

## ABSTRACT

Title of dissertation: A THEORY OF CRAMÉR-RAO  
BOUNDS FOR CONSTRAINED  
PARAMETRIC MODELS

Terrence Joseph Moore, Jr.,  
Doctor of Philosophy, 2010

Dissertation directed by: Professor Benjamin Kedem  
Department of Mathematics

A simple expression for the Cramér-Rao bound (CRB) is presented for the scenario of estimating parameters  $\boldsymbol{\theta}$  that are required to satisfy a differentiable constraint function  $\boldsymbol{f}(\boldsymbol{\theta})$ . A proof of this constrained CRB (CCRB) is provided using the implicit function theorem, and the encompassing theory of the CCRB is proven in a similar manner. This theory includes connecting the CCRB to notions of identifiability of constrained parameters; the linear model under a linear constraint; the constrained maximum likelihood problem, its asymptotic properties and the method of scoring with constraints; and hypothesis testing. The value of the tools developed in this theory are then presented in the communications context for the convolutive mixture model and the calibrated array model.

Report Documentation Page				Form Approved OMB No. 0704-0188	
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.					
1. REPORT DATE <b>2010</b>		2. REPORT TYPE		3. DATES COVERED <b>00-00-2010 to 00-00-2010</b>	
4. TITLE AND SUBTITLE <b>A Theory of Cramer-Rao Bounds for Constrained Parametric Models</b>				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) <b>University of Maryland, College Park, MD, 20742</b>				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT <b>Approved for public release; distribution unlimited</b>					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT <b>A simple expression for the Cramer-Rao bound (CRB) is presented for the scenario of estimating parameters that are required to satisfy a differentiable constraint function <math>f(\cdot)</math>. A proof of this constrained CRB (CCRB) is provided using the implicit function theorem, and the encompassing theory of the CCRB is proven in a similar manner. This theory includes connecting the CCRB to notions of identifiability of constrained parameters; the linear model under a linear constraint the constrained maximum likelihood problem, its asymptotic properties and the method of scoring with constraints; and hypothesis testing. The value of the tools developed in this theory are then presented in the communications context for the convolutive mixture model and the calibrated array model.</b>					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT <b>Same as Report (SAR)</b>	18. NUMBER OF PAGES <b>78</b>	19a. NAME OF RESPONSIBLE PERSON
a. REPORT <b>unclassified</b>	b. ABSTRACT <b>unclassified</b>	c. THIS PAGE <b>unclassified</b>			

A THEORY OF CRAMÉR-RAO BOUNDS FOR CONSTRAINED  
PARAMETRIC MODELS

by

Terrence Joseph Moore, Jr.

Dissertation submitted to the Faculty of the Graduate School of the  
University of Maryland, College Park in partial fulfillment  
of the requirements for the degree of  
Doctor of Philosophy  
2010

Advisory Committee:  
Professor Benjamin Kedem, Chair/Advisor  
Professor Radu Balan  
Professor Paul Smith  
Professor Doron Levy  
Professor Prakash Narayan

© Copyright by  
Terrence Joseph Moore, Jr.  
2010

## Acknowledgments

There are perhaps too many people to thank here without the regret of omitting someone important to me. So many have helped me become a mathematician. To the very talented high school math teachers at West Hills High School in Santee, CA, who encouraged the best students to strive further, to the many remarkable professors at American University who delighted in interacting with their students and no matter how difficult the material ensured the experience was always enjoyable, and to the excellent class of professors at the University of Maryland who challenged and pushed their students to excel. To all these people, I thank you.

I would like to acknowledge my undergraduate and Master's thesis advisor, Dr. Stephen Casey at American University, for his direction in the beginning of my college education. I do need to acknowledge the contribution of my first advisor, Dr. Dennis Healy, who encouraged the construction of this current work from the beginning. He allowed me the freedom to develop the theory using an approach I developed and generated a positive exuberance for whatever results I discovered. I only wish he could have survived to see its completion. I also need to acknowledge my second advisor, Dr. Benjamin Kedem, for his willingness to allow me to complete this document as it was originally conceived as well as for guiding me through the defense process and the selection of the committee.

I also need to acknowledge the financial support I have received from my employer, the U.S. Army Research Lab. Additionally, the professional support and mentoring I received from Dr. Brian Sadler and the advice from team members at

the lab has been invaluable in assisting me succeed through the process. I also need to thank Dr. John Gowens, Mr. Glenn Racine, Dr. Barbara Broome, Dr. Alex Kott, and Dr. Brian Rivera for allowing me the time at and away from work to complete this degree.

# Table of Contents

List of Figures	vii
List of Symbols and Abbreviations	viii
1 INTRODUCTION	1
1.1 A note on the notation . . . . .	3
2 THE CRAMÉR-RAO BOUND	5
2.1 Definition . . . . .	6
2.1.1 Extensions . . . . .	7
2.2 Identifiability . . . . .	8
2.2.1 Local identifiability . . . . .	9
2.2.2 Strong Identifiability . . . . .	10
2.3 Linear Model . . . . .	11
2.3.1 Best Linear Unbiased Estimators . . . . .	12
2.3.2 Gaussian noise . . . . .	12
2.4 Maximum likelihood . . . . .	13
2.4.1 Efficient estimation . . . . .	14
2.4.2 Asymptotic Normality . . . . .	14
2.4.3 Scoring . . . . .	15
2.5 Hypothesis testing . . . . .	15
2.5.1 The Rao statistic . . . . .	15
2.5.2 The Wald statistic . . . . .	17
2.6 Discussion . . . . .	17
3 THE CONSTRAINED CRAMÉR-RAO BOUND	19
3.1 The Constrained CRB . . . . .	22
3.1.1 A proof of the CCRB . . . . .	25
3.1.2 Alternative formulas . . . . .	31
3.1.2.1 Gorman-Hero-Aitchison-Silvey CCRB . . . . .	32
3.1.2.2 Aitchison-Silvey-Crowder CCRB . . . . .	35
3.1.2.3 Xavier’s & Barroso’s Intrinsic Variance Lower Bound	36
3.1.3 Simple Extensions to the CCRB . . . . .	38
3.1.4 Properties of the CCRB . . . . .	41
3.1.5 Derivation of $\mathbf{g}(\boldsymbol{\xi})$ . . . . .	49
3.1.5.1 A Taylor series derivation . . . . .	50
3.1.5.2 A fixed point derivation . . . . .	51
3.2 Identifiability . . . . .	53
3.2.1 Local identifiability . . . . .	53
3.2.1.1 Local identifiability in the Aitchison-Silvey-Crowder	
CCRB formula . . . . .	57
3.2.2 Strong Identifiability . . . . .	58
3.3 Linear Model . . . . .	59

3.3.1	Best Linear Unbiased Estimation . . . . .	60
3.3.2	Uniform Minimum Variance Estimation under Gaussian noise . . . . .	63
3.4	Constrained Maximum Likelihood Estimation . . . . .	66
3.4.1	Efficient estimation . . . . .	68
3.4.2	Asymptotic Normality . . . . .	69
3.4.3	The Method of Scoring Under Parametric Constraints . . . . .	71
3.4.3.1	Convergence Properties . . . . .	75
3.4.3.2	Linear constraints . . . . .	78
3.5	Hypothesis testing . . . . .	81
3.5.1	The Rao statistic . . . . .	82
3.5.2	The Wald statistic . . . . .	83
3.6	Discussion . . . . .	84
4	Applications of the CCRB in Communications Models . . . . .	86
4.1	Convolutional Mixture Model . . . . .	88
4.1.1	Equivalent Convolutional Mixture Models . . . . .	90
4.1.1.1	The Vector-Matrix Model . . . . .	90
4.1.1.2	The Z transform model . . . . .	91
4.1.2	Strict Identifiability . . . . .	95
4.1.3	The Fisher information of the convolutional mixture model . . . . .	98
4.1.3.1	Complex-valued Fisher information . . . . .	98
4.1.3.2	CFIM for the convolutional mixture model . . . . .	99
4.1.3.3	Properties of the CFIM . . . . .	100
4.1.4	Constraints for the convolutional mixture model . . . . .	108
4.1.4.1	Pathways to regularity . . . . .	109
4.1.4.2	Norm channel + real-valued source constraint . . . . .	110
4.1.4.3	Semibind constraints: $s^{(k)}(t) = p(t)$ for $t \in \mathbb{T}$ . . . . .	112
4.1.4.4	Unit Modulus constraint + Semibind constraint . . . . .	113
4.2	Calibrated Array Model . . . . .	117
4.2.1	The Fisher information of the calibrated array model . . . . .	118
4.2.1.1	Indirect derivation of the FIM . . . . .	118
4.2.1.2	Direct derivation of the FIM . . . . .	120
4.2.1.3	Properties of the FIM . . . . .	121
4.2.2	Constraints for the calibrated array model . . . . .	122
4.2.2.1	Constraints on the complex-valued channel gain: $\mathbf{\Gamma} = \mathbf{I}_{K \times K}$ . . . . .	122
4.2.2.2	Semibind constraints: $\mathbf{s}(t) = \mathbf{p}(t)$ for $t \in \mathbb{T}$ . . . . .	124
4.2.2.3	Finite alphabet constraint: $s^{(k)}(n) \in \mathbb{S}$ . . . . .	125
4.2.2.4	Unit modulus constraints: $ \mathbf{s}(n)  = 1$ for all $n$ . . . . .	127
4.2.2.5	Unit modulus constraint; real-valued channel gain: $\text{Im}(\gamma_k) = 0$ for all $k$ . . . . .	129
4.2.2.6	Semi-blind and unit modulus constraint . . . . .	129
4.3	Discussion . . . . .	133



A	Appendices	136
A.1	A proof of the CCRB using the Chapman-Robbins version of the Barankin bound . . . . .	136
A.2	A proof of the CCRB using the method of implicit differentiation . .	138
A.3	Alternative proof of asymptotic normality . . . . .	140
B	Proofs of Convergence Properties of Constrained Scoring	142
C	Proofs of Theorems in Chapter 4	147
	Bibliography	161

## List of Figures

3.1	Reparameterization of $\mathbf{f}(\boldsymbol{\theta}) = 0$ to $\boldsymbol{\theta} = \mathbf{g}(\boldsymbol{\xi})$ . . . . .	26
3.2	Projection of the observations $\mathbf{x}$ onto the linear space $\mathbf{H}\boldsymbol{\Theta}$ and the linear constraint space $\mathbf{H}\boldsymbol{\Theta}_f$ . . . . .	61
3.3	Path created by iterates from the method of scoring with constraints. . . . .	74
4.1	Convolutional Mixture: Finite Impulse Response (tapped delay line) model. . . . .	89
4.2	Convolutional Mixture: Norm-constrained channel estimation performance. . . . .	111
4.3	Convolutional Mixture: Example of multipath channel. . . . .	115
4.4	Convolutional Mixture: Source estimation with varying $\Delta\psi$ . . . . .	116
4.5	Calibrated array model geometries. . . . .	118
4.6	Calibrated Array: CCRBs on AOA for blind, constant modulus, and known signal models. . . . .	131
4.7	Calibrated Array: CCRBs on AOA for semiblind, constant modulus + semiblind, and known signal models. . . . .	132
4.8	Calibrated Array: CCRBs on signal phase for blind, constant modulus, semiblind, and constant modulus + semiblind. . . . .	134

## List of Symbols and Abbreviations

### Symbols

$\mathbf{b}(\boldsymbol{\theta})$	the bias of an estimator of $\boldsymbol{\theta}$
$\mathbb{C}$	the complex numbers
$\mathbb{C}^*$	$\mathbb{C} \cup \{\infty\}$
$\mathcal{CN}(\mathbf{a}, \mathbf{B})$	complex normal with mean $\mathbf{a}$ and (co)variance matrix $\mathbf{B}$
$\chi_r^2$	chi-square distribution with $r$ degrees of freedom
$\xrightarrow{d}$	convergence in distribution
$E_\phi(\cdot)$	the expectation of $(\cdot)$ with respect to the density modeling $\phi$
$\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$	the constraint function
$\mathbf{F}(\boldsymbol{\theta})$	the Jacobian (or gradient) of the constraint function at $\boldsymbol{\theta}$
$\mathbf{I}(\boldsymbol{\theta})$	the Fisher information for the parameter $\boldsymbol{\theta}$
$\mathbf{I}_{m \times m}$ or $\mathbf{I}_m$	the identity matrix of size $m$
$\text{Im}(\cdot)$	the imaginary part of $(\cdot)$
$\mathbf{s}(\mathbf{x}; \boldsymbol{\theta})$	the Fisher score
$\mathcal{N}(\mathbf{a}, \mathbf{B})$	normal with mean $\mathbf{a}$ and variance (covariance) matrix $\mathbf{B}$
$p(\mathbf{x}; \boldsymbol{\theta})$	the probability density function (pdf)
$\mathbb{R}$	the real numbers
$\text{Re}(\cdot)$	the real part of $(\cdot)$
$\mathbf{t}(\mathbf{x})$	an estimator of some parameter (vector)
$\boldsymbol{\theta}$	the parameter
$\Theta$	the parameter space
$\Theta_f$	the constraint manifold (parameter space)
$\hat{\boldsymbol{\theta}}_{\text{ML}}(\mathbf{x})$	the ML estimator of $\boldsymbol{\theta}$ depending on the observation $\mathbf{x}$
$\mathbf{x}$	the observation (sample)
$(\cdot)_+$	$(\cdot)$ if the value is nonnegative, otherwise zero

### Abbreviations

BER	bit-error rate
BLUE	best linear unbiased estimate (or estimator)
CCRB	constrained Cramér-Rao bound
CFIM	complex FIM
CLSE	constrained least squares estimate (or estimator)
CML	constrained maximum likelihood
CMLE	constrained maximum likelihood estimate (or estimator)
CRB	Cramér-Rao bound
FIM	Fisher information matrix
FIR	finite impulse response
iid	independent and identically distributed
LSE	least squares estimate (or estimator)
MIMO	multiple-input, multiple-output
ML	maximum likelihood
MLE	maximum likelihood estimate (or estimator)
MSE	mean-square error
MVUE	minimum variance unbiased estimate (or estimator)

NLB	nullity lower bound of the FIM (or CFIM)
SER	symbol-error rate
SIMO	single-input, multiple-output
SISO	single-input, single-output

## Chapter 1

### INTRODUCTION

The Cramér-Rao bound (CRB), the inverse of the Fisher information, is a limit on the performance of the estimation of parameters under certain conditions. Hence, the variance of any unbiased estimator cannot be lower than this bound, and the CRB may also be interpreted in some sense as a measure of performance potential. For a large number of scenarios, it is often of interest to gauge the performance of estimation under parametric constraints. The traditional approach in deriving CRBs for these cases is to find a reparameterization that represents the constraint. However, such an approach is not always feasible for a large class of constraints and numerical approximations are often restricted to the particular model. An alternative and equivalent option to reparameterization, the constrained Cramér-Rao bound, is presented herein, which is also computationally simple to program.

In communications design research, estimation performance metrics are often interpreted to represent performance potential to study the feasibility of a model to meet a certain measure of desired reliability. This approach is often practical since it avoids the necessity of searching for the best performer over a class of estimators for a particular trial model and since the CRB is analytically or numerically simple to compute. The downside to this approach is that for each model, the CRB needs

to be derived each time. This effectively prohibits the practitioner from studying an overly large class of models or, in the context of this work, a large class of constraints on some base model. This restriction is to a great extent eliminated using the constrained Cramér-Rao bound, as the Fisher information, which involves an integration over many variables, for the base model only needs to be evaluated once.

Chapter 2 offers a quick review of several connections with the Cramér-Rao bound (CRB) within mathematical statistics and serves as a reference point for the study of parameters under parametric equality constraints, discussed in Chapter 3. With the possible exception of the identifiability relationship in section 2.2, much of this section is familiar and well represented in standard mathematical statistics texts.

In Chapter 3, a general theory of the constrained Cramér-Rao bound (CCRB) is presented. In section 3.1, the CCRB is defined and proven, alternative formulas are presented, and several interesting properties of the bound are detailed. In section 3.2, a connection is made between the CCRB and two different notions of identifiability under certain conditions. In section 3.3, the linear model with linear constraints is examined in the context of the CCRB. In section 3.4, connections between the CCRB and constrained maximum likelihood estimation are detailed, including an asymptotic normality result and an adaptation of the method of scoring to the constrained parameter scenario. This chapter concludes with the consideration of hypothesis testing under constraints in section 3.5. Chapters 2 and 3 are designed so that the section numbers correspond directly, i.e., section 3.x relates

a concept for constrained parameters that section 2.x reviews for unconstrained parameters.

In Chapter 4, the analytic tools developed in Chapter 3 are applied in the communications context of the convolutive mixture model (section 4.1) and the calibrated array model (section 4.2). These models are defined and their Fisher information matrices developed in section 4.1.3 and section 4.2.1, respectively. A variety of constraints for these models are considered in sections 4.1.4 and 4.2.2.

## 1.1 A note on the notation

All elements will be denoted in lowercase math font:  $a$ . All vectors will be column vectors and be denoted in a lowercase bold math font:  $\mathbf{a}$ . Hence, the  $i$ th element of the column vector  $\mathbf{a}$  will be denoted as  $a_i$ . (This should not be confused with  $\mathbf{a}_i$ , which is often used as a subvector of the vector  $\mathbf{a}$  and will be defined in context.) All matrices will be denoted in an uppercase bold math font:  $\mathbf{A}$ . All scalars, vectors, and matrices are assumed to have elements with real-valued numbers unless otherwise noted as complex-valued (where the complex number  $i = j = \sqrt{-1}$  should be clear from context).

For vectors and matrices,  $(\cdot)^T$  will denote the transpose operator (do not with confuse  $(\cdot)'$ , which is occasionally used here as a dummy variable),  $(\cdot)^*$  will denote the conjugate operator (do not confuse with  $(\cdot)^\star$ , which is occasionally used as a variant of another vector or matrix), and  $(\cdot)^H$  will denote the Hermitian (or conjugate transpose) operator. When a vector depends on another vector value as in  $\mathbf{a}(\boldsymbol{\theta})$ , then

the Jacobian will be a matrix denoted as  $\mathbf{A}(\boldsymbol{\theta})$  where the  $i$ th row is the transposed of the vector  $\frac{\partial}{\partial \boldsymbol{\theta}} a_i(\boldsymbol{\theta})$  and  $a_i(\boldsymbol{\theta})$  is the  $i$ th row element of  $\mathbf{a}(\boldsymbol{\theta})$ . For square matrices,  $(\cdot)^{-1}$  will denote the inverse of the matrix and  $(\cdot)^\dagger$  will denote the pseudoinverse of the matrix. Of course,  $(\cdot)^2$  will denote the square of the matrix.

For symmetric matrices (or Hermitian matrices in the complex-valued case, the expression  $\mathbf{A} > \mathbf{B}$  will denote that the matrix  $\mathbf{A} - \mathbf{B}$  is positive definite. Similarly,  $\mathbf{A} \geq \mathbf{B}$  will denote that the matrix  $\mathbf{A} - \mathbf{B}$  is positive semidefinite.

Sets will be denoted in an uppercase blackboard math font:  $\mathbb{A}$  or in an uppercase Greek letter math font:  $\Theta$ . Also, all sets will be assumed open sets unless otherwise noted.

For convenience to allow the reader to find referenced items, numbered theorems, corollaries, and examples share numbering. Thus, the theorem immediately following Example x.4 in Chapter x is numbered Theorem x.5 even though the previous theorem is Theorem x.1.



## Chapter 2

### THE CRAMÉR-RAO BOUND

The Cramér-Rao bound (CRB) is a lower bound on the error covariance of any unbiased estimator under certain regularity conditions. As such, it is a measure of the optimal performance of an unbiased estimator under a given model. There are other mean-square error lower bounds, such as the Ziv-Zakai bound [78], the Hammersley-Chapman Robbins bound [26, 16], the Barankin bound [8], the Bhattacharyya bound [11], etc., and indeed, depending on the model, there are numerous other possible performance measures, such as classification bounds like bit-error rate (BER) or symbol-error rate (SER), but the CRB remains a very popular benchmark due to its computational simplicity and its underlying well-developed theory.

This theory has led to numerous connections in areas of mathematical statistics, e.g., identifiability, linear models, maximum likelihood, including asymptotic normality and the method of scoring, and hypothesis testing. This short chapter is a quick review of just a few of these connections. A more complete discussion of the utility and application of the topics discussed in this chapter, as well as proofs of the definitions, theorems and statements herein, may be found in many standard statistical inference texts, such as Shao [62] or Casella and Berger [14] for statisticians, Kay [36] for signal processors, or Van Trees [72] for engineers.

## 2.1 Definition

Suppose we have an observation  $\mathbf{x}$  in  $\mathbb{X} \subset \mathbb{R}^n$  from a probability density function (pdf)  $p(\mathbf{x}; \boldsymbol{\theta})$  where  $\boldsymbol{\theta}$  is in an open set  $\Theta \subset \mathbb{R}^m$  is a vector of deterministic parameters. The Fisher information matrix (FIM) of this model is given by

$$\mathbf{I}(\boldsymbol{\theta}) \triangleq E_{\boldsymbol{\theta}} \{ \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}) \mathbf{s}^T(\mathbf{x}; \boldsymbol{\theta}) \}$$

where  $\mathbf{s}(\mathbf{x}; \boldsymbol{\theta})$  is the Fisher score defined by

$$\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}) \triangleq \left. \frac{\partial \log p(\mathbf{x}; \boldsymbol{\theta}')}{\partial \boldsymbol{\theta}'} \right|_{\boldsymbol{\theta}' = \boldsymbol{\theta}}$$

and the expectation is evaluated at  $\boldsymbol{\theta}$ , i.e.,  $E_{\boldsymbol{\theta}}(\cdot) = \int_{\mathbb{X}} (\cdot) p(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x}$ . And suppose as *regularity conditions*, the pdf is differentiable with respect to  $\boldsymbol{\theta}$  and satisfies

$$\frac{\partial}{\partial \boldsymbol{\theta}^T} E_{\boldsymbol{\theta}}(\mathbf{h}(\mathbf{x})) = E_{\boldsymbol{\theta}}(\mathbf{h}(\mathbf{x}) \mathbf{s}^T(\mathbf{x}; \boldsymbol{\theta})) \quad (2.1)$$

for  $\mathbf{h}(\mathbf{x}) \equiv 1$  and  $\mathbf{h}(\mathbf{x}) \equiv \mathbf{t}(\mathbf{x})$  where  $\mathbf{t}(\mathbf{x})$  is an unbiased estimator of  $\boldsymbol{\theta}$  [62]. These conditions are assured under a number of scenarios, for example, when the Jacobian and Hessian of the density function  $p(\mathbf{x}; \boldsymbol{\theta})$  is absolutely integrable with respect to both  $\mathbf{x}$  and  $\boldsymbol{\theta}$  [72], and essentially permit switching between the order of integration and differentiation. Under these assumptions, we have the following information inequality theorem [14, 62, 36, 72], independently developed by Cramér [17] and Rao [56].

**Theorem 2.1.** The Cramér-Rao bound is the inverse of the FIM,

$$\text{CRB}(\boldsymbol{\theta}) \triangleq \mathbf{I}^{-1}(\boldsymbol{\theta}), \quad (2.2)$$

if it exists, and the variance of any unbiased estimator  $\mathbf{t}(\mathbf{x})$  satisfies the inequality

$$\text{Var}(\mathbf{t}(\mathbf{x})) \geq \text{CRB}(\boldsymbol{\theta})$$

with equality if and only if  $\mathbf{t}(\mathbf{x}) - \boldsymbol{\theta} = \text{CRB}(\boldsymbol{\theta})\mathbf{s}(\mathbf{x}; \boldsymbol{\theta})$  in the mean-square sense.

**Example 2.2.** Let  $\mathbf{x} \sim \mathcal{CN}(\vartheta, \sigma^2)$  with unknown complex-valued mean  $\vartheta$  and known variance  $\sigma^2$ . In terms of real-valued parameters, the equivalent model is

$$\begin{bmatrix} \text{Re}(\mathbf{x}) \\ \text{Im}(\mathbf{x}) \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \text{Re}(\vartheta) \\ \text{Im}(\vartheta) \end{bmatrix}, \frac{\sigma^2}{2} \mathbf{I}_{2 \times 2}\right).$$

From a well-known result on normal distributions [36, equation (3.31)],

$$\mathbf{I}\left(\begin{bmatrix} \text{Re}(\vartheta) \\ \text{Im}(\vartheta) \end{bmatrix}\right) = \frac{2}{\sigma^2} \mathbf{I}_{2 \times 2}$$

and the CRB is  $\frac{\sigma^2}{2} \mathbf{I}_{2 \times 2}$ .

### 2.1.1 Extensions

The performance of a function of parameters, e.g. the transformation  $\boldsymbol{\alpha} = \mathbf{k}(\boldsymbol{\theta})$ , is often of more interest than the performance of the parameters. If the Jacobian of the transformation function is  $\mathbf{K}(\boldsymbol{\theta}) = \left. \frac{\partial \mathbf{k}(\boldsymbol{\theta}')}{\partial \boldsymbol{\theta}'^T} \right|_{\boldsymbol{\theta}' = \boldsymbol{\theta}}$ , then the CRB on the performance of an unbiased estimator of  $\alpha$  is [62, 36]

$$\text{CRB}(\boldsymbol{\alpha}) = \mathbf{K}(\boldsymbol{\theta}) \mathbf{I}^{-1}(\boldsymbol{\theta}) \mathbf{K}^T(\boldsymbol{\theta}),$$

i.e., if  $\mathbf{S}(\mathbf{x})$  is an unbiased estimator of  $\boldsymbol{\alpha}$ , then  $\text{Var}(\mathbf{S}(\mathbf{x})) \geq \mathbf{K}(\boldsymbol{\theta}) \text{CRB}(\boldsymbol{\theta}) \mathbf{K}^T(\boldsymbol{\theta})$ . Implicit in this inequality for the transformation is that  $\boldsymbol{\alpha}$  is differentiable with respect to  $\boldsymbol{\theta}$  and (2.1) must also be satisfied for  $\mathbf{h}(\mathbf{x}) \equiv \mathbf{S}(\mathbf{x})$ . Consequently, if

an estimator  $\mathbf{t}(\mathbf{x})$  is biased with bias  $\mathbf{b}(\boldsymbol{\theta}) = E_{\boldsymbol{\theta}}\mathbf{t}(\mathbf{x}) - \boldsymbol{\theta}$ , then the transformation formula above can be used to attain a bound for  $\boldsymbol{\alpha} = \boldsymbol{\theta} + \mathbf{b}(\boldsymbol{\theta})$ . Then  $\text{Var}(\mathbf{t}(\mathbf{x})) \geq \text{CRB}(\boldsymbol{\theta} + \mathbf{b}(\boldsymbol{\theta}))$  where

$$\text{CRB}(\boldsymbol{\theta} + \mathbf{b}(\boldsymbol{\theta})) = (\mathbf{I}_m + \mathbf{B}(\boldsymbol{\theta})) \text{CRB}(\boldsymbol{\theta}) (\mathbf{I}_m + \mathbf{B}^T(\boldsymbol{\theta}))$$

with  $\mathbf{B}(\boldsymbol{\theta}) = \left. \frac{\partial \mathbf{b}(\boldsymbol{\theta}')}{\partial \boldsymbol{\theta}'^T} \right|_{\boldsymbol{\theta}' = \boldsymbol{\theta}}$ .

Theorem 2.1 requires a nonsingular Fisher information, however, there are a number of interesting cases where this requirement can not be met yet the model is still of interest. For this scenario, the pseudoinverse of the FIM is occasionally used as a bound in place of the CRB, i.e.,

$$\text{Var}(\mathbf{t}(\mathbf{x})) \geq \mathbf{I}^\dagger(\boldsymbol{\theta})$$

for an unbiased estimator  $\mathbf{t}(\mathbf{x})$ . This bound inequality is trivial for nonidentifiable functions of the parameters [66], i.e., the variance is finite only if  $\mathbf{H}(\boldsymbol{\theta}) = \mathbf{H}(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{I}^\dagger(\boldsymbol{\theta})$  where  $\mathbf{H}(\boldsymbol{\theta}) = \mathbf{K}(\boldsymbol{\theta}) + \mathbf{B}(\boldsymbol{\theta})$  and  $\mathbf{t}(\mathbf{x})$  is a biased estimator of  $\boldsymbol{\alpha} = \mathbf{k}(\boldsymbol{\theta})$  with bias  $\mathbf{b}(\boldsymbol{\theta})$ .

## 2.2 Identifiability

The ability to identify parameters determines the validity and utility of certain structural models. Criteria on the identifiability of parameters has numerous connections to parametric statistical measures, such as Kullback-Leibler distance [13] and the Fisher information matrix [58, 29, 69]. In this section, two of these connections are developed to establish conditions under which a particular parametric model is identifiable.

### 2.2.1 Local identifiability

To proceed in examining the identifiability criterion from the CRB, a definition of identifiability is required. A parameter  $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^m$  is *identifiable* in the model  $p(\mathbf{x}; \cdot)$  if there is no other  $\boldsymbol{\theta}' \in \Theta$  such that  $p(\mathbf{x}; \boldsymbol{\theta}) = p(\mathbf{x}; \boldsymbol{\theta}')$  for all  $\mathbf{x} \in \mathbb{R}^n$ . A parameter is *locally identifiable* if there exists an open neighborhood of  $\boldsymbol{\theta}$  such that  $\boldsymbol{\theta}$  is identifiable in that neighborhood.

A parameter is (locally) identifiable in the additive noiseless case if the parameter is solvable (locally). *Estimable* parameters, i.e., expected values of functions of the observations [57], are also identifiable.<sup>1</sup> Hence, a non-identifiable parameter is not estimable regardless of the scheme or the number of observations. These scenarios exist when some inherent ambiguity exists in the model.

**Example 2.3.** The parameter vector  $\boldsymbol{\theta} = [\theta_1, \theta_2]^T$  in the model

$$x = \theta_1 + \theta_2 + w,$$

where  $w$  represents the observation noise, is not identifiable since it is indistinguishable from the parameter vector  $\boldsymbol{\theta}' = [\theta_1 + a, \theta_2 - a]^T$  for any real number  $a$ .

The Fisher information matrix is called *regular* if  $\mathbf{I}(\boldsymbol{\theta})$  is full rank, and  $\boldsymbol{\theta}$  is said to be a *regular point* of  $\mathbf{I}(\boldsymbol{\theta})$ . If the FIM is singular,  $\boldsymbol{\theta}$  is a *singular point*.

**Example 2.4.** The FIM for the model

$$x = \theta^2 + w,$$

---

<sup>1</sup>The converse is not true. While it is possible to develop estimation schemes for any identifiable parameters, there is no guarantee that those estimators will be unbiased.

where  $w \sim \mathcal{N}(0, 1)$ , is  $I(\theta) = 4\theta^2$ . For any  $\theta \neq 0$ ,  $\theta$  is locally identifiable but not identifiable and at the same time a regular point of the FIM. For  $\theta = 0$ ,  $\theta$  is a singular point, yet is identifiable.

Fisher information-regularity implies local identifiability, but as the example demonstrates the converse is not true. Rothenberg [58] found a connection between local identifiability and the FIM under certain conditions.

**Theorem 2.5** (Rothenberg). Assume the FIM  $\mathbf{I}(\boldsymbol{\theta})$  has constant rank locally about  $\boldsymbol{\theta}$ . Then  $\boldsymbol{\theta}$  is locally identifiable if and only if  $\mathbf{I}(\boldsymbol{\theta})$  is regular.

### 2.2.2 Strong Identifiability

Suppose that  $\mathbf{p}(\mathbf{x}; \boldsymbol{\theta})$  is a normal pdf with mean  $\boldsymbol{\mu}(\boldsymbol{\theta}) \in \mathbb{R}^p$  and variance  $\boldsymbol{\Sigma}(\boldsymbol{\theta})$ , whose elements may be explicitly defined by a map  $\boldsymbol{\varphi} : \Theta \rightarrow \mathbb{R}^q$  where  $q \leq p + p(p + 1)/2$  and it is assumed  $m \leq q$ . Then by the given definitions, (local) identifiability holds when  $\boldsymbol{\varphi}$  is injective (locally) and since by a transformation on the FIM [12, p. 157]

$$\mathbf{I}(\boldsymbol{\theta}) = \frac{\partial \boldsymbol{\varphi}^T(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \mathbf{I}(\boldsymbol{\varphi}(\boldsymbol{\theta})) \frac{\partial \boldsymbol{\varphi}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^T}$$

regularity holds when the Jacobian  $\frac{\partial \boldsymbol{\varphi}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^T}$  has full rank  $m$ .

Suppose there exists a set of indices  $i_1, \dots, i_m \in \{1, \dots, q\}$  that make  $\boldsymbol{\varphi}^*(\boldsymbol{\theta}) = [\boldsymbol{\varphi}_{i_1}(\boldsymbol{\theta}), \dots, \boldsymbol{\varphi}_{i_m}(\boldsymbol{\theta})]^T$  injective on  $\Theta$ . Then each  $\boldsymbol{\theta} \in \Theta$  is *strongly identifiable* and  $\boldsymbol{\varphi}^*$  is a *representative mapping*. By this definition, if  $\mathbf{I}(\boldsymbol{\theta})$  is regular at  $\boldsymbol{\theta}$  then  $\boldsymbol{\theta}$  is in a strongly identifiable open neighborhood, and if  $\boldsymbol{\theta}$  is strongly identifiable on  $\Theta$  then it is also identifiable on  $\Theta$ . The converses are not generally true. The following

theorem establishes conditions under which the converses are true [29].

**Theorem 2.6** (Hochwald and Nehorai). Let  $\varphi : \Omega \rightarrow \mathbb{C}^q$  be a holomorphic mapping of  $\mathbf{z} \in \Omega \subset \cup_{\alpha \in A} \Omega_\alpha$ , where  $\Theta \subset \Omega \subset \mathbb{C}^m$  and  $\Omega_\alpha$  is open in  $\mathbb{C}^m$  for each  $\alpha$ . Then

- (a) if  $\mathbf{I}(\mathbf{z})$  is regular, there exists a strongly identifiable open neighborhood about  $\mathbf{z}$ , and
- (b) if there exists a representative mapping  $\varphi_\alpha^* : \Omega_\alpha \rightarrow \mathbb{C}^q$  for each  $\alpha$ ,  $\mathbf{I}(\mathbf{z})$  is regular for every  $\mathbf{z} \in \Omega$ .

