f_1(x) f_2(u) \end{aligned}$$

over the region

$$\{(x,u); -\infty < x < +\infty, 0 \leq u\}.$$

Let  $a_1, a_2, a_3, a_4$  be given constants satisfying

$$-\infty \leq a_1 < a_2 \leq 0 \leq a_3 < a_4 \leq +\infty$$

and define the following subsets of the  $(x,u)$  - sample space:

$$A_i = \{(x,u): u \geq 0, a_i \sqrt{u} < x \leq a_{i+1} \sqrt{u}\}$$

$$B_i = \{(x,u): u \geq 0, (x,u) \notin A_i\}, \quad i = 1, 2, 3.$$

$B_i$  is the complement of  $A_i$  relative to the sample space. Let  $h(x,\theta)$  be a continuous function of  $x$  and a parameter,  $\theta$ . In the applications,  $h(x,\theta) = x$  or  $h(x,\theta) = (x-\theta)^2$ . Let  $k$  be a real constant; in the applications  $k = 0$  or  $k = \theta^2$ . Define the functions

$$g_i(x,u) = h(x,\theta) 1_{B_i}(x,u) + k 1_{A_i}(x,u)$$

$$g_i^*(x,u) = h(x,\theta) 1_{A_i}(x,u) + k 1_{B_i}(x,u)$$

and

$$I(\theta; a, b; c, d) =$$

$$\int_a^b [k - h(t+\theta, \theta)] f_1(t; 0, 1) P\left[\frac{v(t+\theta)^2}{c^2} < \chi_v^2 < \frac{v(t+\theta)^2}{d^2}\right] dt.$$

Then,

$$E[g_1(x,u)] = E[h(x,\theta)] + I(\theta; -\infty, -\theta; a_1, a_2);$$

$$E[g_1^*(x,u)] = k - I(\theta; -\infty, -\theta; a_1, a_2);$$

$$E[g_2(x,u)] = E[h(x,\theta)] + I(\theta; -\infty, -\theta; a_2, 0)$$

$$+ I(\theta; -\theta, \infty; a_3, 0);$$

$$E[g_2^*(x,u)] = k - I(\theta; -\infty, -\theta; a_2, 0) - I(\theta; -\theta, \infty; a_3, 0);$$

$$E[g_3(x,u)] = E[h(x,\theta)] + I(\theta; -\theta, \infty; a_4, a_3);$$

$$E[g_3^*(x,u)] = k - I(\theta; -\theta, \infty; a_4, a_3).$$

Proof. First note that

$$\begin{aligned} E[g_1(x,u)] &= \iint_{B_1} h(x,\theta) f(x,u) du dx + k \iint_{A_1} f(x,u) du dx \\ &= E[h(x,\theta)] + \iint_{A_1} [k - h(x,\theta)] f(x,u) du dx, \end{aligned}$$

and

$$\begin{aligned} E[g_1^*(x,u)] &= \iint_{A_1} h(x,\theta) f(x,u) du dx + k \iint_{B_1} f(x,u) du dx \\ &= k - \iint_{A_1} [k - h(x,\theta)] f(x,u) du dx \\ &= k + E[h(x,\theta)] - E[g_1(x,u)]. \end{aligned}$$

The integrals over the sets  $A_1$  must be evaluated separately. The integral over  $A_1$  is:

$$\begin{aligned}
& \iint_{A_1} [k-h(x,\theta)] f(x,u) \, du \, dx \\
&= \int_{-\infty}^0 [k-h(x,\theta)] f_1(x;\theta,1) \int_{x^2/a_1^2}^{x^2/a_2^2} f_2(u,v) \, du \, dx.
\end{aligned}$$

Change the variable of integration for the outer integral to

$t = x - \theta$ . The inner integral becomes

$$P \left[ \frac{v(t+\theta)^2}{a_1^2} < \chi_v^2 < \frac{v(t+\theta)^2}{a_2^2} \right],$$

the limits of integration change from  $(-\infty, 0)$  to  $(-\infty, -\theta)$ , and

$k - h(x,\theta) = k - h(t+\theta,\theta)$ . Thus, the whole quantity is

$$\begin{aligned}
& \int_{-\infty}^{-\theta} [k - h(t+\theta)] f_1(t; \theta, 1) P \left[ \frac{v(t+\theta)^2}{a_1^2} < \chi_v^2 < \frac{v(t+\theta)^2}{a_2^2} \right] dt \\
&= I(\theta; -\infty, -\theta; a_1, a_2);
\end{aligned}$$

adding  $E[h(x,\theta)]$  produces the result stated in the Lemma.

The other derivations are essentially the same and use the same change of variable:

$$\begin{aligned}
& \iint_{A_2} [k - h(x,\theta)] f(x,u) \, du \, dx \\
&= \int_{-\infty}^0 [k - h(x,\theta)] f_1(x;\theta,1) \int_{x^2/a_2^2}^{\infty} f_2(u,v) \, du \, dx
\end{aligned}$$



$$\begin{aligned}
& + \int_0^{\infty} [k - h(x, \theta)] f_1(x; \theta, 1) \int_{x^2/a_3}^{\infty} f_2(u, v) du dx \\
& = I(\theta; -\infty, -\theta; a_2, 0) + I(\theta; -\theta, \infty, a_3, 0).
\end{aligned}$$

Addition of  $E[h(x, \theta)]$  produces the desired result.

$$\begin{aligned}
& \iint_{A_3} [k - h(x, \theta)] f(x, u) du dx \\
& = \int_0^{\infty} [k - h(x, \theta)] f_1(x) \int_{x^2/a_4}^{x^2/a_3} f_2(u, v) du dx \\
& = I(\theta; -\theta; a_4, a_3).
\end{aligned}$$

Addition of  $E[h(x, \theta)]$  produces the desired result.

The expectations of the functions  $g_i^*(x, u)$  are found from the equations above for  $E[g_i(x, u)]$  and the relations between  $E[g_i(x, u)]$  and  $E[g_i^*(x, u)]$ .

The integrals above with the chi-square probability in the integrand must be evaluated by numerical approximation. In such methods it is important to evaluate the integrand very accurately; Abramowitz and Stegun (1964, Chapter 26) give equations for the direct evaluation of the chi-square distribution function. These formulas were incorporated in a double precision FORTRAN FUNCTION subprogram, QCHI, which was used for calculations in this paper.

Lemma 4.2. Let  $x$  be a random variable with the  $N(\theta, 1)$  distribution, with density denoted  $f_1(x; \theta, 1)$ . Let  $a_1, a_2$  be given constants such that  $-\infty \leq a_1 < a_2 \leq +\infty$  and define the following subsets of the sample space:

$$A = \{x: a_1 \leq x \leq a_2\} = [a_1, a_2].$$

$$B = A^c = \{x: x < a_1 \text{ or } x > a_2\}.$$

Let  $h(x, \theta)$  be a given continuous function and let  $k$  be a given constant. Define

$$g(x) = h(x, \theta) \cdot 1_B(x) + k \cdot 1_A(x)$$

$$g^*(x) = h(x, \theta) \cdot 1_A(x) + k \cdot 1_B(x).$$

Then,

$$\begin{aligned} E[g(x)] &= E[h(x, \theta)] + k[\Phi(a_2 - \theta) - \Phi(a_1 - \theta)] \\ &\quad - \int_{a_1 - \theta}^{a_2 - \theta} h(t + \theta, \theta) f_1(t; 0, 1) dt \end{aligned}$$

and

$$E[g^*(x)] = k + E[h(x, \theta)] - E[g(x)].$$

Proof. First consider

$$\begin{aligned} E[g(x)] &= \int_B h(x, \theta) f_1(x; \theta, 1) dx + k \int_A f_1(x; \theta, 1) dx \\ &= E[h(x, \theta)] + \int_{a_1}^{a_2} [k - h(x, \theta)] f_1(x; \theta, 1) dx. \end{aligned}$$

Let  $t = x - \theta$ ; then

$$E[g(x)] = E[h(x, \theta)] + \int_{a_1 - \theta}^{a_2 - \theta} [k - h(t + \theta, \theta)] f_1(t; 0, 1) dt$$

$$\begin{aligned}
&= E[h(x, \theta)] + k \int_{a_1^{-\theta}}^{a_2^{-\theta}} f_1(t; 0, 1) dt \\
&\quad - \int_{a_1^{-\theta}}^{a_2^{-\theta}} h(t + \theta, \theta) f_1(t; 0, 1) dt
\end{aligned}$$

which is equivalent to the desired result. Now

$$\begin{aligned}
E[g^*(x)] &= \int_A h(x, \theta) f_1(x; \theta, 1) dx + k \int_B f_1(x; \theta, 1) dx \\
&= k + \int_A (h(x, \theta) - k) f_1(x; \theta, 1) dx \\
&= k + E[h(x, \theta)] - E[g(x)].
\end{aligned}$$

as was to be shown.

#### 4.2. Properties of the Estimators Based on Two-Tail Tests

Throughout this section the experimental model is assumed to be as described in section 4.1.1. The distribution function, expected value, and mean square error of the estimator (4.18) will be derived.

##### 4.2.1. Properties with $\sigma^2$ Known

In this subsection it is assumed that  $\sigma_{\epsilon}^2$  is known, and, therefore,  $\sigma_j^2 = \text{Var}(\hat{\beta}_j)$ , can be computed as  $\sigma_{\epsilon}^2$  times the  $j$ -th diagonal element of  $(P'P)^{-1}$ .

Since the integrated mean square error is the criterion under consideration, by Theorem 4.1 the estimators may be considered one at a time. Define:

$$\check{\beta}_j = \begin{cases} \hat{\beta}_j & \text{if } |\hat{\beta}_j| \geq c \sqrt{\text{Var}(\hat{\beta}_j)} \\ 0 & \text{else} \end{cases}, \quad (4.22)$$

The real question of interest is whether  $|\beta_j| > \sqrt{\text{var}(\hat{\beta}_j)}$ , which is equivalent to

$$\frac{|\hat{\beta}_j|}{\sqrt{\text{var}(\hat{\beta}_j)}} \geq 1. \quad (4.23)$$

Define

$$\theta = \frac{\beta_j}{\sqrt{\text{var}(\hat{\beta}_j)}}, \quad (4.24)$$

and its minimum variance unbiased estimator,

$$\hat{\theta} = \frac{\hat{\beta}_j}{\sqrt{\text{var}(\hat{\beta}_j)}} \sim N(\theta, 1).$$

(The symbol " $\sim$ " means "is distributed as".)

Now, inequality (4.23) is equivalent to

$$H_0: |\theta| \geq 1 \quad (4.25)$$

and the procedure (4.22) is motivated by the fact that the UMP level  $\alpha$  test of  $H_0$  (4.25) is (Lehman, 1959, section 3.1):

$$\phi(\hat{\beta}_j) = \begin{cases} 1 & \text{if } |\hat{\beta}_j| \geq c \sqrt{\text{Var}(\hat{\beta}_j)} \\ 0 & \text{else} \end{cases} \quad (4.26)$$

where  $C$  depends on  $\alpha$ , and  $\check{\beta}_j = \hat{\beta}_j \phi(\hat{\beta}_j)$ .

The properties of  $\check{\beta}_j$  are more easily derived in terms of the standardized variables:

$$\begin{aligned}\check{\theta}_j &= \check{\beta}_j / \sqrt{\text{Var}(\hat{\beta}_j)} \\ &= \begin{cases} \hat{\theta} & \text{if } |\hat{\theta}| \geq C \\ 0 & \text{else} \end{cases} \end{aligned} \quad (4.27)$$

$$= \hat{\theta} 1_B(\hat{\theta}) \quad (4.28)$$

where  $1_B$  is the indicator function for the set  $B = (-\infty, -C] \cup [C, +\infty)$ :

$$1_B(x) = \begin{cases} 1 & \text{if } x \in B \\ 0 & \text{if } x \notin B \end{cases}. \quad (4.29)$$

The probability distribution of  $\check{\theta}$  is mixed. Let  $F(\check{\theta}; \theta)$  denote the distribution function for a particular value of the parameter  $\theta$ ; then:

$$F(\check{\theta}; \theta) = \begin{cases} \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\check{\theta}} \exp[-(t-\theta)^2/2] dt, & \text{for } -\infty < \check{\theta} \leq -C \\ F(-C; \theta) & \text{for } -C \leq \check{\theta} < 0 \\ F(+C; \theta) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^C \exp[-(t-\theta)^2/2] dt, & 0 \leq \check{\theta} \leq C \\ \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\check{\theta}} \exp[-(t-\theta)^2/2] dt & \text{for } C \leq \check{\theta} < +\infty. \end{cases} \quad (4.30)$$

$F(\check{\theta}; \theta)$  is continuous except at  $\check{\theta} = 0$ ; it has a continuous derivative,  $f(\check{\theta}; \theta)$ , everywhere except for the set of points,  $\check{\theta} \in \{C, -C, 0\}$ .

Now consider the mean square error of  $\check{\theta}$ , which can be found by application of Lemma 4.2. In that lemma, set  $k = \theta^2$ ,  $x = \hat{\theta}$ ;  $h(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2$  so  $h(t + \theta, \theta) = t^2$ ;  $-a_1 = a_2 = C$  so  $A = [a_1, a_2] = [-C, C]$  and  $B = A^C$  as above. By Lemma 4.2, since  $g(\hat{\theta}) = (\check{\theta} - \theta)^2$ ,

$$\begin{aligned} \text{MSE} &= E[(\check{\theta} - \theta)^2 | \theta, C] = E[g(\hat{\theta})] \\ &= \text{Var}(\hat{\theta}) + \theta^2 [\Phi(C - \theta) - \Phi(-C - \theta) - \int_{-C - \theta}^{C - \theta} t^2 f_1(t; 0, 1) dt] \\ &= 1 + \theta^2 [\Phi(C - \theta) - \Phi(-C - \theta)] - \text{PSM}(-C - \theta, C - \theta) \end{aligned} \quad (4.31)$$

$$= 1 + \int_{-C - \theta}^{C - \theta} \frac{(\theta^2 - t^2) e^{-t^2/2}}{\sqrt{2\pi}} dt; \quad (4.32)$$

where  $\Phi$  denotes the cumulative distribution function and PSM denotes the partial second moment function for the normal distribution, as discussed in section 4.1.3. Although (4.32) is a compact representation of  $\text{MSE}(\theta, C)$ , equation (4.31) is suitable for numerical computations.

A graph of the MSE function for various values of  $C$  and for  $0 < \theta \leq 4$  is displayed in Figure 4.1. It should be noted that the function is symmetric about 0 in  $\theta$ :  $\text{MSE}(\theta, C) = \text{MSE}(-\theta, C)$ , so that only positive values of  $\theta$  need be considered.

It should also be noted that for  $C = 0$ ,  $\check{\theta} = \hat{\theta}$  identically (i.e.,  $\check{\beta}_j = \hat{\beta}_j$ ); this is equivalent to always using the least squares estimator for  $\beta_j$ . Therefore

$$\text{MSE}(\theta, C = 0) = \text{VAR}(\hat{\theta}) = 1, \text{ all } \theta.$$

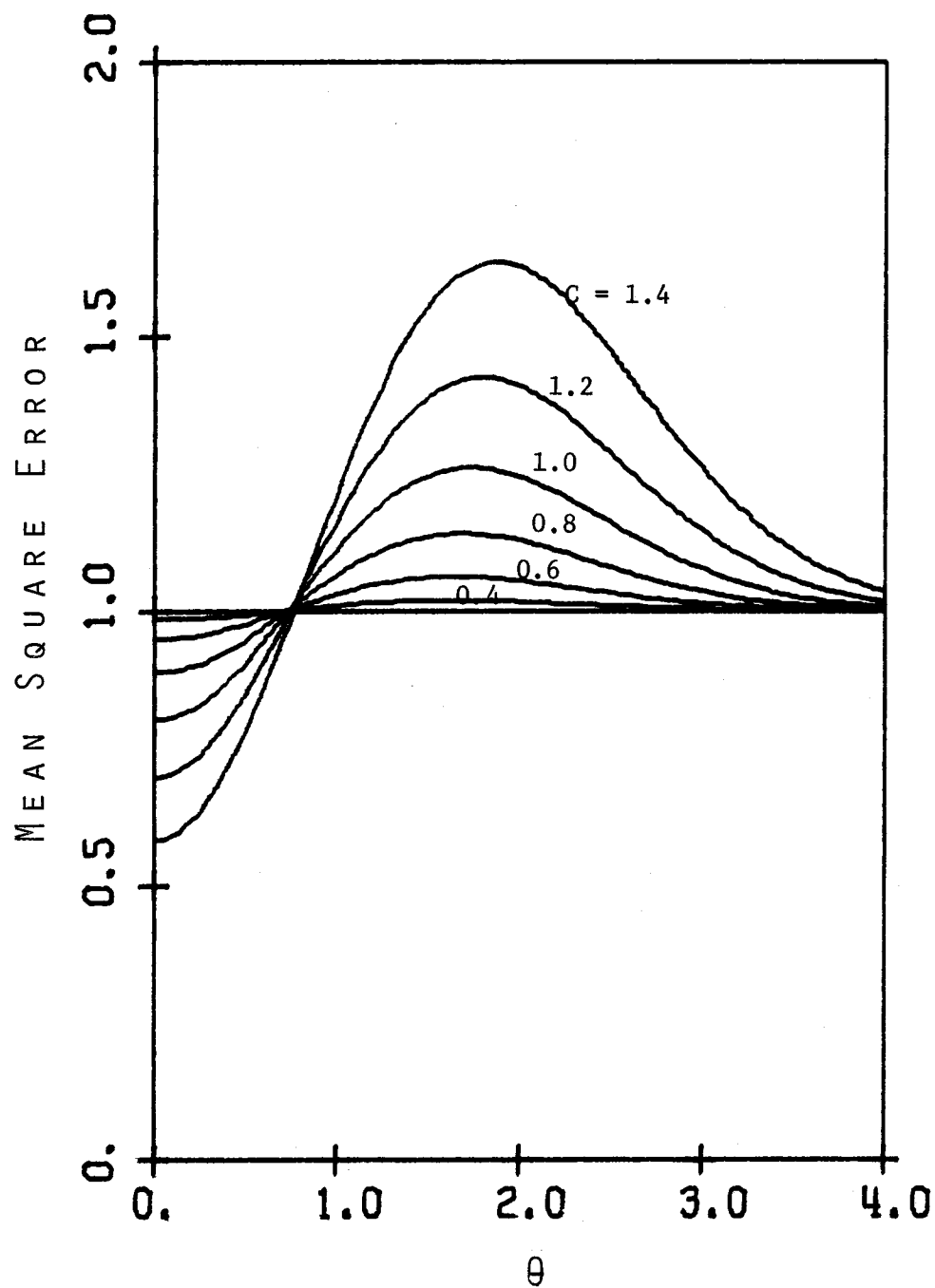


Figure 4.1 Mean Square Error of the two-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$  and  $C$ ;  $\sigma^2$  assumed known.

This straight line is also plotted for reference in Figure 4.1. Discussion of Figure 4.1 is deferred until similar figures are developed for the case where  $\sigma^2$  is unknown.

The mean square error of  $\check{\beta}_j$  can be computed from:

$$\text{MSE}(\check{\beta}_j, \beta_j) = \text{MSE}(\check{\theta}, \theta) \text{var}(\hat{\beta}_j),$$

where  $\beta_j$  is the true value and  $\theta = \beta_j / \sqrt{\text{var}(\hat{\beta}_j)}$ .

The expected value of  $\check{\theta}$  can be found by applying Lemma 4.2 with  $k = 0$ ;  $h(\hat{\theta}, \theta) = \hat{\theta}$ ;  $-a_1 = a_2 = C$ ;  $A = [-C, C] = [a_1, a_2]$ ;  $B = A^C$ , and  $g(\hat{\theta}) = \check{\theta}$ ; by Lemma 4.2,

$$\begin{aligned} E(\check{\theta} | \theta, C) &= E[g(\hat{\theta})] \\ &= \theta + 0[\Phi(C-\theta) - \Phi(-C-\theta)] - \int_{-C-\theta}^{-C-\theta} (t+\theta) f_1(t; 0, 1) dt \\ &= \theta - \theta[\Phi(C-\theta) - \Phi(-C-\theta)] - \text{PFM}(-C-\theta, C-\theta). \end{aligned} \quad (4.35)$$

where  $\Phi$  is the cumulative distribution function and PFM is the partial first moment function for the standard normal distribution, as discussed in section 4.1.3. The bias of  $\check{\theta}$  is:

$$\begin{aligned} \text{BIAS}(\theta, C) &= E(\check{\theta} - \theta) \\ &= \theta[\Phi(-C-\theta) - \Phi(C-\theta)] - \text{PFM}(-C-\theta, C-\theta). \end{aligned} \quad (4.36)$$

This function is plotted in Figure 4.2 for various  $C$  values and for  $0 \leq \theta \leq 4$ . Again it should be noted that if  $C = 0$ ,  $\check{\theta} \equiv \hat{\theta}$ , so the bias is zero. Points are not plotted for  $\theta < 0$  because of the symmetry relationship,



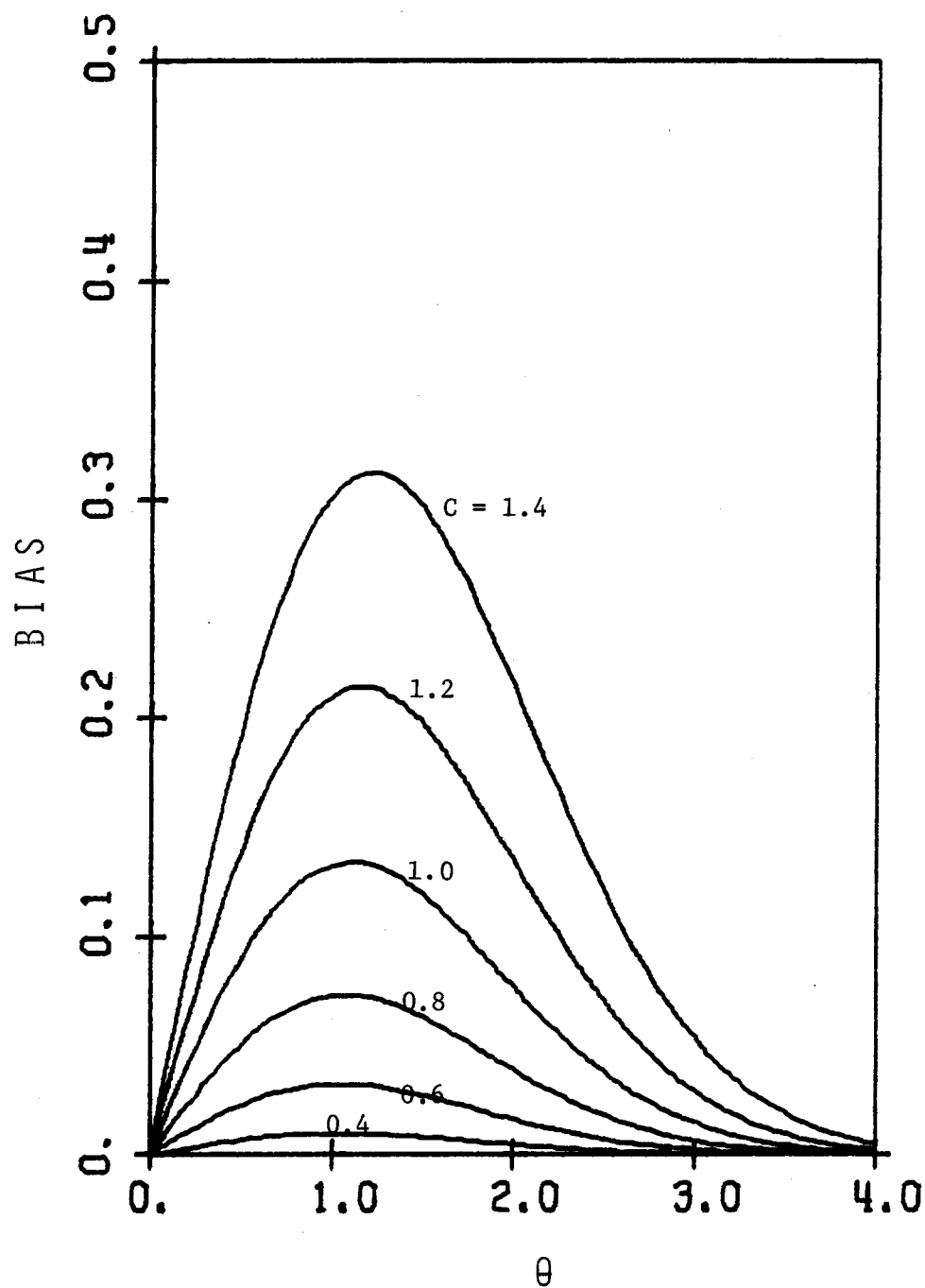


Figure 4.2 Absolute value of the bias of the two-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$  and  $C$ ;  $\sigma^2$  assumed known.

$$\text{BIAS}(-\theta, C) = -\text{BIAS}(\theta, C)$$

which can be inferred from (4.25).

To this point the distribution function, expected value, and mean square error of  $\check{\theta}$  have been displayed as functions of  $\theta$  and  $C$ . The corresponding functions for the coefficient estimators,  $\check{\beta}_j$ , may be found from the defining relation (4.24). In all of these derivations it has been assumed that  $\sigma_{\epsilon}^2$  is known (and, therefore, that  $(\text{var}(\hat{\beta}_j))$  is known), and also that the estimator  $\check{\beta}_j$  is based on a two-tail test. The next section contains an investigation of estimator based on the two-tail test when  $\sigma_{\epsilon}^2$  is estimated. The theory for estimators based on one-tail tests and more general estimators will be developed in later sections.

#### 4.2.2. Properties with $\sigma^2$ Unknown

In this section it is assumed that the experimental model is as described in section 4.1.1;  $\sigma_{\epsilon}^2$  is assumed unknown.

The procedure is essentially the same as in the previous section. For each term in the model

$$\eta(\underline{x}) = \sum_{j=1}^m \beta_j P_j(\underline{x}), \quad (4.37)$$

the hypothesis

$$H_0: |\beta_j| \leq \sqrt{\text{var}(\hat{\beta}_j)}. \quad (4.38)$$

is tested and the estimator  $\check{\beta}_j$  is defined by

$$\check{\beta}_j = \begin{cases} 0 & \text{if } H_0 \text{ is accepted} \\ \hat{\beta}_j & \text{if } H_0 \text{ is rejected} \end{cases}. \quad (4.39)$$

Consider the procedure:

$$\check{\beta}_j = \begin{cases} \hat{\beta}_j & \text{if } |\hat{\beta}_j| \geq C\sqrt{\hat{\text{var}}(\hat{\beta}_j)} \\ 0 & \text{if } |\hat{\beta}_j| < C\sqrt{\hat{\text{var}}(\hat{\beta}_j)} \end{cases}. \quad (4.40)$$

Since the same procedure will be applied to each term in the model (with possibly different values for  $C$ ), the notation can be simplified. Let:

- (a)  $\hat{\beta} \sim N(\beta, \sigma^2)$ ; i.e.,  $\sigma^2$  denotes  $\text{var}(\hat{\beta})$ ;
- (b)  $\theta = \beta/\sigma$ ; (4.41)
- (c)  $\hat{\theta} = \hat{\beta}/\sigma \sim N(\theta, 1)$ ;
- (d)  $\hat{\sigma}^2$  be the estimator of  $\sigma^2 = \text{var}(\hat{\beta})$  based on  $\nu$  degrees of freedom;

thus:

- (e)  $\frac{\nu \hat{\sigma}^2}{\sigma^2} \sim \chi_\nu^2$  independently of  $\hat{\beta}$ . (Graybill, 1961).

Let

$$F = \frac{(\hat{\beta})^2}{\hat{\sigma}^2} \equiv \frac{(\hat{\beta})^2/\sigma^2}{\hat{\sigma}^2/\sigma^2}, \quad (4.42)$$

which has the noncentral  $F$  distribution with 1 and  $\nu$  degrees of freedom and noncentrality parameter  $\theta^2/2$ .  $H_0$  (4.38) is equivalent to

$$H_0: |\theta| \leq 1. \quad (4.43)$$

Toro and Wallace (1968) have shown that the UMP test of  $H_0$  (4.43) is given by:

$$\begin{aligned} \text{Accept } H_0: & \text{ if } F < K \\ \text{Reject } H_0: & \text{ if } F \geq K \end{aligned} \quad (4.44)$$

where the constant  $K$  depends on the degrees of freedom and the level of the test. The acceptance criterion (4.44) is equivalent to:

$$(\hat{\beta})^2 < K \hat{\sigma}^2 \Leftrightarrow |\hat{\beta}| < \sqrt{\hat{\sigma}^2} \quad (4.45)$$

which is equivalent to the procedure (4.40).

In the derivation of the mean square error and expected value of  $\check{\beta}$  it is more convenient to work with the standardized random variables  $\hat{\theta}$  and  $\check{\theta}$ , and the standardized parameter,  $\theta$ . In order to be able to apply Lemma 4.1, define:

$$u = \frac{\hat{\sigma}^2}{\sigma^2}, \quad (4.46)$$

$(vu)$  is a chi-square random variable with  $v$  degrees of freedom, distributed independently of  $\hat{\beta}$ , as previously mentioned. From the relations above,

$$\begin{aligned} E(\check{\beta} | \beta, C, v) &= E(\check{\theta} | \theta, C, v) \sqrt{\text{var}(\hat{\beta})} \\ E[(\check{\beta} - \beta)^2 | \beta, C, v] &= E[(\check{\theta} - \theta)^2 | \theta, C, v] \cdot \text{var}(\hat{\beta}). \end{aligned} \quad (4.47)$$

The mean square error and expected value of  $\check{\theta}$  will now be developed. Consider the following two subsets, A and B, of the two-dimensional sample space of  $(\hat{\theta}, u = \hat{\sigma}^2/\sigma^2)$ :

$$\begin{aligned} B &= \{(\hat{\theta}, u): 0 \leq u, |\hat{\theta}| \sigma \equiv |\hat{\beta}| \geq C \sqrt{\hat{\sigma}^2} \equiv C \sqrt{u\sigma^2}\} \\ &= \{(\hat{\theta}, u): 0 \leq u, \hat{\theta}^2 \geq C^2 u\} \\ A &= \{(\hat{\theta}, u): 0 \leq u, |\hat{\theta}| \sigma \equiv |\hat{\beta}| < C \sqrt{\hat{\sigma}^2} \equiv C \sqrt{u\sigma^2/v}\} \\ &= \{(\hat{\theta}, u): u \geq 0, \hat{\theta}^2 < C^2 u/v\}. \end{aligned} \quad (4.48)$$

Clearly A and B are disjoint and  $A + B$  is the whole sample space.

Also,

$$\check{\theta} = \hat{\theta} 1_B(\hat{\theta}, u)$$

where  $1_B$  is the indicator function of the set B. For the mean square error of  $\check{\theta}$ , at given values of  $\theta$ , C, v:

$$\text{MSE}(\theta, C, v) = E[(\check{\theta} - \theta)^2 | \theta, C, v]$$

apply Lemma 4.1 with the following correspondences:

$$x = \hat{\theta}; u = u; -a_2 = a_3 = C;$$

$$h(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2; h(t+\theta, \theta) = t^2;$$

$$k = \theta^2; g_2(\hat{\theta}, u) = (\check{\theta} - \theta)^2.$$

Define the function

$$I_M(\theta; a, b; c, d) = \int_a^b (\theta^2 - t^2) f_1(t; 0, 1) P \left[ \frac{v(t+\theta)^2}{c^2} < \chi_v^2 < \frac{v(t+\theta)^2}{d^2} \right] dt. \quad (4.49)$$

Then, from Lemma 4.1,

$$\begin{aligned} \text{MSE}(\theta, C, v) &= E[g_2(\hat{\theta}, u)] \\ &= \text{Var}(\hat{\theta}) + I_M(\theta; -\infty, -\theta; -C, 0) + I_M(\theta; -\theta, \infty; C, 0) \\ &= 1 + I_M(\theta; -\infty, +\infty; C, 0). \end{aligned} \quad (4.50)$$

To get the expectation of  $\check{\theta}$ , apply Lemma 4.1 with the correspondences:

$$x = \hat{\theta}; u = u; -a_2 = a_3 = C;$$

$$h(\hat{\theta}, \theta) = \hat{\theta}; h(t+\theta, \theta) = t+\theta;$$

$$k = 0; g_2(\hat{\theta}, u) = \check{\theta};$$

and

$$I_E(\theta; a, b; c, d) = \int_a^b - (t+\theta) f_1(t; 0, 1) P \left[ \frac{v(t+\theta)^2}{c^2} < \chi_v^2 < \frac{v(t+\theta)^2}{d^2} \right] dt. \quad (4.51)$$

From Lemma 4.1,

$$\begin{aligned} E(\check{\theta} | \theta, C, v) &= E[g_2(\hat{\theta}, u)] \\ &= \theta + I_E(\theta; -\infty, -\theta; -C, 0) - I_E(\theta; -\theta, \infty; C, 0) \\ &= \theta + I_E(\theta; -\infty, +\infty; C, 0). \end{aligned}$$

Thus, the bias in  $\check{\theta}$  at the values  $\theta, C, v$  is

$$\text{BIAS}(\theta, C, v) = E(\check{\theta} - \theta | \theta, C, v) = I_E(\theta; -\infty, \infty; C, 0). \quad (4.52)$$

The integrals (4.50) and (4.52) must be approximated by numerical methods. Due to the factor  $e^{-t^2/2}$  and the doubly-infinite range, these integrals would appear to be good candidates for Hermite quadrature (Hildebrand, 1956). However, for large degrees of freedom the cumulative chi-square distribution function is almost a step function at the mean value,  $v$ . The sharpness of the curvature of the integrand near the values

$$\frac{v(t+\theta)^2}{c^2} = v, \text{ i.e., } t = -\theta \pm C \quad (4.53)$$

cause the Hermite quadrature algorithm to fail.

Since the integrands tend to zero very quickly as  $t$  becomes large (due to the  $\exp(-t^2/2)$  and chi-square probability factors) the integrals can be adequately approximated over a relatively short (finite) interval, as follows. Let  $X_0$  be an integer such that

$$P[\chi_v^2/v > X_0^2] < \epsilon = 10^{-10}.$$

The one can use as upper and lower limits of integration:

$$\frac{t_0 + \theta^2}{C} = X_0^2, \text{ i.e., } t_0 = \pm C X_0 - \theta, \quad (4.54)$$

since the integrands are negligible outside the interval  $[-C X_0 - \theta, C X_0 - \theta]$ . Since the integrands have greatest curvature in the neighborhood of the points (4.53), the three subintervals defined by the points

$$-C X_0 - \theta, -C - \theta, C - \theta, C X_0 - \theta \quad (4.55)$$

were used. A 32-point Gaussian quadrature algorithm was applied to the integrands in (4.50) and (4.52) over each of the subintervals defined by the points (4.55). The integrals (4.50), (4.52) were approximated as the sum of the corresponding approximations over the three subintervals. All arithmetic was performed in double precision (about 15 decimal places) on the IBM 360/75. Of course, the final results are not accurate to 15 decimals.

A check on the accuracy of the approximation is available. As  $v$  becomes large the integrands become more ill-conditioned and the accuracy of the approximations decrease. However, for large  $v$ , the

functions  $MSE(\theta, C, \nu)$  and  $BIAS(\theta, C, \nu)$ ,  $\sigma^2$  unknown, tend to  $MSE(\theta, C)$  and  $BIAS(\theta, C)$ ,  $\sigma^2$  known, respectively. Convergence of the corresponding families of plotted curves convinces us that the approximations are adequate.

An additional check on the accuracy is available. The  $n$ -point Gaussian quadrature algorithm is equivalent to fitting an  $n - 1$  degree polynomial to the integrand (the polynomial equals the integrand at  $n$  selected points), and then finding the "exact" integral of the polynomial. If the polynomial very closely approximates the integrand, then the integral of the polynomial will be very close to the integral of the integrand. For several of the parameter value combinations considered, the integrand and corresponding approximating polynomial were plotted. The approximation was very close in each case; in addition the approximating polynomial oscillates about the integrand, which implies that on integration the errors of approximation tend to cancel and produce a very accurate approximate integral.

The curves computed for the mean square error of  $\check{\theta}$ ,  $MSE(\theta, C, \nu)$ , are plotted in Figures 4.3-4.7 for various  $C$ -values and for  $\nu = 2, 5, 10, 25, 50$ .

The curves computed for the bias of  $\check{\theta}$ ,  $BIAS(\theta, C, \nu)$ , are plotted in Figures 4.8-4.12 for various  $C$ -values and for  $\nu = 2, 5, 10, 25, 50$ . The actual values plotted are  $|BIAS(\theta, C, \nu)|$ .

In each case the functions are plotted only for  $\theta \geq 0$ , since the  $MSE$  function and  $|BIAS|$  function are symmetric about zero. However,  $BIAS(-\theta, C, \nu) = -BIAS(\theta, C, \nu)$ .

The mean square error and expectation of the original variable,  $\check{\beta}$ , can be computed from the relations (4.47).



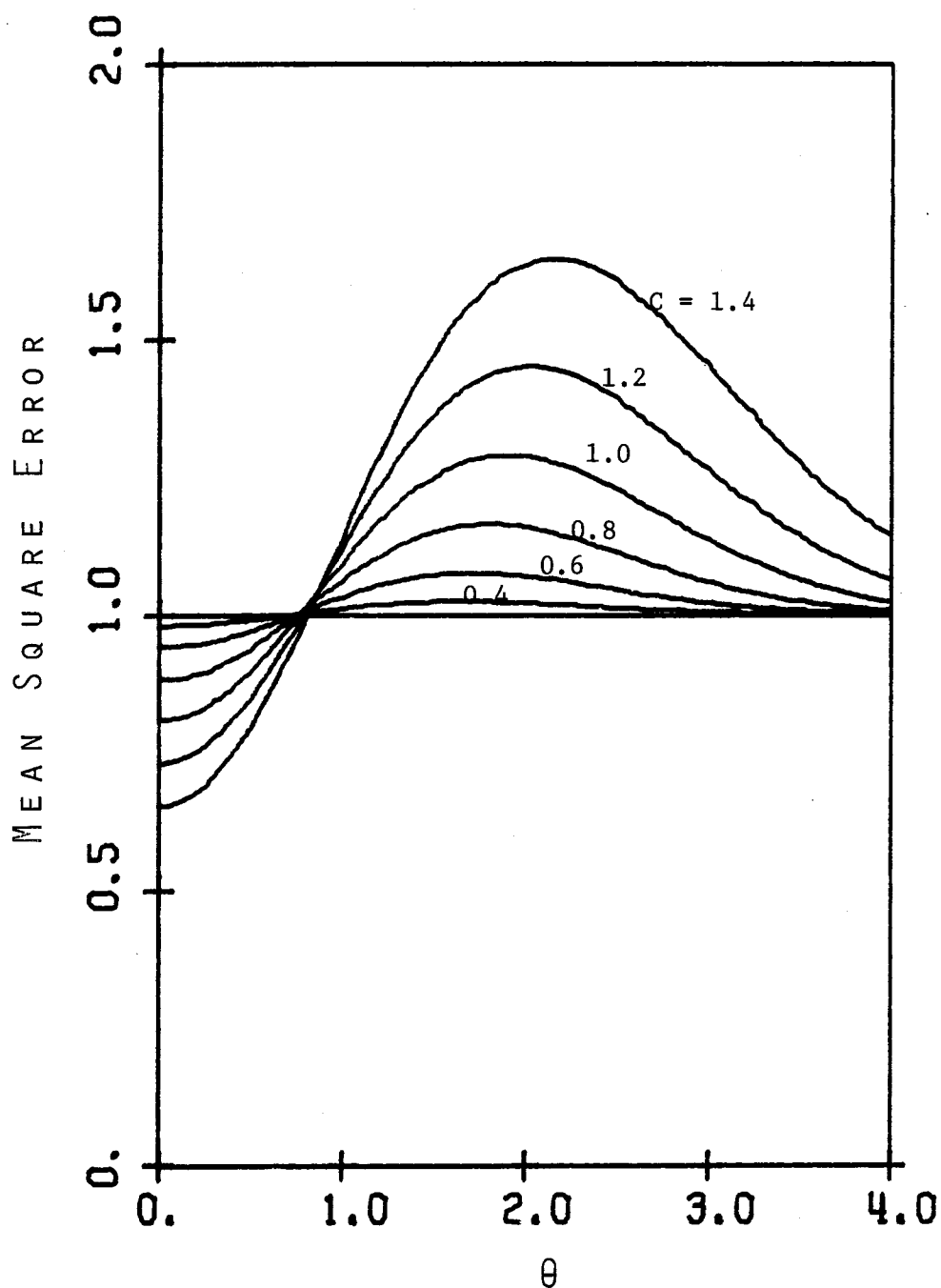


Figure 4.3 Mean Square Error of the two-tail estimator  $\hat{\theta} = \hat{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 2$ ,  $\sigma^2$  estimated.

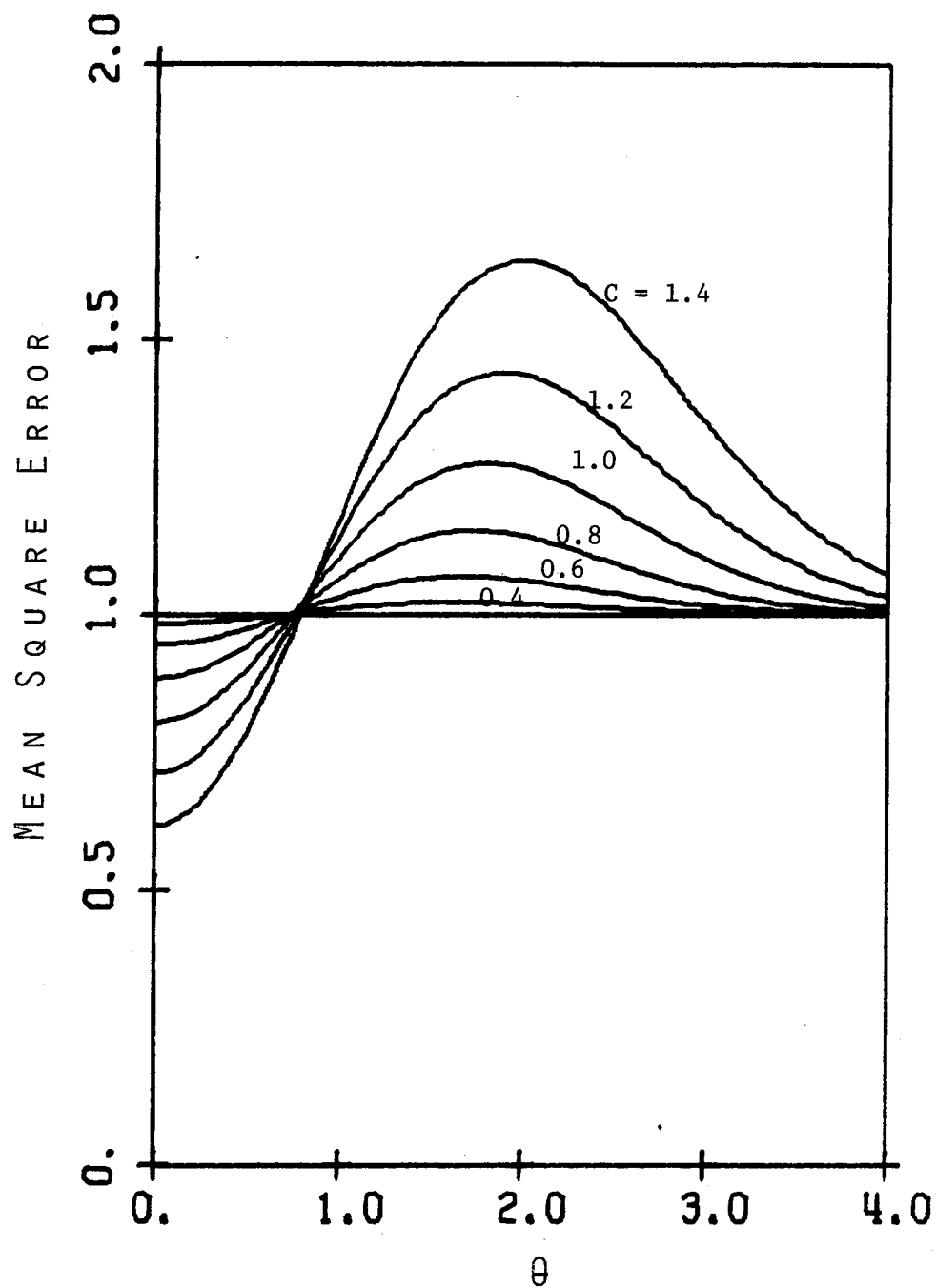


Figure 4.4 Mean Square Error of the two-tail estimator  $\hat{\theta} = \hat{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 5$ ,  $\sigma^2$  estimated.

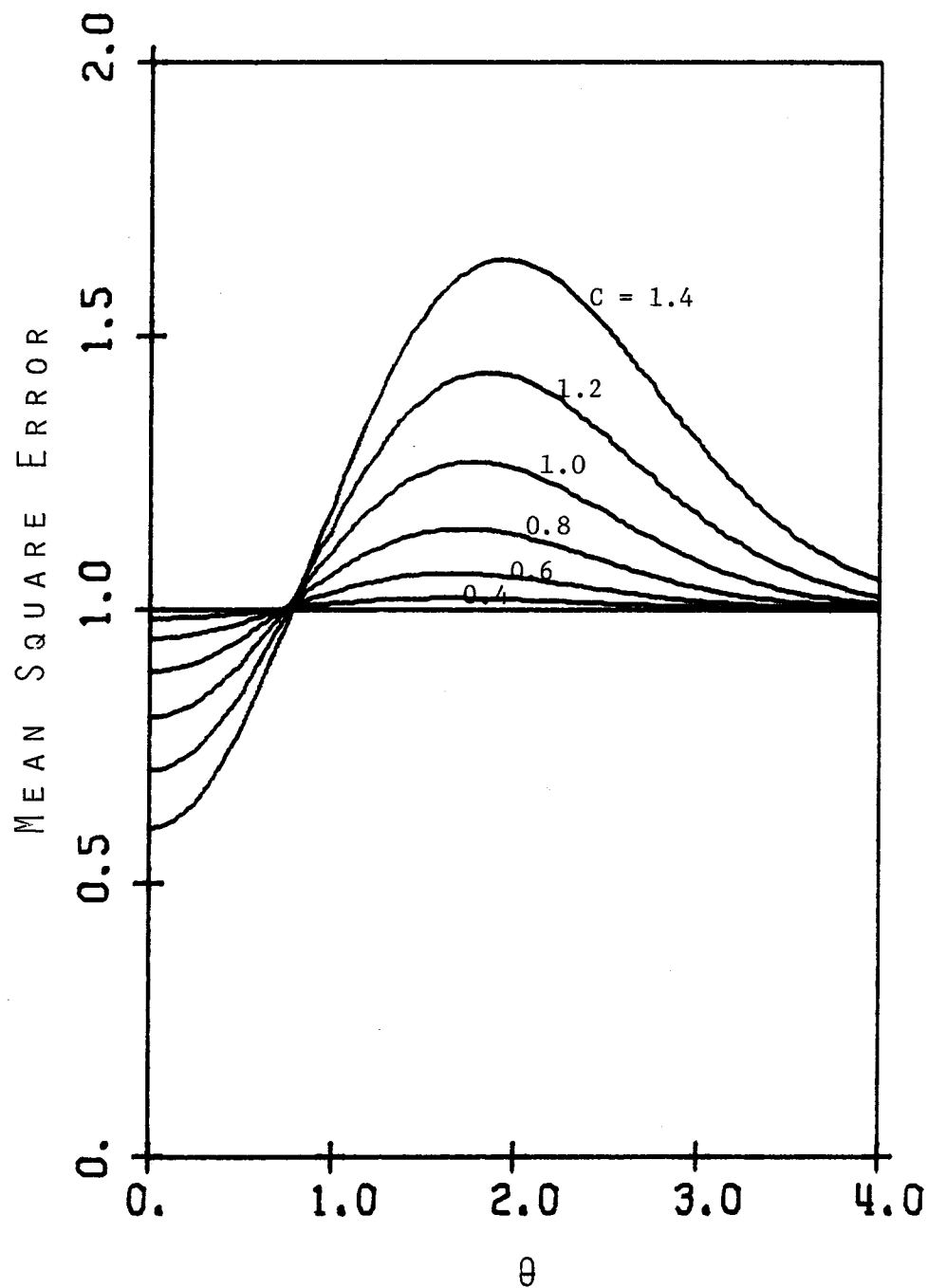


Figure 4.5 Mean Square Error of the two-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 10$ ,  $\sigma^2$  estimated.