Therefore, the existence of a proper holomorphic function(s) equates Fisher information regularity with strong identifiability for normal distributions. And locally constant rank in the FIM equates regularity with local identifiability for arbitrary distributions.

## 2.3 Linear Model

Suppose we have observations  $\mathbf{x}$  from a linear model

$$\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w}, \tag{2.3}$$

on  $\boldsymbol{\theta}$ , where  $\mathbf{H}$  is an observation matrix consisting of known elements and  $\mathbf{w}$  is the noise from the observations with mean zero and variance  $\mathbf{C}$ .

### 2.3.1 Best Linear Unbiased Estimators

The Gauss-Markov theorem [14] states that the best linear unbiased estimator (BLUE) is given by the (weighted) least squares solution

$$\hat{\boldsymbol{\theta}}_{\text{LS}}(\mathbf{x}) = (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}^{-1} \mathbf{x}, \quad (2.4)$$

so called for minimizing the (weighted) least squares  $(\mathbf{x} - \mathbf{H}\boldsymbol{\theta})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H}\boldsymbol{\theta})$ . For any other LUE of  $\boldsymbol{\theta}$ , i.e.  $\mathbf{A}\mathbf{x}$ , then  $\text{Var}(\mathbf{A}\mathbf{x}) \geq \text{Var}(\hat{\boldsymbol{\theta}}_{\text{LS}}(\mathbf{x}))$  with equality if and only if  $\mathbf{A} = (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}^{-1}$ . This assumes a full column rank observation matrix  $\mathbf{H}$ . Otherwise, for any estimable function of the parameters  $\mathbf{d}^T \boldsymbol{\theta}$ , where  $\mathbf{d}$  is in the column space of the transposed observation matrix  $\mathbf{H}^T$ , its least squares solution is  $\mathbf{d}^T \hat{\boldsymbol{\theta}}_{\text{LS}}(\mathbf{x})$  where

$$\hat{\boldsymbol{\theta}}_{\text{LS}}(\mathbf{x}) = (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^\dagger \mathbf{H}^T \mathbf{C}^{-1} \mathbf{x}, \quad (2.5)$$

and  $(\cdot)^\dagger$  is the generalized pseudoinverse of  $(\cdot)$  [57, theorems 11.2B, 11.3A-D]. This solution is also the BLUE with variance  $\mathbf{d}^T (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^\dagger \mathbf{d}$ .

### 2.3.2 Gaussian noise

Additionally, if the noise has a Gaussian distribution, i.e.,  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$ , then the least squares solution is also the maximum likelihood estimator (MLE) and the minimum variance unbiased estimate (MVUE).

**Theorem 2.7.** If the observations obey the linear model in (2.3), where  $\mathbf{H}$  is a known full column rank matrix,  $\boldsymbol{\theta}$  is an unknown parameter vector, and  $\mathbf{w}$  is a



zero-mean normal random vector with variance  $\mathbf{C}$ , then the MVUE is

$$\hat{\boldsymbol{\theta}}_{\text{LS}}(\mathbf{x}) = (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}^{-1} \mathbf{x} \quad (2.6)$$

with estimator covariance equaling the CRB  $\mathbf{I}^{-1}(\boldsymbol{\theta}) = (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1}$ .

Similarly, if  $\mathbf{H}$  is not full rank, and  $\mathbf{d}^T \boldsymbol{\theta}$  is an estimable function, then the MLE is  $\mathbf{d}^T \hat{\boldsymbol{\theta}}_{\text{MLE}}(\mathbf{x})$  where  $\hat{\boldsymbol{\theta}}_{\text{MLE}}(\mathbf{x}) = \hat{\boldsymbol{\theta}}_{\text{LS}}(\mathbf{x})$  from (2.5). This MLE is also the MVUE [57, theorems 11.3F-G].

## 2.4 Maximum likelihood

Given observations  $\mathbf{x}$  from a likelihood (or pdf)  $p(\mathbf{x}; \boldsymbol{\theta})$  depending on an unknown parameter  $\boldsymbol{\theta}$ , a popular method of estimating the parameter is the method of maximum likelihood. This approach chooses as an estimator  $\hat{\boldsymbol{\theta}}_{\text{ML}}(\mathbf{x})$  that, if true, would have the highest probability (the maximum likelihood) of resulting in the given observations  $\mathbf{x}$ , i.e., the optimization problem:

$$\hat{\boldsymbol{\theta}}_{\text{ML}}(\mathbf{x}) = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{x}; \boldsymbol{\theta})$$

where for convenience the log-likelihood is equivalently maximized since  $\log(\cdot)$  is monotone. An analytic solution of the MLE can be found from the first-order conditions on the log-likelihood by considering solutions  $\dot{\boldsymbol{\theta}}(\mathbf{x})$  of

$$\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}') = \mathbf{0}. \quad (2.7)$$

This is the *method of maximum likelihood*. Provided  $\Theta$  is an open set,  $\hat{\boldsymbol{\theta}}_{\text{ML}}(\mathbf{x})$  will satisfy (2.7).

### 2.4.1 Efficient estimation

If an efficient estimator exists, it is well-known that the method of maximum-likelihood finds the estimator [36, exercise 7.12], i.e., such an estimator must be a stationary point of the maximum-likelihood optimization problem. More formally, we have the following theorem.

**Theorem 2.8.** If  $\mathbf{t}(\mathbf{x})$  is an estimator of  $\boldsymbol{\theta}$ , which is also efficient with respect to the CRB, then the estimator is a stationary point of the following optimization problem:

$$\max_{\boldsymbol{\theta}} \log p(\mathbf{x}; \boldsymbol{\theta}).$$

### 2.4.2 Asymptotic Normality

Let the samples  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  be iid as  $\mathbf{x}$  from the pdf  $p(\mathbf{x}; \boldsymbol{\theta})$ . Denote  $\mathbf{y}_n = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$  to be the collection of these samples, so that the likelihood will be  $p(\mathbf{y}_n; \boldsymbol{\theta}) = \prod_{i=1}^n p(\mathbf{x}_i; \boldsymbol{\theta})$ , with the maximum likelihood of these samples denoted  $\hat{\boldsymbol{\theta}}(\mathbf{y}_n)$ .

**Theorem 2.9.** Assuming the regularity conditions stated earlier on the pdf  $p(\mathbf{x}; \boldsymbol{\theta})$ , the MLE of the parameter  $\boldsymbol{\theta}$  is asymptotically distributed according to

$$\sqrt{n} \left( \hat{\boldsymbol{\theta}}(\mathbf{y}_n) - \boldsymbol{\theta} \right) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{I}^{-1}(\boldsymbol{\theta}))$$

where  $\mathbf{I}(\boldsymbol{\theta})$  is derived from the pdf  $p(\mathbf{x}; \boldsymbol{\theta})$ , i.e., it is the Fisher information of a single observation or sample.

### 2.4.3 Scoring

There exists a number of approaches to finding the maximum likelihood, in some cases requiring iterative techniques. One such technique is Fisher's method of scoring. Given an observation  $\mathbf{x}$  from a likelihood or pdf  $p(\mathbf{x}; \boldsymbol{\theta})$  depending on an unknown parameter  $\boldsymbol{\theta}$  and given some initial estimate  $\dot{\boldsymbol{\theta}}^{(1)}$  of  $\boldsymbol{\theta}$ , then iteratively using the update

$$\dot{\boldsymbol{\theta}}^{(k+1)} = \dot{\boldsymbol{\theta}}^{(k)} + \mathbf{I}^{-1}(\dot{\boldsymbol{\theta}}^{(k)})\mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)}) \quad (2.8)$$

will find the MLE under certain conditions, e.g., provided the initial estimate is sufficiently close in a locally convex region.

## 2.5 Hypothesis testing

Given the inclusion of the CRB quantity in the asymptotic normality results in section 2.4.2, it is not surprising that there would also exist connections to some asymptotic hypothesis tests. Assume  $\mathbf{h} : \mathbb{R}^m \rightarrow \mathbb{R}^r$  is a consistent and nonredundant differentiable function, which defines the null hypothesis

$$H_0 : \mathbf{h}(\boldsymbol{\theta}) = \mathbf{0}$$

in the likelihood (or model)  $p(\mathbf{y}_n; \boldsymbol{\theta})$  versus the alternative hypothesis  $H_1 : \mathbf{h}(\boldsymbol{\theta}) \neq \mathbf{0}$ .

### 2.5.1 The Rao statistic

The Rao (or score) test statistic is given by

$$\rho(\mathbf{y}_n) = \frac{1}{n} \mathbf{s}^T(\mathbf{y}_n; \hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \mathbf{I}^{-1}(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \mathbf{s}(\mathbf{y}_n; \hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \quad (2.9)$$

where the Fisher score at  $\mathbf{y}_n$  is  $\mathbf{s}(\mathbf{y}_n, \hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) = \sum_{i=1}^n \mathbf{s}(\mathbf{x}_i, \hat{\boldsymbol{\theta}}_h(\mathbf{y}_n))$  and  $\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)$  is the MLE of  $\boldsymbol{\theta}$  under the null hypothesis  $\mathbf{h}(\boldsymbol{\theta}) = \mathbf{0}$ . A variant of this statistic, the Lagrange multiplier test [3, 63]

$$\hat{\boldsymbol{\lambda}}_n^T \mathbf{H}(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \mathbf{I}^{-1}(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \mathbf{H}^T(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \hat{\boldsymbol{\lambda}}_n,$$

was developed by Silvey [63]. The equivalence between the Rao and Lagrange multiplier test comes from the first order condition to satisfy the constraint [45], i.e.,

$$\mathbf{s}(\mathbf{y}_n; \hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) + \mathbf{H}^T(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \hat{\boldsymbol{\lambda}}_n = \mathbf{0},$$

$$\mathbf{h}(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) = \mathbf{0}$$

where  $\hat{\boldsymbol{\lambda}}_n \in \mathbb{R}^r$  is a vector of Lagrange multiplier estimates. First order Taylor-series expansions of both equations about the true parameter  $\boldsymbol{\theta}$  produces

$$\mathbf{s}(\mathbf{y}_n; \boldsymbol{\theta}) - \mathbf{I}_n(\boldsymbol{\theta})(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n) - \boldsymbol{\theta}) + \mathbf{H}^T(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \hat{\boldsymbol{\lambda}}_n = o(n^{-1/2})$$

$$\mathbf{h}(\boldsymbol{\theta}) + \mathbf{H}(\boldsymbol{\theta})(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n) - \boldsymbol{\theta}) + o(n^{-1/2}) = \mathbf{h}(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n))$$

where  $\mathbf{I}_n(\boldsymbol{\theta}) = n\mathbf{I}(\boldsymbol{\theta})$  ( $n$  times the Fisher information based on a single sample  $\mathbf{x}$ ).

Hence under the null hypothesis, the latter implies  $\mathbf{H}(\boldsymbol{\theta})(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n) - \boldsymbol{\theta}) = o(n^{-1/2})$ ,

and premultiplying the former by  $\mathbf{H}(\boldsymbol{\theta})\mathbf{I}_n^{-1}(\boldsymbol{\theta})$ , then

$$\mathbf{H}(\boldsymbol{\theta})\mathbf{I}_n^{-1}(\boldsymbol{\theta})\mathbf{s}(\mathbf{y}_n; \boldsymbol{\theta}) + \mathbf{H}(\boldsymbol{\theta})\mathbf{I}_n^{-1}(\boldsymbol{\theta})\mathbf{H}^T(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \hat{\boldsymbol{\lambda}}_n = o(n^{-1/2}).$$

Since  $\mathbf{s}(\mathbf{y}_n; \boldsymbol{\theta}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n(\boldsymbol{\theta}))$ , then applying Slutsky's theorem and the continuity of the Fisher information and the hypothesis function,

$$\mathbf{H}(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n))\mathbf{I}_n^{-1}(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n))\mathbf{H}^T(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \hat{\boldsymbol{\lambda}}_n \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{H}(\boldsymbol{\theta})\mathbf{I}_n^{-1}(\boldsymbol{\theta})\mathbf{H}^T(\boldsymbol{\theta})).$$

Therefore,  $\left(\mathbf{H}(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n))\mathbf{I}_n^{-1}(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n))\mathbf{H}^T(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n))\right)^{-1/2}\hat{\boldsymbol{\lambda}}_n$  is a  $r$ -dimensional standard normal variable in distribution, and

$$\rho(\mathbf{y}_n) \xrightarrow{d} \chi_r^2.$$

Hypothesis  $H_0$  is rejected if  $\rho(\mathbf{y}_n) > \chi_{r,\alpha}^2$ , where  $\chi_{r,\alpha}^2$  is the  $(1 - \alpha)$ -quantile of the chi-square distribution with  $r$  degrees of freedom.

## 2.5.2 The Wald statistic

The Wald test statistic is given by

$$\omega(\mathbf{y}_n) = n\mathbf{h}^T(\hat{\boldsymbol{\theta}}(\mathbf{y}_n)) \left( \mathbf{H}(\hat{\boldsymbol{\theta}}(\mathbf{y}_n))\mathbf{I}^{-1}(\hat{\boldsymbol{\theta}}(\mathbf{y}_n))\mathbf{H}^T(\hat{\boldsymbol{\theta}}(\mathbf{y}_n)) \right)^{-1} \mathbf{h}(\hat{\boldsymbol{\theta}}(\mathbf{y}_n)) \quad (2.10)$$

where  $\mathbf{H}(\boldsymbol{\theta}) = \left. \frac{\partial \mathbf{h}(\boldsymbol{\theta}')}{\partial \boldsymbol{\theta}'^T} \right|_{\boldsymbol{\theta}'=\boldsymbol{\theta}}$  and  $\hat{\boldsymbol{\theta}}(\mathbf{y}_n)$  is an MLE of  $\boldsymbol{\theta}$ . (The nonredundancy of  $\mathbf{h}$  implies that  $\mathbf{H}(\cdot)$  is full row rank.) From section 2.1.1, the CRB of  $\mathbf{h}(\boldsymbol{\theta})$  in the model  $p(\mathbf{x}; \boldsymbol{\theta})$  is  $\mathbf{H}(\boldsymbol{\theta})\mathbf{I}^{-1}(\boldsymbol{\theta})\mathbf{H}^T(\boldsymbol{\theta})$  and therefore, using theorem 2.9,

$$\sqrt{n} \left( \mathbf{h}(\hat{\boldsymbol{\theta}}(\mathbf{y}_n)) - \mathbf{h}(\boldsymbol{\theta}) \right) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{H}(\boldsymbol{\theta})\mathbf{I}^{-1}(\boldsymbol{\theta})\mathbf{H}^T(\boldsymbol{\theta})).$$

Hence, under  $H_0$ , using Slutsky's theorem, the convergence in probability of the MLE, and continuity of the FIM and Jacobian of the test function,

$$\omega(\mathbf{y}_n) \xrightarrow{d} \chi_r^2.$$

Therefore, the hypothesis  $H_0$  is rejected if  $\omega(\mathbf{y}_n) > \chi_{r,\alpha}^2$ .

## 2.6 Discussion

In this section, the Cramér-Rao bound (CRB) was defined and in theorem 2.1 it was stated to be a bound on mean-square error performance of an unbiased

estimator. The theory of the CRB's connection to a variety of useful theorems and equations in mathematical statistics was demonstrated. The CRB is included in conditions for local identifiability (theorem 2.5) as well as conditions for strict identifiability (theorem 2.6). The CRB is the projection matrix of the BLUE under a Gaussian model in equation (2.6). There is a connection between the existence of efficiency with respect to the CRB and the method of ML (theorem 2.8). The CRB is also the asymptotic variance of the ML estimator (theorem 2.9) and appears in the update formula in (2.8) for the method of scoring. The CRB also appears in the formulas for the Rao test statistic in (2.9) and for the Wald test statistic (2.10).

These connections are not exhaustive, e.g., the CRB formula can also be useful in defining confidence regions or in useful as a cost function, but these are perhaps the most prevalent general topics in mathematical statistics theory and for that reason serve as a useful comparison for the *constrained* Cramér-Rao bound in chapter 3 and its connections in the theory.

## Chapter 3

### THE CONSTRAINED CRAMÉR-RAO BOUND

While the Cramér-Rao bound (CRB) is a useful measure of parametric estimation, it does not inherently measure the performance of estimators of parameters that satisfy side information in the form of a functional equality constraint

$$\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}. \tag{3.1}$$

The statistical literature is surprisingly somewhat limited in addressing performance measures under this general scenario. The traditional practice is to find some equivalent reparameterization of the particular model and then find the CRB on the parameters of interest using the reparameterized transformation. This approach, however, does not lend itself to theoretical meaning beyond the particular reparameterized model. Typically, works that do examine (3.1) in a general manner are focused on developing methods for decisions (hypothesis testing) instead of measuring estimation performance. Nevertheless, these results have connections to a CRB incorporating the side information in (3.1), or a constrained CRB. A number of papers, including Aitchison and Silvey [3] and Crowder [18], using the method of Lagrangian multipliers, examine the asymptotic normality of the constrained maximum likelihood estimator (CMLE) and as a consequence unintentionally develop a CRB under equality constraints. Under certain conditions, the asymptotic variance of the MLE equaling the CRB lends credence to the claim that the asymptotic vari-

ance of the CMLE should equal to the CRB under equality constraints, although the authors did not always appear cognizant of this fact. Others, including Silvey [63] to develop his Lagrange multiplier test, Osborne [55] with linear constraints to develop a scoring algorithm, and Waldorp, Huizenga, and Grasman [73] to develop a Wald-type test, also use the Lagrange multiplier approach in developing a constrained bound. Again, since these authors were primarily focused on asymptotic properties or hypothesis testing, the nature and perhaps utility of this mathematical quantity in their work is not explicitly stated as a CRB or bound on performance estimation of parameters under constraints. The creation of a constrained bound strictly for the use in performance analysis wasn't achieved until Gorman and Hero [23]. Gorman and Hero derived such a measure by taking the limit of the Hammersley-Chapman-Robbins bound with test points restricted to exist only in the constraint space. This constrained Cramér-Rao bound (CCRB)

$$\mathbf{I}^{-1}(\boldsymbol{\theta}) - \mathbf{I}^{-1}(\boldsymbol{\theta}) \mathbf{F}^T(\boldsymbol{\theta}) (\mathbf{F}(\boldsymbol{\theta}) \mathbf{I}^{-1}(\boldsymbol{\theta}) \mathbf{F}^T(\boldsymbol{\theta}))^{-1} \mathbf{F}(\boldsymbol{\theta}) \mathbf{I}^{-1}(\boldsymbol{\theta}) \quad (3.2)$$

utilizes the Jacobian  $\mathbf{F}(\boldsymbol{\theta})$  of the functional constraint  $\mathbf{f}(\boldsymbol{\theta})$  and the inverse of the Fisher Information matrix (FIM)  $\mathbf{I}(\boldsymbol{\theta})$  (based on the unconstrained model), which must be nonsingular. As with the CRB, there exist a number of alternative derivations. The works of Crowder, Waldorp, et al, Gorman and Hero, and Aitchison [2] include the formula in (3.2) in some manner for which the CRB might be used for the unconstrained scenario in their works, a fact that implicitly proves the validity of the CCRB. With the explicit proof by Gorman and Hero as a guideline, Marzetta [47] provides an elementary proof of this CCRB, which avoids the use of the application of



the Cauchy-Schwarz inequality, and avoids the use of pseudo-inverses, by examining the inequality created from the positive-semidefiniteness property for the variance of a properly defined random variable. A similar construction was used by Stoica and Ng [68] to formulate a more general CCRB

$$\mathbf{U}(\boldsymbol{\theta}) (\mathbf{U}^T(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}))^{-1} \mathbf{U}^T(\boldsymbol{\theta}) \quad (3.3)$$

that incorporates the constraint information without the assumption of a nonsingular FIM. This CCRB utilizes the unconstrained FIM and an orthonormal complement matrix  $\mathbf{U}(\boldsymbol{\theta})$  whose vectors span the null space of the constraint Jacobian matrix. Furthermore, when the FIM is nonsingular, the Stoica-Ng CCRB in (3.3) is equivalent to the Gorman-Hero version of (3.2). Hence, while much of the previous work used the formula in (3.2), the more general formula (3.3) is also applicable. Osborne, independently from much of these other works, developed a method of scoring with constraints that utilizes the Stoica-Ng CCRB formula in (3.3) as the projection matrix in place of the CRB. There are, of course, numerous instances of matrix structures of the same form as (3.3), for example as part of the projection matrix of the generalized least squares estimate of the mean of a linear model.

In this section, we develop a very simple derivation of the CCRB in (3.3). Rather than assuming the parameters satisfy functional constraints, we approach the problem theoretically from an alternative, yet equivalent, perspective and assume the parameters locally fit a reduced parametric model. This approach permits the extension of the existing classical theory underlying the CRB to the case of a model under parametric constraints. While it is true that several of these exten-

sions already exist in the literature, there does not exist a cohesive treatment of these results. However, this chapter should not be viewed as simply a collection of historical results, but a unified and comprehensive development of the theory of the constrained Cramér-Rao bound.

### 3.1 The Constrained CRB

Suppose we observe  $\mathbf{x} \in \mathbb{X} \subset \mathbb{R}^n$  from a probability density function  $p(\mathbf{x}; \boldsymbol{\theta})$  where  $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^m$  is a vector of unknown deterministic parameters and, in addition, suppose these parameters are required to satisfy  $k$  *consistent* and *nonredundant* continuously differentiable parametric equality constraints, i.e.,  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$  for some consistent and nonredundant  $\mathbf{f} : \Theta \rightarrow \mathbb{R}^k$ . We shall denote

$$\Theta_f = \left\{ \boldsymbol{\theta}' \in \Theta : \mathbf{f}(\boldsymbol{\theta}') = \mathbf{0}, \mathbf{f} \text{ consistent, nonredundant} \right\} \quad (3.4)$$

to be the *feasible* set satisfying the constraints. Hence, the constraint can also be stated  $\boldsymbol{\theta} \in \Theta_f$ . The constraints being consistent means that the set  $\Theta_f$  is nonempty. The constraints being nonredundant means that the Jacobian  $\mathbf{F}(\boldsymbol{\theta}') = \frac{\mathbf{f}(\boldsymbol{\theta}')}{\partial \boldsymbol{\theta}'^T}$  has rank  $k$  whenever  $\mathbf{f}(\boldsymbol{\theta}') = \mathbf{0}$ .

As before, the Fisher information matrix (FIM) of this model (ignoring the constraint) is given by

$$\mathbf{I}(\boldsymbol{\theta}) = E_{\boldsymbol{\theta}} \{ \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}) \mathbf{s}^T(\mathbf{x}; \boldsymbol{\theta}) \}$$

where  $\mathbf{s}(\mathbf{x}; \boldsymbol{\theta})$  is the Fisher score defined by

$$\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}) = \left. \frac{\partial \log p(\mathbf{x}; \boldsymbol{\theta}')}{\partial \boldsymbol{\theta}'} \right|_{\boldsymbol{\theta}' = \boldsymbol{\theta}}$$

and the expectation is evaluated at  $\boldsymbol{\theta}$ , i.e.,  $E_{\boldsymbol{\theta}}(\cdot) = \int_{\mathbb{X}} (\cdot) p(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x}$ .

Incorporating this additional *side information* into the information from the observations  $\mathbf{x}$  directly would require an alteration of the pdf's dependence on the unknown parameter. Such an approach can often be analytically impractical or numerically complex. Hence, it is desirable to have a formulaic or prescriptive approach to include the side information, or constraints, indirectly. To meet this need, Stoica and Ng developed a method to incorporate parametric equality constraints into the CRB [68, theorem 1].

**Theorem 3.1** (Stoica & Ng). The constrained Cramér-Rao bound on  $\boldsymbol{\theta} \in \Theta_f$  is given by

$$\text{CCRB}(\boldsymbol{\theta}) = \mathbf{U}(\boldsymbol{\theta}) \left( \mathbf{U}^T(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}) \right)^{-1} \mathbf{U}^T(\boldsymbol{\theta}) \quad (3.5)$$

where  $\mathbf{U}(\boldsymbol{\theta})$  is a matrix whose column vectors form an orthonormal basis for the null space of the Jacobian  $\mathbf{F}(\boldsymbol{\theta})$ , i.e.,

$$\mathbf{F}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}) = \mathbf{0} \text{ , } \mathbf{U}^T(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}) = \mathbf{I}_{(m-k) \times (m-k)} . \quad (3.6)$$

Thus, if  $\mathbf{t}(\mathbf{x})$  is an unbiased estimator of  $\boldsymbol{\theta}$ , which satisfies the constraint (3.4), then

$$\text{Var}(\mathbf{t}(\mathbf{x})) \geq \text{CCRB}(\boldsymbol{\theta})$$

with equality if and only if  $\mathbf{t}(\mathbf{x}) - \boldsymbol{\theta} = \text{CCRB}(\boldsymbol{\theta}) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta})$  in the mean-square sense.<sup>1</sup>

---

<sup>1</sup>The original theorem requires the estimator to satisfy the constraint. In general, the parameter and its unbiased estimator will not simultaneously satisfy the constraint since the implication, mainly that  $\mathbf{f}(E_{\boldsymbol{\theta}} \mathbf{t}(\mathbf{x})) = E_{\boldsymbol{\theta}} \mathbf{f}(\mathbf{t}(\mathbf{x}))$ , is true only under particular conditions. However, the CCRB is the same if either assumption is made exclusively. In this treatise, I assume that the actual parameter  $\boldsymbol{\theta}$  satisfies the constraint and the unbiased estimator  $\mathbf{t}(\mathbf{x})$  does not. (In section 3.4, the constrained maximum likelihood estimator (CMLE) is assumed to satisfy the constraint, but unbiasedness is not assumed.)

The Jacobian  $\mathbf{F}(\boldsymbol{\theta})$  having full row rank is not necessary since the Jacobian does not explicitly appear in the CCRB formula in (3.5). Indeed, the requirement that the column vectors of  $\mathbf{U}(\boldsymbol{\theta})$  form an orthonormal basis is also unnecessary, only that they be linearly independent and that they span the basis of the null space of the row vectors of  $\mathbf{F}(\boldsymbol{\theta})$ , i.e., only that the column space of  $\mathbf{U}(\boldsymbol{\theta})$  is an *orthogonal complement* of the row space of  $\mathbf{F}(\boldsymbol{\theta})$ , is required since it is clear from the structure of (3.5) that the CCRB is invariant to automorphisms in  $\mathbb{R}^{m-\text{rank}(\mathbf{F}(\boldsymbol{\theta}))}$  on  $\mathbf{U}(\boldsymbol{\theta})$ . Regardless, for convenience, and except where otherwise noted, we will assume that the constraints are nonredundant and the columns of  $\mathbf{U}(\boldsymbol{\theta})$  are orthonormal to ensure that  $\text{rank}(\mathbf{U}(\boldsymbol{\theta})) = m - k$  and  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) = \mathbf{I}_{m-k}$ , respectively. The existence of the bound only requires that  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  rather than the FIM be nonsingular.<sup>2</sup> The original proof of this theorem given by Stoica and Ng considers the variance inequality generated by the random variable  $\mathbf{t}(\mathbf{x}) - \boldsymbol{\theta} - \mathbf{W}\mathbf{U}(\boldsymbol{\theta})\mathbf{U}^T(\boldsymbol{\theta})\mathbf{s}(\mathbf{x}; \boldsymbol{\theta})$  and maximizes  $\mathbf{W}$  to attain the tightest bound for  $\text{Var}(\mathbf{t}(\mathbf{x}) - \boldsymbol{\theta})$ .

**Example 3.2** (Unit Modulus Constraint). Suppose  $\vartheta$  in example 2.2 is constrained to be unit modulus, i.e.,  $f(\vartheta) = |\vartheta|^2 - 1$ . Then its gradient in terms of the parameter vector  $\boldsymbol{\theta} = \begin{bmatrix} \text{Re}(\vartheta) \\ \text{Im}(\vartheta) \end{bmatrix}$  is  $\mathbf{F}(\boldsymbol{\theta}) = [2\text{Re}(\vartheta), 2\text{Im}(\vartheta)]$ , which has an orthonormal complement  $\mathbf{U}(\boldsymbol{\theta}) = \begin{bmatrix} -\text{Im}(\vartheta) \\ \text{Re}(\vartheta) \end{bmatrix}$ . The CCRB for this constraint is then

$$\text{CCRB}(\boldsymbol{\theta}) = \frac{\sigma^2}{2} \begin{bmatrix} \text{Im}(\vartheta)^2 & -\text{Im}(\vartheta)\text{Re}(\vartheta) \\ -\text{Re}(\vartheta)\text{Im}(\vartheta) & \text{Re}(\vartheta)^2 \end{bmatrix}.$$

---

<sup>2</sup>A corresponding regularity condition to that mentioned in section 2.1 will be discussed in section 3.1.4.

### 3.1.1 A proof of the CCRB

While the proof given by Stoica and Ng is sufficient to establish the validity of the CCRB, these proofs ignore the existing classical theory encompassing the CRB and FIM, which is already sufficient to prove the CCRB. However, prior to developing the foundation for the CCRB from the existing theory, we require a foray into multivariable calculus and, specifically, the use of the implicit function theorem. The reward for this approach will be a seamless presentation of statistical inference involving the constrained Cramér-Rao bound.

From the perspective of multivariable calculus, the constraint  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$  effectively restricts  $\boldsymbol{\theta}$  to a manifold  $\Theta_f$  of the original vector space  $\Theta$ , with the manifold having dimension  $m - k$  since  $k$  degrees of freedom are lost when  $\text{rank}(\mathbf{F}(\boldsymbol{\theta})) = k$  for all  $\boldsymbol{\theta} \in \Theta_f$ . More formally, the following theorem [65, theorems 5-1 and 5-2] applies.

**Theorem 3.3** (Implicit Function Theorem). Let  $\mathbb{U} \subset \mathbb{R}^m$  be an open set and assume  $\mathbf{f} : \mathbb{U} \rightarrow \mathbb{R}^k$  is a differentiable function such that  $\mathbf{F}(\boldsymbol{\theta})$  has rank  $k$  whenever  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$ . Then  $\Theta_f \cap \mathbb{U}$  is an  $(m - k)$ -dimensional manifold in  $\mathbb{R}^m$ , and for every  $\boldsymbol{\theta} \in \Theta_f \cap \mathbb{U}$  there is an open set  $\mathbb{V} \ni \boldsymbol{\theta}$ , an open set  $\mathbb{W} \subset \mathbb{R}^{m-k}$ , and a 1-1 differentiable function  $\mathbf{g}_{\boldsymbol{\theta}} : \mathbb{W} \rightarrow \mathbb{R}^m$  such that

(a)  $\mathbf{g}_{\boldsymbol{\theta}}(\mathbb{W}) = \Theta_f \cap \mathbb{U} \cap \mathbb{V}$ , and

(b) the Jacobian of  $\mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi}')$  has rank  $m - k$  for each  $\boldsymbol{\xi}' \in \mathbb{W}$ .

Therefore, there exists a function  $\mathbf{g}_{\boldsymbol{\theta}} : \mathbb{R}^{m-k} \rightarrow \mathbb{R}^m$ , and sets  $\mathbb{O} \ni \boldsymbol{\theta}$  and  $\mathbb{P}$  open in  $\Theta_f$  and  $\mathbb{R}^{m-k}$ , respectively, such that  $\mathbf{g}_{\boldsymbol{\theta}} : \mathbb{P} \rightarrow \mathbb{O}$  is a *diffeomorphism* on  $\mathbb{P}$ ,

i.e., a continuously differentiable bijection with a continuously differentiable inverse. A geometric example is shown in Figure 3.1. Note this diffeomorphism depends on the parameter  $\boldsymbol{\theta}$  as the reparameterization is only guaranteed to exist in a local neighborhood of  $\boldsymbol{\theta}$ ; however, for convenience, we will omit this notation in this subsection so that  $\mathbf{g} = \mathbf{g}_{\boldsymbol{\theta}}$ . Thus, we may proceed under the assumption that every  $\boldsymbol{\theta}' \in \mathbb{O} \subset \Theta_f$  is the image of a unique reduced parameter vector  $\boldsymbol{\xi}' \in \mathbb{P} \subset \mathbb{R}^{m-k}$  under  $\mathbf{g}$ , or simply

$$\boldsymbol{\theta}' = \mathbf{g}(\boldsymbol{\xi}'). \quad (3.7)$$

Necessarily, there exists some unique  $\boldsymbol{\xi} \in \mathbb{P}$  for which  $\boldsymbol{\theta} = \mathbf{g}(\boldsymbol{\xi})$ . We will denote the Jacobian of  $\mathbf{g}$  to be  $\mathbf{G}(\boldsymbol{\xi}') = \frac{\partial \mathbf{g}(\boldsymbol{\xi}')}{\partial \boldsymbol{\xi}'^T}$ , which also implicitly depends on  $\boldsymbol{\theta}$ .

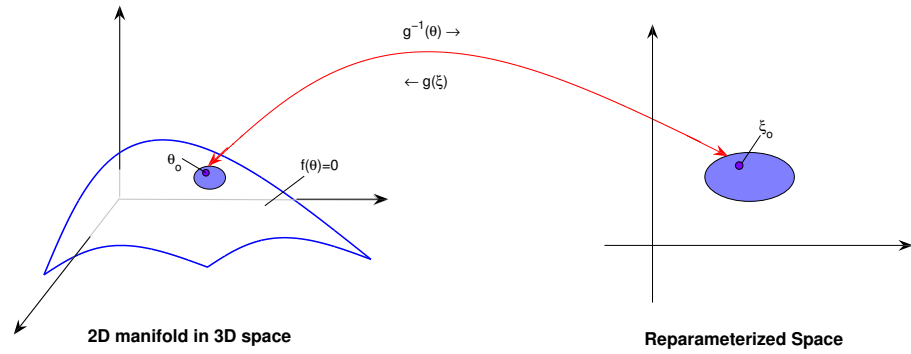


Figure 3.1: Reparameterization of  $\mathbf{f}(\boldsymbol{\theta}) = 0$  to  $\boldsymbol{\theta} = \mathbf{g}(\boldsymbol{\xi})$ .