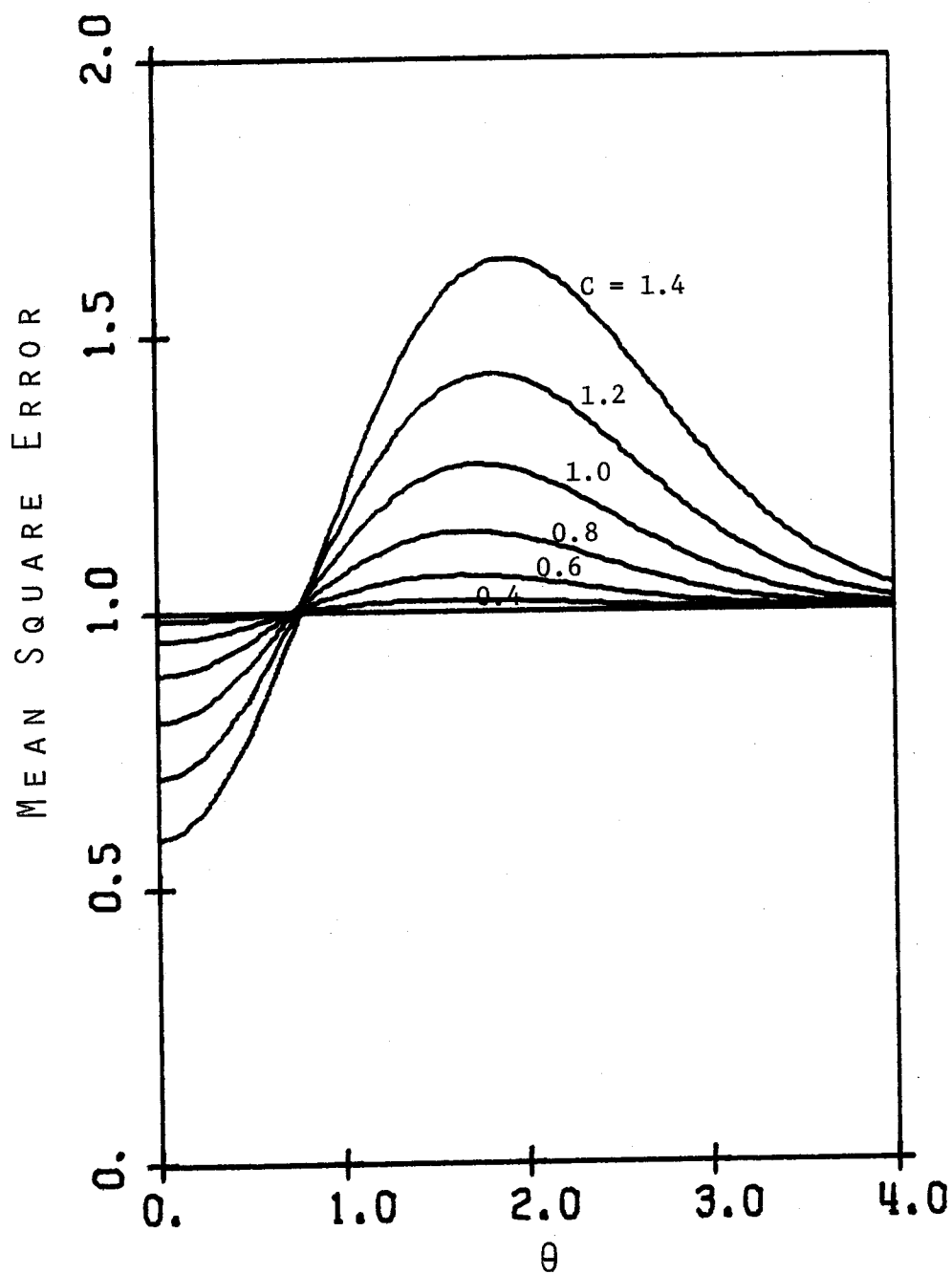


Figure 4.6 Mean Square Error of the two-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $v = 25$ ,  $\sigma^2$  estimated.

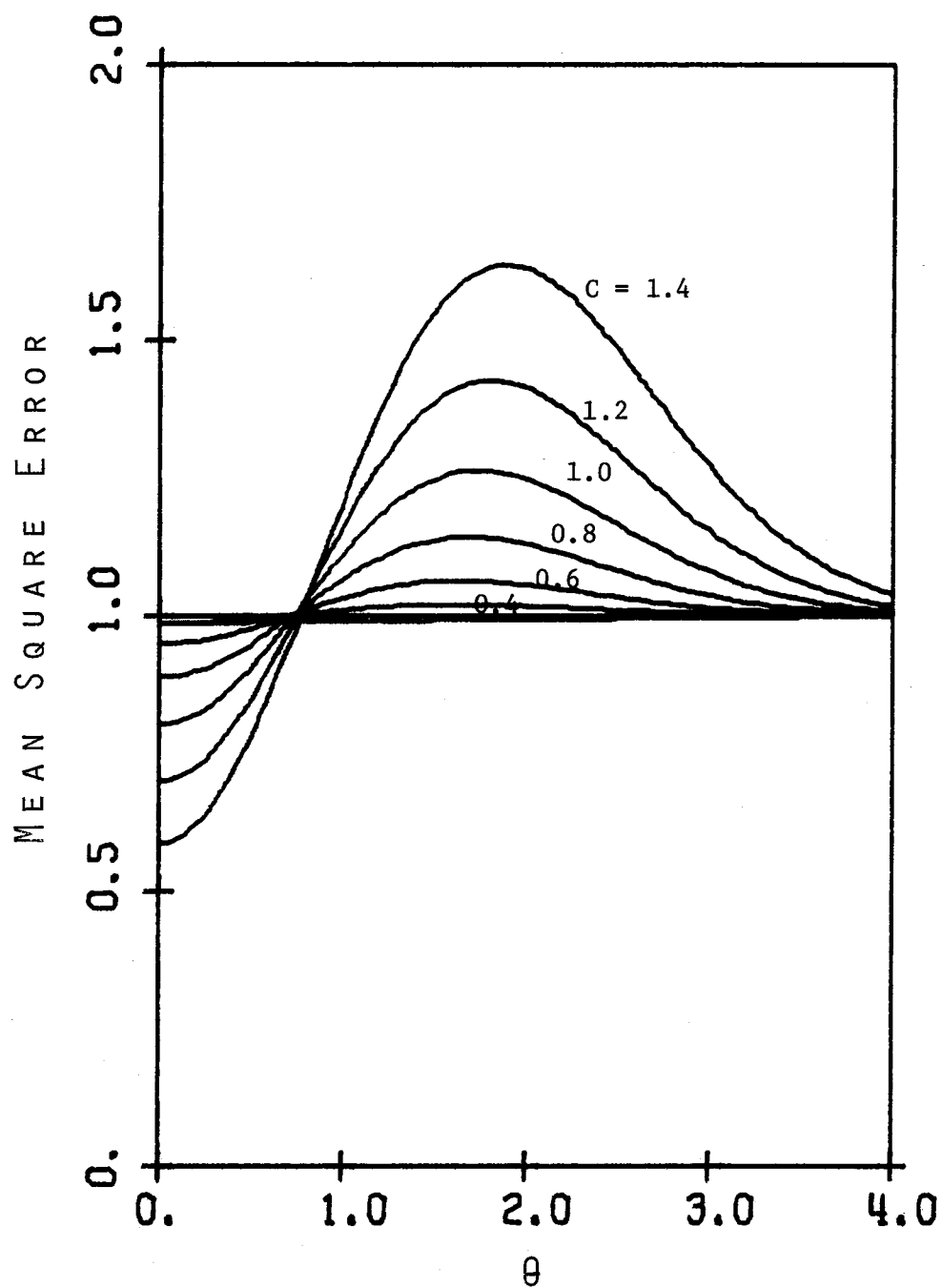


Figure 4.7 Mean Square Error of the two-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 50$ ,  $\sigma^2$  estimated.

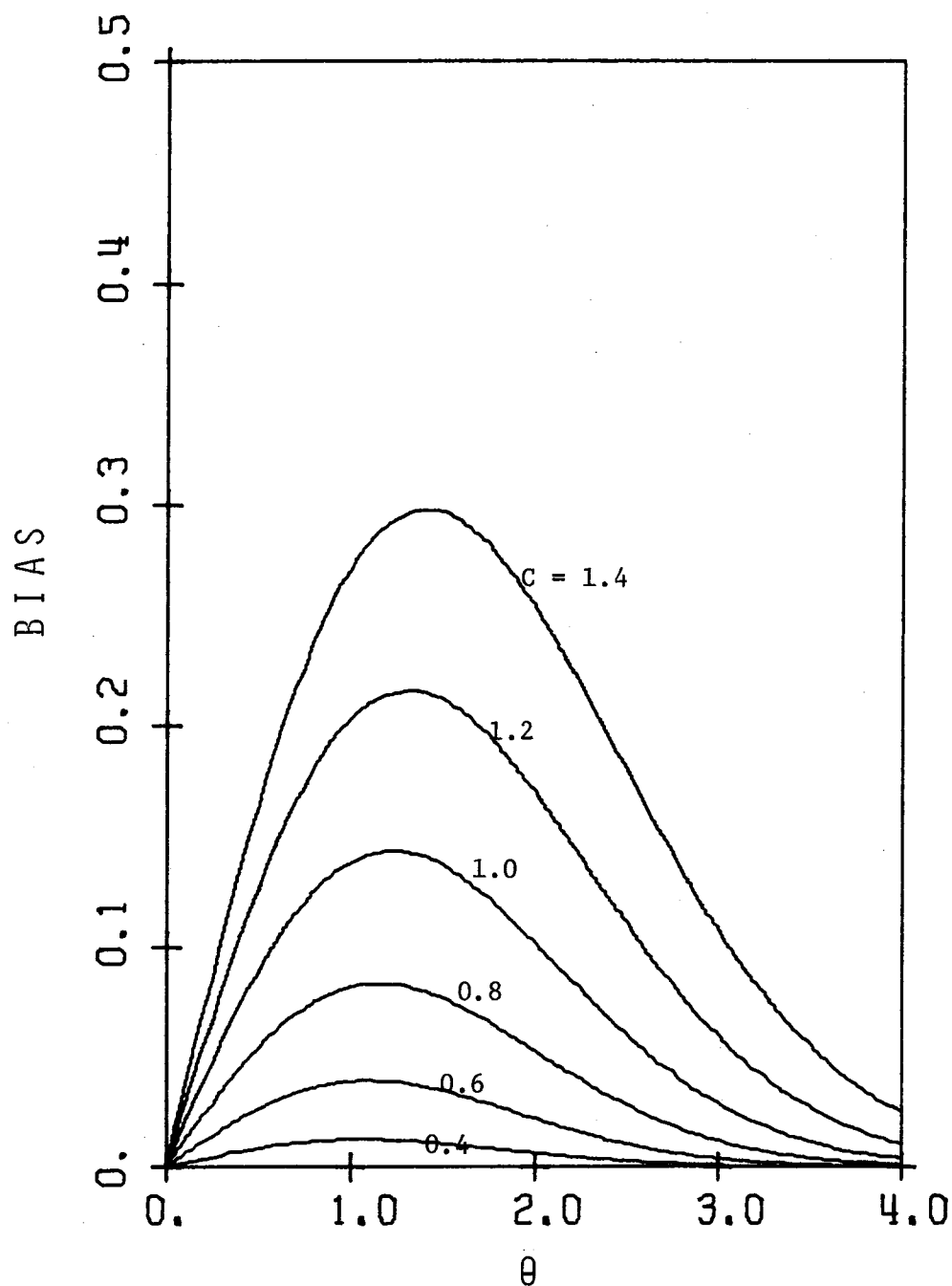


Figure 4.8 Absolute value of the bias of the two-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 2$ ,  $\sigma^2$  estimated.

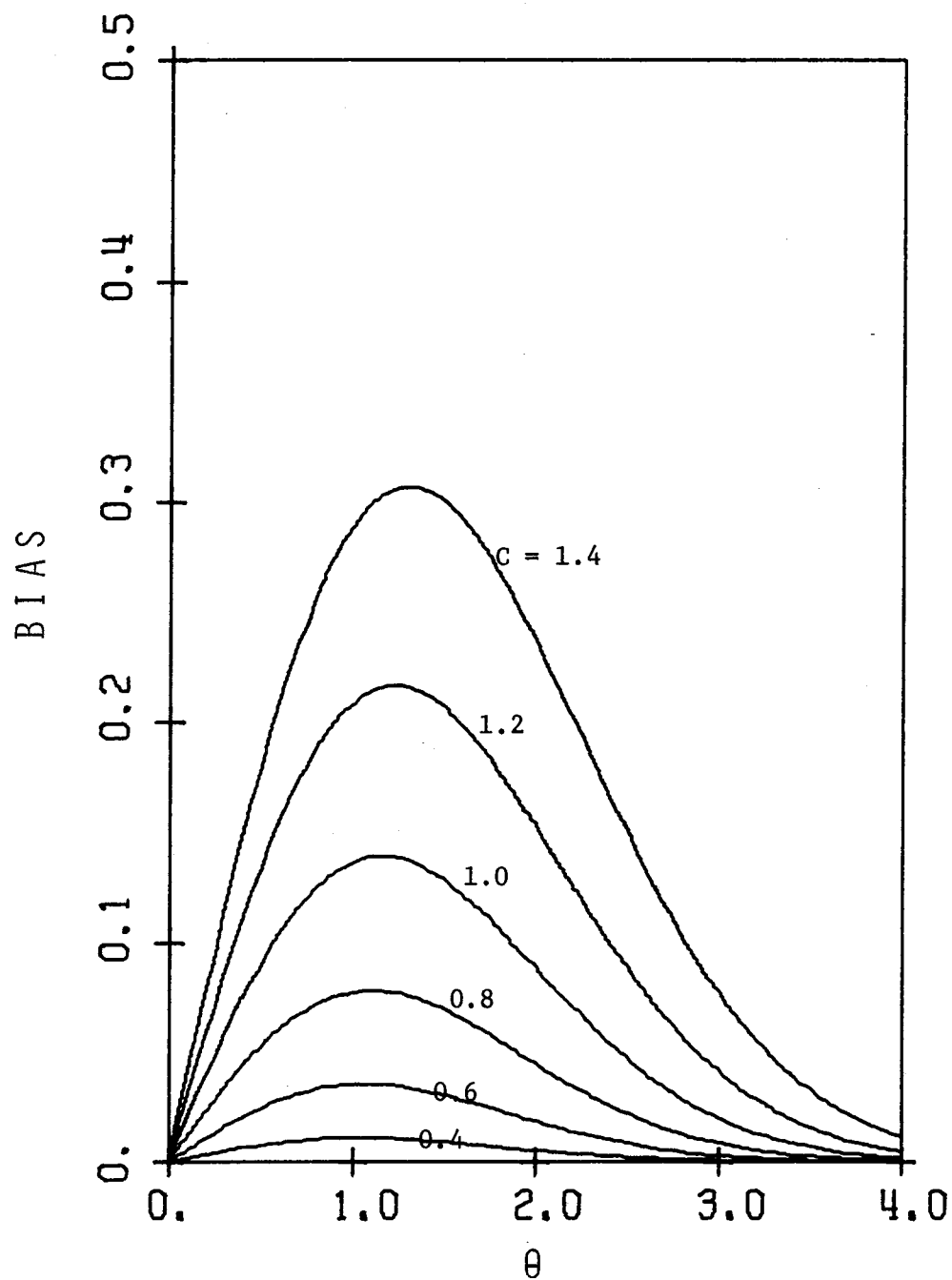


Figure 4.9 Absolute value of the bias of the two-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 5$ ,  $\sigma^2$  estimated.

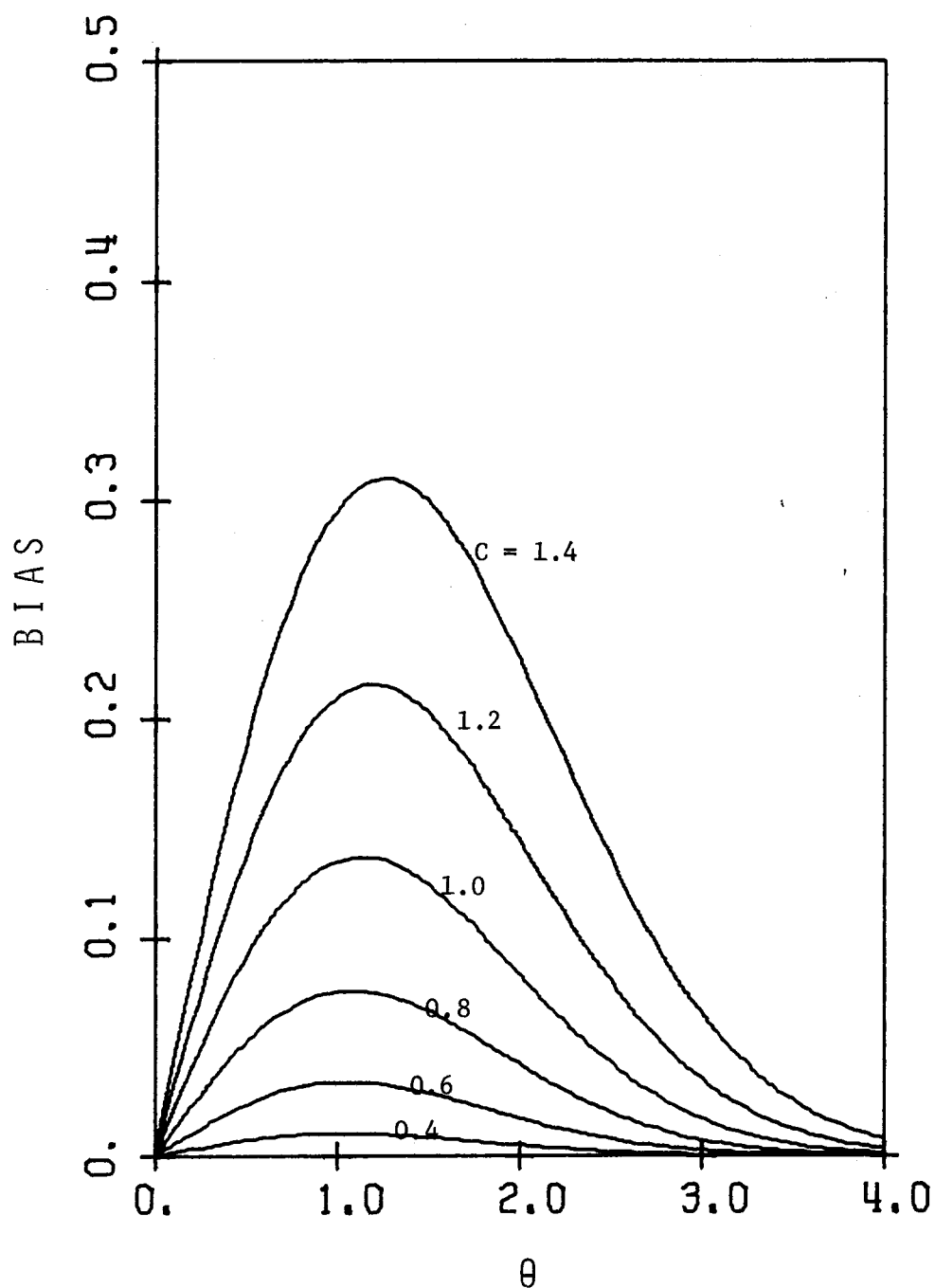


Figure 4.10 Absolute value of the bias of the two-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 10$ ,  $\sigma^2$  estimated.



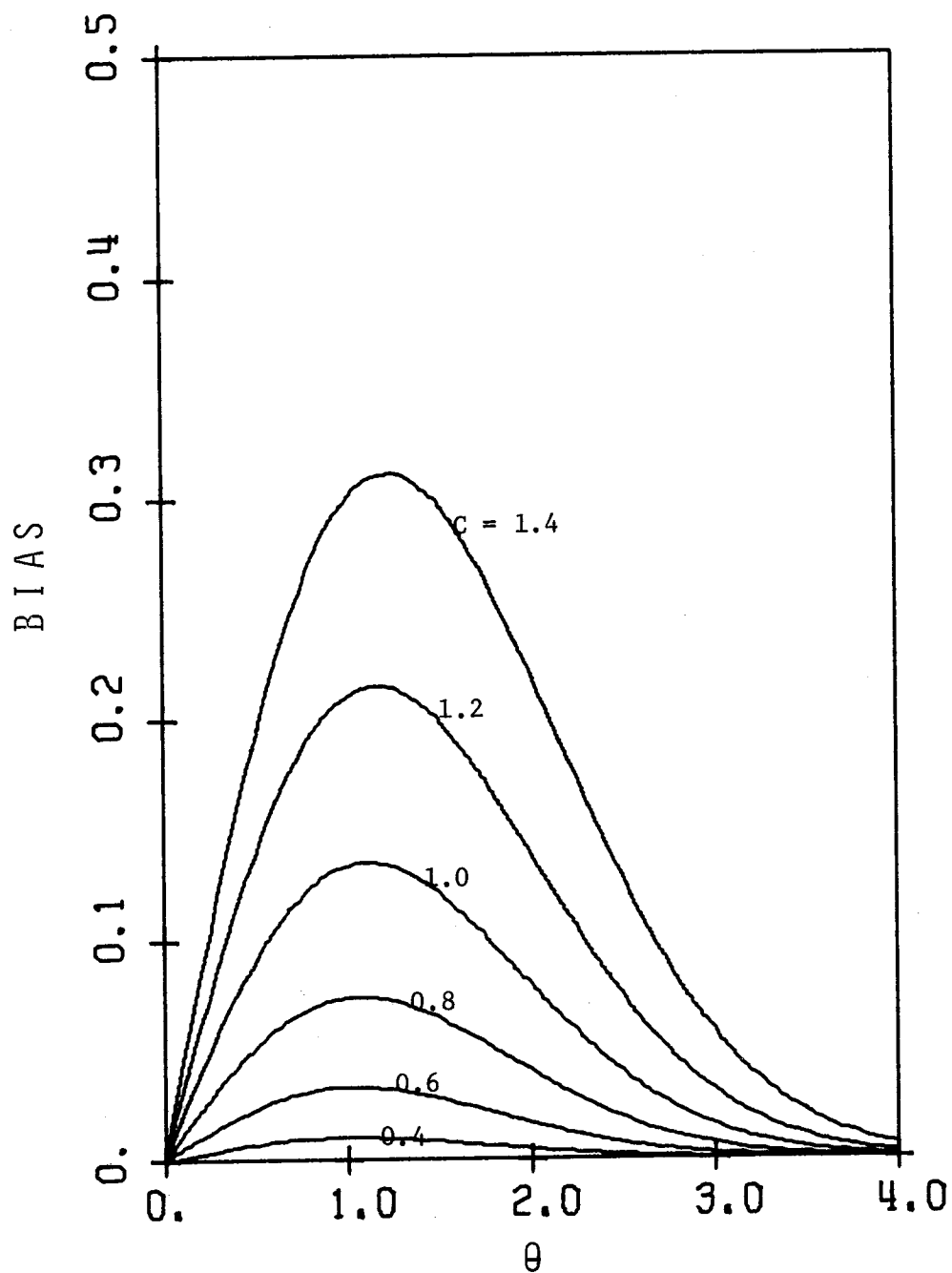


Figure 4.11 Absolute value of the bias of the two-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $v = 25$ ,  $\sigma^2$  estimated.

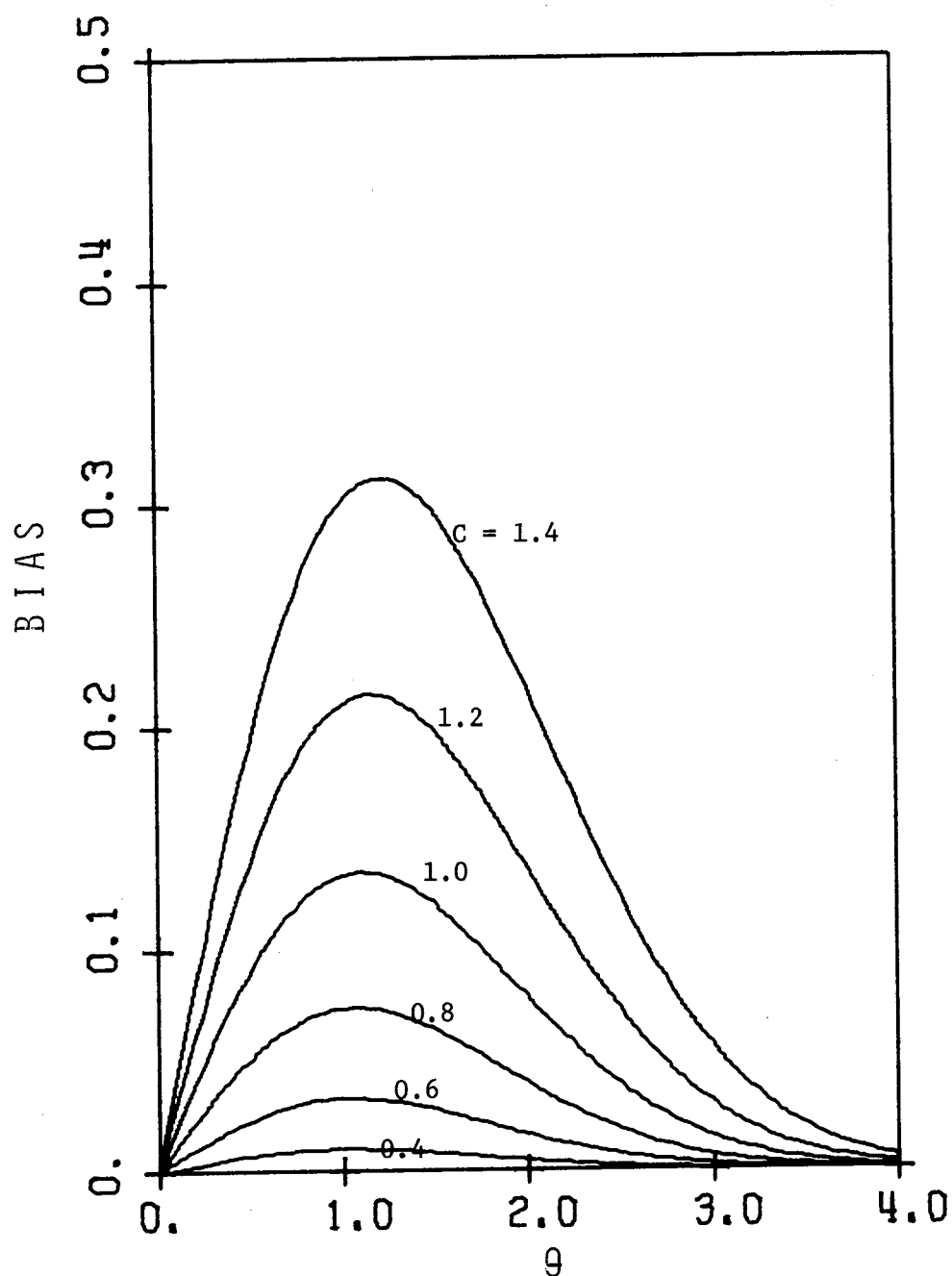


Figure 4.12 Absolute value of the bias of the two-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 50$ ,  $\sigma^2$  estimated.