**Example 3.4** (Unit Modulus Constraint). As an example of this principle, consider a complex parameter  $\vartheta$  with a modulus constraint (as in example 3.2). The parameter vector in this case may be  $\boldsymbol{\theta} = [\theta_1, \theta_2]^T = [\text{Re}(\vartheta), \text{Im}(\vartheta)]^T \in \mathbb{R}^2$  with the constraint being  $f(\boldsymbol{\theta}) = \theta_1^2 + \theta_2^2 - 1 = 0$ . By the implicit function theorem,

constraining  $\vartheta$  to be unit modulus is tantamount to assuming the existence of a  $\xi \in \mathbb{R}$  such that  $\vartheta = e^{-j\xi}$ . Hence,  $\boldsymbol{\theta}$  is a function of  $\xi$ , i.e.,

$$\boldsymbol{\theta} = \begin{bmatrix} \text{Re}(\vartheta) \\ \text{Im}(\vartheta) \end{bmatrix} = \begin{bmatrix} \cos(\xi) \\ -\sin(\xi) \end{bmatrix} = \mathbf{g}(\xi). \quad (3.8)$$

Also,  $\mathbf{g}(\mathbb{R}) = \{\boldsymbol{\theta} \in \mathbb{R}^2 : f(\boldsymbol{\theta}) = 0\}$  and  $\mathbf{G}(\xi) = [-\sin(\xi) \ , \ -\cos(\xi)]^T$  has rank 1 for every  $\xi$ . For the model in example 3.2, then

$$\text{CCRB}(\boldsymbol{\theta}) = \frac{\sigma^2}{2} \begin{bmatrix} \sin^2(\xi) & -\sin(\xi)\cos(\xi) \\ -\cos(\xi)\sin(\xi) & \cos^2(\xi) \end{bmatrix},$$

which is exactly as before.

It must be noted that  $\mathbf{g}$  will not be unique. For the previous example,  $\vartheta = e^{j\xi}$  is another possible reparameterization. Nor is any  $\mathbf{g}$  satisfying the theorem necessarily a 1-1 correspondence between  $\mathbb{R}^{m-k}$  and  $\Theta_f$ ; again, for the previous example,  $\mathbf{g}$  is periodic. Thus, the bijection is only guaranteed locally. Finding a particular  $\mathbf{g}$  for a given  $\mathbf{f}$  and  $\boldsymbol{\theta}$  may not be obvious. Methods for approximating an implicit function will be discussed in section 3.1.5. Regardless, as shall be shown in the context of the CCRB, knowledge of any particular  $\mathbf{g}$  is unnecessary; only its existence, given by the implicit function theorem, is necessary. Why? Using the implicit function theorem to assume a locally equivalent reparameterization for the constraint limits the information from the observations to the density function's local dependence on the unknown parameter. But as the CRB (and hence CCRB) is a local bound that only characterizes the local noise ambiguities in the model, i.e., the average local curvature of the density at the parameter value of interest, this limitation is invariant to the local curvature restricted to a space determined by the constraints.

**Theorem 3.5** ([50]). The CRB on  $\boldsymbol{\theta} \in \mathbb{O}$  under the assumption of (3.7) is given by

$$\mathbf{G}(\boldsymbol{\xi}) \left( \mathbf{G}^T(\boldsymbol{\xi}) \mathbf{I}(\mathbf{g}(\boldsymbol{\xi})) \mathbf{G}(\boldsymbol{\xi}) \right)^{-1} \mathbf{G}^T(\boldsymbol{\xi}) \quad (3.9)$$

and if  $\mathbf{g}$  in (3.7) is an implicit function of  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$  then this bound is equivalent to the constrained Cramér-Rao bound in (3.5).

*Proof.* This CRB can be developed from a transformation of parameters on the FIM and the CRB. From the CRB transformation of parameters on the CRB (see section 2.1.1), if  $\mathbf{t}(\mathbf{x})$  is an unbiased estimator<sup>3</sup> of  $\boldsymbol{\theta} = \mathbf{g}(\boldsymbol{\xi})$  then

$$\text{Var}(\mathbf{t}(\mathbf{x})) \geq \text{CRB}(\mathbf{g}(\boldsymbol{\xi})) = \mathbf{G}(\boldsymbol{\xi}) \tilde{\mathbf{I}}^{-1}(\boldsymbol{\xi}) \mathbf{G}^T(\boldsymbol{\xi}) \quad (3.10)$$

with equality if and only if  $\mathbf{t}(\mathbf{x}) - \mathbf{g}(\boldsymbol{\xi}) = \mathbf{G}(\boldsymbol{\xi}) \tilde{\mathbf{I}}^{-1}(\boldsymbol{\xi}) \tilde{\mathbf{s}}(\mathbf{x}; \boldsymbol{\xi})$  in the mean-square sense [36, Appendix 3B], where  $\tilde{\mathbf{I}}(\boldsymbol{\xi}) = E_{\boldsymbol{\xi}} \tilde{\mathbf{s}}(\mathbf{x}; \boldsymbol{\xi}) \tilde{\mathbf{s}}^T(\mathbf{x}; \boldsymbol{\xi})$  is the FIM on  $\boldsymbol{\xi}$  and  $\tilde{\mathbf{s}}(\mathbf{x}; \boldsymbol{\xi}) = \left. \frac{\partial \log q(\mathbf{x}; \boldsymbol{\xi}')}{\partial \boldsymbol{\xi}'} \right|_{\boldsymbol{\xi}' = \boldsymbol{\xi}}$  is the Fisher score of the pdf with respect to  $\boldsymbol{\xi}$ , this pdf being  $q(\mathbf{x}; \boldsymbol{\xi}) = p(\mathbf{x}; \mathbf{g}(\boldsymbol{\xi}))$ . By application of the derivative chain rule the Fisher score of  $\mathbf{x}$  with respect to  $\boldsymbol{\xi}$  is  $\tilde{\mathbf{s}}(\mathbf{x}; \boldsymbol{\xi}) = \mathbf{G}^T(\boldsymbol{\xi}) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta})$ . Hence, from the transformation of parameters on the FIM [12, p. 157], the FIM on  $\boldsymbol{\xi}$  is<sup>4</sup>

$$\begin{aligned} \tilde{\mathbf{I}}(\boldsymbol{\xi}) &= E \{ \tilde{\mathbf{s}}(\mathbf{x}; \boldsymbol{\xi}) \tilde{\mathbf{s}}^T(\mathbf{x}; \boldsymbol{\xi}) \} \\ &= E \{ \mathbf{G}^T(\boldsymbol{\xi}) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}) \mathbf{s}^T(\mathbf{x}; \boldsymbol{\theta}) \mathbf{G}(\boldsymbol{\xi}) \} \\ &= \mathbf{G}^T(\boldsymbol{\xi}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{G}(\boldsymbol{\xi}) \\ &= \mathbf{G}^T(\boldsymbol{\xi}) \mathbf{I}(\mathbf{g}(\boldsymbol{\xi})) \mathbf{G}(\boldsymbol{\xi}). \end{aligned} \quad (3.11)$$

---

<sup>3</sup>Again, there is no actual use of the assumption here that  $\mathbf{t}(\mathbf{x}) \in \Theta_f$  although if  $\mathbf{t}(\mathbf{x}) \in g(\mathcal{P})$  then  $\mathbf{t}(\mathbf{x})$  does indeed satisfy the constraint. The theorem result for unbiased estimators holds regardless and this will depend on the regularity condition detailed in section 3.1.4.

<sup>4</sup>Implicitly, it is assumed that the regularity conditions of section 2.1 hold with respect to  $\boldsymbol{\xi}$ . For how this applies to  $\boldsymbol{\theta}$ , see section 3.1.4.



Substituting (3.11) into (3.10) establishes the CRB under the assumption  $\boldsymbol{\theta} = \mathbf{g}(\boldsymbol{\xi})$ .

To establish the equivalence with the CCRB, note since  $\mathbf{f} \circ \mathbf{g}(\boldsymbol{\xi}) = \mathbf{0}$ , then by taking the Jacobian with respect to  $\boldsymbol{\xi}$  we have

$$\mathbf{0} = \frac{\partial}{\partial \boldsymbol{\xi}^T} \mathbf{0} = \frac{\partial}{\partial \boldsymbol{\xi}} \mathbf{f}(\mathbf{g}(\boldsymbol{\xi})) = \mathbf{F}(\mathbf{g}(\boldsymbol{\xi})) \mathbf{G}(\boldsymbol{\xi}) = \mathbf{F}(\boldsymbol{\theta}) \mathbf{G}(\mathbf{g}^{-1}(\boldsymbol{\theta})).$$

Hence, the columns of  $\mathbf{G}(\boldsymbol{\xi})$  reside in the null space of the row vectors of  $\mathbf{F}(\boldsymbol{\theta})$ .

And since  $\mathbf{g}$  has an inverse locally at  $\boldsymbol{\xi} = \mathbf{g}^{-1}(\boldsymbol{\theta})$ , then  $\mathbf{G}(\boldsymbol{\xi})$  has full column rank  $m - k$  locally about  $\boldsymbol{\xi}$ . (This is true on the whole set  $\mathbb{P} \subset \mathbb{R}^{m-k}$ .) Therefore,  $\text{span}(\mathbf{G}(\mathbf{g}^{-1}(\boldsymbol{\theta}))) = \text{span}(\mathbf{U}(\boldsymbol{\theta}))$  and there exists an  $m - k \times m - k$  full rank transformation matrix  $\mathbf{S}(\boldsymbol{\theta})$  such that  $\mathbf{G}(\mathbf{g}^{-1}(\boldsymbol{\theta})) \mathbf{S}(\boldsymbol{\theta}) = \mathbf{U}(\boldsymbol{\theta})$ . (This is true on the whole set  $\mathbb{O} \subset \Theta_f$ .) This matrix  $\mathbf{S}(\boldsymbol{\theta})$  is merely an orthonormalizing change of basis on the columns of  $\mathbf{G}(\mathbf{g}^{-1}(\boldsymbol{\theta}))$ . Therefore,

$$\begin{aligned} & \mathbf{G}(\boldsymbol{\xi}) (\mathbf{G}^T(\boldsymbol{\xi}) \mathbf{I}(\mathbf{g}(\boldsymbol{\xi})) \mathbf{G}(\boldsymbol{\xi}))^{-1} \mathbf{G}^T(\boldsymbol{\xi}) \\ &= \mathbf{G}(\mathbf{g}^{-1}(\boldsymbol{\theta})) (\mathbf{S}^{-T}(\boldsymbol{\theta}) \mathbf{U}^T(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}) \mathbf{S}^{-1}(\boldsymbol{\theta}))^{-1} \mathbf{G}^T(\mathbf{g}^{-1}(\boldsymbol{\theta})) \\ &= \mathbf{G}(\mathbf{g}^{-1}(\boldsymbol{\theta})) \mathbf{S}(\boldsymbol{\theta}) (\mathbf{U}^T(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}))^{-1} \mathbf{S}^T(\boldsymbol{\theta}) \mathbf{G}(\mathbf{g}^{-1}(\boldsymbol{\theta})) \\ &= \mathbf{U}(\boldsymbol{\theta}) (\mathbf{U}^T(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}))^{-1} \mathbf{U}^T(\boldsymbol{\theta}). \end{aligned}$$

Also,  $\text{Var}(\mathbf{t}(\boldsymbol{\xi})) \geq \text{CCRB}(\boldsymbol{\theta})$  with equality if and only if

$$\begin{aligned} \mathbf{t}(\mathbf{x}) - \mathbf{g}(\boldsymbol{\xi}) &= \mathbf{G}(\boldsymbol{\xi}) \tilde{\mathbf{I}}^{-1}(\boldsymbol{\xi}) \tilde{\mathbf{s}}(\mathbf{x}; \boldsymbol{\xi}) \\ \mathbf{t}(\mathbf{x}) - \boldsymbol{\theta} &= \mathbf{G}(\boldsymbol{\xi}) \tilde{\mathbf{I}}^{-1}(\boldsymbol{\xi}) \mathbf{G}^T(\boldsymbol{\xi}) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}) \\ &= \text{CCRB}(\boldsymbol{\theta}) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}). \end{aligned}$$

□

This proof of the CCRB uses the implicit function theorem, the CRB transformation formula, the FIM transformation formula, as well as well-known properties of rank and the derivative chain rule. An advantage of establishing the CCRB from these classical results will become clear as we establish the connection throughout statistical information theory. An example of this is immediately evident in an alternative proof of a proposition of Stoica and Ng [68, proposition 1].

**Corollary 3.6** (Stoica & Ng). Given the regularity conditions on  $\boldsymbol{\xi}$ , a necessary and sufficient condition for the existence of a finite CCRB of  $\boldsymbol{\theta}$  is

$$|\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})| \neq 0,$$

i.e.,  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  is nonsingular.

*Proof.* From the prior theorem, it is clear that  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  is nonsingular if and only if  $\mathbf{G}^T(\boldsymbol{\xi})\mathbf{I}(\mathbf{g}(\boldsymbol{\xi}))\mathbf{G}(\boldsymbol{\xi})$  is nonsingular if and only if  $\tilde{\mathbf{I}}(\boldsymbol{\xi})$  is nonsingular. Since a necessary and sufficient condition for the existence of a finite CRB of  $\mathbf{g}(\boldsymbol{\xi})$  is that  $\tilde{\mathbf{I}}(\boldsymbol{\xi})$  is nonsingular, the corollary is proven.  $\square$

Thus, with the usual regularity conditions (see section 2.1) being maintained for  $\boldsymbol{\xi}$ , the conditions for the existence of the CCRB with respect to  $\boldsymbol{\theta}$  are equivalent to the conditions for the existence of the CRB in the reduced parameter space of  $\boldsymbol{\xi}$ .

Moreover, the matrix  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  is nonsingular if and only if  $\mathbf{U}(\boldsymbol{\theta})$  has full column rank  $m-k$  and no linear combination of the columns of  $\mathbf{U}(\boldsymbol{\theta})$  reside in the null space of  $\mathbf{I}(\boldsymbol{\theta})$ . The first condition is satisfied always by definition. The second condition implies, for example, that if  $\mathbf{L}(\boldsymbol{\theta})\mathbf{L}^T(\boldsymbol{\theta})$  is the Cholesky decomposition [20,

p.194] of the FIM, where  $\mathbf{L}(\boldsymbol{\theta}) \in \mathbb{R}^{m \times \text{rank}(\mathbf{I}(\boldsymbol{\theta}))}$  is a full column rank lower triangular matrix with strictly positive values on the diagonal, then  $\mathbf{U}(\boldsymbol{\theta}) = \mathbf{L}(\boldsymbol{\theta})\mathbf{A}(\boldsymbol{\theta}) + \mathbf{K}(\boldsymbol{\theta})\mathbf{B}(\boldsymbol{\theta})$  where  $\mathbf{K}(\boldsymbol{\theta})$  is an orthogonal complement of  $\mathbf{L}(\boldsymbol{\theta})$  and for some full column rank  $\mathbf{A}(\boldsymbol{\theta}) \in \mathbb{R}^{\text{rank}(\mathbf{I}(\boldsymbol{\theta})) \times m-k}$ .

Further properties of the CCRB will be discussed in section 3.1.4.

### 3.1.2 Alternative formulas

The formula Stoica and Ng used to express the constrained Cramér-Rao bound is a generalization of an expression developed earlier first by Gorman and Hero [23, theorem 1] and later by Marzetta [47, theorem 2]. This is the CCRB formula in (3.12). Although Gorman and Hero's formula requires a nonsingular Fisher information, the version developed by Stoica and Ng appears to be inspired by a result [23, (19) in lemma 2] in Gorman and Hero's paper that unnecessarily assumes a positive definite FIM. However, this result was, in essence, not unknown in the literature. Gorman and Hero were aware of the prior work of Aitchison and Silvey [3, theorem 2 and  $\mathbf{P}$  on p.823], which is concerned with the asymptotic variance of the maximum likelihood estimator subject to restraints. But they were perhaps unaware (by lack of citation) of a later paper on hypothesis tests associated with the maximum likelihood, in which Aitchison and Silvey suggest a solution to the problem of singular information matrices [4, section 6]. This is the CCRB formula in (3.14). This version of the CCRB was also proven by Crowder [18, theorem 3]. Concurrent to the Gorman and Hero effort, Hendriks [27] and later Oller and Corcuera [53] developed

an extension of the Cramér-Rao bound intrinsic to the manifold using Riemannian geometry. More recently, Xavier and Barroso [74, 75] specified the lower bound on the geodesic of estimators to the true parameter. The original and latest versions of their bound are expressed in (3.16) and (3.17), respectively.

While these CCRB expressions are not the focus of the current treatise, they are still important for possible insights into the constrained Cramér-Rao bound. In this section, aspects of these insights will be briefly discussed as well as conditions for equality with the CCRB expression in (3.5).

### 3.1.2.1 Gorman-Hero-Aitchison-Silvey CCRB

Aitchison and Silvey used the method of Lagrange multipliers to show that the weighted asymptotic variance of the constrained maximum likelihood estimator, which should be the CCRB (implicitly), tends to

$$\text{CCRB}_2(\boldsymbol{\theta}) = \mathbf{I}^{-1}(\boldsymbol{\theta}) - \mathbf{I}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}) \left( \mathbf{F}(\boldsymbol{\theta})\mathbf{I}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}) \right)^{-1} \mathbf{F}(\boldsymbol{\theta})\mathbf{I}^{-1}(\boldsymbol{\theta}) \quad (3.12)$$

when the Fisher information is nonsingular. Alternatively, Gorman and Hero developed this same CCRB by restricting *test points* in the Chapman-Robbins bound (a Barankin-type bound) to be in the constraint space  $\Theta_f$  and then finding the derivatives as the limit of the finite difference expressions in the Chapman-Robbins bound.<sup>5</sup> A simpler proof was provided by Marzetta by considering the positive semidefiniteness of a properly chosen random variable.

A particular advantage to this presentation of the CCRB is the explicit quan-

---

<sup>5</sup>A definition of the Chapman-Robbins bound as well as a variant of the Gorman and Hero proof, which allows for a singular FIM, can be found in Appendix A.1.

tification of the gain in performance potential. Imposing a constraint  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$  on a set of parameters improves (lowers) the unconstrained bound from  $\mathbf{I}^{-1}(\boldsymbol{\theta})$  by exactly  $\mathbf{I}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}) (\mathbf{F}(\boldsymbol{\theta})\mathbf{I}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}))^{-1} \mathbf{F}(\boldsymbol{\theta})\mathbf{I}^{-1}(\boldsymbol{\theta})$ . Since the CRB of  $\mathbf{f}(\boldsymbol{\theta})$  is  $\mathbf{F}(\boldsymbol{\theta})\mathbf{I}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta})$  then its inverse is the Fisher information of  $\mathbf{f}(\boldsymbol{\theta})$  or  $\mathbf{0}$ . And  $\mathbf{F}^T(\boldsymbol{\theta}) (\mathbf{F}(\boldsymbol{\theta})\mathbf{I}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}))^{-1} \mathbf{F}(\boldsymbol{\theta})$  is the Fisher information of  $\boldsymbol{\theta}$  generated from the constraint  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$ . A disadvantage is the requirement of a nonsingular FIM. There exist numerous scenarios that require constraints for the original model to be identifiable (see Chapter 4). Additionally, this CCRB formula requires nonredundant constraints, i.e., the Jacobian  $\mathbf{F}(\boldsymbol{\theta})$  must be full row rank.

Similarities include the computational complexity of both formulas, which is  $O(m^3)$ . And when the Fisher information is nonsingular (and the constraints nonredundant), both formulas are equivalent.

**Theorem 3.7.** When the Fisher information is nonsingular and the constraints nonredundant, then an equivalent formula for the CCRB in (3.5) is  $\text{CCRB}_2(\boldsymbol{\theta})$ .

*Proof.* This is a different proof than the one provided in [68, corollary 1].<sup>6</sup> The existence of the Gorman-Hero-Marzetta formula assumes that the FIM  $\mathbf{I}(\boldsymbol{\theta})$  and the CRB of the constraint  $\mathbf{F}(\boldsymbol{\theta})\mathbf{I}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta})$  are regular (non-singular) [23, 47].

Correspondingly, the existence of the Stoica and Ng CCRB formula assumes that

---

<sup>6</sup>In addition to this alternate proof, they reference an algebraic identity from [37] that is useful in establishing the result.

**Lemma 3.8** (Khatri). Suppose  $\mathbf{A}$  is  $p \times q$  and  $\mathbf{B}$  is  $p \times p - q$  have ranks  $q$  and  $p - q$  respectively such that  $\mathbf{B}^T \mathbf{A} = \mathbf{0}$ . Then for any symmetric positive definite matrix  $\mathbf{S}$ ,

$$\mathbf{S}^{-1} - \mathbf{S}^{-1} \mathbf{A} (\mathbf{A}^T \mathbf{S}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{S}^{-1} = \mathbf{B} (\mathbf{B}^T \mathbf{S} \mathbf{B})^{-1} \mathbf{B}^T.$$

Substituting  $\mathbf{I}(\boldsymbol{\theta})$  for  $\mathbf{S}$ ,  $\mathbf{F}(\boldsymbol{\theta})$  for  $\mathbf{A}^T$ , and  $\mathbf{U}(\boldsymbol{\theta})$  for  $\mathbf{B}$  shows the equivalence between the two CCRBs for when the FIM is nonsingular.

$\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  is non-singular. Now, right-multiplying both formulas by  $\mathbf{F}^T(\boldsymbol{\theta})$  returns the results

$$\text{CCRB}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}) = \mathbf{U}(\boldsymbol{\theta}) (\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}))^{-1} \mathbf{U}^T(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}) = \mathbf{0},$$

$$\text{CCRB}_2(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}) = \mathbf{I}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}) - \mathbf{I}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}) = \mathbf{0}$$

since  $\mathbf{F}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) = \mathbf{0}$  in the first equation and the elimination of the inverse in  $\text{CRB}(\mathbf{f}(\boldsymbol{\theta}))$  in the second. Alternatively, right-multiplying both formulas by the matrix  $\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  returns the results

$$\text{CCRB}(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) = \mathbf{U}(\boldsymbol{\theta}),$$

$$\text{CCRB}_2(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) = \mathbf{U}(\boldsymbol{\theta})$$

again by eliminating the inverse in the first equation and since  $\mathbf{F}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) = \mathbf{0}$  in the second. Hence we have the equality

$$\text{CCRB}(\boldsymbol{\theta}) [\mathbf{F}^T(\boldsymbol{\theta}) \quad \mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})] = \text{CCRB}_2(\boldsymbol{\theta}) [\mathbf{F}^T(\boldsymbol{\theta}) \quad \mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})]$$

and if it can be shown that the matrix  $[\mathbf{F}^T(\boldsymbol{\theta}), \mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})]$  is regular, then the two CCRB formulas are shown to be equivalent. Suppose there exists vectors  $\boldsymbol{\alpha} \in \mathbb{R}^k$  and  $\boldsymbol{\beta} \in \mathbb{R}^{m-k}$  such that

$$\mathbf{F}^T(\boldsymbol{\theta})\boldsymbol{\alpha} + \mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})\boldsymbol{\beta} = \mathbf{0}. \tag{3.13}$$

Pre-multiplying (3.13) by  $\mathbf{U}^T(\boldsymbol{\theta})$  implies that

$$\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})\boldsymbol{\beta} = \mathbf{0}.$$

Since  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  is regular then  $\boldsymbol{\beta} = \mathbf{0}$ . Likewise, premultiplying (3.13) by  $\mathbf{F}(\boldsymbol{\theta})\mathbf{I}^{-1}(\boldsymbol{\theta})$  implies that

$$\mathbf{F}(\boldsymbol{\theta})\mathbf{I}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta})\boldsymbol{\alpha} = \mathbf{0}.$$

Since  $\mathbf{F}(\boldsymbol{\theta})\mathbf{I}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta})$  is regular then  $\boldsymbol{\alpha} = \mathbf{0}$ . Hence  $\mathbf{F}^T(\boldsymbol{\theta})\boldsymbol{\alpha} + \mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})\boldsymbol{\beta} = \mathbf{0}$  implies  $\boldsymbol{\alpha} = \mathbf{0}$  and  $\boldsymbol{\beta} = \mathbf{0}$ , which proves that  $[\mathbf{F}^T(\boldsymbol{\theta}) \quad \mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})]$  is full rank.  $\square$

### 3.1.2.2 Aitchison-Silvey-Crowder CCRB

The solution for resolving invertibility of singular FIMs in the variance and test results of Aitchison and Silvey [4] was to load the Fisher information with a matrix of the form  $\mathbf{F}^T(\boldsymbol{\theta})\mathbf{F}(\boldsymbol{\theta})$ . This was made more rigorous by Crowder [18] by replacing the Fisher information  $\mathbf{I}(\boldsymbol{\theta})$  with a *loaded FIM*  $\mathbf{D}(\boldsymbol{\theta}) = \mathbf{I}(\boldsymbol{\theta}) + \mathbf{F}^T(\boldsymbol{\theta})\mathbf{K}\mathbf{F}(\boldsymbol{\theta})$ , where  $\mathbf{K}$  is any positive semidefinite matrix such that  $\mathbf{D}(\boldsymbol{\theta})$  is regular. Hence, a generalization of (3.12) is

$$\text{CCRB}_3(\boldsymbol{\theta}) = \mathbf{D}^{-1}(\boldsymbol{\theta}) - \mathbf{D}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}) (\mathbf{F}(\boldsymbol{\theta})\mathbf{D}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}))^{-1} \mathbf{F}(\boldsymbol{\theta})\mathbf{D}^{-1}(\boldsymbol{\theta}). \quad (3.14)$$

This extension permits a singular FIM. This formulation is also independent of the choice of  $\mathbf{K}$ .

**Theorem 3.9.** An equivalent formula for the CCRB in (3.5) is  $\text{CCRB}_3(\boldsymbol{\theta})$ .

*Proof.* Replace  $\mathbf{I}(\boldsymbol{\theta})$  with  $\mathbf{D}(\boldsymbol{\theta})$  in the proof of theorem 3.7.<sup>7</sup>  $\square$

---

<sup>7</sup>In appendix A.3, this Crowder formula for the asymptotic variance of the constrained maximum likelihood estimate is shown to be equivalent to the CCRB using lemma 3.8 (also see section 3.4.2).

As the computational complexity of the two formulas are  $O(m^3)$ , there is no direct advantage of one CCRB version over the other. Certainly, in the nonsingular FIM case,  $\text{CCRB}_3(\boldsymbol{\theta})$  is the same as the  $\text{CCRB}_2(\boldsymbol{\theta})$  (with  $\mathbf{K} = \mathbf{0}$ ). Unfortunately, this CCRB formula does not appear to have simple connections to other areas of statistical inference in an inherent manner. Nevertheless, this possible inefficacy in theoretical applications does not effect its practical use.

### 3.1.2.3 Xavier’s & Barroso’s Intrinsic Variance Lower Bound

The prior CCRB metrics were in Euclidean  $\mathbb{R}^m$  space, i.e., the lower bound is on the measurement of the distance (in some direction or dimension) of the estimator to the true value of the parameter measured by “cutting through” the manifold. In some scenarios, it may be of interest to know what the bound is on the measurement of the distance “over the surface” of the manifold. Since dimensional directions can be somewhat ambiguous depending on the manifold, of particular interest is the geodesic, or shortest distance.

For this scenario, Xavier and Barroso [74, 75] formulated an inequality, the intrinsic variance lower bound (IVLB), for the variance of the geodesic to an unbiased estimator intrinsic to the manifold. Their results are derived from those of Hendriks [27] and Oller and Corcuera [53]. Ignoring elements of Riemannian geometric theory that are beyond the scope of this presentation, their IVLB result essentially relies on the inequality

$$\sqrt{C}\sqrt{\text{var}(\vartheta)}\cot(\sqrt{C}\sqrt{\text{var}(\vartheta)}) \leq \frac{\sqrt{\text{var}(\vartheta)}}{\sqrt{\lambda_{\boldsymbol{\theta}}}}, \quad (3.15)$$



where  $C$  is an upper bound on the sectional curvature of the manifold,  $\lambda_{\boldsymbol{\theta}}$  is a bound on variance of the Euclidean estimator error, and  $\text{var}(\vartheta)$  is the variance of the geodesic. Precisely,  $C = \max_{\boldsymbol{\theta} \in \Theta_f} K_{\boldsymbol{\theta}}$  where

$$K_{\boldsymbol{\theta}} = \max_{\substack{\mathbf{v}_1, \mathbf{v}_2 \text{ orthonormal} \\ \mathbf{F}(\boldsymbol{\theta})\mathbf{v}_i = \mathbf{0}}} \langle \Pi(\mathbf{v}_1, \mathbf{v}_1), \Pi(\mathbf{v}_2, \mathbf{v}_2) \rangle - \langle \Pi(\mathbf{v}_1, \mathbf{v}_2), \Pi(\mathbf{v}_1, \mathbf{v}_2) \rangle$$

and  $\Pi(\cdot, \cdot)$  is the second fundamental form [34] of  $\Theta_f$  defined by

$$\Pi(\mathbf{a}, \mathbf{b}) = -\mathbf{F}^T(\boldsymbol{\theta}) (\mathbf{F}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}))^{-1} \left[ \mathbf{a}^T \frac{\partial \mathbf{F}_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \mathbf{b} \right]_{k \times 1}$$

on  $\mathcal{U} \times \mathcal{U}$  where  $\mathcal{U} = \text{span}\{\mathbf{U}(\boldsymbol{\theta})\}$ .

Xavier and Barroso use a polynomial bound on the cotangent to solve the lower bound of a quadratic in terms of  $\text{var}(\vartheta)$ . In an earlier variant of the IVLB [74], Xavier and Barroso chose  $\lambda_{\boldsymbol{\theta}}^{-1} = \max_{\mathbf{v} \in \mathcal{U}, \|\mathbf{v}\|=1} \mathbf{v}^T \mathbf{I}(\boldsymbol{\theta}) \mathbf{v}$  and bounded  $t \cot(t) \geq 1 - \frac{2}{3}t^2$ , for  $0 \leq t \leq T \equiv 1.35$ , in (3.15) where  $t = \sqrt{C} \sqrt{\text{var}(\vartheta)}$ , to bound on the variance of the estimator's geodesic to the mean by

$$\text{var}(\vartheta) \geq \frac{4C + 3\lambda_{\boldsymbol{\theta}} - \sqrt{\lambda_{\boldsymbol{\theta}}(9\lambda_{\boldsymbol{\theta}} + 24C)}}{\frac{8}{3}C^2}. \quad (3.16)$$

Unfortunately, this bound was optimistic in the limit for the simple Euclidean case ( $C = 0$ ). In a more recent paper [75], the authors improved the bound by choosing  $\lambda_{\boldsymbol{\theta}} = \text{tr}(\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}))^{-1}$  and an alternative lower bound for  $t \cot(t)$  to obtain

$$\text{var}(\vartheta) \geq \frac{\lambda_{\boldsymbol{\theta}}C + 1 - \sqrt{2\lambda_{\boldsymbol{\theta}}C + 1}}{\frac{1}{2}C^2\lambda_{\boldsymbol{\theta}}}. \quad (3.17)$$

Although the authors omitted a proof<sup>8</sup>, the alternative lower bound for  $t \cot(t)$  appears to be  $t \cot(t) \geq 1 - \frac{1}{2}t^2$  for  $0 \leq t \leq T$ . An immediate benefit from this

---

<sup>8</sup>Xavier and Barroso stated in [75] that the proof would “be found in the companion paper [14]”, but this companion paper appears to have never been published.

improved bound is the agreement with the simple ( $C = 0$ ) Euclidean case, i.e.,  $\text{var}(\vartheta) \geq \lambda_{\boldsymbol{\theta}} = \text{tr}(\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}))^{-1}$ . This consistency is entirely due to the choice of  $\lambda_{\boldsymbol{\theta}}$ . Even so, greater improvement of the bound, not found in the literature, is possible in the curved scenarios ( $C > 0$ ) at least in the bound of  $t \cot(t)$ . Note, for  $0 \leq t < \pi$

$$\begin{aligned} t \cot(t) &= 1 + \sum_{s=1}^{\infty} \frac{(-1)^s 2^{2s} B_{2s} t^{2s}}{(2s)!} = 1 - \sum_{s=1}^{\infty} 2 \left( \frac{t}{\pi} \right)^{2s} \zeta(2s) \\ &\geq 1 - \sum_{s=1}^{\infty} 2 \left( \frac{t}{\pi} \right)^{2s} \frac{\pi^2}{6} = 1 - \frac{\pi^2}{3} \frac{\left( \frac{t}{\pi} \right)^2}{1 - \left( \frac{t}{\pi} \right)^2} \\ &\geq 1 - \frac{b}{2} t^2 \end{aligned}$$

where  $B_{2s}$  are Bernoulli numbers,  $\zeta(\cdot)$  is the Riemann zeta function, and  $b = \frac{2/3}{1 - \left( \frac{\pi}{\pi} \right)^2} \approx 0.8177$ . Then

$$\text{var}(\vartheta) \geq \frac{\lambda_{\boldsymbol{\theta}} C b + 1 - \sqrt{2\lambda_{\boldsymbol{\theta}} C b + 1}}{\frac{1}{2} C^2 b^2 \lambda_{\boldsymbol{\theta}}}. \quad (3.18)$$

### 3.1.3 Simple Extensions to the CCRB

As with unconstrained parameters, the performance of a function of parameters often may be of greater interest. Consider a continuously differentiable function  $\mathbf{k} : \Theta_f \rightarrow \mathbb{R}^q$ . Denote the Jacobian of this transfer function to be  $\mathbf{K}(\boldsymbol{\theta}) = \left. \frac{\partial \mathbf{k}(\boldsymbol{\theta}')}{\partial \boldsymbol{\theta}'^T} \right|_{\boldsymbol{\theta}' = \boldsymbol{\theta}}$ . We have a simple extension of the classical transformation of parameters in section 2.1.1.

**Corollary 3.10.** If  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$ , then the variance of any unbiased estimator  $\mathbf{S}(\mathbf{x})$  of

$\boldsymbol{\alpha} = \mathbf{k}(\boldsymbol{\theta})$  satisfies the inequality

$$\text{Var}(\mathbf{S}(\mathbf{x})) \geq \text{CCRB}(\boldsymbol{\alpha}) \triangleq \mathbf{K}(\boldsymbol{\theta})\text{CCRB}(\boldsymbol{\theta})\mathbf{K}^T(\boldsymbol{\theta}).$$

*Proof.* Let  $\mathbf{g}_{\boldsymbol{\theta}}$  be the implicit function defined by  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$ . Then  $\boldsymbol{\alpha} = \mathbf{k}(\mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi}))$  has a Jacobian  $\mathbf{K}(\boldsymbol{\theta})\mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\xi})$ . From the classical transformation of parameters, the inequality  $\text{Var}(\mathbf{S}(\mathbf{x})) \geq \text{CRB}(\boldsymbol{\alpha})$  holds where

$$\begin{aligned} \text{CRB}(\boldsymbol{\alpha}) &= \mathbf{K}(\boldsymbol{\theta})\mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\xi})\tilde{\mathbf{I}}^{-1}(\boldsymbol{\xi})\mathbf{G}_{\boldsymbol{\theta}}^T(\boldsymbol{\xi})\mathbf{K}^T(\boldsymbol{\theta}) \\ &= \mathbf{K}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})\left(\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})\right)^{-1}\mathbf{U}^T(\boldsymbol{\theta})\mathbf{K}^T(\boldsymbol{\theta}). \end{aligned}$$

□

This transformation property is useful in extending the constrained Cramér-Rao bound to biased estimation.