### 4.3. Estimators Based on One-Tail Tests

To this point estimators based on two-tail tests have been considered. Since an investigator often knows the algebraic sign of certain of the coefficients, procedures which utilize this information will be investigated. Throughout this chapter the experimental model is assumed to be as described in section 4.1.1.

The motivating idea continues to be that setting the estimator to  $\check{\beta}_j = 0$  whenever  $|\beta_j|^2 < \text{Var}(\hat{\beta}_j)$ , and to  $\check{\beta}_j = \hat{\beta}_j$ , the least squares estimator, when  $\beta_j^2 > \text{Var}(\hat{\beta}_j)$ . Since the sign of the coefficient is known, it can be assumed without loss of generality that  $\beta_j \geq 0$ , for if  $\beta_j \leq 0$ , one can consider  $-\beta_j$ .

The procedure is to test the hypothesis

$$\begin{aligned} H_0: \beta_j &\leq \sqrt{\text{Var}(\hat{\beta}_j)} \\ \text{vs. } H_a: \beta_j &> \sqrt{\text{Var}(\hat{\beta}_j)} \end{aligned} \quad (4.59)$$

and use the estimator:

$$\check{\beta}_j = \begin{cases} \hat{\beta}_j & \text{if } H_0 \text{ is rejected} \\ 0 & \text{if } H_0 \text{ is accepted} \end{cases}$$

For notational simplicity in the following derivations, let:

$$\hat{\beta} \text{ denote } \hat{\beta}_j$$

$$\sigma^2 \text{ denote } \text{Var}(\hat{\beta}_j)$$

so that:

$$\hat{\beta} \sim N(\beta, \sigma^2),$$

$$\theta = \beta/\sigma,$$

$$\hat{\theta} = \hat{\beta}/\sigma \sim N(\theta, 1),$$

$\hat{\sigma}^2$  is an unbiased estimator of  $\sigma^2$  based on  $\nu$  degrees of freedom,

$$\frac{\nu \hat{\sigma}^2}{\sigma^2} \sim \chi_{\nu}^2, \text{ independently of } \hat{\beta}.$$

The hypothesis (4.59), in terms of the above notation, becomes:

$$H_0: \theta \leq 1$$

$$H_a: \theta > 1 \quad (4.60)$$

There are two situations, according to whether  $\sigma^2$  is known. When  $\sigma^2$  is known, the test is based on the statistic  $\hat{\theta}$ . It is well known that the normal probability density has monotone likelihood ratio (for variation in the mean,  $\theta$ ); application of Lehman's Theorem 2 (Lehman, 1959, p. 68) yields the UMP test (for a given  $\alpha$  level) of the hypothesis (4.60):

$$\phi(\hat{\theta}) = \begin{cases} 1 & \text{if } \hat{\theta} > c \\ 0 & \text{if } \hat{\theta} < c \end{cases}$$

where  $C$  is determined from  $\alpha$ .

If  $\sigma^2$  is unknown, consider the statistic:

$$t = \frac{\hat{\beta}}{\sqrt{\hat{\sigma}^2}} \equiv \frac{\hat{\beta}/\sigma}{\sqrt{\hat{\sigma}^2/\sigma^2}} = \frac{\hat{\theta}}{\sqrt{\hat{\sigma}^2/\sigma^2}}$$

which has the noncentral  $t$  distribution with noncentrality parameter  $\theta$ . Lehman (1959, p. 233) shows that this distribution has the monotone likelihood ratio property (for variation in  $\theta$ ), so another application of Lehmann's Theorem 2 (Lehmann, 1959, p. 68) yields the UMP test of the hypothesis (4.60) as:

$$\phi(t) = \begin{cases} 1 & \text{if } t > C \\ 0 & \text{if } t < C \end{cases}$$

where  $C$  is determined from  $\alpha$  and the noncentral  $t$  distribution. For simplicity, rewrite the test above in the form:

$$\phi(t) = \begin{cases} 1 & \text{if } \hat{\theta} > C \sqrt{\hat{\sigma}^2/\sigma^2} \\ 0 & \text{if } \hat{\theta} < C \sqrt{\hat{\sigma}^2/\sigma^2} \end{cases}$$

Thus, in either case ( $\sigma^2$  known or not), the estimator for  $\theta$  based on the UMP one-tail test can be written:

$$\begin{aligned} \check{\theta} &= \begin{cases} \hat{\theta} & \text{if } \hat{\theta} > C \sqrt{\hat{\sigma}^2/\sigma^2} \\ 0 & \text{if } \hat{\theta} < C \sqrt{\hat{\sigma}^2/\sigma^2} \end{cases} \\ &= \hat{\theta} 1_B(\hat{\theta}) \end{aligned}$$

where

$$\begin{aligned} B &= \{\hat{\theta}, \hat{\sigma}^2 > 0: \hat{\theta} > C \sqrt{\hat{\sigma}^2/\sigma^2}\} \\ &= \{\hat{\theta} > 0, \hat{\sigma}^2 > 0: \frac{\sqrt{\hat{\theta}^2}}{C^2} > \frac{\sqrt{\hat{\sigma}^2}}{\sigma^2}\} \end{aligned}$$

and where  $\hat{\sigma}^2$  is replaced by  $\sigma^2$ , when known. This also yields the estimator

$$\check{\beta} = \hat{\beta} 1_B(\hat{\beta})$$

where

$$B = \{\hat{\beta}, \hat{\sigma}^2 > 0: \hat{\beta} > c \sqrt{\hat{\sigma}^2}\}$$

with  $\hat{\sigma}^2$  replaced by  $\sigma^2$ , if known.

The properties of the estimators may now be derived, depending on whether  $\sigma^2$  is known.

#### 4.3.1. Properties of the One-Tail Estimators When $\sigma^2$ is Known

For notational convenience the expected value and mean square error of  $\check{\theta}$  will be evaluated. Let  $f(\hat{\theta}; \theta)$  denote the  $N(\theta, 1)$  density and consider:

$$B = \{\hat{\theta}: \hat{\theta} > c\}$$

$$\text{MSE}(\theta, c) = E[(\check{\theta} - \theta)^2 | \theta, c]$$

$$= E\{[\hat{\theta} 1_B(\hat{\theta}) - \theta]^2 | \theta, c\}$$

$$= \int_{-\infty}^{\infty} [\hat{\theta} 1_B(\hat{\theta}) - \theta]^2 f(\hat{\theta}; \theta) d\hat{\theta}$$

$$= \int_c^{\infty} (\hat{\theta} - \theta)^2 f(\hat{\theta}; \theta) d\hat{\theta} + \theta^2 \int_{-\infty}^c f(\hat{\theta}; \theta) d\hat{\theta}$$

Application of the change of variable  $t = \hat{\theta} - \theta$  yields:

$$\text{MSE}(\theta, c) = \text{PSM}(c - \theta, +\infty) + \theta^2 \Phi(c - \theta) \quad (4.61)$$

where  $\Phi$  denotes the cumulative distribution function and PSM denotes the partial second moment function, discussed in section 4.1.3.

Similarly,

$$\begin{aligned}
 E[\check{\theta}|\theta, C] &= \int_{-\infty}^{\infty} \hat{\theta} 1_B(\hat{\theta}) f(\hat{\theta}; \theta) d\hat{\theta} = \int_C^{\infty} \hat{\theta} f(\hat{\theta}) d\hat{\theta} \\
 &= \theta - \int_{-\infty}^C \hat{\theta} f(\hat{\theta}; \theta) d\hat{\theta} = \theta - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{C-\theta} (t+\theta) \exp(-t^2/2) dt \\
 &= \theta - \theta \Phi(C - \theta) - \text{PFM}(-\infty, C - \theta) \quad (4.62)
 \end{aligned}$$

where PFM denotes the partial first moment function for the standard normal density, discussed in section 4.1.3. Thus the bias of  $\check{\theta}$  at  $\theta$ ,  $C$  is

$$\text{BIAS}(\theta, C) = -\text{PFM}(-\infty, C - \theta) - \theta \Phi(C - \theta). \quad (4.63)$$

The MSE and BIAS functions are plotted for various  $C$ -values in Figures 4.13 and 4.14 respectively. Discussion of these Figures is deferred until corresponding results are developed for the unknown variance case.

#### 4.3.2. Properties of One-Tail Estimators When Variance is Unknown

The mean square error of  $\check{\theta}$  can be found by applying Lemma 4.1 with the following correspondences:

$$x = \hat{\theta}; u = u; a_2 = -\infty; a_3 = C;$$

$$h(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2; h(t + \theta, \theta) = t^2;$$

$$k = \theta^2; \text{ and } g_2(\hat{\theta}, u) = (\check{\theta} - \theta)^2.$$

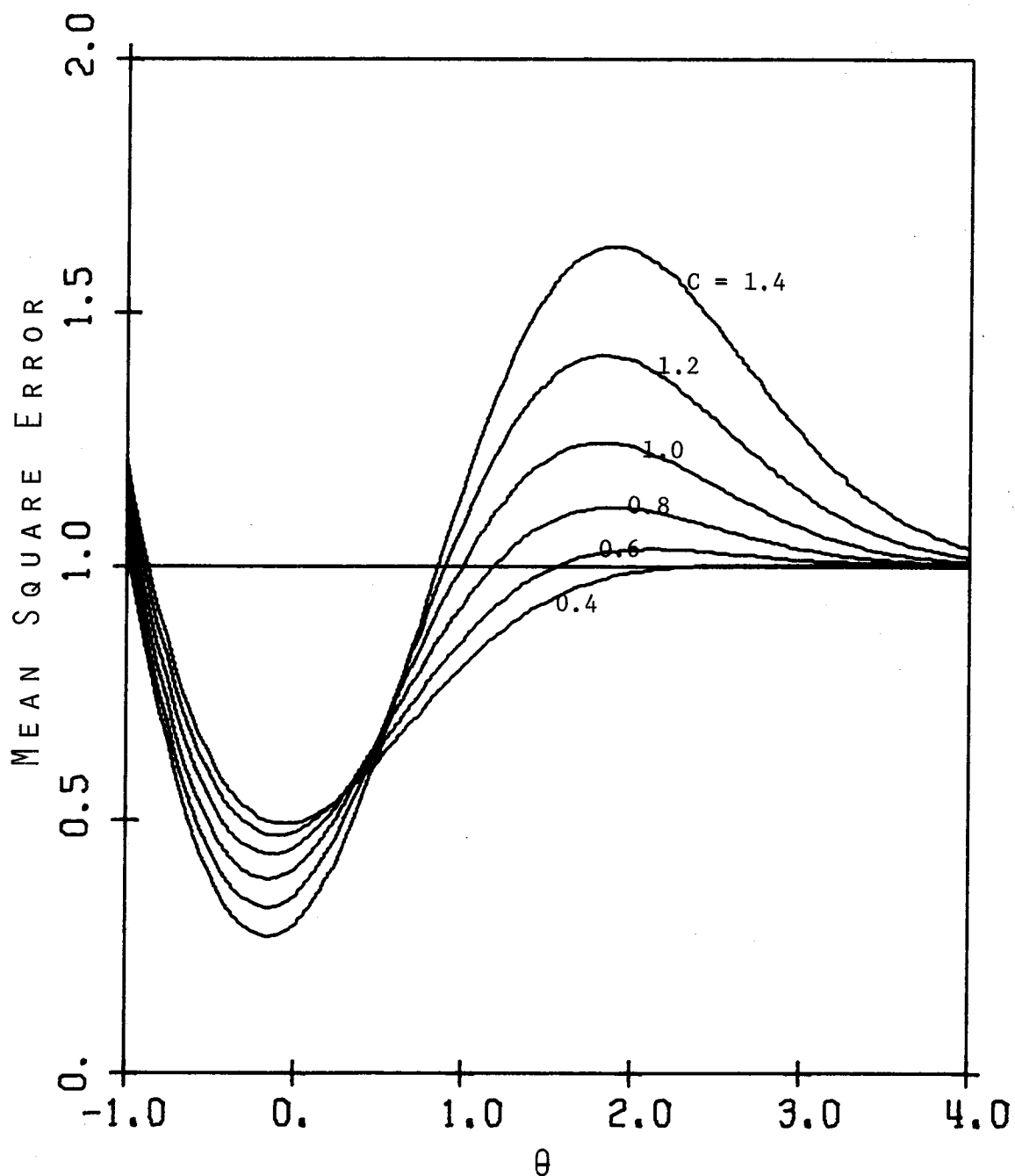


Figure 4.13 Mean Square Error of the one-tail estimator  
 $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$  and  $C$ ;  
 $\sigma^2$  assumed known.



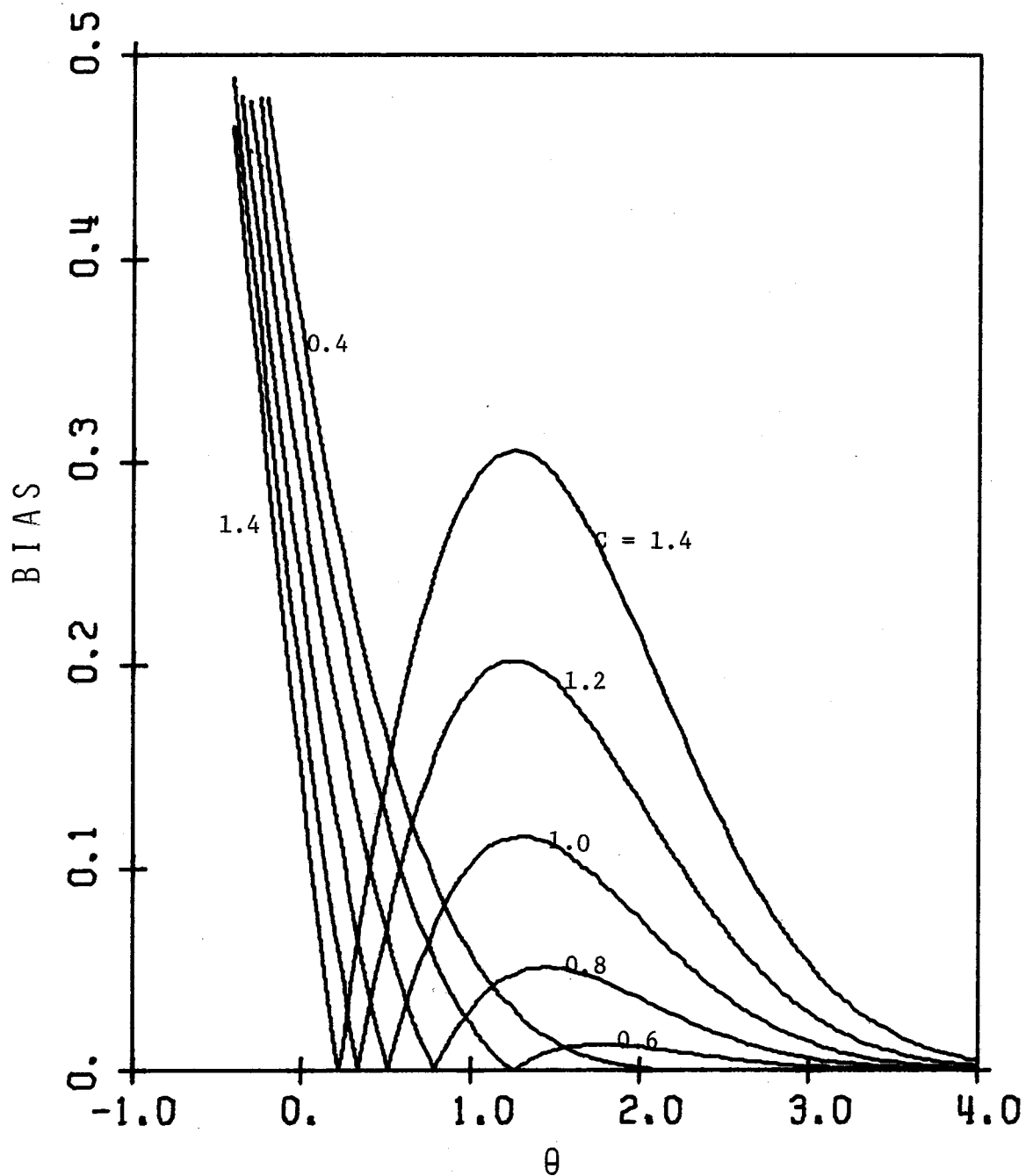


Figure 4.14 Absolute value of the bias of the one-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$  and  $C$ ;  $\sigma^2$  assumed known.

Let  $I_M(\theta; a, b; c, d)$  be defined as in the derivation of the MSE of the two-tail estimator (equation 4.49). Then, by Lemma 4.1,

$$\begin{aligned} \text{MSE}(\theta, C, v) &= E[(\check{\theta} - \theta)^2 | \theta, C, v] = E[g_2(\hat{\theta}, u)] \\ &= 1 + I_M(\theta; -\infty, -\theta; -\infty, 0) + I_M(\theta; -\theta, \infty; C, 0). \end{aligned}$$

The chi-square probability in the expression  $I_M(\theta; -\infty, -\theta; -\infty, 0)$  is identically 1.0; thus the term has value

$$\theta^2 [\Phi(-\theta) - \Phi(-\infty)] - \text{PSM}(-\infty, -\theta).$$

Hence,

$$\begin{aligned} \text{MSE}(\theta, C, v) &= 1 + \theta^2 \Phi(-\theta) - \text{PSM}(-\infty, -\theta) \\ &\quad + I_M(\theta; -\theta, \infty; C, 0). \end{aligned} \tag{4.63}$$

The expected value of  $\check{\theta}$  can also be found by application of Lemma 4.1, with the following correspondences:

$$\begin{aligned} x &= \hat{\theta}; u = u; a_2 = -\infty; a_3 = C; \\ h(\hat{\theta}, \theta) &= \hat{\theta}; h(t+\theta, \theta) = t+\theta; k = 0; \\ g_2(\hat{\theta}, u) &= \check{\theta}. \end{aligned}$$

Let  $I_E(\theta; a, b; c, d)$  be defined as at equation (4.51) in the derivation of the expected value of the two-tail estimator; then by Lemma 4.1,

$$\begin{aligned} E(\check{\theta} | \theta, C, v) &= E[g_2(\hat{\theta}, u)] \\ &= \theta + I_E(\theta; -\infty, -\theta; -\infty, 0) + I_E(\theta; -\theta, \infty; C, 0). \end{aligned}$$

Because the chi-square probability is identically 1.0 in the second term,

$$I_E(\theta; -\infty, -\theta; -\infty, 0) = -\theta \Phi(-\theta) - \text{PFM}(-\infty, -\theta).$$

Thus,

$$\begin{aligned} E(\check{\theta} | \theta, C, \nu) &= \theta - \theta \Phi(-\theta) - \text{PFM}(-\infty, -\theta) \\ &\quad + I_E(\theta; -\theta, \infty; C, 0). \end{aligned} \quad (4.64)$$

The last terms in (4.63) and (4.64) must be approximated numerically. Since the integrands of these terms are identical to the corresponding integrands in the two-tail case similar techniques were used for evaluation and similar remarks on the accuracy of the approximation apply.

Specifically, the integrals were approximated over the interval  $[-\theta, C X_0 - \theta]$ , where  $X_0$  is an integer satisfying.

$$P[\chi_{\nu}^2/\nu > X_0] \leq \epsilon = 10^{-10}.$$

As before, for large  $\nu$  the integrands have greatest curvature in the neighborhood of:

$$\nu \left( \frac{\theta + t}{C} \right)^2 = \nu, \text{ i.e., } t = \pm C - \theta.$$

The integrals were approximated by a 32-point Gaussian quadrature technique (double precision arithmetic) over each of the two intervals  $[-\theta, C - \theta]$ ,  $[C - \theta, C X_0 - \theta]$ ; the sum of the integrals over the subintervals was taken as the approximation for the integral over the whole interval  $[-\theta, +\infty]$ .

The remarks on the accuracy of the approximation in the two-tail case apply here also.

Graphs of the MSE and  $|\text{BIAS}|$  functions for various C-values, for  $v = 2, 5, 10, 25, 50$  and for  $-1.0 \leq \theta \leq 4.0$  are given in Figures 4.15-4.19 and 4.20-4.24 respectively.

#### 4.4. Generalized Estimators

In some situations it may be desirable to use the generalized estimator,

$$\check{\beta}_j = \begin{cases} 0 & \text{if } \beta_j \in A^* \\ \beta_j & \text{if } \hat{\beta}_j \in B^* = A^{*c} \end{cases}$$

where  $A^*$  is a subset of the sample space and  $B^*$  is the complement of  $A^*$  relative to the sample space. The properties of this estimator will be developed for the cases in which  $A^*$  is a finite interval, the complement of a finite interval, or a half-line. If the variance is assumed unknown, the endpoints of  $A^*$  will be random.

The properties of the estimator will be developed in terms of the standardized variables

$$\hat{\theta} = \hat{\beta}_j / \sqrt{\text{Var}(\hat{\beta}_j)}$$

and

$$\check{\theta} = \check{\beta}_j / \sqrt{\text{Var}(\check{\beta}_j)}.$$

The properties of this  $\check{\theta}$  will be developed along the same lines as derivations in previous sections.

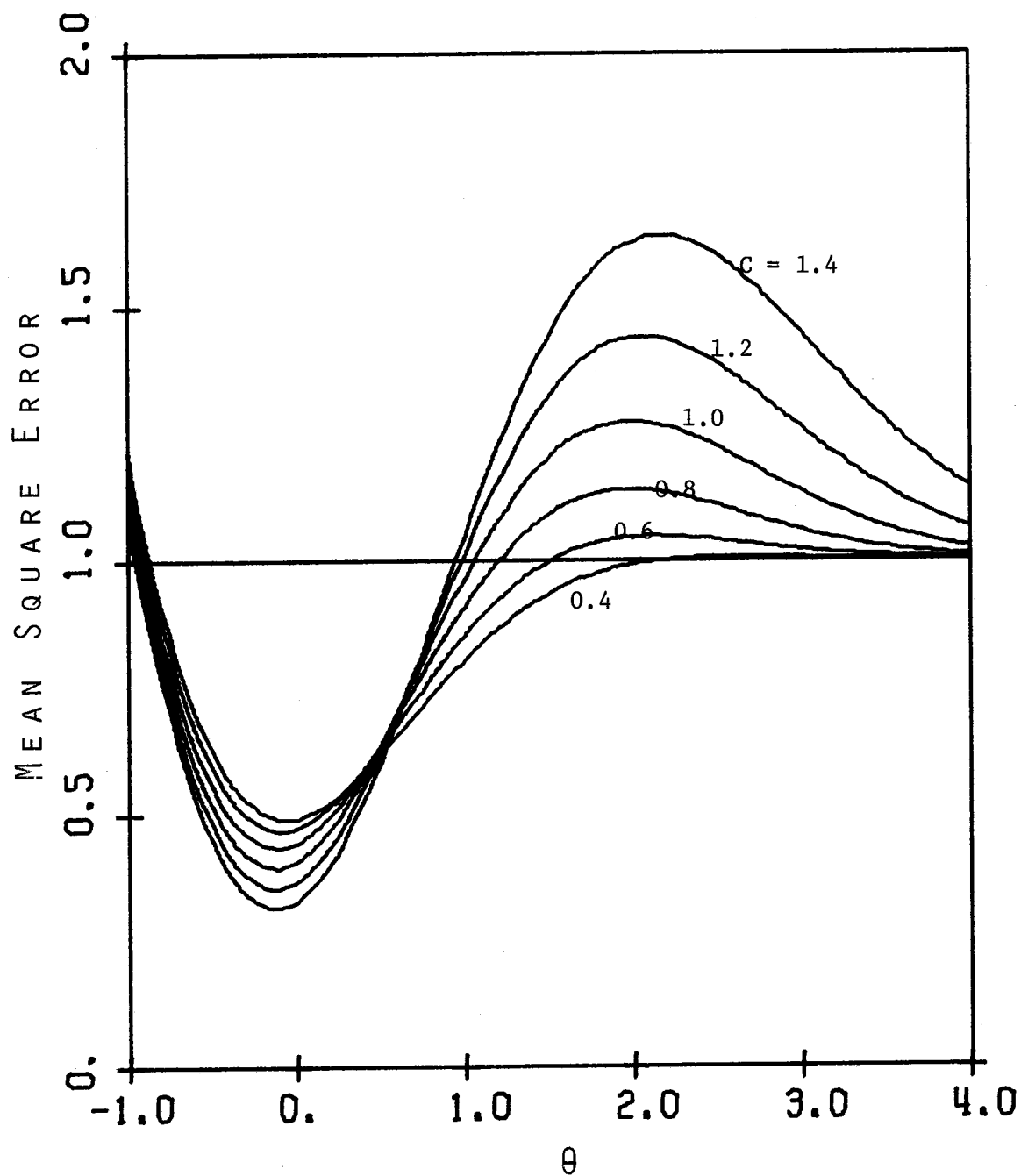


Figure 4.15 Mean Square Error of the one-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 2$ ,  $\sigma^2$  estimated.