**Example 3.11** (Biased Estimation). Assume  $\mathbf{t}(\mathbf{x})$  is a biased estimator of a constrained parameter  $\boldsymbol{\theta}$  with bias  $\mathbf{b}(\boldsymbol{\theta}) = E_{\boldsymbol{\theta}}\mathbf{t}(\mathbf{x}) - \boldsymbol{\theta}$  and constraint  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$ . Define  $\mathbf{k}(\boldsymbol{\theta}) = \boldsymbol{\theta} + \mathbf{b}(\boldsymbol{\theta})$ . Then  $\mathbf{t}(\mathbf{x})$  is an unbiased estimator of  $\boldsymbol{\alpha} = \mathbf{k}(\boldsymbol{\theta})$ . Since  $\mathbf{K}(\boldsymbol{\theta}) = \mathbf{I}_{m \times m} + \mathbf{B}(\boldsymbol{\theta})$  where  $\mathbf{B}(\boldsymbol{\theta}) = \left. \frac{\partial \mathbf{b}(\boldsymbol{\theta}')}{\partial \boldsymbol{\theta}'^T} \right|_{\boldsymbol{\theta}' = \boldsymbol{\theta}}$ , then we have the inequality

$$\text{Var}(\mathbf{t}(\mathbf{x})) \geq (\mathbf{I}_{m \times m} + \mathbf{B}(\boldsymbol{\theta}))\text{CCRB}(\boldsymbol{\theta})(\mathbf{I}_{m \times m} + \mathbf{B}^T(\boldsymbol{\theta})) \triangleq \text{CCRB}(\boldsymbol{\theta} + \mathbf{b}(\boldsymbol{\theta})).$$

Often, when the Fisher information matrix is singular, its pseudoinverse is used as a bound on the variance of an estimator. As mentioned in section 2.1.1, the bound is trivially true for some component of the estimator except under certain conditions.

**Corollary 3.12.** If  $\mathbf{t}(\mathbf{x})$  is an estimator of  $\mathbf{k}(\boldsymbol{\theta})$  having bias  $\mathbf{b}(\boldsymbol{\theta})$ , where  $\boldsymbol{\theta}$  satisfies the constraint  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$ , then the inequality

$$\text{Var}(\mathbf{t}(\mathbf{x})) \geq \mathbf{H}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) \left( \mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) \right)^\dagger \mathbf{U}^T(\boldsymbol{\theta})\mathbf{H}^T(\boldsymbol{\theta}),$$

where  $\mathbf{H}(\boldsymbol{\theta}) = \mathbf{K}(\boldsymbol{\theta}) + \mathbf{B}(\boldsymbol{\theta})$ , is nontrivially satisfied, i.e., all components of  $\mathbf{t}(\mathbf{x})$  have finite variance, if and only if

$$\mathbf{H}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) = \mathbf{H}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) \left( \mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) \right)^\dagger \mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}).$$

*Proof.* Suppose  $\mathbf{t}(\mathbf{x})$  is an estimator of  $\boldsymbol{\alpha} = \mathbf{k}(\boldsymbol{\theta})$  with bias  $\mathbf{b}(\boldsymbol{\theta})$  and under the constraint  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$ . Define  $\mathbf{H}(\boldsymbol{\theta}) = \left. \frac{\partial}{\partial \boldsymbol{\theta}^T} (\mathbf{k}(\boldsymbol{\theta}) + \mathbf{b}(\boldsymbol{\theta})) \right|_{\boldsymbol{\theta}' = \boldsymbol{\theta}} = \mathbf{K}(\boldsymbol{\theta}) + \mathbf{B}(\boldsymbol{\theta})$ . If  $\mathbf{g}_\theta$  is the implicit function defined by  $\mathbf{f}$ , then we have the inequality

$$\text{Var}(\mathbf{t}(\mathbf{x})) \geq \tilde{\mathbf{H}}(\boldsymbol{\xi})\tilde{\mathbf{I}}^\dagger(\boldsymbol{\xi})\tilde{\mathbf{H}}^T(\boldsymbol{\xi})$$

where  $\tilde{\mathbf{H}}(\boldsymbol{\xi}) = \left. \frac{\partial}{\partial \boldsymbol{\xi}^T} (\mathbf{k}(\mathbf{g}_\theta(\boldsymbol{\xi}')) + \mathbf{b}(\mathbf{g}_\theta(\boldsymbol{\xi}')) \right|_{\boldsymbol{\xi}' = \boldsymbol{\xi}} = \mathbf{H}(\boldsymbol{\theta})\mathbf{G}_\theta(\boldsymbol{\xi})$ . (This can also be inferred from [23, lemma 2], although no assumption is made here about the singularity of the Fisher information.) Hence, all the components of  $\mathbf{t}(\mathbf{x})$  can have finite variance if and only if [66]

$$\begin{aligned} \tilde{\mathbf{H}}(\boldsymbol{\xi}) &= \tilde{\mathbf{H}}(\boldsymbol{\xi})\tilde{\mathbf{I}}(\boldsymbol{\xi})\tilde{\mathbf{I}}^\dagger(\boldsymbol{\xi}) \\ \mathbf{H}(\boldsymbol{\theta})\mathbf{G}_\theta(\boldsymbol{\xi}) &= \mathbf{H}(\boldsymbol{\theta})\mathbf{G}_\theta(\boldsymbol{\xi})\tilde{\mathbf{I}}^\dagger(\boldsymbol{\xi})\tilde{\mathbf{I}}(\boldsymbol{\xi}) \\ &= \mathbf{H}(\boldsymbol{\theta})\mathbf{G}_\theta(\boldsymbol{\xi})\tilde{\mathbf{I}}^\dagger(\boldsymbol{\xi})\mathbf{G}_\theta^T(\boldsymbol{\xi})\mathbf{I}(\boldsymbol{\theta})\mathbf{G}_\theta(\boldsymbol{\xi}) \\ &= \mathbf{H}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) \left( \mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) \right)^\dagger \mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{G}_\theta(\boldsymbol{\xi}) \\ \mathbf{H}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) &= \mathbf{H}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) \left( \mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) \right)^\dagger \mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}). \end{aligned}$$

□

By definition  $\mathbf{U}(\boldsymbol{\theta})$  will always be full column rank, so singularity of the matrix  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  implies singularity in the FIM  $\mathbf{I}(\boldsymbol{\theta})$  and insufficiency in the constraints to resolve the inherent ambiguities in the model. This scenario will be discussed in greater detail in the next section.

### 3.1.4 Properties of the CCRB

The constrained Cramér-Rao bound for the constraint  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$  is equivalent to the CRB for any reparameterization of the parameters satisfying the constraint. Implicit in that equivalence and in the proof of the CCRB presented in section 3.1.1 are the regularity conditions for the CRB on the implicit parameter  $\boldsymbol{\xi}$ , i.e.,  $\frac{\partial}{\partial \boldsymbol{\xi}^T} E_{\boldsymbol{\xi}}(\mathbf{h}(x)) = E_{\boldsymbol{\xi}}(\mathbf{h}(x)\tilde{\mathbf{s}}^T(\mathbf{x}; \boldsymbol{\xi}))$  for  $\mathbf{h}(x) \equiv 1$  and  $\mathbf{h}(x) \equiv \mathbf{t}(x)$ , where  $\mathbf{t}(x)$  in this case is an unbiased estimator of  $\boldsymbol{\theta} = \mathbf{g}(\boldsymbol{\xi})$ . This condition translates to  $\frac{\partial}{\partial \boldsymbol{\theta}^T} E_{\boldsymbol{\xi}}(\mathbf{h}(x))\mathbf{G}(\boldsymbol{\xi}) = E_{\boldsymbol{\xi}}(\mathbf{h}(x)\mathbf{s}^T(\mathbf{x}; \boldsymbol{\theta})\mathbf{G}(\boldsymbol{\xi}))$  or, strictly in terms of  $\boldsymbol{\theta}$ ,

$$\frac{\partial}{\partial \boldsymbol{\theta}^T} E_{\boldsymbol{\theta}}(\mathbf{h}(x))\mathbf{U}(\boldsymbol{\theta}) = E_{\boldsymbol{\theta}}(\mathbf{h}(x)\mathbf{s}^T(\mathbf{x}; \boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})) \quad (3.19)$$

for  $\mathbf{h}(x) \equiv 1$  and  $\mathbf{h}(x) \equiv \mathbf{t}(x)$ . From this condition, it can be shown that

$$E_{\boldsymbol{\theta}}(\mathbf{t}(x) - \boldsymbol{\theta})\mathbf{s}^T(\mathbf{x}; \boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})\mathbf{U}^T(\boldsymbol{\theta}) = \mathbf{U}(\boldsymbol{\theta})\mathbf{U}^T(\boldsymbol{\theta}), \quad (3.20)$$

which is a variant of the regularity condition required and proven by Marzetta [47] and the same regularity condition simply stated (but not proven) by Stoica and Ng [68] in their proofs of the CCRB. Thus, as with the CCRB, the regularity condition for the constraint is equivalent to the regularity condition for the reparameterization of parameters satisfying the constraint. This general fact, a restatement of theorem

3.5, is the quintessential property of the CCRB formula. More explicitly, the CCRB is a generalization of both the degenerate and the determinate constraint cases, as examples 3.13 and 3.14 demonstrate.

**Example 3.13** (Degenerate Case). The scenario with no constraint is equivalent to the statement that the function which describes the constraint is null. That is,  $\mathbf{f} : \Theta \rightarrow \mathbb{R}^0$  and  $\mathbf{f}(\boldsymbol{\theta}) = []$ . Then  $\mathbf{F} : \Theta \rightarrow \mathbb{R}^{0 \times m}$  is also a null gradient row vector having rank 0. Any nonsingular  $m \times m$  matrix  $\mathbf{U}(\boldsymbol{\theta})$  satisfies (3.6), and thus

$$\mathbf{U}(\boldsymbol{\theta}) (\mathbf{U}^T(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}))^{-1} \mathbf{U}^T(\boldsymbol{\theta}) = \mathbf{I}^{-1}(\boldsymbol{\theta}).$$

Therefore the CCRB formula also incorporates the unconstrained scenario.

**Example 3.14** (Determinate Case). Suppose the constraint  $\mathbf{f}$  completely determines the parameter. Then necessarily, since we have  $m$  unknowns in the parameter vector  $\boldsymbol{\theta}$ , there must be at least  $k \geq m$  equations in the constraint equation and the Jacobian  $\mathbf{F}(\boldsymbol{\theta})$  must have rank  $m$ . Since  $\mathbf{f}$  is assumed to have nonredundant constraints,  $\mathbf{F}(\boldsymbol{\theta})$  is actually a nonsingular square matrix. Only the null vector  $\mathbf{U} : \Theta \rightarrow \mathbb{R}^{m \times 0}$  satisfies (3.6) ( $\mathbf{U}^T(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}) = \mathbf{I}_{0 \times 0}$ ), therefore  $\mathbf{U}(\boldsymbol{\theta}) (\mathbf{U}^T(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}))^{-1} \mathbf{U}^T(\boldsymbol{\theta})$  is a null element, and the CCRB does not exist for a completely determined parameter set.

Thus far, only parametric equality constraints have been considered. Gorman and Hero show that only *active* constraints results in a reduction in the bound [23, lemma 4], i.e., *inactive* (or strict) inequality constraints do not contribute information to the model in the CCRB sense.<sup>9</sup> As the test points approach the parameter

---

<sup>9</sup>Because the CCRB (CRB) is a local bound that only accounts for local fluctuations of the

in the Chapman-Robbins bound, they become interior points of the constraint set and the inequalities have no impact on the information metric. This can also be shown without resorting to the Chapman-Robbins approach.

**Example 3.15** (Inequality Constraints Are Non-Informative). Assume, in addition to the constraint set  $\Theta_f$ , the parameters are required to satisfy the strict inequality  $h(\boldsymbol{\theta}) < 0$  where  $h : \Theta \rightarrow \mathbb{R}$  is a continuously differentiable function. To incorporate this constraint, we introduce a *dummy* parameter  $\vartheta$  and add a new equality constraint  $\vartheta^2 = -h(\boldsymbol{\theta})$ , which is equivalent to the strict inequality constraint (whenever  $\vartheta \neq 0$ ). In addition we create an extended parameterization  $\boldsymbol{\phi}' = \begin{bmatrix} \boldsymbol{\theta}' \\ \vartheta' \end{bmatrix} \in \mathbb{R}^{m+1}$  and constraint function  $\mathbf{f}^*(\boldsymbol{\phi}') = \begin{bmatrix} \mathbf{f}(\boldsymbol{\theta}') \\ \vartheta'^2 + h(\boldsymbol{\theta}') \end{bmatrix}$ . This will generate a Fisher information and a Jacobian matrix defined by

$$\mathbf{I}^*(\boldsymbol{\phi}) = \begin{bmatrix} \mathbf{I}(\boldsymbol{\theta}) & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix}, \quad \mathbf{F}^*(\boldsymbol{\phi}) = \begin{bmatrix} \mathbf{F}(\boldsymbol{\theta}) & \mathbf{0} \\ \mathbf{H}^T(\boldsymbol{\theta}) & 2\vartheta \end{bmatrix},$$

respectively, where  $\mathbf{H}(\boldsymbol{\theta}) = \left. \frac{\partial h(\boldsymbol{\theta}')}{\partial \boldsymbol{\theta}'} \right|_{\boldsymbol{\theta}' = \boldsymbol{\theta}}$ . Note that  $\mathbf{F}^*(\boldsymbol{\phi})$  will be full row rank as long as  $\vartheta \neq 0$ . If  $\mathbf{U}(\boldsymbol{\theta})$  is defined as in (3.6), then  $\mathbf{U}^*(\boldsymbol{\phi}) = \begin{bmatrix} \mathbf{U}(\boldsymbol{\theta}) \\ \mathbf{v}^T(\boldsymbol{\theta}, \vartheta) \end{bmatrix}$  will also satisfy (3.6) with respect to  $\mathbf{F}^*(\boldsymbol{\phi})$ , where  $\mathbf{v}(\boldsymbol{\theta}, \vartheta) = \frac{1}{2\vartheta} \mathbf{H}^T(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta})$ . By theorem

---

true parameter, inactive constraints have no impact on performance potential. As such the CCRB (CRB) only provides information on the pdf for the mainlobe of the density function, which typically corresponds to the true parameter. For a number of scenarios, e.g., when there is sufficiently large variance in the pdf, sidelobes of the distribution impact the performance, thereby making the CCRB (CRB) overly optimistic. This occurs frequently in communications when the signal-to-noise ratio (SNR) or data transmission size decreases.

3.1, then

$$\begin{aligned}
& \mathbf{U}^*(\boldsymbol{\phi}) \left( \mathbf{U}^{*T}(\boldsymbol{\phi}) \mathbf{I}^*(\boldsymbol{\phi}) \mathbf{U}^*(\boldsymbol{\phi}) \right)^{-1} \mathbf{U}^{*T}(\boldsymbol{\phi}) \\
&= \mathbf{U}^*(\boldsymbol{\phi}) \left( \mathbf{U}(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}) \right)^{-1} \mathbf{U}^{*T}(\boldsymbol{\phi}) \\
&= \begin{bmatrix} \mathbf{U}(\boldsymbol{\theta}) \left( \mathbf{U}(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}) \right)^{-1} \mathbf{U}(\boldsymbol{\theta}) & \mathbf{U}(\boldsymbol{\theta}) \left( \mathbf{U}(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}) \right)^{-1} \mathbf{v}(\boldsymbol{\theta}, \vartheta) \\ \mathbf{v}^T(\boldsymbol{\theta}, \vartheta) \left( \mathbf{U}(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}) \right)^{-1} \mathbf{U}^T(\boldsymbol{\theta}) & \mathbf{v}^T(\boldsymbol{\theta}, \vartheta) \left( \mathbf{U}(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}) \right)^{-1} \mathbf{v}(\boldsymbol{\theta}, \vartheta) \end{bmatrix}.
\end{aligned}$$

Hence, the CCRB on the  $\boldsymbol{\theta}$  components of  $\boldsymbol{\phi}$ , the upper-left submatrix, remains unchanged.

Equality constraints, however, do add side information to the model. Thus, it is intuitive to expect that the constrained model should result in a lower bound compared to the bound for the model without constraints. This statement was made in Gorman and Hero [23, p.1292], but not in Marzetta [47] nor Stoica and Ng [68]. In the latter case, the statement is only true under certain conditions. Prior to establishing when the bound is lowered, a powerful lemma will be proven.

**Lemma 3.16.** For an arbitrary full column rank matrix  $\mathbf{A}$ , and an arbitrary symmetric positive semidefinite matrix  $\mathbf{B}$ , the inequality

$$\mathbf{A} \left( \mathbf{A}^T \mathbf{B} \mathbf{A} \right)^\dagger \mathbf{A}^T \leq \mathbf{B}^\dagger \quad (3.21)$$

holds over the projection subspace of  $\mathbf{B}^\dagger \mathbf{B}$  with equality if and only if  $\text{rank}(\mathbf{A}^T \mathbf{B} \mathbf{A}) = \text{rank}(\mathbf{B})$ .

*Proof.* Let  $\mathbf{L}\mathbf{L}^T$  be the Cholesky decomposition of  $\mathbf{B}$  [20, p. 194]. Then  $\mathbf{L} \in \mathbb{R}^{m \times \text{rank}(\mathbf{B})}$  is a full column rank lower triangular matrix with strictly positive values on the diagonal. To show the inequality, consider linear unbiased estimates of the



mean in the model  $\mathbf{y} = \mathbf{D}\boldsymbol{\beta} + \boldsymbol{\epsilon}$  where  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, (\mathbf{L}^T \mathbf{L})^{-2})$ , and  $\boldsymbol{\beta}$  is treated as the unknown parameter. In particular,  $\mathbf{y}$  is such an estimate with variance equaling  $(\mathbf{L}^T \mathbf{L})^{-2}$  where as the best linear unbiased estimate (see section 2.3.1)

$$\mathbf{D}\hat{\boldsymbol{\beta}} = \mathbf{D} (\mathbf{D}^T (\mathbf{L}^T \mathbf{L})^2 \mathbf{D})^\dagger \mathbf{D}^T (\mathbf{L}^T \mathbf{L})^2 \mathbf{y},$$

with variance equal to  $\mathbf{D} (\mathbf{D}^T (\mathbf{L}^T \mathbf{L})^2 \mathbf{D})^\dagger \mathbf{D}^T$ . By the Gauss-Markov theorem we have the inequality with the BLUE's variance

$$\mathbf{D} (\mathbf{D}^T (\mathbf{L}^T \mathbf{L})^2 \mathbf{D})^\dagger \mathbf{D}^T \leq (\mathbf{L}^T \mathbf{L})^{-2}$$

with equality if and only if  $(\mathbf{L}^T \mathbf{L}) \mathbf{D} (\mathbf{D}^T (\mathbf{L}^T \mathbf{L})^2 \mathbf{D})^\dagger \mathbf{D}^T (\mathbf{L}^T \mathbf{L}) = \mathbf{I}_{\text{rank}(\mathbf{L}) \times \text{rank}(\mathbf{L})}$ . Substituting in  $\mathbf{D} = (\mathbf{L}^T \mathbf{L})^{-1} \mathbf{L}^T \mathbf{A}$  in 3.1.4 and pre- and post-multiplying both sides by  $\mathbf{L}$  and  $\mathbf{L}^T$ , respectively, we have the inequality

$$\mathbf{B}^\dagger \mathbf{B} \mathbf{A} (\mathbf{A}^T \mathbf{L} \mathbf{L}^T \mathbf{A})^\dagger \mathbf{A}^T \mathbf{B} \mathbf{B}^\dagger \leq \mathbf{B}^\dagger.$$

Considering quadratic forms in  $\mathbf{B}^\dagger \mathbf{B} \mathbf{v}$  proves the result since for any vector  $\mathbf{v} \in \mathbb{R}^m$ ,  $\mathbf{v}^T \mathbf{B}^\dagger \mathbf{B} \mathbf{B}^\dagger \mathbf{B} \mathbf{B}^\dagger \mathbf{v} = \mathbf{v}^T \mathbf{B}^\dagger \mathbf{v}$ . Moreover, from the definition of  $\mathbf{D}$ , we have  $(\mathbf{L}^T \mathbf{L}) \mathbf{D} (\mathbf{D}^T (\mathbf{L}^T \mathbf{L})^2 \mathbf{D})^\dagger \mathbf{D}^T (\mathbf{L}^T \mathbf{L}) = \mathbf{I}_{\text{rank}(\mathbf{L}) \times \text{rank}(\mathbf{L})}$  if and only if it can be shown that  $\mathbf{L}^T \mathbf{A} (\mathbf{A}^T \mathbf{L} \mathbf{L}^T \mathbf{A})^\dagger \mathbf{A}^T \mathbf{L} = \mathbf{I}_{\text{rank}(\mathbf{L}) \times \text{rank}(\mathbf{L})}$ , which is so if and only if  $\mathbf{A}$  and  $\mathbf{L}$  satisfy  $\text{rank}(\mathbf{A}^T \mathbf{L}) = \text{rank}(\mathbf{B})$ . □

This lemma proves the following results, including that in a linear subspace the  $\text{CCRB}(\boldsymbol{\theta})$  is lesser than or equal to the  $\text{CRB}(\boldsymbol{\theta})$  in the matrix sense and that constraints strictly on the ambiguous information have no impact on the performance of the unambiguous information.

**Theorem 3.17.** Let  $\mathbf{U}(\boldsymbol{\theta})$  be defined as in (3.6) from a constraint  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$ . Then

$$\mathbf{U}(\boldsymbol{\theta}) (\mathbf{U}^T(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}))^\dagger \mathbf{U}^T(\boldsymbol{\theta}) \leq \mathbf{I}^\dagger(\boldsymbol{\theta})$$

over the linear projection subspace defined by  $\mathbf{I}^\dagger(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta})$ .

*Proof.* From the lemma, defining  $\mathbf{A} = \mathbf{U}(\boldsymbol{\theta})$  and  $\mathbf{B} = \mathbf{I}(\boldsymbol{\theta})$  gives the inequality.

□

The conditions under which  $\mathbf{U}(\boldsymbol{\theta}) (\mathbf{U}^T(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}))^\dagger \mathbf{U}^T(\boldsymbol{\theta})$  and  $\mathbf{I}^\dagger(\boldsymbol{\theta})$  are nontrivial bounds are detailed in corollary 3.12 and section 2.1.1, respectively. The projection space of  $\mathbf{I}^\dagger(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta})$  corresponds to the identifiable components of  $\boldsymbol{\theta}$  without the constraints (e.g., see [18, section 4]). The theorem establishes the result that the constraints can only lower the bound and thereby increase performance potential for (functions of) parameters that are already identifiable. This is regardless of whether the Fisher information is singular or whether the parameters are identifiable under constraints (whether  $\mathbf{U}^T(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta})$  is singular, see Theorem 3.24). Theorem 3.17 is more general than the result shown by Ash<sup>10</sup>, which is given by the following corollary.

**Corollary 3.18** (Ash). Provided no linear combination of  $\mathbf{U}(\boldsymbol{\theta})$  lies in the null space of the Fisher information matrix, then

$$\mathbf{U}(\boldsymbol{\theta}) (\mathbf{U}^T(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}))^{-1} \mathbf{U}^T(\boldsymbol{\theta}) \leq \mathbf{I}^\dagger(\boldsymbol{\theta})$$

over the linear projection subspace defined by  $\mathbf{I}^\dagger(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta})$ .

---

<sup>10</sup>Although not stated in this manner, this result appears in Ash's thesis [5, equation (3.63)] as well as in a publication of his third chapter in Ash and Moses [6, equation (63)]. The results from Lemma 3.16 are more general than those in [5] or [6].

*Proof.* Since  $\mathbf{U}(\boldsymbol{\theta})$  is full column rank and the range space of its columns does not exist in the null space of the FIM, then  $(\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}))$  is nonsingular.  $\square$

Furthermore, if the Fisher information matrix is nonsingular, then  $\mathbf{I}^\dagger(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})$  is the identity matrix, so the inequality holds as a quadratic form over  $\mathbb{R}^m$  and

$$\mathbf{U}(\boldsymbol{\theta}) (\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}))^{-1} \mathbf{U}^T(\boldsymbol{\theta}) \leq \mathbf{I}^{-1}(\boldsymbol{\theta}).$$

This was also cleverly shown by Gorman and Hero [23, (44) in remark 4], again, with the unnecessary assumption that  $\mathbf{I}(\boldsymbol{\theta})$  is regular (non-singular). Next, we use the lemma to observe the existence of *non-informative* constraints.

**Corollary 3.19.** Assume the row vectors of  $\mathbf{F}(\boldsymbol{\theta}) \in \mathbb{R}^{k \times m}$  form a linearly independent basis for the null space of  $\mathbf{I}(\boldsymbol{\theta})$ . Then  $\mathbf{F}_1(\boldsymbol{\theta}) \in \mathbb{R}^{k' \times m}$  is a linear combination of a submatrix of  $\mathbf{F}(\boldsymbol{\theta})$  if and only if

$$\mathbf{U}_1(\boldsymbol{\theta}) (\mathbf{U}_1^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}_1(\boldsymbol{\theta}))^\dagger \mathbf{U}_1^T(\boldsymbol{\theta}) = \mathbf{I}^\dagger(\boldsymbol{\theta}),$$

where  $\mathbf{U}_1(\boldsymbol{\theta})$  is defined as in (3.6) relative to  $\mathbf{F}_1(\boldsymbol{\theta})$ .

*Proof.* Without loss of generality, partition  $\mathbf{F}$  as  $\begin{bmatrix} \mathbf{F}_1(\boldsymbol{\theta}) \\ \mathbf{F}_2(\boldsymbol{\theta}) \end{bmatrix}$ . Then if  $\mathbf{U}(\boldsymbol{\theta})$  is defined as in (3.6), define

$$\mathbf{U}_1(\boldsymbol{\theta}) = \left[ \mathbf{U}(\boldsymbol{\theta}), \left( I_{m \times m} - \mathbf{F}_1^T (\mathbf{F}_1 \mathbf{F}_1^T)^{-1} \mathbf{F}_1 \right) \mathbf{F}_2^T(\boldsymbol{\theta}) \mathbf{D}(\boldsymbol{\theta}) \right],$$

where  $\mathbf{D}(\boldsymbol{\theta})$  represents a Gram-Schmidt processing matrix which orthonormalizes the column vectors of  $\left( I_{m \times m} - \mathbf{F}_1^T (\mathbf{F}_1 \mathbf{F}_1^T)^{-1} \mathbf{F}_1 \right) \mathbf{F}_2^T(\boldsymbol{\theta})$  (these are already orthogonal to the column vectors of  $\mathbf{U}(\boldsymbol{\theta})$ ). This satisfies (3.6). Since  $\text{span}(\mathbf{I}(\boldsymbol{\theta})) \subset$

$\text{span}(\mathbf{U}_1(\boldsymbol{\theta}))$  then  $\text{rank}(\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})) = \text{rank}(\mathbf{I}(\boldsymbol{\theta}))$ , and thus, by lemma 3.16,

$$\mathbf{U}_1(\boldsymbol{\theta}) (\mathbf{U}_1^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}_1(\boldsymbol{\theta}))^\dagger \mathbf{U}_1^T(\boldsymbol{\theta}) = \mathbf{I}^\dagger(\boldsymbol{\theta}).$$

□

This corollary details the necessary and sufficient characteristics for a constraint to be non-informative. Implicit from the lemma and above corollary is the CCRB's invariance to linear combinations of columns of  $\mathbf{U}(\boldsymbol{\theta})$  that exist in the null space of  $\mathbf{I}(\boldsymbol{\theta})$ . The specific case when  $\text{span}(\mathbf{I}(\boldsymbol{\theta})) = \text{span}(\mathbf{U}(\boldsymbol{\theta}))$  is also shown by Ash in [5, equation (3.73)] and [6, equation (73)].<sup>11</sup>

The CCRB can be interpreted geometrically, as in Figure 3.1 as a contraction of the information ambiguity to find the bound and then an expansion. The column vectors of  $\mathbf{U}(\boldsymbol{\theta})$ , being in the null space of the Jacobian of the constraints, restrict the information in  $\Theta$  into the constraint space  $\Theta_f$ . The bound can be found from (the inverse of) the information in  $\Theta_f$ , and is then projected back into the original space of the parameters of interest. This down-and-up projection is easiest to observe for linear constraints, where any local reparameterization is also a global reparameterization.

**Example 3.20.** Assume the parameters satisfy the linear constraint

$$\mathbf{f}(\boldsymbol{\theta}) = \mathbf{F}\boldsymbol{\theta} + \mathbf{v} = \mathbf{0}.$$

In this case, solutions of  $\boldsymbol{\theta}$  are of the form  $\boldsymbol{\theta} = -\mathbf{F}^T (\mathbf{F}\mathbf{F}^T)^{-1} \mathbf{v} + \mathbf{U}\boldsymbol{\xi} = \mathbf{g}(\boldsymbol{\xi})$  where

---

<sup>11</sup>In Ash's thesis and paper, this scenario was referred to as the *minimally constrained* case, where only all the unknown information ambiguities are constrained. However, in this current work, the minimal number of constraints is zero, the degenerate constraint case of example 3.13, thus it makes more sense to refer to this scenario as a non-informative constraint.

$\mathbf{U}$  is defined as in (3.6). That this is a one-to-one correspondence locally (or in this case globally) is made clear by the existence of the inverse  $\mathbf{g}^{-1}(\boldsymbol{\theta}) = \mathbf{U}^T \boldsymbol{\theta} = \boldsymbol{\xi}$ . Then, since  $\mathbf{G}(\boldsymbol{\xi}) = \mathbf{U}$ , we have from (3.9)

$$\text{CCRB}(\boldsymbol{\theta}) = \mathbf{U} (\mathbf{U}^T \mathbf{I}(\boldsymbol{\theta}) \mathbf{U})^{-1} \mathbf{U}^T$$

which is exactly the CCRB provided in (3.5).

Alternatively, another interpretation [9] is that the CCRB is less than (in a matrix sense) the CRB because the performance bound is over an expanded class of estimators. The (non-biased) CRB is a bound on the mean-square error of unbiased estimators in  $\Theta$ , whereas the CCRB is a bound for estimators that only need to be unbiased on  $\Theta_f$  and not on the whole set  $\Theta$ .

### 3.1.5 Derivation of $\mathbf{g}(\boldsymbol{\xi})$

In the proof of the CCRB as well as its application, the need for an explicit reparameterization  $\mathbf{g}_{\boldsymbol{\theta}}$  proved unnecessary. However, in examples 3.4 and 3.20 of the previous section, we presented scenarios where, given the constraint function  $\mathbf{f}$ , we were able to define a locally equivalent continuously differentiable  $\mathbf{g}_{\boldsymbol{\theta}}$ . There may exist other scenarios where it is desirable to obtain a  $\mathbf{g}_{\boldsymbol{\theta}}$  explicitly.

In this section, we present two procedures to do so. First, we detail an approach using the Taylor expansion of  $\mathbf{g}_{\boldsymbol{\theta}}$  given explicit knowledge of  $\mathbf{U}(\boldsymbol{\theta}')$ . Then, we demonstrate an approach based on fixed point methods.

### 3.1.5.1 A Taylor series derivation

The Taylor series expansion of  $\boldsymbol{\theta}' = \mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi}')$  about  $\boldsymbol{\xi}$  is given by

$$\mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi}') = \boldsymbol{\theta} + \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\xi})(\boldsymbol{\xi}' - \boldsymbol{\xi}) + \cdots$$

where  $\boldsymbol{\theta} = \mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi})$ . For any differentiable function  $h : \mathbb{R}^m \rightarrow \mathbb{R}$  we have

$$\frac{\partial}{\partial \boldsymbol{\xi}'^T} h(\boldsymbol{\theta}') = \frac{\partial}{\partial \boldsymbol{\theta}'^T} h(\boldsymbol{\theta}') \cdot \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\xi}').$$

Since the implicit function from  $\mathbb{R}^{m-k}$  into  $\Theta_f$  is not unique, we can choose a reparameterization which uses the transformation matrix  $\mathbf{S}(\boldsymbol{\theta}) = \mathbf{I}_{m-k}$ , i.e., we choose a reparameterization with Jacobian  $\mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\xi}') = \mathbf{U}(\mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi}'))$  for any null space matrix  $\mathbf{U}(\boldsymbol{\theta})$  that satisfies (3.6) and choose  $\boldsymbol{\xi} = \mathbf{0}$ . With this selection of the reparameterization, the  $r$ th order derivatives of  $\mathbf{g}_{\boldsymbol{\theta}}$  are the  $(r-1)$ st order derivatives of the elements of  $\mathbf{U}(\boldsymbol{\theta}')$  with respect to  $\boldsymbol{\xi}'$ , which are to be evaluated at  $\boldsymbol{\theta}$  for the coefficients in the Taylor series.