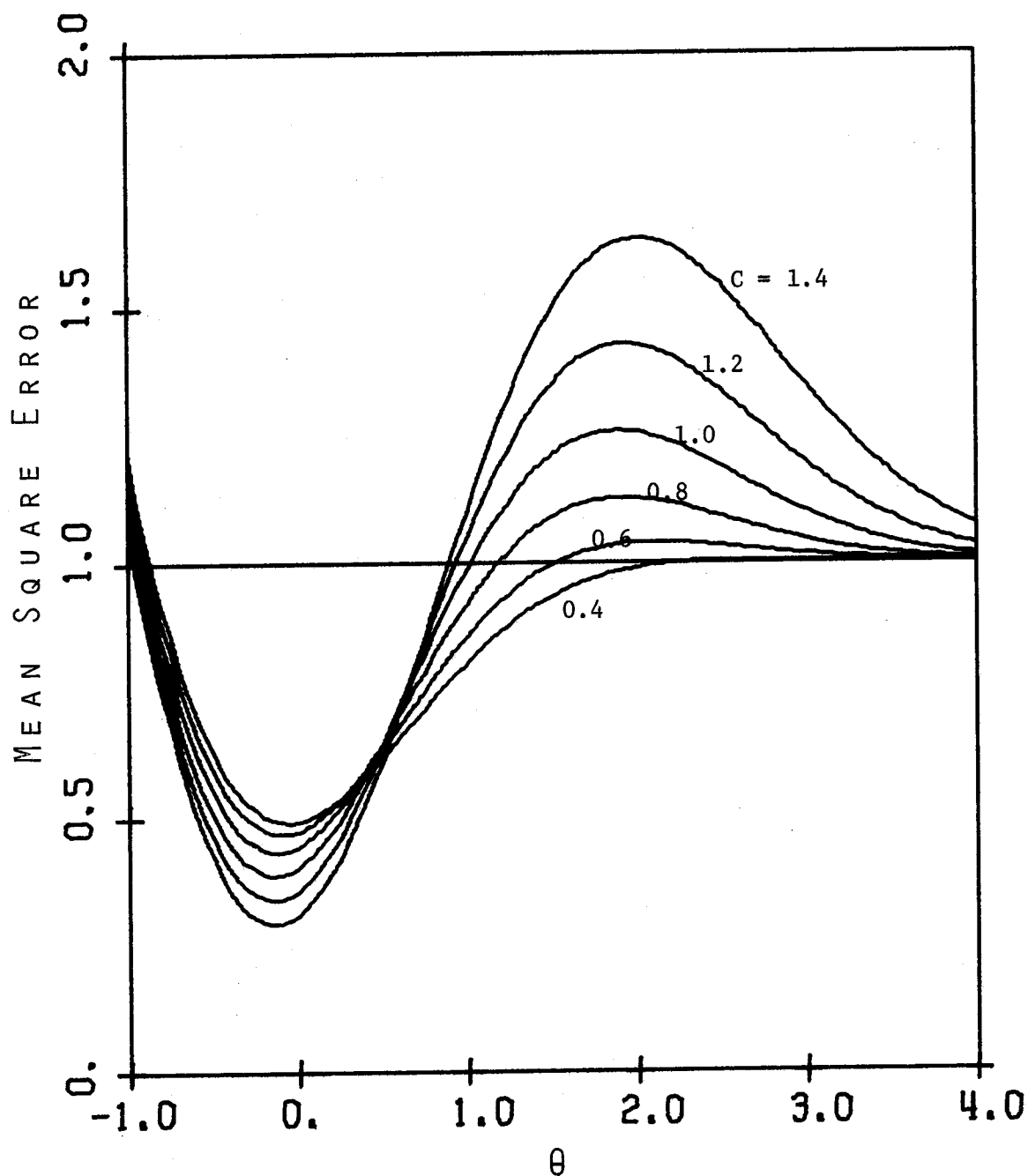


Figure 4.16 Mean Square Error of the one-tail estimator  $\hat{\theta} = \hat{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 5$ ,  $\sigma^2$  estimated.

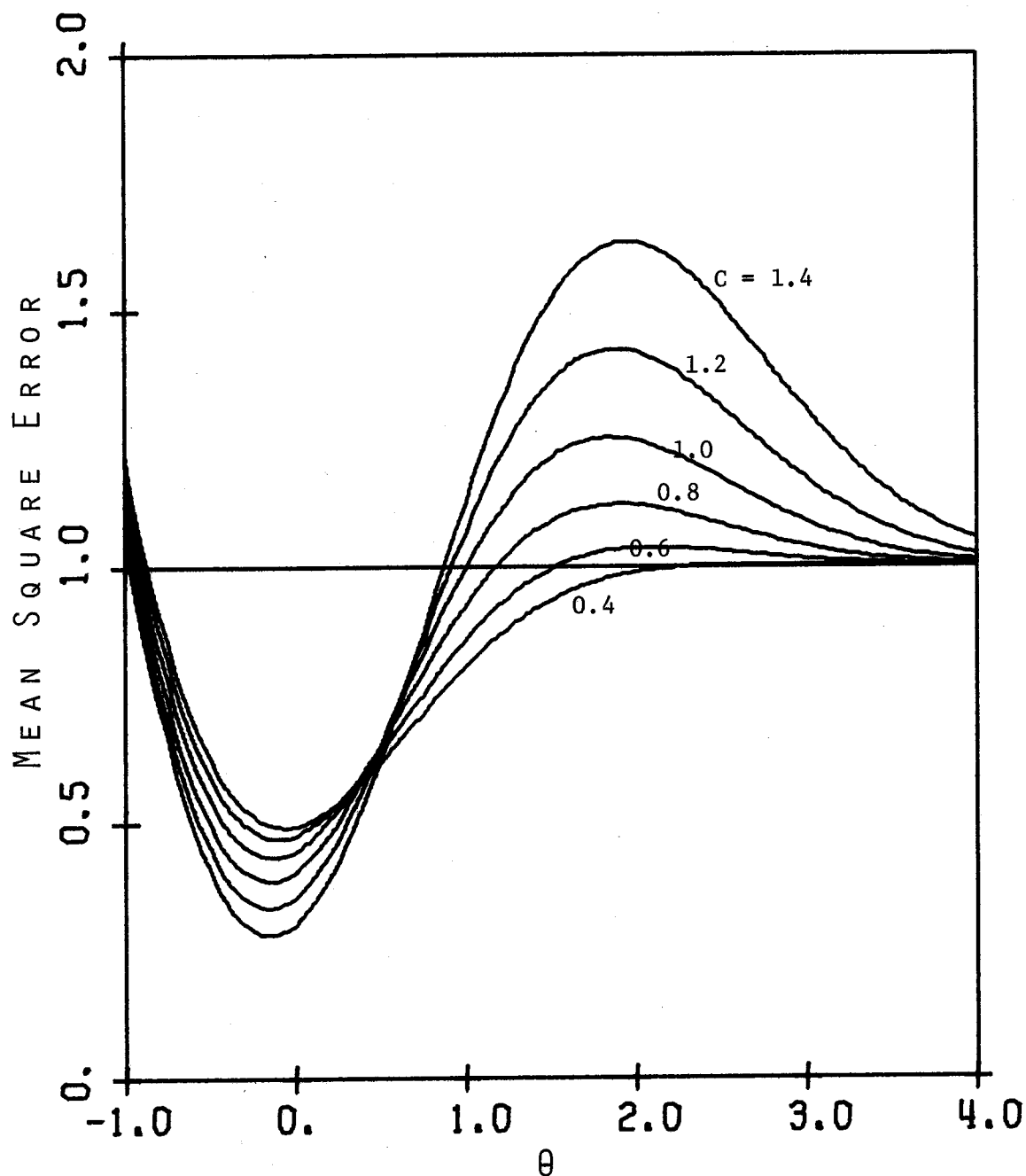


Figure 4.17 Mean Square Error of the one-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 10$ ,  $\sigma^2$  estimated.

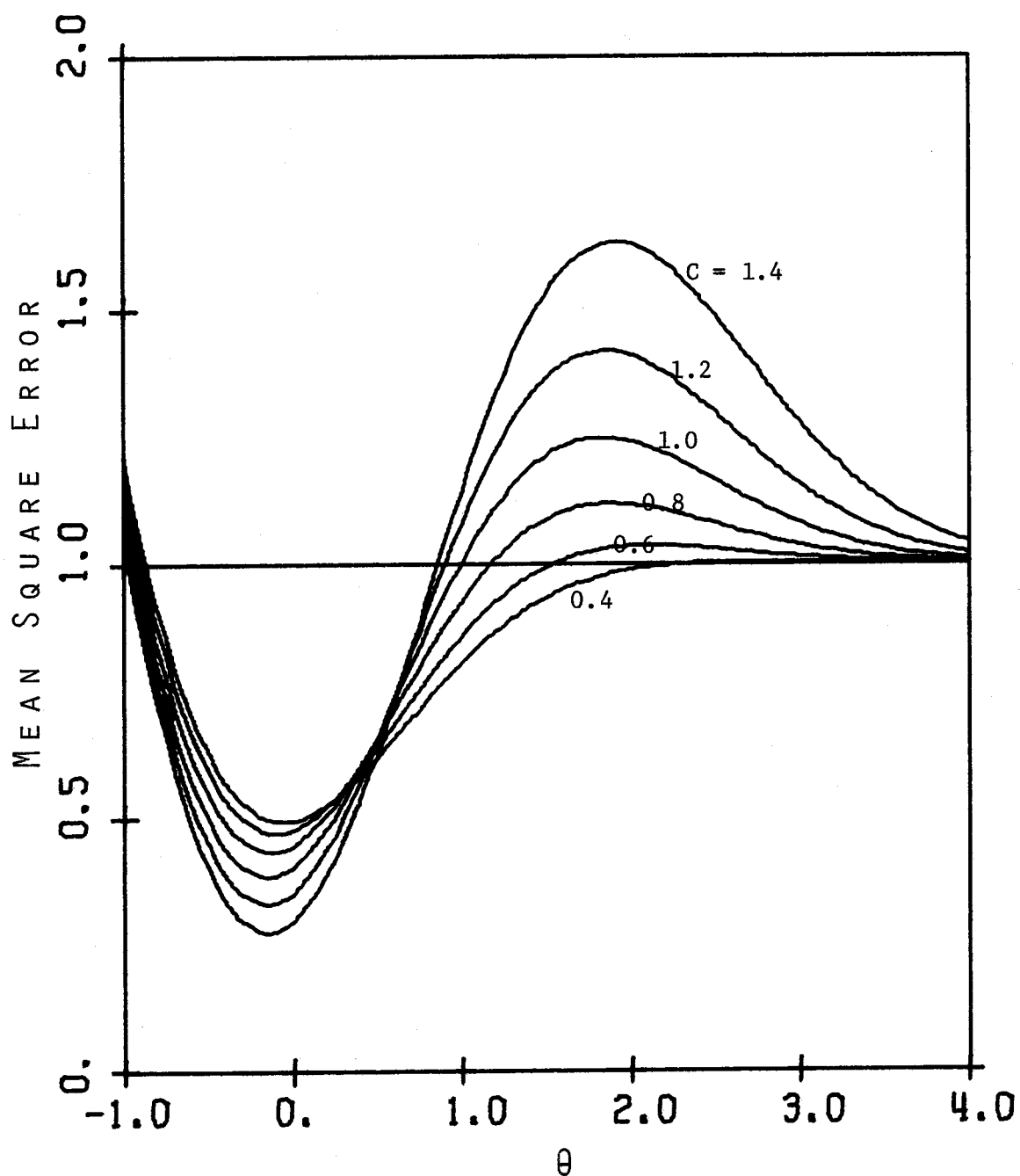


Figure 4.18 Mean Square Error of the one-tail estimator  $\hat{\theta} = \hat{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 25$ ,  $\sigma^2$  estimated.



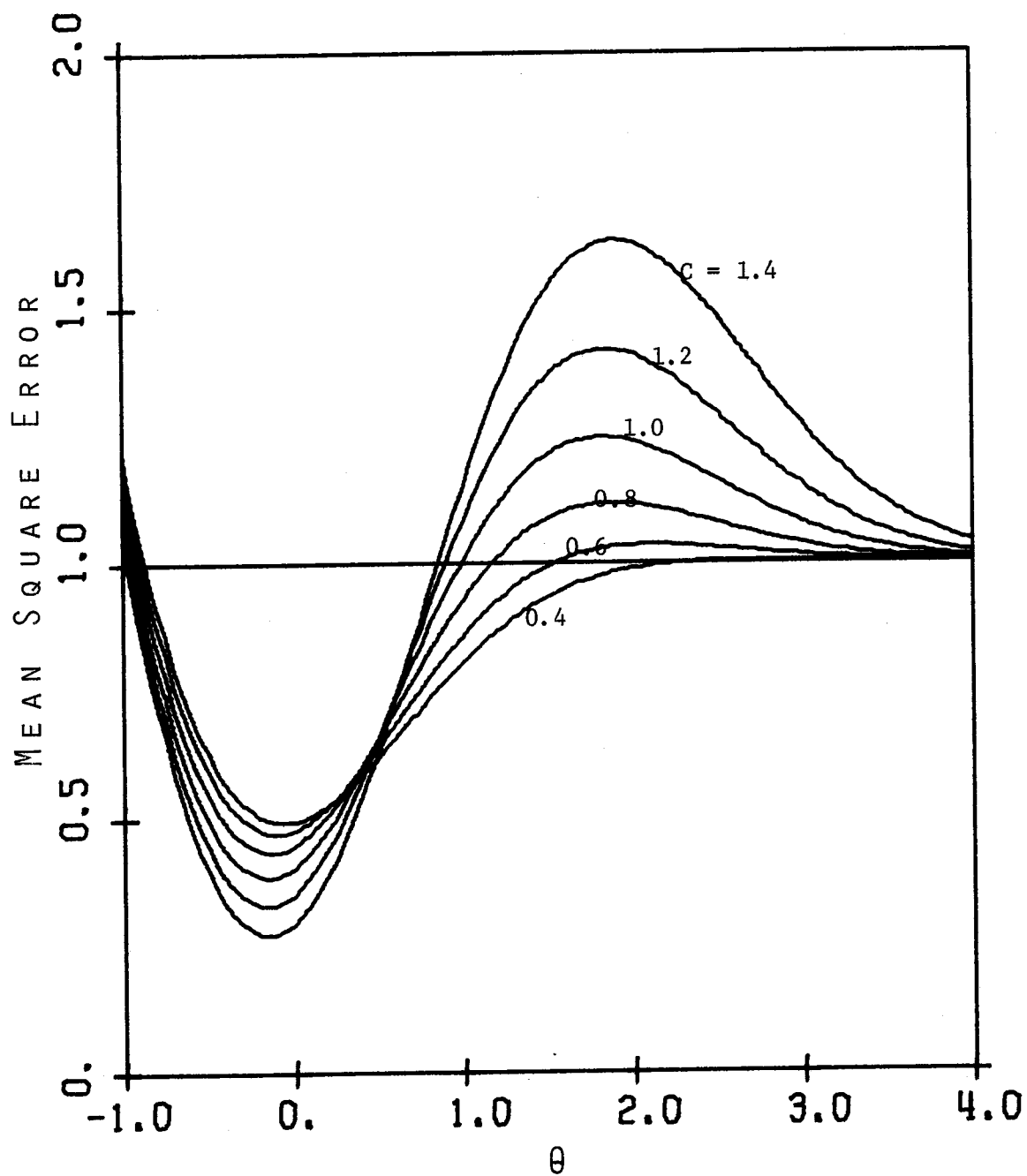


Figure 4.19 Mean Square Error of the one-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 50$ ,  $\sigma^2$  estimated.

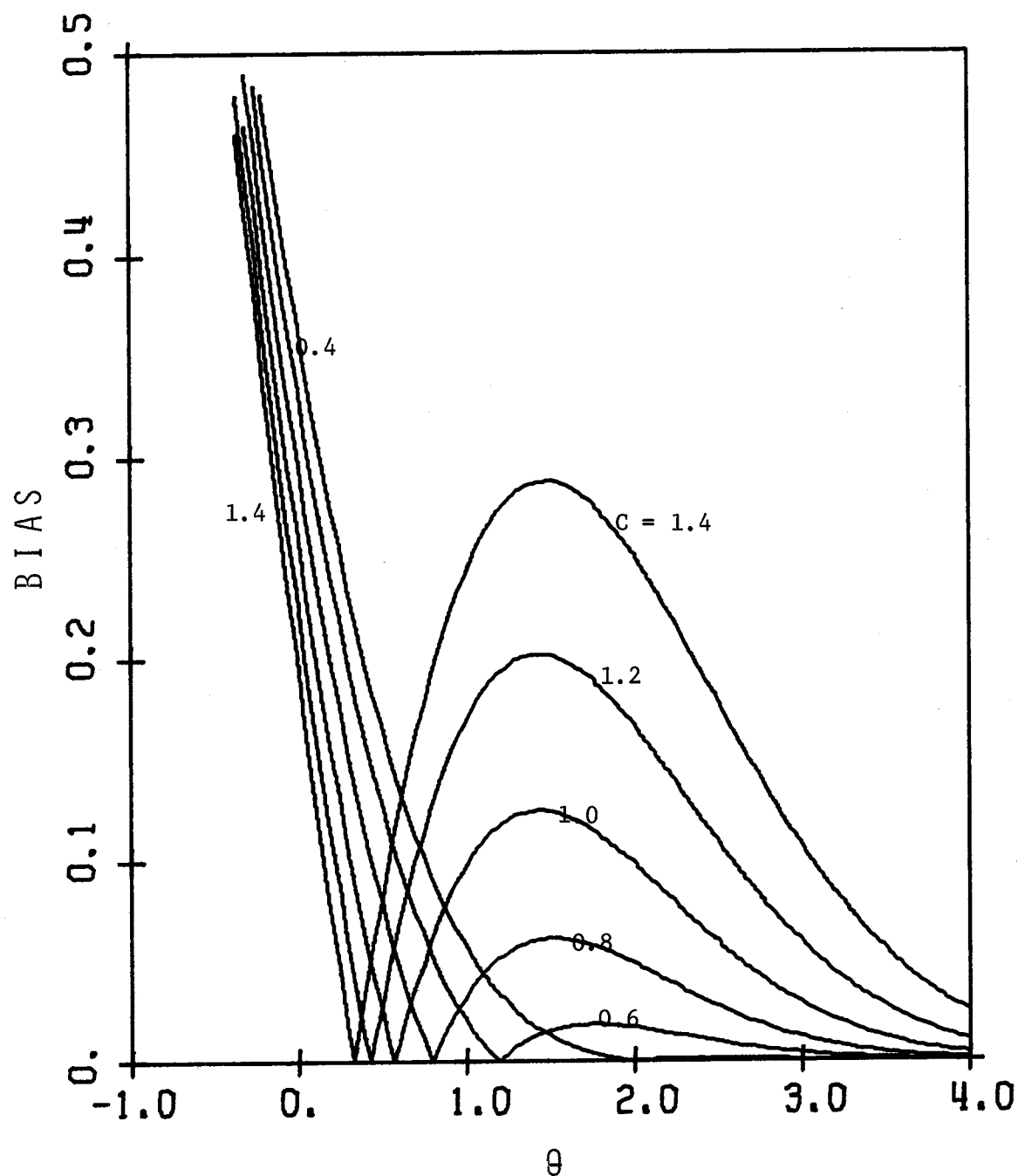


Figure 4.20. Absolute value of the bias of the one-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 2$ ,  $\sigma^2$  estimated.

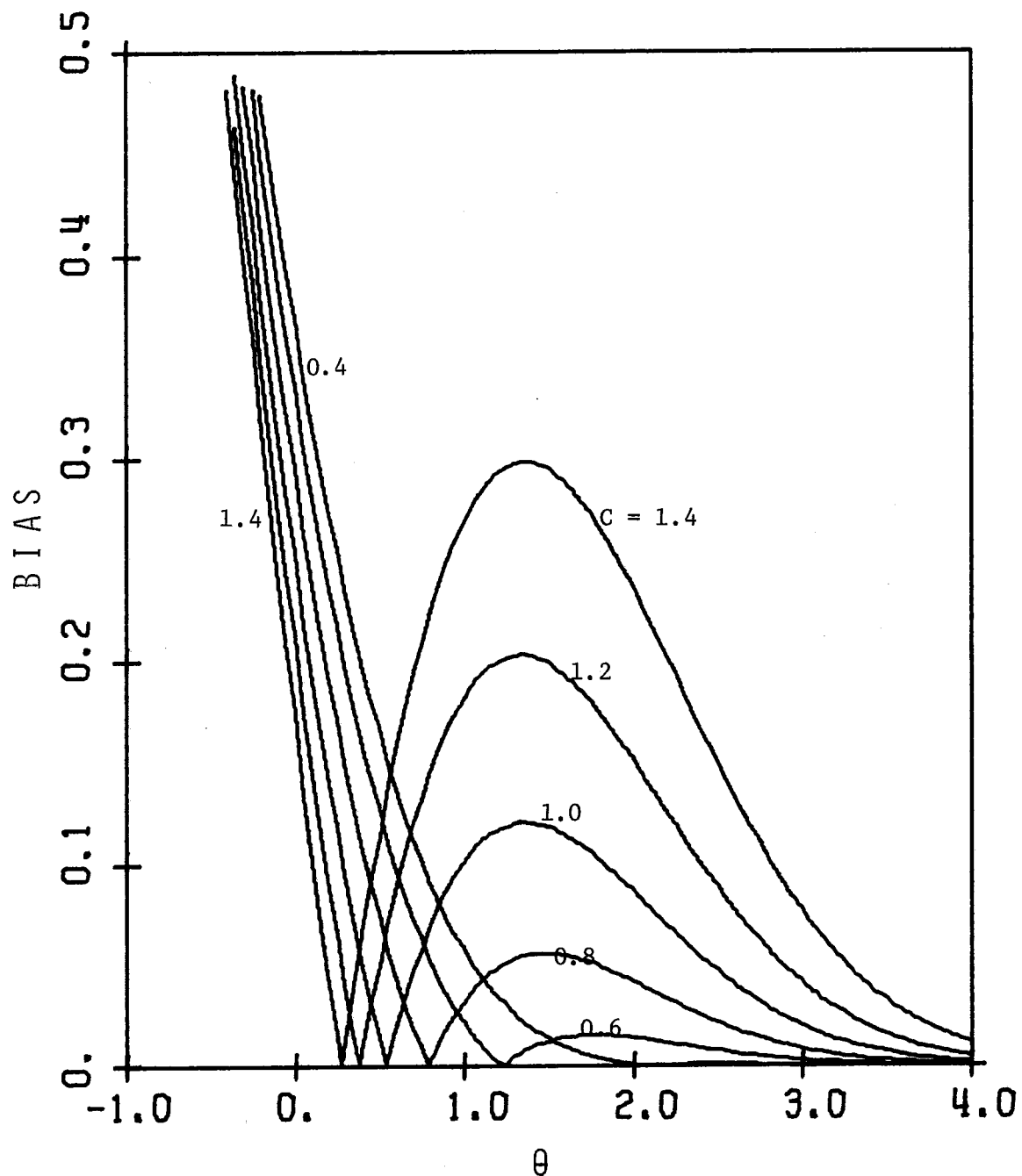


Figure 4.21 Absolute value of the bias of the one-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 5$ ,  $\sigma^2$  estimated.