**Example 3.21.** Reviewing example 3.4, note

$$\mathbf{U}(\boldsymbol{\theta}') = \begin{bmatrix} \theta'_2 \\ -\theta'_1 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} \theta'_1 \\ \theta'_2 \end{bmatrix}$$

and hence

$$\frac{\partial \mathbf{U}(\boldsymbol{\theta}')}{\partial \boldsymbol{\theta}'^T} = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \quad (3.22)$$

is independent of  $\boldsymbol{\theta}'$ , so using (3.22) as a reference, we have

$$\frac{\partial^r \mathbf{U}(\boldsymbol{\theta}')}{\partial \boldsymbol{\xi}'^r} = \frac{\partial^{r-1}}{\partial \boldsymbol{\xi}'^{r-1}} \left( \frac{\partial \mathbf{U}(\boldsymbol{\theta}')}{\partial \boldsymbol{\theta}'^T} \cdot \mathbf{U}(\boldsymbol{\theta}') \right) = \cdots = \left( \frac{\partial \mathbf{U}(\boldsymbol{\theta}')}{\partial \boldsymbol{\theta}'^T} \right)^r \mathbf{U}(\boldsymbol{\theta}').$$

With a natural selection of  $\boldsymbol{\theta} = \mathbf{g}_{\boldsymbol{\theta}}(0) = (1, 0)^T$  as an initial value for the Taylor series, the reparameterized function is found to be

$$\begin{aligned}
\boldsymbol{\theta}' = \mathbf{g}(\xi') &= \sum_{r=0}^{\infty} \left( \frac{\partial \mathbf{U}(\boldsymbol{\theta}')}{\partial \boldsymbol{\theta}'^T} \right)^r \boldsymbol{\theta} \frac{\xi'^r}{r!} \\
&= \sum_{r=0}^{\infty} \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}^r \begin{bmatrix} 1 \\ 0 \end{bmatrix} \frac{\xi'^r}{r!} \\
&= \begin{bmatrix} \sum_{r=0}^{\infty} (-1)^r \frac{\xi'^{2r}}{(2r)!} \\ \sum_{r=0}^{\infty} (-1)^{r+1} \frac{\xi'^{2r+1}}{(2r+1)!} \end{bmatrix} \\
&= \begin{bmatrix} \cos(\xi') \\ -\sin(\xi') \end{bmatrix}.
\end{aligned}$$

This particular choice of  $\boldsymbol{\theta}$  produces a reparameterization  $\mathbf{g}_{\boldsymbol{\theta}}$  in agreement with the one chosen in (3.8). Any alternative reparameterization can be found utilizing an alternative transformation matrix  $\mathbf{S}(\boldsymbol{\theta}')$  or an alternative initialization.

Such an approach will only derive a local bijective map and convergence of the Taylor series may not result in a known functional form for any given constraint  $\mathbf{f}$ . When  $\mathbf{U}(\boldsymbol{\theta}')$  is not known as a function of  $\boldsymbol{\theta}'$ , numerical techniques are available to find the derivatives of  $\mathbf{U}(\boldsymbol{\theta}')$  with respect to  $\boldsymbol{\xi}'$  using the equation  $\mathbf{F}(\boldsymbol{\theta}')\mathbf{U}(\boldsymbol{\theta}') = \mathbf{0}$ .

### 3.1.5.2 A fixed point derivation

More commonly, a fixed point approach is taken to the derivation of an implicit function. As in appendix A.2, the parameter vector is partitioned  $\boldsymbol{\theta} = \begin{bmatrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \end{bmatrix}$  and the constraint function rewritten as  $\mathbf{f} : \mathbb{R}^{m-k} \times \mathbb{R}^k \rightarrow \mathbb{R}^k$  defined as  $\mathbf{f}^*(\boldsymbol{\theta}_1', \boldsymbol{\theta}_2') = \mathbf{f}(\begin{bmatrix} \boldsymbol{\theta}_1' \\ \boldsymbol{\theta}_2' \end{bmatrix})$ . If  $\mathbf{f}_{\boldsymbol{\theta}_2'}^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \frac{\partial}{\partial \boldsymbol{\theta}_2'^T} \mathbf{f}^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2') \Big|_{\boldsymbol{\theta}_2' = \boldsymbol{\theta}_2}$  is nonsingular, then there exist a

unique continuous function  $\boldsymbol{\theta}_2 = \boldsymbol{\theta}_2(\boldsymbol{\theta}_1)$  about  $\boldsymbol{\theta}_1$  such that  $\boldsymbol{f}^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2(\boldsymbol{\theta}_1)) = \mathbf{0}$ .

This is merely a particular variation of the implicit function theorem. One proof of this version proves the existence of a fixed point in the contraction map

$$\boldsymbol{d}(\boldsymbol{\theta}_2(\boldsymbol{\theta}'_1)) = \boldsymbol{\theta}_2(\boldsymbol{\theta}'_1) - \boldsymbol{D}^{-1}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \boldsymbol{f}(\boldsymbol{\theta}'_1, \boldsymbol{\theta}_2(\boldsymbol{\theta}'_1))$$

where  $\boldsymbol{D}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \boldsymbol{f}_{\boldsymbol{\theta}'_2}^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$ . This motivates the use of the iteration

$$\boldsymbol{\theta}_2^{(i+1)}(\boldsymbol{\theta}'_1) = \boldsymbol{\theta}_2^{(i)}(\boldsymbol{\theta}'_1) - \boldsymbol{D}^{-1}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \boldsymbol{f}(\boldsymbol{\theta}'_1, \boldsymbol{\theta}_2^{(i)}(\boldsymbol{\theta}'_1))$$

to generate the (fixed point) implicit function. The iteration is essentially an application of Newton's method [64].

**Example 3.22.** Reviewing example 3.4 again, we have the constraint  $f^*(\theta_1, \theta_2) = \theta_1^2 + \theta_2^2 - 1 = 0$  and wish to find a local reparameterization at  $(\theta_1, \theta_2) = (1, 0)$ . First note,

$$f_{\theta_1}^*(\theta_1, \theta_2) = 2\theta_1 = 2 \quad , \quad f_{\theta_2}^*(\theta_1, \theta_2) = 2\theta_2 = 0,$$

and we only have a nonsingularity with respect to  $\theta_1$ . Hence, we cannot find a function for  $\theta_2$  in terms of  $\theta_1$  using this approach at  $(1, 0)$ , but we can find  $\theta'_1(\theta'_2)$  there. Defining  $D(\theta_1, \theta_2) = f_{\theta_1}^*(\theta_1, \theta_2) = 2$  and initializing with  $\theta_1^{'(1)}(\theta'_2) = 1$ , the fixed point method dictates the next iterate to be

$$\begin{aligned} \theta_1^{'(2)}(\theta'_2) &= \theta_1^{'(1)} - D^{-1}(\theta_1, \theta_2) \left( (\theta_1^{'(1)}(\theta'_2))^2 + \theta_2'^2 - 1 \right) \\ &= 1 - \frac{1}{2} (1 + \theta_2'^2 - 1) = 1 - \frac{1}{2} \theta_2'^2. \end{aligned}$$

Continuing, the third iterate is

$$\begin{aligned} \theta_1^{'(3)}(\theta'_2) &= 1 - \frac{1}{2} \theta_2'^2 - \frac{1}{2} \left( \left( 1 - \frac{1}{2} \theta_2'^2 \right)^2 + \theta_2'^2 - 1 \right) \\ &= 1 - \frac{1}{2} \theta_2'^2 - \frac{1}{8} \theta_2'^4, \end{aligned}$$



and the fourth iterate is

$$\begin{aligned}\theta_1^{(4)}(\theta'_2) &= 1 - \frac{1}{2}\theta_2'^2 - \frac{1}{8}\theta_2'^4 - \frac{1}{2} \left( \left(1 - \frac{1}{2}\theta_2'^2 - \frac{1}{8}\theta_2'^4\right)^2 + \theta_2'^2 - 1 \right) \\ &= 1 - \frac{1}{2}\theta_2'^2 - \frac{1}{8}\theta_2'^4 - \frac{1}{16}\theta_2'^6 - \frac{1}{128}\theta_2'^8.\end{aligned}$$

As expected, as  $r \rightarrow \infty$  then  $\theta_1^{(r)}(\theta'_2)$  approaches a limit function

$$\theta_1'(\theta'_2) = \sqrt{1 - \theta_2'^2} = 1 - \frac{1}{2}\theta_2'^2 - \frac{1}{8}\theta_2'^4 - \frac{1}{16}\theta_2'^6 - \frac{5}{128}\theta_2'^8 + O(\theta_2'^9)$$

near  $\theta_2 = 0$ .

## 3.2 Identifiability

Identifiability conditions based on the CRB (or FIM) were detailed in Section 2.2, and the definitions of local and strong identifiability are given therein. In this section, the identifiability of parameters under functional equality constraints is considered.

### 3.2.1 Local identifiability

To establish a new identifiability criterion from the CCRB, we will first examine an existing criterion. Rothenberg [58, theorem 6] developed conditions for identifiability of a parameter vector  $\boldsymbol{\theta}$  under the constraints  $\mathbf{f}(\boldsymbol{\theta}) = 0$ , which was later partially re-derived by Crowder [18, lemma 1 is the “only if” portion of the theorem statement].

**Theorem 3.23** (Rothenberg-Crowder). Assume both  $\mathbf{F}(\boldsymbol{\theta}')$  and

$$\mathbf{M}(\boldsymbol{\theta}') = \begin{bmatrix} \mathbf{I}(\boldsymbol{\theta}') \\ \mathbf{F}(\boldsymbol{\theta}') \end{bmatrix}$$

have constant rank in a local neighborhood about  $\boldsymbol{\theta}$ . Then  $\boldsymbol{\theta}$  is locally identifiable if and only if  $\mathbf{M}(\boldsymbol{\theta})$  has full column rank  $m$ .

The proof of this theorem is based on the unconstrained proof of Theorem 2.5. An immediate implication is that if the FIM  $\mathbf{I}(\boldsymbol{\theta})$  is regular, then not only is the unconstrained model locally identifiable, but the constrained model is as well. Otherwise, if the FIM is singular, then the constraints must be such that the row vectors of its Jacobian  $\mathbf{F}(\boldsymbol{\theta})$  eliminate the null space of the FIM. In doing so, the constraints eliminate whatever inherent ambiguity was in the model that led to local unidentifiability and a singular FIM (Theorem 2.5). If row vectors of the Jacobian did not eliminate the null space of the Fisher information then there would exist a linear combination of column vectors of  $\mathbf{U}(\boldsymbol{\theta})$  such that  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  is singular, as shall be shown shortly.

Additionally, regardless of the information, in the trivial case where the constraints are such that the Jacobian is full column rank, then the theorem says the model is locally identifiable. Indeed, when  $\text{rank}(\mathbf{F}(\boldsymbol{\theta})) = m$  the constraints completely determine the parameter (e.g., see example 3.14).

Rothenberg's theorem for unconstrained identifiability is useful in establishing an alternative criterion for identifiability in relation to a component of the CCRB.

**Theorem 3.24.** Let  $\boldsymbol{\theta} \in \Theta_f$  and assume  $\mathbf{U}^T(\boldsymbol{\theta}')\mathbf{I}(\boldsymbol{\theta}')\mathbf{U}(\boldsymbol{\theta}')$  has constant rank in a neighborhood of  $\boldsymbol{\theta}$ . Then  $\boldsymbol{\theta}$  is locally identifiable if and only if  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  is regular.

*Proof.* Let  $\mathbf{g}_{\boldsymbol{\theta}}$  satisfy Theorem 3.3. If  $\mathbf{U}^T(\boldsymbol{\theta}')\mathbf{I}(\boldsymbol{\theta}')\mathbf{U}(\boldsymbol{\theta}')$  has constant rank in a

local neighborhood about  $\boldsymbol{\theta} = \mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi})$ . Then  $\tilde{\mathbf{I}}(\boldsymbol{\xi}')$  has constant rank in a local neighborhood about  $\boldsymbol{\xi}$ . And since  $\mathbf{g}_{\boldsymbol{\theta}}$  is injective,  $\boldsymbol{\theta}$  is locally identifiable in  $\Theta_f$  if and only if  $\boldsymbol{\xi}$  is locally identifiable in  $R^{m-k}$ . By Theorem 2.5,  $\boldsymbol{\xi}$  is locally identifiable if and only if  $\tilde{\mathbf{I}}(\boldsymbol{\xi})$  is regular. And  $\tilde{\mathbf{I}}(\boldsymbol{\xi})$  is regular if and only if  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  is regular.  $\square$

This is the corresponding theorem to Rothenberg's Theorem 2.5 and agrees with his Theorem 3.23, although the proof does not rely on this latter theorem's result because the implicit function  $\mathbf{g}_{\boldsymbol{\theta}}$  simplifies the approach. It is, however, possible to prove a more general result connecting the rank of the  $\mathbf{M}(\cdot)$  matrix of Theorem 3.23 to the implicit Fisher information  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$ .

**Theorem 3.25.** Assume  $\mathbf{U}^T(\boldsymbol{\theta}')\mathbf{I}(\boldsymbol{\theta}')\mathbf{U}(\boldsymbol{\theta}')$  has constant rank in a neighborhood of  $\boldsymbol{\theta}$ . Then

$$\text{nullity}(\mathbf{M}(\boldsymbol{\theta}')) = \text{nullity}(\mathbf{U}^T(\boldsymbol{\theta}')\mathbf{I}(\boldsymbol{\theta}')\mathbf{U}(\boldsymbol{\theta}')).$$

*Proof.* First, for some fixed  $\boldsymbol{\theta}'$ , assume  $\mathbf{M}(\boldsymbol{\theta}')$  is not full column rank and the vectors  $\mathbf{v}_1, \dots, \mathbf{v}_r$  are linearly independent and span the null space of  $\mathbf{M}(\boldsymbol{\theta}')$ , i.e., for  $l = 1, \dots, r$ ,

$$\mathbf{I}(\boldsymbol{\theta}')\mathbf{v}_l = \mathbf{0}$$

$$\mathbf{F}(\boldsymbol{\theta}')\mathbf{v}_l = \mathbf{0}.$$

Let  $\mathbf{U}(\boldsymbol{\theta}')$  be a matrix defined as in (3.6). Since  $\mathbf{F}(\boldsymbol{\theta}')\mathbf{v}_l = \mathbf{0}$  then each  $\mathbf{v}_l = \mathbf{U}(\boldsymbol{\theta}')\mathbf{w}_l$  for some  $\mathbf{w}_1, \dots, \mathbf{w}_r \in \mathbb{R}^{m-k}$ . Now if  $\sum_{l=1}^r \gamma_l \mathbf{w}_l = \mathbf{0}$ , then  $\sum_{l=1}^r \gamma_l \mathbf{v}_l = \sum_{l=1}^r \mathbf{U}(\boldsymbol{\theta}')\gamma_l \mathbf{w}_l = \mathbf{0}$ , which implies  $\gamma_1 = \dots = \gamma_r = 0$  since the  $\mathbf{v}_i$  are linearly independent. Hence,

the  $\mathbf{w}_1, \dots, \mathbf{w}_r$  are also linearly independent and span an  $r$ -dimensional subspace of  $\mathbb{R}^{m-k}$ . Thus,  $\mathbf{I}(\boldsymbol{\theta}')\mathbf{U}(\boldsymbol{\theta}')\mathbf{w}_l = \mathbf{0}$  for  $l = 1, \dots, r$ . This proves the implication, namely, that  $\text{nullity}(\mathbf{U}^T(\boldsymbol{\theta}')\mathbf{I}(\boldsymbol{\theta}')\mathbf{U}(\boldsymbol{\theta}')) \geq \text{nullity}(\mathbf{M}(\boldsymbol{\theta}'))$ .

To show the inverse, assume the vectors  $\mathbf{w}_1, \dots, \mathbf{w}_r$  are linearly independent and span the null space of the  $\mathbf{U}^T(\boldsymbol{\theta}')\mathbf{I}(\boldsymbol{\theta}')\mathbf{U}(\boldsymbol{\theta}')$ . Define  $\mathbf{v}_l = \mathbf{U}(\boldsymbol{\theta}')\mathbf{w}_l$  for  $l = 1, \dots, r$ . Since  $\mathbf{I}(\boldsymbol{\theta}')\mathbf{v}_l = \mathbf{0}$  and  $\mathbf{F}(\boldsymbol{\theta}')\mathbf{v}_l = \mathbf{0}$ , then  $\mathbf{M}(\boldsymbol{\theta}')\mathbf{v}_l = \mathbf{0}$  for  $l = 1, \dots, r$ . Note for any  $\gamma_1, \dots, \gamma_r$ , then  $\sum_{l=1}^r \gamma_l \mathbf{v}_l = \sum_{l=1}^r \mathbf{U}(\boldsymbol{\theta}')\gamma_l \mathbf{w}_l = \mathbf{0}$  if and only if  $\sum_{l=1}^r \gamma_l \mathbf{w}_l = \mathbf{0}$  (since  $\mathbf{U}(\boldsymbol{\theta}')$  is full column rank) which is true if and only if  $\gamma_1 = \dots = \gamma_r = 0$ . This proves the converse, i.e.,  $\text{nullity}(\mathbf{M}(\boldsymbol{\theta}')) \geq \text{nullity}(\mathbf{U}^T(\boldsymbol{\theta}')\mathbf{I}(\boldsymbol{\theta}')\mathbf{U}(\boldsymbol{\theta}'))$ .  $\square$

The theorem essentially states that in the local neighborhood where the implicit function  $\mathbf{g}$  is defined,  $\mathbf{M}(\boldsymbol{\theta})$  has full column rank if and only if  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  does, for any  $\boldsymbol{\theta} \in \Theta_f$ . As a consequence, we have the following corollary.

**Corollary 3.26.** If  $\mathbf{I}(\boldsymbol{\theta})$  is regular, then  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  is also regular. And if  $\boldsymbol{\theta}$  is locally identifiable in  $\Theta$ , then  $\boldsymbol{\theta}$  is locally identifiable in  $\Theta_f$ .

*Proof.* If  $\boldsymbol{\theta}$  is locally identifiable, then  $\text{nullity}(\mathbf{I}(\boldsymbol{\theta})) = 0$  by theorem 2.5. And if  $\text{nullity}(\mathbf{I}(\boldsymbol{\theta})) = 0$ , then  $\text{nullity}(\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})) = 0$  by theorem 3.25. Finally, if  $\text{nullity}(\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})) = 0$ , then  $\boldsymbol{\theta}$  is locally identifiable in  $\Theta_f$  by theorem 3.24.  $\square$

### 3.2.1.1 Local identifiability in the Aitchison-Silvey-Crowder CCRB formula

In subsection 3.1.2.2, an alternative form of the CCRB is presented where a loaded Fisher information is used in place of the FIM to resolve issues of singularity of the FIM. The Aitchison-Silvey loaded FIM is  $\mathbf{I}(\boldsymbol{\theta}) + \mathbf{F}^T(\boldsymbol{\theta})\mathbf{F}(\boldsymbol{\theta})$  [4], whereas the Crowder loaded FIM  $\mathbf{I}(\boldsymbol{\theta}) + \mathbf{F}^T(\boldsymbol{\theta})\mathbf{K}\mathbf{F}(\boldsymbol{\theta})$  [18], where  $\mathbf{K}$  is chosen such that the loaded FIM is full rank. The following theorem, which connects local identifiability and the Aitchison-Silvey-Crowder CCRB, shows that under certain conditions, the matrix  $\mathbf{K}$  is unnecessary. This result is hinted at, but not clearly stated in [18, lemma 6].

**Theorem 3.27.** The Aitchison-Silvey-Crowder loaded FIM  $\mathbf{I}(\boldsymbol{\theta}) + \mathbf{F}^T(\boldsymbol{\theta})\mathbf{F}(\boldsymbol{\theta})$  is nonsingular if and only if  $\mathbf{M}(\boldsymbol{\theta})$  is full column rank. Hence, if  $\mathbf{I}(\boldsymbol{\theta}') + \mathbf{F}^T(\boldsymbol{\theta}')\mathbf{F}(\boldsymbol{\theta}')$  is constant locally about  $\boldsymbol{\theta}$ , then  $\boldsymbol{\theta}$  is identifiable if and only if  $\mathbf{I}(\boldsymbol{\theta}) + \mathbf{F}^T(\boldsymbol{\theta})\mathbf{F}(\boldsymbol{\theta})$  is nonsingular.

*Proof.* To show the contrapositive, assume  $\mathbf{M}(\boldsymbol{\theta})$  is not full column rank and  $\mathbf{v}$  is a nontrivial vector such that  $\mathbf{M}(\boldsymbol{\theta})\mathbf{v} = \mathbf{0}$ . Then  $\mathbf{I}(\boldsymbol{\theta})\mathbf{v} = \mathbf{0}$  and  $\mathbf{F}(\boldsymbol{\theta})\mathbf{v} = \mathbf{0}$ . Therefore,  $\mathbf{v}$  is in the null space of the Gram matrix  $\mathbf{I}(\boldsymbol{\theta}) + \mathbf{F}^T(\boldsymbol{\theta})\mathbf{F}(\boldsymbol{\theta})$ . To show the inverse, assume  $(\mathbf{I}(\boldsymbol{\theta}) + \mathbf{F}^T(\boldsymbol{\theta})\mathbf{F}(\boldsymbol{\theta}))\mathbf{v} = \mathbf{0}$  for some nontrivial vector  $\mathbf{v} \in \mathbb{R}^m$ . Then since

$$\mathbf{I}(\boldsymbol{\theta}) + \mathbf{F}^T(\boldsymbol{\theta})\mathbf{F}(\boldsymbol{\theta}) = [\mathbf{I}^{1/2}(\boldsymbol{\theta}) \quad \mathbf{F}^T(\boldsymbol{\theta})] \cdot \begin{bmatrix} \mathbf{I}^{1/2}(\boldsymbol{\theta}) \\ \mathbf{F}(\boldsymbol{\theta}) \end{bmatrix},$$

we must have that  $\mathbf{I}^{1/2}(\boldsymbol{\theta})\mathbf{v} = \mathbf{0}$  and  $\mathbf{F}(\boldsymbol{\theta})\mathbf{v} = \mathbf{0}$ . Hence,  $\mathbf{I}(\boldsymbol{\theta})\mathbf{v} = \mathbf{I}^{1/2}(\boldsymbol{\theta})\mathbf{I}^{1/2}(\boldsymbol{\theta})\mathbf{v} = \mathbf{0}$ , which implies  $\mathbf{M}(\boldsymbol{\theta})\mathbf{v} = \mathbf{0}$  and  $\mathbf{M}(\boldsymbol{\theta})$  is not full column rank. □

### 3.2.2 Strong Identifiability

Restricting the discussion to normal distributions in this section we can extend the equivalence criterion between regularity and strong identifiability as defined in section 2.2.2. That is, assume  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}(\boldsymbol{\theta}), \boldsymbol{\Sigma}(\boldsymbol{\theta}))$  with  $\boldsymbol{\mu}(\boldsymbol{\theta}) \in \mathbb{R}^p$ , with the elements of the mean and variance explicitly defined by a map  $\boldsymbol{\varphi} : \Theta \rightarrow \mathbb{R}^q$  where  $q \leq p + p(p+1)/2$  and assume  $m \leq q$ .

**Theorem 3.28.** If  $\boldsymbol{\theta}$  is strongly identifiable on  $\Theta$  and  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$ , then  $\boldsymbol{\theta}$  is also strongly identifiable on  $\Theta_f$ .

*Proof.* If  $\boldsymbol{\theta}$  is strongly identifiable, then there exists a representative mapping  $\boldsymbol{\varphi}^*$ , which is injective on  $\Theta$ . Since  $\Theta_f \subset \Theta$ ,  $\boldsymbol{\varphi}^*$  is still injective and a representative mapping on  $\Theta_f$ . □

This theorem is complementary to Corollary 3.26. Essentially, the imposition of constraints does not take away existing identifiability (local or strong) or Fisher information regularity that already exists in a model. However, it is not always the case that the original (unconstrained) model is information regular or identifiable. The following theorem, an extension of Theorem 2.6, connects the notion of strong identifiability with regularity of  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$ .

**Theorem 3.29.** Assume  $\boldsymbol{\varphi}$  is a holomorphic mapping of  $\mathbf{z} \in \Omega \subset \cup_{\alpha \in A} \Omega_\alpha$  into  $\mathbb{C}^q$ , where  $\Theta_f \subset \Omega \subset \mathbb{C}^m$  and  $\Omega_\alpha$  is open in  $\mathbb{C}^m$  for each  $\alpha$ . Then

- (a) if  $\mathbf{U}^T(\mathbf{z})\mathbf{I}(\mathbf{z})\mathbf{U}(\mathbf{z})$  is regular, there exists a strongly identifiable open neighborhood about  $\mathbf{z}$ , and

(b) if there exists a representative mapping  $\varphi_\alpha^* : \Omega_\alpha \rightarrow \mathbb{C}^q$  for each  $\alpha$ , then the matrix  $\mathbf{U}^T(\mathbf{z})\mathbf{I}(\mathbf{z})\mathbf{U}(\mathbf{z})$  is regular for every  $\mathbf{z} \in \Omega$ .

*Proof.* By the implicit function theorem (Theorem 3.3), then for any matrix  $\mathbf{U}(\boldsymbol{\theta})$  whose columns form a basis for the null space of the Jacobian of  $\mathbf{f}(\boldsymbol{\theta})$ , there exists an open set  $\mathbb{V} \ni \boldsymbol{\theta}$ , an open set  $\mathbb{W} \subset \mathbb{R}^{m-k}$ , and some transformation  $\mathbf{g}_\theta : \mathbb{W} \rightarrow \mathbb{R}^m$  such that  $\boldsymbol{\theta} = \mathbf{g}_\theta(\boldsymbol{\xi})$  for some  $\boldsymbol{\xi} \in \mathbb{W}$ . In particular, in this reduced parameter space, there exists a FIM such that  $\tilde{\mathbf{I}}(\boldsymbol{\xi}) = \mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$ . Since Theorem 2.6 applies to  $\tilde{\mathbf{I}}(\boldsymbol{\xi})$ , the result for  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  is proven.  $\square$

Regardless of the regularity of the FIM, only regularity of  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  determines strict identifiability under constraints for normal distributions, given a proper holomorphic function(s).

### 3.3 Linear Model

Assume the observations  $\mathbf{x}$  model a linear function of the parameters  $\boldsymbol{\theta}$ , as in

$$\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w}, \quad (3.23)$$

where  $\mathbf{H}$  is a full column rank  $n \times m$  observation matrix consisting of known elements and  $\mathbf{w}$  is a random noise vector with mean zero and known variance  $\mathbf{C}$ . As noted in section 2.3.1, the (weighted) LSE

$$\hat{\boldsymbol{\theta}}_{\text{LS}}(\mathbf{x}) = (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}^{-1} \mathbf{x}$$

is the BLUE. The LSE has a variance of  $\mathbf{Q}^{-1}$  where  $\mathbf{Q} = (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})$ . In addition, we assume a linear constraint

$$\mathbf{f}(\boldsymbol{\theta}) = \mathbf{F}\boldsymbol{\theta} + \mathbf{v} = \mathbf{0},$$

where  $\mathbf{F}$  is a known full row rank  $k \times m$  projection matrix and  $\mathbf{v}$  is a known shift vector. For a linear constraint, the Jacobian  $\mathbf{F}(\boldsymbol{\theta}) = \mathbf{F}$  does not depend on the parameter. As the linear problem is well-studied, many of the results in this section are known (e.g., see [61, section 3.8] and [57, section 11.3.3]) but are presented here from a different perspective.

### 3.3.1 Best Linear Unbiased Estimation

The constrained (weighted) LSE (CLSE) is most often given by [36, p. 252]

$$\hat{\boldsymbol{\theta}}_{\text{CLS}}(\mathbf{x}) = \hat{\boldsymbol{\theta}}_{\text{LS}}(\mathbf{x}) - \mathbf{Q}^{-1} \mathbf{F}^T (\mathbf{F} \mathbf{Q}^{-1} \mathbf{F}^T)^{-1} (\mathbf{F} \hat{\boldsymbol{\theta}}_{\text{LS}}(\mathbf{x}) + \mathbf{v}). \quad (3.24)$$

Simple calculation confirms that this CLSE exists in  $\ker(\mathbf{f}) = \Theta_f$  (i.e., it satisfies the constraint), is unbiased, and has variance

$$\text{Var}(\hat{\boldsymbol{\theta}}_{\text{CLS}}(\mathbf{x})) = \mathbf{Q}^{-1} - \mathbf{Q}^{-1} \mathbf{F}^T (\mathbf{F} \mathbf{Q}^{-1} \mathbf{F}^T)^{-1} \mathbf{F} \mathbf{Q}^{-1}. \quad (3.25)$$

The BLUE property of the LSE is preserved for the CLSE. This is not a surprising result since a linear constraint on a linear model is still a linear model, as seen in figure 3.2.

In the context of the null space approach of the CCRB, an alternative CLSE may be developed. One particular advantage of this approach is the avoidance of



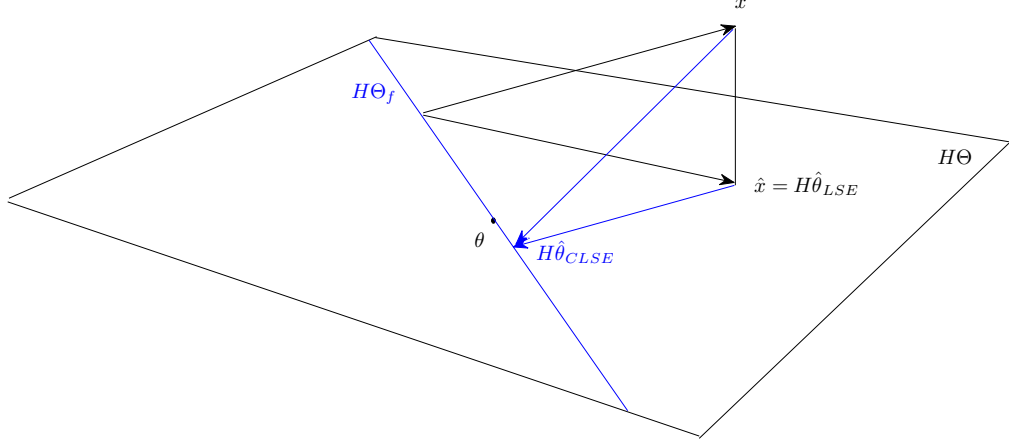


Figure 3.2: Projection of the observations  $\mathbf{x}$  onto the linear space  $\mathbf{H}\Theta$  and the linear constraint space  $\mathbf{H}\Theta_f$ .

the need for directly using and solving for Lagrange multipliers. Note that zero solutions of  $\mathbf{f}(\boldsymbol{\theta})$  are of the form

$$\boldsymbol{\theta} = -\mathbf{F}^T (\mathbf{F}\mathbf{F}^T)^{-1} \mathbf{v} + \mathbf{U}\boldsymbol{\xi} = \mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi})$$

where  $\mathbf{U}$  satisfies the equations in (3.6), i.e., the columns of  $\mathbf{U}$  form an orthonormal basis for the null space of the row vectors of  $\mathbf{F}$ , and  $\boldsymbol{\xi} \in \mathbb{R}^{m-k}$  is a parameter representing the projection of  $\boldsymbol{\theta}$  to the constraint space  $\Theta_f$ . Since the Jacobian  $\mathbf{F}$  is independent of the parameter, then  $\mathbf{U}(\boldsymbol{\theta}) = \mathbf{U}$  and, hence,  $\mathbf{g}_{\boldsymbol{\theta}}$  are as well. Moreover, the typical local properties for the implicit function hold globally. Substituting this solution for  $\boldsymbol{\theta}$  the linear model is reformulated

$$\mathbf{y} = \mathbf{H}\mathbf{U}\boldsymbol{\xi} + \mathbf{w} \quad (3.26)$$

where  $\mathbf{y} = \mathbf{x} + \mathbf{H}\mathbf{F}^T (\mathbf{F}\mathbf{F}^T)^{-1} \mathbf{v}$ . Following the least squares result for the CRB,

we desire a solution that minimizes a quadratic objective function

$$\hat{\boldsymbol{\xi}}_{\text{LS}} = \arg \min_{\boldsymbol{\xi}} (\mathbf{y} - \mathbf{H}\mathbf{U}\boldsymbol{\xi})^T \mathbf{C}^{-1} (\mathbf{y} - \mathbf{H}\mathbf{U}\boldsymbol{\xi}).$$

This solution must satisfy the normal equations given by

$$(\mathbf{U}^T \mathbf{Q} \mathbf{U}) \hat{\boldsymbol{\xi}}_{\text{LS}}(\mathbf{y}) = \mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} \mathbf{y}.$$

Provided the normalizing matrix  $(\mathbf{U}^T \mathbf{Q} \mathbf{U})$  is full rank, the LSE of  $\boldsymbol{\xi}$  is given by

$$\hat{\boldsymbol{\xi}}_{\text{LS}}(\mathbf{y}) = (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} \mathbf{y}$$

and is the BLUE ((2.4) in section 2.3.1). The corresponding LSE of  $\boldsymbol{\theta}$  based on this null space approach is

$$\begin{aligned} \hat{\boldsymbol{\theta}}_{\text{CLS}}(\mathbf{x}) &= -\mathbf{F}^T (\mathbf{F} \mathbf{F}^T)^{-1} \mathbf{v} + \mathbf{U} \hat{\boldsymbol{\xi}}_{\text{LS}}(\mathbf{y}) \\ &= -\mathbf{F}^T (\mathbf{F} \mathbf{F}^T)^{-1} \mathbf{v} + \mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} + \mathbf{H} \mathbf{F}^T (\mathbf{F} \mathbf{F}^T)^{-1} \mathbf{v}) \\ &= \mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{Q} (\boldsymbol{\theta}_1 + \mathbf{F}^T (\mathbf{F} \mathbf{F}^T)^{-1} \mathbf{v}) - \mathbf{F}^T (\mathbf{F} \mathbf{F}^T)^{-1} \mathbf{v} \\ &\quad + \mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H} \boldsymbol{\theta}_1) \\ &= \boldsymbol{\theta}_1 + \mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H} \boldsymbol{\theta}_1) \end{aligned} \tag{3.27}$$

where  $\boldsymbol{\theta}_1 = -\mathbf{F}^T (\mathbf{F} \mathbf{F}^T)^{-1} \mathbf{v} + \mathbf{U} \boldsymbol{\xi}_1$  can be any arbitrary point satisfying the linear constraint ( $\boldsymbol{\xi}_1$  is unrestricted). This alternative CLSE is more general than the prior formula as it is applicable in scenarios when  $\mathbf{H}$  and  $\mathbf{F}$  are not necessarily full column rank and full row rank, respectively. As such, it has the more general expression for its variance

$$\mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T,$$

which is equivalent to (3.25) when  $\mathbf{H}$  and  $\mathbf{F}$  are full column rank and full row rank, respectively. Replacing the rank conditions on  $\mathbf{H}$  and  $\mathbf{F}$ , the weaker necessary conditions for this CLSE are that  $\mathbf{H}\mathbf{U}$  be full column rank. If the stronger necessary conditions on  $\mathbf{H}$  and  $\mathbf{F}$  exist, then the CLSE can be reformulated in terms of the LSE as before, e.g.,

$$\hat{\boldsymbol{\theta}}_{\text{CLS}}(\mathbf{x}) = \boldsymbol{\theta}_1 + \mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{Q} (\hat{\boldsymbol{\theta}}_{\text{LS}}(\mathbf{x}) - \boldsymbol{\theta}_1).$$

It is easy to confirm that  $\hat{\boldsymbol{\theta}}_{\text{CLS}}(\mathbf{x}) \in \ker(\mathbf{f})$  by recalling that  $\mathbf{F}\mathbf{U} = \mathbf{0}$  and  $\boldsymbol{\theta}_1 \in \ker(\mathbf{f})$ .