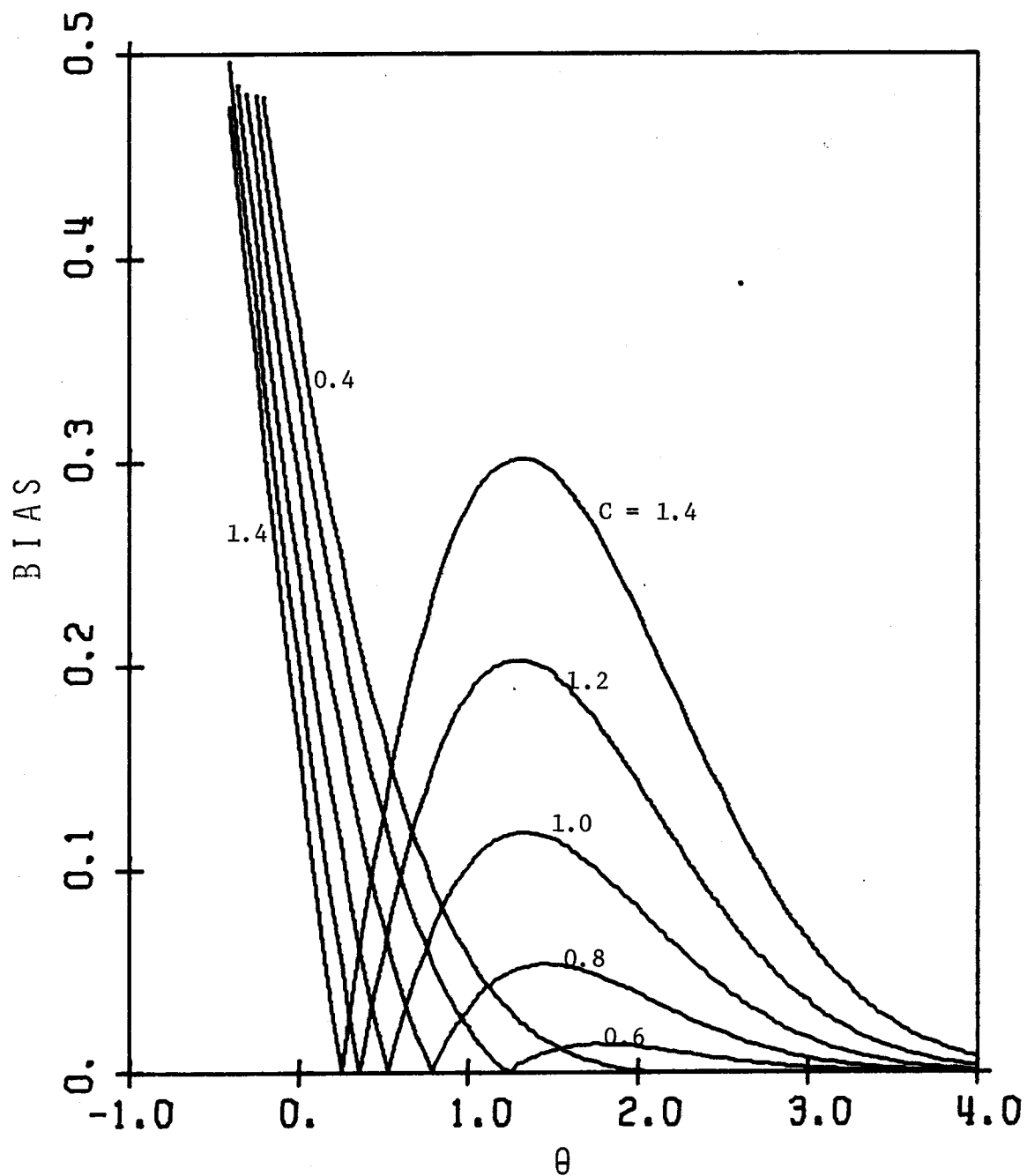


Figure 4.22 Absolute value of the bias of the one-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 10$ ,  $\sigma^2$  estimated.

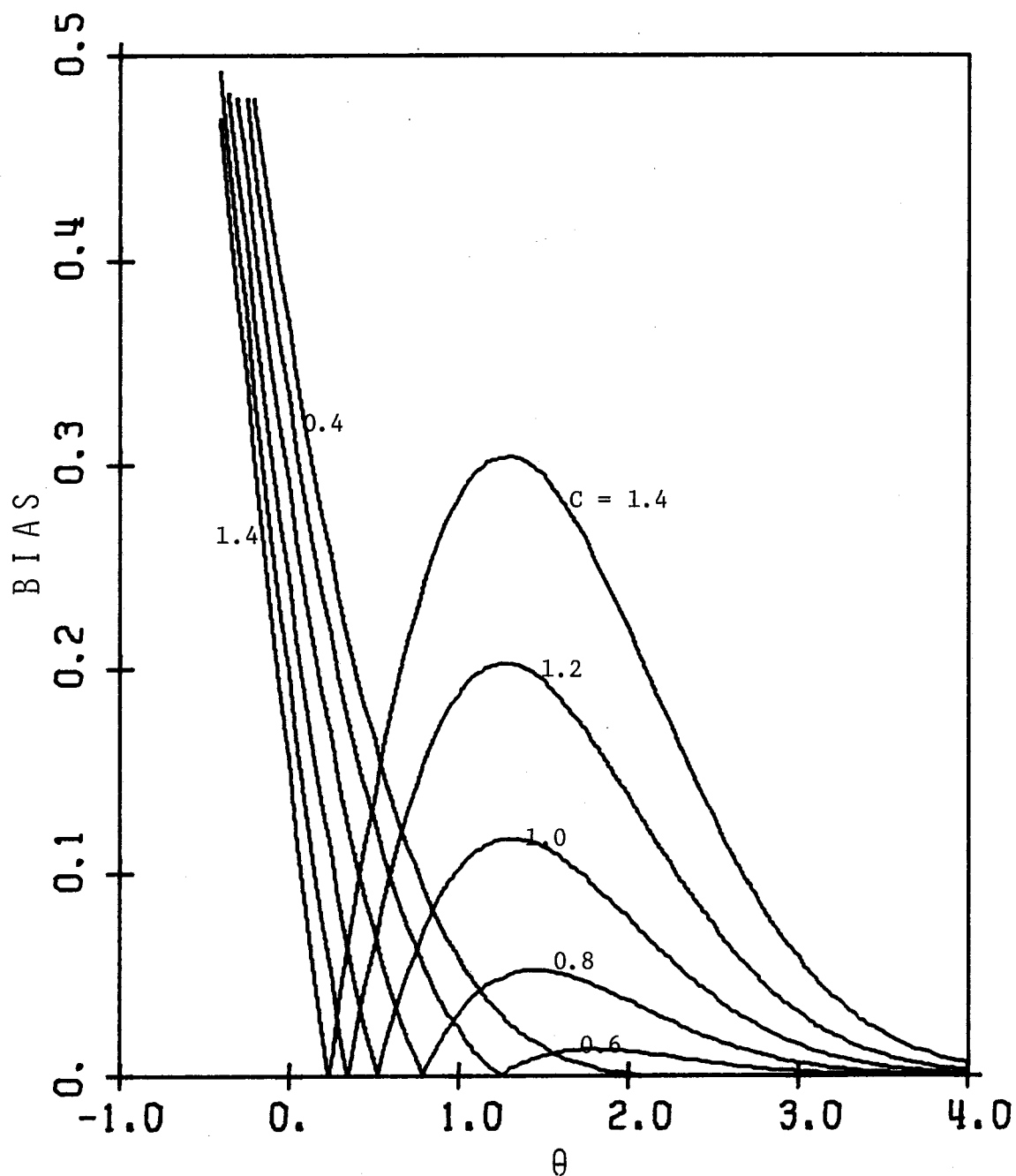


Figure 4.23 Absolute value of the bias of the one-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 25$ ,  $\sigma^2$  estimated.

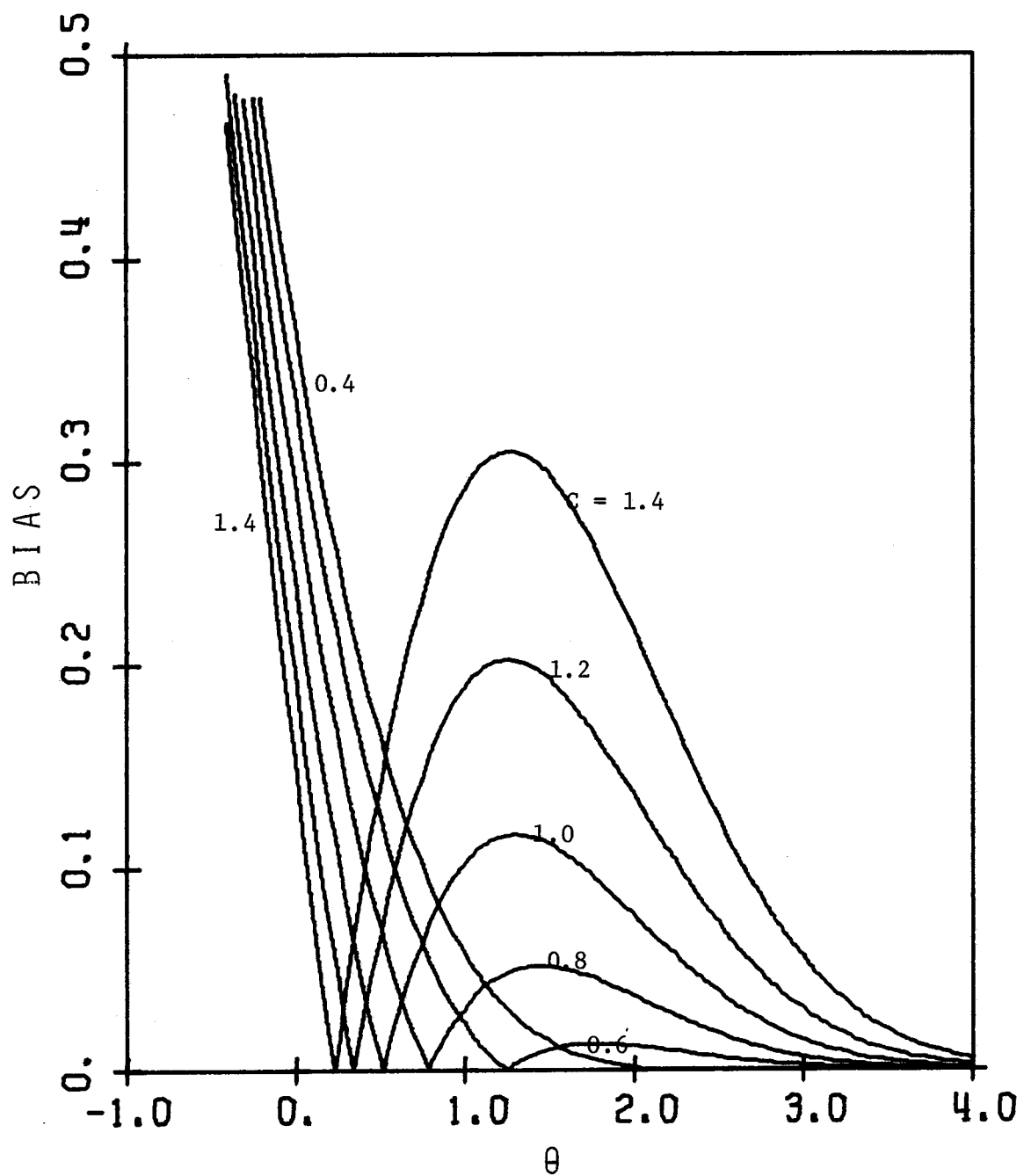


Figure 4.24 Absolute value of the bias of the one-tail estimator  $\check{\theta} = \check{\beta}/\text{s.d.}(\hat{\beta})$  as a function of  $\theta$ ,  $C$  and  $\nu = 50$ ,  $\sigma^2$  estimated.

Consider the case in which the variance is known. Let the constants  $C_1, C_2 (C_1 < C_2)$  be given, with  $C_1 = -\infty$  and/or  $C_2 = +\infty$  allowed. Define  $A = [C_1, C_2]$   $B = A^c$ . Let

$$\check{\theta}_1 = \begin{cases} 0 & \text{if } \hat{\theta} \in A \\ \hat{\theta} & \text{if } \hat{\theta} \in B \end{cases} = \hat{\theta} 1_B(\hat{\theta})$$

$$\check{\theta}_2 = \begin{cases} 0 & \text{if } \hat{\theta} \in B \\ \hat{\theta} & \text{if } \hat{\theta} \in A \end{cases} = \hat{\theta} 1_A(\hat{\theta})$$

The expected value, bias, and MSE of each of these estimators are found by simple applications of Lemma 4.2. For expected values, define (in Lemma 4.2)  $k = 0$ ,  $h(\hat{\theta}, \theta) = \hat{\theta}$ ,  $\check{\theta}_1 = g(\hat{\theta})$ , and  $\check{\theta}_2 = g^*(\hat{\theta})$ . Then, by Lemma 4.2,

$$\begin{aligned} E(\check{\theta}_1 | \theta, A) &= E(\hat{\theta}) + 0[\Phi(C_2 - \theta) - \Phi(C_1 - \theta)] - \int_{C_1 - \theta}^{C_2 - \theta} (t + \theta) f_1(t; 0, 1) dt \\ &= \theta[1. - \Phi(C_2 - \theta) + \Phi(C_1 - \theta)] - \text{PFM}(C_1 - \theta, C_2 - \theta), \end{aligned}$$

and

$$\begin{aligned} E(\check{\theta}_2 | \theta, B) &= 0 + \theta - E(\check{\theta}_1 | \theta, A) \\ &= \Phi(C_2 - \theta) - \Phi(C_1 - \theta) + \text{PFM}(C_1 - \theta, C_2 - \theta). \end{aligned}$$

Thus,

$$\begin{aligned} \text{BIAS}(\theta, A) &= E(\check{\theta}_1 - \theta | \theta, A) \\ &= \Phi(C_1 - \theta) - \Phi(C_2 - \theta) - \text{PFM}(C_1 - \theta, C_2 - \theta), \end{aligned}$$

and

$$\begin{aligned} \text{BIAS}(\theta, B) &= E(\check{\theta}_2 - \theta | \theta, A) \\ &= -\theta + \Phi(C_2 - \theta) - \Phi(C_1 - \theta) + \text{PFM}(C_1 - \theta, C_2 - \theta). \end{aligned}$$

For the MSE functions, define  $k = \theta^2$ ,  $h(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2$ ,  $g(\hat{\theta}) = (\check{\theta}_1 - \theta)^2$ ,  $g^*(\hat{\theta}) = (\check{\theta}_2 - \theta)^2$ , and note that  $h(t+\theta, \theta) = [(t+\theta) - \theta]^2 = t^2$ . By Lemma 4.2,

$$\begin{aligned} \text{MSE}(\theta, A) &= E[(\check{\theta}_1 - \theta)^2 | \theta, A] = E[g(\hat{\theta})] \\ &= 1 + \theta^2[\Phi(C_2 - \theta) - \Phi(C_1 - \theta)] - \int_{C_1 - \theta}^{C_2 - \theta} t^2 f_1(t; 0, 1) dt \\ &= 1 + \theta^2[\Phi(C_2 - \theta) - \Phi(C_1 - \theta)] - \text{PSM}(C_1 - \theta, C_2 - \theta); \end{aligned}$$

$$\begin{aligned} \text{MSE}(\theta, B) &= E[(\check{\theta}_2 - \theta)^2 | \theta, A] = E[g^*(\hat{\theta})] = \theta^2 + 1 - \text{MSE}(\theta, A) \\ &= \theta^2[1 - \Phi(C_2 - \theta) + \Phi(C_1 - \theta)] + \text{PSM}(C_1 - \theta, C_2 - \theta). \end{aligned}$$

Consider the case in which  $\text{Var}(\hat{\beta}_j)$  is unknown. Let constants  $C_1, C_2, C_3$ , and  $C_4$  be given such that

$$-\infty \leq C_1 < C_2 \leq 0 \leq C_3 < C_4 \leq +\infty.$$

Define the three subsets of the sample space:

$$A_i = \{(\hat{\theta}, u): u \geq 0, C_i \sqrt{u} < \hat{\theta} \leq C_{i+1} \sqrt{u}\}, i = 1, 2, 3,$$

and the complements (relative to the sample space),

$$B_i = \{(\hat{\theta}, u): u \geq 0, (\hat{\theta}, u) \notin A_i\}, i = 1, 2, 3,$$



where  $u = \hat{\sigma}^2 / \sigma^2$  as in previous sections. Define the estimators:

$$\check{\theta}_i = \begin{cases} 0 & \text{if } (\hat{\theta}, u) \in A_i \\ \hat{\theta} & \text{if } (\hat{\theta}, u) \in B_i \end{cases} = \hat{\theta} 1_{B_i}(\hat{\theta}, u), \quad i = 1, 2, 3;$$

$$\check{\theta}_{ci} = \begin{cases} 0 & \text{if } (\hat{\theta}, u) \in B_i \\ \hat{\theta} & \text{if } (\hat{\theta}, u) \in A_i \end{cases} = \hat{\theta} 1_{A_i}(\hat{\theta}, u), \quad i = 1, 2, 3.$$

The expected values and BIAS functions for these estimators are found from Lemma 4.1 by setting  $k = 0$  and  $h(\hat{\theta}, \theta) = \hat{\theta}$ , so that  $g_i(\hat{\theta}, u) = \check{\theta}_i$ ,  $g_i^*(\hat{\theta}, u) = \check{\theta}_{ci}$ . Temporarily define

$$I_E(\theta; a, b; c, d) = - \int_a^b (t+\theta) f_1(t; 0, 1) P \left[ \frac{v(t+\theta)^2}{c^2} < \chi_v^2 < \frac{v(t+\theta)^2}{d^2} \right] dt.$$

Then by Lemma 4.1,

$$E(\check{\theta}_1 | \theta, A_1) = \theta + I_E(\theta; -\infty, -\theta; C_1, C_2);$$

$$E(\check{\theta}_{c1} | \theta, A_1) = -I_E(\theta; -\infty, -\theta; C_1, C_2);$$

$$E(\check{\theta}_2 | \theta, A_2) = \theta + I_E(\theta; -\infty, -\theta; C_2, 0) + I_E(\theta; -\theta, \infty; C_3, 0);$$

$$E(\check{\theta}_{c2} | \theta, A_2) = -I(\theta; -\infty, -\theta; C_2, 0) - I(\theta; -\theta, \infty; C_3, 0);$$

$$E(\check{\theta}_3 | \theta, A_3) = \theta + I(\theta; -\theta, \infty; C_4, C_3);$$

$$E(\check{\theta}_{c3} | \theta, A_3) = -I(\theta; -\theta, \infty; C_4, C_3).$$

The BIAS functions for the various estimators are easily written from the equations above.

The MSE functions can be found by applying Lemma 4.1 with  $k = \theta^2$ ,  $h(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2$ , so that  $h(t+\theta, \theta) = t^2$ ,  $g_1(\hat{\theta}, u) = (\check{\theta}_1 - \theta)^2$  and  $g_1^*(\hat{\theta}, u) = (\check{\theta}_{c1} - \theta)^2$ . Temporarily define the function

$$I_M(\theta; a, b; c, d) = \int_a^b (\theta^2 - t^2) f_1(t; 0, 1) P \left[ \frac{v(t+\theta)^2}{c^2} < \chi^2 < \frac{v(t+\theta)^2}{d^2} \right] dt.$$

Then, by Lemma 4.1,

$$\text{MSE}(\check{\theta}_1 | \theta, A_1) = 1 + I_M(\theta; -\infty, -\theta; C_1, C_2);$$

$$\text{MSE}(\check{\theta}_{c1} | \theta, A_1) = \theta^2 - I_M(\theta; -\infty, -\theta; C_1, C_2);$$

$$\text{MSE}(\check{\theta}_2 | \theta, A_2) = 1 + I_M(\theta; -\infty, -\theta; C_2, 0) + I_M(\theta; -\theta, \infty; C_3, 0);$$

$$\text{MSE}(\check{\theta}_{c2} | \theta, A_2) = \theta^2 - I_M(\theta, -\infty, -\theta; C_2, 0) - I_M(\theta; -\theta, \infty; C_3, 0);$$

$$\text{MSE}(\check{\theta}_3 | \theta, A_3) = 1 + I_M(\theta; -\theta, \infty; C_4, C_3);$$

$$\text{MSE}(\check{\theta}_{c3} | \theta, A_3) = \theta^2 - I_M(\theta, -\theta; \infty; C_4, C_3).$$

It is easily seen that these results include the results of the previous sections as special cases. The symmetric two-tail estimator is obtained by setting  $C_2 = -C_3$  and using  $\check{\theta}_2$ . The one-tail estimator is obtained by setting  $C_2 = -\infty$  and using  $\check{\theta}_2$ .

The functions above are evaluated numerically in the same manner as the corresponding functions for estimators discussed in previous sections. The computer subroutines used for the previous estimators were altered to compute the properties of these general estimators. Previous comments on accuracy apply here, also.

It is important to note that the MSE and BIAS functions in this case are not symmetric. The figures on the following pages illustrate MSE and BIAS curves for various  $(C_1, C_2)$  values for the known variance case.

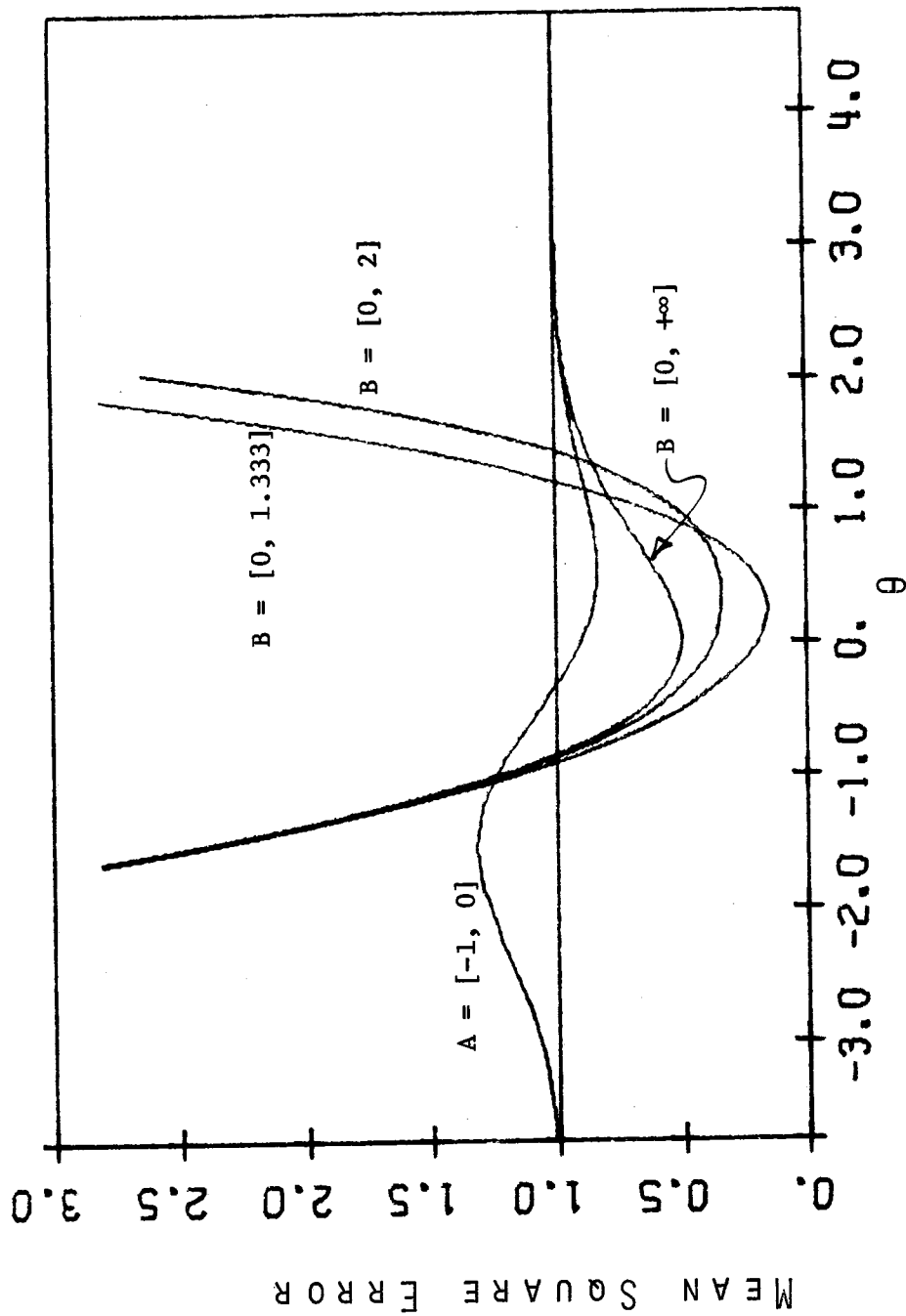


Figure 4.25 Mean Square Error of the generalized estimator  $\hat{\theta} = 0 \cdot I_A(\theta) + \theta \cdot I_B(\theta)$ , variance assumed known.