Also, as in section 2.3.1, if  $\mathbf{H}\mathbf{U}$  is not full column rank, then for estimable functions  $\mathbf{d}^T \boldsymbol{\theta}$ , i.e., for vectors  $\mathbf{d}$  in the column space of  $\mathbf{U}^T \mathbf{H}^T$ , the LSE is BLUE and is given by

$$\mathbf{d}^T \hat{\boldsymbol{\theta}}_{\text{CLS}}(\mathbf{x}) = \mathbf{d}^T \boldsymbol{\theta}_1 + \mathbf{d}^T \mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^\dagger \mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H} \boldsymbol{\theta}_1) \quad (3.28)$$

similar to (3.27) with variance  $\mathbf{d}^T \mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^\dagger \mathbf{U}^T \mathbf{d}$ .

### 3.3.2 Uniform Minimum Variance Estimation under Gaussian noise

Under the assumption that the noise is normally distributed, the (unconstrained) LSE is also the MLE. The FIM<sup>12</sup> is  $\mathbf{I}(\boldsymbol{\theta}) = \mathbf{Q} = (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})$ , and the LSE/MLE is the MVUE being efficient with respect to  $\text{CRB} = \mathbf{I}^{-1}(\boldsymbol{\theta})$ . Given this general principle that for linear models with additive Gaussian noise, the LSE is the MLE, then since a linear constraint is essentially a reduced dimensional linear

---

<sup>12</sup>Note that neither the Fisher information, the Jacobian of the constraints, nor the null space matrix depend on the parameters in the linear model with linear constraints.

model as evidenced in (3.26), the CLSE should be the constrained MLE (CMLE).

As we shall see, this is indeed the case.

First, for the Gaussian linear model in (3.26), the pdf is

$$q(\mathbf{y}; \boldsymbol{\xi}) = \frac{1}{(2\pi)^{(m-k)/2}(\det \mathbf{C})^{(m-k)/2}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{H}\mathbf{U}\boldsymbol{\xi})^T \mathbf{C}^{-1} (\mathbf{y} - \mathbf{H}\mathbf{U}\boldsymbol{\xi}) \right\}.$$

The Fisher score  $\tilde{\mathbf{s}}(\mathbf{y}; \boldsymbol{\xi}) = -\mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{y} - \mathbf{H}\mathbf{U}\boldsymbol{\xi})$  has a variance, or Fisher information, of  $\tilde{\mathbf{I}}(\boldsymbol{\xi}) = \mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} \mathbf{H}\mathbf{U}$ . Hence the CCRB is

$$\begin{aligned} \text{CCRB}(\boldsymbol{\theta}) &= \mathbf{G}(\boldsymbol{\xi}) \tilde{\mathbf{I}}^{-1}(\boldsymbol{\xi}) \mathbf{G}^T(\boldsymbol{\xi}) \\ &= \mathbf{U} (\mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} \mathbf{H}\mathbf{U})^{-1} \mathbf{U}^T \end{aligned}$$

(see example 3.20).

Maximizing the likelihood (or pdf) is equivalent to minimizing the quadratic, therefore the LSE of  $\boldsymbol{\xi}$  is also the MLE of  $\boldsymbol{\xi}$ . And by the invariance property [36, 14] of the MLE, then the CMLE of  $\boldsymbol{\theta}$  is

$$\hat{\boldsymbol{\theta}}_{\text{CML}}(\mathbf{x}) = \mathbf{g}(\hat{\boldsymbol{\xi}}_{\text{ML}}(\mathbf{y})) = \hat{\boldsymbol{\theta}}_{\text{CLS}}(\mathbf{x}).$$

**Theorem 3.30.** The CMLE is optimal for the linear model under linear constraints.

That is, if the observations obey the linear model in (3.23), where  $\mathbf{H}$  is a known matrix,  $\boldsymbol{\theta}$  is an unknown parameter vector subject to the linear constraint  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{F}\boldsymbol{\theta} + \mathbf{v} = \mathbf{0}$ , and  $\mathbf{w}$  is a zero-mean normal random vector with known variance  $\mathbf{C}$ , then provided  $\mathbf{H}\mathbf{U}$  is full column rank, where  $\mathbf{U}$  is defined by (3.6), the CMLE

$$\hat{\boldsymbol{\theta}}_{\text{CML}}(\mathbf{x}) = \boldsymbol{\theta}_1 + \mathbf{U} (\mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} \mathbf{H}\mathbf{U})^{-1} \mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H}\boldsymbol{\theta}_1) \quad (3.29)$$

is unbiased and efficient.

While the formula for the CMLE based on the MLE might seem preferable, analogous to the formula for the CLSE based on the LSE, i.e.,

$$\hat{\boldsymbol{\theta}}_{\text{CML}}(\boldsymbol{x}) = \boldsymbol{\theta}_1 + \boldsymbol{U} (\boldsymbol{U}^T \boldsymbol{Q} \boldsymbol{U})^{-1} \boldsymbol{U}^T \boldsymbol{Q} (\hat{\boldsymbol{\theta}}_{\text{ML}}(\boldsymbol{x}) - \boldsymbol{\theta}_1), \quad (3.30)$$

with  $\boldsymbol{Q} = \boldsymbol{H}^T \boldsymbol{C}^{-1} \boldsymbol{H}$ , this formulation requires the existence of a full rank FIM  $\boldsymbol{I}(\boldsymbol{\theta}) = \boldsymbol{Q}$  in the MLE. The benefit of using the CMLE in (3.29) versus the CMLE in (3.30) is that the following proof does not require this assumption.

*Proof.* First note for any  $\boldsymbol{\theta}, \boldsymbol{\theta}_1$  satisfying the constraints

$$\begin{aligned} \boldsymbol{\theta} - \boldsymbol{\theta}_1 &= -\boldsymbol{F}^T (\boldsymbol{F} \boldsymbol{F}^T)^{-1} \boldsymbol{v} + \boldsymbol{U} \boldsymbol{\xi} + \boldsymbol{F}^T (\boldsymbol{F} \boldsymbol{F}^T)^{-1} \boldsymbol{v} - \boldsymbol{U} \boldsymbol{\xi}_1 \\ &= \boldsymbol{U} (\boldsymbol{\xi} - \boldsymbol{\xi}_1) \end{aligned}$$

for some  $\boldsymbol{\xi}, \boldsymbol{\xi}_1 \in \mathbb{R}^{m-k}$ . Therefore, the expected value of the CMLE is

$$\begin{aligned} E_{\boldsymbol{\theta}} \hat{\boldsymbol{\theta}}_{\text{CML}}(\boldsymbol{x}) &= \boldsymbol{\theta}_1 + \boldsymbol{U} (\boldsymbol{U}^T \boldsymbol{Q} \boldsymbol{U})^{-1} \boldsymbol{U}^T \boldsymbol{Q} (\boldsymbol{\theta} - \boldsymbol{\theta}_1) \\ &= \boldsymbol{\theta}_1 + \boldsymbol{U} (\boldsymbol{U}^T \boldsymbol{Q} \boldsymbol{U})^{-1} \boldsymbol{U}^T \boldsymbol{Q} \boldsymbol{U} (\boldsymbol{\xi} - \boldsymbol{\xi}_1) \\ &= \boldsymbol{\theta}_1 + \boldsymbol{U} (\boldsymbol{\xi} - \boldsymbol{\xi}_1) \\ &= \boldsymbol{\theta}_1 + \boldsymbol{\theta} - \boldsymbol{\theta}_1 \\ &= \boldsymbol{\theta}. \end{aligned}$$

Finally, the variance of the CMLE is

$$\begin{aligned}
\text{Var}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}_{\text{CML}}(\mathbf{x})) &= \text{Var}_{\boldsymbol{\theta}}\left(\boldsymbol{\theta}_1 + \mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H} \boldsymbol{\theta}_1)\right) \\
&= \text{Var}_{\boldsymbol{\theta}}\left(\mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H} \boldsymbol{\theta} + \mathbf{H} \boldsymbol{\theta} - \mathbf{H} \boldsymbol{\theta}_1)\right) \\
&= \text{Var}_{\boldsymbol{\theta}}\left(\mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H} \boldsymbol{\theta})\right) \\
&= \mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} \mathbf{C} \mathbf{C}^{-1} \mathbf{H} \mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T \\
&= \mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{Q} \mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T \\
&= \mathbf{U} (\mathbf{U}^T \mathbf{Q} \mathbf{U})^{-1} \mathbf{U}^T,
\end{aligned}$$

i.e., the CCRB of  $\boldsymbol{\theta}$ . □

Thus, the CMLE is the MVUE for the linear model with linear constraints under a Gaussian assumption.

Additionally, when  $\mathbf{H} \mathbf{U}$  is not full column rank, then when  $\mathbf{d}$  is in the column space of  $\mathbf{U}^T \mathbf{H}^T$ , the MLE of  $\mathbf{d}^T \boldsymbol{\theta}$  is still the MVUE and is given by  $\mathbf{d}^T \hat{\boldsymbol{\theta}}_{\text{CML}}(\mathbf{x}) = \mathbf{d}^T \hat{\boldsymbol{\theta}}_{\text{CLS}}(\mathbf{x})$  from (3.28), with variance  $\mathbf{d}^T \text{CCRB}(\boldsymbol{\theta}) \mathbf{d}$ .

### 3.4 Constrained Maximum Likelihood Estimation

The constrained MLE (CMLE) of the parameter vector  $\boldsymbol{\theta}$  constrained to the manifold  $\Theta_f = \{\boldsymbol{\theta}' : \mathbf{f}(\boldsymbol{\theta}') = \mathbf{0}\}$  is the estimator in  $\Theta_f$ , which given the observations  $\mathbf{x}$ , that maximizes the likelihood distribution  $p(\mathbf{x}; \cdot)$ , i.e. it is the maximum likelihood in  $\Theta_f$ . Since  $\log(\cdot)$  is concave, it is convenient to equivalently maximize the log-likelihood  $\log p(\mathbf{x}; \cdot)$  since then the Jacobian of the objective is the Fisher score. In an optimization context, the CMLE, which will be denoted  $\hat{\boldsymbol{\theta}}_{\text{CML}}(\mathbf{x})$ , is

the solution to the following constrained optimization problem:

$$\begin{aligned} \max_{\boldsymbol{\theta}'} \quad & \log p(\mathbf{x}; \boldsymbol{\theta}') \\ \text{s.t.} \quad & \mathbf{f}(\boldsymbol{\theta}') = \mathbf{0}. \end{aligned} \quad (3.31)$$

Analogous to the method of maximum likelihood approach of (2.7), solutions  $\hat{\boldsymbol{\theta}}(\mathbf{x}) = \mathbf{g}_{\hat{\boldsymbol{\theta}}(\mathbf{x})}(\hat{\boldsymbol{\xi}}(\mathbf{x}))$  satisfying

$$\left. \frac{\partial}{\partial \boldsymbol{\xi}'} \log p(\mathbf{x}, \mathbf{g}_{\boldsymbol{\theta}'}(\boldsymbol{\xi}')) \right|_{\boldsymbol{\xi}' = \hat{\boldsymbol{\xi}}(\mathbf{x})} = \mathbf{0},$$

where  $\mathbf{g}_{\boldsymbol{\theta}'}$  is defined by (3.7), are candidates to be the CMLE. More formally, a solution to this optimization problem must satisfy the Karush-Kuhn-Tucker conditions [45], i.e.,

$$\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}') - \boldsymbol{\lambda}^T \mathbf{F}(\boldsymbol{\theta}') = \mathbf{0} \quad (\text{stationarity}) \quad (3.32)$$

$$\mathbf{f}(\boldsymbol{\theta}') = \mathbf{0}. \quad (\text{feasibility}) \quad (3.33)$$

Any point satisfying these conditions is a *stationary, feasible point*.

Since  $\hat{\boldsymbol{\theta}}_{\text{CML}}(\mathbf{x}) \in \Theta_f$ , then the implicit function theorem implies there exists an open set  $\mathbb{O} \subset \Theta_f$  containing  $\hat{\boldsymbol{\theta}}_{\text{CML}}(\mathbf{x})$ , an open set  $\mathbb{P} \subset \mathbb{R}^{m-k}$ , and a continuously differentiable bijection  $\mathbf{g}_{\hat{\boldsymbol{\theta}}(\mathbf{x})} : \mathbb{P} \rightarrow \mathbb{O}$  such that  $\hat{\boldsymbol{\theta}}(\mathbf{x}) = \mathbf{g}_{\hat{\boldsymbol{\theta}}(\mathbf{x})}(\hat{\boldsymbol{\xi}}(\mathbf{x}))$  for some  $\hat{\boldsymbol{\xi}}(\mathbf{x}) \in \mathbb{P}$ . If  $\hat{\boldsymbol{\theta}}(\mathbf{x})$  is a maximizer of the likelihood function  $p(\mathbf{x}; \boldsymbol{\theta})$  in the constraint set  $\Theta_f$ , then the likelihood  $q(\mathbf{x}; \boldsymbol{\xi}') = p(\mathbf{x}; \mathbf{g}_{\hat{\boldsymbol{\theta}}(\mathbf{x})}(\boldsymbol{\xi}'))$  has a maximum at  $\hat{\boldsymbol{\xi}}(\mathbf{x})$  in  $\mathbb{P}$  (i.e., a local maximum at  $\hat{\boldsymbol{\xi}}(\mathbf{x})$  in  $\mathbb{R}^{m-k}$ ). It cannot be said that  $q(\mathbf{x}; \boldsymbol{\xi})$  has a global maximum at  $\hat{\boldsymbol{\xi}}(\mathbf{x})$ , since  $\boldsymbol{\theta}' = \mathbf{g}_{\hat{\boldsymbol{\theta}}(\mathbf{x})}(\boldsymbol{\xi}')$  is only guaranteed to exist in  $\Theta_f$  when  $\boldsymbol{\xi}' \in \mathbb{P}$ , i.e., there may exist a point  $\boldsymbol{\xi}' \in \mathbb{R}^{m-k} \setminus \mathbb{P}$  such that  $q(\mathbf{x}; \boldsymbol{\xi}') > q(\mathbf{x}; \hat{\boldsymbol{\xi}}(\mathbf{x}))$  and  $\mathbf{g}_{\hat{\boldsymbol{\theta}}(\mathbf{x})}(\boldsymbol{\xi}') \notin \Theta_f$ .

### 3.4.1 Efficient estimation

In Section 2.4.1, it was explained that when an efficient estimator exists, the method of maximum likelihood finds the estimator [36, 62]. It is useful to note the connection between efficiency and the method of constrained maximum likelihood since Stoica and Ng ignored this extension in their paper [68], despite Marzetta having showed that this result extends to the constrained case when the FIM is non-singular [47, theorem 3]. What follows is the general extension of this result, including the case for singular FIMs.

**Theorem 3.31.** If  $\mathbf{t}(\mathbf{x})$  is a constrained estimator of  $\boldsymbol{\theta}$ , required to satisfy the constraint  $\mathbf{f}(\mathbf{t}(\mathbf{x})) = \mathbf{0}$ , which is also efficient with respect to the CCRB, then the estimator is a stationary point for the constrained optimization problem in (3.31).

*Proof.* This is perhaps more easily proven strictly from the constrained parameter perspective, since the global maximum of the likelihood relative to the implicit reparameterization may not correspond to global maximum in  $\Theta_f$  relative to the constrained parameterization. Since  $\mathbf{t}(\mathbf{x})$  is efficient then in the mean-square sense we have  $\mathbf{t}(\mathbf{x}) - \boldsymbol{\theta} = \text{CCRB}(\boldsymbol{\theta})\mathbf{s}(\mathbf{x}; \boldsymbol{\theta})$  as a function of  $\boldsymbol{\theta}$ . Then as  $\boldsymbol{\theta} \rightarrow \mathbf{t}(\mathbf{x})$  (this assumes the observations are consistent with  $\boldsymbol{\theta}$ ) we have

$$\mathbf{0} \leftarrow \mathbf{t}(\mathbf{x}) - \boldsymbol{\theta} = \text{CCRB}(\boldsymbol{\theta})\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}).$$

The continuity of the CCRB and the Fisher score implies  $\mathbf{s}(\mathbf{x}; \mathbf{t}(\mathbf{x})) = \mathbf{F}^T(\mathbf{t}(\mathbf{x})) \cdot \boldsymbol{\lambda}$  for some  $\boldsymbol{\lambda} \in \mathbb{R}^k$  or

$$\mathbf{s}(\mathbf{x}; \mathbf{t}(\mathbf{x})) - \boldsymbol{\lambda}^T \mathbf{F}(\mathbf{t}(\mathbf{x})) = \mathbf{0},$$



which defines the stationarity condition (3.32) of the constrained optimization problem with  $\boldsymbol{\lambda}$  being the vector of Lagrange multipliers.  $\square$

### 3.4.2 Asymptotic Normality

The asymptotic properties of the MLE can be found in section 2.4.2. Therein, it was mentioned that the maximum likelihood estimator was asymptotically unbiased and efficient with variance asymptotically equivalent to the CRB. A corresponding relationship exists between the CMLE and the CCRB. As before, let the samples  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ , be iid as  $\mathbf{x}$  from the likelihood  $p(\mathbf{x}; \boldsymbol{\theta})$ , where  $\boldsymbol{\theta}$  is assumed to exist in  $\Theta_f$ . Denote  $\mathbf{y}_n = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$  to be the collection of these samples, so that the likelihood will be  $p(\mathbf{y}_n; \boldsymbol{\theta}) = \prod_{i=1}^n p(\mathbf{x}_i; \boldsymbol{\theta})$ . Hence, the asymptotic CMLE will be denoted  $\hat{\boldsymbol{\theta}}(\mathbf{y}_n)$ .

**Theorem 3.32.** Assuming the pdf satisfies the regularity conditions (see (3.19) and discussion after proof), then the CMLE is asymptotically distributed according to

$$\sqrt{n} \left( \hat{\boldsymbol{\theta}}(\mathbf{y}_n) - \boldsymbol{\theta} \right) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \text{CCRB}(\boldsymbol{\theta})) .$$

There exists a number of results in the literature regarding the asymptotic characteristics of the CMLE (e.g., the works of Aitchison and Silvey [3, 63, 4, 2], and of Crowder [18]). For example, Crowder shows that

$$\sqrt{n} \left( \hat{\boldsymbol{\theta}}(\mathbf{y}_n) - \boldsymbol{\theta} \right) \xrightarrow{d} \mathcal{N} \left( \mathbf{0}, \mathbf{D}^{-1}(\boldsymbol{\theta}) - \mathbf{D}^{-1}(\boldsymbol{\theta}) \mathbf{F}^T(\boldsymbol{\theta}) \left( \mathbf{F}(\boldsymbol{\theta}) \mathbf{D}^{-1}(\boldsymbol{\theta}) \mathbf{F}^T(\boldsymbol{\theta}) \right)^{-1} \mathbf{F}(\boldsymbol{\theta}) \mathbf{D}^{-1}(\boldsymbol{\theta}) \right)$$

where  $\mathbf{D}(\boldsymbol{\theta}) = \mathbf{I}(\boldsymbol{\theta}) + \mathbf{F}^T(\boldsymbol{\theta}) \mathbf{K} \mathbf{F}(\boldsymbol{\theta})$  for any positive semi-definite matrix  $\mathbf{K}$  that ensures the nonsingularity of  $\mathbf{D}(\boldsymbol{\theta})$ . And while it would be sufficient to use these

existing results to verify the connection between the CMLE and the CCRB, it is also insightful (and the point of this treatise) to examine the problem entirely from the perspective of the reduced parameter space, i.e., using the implicit function or a null space approach.<sup>13</sup>

*Proof.* By the implicit function theorem (Theorem 3.3), there exists an open set  $\mathbb{O} \subset \Theta_f$  containing  $\boldsymbol{\theta}$ , an open set  $\mathbb{P} \subset \mathbb{R}^{m-k}$ , and a continuously differentiable bijection  $\mathbf{g}_{\boldsymbol{\theta}} : \mathbb{P} \rightarrow \mathbb{O}$  such that  $\boldsymbol{\theta} = \mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi})$  for some  $\boldsymbol{\xi} \in \mathbb{P}$ . The likelihood for  $\boldsymbol{\xi}$  is given by  $q(\mathbf{y}_n; \boldsymbol{\xi}) = p(\mathbf{y}_n; \mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi}))$ .

Let  $\hat{\boldsymbol{\xi}}(\mathbf{y}_n)$  be the MLE of  $\boldsymbol{\xi}$  based on the likelihood  $q(\mathbf{y}_n; \boldsymbol{\xi})$ . Since the MLE is consistent and asymptotically efficient, then

$$\sqrt{n} \left( \hat{\boldsymbol{\xi}}(\mathbf{y}_n) - \boldsymbol{\xi} \right) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \tilde{\mathbf{I}}^{-1}(\boldsymbol{\xi})). \quad (3.34)$$

In particular, since  $\hat{\boldsymbol{\xi}}(\mathbf{y}_n) \rightarrow \boldsymbol{\xi}$  as  $n \rightarrow \infty$ , then for  $n$  sufficiently large, say  $n > N$ ,  $\hat{\boldsymbol{\xi}}(\mathbf{y}_n) \in \mathbb{P}$ . Let  $\hat{\boldsymbol{\theta}}(\mathbf{y}_n)$  be the CMLE of  $\boldsymbol{\theta}$  based on the likelihood  $p(\mathbf{y}_n; \boldsymbol{\theta})$  and the constraint  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$ . By the invariance property [36, 62], for  $n > N$ ,  $\hat{\boldsymbol{\theta}}(\mathbf{y}_n) = \mathbf{g}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\xi}}(\mathbf{y}_n))$  and  $\hat{\boldsymbol{\theta}}(\mathbf{y}_n) \in \mathbb{O}$ . Therefore, since  $\hat{\boldsymbol{\xi}}(\mathbf{y}_n) \rightarrow \boldsymbol{\xi}$  as  $n \rightarrow \infty$ ,  $\hat{\boldsymbol{\theta}}(\mathbf{y}_n) \rightarrow \boldsymbol{\theta}$  also and the CMLE is consistent.

The Taylor series expansion (see section 3.1.5) of  $\mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi}')$  can be truncated using a Lagrange remainder term [38, p. 232] as in

$$\mathbf{g}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\xi}}(\mathbf{y}_n)) = \mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi}) + \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\xi}, \hat{\boldsymbol{\xi}}(\mathbf{y}_n)) \left( \hat{\boldsymbol{\xi}}(\mathbf{y}_n) - \boldsymbol{\xi} \right)$$

where  $\mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\xi}, \hat{\boldsymbol{\xi}}(\mathbf{y}_n))$  is meant to represent a matrix of the form of  $\mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\xi})$  where each

---

<sup>13</sup>Nevertheless, a proof directly from Crowder's asymptotic normality result is detailed in appendix A.3.

row is evaluated at possibly different points  $\boldsymbol{\xi}'^{(i)}$ ,  $i = 1, \dots, m - k$ , each existing on the line segment starting at  $\boldsymbol{\xi}$  and ending with  $\hat{\boldsymbol{\xi}}(\mathbf{y}_n)$ . From the invariance property of the MLE, this can be rewritten as

$$\hat{\boldsymbol{\theta}}(\mathbf{y}_n) - \boldsymbol{\theta} = \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\xi}, \hat{\boldsymbol{\xi}}(\mathbf{y}_n)) \left( \hat{\boldsymbol{\xi}}(\mathbf{y}_n) - \boldsymbol{\xi} \right).$$

Since the MLE  $\hat{\boldsymbol{\xi}}(\mathbf{y}_n)$  is consistent, then  $\mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\xi}, \hat{\boldsymbol{\xi}}(\mathbf{y}_n)) \xrightarrow{p} \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\xi})$ . Given this and (3.34), then by Slutsky's theorem [62, p. 60]

$$\sqrt{n} \left( \hat{\boldsymbol{\theta}}(\mathbf{y}_n) - \boldsymbol{\theta} \right) \xrightarrow{d} \mathcal{N} \left( \mathbf{0}, \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\xi}) \tilde{\mathbf{I}}^{-1}(\boldsymbol{\xi}) \mathbf{G}_{\boldsymbol{\theta}}^T(\boldsymbol{\xi}) \right),$$

which by theorem 3.5 shows  $\hat{\boldsymbol{\theta}}(\mathbf{y}_n)$  is asymptotically efficient with respect to the CCRB. □

The conditions for asymptotic normality with respect to the CCRB are the conditions that  $\hat{\boldsymbol{\xi}}(\mathbf{y}_n)$  be asymptotically normal [14, p. 516]. For the MLE these include (a) differentiability of the Fisher score, (b) the Fisher information continuous with respect to the parameter and nonzero at  $\boldsymbol{\xi}$ , and (c) consistency. For the CMLE and theorem 3.32, this translates to (a) differentiability of the Fisher score and the existence of first and second derivatives of any implicit function (or equivalently, the constraint  $\mathbf{f}$ ), (b)  $\mathbf{U}^T(\boldsymbol{\theta}') \mathbf{I}(\boldsymbol{\theta}') \mathbf{U}(\boldsymbol{\theta}')$  continuous with respect to  $\boldsymbol{\theta}'$  and regular at  $\boldsymbol{\theta}' = \boldsymbol{\theta}$ , and (c) consistency of the CMLE.

### 3.4.3 The Method of Scoring Under Parametric Constraints

The method of scoring for unconstrained parameters is detailed in section 2.4.3. Here, we examine scoring with constraints. Assume we have an iterate  $\dot{\boldsymbol{\theta}}^{(k)} \in \Theta_f$  and

we wish to improve this iterate in the sense of the optimization problem expressed in (3.31). The method of scoring does not directly apply, since any projection step will not take into account the constraint, i.e., it is likely the direction of steepest ascent is not the appropriate path in terms of maximizing the likelihood subject to the functional equality constraints. Thus, it is desirable to have projected direction and restoration steps that take the constraints into consideration.

Given an initial estimate  $\dot{\boldsymbol{\theta}}^{(k)}$ , there exists a set  $\mathbb{O} \ni \dot{\boldsymbol{\theta}}^{(1)}$  open in  $\Theta_f$ , a set  $\mathbb{P}$  open in  $\mathbb{R}^{m-k}$ , and a continuously differentiable function  $\mathbf{g}_{\dot{\boldsymbol{\theta}}^{(k)}} : \mathbb{R}^{m-k} \rightarrow \mathbb{R}^m$  such that  $\mathbf{g}_{\dot{\boldsymbol{\theta}}^{(k)}}$  is a diffeomorphism on  $\mathbb{P}$ ,  $\mathbf{g}_{\dot{\boldsymbol{\theta}}^{(k)}}(\mathbb{P}) = \mathbb{O}$ , and in particular there exists a  $\dot{\boldsymbol{\xi}}^{(k)} \in \mathbb{P}$  such that  $\mathbf{g}_{\dot{\boldsymbol{\theta}}^{(k)}}(\dot{\boldsymbol{\xi}}^{(k)}) = \dot{\boldsymbol{\theta}}^{(k)}$ . Scoring can now be applied in the reduced parameter space of  $\mathbb{R}^{m-k}$ .

For the given set of observations  $\mathbf{x}$  and this corresponding initial estimate  $\dot{\boldsymbol{\xi}}^{(k)}$ , the method of scoring suggests the projection step

$$\dot{\boldsymbol{\xi}}^{(k+1)} = \dot{\boldsymbol{\xi}}^{(k)} + \tilde{\mathbf{I}}^{-1}(\dot{\boldsymbol{\xi}}^{(k)})\tilde{\mathbf{s}}(\mathbf{x}; \dot{\boldsymbol{\xi}}^{(k)})$$

to generate a better estimate  $\dot{\boldsymbol{\xi}}^{(k+1)}$  in the sense of maximizing the likelihood  $q(\mathbf{x}; \dot{\boldsymbol{\xi}}') = p(\mathbf{x}; \mathbf{g}_{\dot{\boldsymbol{\theta}}^{(k)}}(\dot{\boldsymbol{\xi}}'))$ . As with many iterative procedures, convergence is only guaranteed under certain initial conditions. If the projection step or shift is too large, then  $\dot{\boldsymbol{\theta}}^{(k+1)}$  may not be a *usable* point, i.e., an iterate that increases the value of the likelihood function. To add stability to the procedure, often a step-size rule or shift-cutting is employed. This amounts to the inclusion of a multiplicative factor  $\alpha^{(k)} \in [0, 1]$ , modifying the projection step to

$$\dot{\boldsymbol{\xi}}^{(k+1)} = \dot{\boldsymbol{\xi}}^{(k)} + \alpha^{(k)} \tilde{\mathbf{I}}^{-1}(\dot{\boldsymbol{\xi}}^{(k)})\tilde{\mathbf{s}}(\mathbf{x}; \dot{\boldsymbol{\xi}}^{(k)}).$$

Choosing an appropriate *step-size rule* for  $\alpha^{(k)}$  will guarantee convergence, although typically at a cost to the rate of convergence.

The Taylor series expansion (see section 3.1.5) of  $\mathbf{g}_{\dot{\boldsymbol{\theta}}^{(k)}}$  about  $\dot{\boldsymbol{\xi}}^{(k)}$  and evaluated at  $\dot{\boldsymbol{\xi}}^{(k+1)}$  is given by

$$\mathbf{g}_{\dot{\boldsymbol{\theta}}^{(k)}}(\dot{\boldsymbol{\xi}}^{(k+1)}) = \mathbf{g}_{\dot{\boldsymbol{\theta}}^{(k)}}(\dot{\boldsymbol{\xi}}^{(k)}) + \mathbf{G}_{\dot{\boldsymbol{\theta}}^{(k)}}(\dot{\boldsymbol{\xi}}^{(k)}) \cdot (\dot{\boldsymbol{\xi}}^{(k+1)} - \dot{\boldsymbol{\xi}}^{(k)}) + o(\|\dot{\boldsymbol{\xi}}^{(k+1)} - \dot{\boldsymbol{\xi}}^{(k)}\|)$$

where  $o(\|\dot{\boldsymbol{\xi}}^{(k+1)} - \dot{\boldsymbol{\xi}}^{(k)}\|)$  is a term that shrinks faster than  $\|\dot{\boldsymbol{\xi}}^{(k+1)} - \dot{\boldsymbol{\xi}}^{(k)}\|$  as  $k \rightarrow \infty$ . Ignoring this error term, this generates an iteration in the larger dimensional parameter space  $\Theta \subset \mathbb{R}^m$  by defining the next iterate  $\dot{\boldsymbol{\theta}}^{(k+1)} = \mathbf{g}_{\dot{\boldsymbol{\theta}}^{(k)}}(\dot{\boldsymbol{\xi}}^{(k+1)})$ . That is,

$$\begin{aligned} \dot{\boldsymbol{\theta}}^{(k+1)} &= \mathbf{g}_{\dot{\boldsymbol{\theta}}^{(k)}}(\dot{\boldsymbol{\xi}}^{(k)}) + \mathbf{G}_{\dot{\boldsymbol{\theta}}^{(k)}}(\dot{\boldsymbol{\xi}}^{(k)}) \cdot (\dot{\boldsymbol{\xi}}^{(k+1)} - \dot{\boldsymbol{\xi}}^{(k)}) \\ &= \dot{\boldsymbol{\theta}}^{(k)} + \alpha^{(k)} \mathbf{G}_{\dot{\boldsymbol{\theta}}^{(k)}}(\dot{\boldsymbol{\xi}}^{(k)}) \tilde{\mathbf{I}}^{-1}(\dot{\boldsymbol{\xi}}^{(k)}) \tilde{\mathbf{s}}(\mathbf{x}; \dot{\boldsymbol{\xi}}^{(k)}) \\ &= \dot{\boldsymbol{\theta}}^{(k)} + \alpha^{(k)} \mathbf{G}_{\dot{\boldsymbol{\theta}}^{(k)}}(\dot{\boldsymbol{\xi}}^{(k)}) (\mathbf{G}_{\dot{\boldsymbol{\theta}}^{(k)}}^T(\dot{\boldsymbol{\xi}}^{(k)}) \mathbf{I}(\dot{\boldsymbol{\theta}}^{(k)}) \mathbf{G}_{\dot{\boldsymbol{\theta}}^{(k)}}(\dot{\boldsymbol{\xi}}^{(k)}))^{-1} \mathbf{G}_{\dot{\boldsymbol{\theta}}^{(k)}}^T(\dot{\boldsymbol{\xi}}^{(k)}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)}) \\ &= \dot{\boldsymbol{\theta}}^{(k)} + \alpha^{(k)} \mathbf{U}(\dot{\boldsymbol{\theta}}^{(k)}) (\mathbf{U}^T(\dot{\boldsymbol{\theta}}^{(k)}) \mathbf{I}(\dot{\boldsymbol{\theta}}^{(k)}) \mathbf{U}(\dot{\boldsymbol{\theta}}^{(k)}))^{-1} \mathbf{U}^T(\dot{\boldsymbol{\theta}}^{(k)}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)}). \end{aligned}$$

In comparison with the classical method of scoring, this iteration

$$\dot{\boldsymbol{\theta}}^{(k+1)} = \dot{\boldsymbol{\theta}}^{(k)} + \alpha^{(k)} \text{CCRB}(\dot{\boldsymbol{\theta}}^{(k)}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)}) \quad (3.35)$$

is essentially a replacement of the CRB with the CCRB. This is the projection step of the method of scoring with parametric equality constraints.<sup>14</sup> Intuitively, this should seem appropriate since the CCRB is a generalization of the CRB. However,

---

<sup>14</sup>Osborne [55] used a Lagrangian multiplier approach to develop the method of scoring. But his scenario was restricted to linear constraints and, hence, lacked the restoration step. Additionally, he makes no mention that the matrix projecting the negative Jacobian of the objective is a constrained Cramér-Rao bound. Note the structure of this projection step is well-known as a nonstatistical formulation exists for the conventional optimization problem in [25, p. 178].

even with an appropriate step-size rule to generate usable iterates, since there is no certainty that  $\dot{\boldsymbol{\xi}}^{(k+1)} \in \mathbb{P}$ , then it is likely that  $\dot{\boldsymbol{\theta}}^{(k+1)}$  will not be a feasible point. To correct this, an encompassing restoration step is required to produce the next iterate, i.e.,

$$\dot{\boldsymbol{\theta}}^{(k+1)} = \boldsymbol{\pi} [\dot{\boldsymbol{\theta}}^{(k)} + \alpha^{(k)} \text{CCRB}(\dot{\boldsymbol{\theta}}^{(k)}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})] \quad (3.36)$$

where  $\boldsymbol{\pi}[\cdot]$  is the natural projection of  $\mathbb{R}^m$  onto  $\Theta_f$ . This is the method of scoring with parametric equality constraints. With this additional restoration step, the

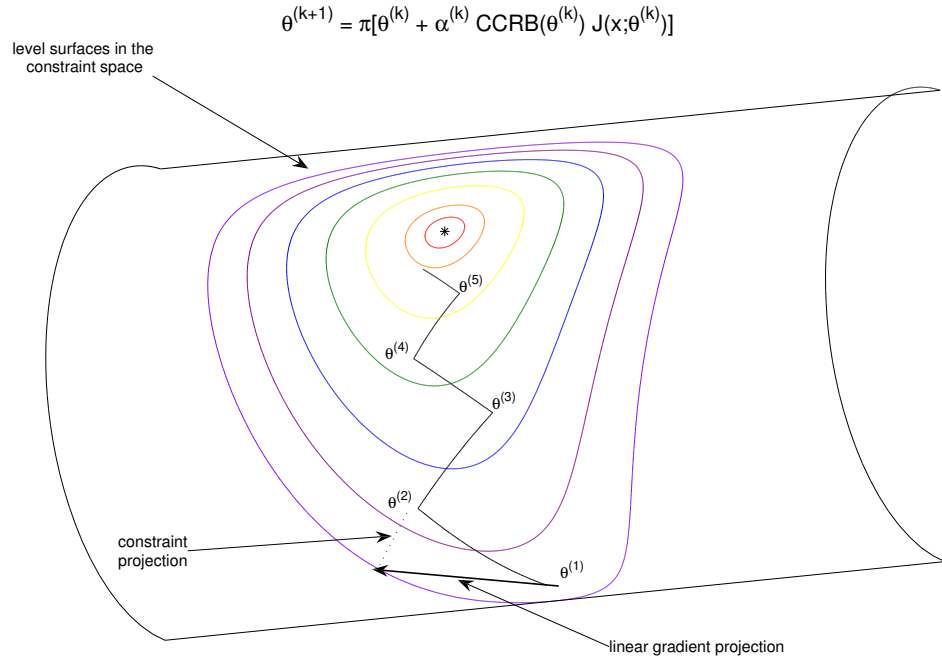


Figure 3.3: Path created by iterates from the method of scoring with constraints.

usability of the iterate would be tested and accepted or rejected after (3.36) as opposed to after (3.35). For a convex set, the natural projection is well-defined. In general, though, some other rule will likely need to be applied in cases for which there does not exist a unique shortest distance to  $\Theta_f$ , e.g., reducing the step size

$\alpha^{(k)}$ . Simple projections, e.g., onto planar or spherical constraints, are relatively simple operations, however, it might be more commonly the case that the restoration cannot be expressed analytically. To ensure the iterates satisfy the constraints approximately, one approach is to apply an additional iterative process [25]

$$\dot{\boldsymbol{\theta}}^{(k,l+1)} = \dot{\boldsymbol{\theta}}^{(k,l)} - \mathbf{F}^T(\dot{\boldsymbol{\theta}}^{(k,l)}) (\mathbf{F}(\dot{\boldsymbol{\theta}}^{(k,l)}) \mathbf{F}^T(\dot{\boldsymbol{\theta}}^{(k,l)}))^{-1} \mathbf{f}(\dot{\boldsymbol{\theta}}^{(k,l)}),$$

where  $\dot{\boldsymbol{\theta}}^{(k,1)} = \dot{\boldsymbol{\theta}}^{(k)}$  and  $\dot{\boldsymbol{\theta}}^{(k+1)} = \boldsymbol{\pi} [\dot{\boldsymbol{\theta}}^{(k)}] = \dot{\boldsymbol{\theta}}^{(k,l)}$  when  $\mathbf{f}(\dot{\boldsymbol{\theta}}^{(k,l)}) \approx \mathbf{0}$  to a desired degree of accuracy and provided the iterate is still usable. Alternatively, a penalty can be added to the cost (objective) function, e.g., as in

$$\log p(\mathbf{x}; \boldsymbol{\theta}') + \eta \sum_{i=1}^k |f_i(\boldsymbol{\theta}')|$$

for some positive  $\eta$  and where  $f_i$  is the  $i$ th constraint equation, to limit the divergence of the iterations away from  $\Theta_f$ .

### 3.4.3.1 Convergence Properties

There is a large class of conditions that guarantee convergence in fixed point theorems, some of which can be found in [64, 25, 10]. The most general statement is that given an initialization “sufficiently close” to the maximum value  $\hat{\boldsymbol{\theta}}(\mathbf{x})$ , the sequence  $\{\dot{\boldsymbol{\theta}}^{(k)}\}$  generated by the algorithm will converge to this CMLE. As the method of scoring with parametric equality constraints is a Newton-type method, convergence properties that already exist for these methods can be adapted here. As it is impossible to cover all the potential approaches to developing properties for this constrained scoring algorithm, this section focuses on properties similar to those

found in Goldstein [22]. To ease reading of this section, the proofs of the theorems in this section are presented in appendix B.

First, define  $\Theta_{\dot{\boldsymbol{\theta}}^{(k)}} = \{\boldsymbol{\theta}' \in \Theta_f : p(\mathbf{x}; \boldsymbol{\theta}') \geq p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})\}$  as the set of all feasible and usable iterates after the  $k$ th iterate  $\dot{\boldsymbol{\theta}}^{(k)}$ . The step rule for the properties in this section is as follows: for a fixed  $\beta \in (0, 1)$  choose the least positive integer  $m^{(k)}$  such that  $\alpha^{(k)} = \beta^{m^{(k)}}$  satisfies the inequality

$$\alpha^{(k)} (\log p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k+1)}) - \log p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})) \geq \kappa \left\| \dot{\boldsymbol{\theta}}^{(k+1)} - \dot{\boldsymbol{\theta}}^{(k)} \right\|_{\mathbf{I}(\dot{\boldsymbol{\theta}}^{(k)})}^2 \quad (3.37)$$

where  $\dot{\boldsymbol{\theta}}^{(k+1)}$  is defined by (3.36). If no finite  $m^{(k)}$  exists, then choose  $\alpha^{(k)} = 0$ . This type of step-size rule enforces a stepwise Lipschitz condition. For theorems 3.33 and 3.37, we require that  $\Theta_f$  be convex.<sup>15</sup>

**Theorem 3.33.** If for any iterate  $\dot{\boldsymbol{\theta}}^{(k)} \in \Theta_f$  there does not exist an  $\alpha^{(k)} > 0$  that satisfies (3.37), then  $\dot{\boldsymbol{\theta}}^{(k)}$  is a stationary point.

Therefore, when the step rule forces the choice of  $\alpha^{(k)} = 0$  then the method of scoring with parametric equality constraints has converged. The next theorem details a property on the sequence of likelihood functions generated by the iterates.

**Theorem 3.34.** The sequence  $\{p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})\}$  is a monotone increasing sequence. Furthermore, if  $p(\mathbf{x}; \cdot)$  is bounded above, then  $\{p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})\}$  converges.

---

<sup>15</sup>For any nonlinear equality constraint,  $\Theta_f$  will not be convex. However, locally, the restoration of the linear projection onto the tangent space of the CCRB has the appearance of the restoration onto a convex set for sufficiently small  $\alpha^{(k)}$ , i.e.,  $\Theta_f$  appears locally convex. While it may be possible to restrict the step size to this local convexity, such an enhancement is beyond the scope of the work presented here. For inequality constraints,  $\Theta_f$  may be convex; although such constraints, as in example 3.15, will not inform the projection update in (3.35), they might inform the restoration update in (3.36).



Thus,  $\Theta_{\dot{\boldsymbol{\theta}}^{(k)}}$  is a decreasing sequence of closed (nested) sets, i.e.,  $\Theta_{\dot{\boldsymbol{\theta}}^{(k+1)}} \subset \Theta_{\dot{\boldsymbol{\theta}}^{(k)}}$  or for any given sequence,  $\dot{\boldsymbol{\theta}}^{(i)} \in \Theta_{\dot{\boldsymbol{\theta}}^{(j)}}$  provided  $i \geq j$ . That is, using a proper step size rule will guarantee usable iterates. The monotonicity of  $\{p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})\}$ , even if bounded above, does not imply monotonicity in the sequence  $\{\log p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k+1)}) - \log p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})\}$ . However, this does guarantee convergence.

**Theorem 3.35.** If the likelihood  $p(\mathbf{x}; \cdot)$  is bounded above, then the sequence

$$\{\log p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k+1)}) - \log p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})\}$$

vanishes.

Hence, a bounded likelihood function guarantees the existence of a maximum likelihood solution(s). This can also be shown by the nested interval theorem [38, problem 2-1-12]. These previous properties would be a consequence of any rule that chooses feasible, usable iterates. The value of the rule in (3.37) is that it allows for statements to be made on the sequence  $\{\dot{\boldsymbol{\theta}}^{(k)}\}$ .

**Theorem 3.36.** If the likelihood  $p(\mathbf{x}; \cdot)$  is bounded above, then the sequence

$$\{\|\dot{\boldsymbol{\theta}}^{(k+1)} - \dot{\boldsymbol{\theta}}^{(k)}\|_{\mathbf{I}(\dot{\boldsymbol{\theta}}^{(k)})}\}$$

vanishes as  $k \rightarrow \infty$ .

This theorem does not guarantee that the sequence  $\{\dot{\boldsymbol{\theta}}^{(k)}\}$  will converge or even have a *limit point*.<sup>16</sup> (The series  $\sum_{j=1}^k \frac{1}{j}$  is an example of a sequence satisfying the above theorem with no real-valued limit points.) Although, if the set  $\Theta_{\dot{\boldsymbol{\theta}}^{(k)}}$  is

---

<sup>16</sup>A point  $a$  is a *limit point* of a sequence  $\{a_n\}$  if for any integer  $K$  and any  $\epsilon > 0$ , there exists an  $k > K$  such that  $|a_k - a| < \epsilon$ .

bounded for some  $k$ , then the Bolzano-Weierstrass theorem [38, p. 52, theorem 2-12] implies the existence of a limit point of the sequence.

**Theorem 3.37.** If  $\Theta_{\dot{\theta}^{(1)}}$  is compact and convex, then limit points of the sequence  $\{\dot{\theta}^{(k)}\}$  are also stationary points.

The theorem remains true if  $\Theta_{\dot{\theta}^{(k)}}$  is compact and convex for any  $k$ .

**Theorem 3.38.** If  $\Theta_{\dot{\theta}^{(1)}}$  is compact for all sequences in a closed set of  $\Theta_f$  and if there is a unique limit point  $\theta^*$  for all such sequences then  $\lim_{k \rightarrow \infty} \dot{\theta}^{(k)} = \theta^*$  for every sequence  $\{\dot{\theta}^{(k)}\}$ . Also,  $\theta^*$  is the maximum of  $p(\mathbf{x}; \cdot)$ .

### 3.4.3.2 Linear constraints

Linear constraints on the parameter typically restrict the parameter to a set of the form  $\Theta_f = \{\theta' \in \Theta : \mathbf{F}\theta' + \mathbf{v} = \mathbf{0}\}$  (see section 3.3). Under this linear constraint, the restoration operation  $\pi[\cdot]$  is redundant since any step remains in the constraint space, i.e., since  $\dot{\theta}^{(k)} \in \Theta_f$  and  $\mathbf{F} \cdot \text{CCRB}(\dot{\theta}^{(k)}) = \mathbf{0}$ . Thus the method of scoring with parametric equality constraints in (3.36) simplifies to the iteration

$$\dot{\theta}^{(k+1)} = \dot{\theta}^{(k)} + \alpha^{(k)} \text{CCRB}(\dot{\theta}^{(k)}) \mathbf{s}(\mathbf{x}; \dot{\theta}^{(k)}).$$

**Example 3.39** (Linear model with linear constraints). In the linear model with normal noise case of section 3.3.2, we have  $\mathbf{x} = \mathbf{H}\theta + \mathbf{w}$  with  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$ . In this case, the negative Hessian is the FIM and the optimization problem becomes a null space quadratic exercise [21], i.e., the minimization of a quadratic objective subject to a linear constraint. The  $\text{CCRB} = \mathbf{U} (\mathbf{U}^T \mathbf{H}^T \mathbf{C}^{-1} \mathbf{H} \mathbf{U})^{-1} \mathbf{U}^T$ , which is constant with

respect to the parameter; and the Fisher score is  $\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}') = \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H}\boldsymbol{\theta}')$ . Hence, if  $\dot{\boldsymbol{\theta}}^{(1)}$  is any feasible vector, e.g.,  $\dot{\boldsymbol{\theta}}^{(1)} = -\mathbf{F}^T (\mathbf{F}\mathbf{F}^T)^{-1} \mathbf{v}$ , then the method of scoring with constraints finds the CMLE in one step to be

$$\dot{\boldsymbol{\theta}}^{(2)} = \dot{\boldsymbol{\theta}}^{(1)} + \text{CCRB} \cdot \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H}\dot{\boldsymbol{\theta}}^{(1)})$$

which is exactly the formula in (3.29) from theorem 3.30. The next iterate  $\dot{\boldsymbol{\theta}}^{(3)} = \dot{\boldsymbol{\theta}}^{(2)} + \text{CCRB} \cdot \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H}\dot{\boldsymbol{\theta}}^{(2)})$  reveals that the procedure reaches a fixed point since

$$\begin{aligned} & \text{CCRB} \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H}\dot{\boldsymbol{\theta}}^{(2)}) \\ &= \text{CCRB} \mathbf{H}^T \mathbf{C}^{-1} \left( \mathbf{x} - \mathbf{H} (\dot{\boldsymbol{\theta}}^{(1)} + \text{CCRB} \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H}\dot{\boldsymbol{\theta}}^{(1)})) \right) \\ &= \text{CCRB} \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H}\dot{\boldsymbol{\theta}}^{(1)}) - \text{CCRB} \mathbf{H}^T \mathbf{C}^{-1} \mathbf{H} \text{CCRB} \mathbf{H}^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H}\dot{\boldsymbol{\theta}}^{(1)}) \\ &= \mathbf{0} \end{aligned}$$

Therefore,  $\hat{\boldsymbol{\theta}}_{\text{CML}}(\mathbf{x}) = \dot{\boldsymbol{\theta}}^{(2)}$ .

**Example 3.40** (Jamshidian's GP algorithm). Jamshidian [33] developed a Gradient Projection (GP) algorithm for maximizing the likelihood subject to linear parameter constraints using the iteration

$$\dot{\boldsymbol{\theta}}^{(k+1)} = \dot{\boldsymbol{\theta}}^{(k)} + \alpha^{(k)} \left( \mathbf{W}^{-1} - \mathbf{W}^{-1} \mathbf{F}^T (\mathbf{F} \mathbf{W}^{-1} \mathbf{F}^T)^{-1} \mathbf{F} \mathbf{W}^{-1} \right) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)}) \quad (3.38)$$

for some positive definite matrix  $\mathbf{W}$ . An optimal choice with regard to the algorithm's rate of convergence, Jamshidian suggests, is a possibly diagonally loaded Hessian of the log-likelihood

$$\mathbf{W}(\mathbf{x}, \dot{\boldsymbol{\theta}}^{(k)}) = \frac{\partial}{\partial \boldsymbol{\theta}^T} \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)}) + \gamma^{(k)} \mathbf{I}_{m \times m},$$

where  $\gamma^{(k)} \geq 0$  is chosen to be sufficiently large enough to ensure the positive definiteness of the matrix  $\mathbf{W}$ . This formulation is closely connected to the method of scoring with constraints. The GP iteration is equivalent to scoring by choosing  $\mathbf{W}(\mathbf{x}, \dot{\boldsymbol{\theta}}^{(k)})$  to be the FIM, when it is nonsingular. Indeed, the projecting matrix is similar to the Marzetta form of the CCRB in [47] with  $\mathbf{W}(\mathbf{x}, \dot{\boldsymbol{\theta}}^{(k)})$  replacing the FIM. This fact produces a slight generalization of the GP iteration, given by

$$\dot{\boldsymbol{\theta}}^{(k+1)} = \dot{\boldsymbol{\theta}}^{(k)} + \alpha^{(k)} \mathbf{U} (\mathbf{U}^T \mathbf{W}(\mathbf{x}, \dot{\boldsymbol{\theta}}^{(k)}) \mathbf{U})^{-1} \mathbf{U}^T \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})$$

where  $\mathbf{U}$  is defined as in (3.6). In this formulation, the projecting metric  $\mathbf{W}(\mathbf{x}, \dot{\boldsymbol{\theta}}^{(k)})$  only needs to be positive semidefinite. Alternatively, in (3.38) the Aitchison and Silvey [4] substitution for the FIM,  $\mathbf{I}(\dot{\boldsymbol{\theta}}^{(k)}) + \mathbf{F}^T(\dot{\boldsymbol{\theta}}^{(k)}) \mathbf{K} \mathbf{F}(\dot{\boldsymbol{\theta}}^{(k)})$ , instead of a diagonally loaded Hessian of the log-likelihood (or even a diagonally loaded Fisher information matrix).

In this sense, the two iterations are equivalent for the linear model, when the FIM is simply the negative Hessian. This occurs when the log-likelihood is quadratic (normal). This also suggests the adaption of Jamshidian's GP algorithm to cases of nonlinear constraints. Likewise, asymptotically, where  $\mathbf{y}_n$  is denoted as in section 3.4.2, then by the law of large numbers the Jamshidian projection matrix

$$\mathbf{W}(\mathbf{y}_n, \dot{\boldsymbol{\theta}}^{(k)}) \rightarrow n \mathbf{I}(\dot{\boldsymbol{\theta}}^{(k)}) + \gamma^{(k)} \mathbf{I}_{m \times m}$$

for an arbitrarily small (possibly zero)  $\gamma^{(k)} > 0$ , as  $\mathbf{I}(\dot{\boldsymbol{\theta}}^{(k)})$  is positive semidefinite.

The projection step update then becomes

$$\dot{\boldsymbol{\theta}}^{(k+1)} = \dot{\boldsymbol{\theta}}^{(k)} + \alpha^{(k)} \text{CCRB}(\dot{\boldsymbol{\theta}}^{(k)}) \mathbf{s}(\mathbf{y}_n; \dot{\boldsymbol{\theta}}^{(k)}) - \alpha^{(k)} \gamma^{(k)} \mathbf{s}(\mathbf{y}_n; \dot{\boldsymbol{\theta}}^{(k)}),$$

i.e., essentially equivalent to the method of scoring with parametric equality constraints in (3.35).

### 3.5 Hypothesis testing

In section 2.5, hypothesis testing (the Rao and Wald tests) using the CRB was reviewed. In this section, hypothesis testing is considered under a constrained alternative. Assume  $\mathbf{h} : \mathbb{R}^m \rightarrow \mathbb{R}^r$  is a consistent and nonredundant differentiable function, which is also consistent and nonredundant with the differentiable function  $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^k$ . Hence,  $\Theta_h = \{\boldsymbol{\theta}' : \mathbf{h}(\boldsymbol{\theta}') = \mathbf{0}\} \subset \Theta_f$ , where  $\Theta_f$  is defined as in (3.4), and also  $\text{rank}\left(\begin{bmatrix} \mathbf{H}(\boldsymbol{\theta}) \\ \mathbf{F}(\boldsymbol{\theta}) \end{bmatrix}\right) = r + k \leq m$  and  $r < m - k$ . Then the hypothesis test can be stated as

$$H_0 : \mathbf{h}(\boldsymbol{\theta}) = \mathbf{0} \quad \text{vs.} \quad H_1 : \mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}. \quad (3.39)$$

Naturally,  $\mathbf{f}(\boldsymbol{\theta}') = \mathbf{0}$  under these conditions defines an implicit function locally, so assume  $\mathbf{g}_{\boldsymbol{\theta}'} : \mathbb{R}^{m-k} \rightarrow \mathbb{R}^k$  is such a function satisfying theorem 3.3 for any  $\boldsymbol{\theta}' \in \Theta_f$ . Then, a locally (or asymptotically) equivalent hypothesis can be stated as

$$H_0 : \mathbf{h}(\mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi})) = \mathbf{0} \quad \text{vs.} \quad H_1 : \mathbf{h}(\mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi})) \neq \mathbf{0}. \quad (3.40)$$

In this formulation, the well-known Rao and Wald statistics were shown in section 2.5.

### 3.5.1 The Rao statistic

For the hypothesis testing scenario in (3.40), the Rao test statistic presented in section 2.5.1 is given by

$$\rho(\mathbf{y}_n) = \tilde{\mathbf{s}}^T(\mathbf{y}_n; \hat{\boldsymbol{\xi}}_{h(g)}(\mathbf{y}_n)) \tilde{\mathbf{I}}_n^{-1}(\hat{\boldsymbol{\xi}}_{h(g)}(\mathbf{y}_n)) \tilde{\mathbf{s}}(\mathbf{y}_n; \hat{\boldsymbol{\xi}}_{h(g)}(\mathbf{y}_n))$$

where  $\tilde{\mathbf{s}}(\mathbf{y}_n; \hat{\boldsymbol{\xi}}_{h(g)}(\mathbf{y}_n))$  is the Fisher score of the observations  $\mathbf{y}_n$  (as defined in section 3.4.2) and evaluated at the  $H_0$ -constrained maximum likelihood estimate or constrained root of the likelihood estimate (CRLE)  $\hat{\boldsymbol{\xi}}_{h(g)}(\mathbf{y}_n)$  of the likelihood  $q(\mathbf{y}_n; \boldsymbol{\xi}) = p(\mathbf{y}_n; \mathbf{g}_{\boldsymbol{\theta}}(\boldsymbol{\xi}))$ , and  $\tilde{\mathbf{I}}_n^{-1}(\hat{\boldsymbol{\xi}}_{h(g)}(\mathbf{y}_n))$  is the  $n$ -sample Fisher information evaluated at the CMLE. In this context, the CMLE is the solution to the optimization problem  $\max_{\boldsymbol{\xi}'} \log q(\mathbf{x}; \boldsymbol{\xi}') \text{ s.t. } \mathbf{h} \circ \mathbf{g}_{\boldsymbol{\theta}'}(\boldsymbol{\xi}') = \mathbf{0}$ . As in theorem 3.5,  $\tilde{\mathbf{s}}(\mathbf{y}_n; \boldsymbol{\xi}') = \mathbf{G}_{\boldsymbol{\theta}'}(\boldsymbol{\xi}') \mathbf{s}(\mathbf{y}_n; \boldsymbol{\theta}')$  and  $\tilde{\mathbf{I}}_n(\boldsymbol{\xi}') = n \mathbf{G}^T(\boldsymbol{\xi}') \mathbf{I}(\mathbf{g}_{\boldsymbol{\theta}'}(\boldsymbol{\xi}')) \mathbf{G}(\boldsymbol{\xi}')$ . Also, recall that for sufficiently large  $n$ ,  $\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n) = \mathbf{g}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\xi}}_{h(g)}(\mathbf{y}_n))$ . Therefore, the locally (or asymptotically) equivalent Rao test statistic for the hypothesis in (3.39) is

$$\rho(\mathbf{y}_n) = \frac{1}{n} \mathbf{s}^T(\mathbf{y}_n; \hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \text{CCRB}(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \mathbf{s}(\mathbf{y}_n; \hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)), \quad (3.41)$$

which is analogous to (2.9) with the CCRB replacing the CRB. Under  $H_0$ ,  $\rho(\mathbf{y}_n)$  is still asymptotically  $\chi_r^2$  in distribution. The corresponding Lagrange-multiplier variant of this statistic is given by

$$\rho(\mathbf{y}_n) = \frac{1}{n} \hat{\boldsymbol{\lambda}}_h(\mathbf{y}_n) \mathbf{H}(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \text{CCRB}(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \mathbf{H}^T(\hat{\boldsymbol{\theta}}_h(\mathbf{y}_n)) \hat{\boldsymbol{\lambda}}_h(\mathbf{y}_n), \quad (3.42)$$

where the Lagrange multiplier estimates  $\hat{\boldsymbol{\lambda}}_h(\mathbf{y}_n)$  are based on the first order conditions relating to the constraint  $\mathbf{h}$  (not  $\mathbf{h}$  and  $\mathbf{f}$ ).

The result in (3.41) is consistent with the classical results in [63], although not explicitly in this form. For the hypothesis scenario in (3.39) the Lagrange multiplier statistic should be  $\hat{\boldsymbol{\lambda}}_h^T \mathbf{R}_{h,\hat{\boldsymbol{\theta}}}^{-1} \hat{\boldsymbol{\lambda}}_h$ , where  $\mathbf{R}_{h,\hat{\boldsymbol{\theta}}}$  is defined by

$$\begin{bmatrix} \mathbf{P}_{\boldsymbol{\theta}} & * & * \\ * & * & * \\ * & * & \mathbf{R}_{h,\boldsymbol{\theta}} \end{bmatrix} = \begin{bmatrix} \mathbf{I}(\boldsymbol{\theta}) + \mathbf{F}^T(\boldsymbol{\theta})\mathbf{F}(\boldsymbol{\theta}) & -\mathbf{F}^T(\boldsymbol{\theta}) & \mathbf{H}^T(\boldsymbol{\theta}) \\ \mathbf{F}(\boldsymbol{\theta}) & \mathbf{0} & \mathbf{0} \\ \mathbf{H}(\boldsymbol{\theta}) & \mathbf{0} & \mathbf{0} \end{bmatrix}^{-1},$$

which is a variant of what appears in [63, equation (6.5)]. Finding the inverse using the Schur complement, it is clear that

$$\mathbf{R}_{h,\boldsymbol{\theta}} = \mathbf{H}^T(\boldsymbol{\theta}) \left[ \mathbf{D}^{-1}(\boldsymbol{\theta}) - \mathbf{D}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}) (\mathbf{F}(\boldsymbol{\theta})\mathbf{D}^{-1}(\boldsymbol{\theta})\mathbf{F}^T(\boldsymbol{\theta}))^{-1} \mathbf{F}(\boldsymbol{\theta})\mathbf{D}^{-1}(\boldsymbol{\theta}) \right] \mathbf{H}(\boldsymbol{\theta})$$

with  $\mathbf{D}(\boldsymbol{\theta}) = \mathbf{I}(\boldsymbol{\theta}) + \mathbf{F}^T(\boldsymbol{\theta})\mathbf{F}(\boldsymbol{\theta})$ . Recognizing the inner matrix as the Aitchison-Silvey-Crowder variant of the CCRB formula in (3.14) and substituting the CMLE for the parameter obtains (3.42).

### 3.5.2 The Wald statistic

Similarly, the Wald test statistic presented in section 2.5.2 is

$$\mathbf{h}^T(\mathbf{g}(\hat{\boldsymbol{\xi}}(\mathbf{y}_n))) \left( \mathbf{H}(\mathbf{g}(\hat{\boldsymbol{\xi}}(\mathbf{y}_n))) \mathbf{G}^T(\hat{\boldsymbol{\xi}}(\mathbf{y}_n)) \tilde{\mathbf{I}}_n^{-1}(\hat{\boldsymbol{\xi}}(\mathbf{y}_n)) \mathbf{G}(\hat{\boldsymbol{\xi}}(\mathbf{y}_n)) \mathbf{H}^T(\mathbf{g}(\hat{\boldsymbol{\xi}}(\mathbf{y}_n))) \right)^{-1} \mathbf{h}(\mathbf{g}(\hat{\boldsymbol{\xi}}(\mathbf{y}_n)))$$

for the testing problem in (3.40), where  $\mathbf{g}$  is localized about  $\boldsymbol{\theta}$ . Following the steps in section 3.5.1, then to test (3.39), the corresponding Wald test statistic is

$$\omega(\mathbf{y}_n) = n \mathbf{h}^T(\hat{\boldsymbol{\theta}}(\mathbf{y}_n)) \left( \mathbf{H}(\hat{\boldsymbol{\theta}}(\mathbf{y}_n)) \text{CCRB}(\hat{\boldsymbol{\theta}}(\mathbf{y}_n)) \mathbf{H}^T(\hat{\boldsymbol{\theta}}(\mathbf{y}_n)) \right)^{-1} \mathbf{h}(\hat{\boldsymbol{\theta}}(\mathbf{y}_n)).$$

As with the Rao statistic, this general Wald statistic replaces the CRB in (2.10) with the CCRB, and under  $H_0$ ,  $\omega(\mathbf{y}_n)$  is asymptotically  $\chi_r^2$  in distribution.

This agrees with the classical result in [2, ‘ $\lambda_{21}(\theta)$ ’ on p. 240] where the Gorman-Hero-Aitchison-Silvey variant of the CCRB formula in (3.12) is used instead. (The general scenario when the FIM is singular is discussed in [2, section 3.9].)

A requirement for the existence of this statistic is that  $\mathbf{H}(\boldsymbol{\theta})\text{CCRB}(\boldsymbol{\theta})\mathbf{H}^T(\boldsymbol{\theta})$  be regular. This is not an additional requirement, but a necessity in testing that the hypothesis testing function itself be identifiable for the hypothesis to be valid.

### 3.6 Discussion

The previous sections have established that the theory of the constrained CRB is equivalent to that of the CRB. The majority of the proofs, besides their novelty in the recent literature on the CCRB, essentially rely on the generation of an implicit function that translates the constrained parametric problem into an unconstrained parametric problem, for which the theory of the CRB is well-established.

This theory has been extended to the CCRB, in particular, for identifiability under constraints; for the linear model with linear constraints (already well-documented in the literature but perhaps not with reference to this CCRB); for the constrained maximum likelihood problem, its asymptotic normality and the method of scoring; as well as for the Wald and Rao hypothesis tests. This list is by no means exhaustive. There exists recent research, for example, related to biased estimation with the CCRB [9], and an extension of the CCRB for complex-valued parameters [32]. In addition, there are other areas in mathematical statistics that, as far as I know, have not yet been connected to the CCRB, including a geometric interpreta-



tion of the CCRB (possibly in the manner of [60] or an extension of [5]), biasedness issues with constrained estimation, confidence intervals or sets, and the plausibility of a Bayesian version of the CCRB.

While the primary contributions of this thesis are the theoretical results and their proofs existing in this chapter, this should not discount the practical applications of these ideas. From the practitioners viewpoint, this research has produced a number of useful tools. For example, to test local identifiability, in addition to the classic result of Rothenberg (theorem 3.23), theorem 3.24 may be used. Likewise, for strict identifiability, theorem 3.29 is useful. And the method of scoring in (3.36) adds to the list of constrained maximum likelihood methods. These and others will be applied in a communications context, in the next chapter.

## Chapter 4

### Applications of the CCRB in Communications Models

Communications models in statistical signal processing often have the structure

$$y(n) = \mu(n, \boldsymbol{\theta}_n) + w(n) \text{ , } n = 0, \dots, N, \quad (4.1)$$

where  $y(n)$  are the received observations of some model  $\mu(n, \boldsymbol{\theta}_n)$  affected (additively) by noise  $w(n)$  over a series of time samples  $n$ . The parameters may or may not be dependent on the number of time samples and the noise may or may not be independent (or even normally distributed). This general model encompasses a number of areas of signal processing, e.g., communications, sonar, radar, speech, imaging, control, sensing, networks, etc. In every one of these areas, there exist models where parametric equality constraints are of interest to practitioners in the field. The CCRB has proven useful as a performance analysis tool in localization [6], watermarking security [15], tomography [28], source bearing and symbol estimation [59], space-time block coding [42], and even a variant of least squares [46].

In this chapter, we shall detail just two rather general signal processing communications models: the convolutive channel and the calibrated array, each with unknown signal and channel components. Because the source and channel interact multiplicatively in the models, numerous variations of these two basic models are possible (and necessary). We shall formulate constraints for just a few of these vari-

ations and connect the results of this chapter to the theory developed in Chapter 3. For the sake of coherence in the presentation, the lengthier proofs of theorems in this chapter are relegated to appendix C.

In section 4.1, the convolutive mixture model with deterministic parameters is presented using a variety of descriptions. In addition, several useful terms are defined in section 4.1.1.2 that are useful to characterize conditions on the a notion of near-identifiability presented in section 4.1.2 and conditions on the Fisher information derived in 4.1.3. The corresponding complex-valued FIM (CFIM) is defined in 4.1.3.2 and important properties are given in 4.1.3.3. These properties are critical to understanding the particular inherent ambiguities in the basic convolutive mixture model and determining the class of constraints are necessary to obtain a CCRB, as discussed in 4.1.4.1. Two constraint models, the norm constraint (in 4.1.4.2) and the semiblind constraint (in 4.1.4.3) serve as further validation of the CCRB method by obtaining prior results in the literature; while another constraint model, the combination of the semiblind and unit modulus constraints in 4.1.4.4, demonstrates an important constraint model that did not previously exist in the literature.

In section 4.2, a special case of the convolutive mixture model, called the calibrated array model, is considered. The FIM for this (sub)model and its properties are detailed in section 4.2.1 and various constraints are considered in section 4.2.2. As before, a number of the constraint models are validations of the CCRB method compared with previous results in the literature (in 4.2.2.1, 4.2.2.2 and 4.2.2.5). The CCRB approach is also used to detail constraint models where the constraints are inappropriate either because the constraints are over-determined (in 4.2.2.3) or

because they are under-determined (in 4.2.2.4). As with the more general convolutive mixture model, a constraint model case is presented in 4.2.2.6 that did not previously exist in the literature.