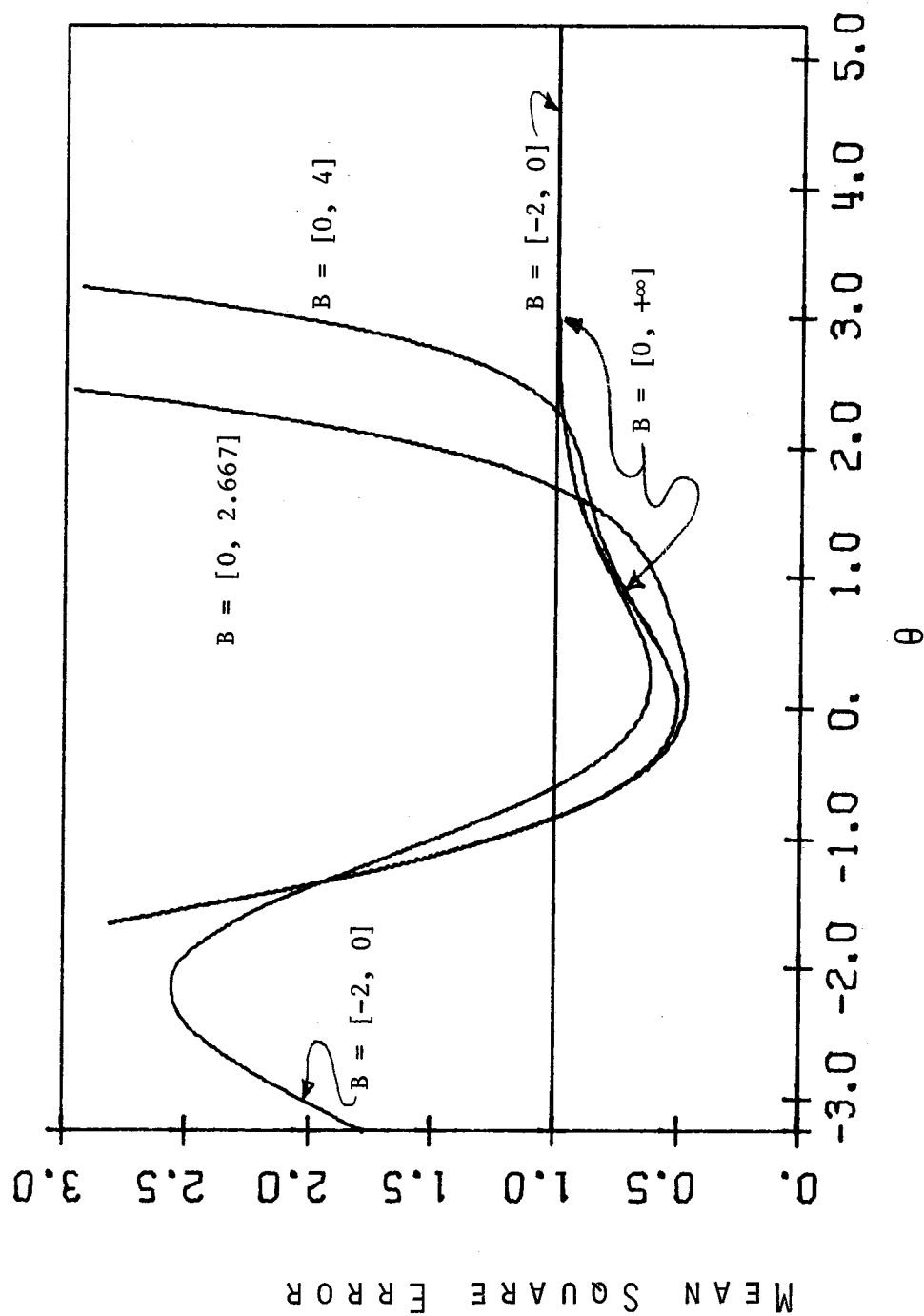


Figure 4.26 Mean Square Error of the generalized estimator  
 $\check{\theta} = 0 \quad I_A(\check{\theta}) + \theta \quad I_B(\check{\theta})$ , variance assumed known.

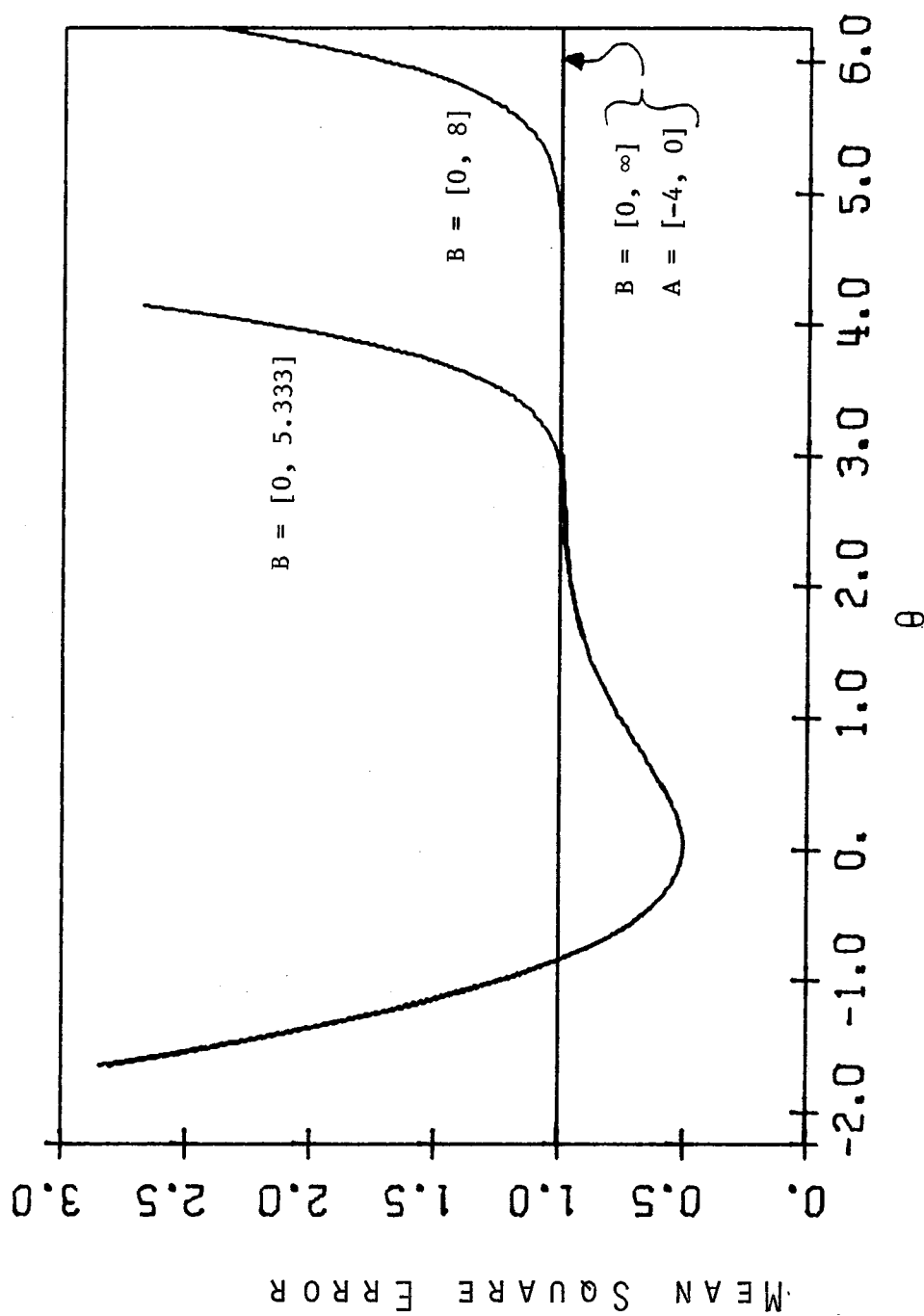


Figure 4.27 Mean Square Error of the generalized estimator  $\hat{\theta} = 0 \mathbf{1}_A(\hat{\theta}) + \hat{\theta} \mathbf{1}_B(\hat{\theta})$ , variance assumed known.

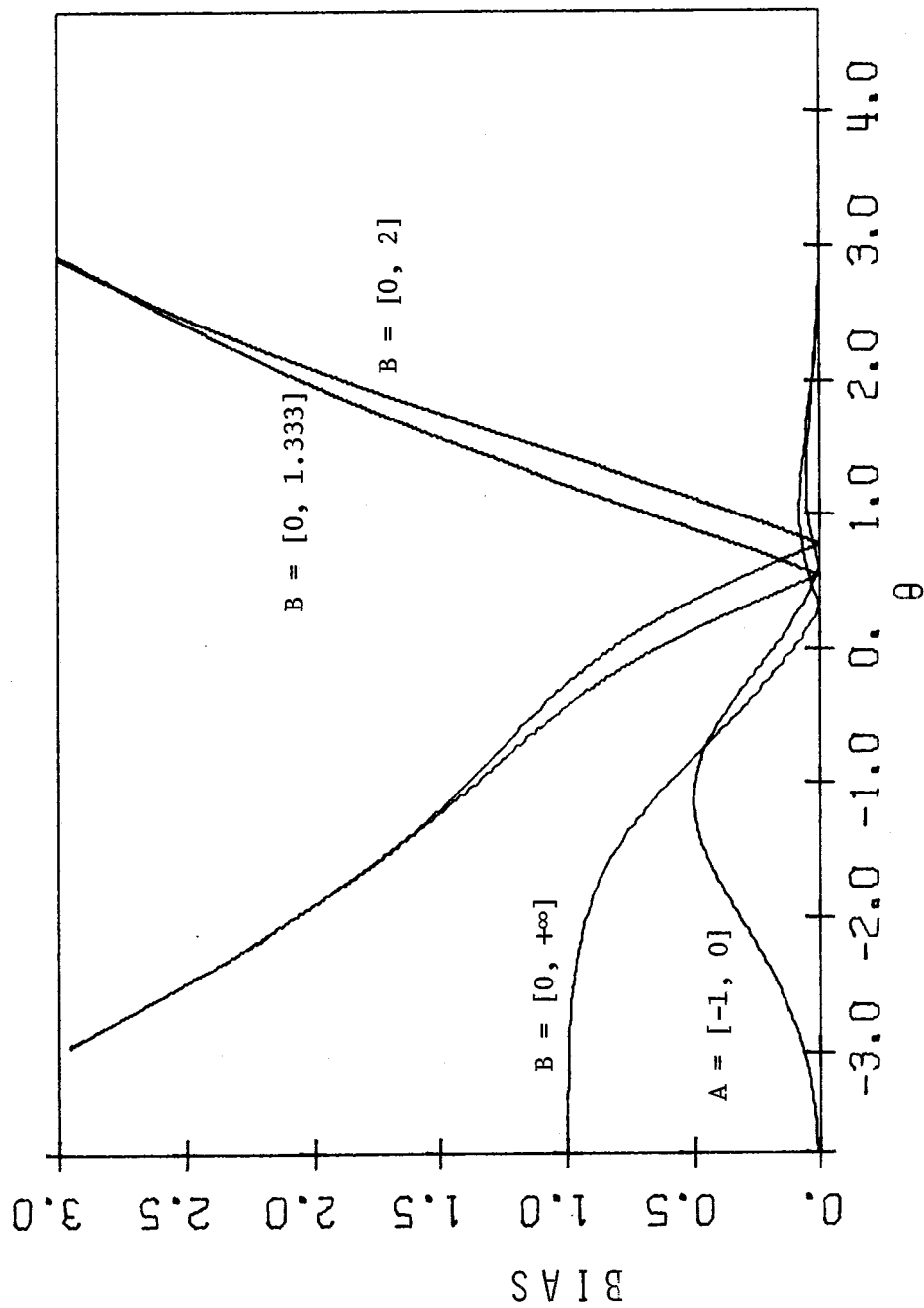


Figure 4.28 Absolute value of the bias of the generalized estimator  $\hat{\theta} = 0 \cdot 1_A(\theta) + \theta \cdot 1_B(\theta)$ , variance assumed known.

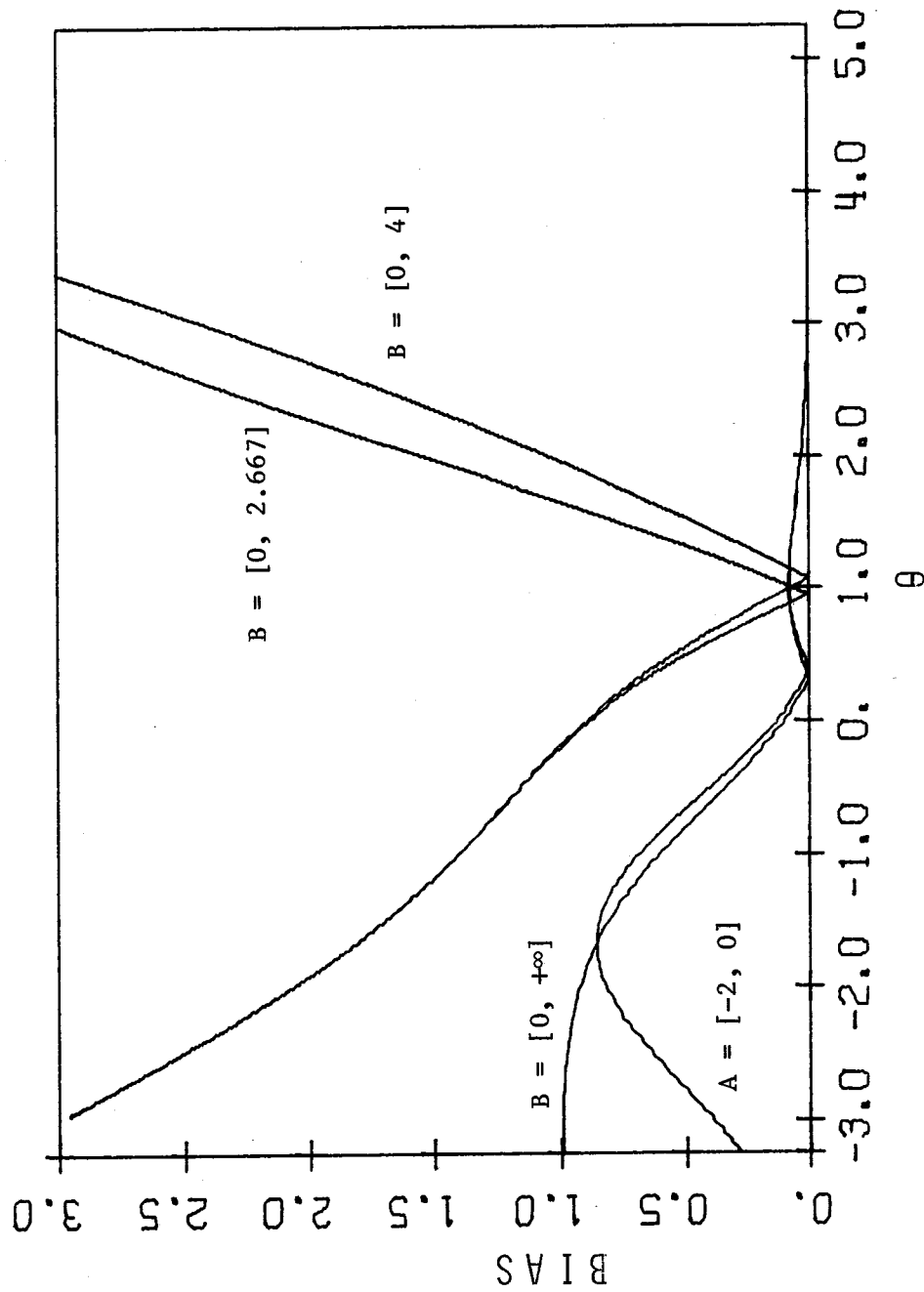


Figure 4.29 Absolute value of the bias of the generalized estimator  $\hat{\theta} = 0 \quad l_A(\hat{\theta}) + \hat{\theta} \quad l_B(\hat{\theta})$ , variance assumed known.



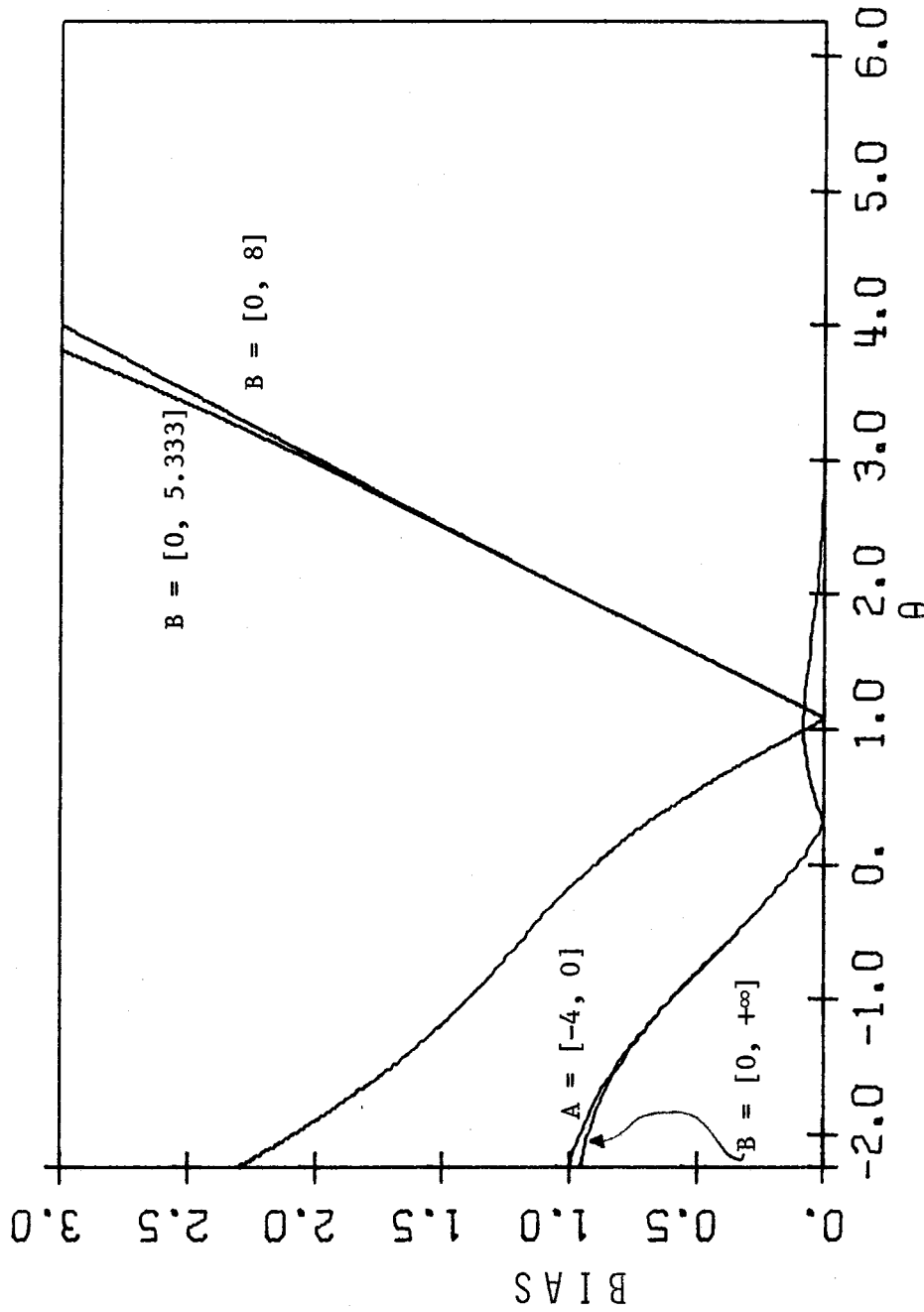


Figure 4.30 Absolute value of the bias of the generalized estimator  $\hat{\theta} = 0 \quad l_A(\theta) + \theta \quad l_B(\theta)$ , variance assumed known.

## 5. SOME TECHNIQUES FOR THE SELECTION OF CUTOFF POINTS

In order to apply the procedures discussed in the previous sections one must select the cutoff points ( $C$ , or  $C_1$  and  $C_2$ ). Presumably, if there is prior information about a particular parameter, one would like to use this information in the selection of cutoff points in an attempt to decrease the mean square errors of an estimator,  $\check{\beta}_j$ . For example, if one were absolutely certain that  $|\beta_j| < \sqrt{\text{var}(\hat{\beta}_j)}$ , one would always delete the  $j$ -th term from the estimation equation ( $\check{\beta}_j = 0$ ); i.e., in effect one would use a symmetric two-tail procedure with  $C = +\infty$ . On the other hand, if one were absolutely certain that  $|\beta_j| > \sqrt{\text{var}(\hat{\beta}_j)}$  one would use the symmetric two-tail procedure with  $C = 0$ ; i.e., always include the  $j$ -th term in the estimation equation, setting  $\check{\beta}_j = \hat{\beta}_j$ . These two extremes help explain the relationship between  $\theta$ -values and optimum  $C$ -values. In general, if prior information indicates that  $\beta_j^2 < \text{Var}(\hat{\beta}_j)$ , one would tend to use large  $C$ -values with either the one-tail or symmetric two-tail procedures. If prior information indicates that  $\beta_j^2 > \text{var}(\hat{\beta}_j)$ , one would tend to use small  $C$ -values (use  $\check{\beta}_j = \hat{\beta}_j$  a larger proportion of the time, or increase the probability that  $\check{\beta}_j = \hat{\beta}_j$ ).

In this section two techniques are discussed and illustrated for the selection of cutoff points when prior information is available in the form of a "prior distribution" for a particular parameter,  $\theta$ .

It is assumed throughout this discussion that the variance is known. Although the techniques can also be applied to the unknown variance situation, the MSE curves for the known and unknown variance situations are so similar that it is expected that the selected cutoff points for the two situations would not be very different.

The following notation will be useful. Let:

$g(\theta)$  denote the prior density of  $\theta$ ;

$f(\hat{\theta}|\theta)$  denote the conditional density of  $\hat{\theta}$  given  $\theta$ ,

i.e., the  $N(\theta, 1)$  density;

$f(\hat{\theta}, \theta) = g(\theta)f(\hat{\theta}|\theta)$  denote the joint density of  $\hat{\theta}$ ,  $\theta$ ;

$f(\hat{\theta}) = \int f(\hat{\theta}, \theta)d\theta$  denote the marginal density of  $\theta$ ;

$h(\theta|\hat{\theta}) = f(\hat{\theta}, \theta)/f(\hat{\theta})$  denote the conditional density of  $\theta$  given  $\hat{\theta}$ ,

i.e., the posterior density of  $\theta$ .

It is assumed that all of the functions above exist and are valid probability density functions. In Bayesian terminology, the "decision function" is

$$d(\hat{\theta}) = \check{\theta} = 0 \cdot 1_A(\hat{\theta}) + \hat{\theta} \cdot 1_B(\hat{\theta}),$$

where the set A is the set on which  $\check{\theta} = 0$ , and B is the set on which  $\check{\theta} = \hat{\theta}$ . The loss function is:

$$L(\theta, d(\hat{\theta})) = (\check{\theta} - \theta)^2 = \theta^2 \cdot 1_A(\hat{\theta}) + (\hat{\theta} - \theta)^2 \cdot 1_B(\hat{\theta});$$

the risk (expected loss) is:

$$R(\theta, d) = \int (\hat{\theta} - \theta)^2 f(\hat{\theta}|\theta) d\hat{\theta} = \text{MSE}(\theta, C_1, C_2).$$

### 5.1. The Bayes Decision Procedure

One strategy for the selection of the cutoff points (equivalent to the selection of the decision function,  $d(\hat{\theta})$ ) is to use the Bayes decision procedure (Lindgren, 1962; Michaels, 1969), which consists of choosing  $d$  (i.e.,  $C_1$  and  $C_2$ ) to minimize the expected risk with respect to the prior of  $\theta$ , i.e., choose  $C_1$  and  $C_2$  to minimize

$$E_g[\text{MSE}(\theta, C_1, C_2)] = \int g(\theta) \text{MSE}(\theta, C_1, C_2) d\theta.$$

Lindgren (1962, p. 285) shows this strategy is equivalent to minimization of the expected posterior loss, i.e., after  $\hat{\theta}$  is observed, compute  $h(\theta|\hat{\theta})$  and choose the sets  $A$  and  $B$  to minimize:

$$\begin{aligned} E_h[L(\theta, d(\hat{\theta}))] &= \int [\theta^2 1_A(\hat{\theta}) + (\hat{\theta} - \theta)^2 1_B(\hat{\theta})] h(\theta|\hat{\theta}) d\theta \\ &= \begin{cases} E_h(\theta^2) & \text{if } \hat{\theta} \in A \\ E_h(\hat{\theta} - \theta)^2 = \hat{\theta}^2 - 2\hat{\theta} E_h(\theta) + E_h(\theta^2) & \text{if } \hat{\theta} \in B. \end{cases} \end{aligned}$$

Thus, the set  $A$  should be chosen so that  $\hat{\theta} \in A$  ( $\hat{\theta} = 0$ ) if

$$E_h(\theta^2) < \hat{\theta}^2 - 2\hat{\theta} E_h(\theta) + E_h(\theta^2)$$

$$\Leftrightarrow \hat{\theta}[2E_h(\theta) - \hat{\theta}] < 0.$$

The Bayes decision procedure (for selecting A and B or  $C_1$  and  $C_2$ ) can be summarized as follows:

- (1) Compute the posterior density of  $\theta$  for the observed value of  $\hat{\theta} = \hat{\beta}_j \sqrt{\text{Var}(\hat{\beta}_j)}$ ;
- (2) Compute the mean of the posterior density,  $E_h(\theta)$ . Then compute  $q(\hat{\theta}) = \hat{\theta}[2E_h(\theta) - \hat{\theta}]$ .
- (3) If  $q(\hat{\theta}) < 0$ , set  $\check{\theta} = 0$ ; otherwise  $\check{\theta} = \hat{\theta}$ .

The procedure does not yield the sets A and B directly; one method for finding the "effective" cutoff points for the above procedure is to "search" for the cutoff points by repeating the procedure above for various  $\hat{\theta}$  values. For the normal prior distribution it is possible to find  $E_h(\theta)$  and the sets A, B explicitly in terms of the parameters of the prior and  $\hat{\theta}$  (see section 5.3).

The Bayes decision procedure is optimal in the sense that the sets A and B are chosen so as to minimize the expected value of  $\text{MSE}(\theta, C_1, C_2)$  with respect to the prior distribution. However, in some cases the procedure is very sensitive to small changes in the parameters of the prior distribution.

## 5.2. The Posterior $C=1$ Procedure

If the parameter  $\theta$  (or  $\beta_j$ ) were known exactly, the appropriate cutoff points would be -1 and 1, which motivates the following procedure:

- (1) Compute the mean,  $E_h(\theta|\hat{\theta})$ , of the posterior distribution;

(2) Apply the rule:

$$\check{\theta} = \begin{cases} 0 & \text{if } |E_h(\theta|\hat{\theta})| < 1 \\ \hat{\theta} & \text{otherwise} \end{cases}.$$

While this rule is not optimum in the sense that the Bayes decision procedure is optimum, it does have intuitive appeal. The mean of the posterior distribution in a sense summarizes the prior information and the information from the sample regarding the value of  $\theta$ ; the same procedure is applied to the posterior mean that would be applied to the parameter itself, if it were known.

Aside from any intuitive appeal, the cutoff points produced by the "Posterior C=1 Procedure" for a normal prior are not as sensitive to small changes in the prior parameters as the Bayes procedure, and may be more appealing in other ways than the cutoff points produced by the Bayes decision procedure.

The Posterior C=1 Procedure does not produce explicit algebraic formulas for cutoff points. For the normal  $N(\mu, t^2)$  prior, the cutoff points can be found explicitly in terms of  $\mu$  and  $t^2$  (see section 5.3); for other priors a search procedure may be necessary.

### 5.3. Cutoff Point Selection with a Normal Prior

Hogg and Craig (1965, p. 165) show that for the  $N(\mu, t^2)$  prior for  $\theta$ , and  $N(\theta, 1)$  density for  $\hat{\theta}$ , the posterior distribution of  $\theta$  given  $\hat{\theta}$  is the normal distribution with mean

$$\mu^* = \frac{\mu + t^2 \hat{\theta}}{1 + t^2} = E_h(\theta|\hat{\theta}) = \left(\frac{1}{1+t^2}\right) \mu + \left(\frac{t^2}{1+t^2}\right) \hat{\theta},$$

and variance

$$t^{*2} = \frac{t^2}{1+t^2}.$$

Note that the posterior mean is a weighted average of the prior mean and the observed  $\hat{\theta}$ .

From the results of section 5.1, the Bayes decision procedure depends on the sign of the function

$$q(\hat{\theta}) = \hat{\theta}[2E_h(\theta) - \hat{\theta}] = \hat{\theta}[2\mu^* - \hat{\theta}]$$

which simplifies to:

$$q(\hat{\theta}) = \hat{\theta}^2 \left( \frac{t^2-1}{t^2+1} \right) + \hat{\theta} \left( \frac{2\mu}{t^2+1} \right),$$

with roots at  $\hat{\theta} = 0$  and  $\hat{\theta} = -2\mu/(t^2-1)$  if  $t^2 \neq 1$ , or just  $\hat{\theta} = 0$  if  $t^2 = 1$ . Typical graphs of  $q(\hat{\theta})$  for various combinations of  $\mu, t^2$  values are presented in Figure 5.1. Remember that  $\check{\theta} = 0$  wherever  $q(\hat{\theta}) < 0$ , and  $\theta = \hat{\theta}$  wherever  $q(\hat{\theta}) \geq 0$ .

The application of the Posterior C=1 Procedure is as follows:

$$\check{\theta} = \begin{cases} 0 & \text{if } |\mu^*| = \left| \frac{\mu + \hat{\theta}t^2}{t^2+1} \right| < 1 \\ \hat{\theta} & \text{otherwise} \end{cases},$$

which is equivalent to:

$$\check{\theta} = \begin{cases} 0 & \text{if } c_1 \leq \hat{\theta} \leq c_2 \\ \hat{\theta} & \text{otherwise} \end{cases},$$

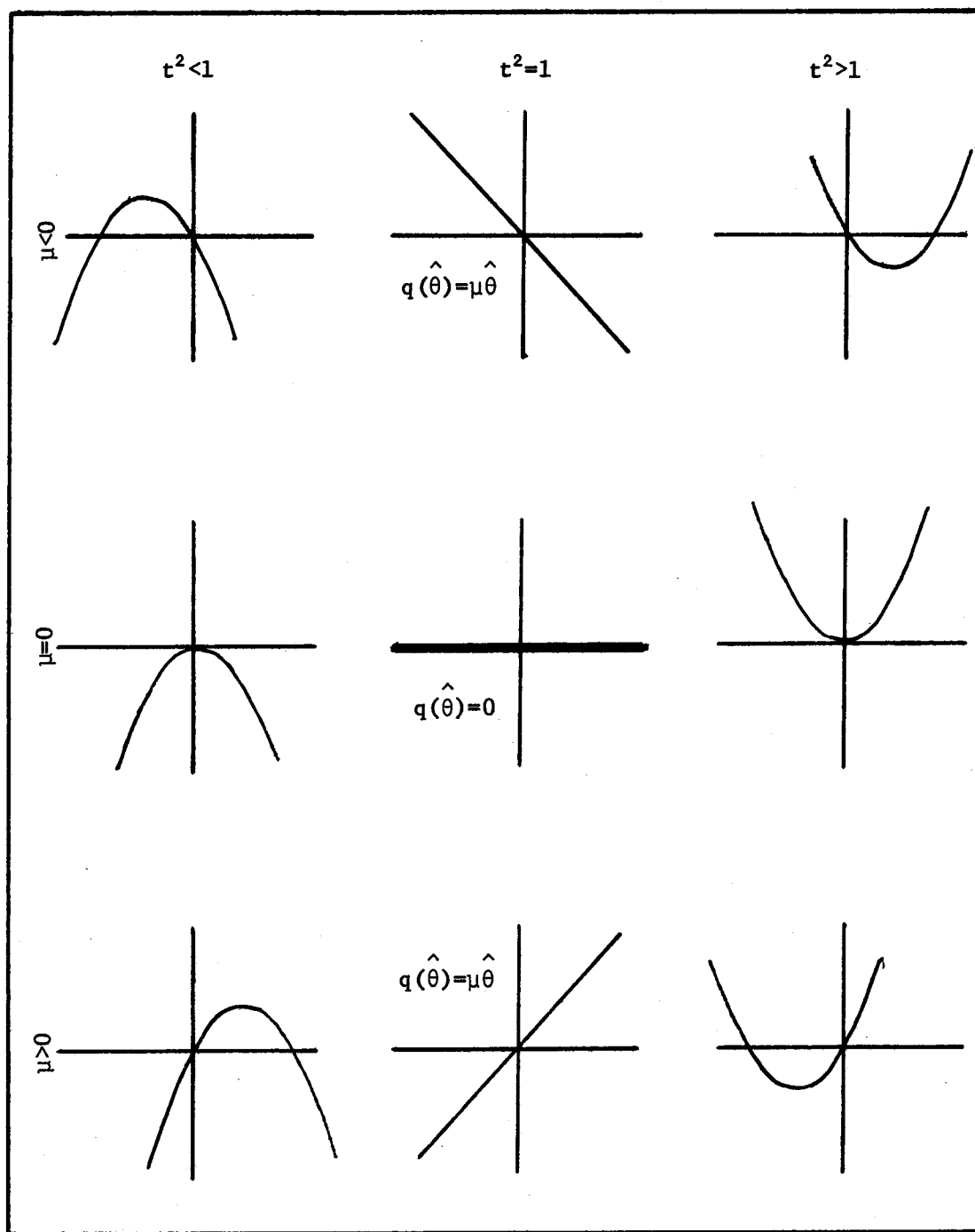


Figure 5.1 Configurations of  $q(\hat{\theta})$  for various combinations of prior distribution parameters.



where:

$$C_1 = -\frac{\mu}{t^2} - \left( \frac{1+t^2}{t^2} \right) = -\left( \frac{1+\mu}{t^2} \right) - 1$$

$$C_2 = -\frac{\mu}{t^2} + \left( \frac{1+t^2}{t^2} \right) = \left( \frac{1-\mu}{t^2} \right) + 1.$$

Various properties of the two procedures may be compared, the most important being the expected value of the MSE function with respect to the prior distribution. The Bayes procedure is, by definition, optimum for this criterion. A comparison of the expected MSE functions is given in Table 5.1 for several  $(\mu, t^2)$  -values in the ranges likely to be encountered in practice. In examining the table one should note that if  $\check{\theta} \equiv \hat{\theta}$  the expected MSE is 1.0; if  $\check{\theta} \equiv 0$  the expected MSE is  $E_g(\theta^2) = \mu^2 + t^2$ . The expected MSE for the optimal (Bayes) procedure must lie between 0.0 and 1.0.

The Bayes procedure is sharply superior for small  $t^2 (< 1)$  and  $\mu \approx 1$  values. Note that in some cases the expected MSE is greater than 1.0 for the Posterior  $C=1$  Procedure; in such cases the use of  $\check{\theta} \equiv \hat{\theta}$  would be a better procedure.

There are other properties of the procedures which can be compared. The Bayes procedure is "sensitive" to certain prior parameter values; at certain points slight changes in parameter values can make striking changes in the cutoff region A (over which  $\check{\theta} \equiv 0$ ) and its complement, B. The procedure is sensitive at  $\mu = 0$  and any value of  $t^2$ , at  $t^2 = 1$  and any value of  $\mu$ , and is particularly sensitive at  $\mu = 0$  and  $t^2 = 1$ . For example, if  $\mu = -0.01$ ,  $t^2 = 0.99$ , then



Table 5.1.--Continued

$\mu$	$t^2$	Posterior C=1 Procedure			Bayes Procedure		Ratio = $\frac{E(MSE:P \ C=1)}{E(MSE:Bayes)}$	
		$C_1$	$C_2$	E(MSE)	$C_1$	$C_2$	E(MSE)	
2.0	0.25	-13.000	-3.000	0.9999	0.0	5.3333	0.9379	1.0661
2.0	0.50	-7.000	-1.000	0.9683	0.0	8.0000	0.9217	1.0505
2.0	1.00	-4.000	0.0	0.8996	0.0	0.0	0.8995	1.0001
2.0	2.00	-2.500	0.500	0.9414	-4.0000	0.0	0.9108	1.0335
2.0	4.00	-1.750	0.750	1.0274	-1.3333	0.0	0.9792	1.0493
4.0	0.25	-21.000	-11.000	1.0000	0.0	10.6667	0.9997	1.0003
4.0	0.50	-11.000	-5.000	1.0000	0.0	16.0000	0.9990	1.0010
4.0	1.00	-6.000	-2.000	0.9999	0.0	0.0	0.9961	1.0038
4.0	2.00	-3.500	-0.500	0.9891	-8.0000	0.0	0.9858	1.0034
4.0	4.00	-2.250	0.250	0.9796	-2.6667	0.0	0.9770	1.0026

NOTE: For  $t^2 < 1$ , the cutoff region (on which  $\hat{\theta} = 0$ ) for the Bayes procedure is the complement of  $[C_1, C_2]$ ; for  $t^2 = 1$ , the cutoff region is  $(-\infty, 0)$ ; for  $t^2 > 1$  the region is  $[C_1, C_2]$ . The cutoff region for the Posterior C=1 Procedure is  $[C_1, C_2]$ .

$A = (-\infty, -2) + (0, +\infty)$ ,  $B = [-2, 0]$ . But if  $\mu = +0.01$ ,  $t^2 = 1.01$ , then  $A = [-2, 0]$  and  $B = (-\infty, -2) + (0, +\infty)$ . A change of only 0.02 in each of the parameters completely reverses the procedure! Table 5.1 reveals why the Bayes procedure is so sensitive; in the neighborhood of  $\mu = 0$ ,  $t^2 = 1$  it makes little difference whether one sets  $\check{\theta} = \hat{\theta}$  or  $\check{\theta} = 0$ ; the expected MSE's are nearly identical, and approximately equal to 1.0. A slight departure from the point  $\mu = 0$ ,  $t^2 = 1$  produces a very slight change in the expected MSE, which makes one of  $\check{\theta} = \hat{\theta}$  or  $\check{\theta} = 0$  slightly superior to the other; slight changes in opposite directions would produce opposite cutoff regions but nearly the same expected MSE. The Posterior C=1 Procedure is insensitive to small parameter value changes.

All the results discussed in this section hold for the normal prior distribution; other prior distributions may produce quite different results.

#### 5.4. Cutoff Point Selection with a Uniform Prior

This section is included as an example of the application of the Bayes decision and Posterior C=1 procedures for a non-normal prior. The posterior distribution is messy; no detailed analyses or comparisons of the two procedures will be attempted. Let

$$g(\theta; a, b) = \begin{cases} \frac{1}{b-a} & a \leq \theta \leq b \\ 0 & \text{elsewhere} \end{cases}$$

be the prior of  $\theta$ , and let  $f(\hat{\theta}|\theta)$  denote the  $N(\theta,1)$  density. Then the joint, marginal, and posterior densities are found as follows:

$$f(\hat{\theta}, \theta) = \frac{\exp\{-(\hat{\theta}-\theta)^2/2\}}{(b-a)\sqrt{2\pi}}, \quad \begin{matrix} -\infty < \hat{\theta} < +\infty \\ a \leq \theta \leq b \end{matrix}$$

$$f(\hat{\theta}) = \int_a^b f(\hat{\theta}, \theta) d\theta = \frac{1}{b-a} [\Phi(b-\hat{\theta}) - \Phi(a-\hat{\theta})], \quad -\infty < \hat{\theta} < +\infty$$

$$h(\theta|\hat{\theta}) = \frac{\exp\{-(\hat{\theta}-\theta)^2/2\}}{\sqrt{2\pi} [\Phi(b-\hat{\theta}) - \Phi(a-\hat{\theta})]}, \quad a \leq \theta \leq b.$$

The posterior is essentially the  $N(\theta,1)$  distribution truncated at  $a$  and  $b$ . The expected value of the posterior is:

$$\begin{aligned} E_h[\theta|\hat{\theta}] &= \frac{\int_a^b \theta h(\theta|\hat{\theta}) d\theta}{[\Phi(b-\hat{\theta}) - \Phi(a-\hat{\theta})]} \\ &= \frac{\int_{a-\hat{\theta}}^{b-\hat{\theta}} (t+\hat{\theta}) \exp\{-t^2/2\} dt}{\sqrt{2\pi} [\Phi(b-\hat{\theta}) - \Phi(a-\hat{\theta})]} \\ &= \hat{\theta} + \frac{\text{PFM}(a-\hat{\theta}, b-\hat{\theta})}{[\Phi(b-\hat{\theta}) - \Phi(a-\hat{\theta})]} = \mu^*. \end{aligned}$$

This function is easily evaluated on a computer for a particular  $(\hat{\theta}, a, b)$  value. Application of the Posterior C=1 Procedure is trivial once  $E_h(\theta)$  is available. The Bayes decision procedure is also trivial:  $\checkmark$   
 $\theta = 0$  if  $\hat{\theta}[2\mu^* - \hat{\theta}] < 0$ . It appears to be difficult to find explicit formulas for the C-values for these procedures. For particular  $(a, b)$  values the effective C-values may be found by a search procedure on a computer.

Estimator of  $\eta$  is simply the BLUE of the best linear approximating function of  $\eta$ .

In Section 4 a procedure was introduced for using one set of data both for determination of terms to be included in an estimation model and for estimation of the parameters in the resulting model. Due to the fact that when non-zero terms are deleted from an estimation model the resulting estimator is biased, the variance criterion was replaced by the Integrated Mean Square Error (IMSE) criterion for the comparison of estimators. It was shown in Section 4 that the IMSE of an estimator,  $\hat{\eta}$ , for  $\eta$  (both expressed in terms of integration-orthonormal functions) can be written as the termwise sum of the mean square errors of each of the coefficients. Because there are no covariance terms in this expression for the IMSE one can select terms to be included in the estimation model on a term-by-term basis. The procedure consists of examining each term in the model and setting the estimator of the corresponding coefficient equal to either zero (to delete the term from the model) or to the least squares estimate of the coefficient.

Three estimation procedures were considered. The first is called a "two-tail" estimator and is based on the UMP test of the hypothesis  $H_0: |\hat{\beta}_j| \geq \text{s.d.}(\hat{\beta}_j)$  versus  $H_a: |\hat{\beta}_j| < \text{s.d.}(\hat{\beta}_j)$ . The estimator is defined as:

$$\hat{\beta}_j \begin{cases} = \hat{\beta}_j & \text{if } |\hat{\beta}_j| > C \text{ s.d.}(\hat{\beta}_j) \\ = 0 & \text{if } |\hat{\beta}_j| \leq C \text{ s.d.}(\hat{\beta}_j), \end{cases}$$

## 6. SUMMARY

Techniques have been presented for attack on one of the very difficult problems facing applied statisticians, that of using one set of data both for determining which terms to include in a general linear model and for estimating the parameters of the resulting model.

The major result of Section 3 was the extension of the concept of Minimum Bias Estimation to a very general setting. First, some linear approximation theory was presented leading to the definition of the best linear approximating function of a given function  $\eta$  with respect to a given region of interest,  $R$ , weight measure  $W$ , and set of linearly independent functions,  $F$ . The best linear approximating function was given in terms of integration-orthogonal functions; algorithms were given for computing the set of integration-orthogonal functions from the set  $F$  of linearly independent functions, and for transformations of representations of  $\eta$  in terms of the linearly independent functions in  $F$  to representations in terms of the set of integration-orthogonal functions, and vice versa. It was shown that the best linear approximating function of  $\eta$  is easily computed; one simply deletes unwanted terms from the representation in terms of integration-orthogonal functions. Similarly, one computes the Best Linear Unbiased Estimator (BLUE) of the best linear approximating function by simply deleting terms from the BLUE of the integration-orthogonal-function representation of  $\eta$ . Finally, it was shown that the Minimum Bias

where the constant  $C$  is determined from other considerations (see Section 5) and where  $\text{s.d.}(\hat{\beta}_j)$  denotes the known or estimated standard deviation of  $\hat{\beta}_j$ .

The distribution, expected value (and bias) and mean square error of the estimator  $\check{\beta}_j$  were derived and the bias and MSE plotted for various  $C$ -values and various values of  $\nu$ , the degrees of freedom for the estimate of the variance of  $\hat{\beta}_j$ . This estimator was called the "two-tail" estimator because the resulting estimator can take positive or negative (or zero) values.

For the case in which the sign of a particular coefficient,  $\beta_j$ , is known from prior consideration, a "one-tail" estimation procedure was presented. Again, the BLUE,  $\hat{\beta}_j$ , of  $\beta_j$  is compared with the (known or estimated) standard deviation of  $\hat{\beta}_j$ ; the resulting estimator,  $\check{\beta}_j$ , is defined as:

$$\check{\beta}_j = \begin{cases} = \hat{\beta}_j & \text{if } \hat{\beta}_j > C \text{ s.d.}(\hat{\beta}_j) \\ = 0 & \text{if } \hat{\beta}_j \leq C \text{ s.d.}(\hat{\beta}_j). \end{cases}$$

(The above procedure is for  $\beta_j$  assumed positive; if  $\beta_j$  is assumed negative,  $\check{\beta}_j = \hat{\beta}_j$  if  $-\hat{\beta}_j > C \text{ s.d.}(\hat{\beta}_j)$ .) The distribution, expected value (and bias) and mean square error of the estimator  $\check{\beta}_j$  were derived, and the bias and mean square error were plotted for various  $C$ -values and various values of  $\nu$ , the degrees of freedom for the estimate of the variance of  $\hat{\beta}_j$ .

Finally, a generalized estimator of the form



$$\check{\beta}_j = \begin{cases} \hat{\beta}_j & \text{if } \hat{\beta}_j \in B \\ 0 & \text{if } \hat{\beta}_j \in A \end{cases} = \hat{\beta}_j \cdot 1_B(\hat{\beta}_j)$$

was considered, where  $A$  is a subset of the sample space and  $B$  is the complement of  $A$  relative to the sample space. The expected value (and bias) and mean square error of  $\check{\beta}$  were derived for  $A$ -sets which are finite intervals, complements of finite intervals, or half-lines. For the unknown variance case, the endpoints of  $A$  and  $B$  depend on the estimate of the variance.

Two techniques were presented in Section 5 for selection of the  $A$  and  $B$  sets for the generalized estimator. By specifying a prior distribution for a particular coefficient, one can select those sets which give the smallest expected value (with respect to the prior) of the mean square error of the estimator. A formula was given for computation of the estimator as a function of the mean of the posterior distribution. For the normal prior the boundaries of  $A$  and  $B$  were expressed in terms of the mean and variance of the prior. A second technique was presented and compared with the optimum technique for the case in which the prior is normal.

All of the results above hold for the response function expressed in terms of integration-orthonormal functions; it was shown in Section 4 that the IMSE of the estimator resulting from the above procedures is simply the sum of the mean square errors of the individual coefficients. Thus, the procedures can be applied term-by-term to each term in the "full model" (each term of "potential significance"). Since each term may or may not enter the final estimation model, the procedures, in

effect, allow consideration of "all possible regressions" subject to the restriction that terms are deleted by setting corresponding coefficients to zero. The procedure is remarkably easy to apply (compared with stepwise regression, for example): one simply compares each coefficient estimate with its standard deviation times a constant ( $C$ , or  $C_1$  and  $C_2$ ). Moreover, the properties of the final estimation model are known. (The properties of the final estimation model in stepwise regression are not known.)

In summary, a procedure has been presented which allows one to determine an estimation model and estimate the resulting model with the same data. The procedure is applied to an expression of a response function in terms of integration-orthogonal functions. (Transformations between representation in terms of integration-orthogonal functions and representation in terms of "standard" (linearly independent) functions consist of multiplication of a vector by an upper triangular matrix.) The procedure is easy to apply, is applied on a term-by-term basis to each potential term in the model, has known properties, and is closely related to Minimum Bias Estimation.

There are at least three areas in which further study is to be recommended:

- (1) A comparison of the estimators above with those produced by least squares estimation of the reduced model. The properties of the least squares estimator are design-dependent; comparisons would have to be made for various standard design.

(2) A study should be made to determine efficient designs for the estimator  $\hat{\eta}^v$  in various situations; the designs would depend on the region of interest and the weight function.

(3) The procedures given here should be extended to multivariate general linear model settings.

## LIST OF REFERENCES

- Abramowitz, M. and I. A. Stegun. 1964. Handbook of Mathematical Functions. National Bureau of Standards Applied Mathematics Series, No. 55. U. S. Government Printing Office, Washington, D. C.
- Bancroft, T. A. 1944. On Biases in Estimation Due to the Use of Preliminary Tests of Significance. Annals of Mathematical Statistics. 15: 190-204.
- Box, G. E. P. and N. R. Draper. 1959. A Basis for the Selection of a Response Surface Design. Jour. Amer. Stat. Assoc. 54: 622-654.
- Box, G. E. P. and N. R. Draper. 1963. The Choice of a Second Order Rotatable Design. Biometrika 50: 335-352.
- Davis, Phillip J. 1962. Orthonormalizing Codes in Numerical Analysis, pp. 347-379. In John Todd (ed.), Survey of Numerical Analysis. McGraw-Hill Book Company, Inc., New York.
- Draper, N. R. and H. Smith. 1967. Applied Regression Analysis. John Wiley and Sons, Inc., New York.
- Faddeeva, V. N. 1959. Computational Methods of Linear Algebra. Translated from the Russian by Curtis D. Benster, Dover Publications, Inc., New York.
- Gorman, J. W. and R. J. Toman. 1966. Selection of Variables for Fitting Equations to Data. Technometrics. 8: 27-51.
- Graybill, F. A. 1961. An Introduction to Linear Statistical Models. Volume 1. McGraw-Hill Book Company, Inc., New York.
- Halmos, Paul R. 1950. Measure Theory. D. Van Nostrand Company, Inc., New York.
- Handscomb, D. C. 1965. Methods of Numerical Approximation. Pergamon Press, New York.
- Hildebrand, F. B. 1956. Introduction to Numerical Analysis. McGraw-Hill Book Company, Inc., New York.

- Hocking, R. R. and N. N. Leslie. 1967. Selection of the Best Subset in Regression Analysis. *Technometrics* 9: 531-540.
- Hogg, Robert V. and Allen T. Craig. 1965. Introduction to Mathematical Statistics. The MacMillan Company. New York.
- Karson, M. J., A. R. Manson, and R. J. Hader, 1969. Minimum Bias Estimation and Experimental Design for Response Surfaces. (To appear in *Technometrics*, September, 1969.)
- Larson, Harold J. and T. A. Bancroft. 1963a. Sequential Model Building for Prediction in Regression Analysis, I. *The Annals of Mathematical Statistics*. 34: 462-479.
- Larson, Harold J. and T. A. Bancroft. 1963b. Biases in Prediction by Regression for Certain Incompletely Specified Models. *Biometrika* 50: 391-402.
- Lehmann, E. L. 1959. Testing Statistical Hypotheses. John Wiley & Sons, Inc., New York.
- Lindgren, B. W. 1962. Statistical Theory. The MacMillan Company. New York.
- Lindley, D. V. 1968. The Choice of Variables in Multiple Regression. *Journal of the Royal Statistical Society B*. 30: 31-66.
- Loeve, Michel. 1963. Probability Theory. D. Van Nostrand Company, Inc., New York.
- Michaels, Scott E. 1969. Optimum Design and Test/Estimation Procedures for Regression Models. Unpublished Ph. D. thesis, Department of Experimental Statistics, North Carolina State University, Raleigh, North Carolina.
- Schatzoff, R. Tsao and S. Feinberg. 1968. Efficient Calculation of all Possible Regressions. *Technometrics*. 10: 769-779.
- Sclove, S. L. 1968. Improved Estimators for Coefficients in Linear Regression. *Jour. Amer. Stat. Assoc.* 63: 596-606.
- Todd, John. 1962. Survey of Numerical Analysis. McGraw-Hill Book Company, Inc., New York.
- Toro, Carlos and T. D. Wallace. 1969. A Test of the Mean Square Error Criterion for Restrictions in Linear Regression. *Jour. Amer. Stat. Assoc.* 63: 558-572.