## 4.1 Convolutive Mixture Model

The complex baseband representation of a multi-input, multi-output (MIMO) finite impulse response (FIR) system, or the convolutive mixture model, may be written as

$$\begin{aligned} y_m(n) &= \sum_{k=1}^K s^{(k)}(n) * h_m^{(k)}(n) + w_m(n) \\ &= \sum_{k=1}^K \sum_{l=0}^{L_k} s^{(k)}(n-l) h_m^{(k)}(l) + w_m(n) \end{aligned} \quad (4.2)$$

for the  $n$ th observation of the  $m$ th channel where the model consists of  $K$  sources,  $M$  channels ( $1 \leq m \leq M$ ) with maximal order  $L_k$  for the  $k$ th source (or  $L_k + 1$  taps for the  $k$ th source),  $N$  output samples ( $0 \leq n \leq N - 1$ ) per channel and  $N + L_k$  input samples per channel per source. For each source, the model can be seen in figure 4.1. The dimensions  $K$ ,  $M$ ,  $L_k$  for  $k = 1, \dots, K$ , and  $N$  are all assumed known. The elements  $s^{(k)}(n)$  denote the scalar, complex value of the  $k$ th input source at time  $n$ ; the elements of  $h_m^{(k)}(l)$  denote the scalar, complex value of the  $l$ th filter coefficient (or lag) of the  $k$ th source processed by the  $m$ th channel. Both the signal inputs and channel coefficients are treated as deterministic unknowns, i.e., parameters having true values that can be estimated. The noise elements  $w_m(n)$  are commonly assumed to be zero-mean, complex-valued, circular Gaussian iid over

space and time (and channel) with known variance  $\sigma^2$ .<sup>1</sup>

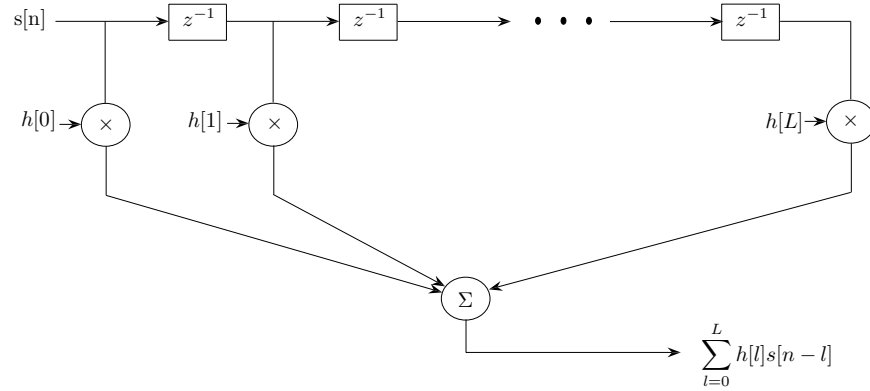


Figure 4.1: Finite Impulse Response (tapped delay line) model.

The convolution aspect of this model is often used to characterize the intersymbol interference in direct sequence CDMA (code-division multiple-access) over dispersive channels [44]. This occurs, for example, when the propagation time of the signal is shorter than the coherence time of the channel in wideband signals. The varying channel order lengths  $L_k$  represent the different coherence times of the different source-channel links. This convolution also applies to scenarios where reverberation of the transmitted signal is present, which occurs when the communications is echoed. Additionally, the additive aspect over the sources represents the multiuser interference (e.g., the cocktail problem) common in communications systems today. Moreover, this convolutive mixture model incorporates a number of important model subclasses:

1. the convolutive single-input, multi-output (SIMO) model when  $K = 1$ ,

---

<sup>1</sup>Under the Gaussian assumption, an unknown variance parameter decouples from the unknown mean parameters in the Fisher information. So, while the noise power will affect the performance potential, whether  $\sigma^2$  is known or not does not affect how the CRB (CCRB) depends on the parameters. It does, of course, affect estimation.

2. the convolutive single-input, single-output (SISO) model when  $K = M = 1$ ,
3. the memoryless, instantaneous mixing model when  $L_k = 0$  for all  $k = 1, \dots, K$ ,  
and
4. the calibrated model (with constraints and a transformation of parameters).

For generality, we consider the full (MIMO) model, but the results contained in this section also apply to the preceding subclasses.

#### 4.1.1 Equivalent Convolutive Mixture Models

##### 4.1.1.1 The Vector-Matrix Model

In vector-matrix notation, the model may be written as

$$\mathbf{y} = \sum_{k=1}^K \mathbf{H}_M^{(k)} \mathbf{s}^{(k)} + \mathbf{w} = \mathbf{H}_M \mathbf{s} + \mathbf{w} \quad (4.3)$$

where the observations are contained in  $\mathbf{y}^T = [\mathbf{y}_1^T, \dots, \mathbf{y}_M^T] \in \mathbb{C}^{NM}$  and the observations for each channel contained in  $\mathbf{y}_m^T = [y_m(0), \dots, y_m(N-1)] \in \mathbb{C}^N$  for  $m = 1, \dots, M$ ,  $\mathbf{s}^{(k)T} = [s^{(k)}(-L_k), \dots, s^{(k)}(N-1)] \in \mathbb{C}^{N+L_k}$  represents the input sequence of the  $k$ th source, the channel matrix is represented by the  $NM \times N + L_k$  matrix

$$\mathbf{H}_M^{(k)} = \begin{bmatrix} \mathbf{H}_{(1)}^{(k)} \\ \vdots \\ \mathbf{H}_{(M)}^{(k)} \end{bmatrix} \quad (4.4)$$

where the  $m$ th channel submatrix over the  $k$ th source is given by the  $N \times N + L_k$  matrix

$$\mathbf{H}_{(m)}^{(k)} = \begin{bmatrix} h_m^{(k)}(L_k) & h_m^{(k)}(L_k - 1) & \cdots & h_m^{(k)}(0) & & \\ & h_m^{(k)}(L_k) & h_m^{(k)}(L_k - 1) & \cdots & h_m^{(k)}(0) & \\ & & \ddots & & & \ddots \\ & & & h_m^{(k)}(L_k) & h_m^{(k)}(L_k - 1) & \cdots & h_m^{(k)}(0) \end{bmatrix} \quad (4.5)$$

for  $m = 1, \dots, M$ , and the noise vector is  $\mathbf{w}^T = [\mathbf{w}_1^T, \dots, \mathbf{w}_M^T] \in \mathbb{C}^{MN}$  with the noise for each channel given by  $\mathbf{w}_m^T = [w_m(0), \dots, w_m(N - 1)] \in \mathbb{C}^N$ . This particular vector-matrix ordering of the model can also be represented as

$$\mathbf{y} = \sum_{k=1}^K (\mathbf{I}_{M \times M} \otimes \mathbf{S}^{(k)}) \mathbf{h}^{(k)} + \mathbf{w} = \mathbf{S}_M \mathbf{h} + \mathbf{w} \quad (4.6)$$

where the input samples are now organized into an  $N \times L_k + 1$  Toeplitz matrix

$$\mathbf{S}^{(k)} = \begin{bmatrix} s^{(k)}(0) & s^{(k)}(-1) & \cdots & s^{(k)}(-L_k) \\ s^{(k)}(1) & s^{(k)}(0) & \cdots & s^{(k)}(-L_k + 1) \\ \vdots & \vdots & \ddots & \vdots \\ s^{(k)}(N - 1) & s^{(k)}(N - 2) & \cdots & s^{(k)}(N - L_k - 1) \end{bmatrix}$$

with  $\otimes$  being the Kronecker product, and the channel elements are vectorized as  $\mathbf{h}^{(k)T} = [\mathbf{h}_1^{(k)T}, \dots, \mathbf{h}_M^{(k)T}] \in \mathbb{C}^{M(L_k+1)}$  with  $\mathbf{h}_m^{(k)T} = [h_m^{(k)}(0), \dots, h_m^{(k)}(L_k)] \in \mathbb{C}^{L_k+1}$  for each  $k = 1, \dots, K$  and  $m = 1, \dots, M$ . The purpose of these alternative vector-matrix methods will become evident in the development of the Fisher information (see section 4.1.3) for the original model in (4.2).

#### 4.1.1.2 The Z transform model

Yet another alternative representation of this model, necessary to introduce relevant characteristics of the model parameters, and often referred to as being in

the  $Z$ -transform domain, is the  $M$ -variate stationary process

$$\mathbf{y}(n) = [\mathbf{H}(z)] * \mathbf{s}(n) + \mathbf{w}(n) \quad (4.7)$$

where  $\mathbf{H}(z)$  is a  $M \times K$  (global) transfer function (polynomial) defined by

$$\mathbf{H}(z) = [\mathbf{H}^{(1)}(z), \dots, \mathbf{H}^{(K)}(z)]$$

for nonzero  $z \in \mathbb{C}^*$  (including  $\infty$ ) with  $\mathbf{H}^{(k)}(z)$  being the  $k$ th source transfer function

$$\mathbf{H}^{(k)}(z) = \sum_{l=0}^{L_k} \mathbf{h}^{(k)}(l) z^{-l}$$

and with  $\mathbf{h}^{(k)}(l)^T = [h_1^{(k)}(l), h_2^{(k)}(l), \dots, h_M^{(k)}(l)] \in \mathbb{C}^M$  for  $l = 0, \dots, L_k$  and  $k = 1, \dots, K$ ,  $\mathbf{s}(n)^T = [s^{(1)}(n), s^{(2)}(n), \dots, s^{(K)}(n)] \in \mathbb{C}^K$  for  $n = 0, \dots, N-1$ , and  $\mathbf{y}(n)$  and  $\mathbf{w}(n)$  correspond to match the models in (4.3) and (4.6).

Hence,  $\mathbf{H}(z)$  is a polynomial matrix of the backward shift  $z^{-1}$ . The  $k$ th source transfer function is said to have a *common zero* if there exists a nonzero  $z_0 \in \mathbb{C}^*$  such that  $\mathbf{H}^{(k)}(z_0) = \mathbf{0}$ . If the transfer function does not have a common zero, then the polynomials  $\mathbf{H}_{(m)}^{(k)}(z) = \sum_{l=0}^{L_k} h_m^{(k)}(l) z^{-l}$ , for  $m = 1, \dots, M$ , are said to be *coprime*. The global transfer function is said to be *reducible* if there exists a nonzero  $z_0 \in \mathbb{C}^*$  such that  $\text{rank}(\mathbf{H}(z_0)) < K$ . If not, it is said to be *irreducible*. The global transfer function is said to be *column-reduced* if  $\lim_{z \rightarrow 0} \text{rank}(\mathbf{H}(z)\mathbf{Z}) = K$  where  $\mathbf{Z} = \text{diag}\{z^{L_1}, \dots, z^{L_K}\}$ . This necessarily implies for an irreducible and column-reduced transfer function, called a *minimum polynomial basis* [35], that  $M \geq K$ . The connection between irreducibility in the multi-source case and not having common zeros in the single source case is made clear in the following result.



**Theorem 4.1.**  $\mathbf{H}(z)$  is irreducible if and only if  $\sum_{k=1}^K \alpha_k \mathbf{H}^{(k)}(z)$  has no common zeros for any nontrivial complex-valued collection  $\alpha_1, \dots, \alpha_K$ .

*Proof.*  $\mathbf{H}(z)$  is reducible if and only if there exists some point  $z_0 \in \mathbb{C}^*$  such that  $\text{rank}(\mathbf{H}(z_0)) < K$ , which follows if and only if there exists nontrivial  $\beta_1, \dots, \beta_K \in \mathbb{C}$  such that  $\sum_{k=1}^K \beta_k \mathbf{H}^{(k)}(z_0) = \mathbf{0}$ , i.e., if and only if  $\sum_{k=1}^K \beta_k \mathbf{H}^{(k)}(z)$  has a common zero.  $\square$

The connection between the model in (4.3) and (4.7) is that the  $K$ -source,  $M$ -channel matrix  $\mathbf{H}_M$  is a column-rotated, generalized Sylvester matrix of the block Toeplitz matrix

$$\begin{bmatrix} \mathbf{H}(L) & \mathbf{H}(L-1) & \cdots & \mathbf{H}(0) & & & \\ & \mathbf{H}(L) & \mathbf{H}(L-1) & \cdots & \mathbf{H}(0) & & \\ & & \ddots & & & \ddots & \\ & & & \mathbf{H}(L) & \mathbf{H}(L-1) & \cdots & \mathbf{H}(0) \end{bmatrix}$$

where  $\mathbf{H}(l) = [\mathbf{h}^{(1)}(l), \dots, \mathbf{h}^{(K)}(l)]$  with  $\mathbf{h}^{(k)}(l)$  a null vector when  $l > L_k$  and  $L = \max_k L_k$ . Furthermore, under certain conditions, the rank of  $\mathbf{H}_M$  is determined by the characterization of the basis  $\mathbf{H}(z)$ .

**Theorem 4.2.** Assume  $M > K$  and  $N \geq \sum_{k=1}^K L_k$ . Then  $\mathbf{H}_M$  is full column rank if and only if  $\mathbf{H}(z)$  is a minimum polynomial basis.

*Proof.* This can be shown by a variant of the proof in Loubaton and Moulines [44, theorem 1].  $\square$

The input sequence  $\mathbf{s}^{(k)}$  is said to have  $p_k$  modes<sup>2</sup> if it can be written as a

---

<sup>2</sup>There are a number of alternative definition of modes [30, 43], [35, p. 168].

linear combination

$$s^{(k)}(n) = \sum_{i=1}^{p_k} c_i m_i^{n+L_k},$$

where  $c_i$ ,  $i = 1, \dots, p_k$  are complex-valued weights and  $m_i$ ,  $i = 1, \dots, p_k$  are the  $\mathbb{C}^*$ -valued roots of the polynomial

$$a(0) + a(1)z^{-1} + a(2)z^{-2} + \dots + a(p_k)z^{-p_k},$$

whose coefficients satisfy

$$\sum_{j=0}^{p_k} s^{(k)}(i+j)a(j) = 0$$

for  $i = -L_k, \dots, N - p_k - 1$ . Hence, the Toeplitz matrix

$$\mathbf{S}^{(k)}(n) = \begin{bmatrix} s^{(k)}(n) & s^{(k)}(n-1) & \dots & s^{(k)}(-L_k) \\ s^{(k)}(n+1) & s^{(k)}(n) & \dots & s^{(k)}(-L_k+1) \\ \vdots & & \ddots & \vdots \\ s^{(k)}(N-1) & s^{(k)}(N-2) & \dots & s^{(k)}(N-n-L_k-1) \end{bmatrix} \quad (4.8)$$

has rank  $\min\{N-n, p, n+L_k+1\}$  [76, lemma 1]. Note,  $\mathbf{S}^{(k)}(0) = \mathbf{S}^{(k)}$  and  $\mathbf{S}^{(k)}(-L_k) = \mathbf{s}^{(k)}$  from section 4.1.1.1. If one of the modes of  $\mathbf{s}^{(k)}$  is a common zero of the channel transfer function  $\mathbf{H}^{(k)}(z)$ , say  $m_v$ , then the channel makes no distinction between  $\mathbf{s}^{(k)}$  and  $\mathbf{s}^{(k)'}$ , defined by

$$s^{(k)'}(n) = \sum_{\substack{i=1 \\ i \neq v}}^{p_k} c_i m_i^{n+L_k},$$

even though  $\mathbf{s}^{(k)} \neq \mathbf{s}^{(k)'}$ . Ironically, if a channel lacks sufficient diversity, the more modes an input has leads to greater risk of lost information on the input but potentially less meaningful loss depending on the weights.

**Theorem 4.3.** The matrix

$$\mathbf{S}(n) = [\mathbf{S}^{(1)}(n) \quad \mathbf{S}^{(2)}(n) \quad \dots \quad \mathbf{S}^{(K)}(n)] \quad (4.9)$$

is full column rank only if  $N \geq (K+1)n + K + \sum_{j=1}^K L_j$ ,  $p_{\text{total}} \geq Kn + K + \sum_{k=1}^K L_k$ , and  $p_k \geq L_k + 1 + n$  for each  $k = 1, \dots, K$  and if  $N \geq (K+1)n + K + (K+2) \sum_{j=1}^K L_j$ ,  $p_{\text{total}} \geq K * n + K + (K+1) \sum_{k=1}^K L_k$ , and  $p_k \geq L_k + 1 + n + \sum_{j=1}^K L_j$ .

*Proof.* This is a variation of the results in [76, 1, 48, 49].  $\square$

In particular, for  $n = 0$  then  $\mathbf{S}(0)$  requires  $N \geq K + \sum_{j=1}^K L_j$ ,  $p_{\text{total}} \geq K + \sum_{j=1}^K L_j$ , and  $p_k \geq L_k + 1$ , for each  $k = 1, \dots, K$ , to be full rank. Or conversely,  $\mathbf{S}(0)$  is full rank if  $N \geq K + (K+2) \sum_{j=1}^K L_j$ ,  $p_{\text{total}} \geq K + (K+1) \sum_{k=1}^K L_k$ , and  $p_k \geq L_k + 1 + \sum_{j=1}^K L_j$ .

Before continuing to the development of the Fisher information for this convolutive mixture model, it is relevant to consider a notion of identifiability using the concepts of the Z-transform model in (4.7). However, this next section is a brief aside that may be skipped if not of interest to the reader.

#### 4.1.2 Strict Identifiability

The  $K$ -user,  $M$ -channel FIR system in (4.2) or (4.3) is said to be *strictly identifiable (SI)* if and only if

$$\mathbf{H}_M \mathbf{s} = \mathbf{H}'_M \mathbf{s}' \iff \mathbf{H}'(z) = \mathbf{H}(z) \mathbf{A} \text{ and } \mathbf{s}'(n) = \mathbf{A}^{-1} \mathbf{s}(n)$$

for some nonsingular matrix  $\mathbf{A}$ . The converse statement is always true, so strict identifiability depends on the conditional statement. The term strict identifiability is a misnomer. As is clear from the definition, the deterministic parameters are not “identifiable” when they are strictly identifiable. Instead, the parameters are

identifiable up to some minimal ambiguity. When this situation occurs, then it is possible to treat the channel (signal) components statistically to truly identify the deterministic signal (channel) parameters using stochastic approaches, e.g., the subspace method [24, 44]. In the context of this thesis, where the parameters are not treated as random but as deterministic, the reduction to a minimal ambiguity set also reduces the number of necessary constraints to eliminate it.

Strict identifiability in the convolutive channel model, it shall be shown, has some interesting connections with the Fisher information matrix. This is not surprising, given the results in section 3.2. It is perhaps also intuitive to expect that as strict identifiability is a notion of *near-identifiability*, then the corresponding Fisher information will satisfy some notion of *near-regularity*.

Abed-Meraim and Hua [1] developed necessary and sufficient conditions for strict identifiability in terms of characteristics in the Z-transform model, i.e., no channel zeros, the number of signal modes, more sources than sensors, etc. The following two theorems will be presented without proof, as they were in [1]. Proofs I developed can be found in [48].

**Theorem 4.4** (SI necessary conditions). The  $M$ -channel  $K$ -source FIR system is strictly identifiable only if<sup>3</sup>

(a)  $\mathbf{H}(z)$  is irreducible and column-reduced,

$$(b) \ p_{\text{total}} \geq K + \sum_{j=1}^K L_j,$$

---

<sup>3</sup>In [1], condition (a) omits the column-reducedness requirement, condition (c) did not include the special case, and condition (d) is originally (and I think erroneously) stated  $N \geq 2K + \sum_{j=1}^K L_j$ .

(c)  $p_k \geq L_k + 2$  for  $k = 1, \dots, K$  of  $p_k \geq 1$  if  $L_k = 0$ , and

(d)  $N \geq K + \sum_{j=1}^K L_j.$

**Theorem 4.5** (SI sufficiency conditions). The  $M$ -channel  $K$ -source FIR system is strictly identifiable if

(a)  $\mathbf{H}(z)$  is irreducible and column-reduced,

(b)  $p_{\text{total}} \geq K + (K + 1) \sum_{j=1}^K L_j,$

(c)  $p_k \geq L_k + 1 + \sum_{j=1}^K L_j$  for  $k = 1, \dots, K$ , and

(d)  $N \geq K + (K + 2) \sum_{j=1}^K L_j.$

Yet other notions of near-identifiability in the convolutive channel model exist in the literature, e.g., *cross-relation-based identifiability* [43], which have an equivalence to strict identifiability in the SIMO case [30]. The necessary and sufficient conditions presented here are independent of the number of channels, although one should expect that increasing the number of channels would increase channel diversity thereby weakening the requirements on the other channel or source characteristics. Hence, theorems 4.4 and 4.5 are in some sense conditions for the least diversified, strictly identifiable scenario  $M = K + 1$ . Table 4.1 shows the growth in necessary and sufficient conditions as the number of users increases for a fixed channel size and for fixed channel orders.

Table 4.1: Necessary and sufficient data sizes for strict identifiability

# of sources $K$	1	2	3	4	5
necessary data size $N$	6	12	18	24	30
sufficient data size $N$	16	42	78	124	180

( $M = 6$  channels,  $L_k = 5$  for all sources)

### 4.1.3 The Fisher information of the convolutive mixture model

#### 4.1.3.1 Complex-valued Fisher information

Before presenting the Fisher information on the model in (4.2), certain details about Fisher information matrices on complex-valued parameters are necessary. The structure of the FIM depends on the ordering of the parameters. For complex-valued parameters, the FIM parameter vector consists of the real and imaginary parts of the parameters. If the complex-valued parameters are collected in the vector  $\boldsymbol{\vartheta}$  and the (FIM) real parameter vector is defined to be  $\boldsymbol{\theta}^T = [\text{Re}(\boldsymbol{\vartheta}^T), \text{Im}(\boldsymbol{\vartheta}^T)]$  and the real-valued parameter FIM has the structure

$$\mathbf{I}(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{E}_1(\boldsymbol{\theta}) & -\mathbf{E}_2(\boldsymbol{\theta}) \\ \mathbf{E}_2(\boldsymbol{\theta}) & \mathbf{E}_1(\boldsymbol{\theta}) \end{bmatrix}, \quad (4.10)$$

then the *complex-valued parameter Fisher information matrix* (CFIM) may be defined as<sup>4</sup>

$$\mathcal{I}(\boldsymbol{\vartheta}) = E_{\boldsymbol{\vartheta}} \mathbf{f}(\mathbf{x}; \boldsymbol{\vartheta}) \mathbf{f}^H(\mathbf{x}; \boldsymbol{\vartheta}) = \frac{1}{2} (\mathbf{E}_1(\boldsymbol{\theta}) + j \cdot \mathbf{E}_2(\boldsymbol{\theta}))$$

where  $\mathbf{f}(\mathbf{x}; \boldsymbol{\vartheta}) = \frac{\partial \log p(\mathbf{x}; \boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}^*}$ .<sup>5</sup> A number of properties for the real-valued parameter FIM can be gleaned from properties on the CFIM.

<sup>4</sup>This CFIM is a submatrix of the complex-valued parameter FIM developed by van den Bos [70], which would be the preferred FIM for use in a performance metric or in applying constraints [32]. However, in this and the next sections, the CFIM presented here is only used to obtain properties relevant to the real-valued parameter FIM.

<sup>5</sup>The complex derivative is defined to be  $\frac{\partial}{\partial \boldsymbol{\vartheta}} = \frac{\partial}{\partial \text{Re} \boldsymbol{\vartheta}} - j \cdot \frac{\partial}{\partial \text{Im} \boldsymbol{\vartheta}}$  in this thesis.

**Theorem 4.6.** The null space of the FIM of the form (4.10) has dimension exactly twice the dimension of the null space of the corresponding CFIM, i.e.,

$$\text{nullity}(\mathbf{I}(\boldsymbol{\theta})) = 2 \cdot \text{nullity}(\mathcal{I}(\boldsymbol{\vartheta})).$$

*Proof.* Note  $\begin{bmatrix} \mathbf{E}_1 & -\mathbf{E}_2 \\ \mathbf{E}_2 & \mathbf{E}_1 \end{bmatrix} \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix} = \mathbf{0}$  if and only if  $\mathbf{E}_1 \mathbf{a} - \mathbf{E}_2 \mathbf{b} = \mathbf{0}$  and  $\mathbf{E}_2 \mathbf{a} + \mathbf{E}_1 \mathbf{b} = \mathbf{0}$  if and only if  $(\mathbf{E}_1 + j\mathbf{E}_2)(\mathbf{a} + j\mathbf{b}) = (\mathbf{E}_1 \mathbf{a} - \mathbf{E}_2 \mathbf{b}) + j(\mathbf{E}_2 \mathbf{a} + \mathbf{E}_1 \mathbf{b}) = \mathbf{0}$ . Also,  $\begin{bmatrix} \mathbf{E}_1 & -\mathbf{E}_2 \\ \mathbf{E}_2 & \mathbf{E}_1 \end{bmatrix} \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix} = \mathbf{0}$  if and only if  $\begin{bmatrix} \mathbf{E}_1 & -\mathbf{E}_2 \\ \mathbf{E}_2 & \mathbf{E}_1 \end{bmatrix} \begin{bmatrix} -\mathbf{b} \\ \mathbf{a} \end{bmatrix} = \mathbf{0}$ . Finally, the dimension of  $\begin{bmatrix} \mathbf{a} & -\mathbf{b} \\ \mathbf{b} & \mathbf{a} \end{bmatrix}$  is 2 unless  $\mathbf{a} = \mathbf{b} = \mathbf{0}$ ; however, the dimension of  $[\mathbf{a} + j\mathbf{b}, -\mathbf{b} + j\mathbf{a}]$  is only 1 since  $(-\mathbf{b} + j\mathbf{a}) = -j \cdot (\mathbf{a} + j\mathbf{b})$ .  $\square$

A version of Bang's formula [36, a variant of (15.52) on p. 525] can be developed for complex-valued parameters, i.e.,

$$\begin{aligned} \mathcal{I}_{ij}(\boldsymbol{\vartheta}) &= \text{tr} \left[ \mathbf{C}^{-1}(\boldsymbol{\vartheta}) \frac{\partial \mathbf{C}(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}_i^*} \mathbf{C}^{-1}(\boldsymbol{\vartheta}) \frac{\partial \mathbf{C}(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}_j} \right] \\ &\quad + \frac{\partial \boldsymbol{\mu}^H(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}_i^*} \mathbf{C}^{-1}(\boldsymbol{\vartheta}) \frac{\partial \boldsymbol{\mu}(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}_j} + \frac{\partial \boldsymbol{\mu}^T(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}_i^*} \mathbf{C}^{-1}(\boldsymbol{\vartheta}) \frac{\partial \boldsymbol{\mu}^*(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}_j} \end{aligned}$$

for any observations  $\mathbf{y} \sim \mathcal{CN}(\boldsymbol{\mu}(\boldsymbol{\vartheta}), \mathbf{C}(\boldsymbol{\vartheta}))$ .

#### 4.1.3.2 CFIM for the convolutive mixture model

For the convolutive mixture model in (4.3) and (4.6), only the mean vector depends on the unknown parameters, hence

$$\mathcal{I}(\boldsymbol{\vartheta}) = \frac{\partial \boldsymbol{\mu}^H(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}^*} \mathbf{C}^{-1} \frac{\partial \boldsymbol{\mu}(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}^T} + \frac{\partial \boldsymbol{\mu}^T(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}^*} \mathbf{C}^{-1} \frac{\partial \boldsymbol{\mu}^*(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}^T}.$$

The mean vector is

$$\boldsymbol{\mu}(\boldsymbol{\vartheta}) = \sum_{k=1}^K \mathbf{H}_M^{(k)} \mathbf{s}^{(k)} = \sum_{k=1}^K (\mathbf{I}_M \otimes \mathbf{S}^{(k)}) \mathbf{h}^{(k)}.$$

Therefore, if the complex-valued parameter vector in the model in (4.3) is defined by

$$\boldsymbol{\vartheta}^T = \begin{bmatrix} \mathbf{h}^{(1)T} \\ \mathbf{s}^{(1)T} \\ \vdots \\ \mathbf{h}^{(K)T} \\ \mathbf{s}^{(K)T} \end{bmatrix} \quad (4.11)$$

then the complex-valued FIM of the model can be shown to be

$$\mathcal{I}(\boldsymbol{\vartheta}) = \frac{1}{\sigma^2} \boldsymbol{\mathcal{Q}}^H \boldsymbol{\mathcal{Q}} = \frac{1}{\sigma^2} \begin{bmatrix} \mathbf{Q}_1^H \mathbf{Q}_1 & \mathbf{Q}_1^H \mathbf{Q}_2 & \cdots & \mathbf{Q}_1^H \mathbf{Q}_K \\ \mathbf{Q}_2^H \mathbf{Q}_1 & \mathbf{Q}_2^H \mathbf{Q}_2 & \cdots & \mathbf{Q}_2^H \mathbf{Q}_K \\ \vdots & \vdots & \cdots & \vdots \\ \mathbf{Q}_K^H \mathbf{Q}_1 & \mathbf{Q}_K^H \mathbf{Q}_2 & \cdots & \mathbf{Q}_K^H \mathbf{Q}_K \end{bmatrix} \quad (4.12)$$

where  $\boldsymbol{\mathcal{Q}} = [\mathbf{Q}_1, \dots, \mathbf{Q}_K]$  and  $\mathbf{Q}_k = \begin{bmatrix} \mathbf{I}_M \otimes \mathbf{S}^{(k)} & \mathbf{H}_M^{(k)} \end{bmatrix}$ .

#### 4.1.3.3 Properties of the CFIM

In this section, I develop properties of this CFIM. In particular, the singularity of the CFIM is proven and a limit on the dimension of this singularity is detailed, as well as necessary and sufficient conditions on the parameter characteristics to attain this limit.

Given the inherent relationship between regularity of the FIM and identifiability as detailed in section 3.2, it should not be surprising that the FIM, and hence CFIM, is singular. The model presented in (4.3) or (4.6) has a multiplicative ambiguity with any source interacting with its corresponding channel, i.e., for any nonsingular matrix  $\mathbf{A} \in \mathbb{R}^{L_k+1 \times L_k+1}$ , the input source matrix  $\mathbf{S}^{(k)}$  and the channel vector  $\mathbf{h}_m^{(k)}$  are indistinguishable from  $\mathbf{S}^{(k)} \mathbf{A}$  and  $\mathbf{A}^{-1} \mathbf{h}_m^{(k)}$ , respectively. Over all the channels, this presumes a complex-valued multiplicative ambiguity of at least  $ML_k$ . Additionally, cross-source ambiguities exist between a source input and a different



source channel. The limit for the minimal degrees of freedom or the rank of the ambiguity is given in the following theorem.

**Theorem 4.7.** The CFIM is singular and the dimension of its null space is lower bounded as

$$\text{nullity}(\mathcal{I}(\boldsymbol{\vartheta})) \geq \sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+, \quad (4.13)$$

where  $(a)_+ = a$  for  $a \geq 0$  and  $(a)_+ = 0$  for  $a < 0$ . This limit quantity is the *nullity lower bound (NLB)*.

The proof can be found in the appendix C. The proof constructs the following matrix, whose linearly independent columns are a basis for the null space of  $\mathcal{I}(\boldsymbol{\vartheta})$ ,

$$\mathcal{N} = \begin{bmatrix} \mathbf{h}^{(1)} & \boldsymbol{\mathcal{H}}_{(1)}^{(2)} & \mathbf{0} & \mathbf{0} & \cdots & \boldsymbol{\mathcal{H}}_{(1)}^{(K)} & \mathbf{0} & \cdots & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ -\mathbf{s}^{(1)} & \mathbf{0} & \mathbf{0} & -\boldsymbol{\mathcal{S}}_{(1)}^{(2)} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \cdots & \mathbf{0} & -\boldsymbol{\mathcal{S}}_{(1)}^{(K)} \\ \mathbf{0} & \mathbf{0} & \mathbf{h}^{(2)} & \boldsymbol{\mathcal{H}}_{(2)}^{(1)} & \cdots & \mathbf{0} & \boldsymbol{\mathcal{H}}_{(2)}^{(K)} & \cdots & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & -\boldsymbol{\mathcal{S}}_{(2)}^{(1)} & -\mathbf{s}^{(2)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \cdots & -\boldsymbol{\mathcal{S}}_{(2)}^{(K)} & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \cdots & \vdots & \cdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{h}^{(K)} & \cdots & \boldsymbol{\mathcal{H}}_{(K)}^{(2)} & \boldsymbol{\mathcal{H}}_{(K)}^{(1)} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & -\boldsymbol{\mathcal{S}}_{(K)}^{(1)} & -\boldsymbol{\mathcal{S}}_{(K)}^{(2)} & \cdots & -\mathbf{s}^{(K)} & \cdots & \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad (4.14)$$

where  $\boldsymbol{\mathcal{H}}_{(j)}^{(i)}$  and  $\boldsymbol{\mathcal{S}}_{(j)}^{(i)}$  are defined in (C.2) and (C.4), respectively.

The nullity lower bound (NLB) seems to be an unusual quantity. The term  $(L_i - L_j + 1)_+$  represents the ambiguity for the “tap window” of the  $i$ th channel masking the interaction of the  $j$ th channel with its corresponding source. This quantity also reveals that an ambiguity exists only if the coherence-propagation delay of the  $i$ th channel is at least as great as for the  $j$ th channel. Intuitively, the greater the spread between the channel orders increases the overlapping window for the  $j$ th channel’s interaction to be masked and thereby increases the ambiguity. Conversely, when the channel orders have narrow differences, this increases the

diversity of the channel by limiting the window where one channel can cover that of another.

A simple corollary follows stating that the NLB or the dimension of the ambiguity set is at least the square of the number of sources.

**Corollary 4.8.**  $\text{nullity}(\mathcal{I}(\boldsymbol{\vartheta})) \geq K^2$ .

*Proof.* Note  $(L_i - L_j + 1)_+ + (L_j - L_i + 1)_+ \geq 2$  with equality if and only if  $|L_i - L_j| \leq 1$ . Thus,

$$\begin{aligned} \sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+ &= \sum_{i=1}^K \sum_{\substack{j=1 \\ j \neq i}}^K (L_i - L_j + 1)_+ + \sum_{i=1}^K (L_i - L_i + 1)_+ \\ &= \sum_{i=1}^K \sum_{j=i+1}^K (L_i - L_j + 1)_+ + (L_j - L_i + 1)_+ + K \\ &\geq \frac{K(K-1)}{2} \cdot 2 + K = K^2. \end{aligned}$$

□

The proof also reveals that this  $K^2$  degree of uncertainty can only be attained if any two orders differ by at most 1 tap, which agrees with the intuition that channel diversity is enhanced when the channel orders are not widely spread.

In addition to the ambiguity due to the channel order spread, the degrees of freedom in the model depends on the number of parameters and the number of observations.

**Theorem 4.9.**  $\text{nullity}(\mathcal{I}(\boldsymbol{\vartheta})) \geq \sum_{k=1}^K (N + L_k + M(L_k + 1)) - MN$ .

*Proof.* Since  $\mathcal{I}(\boldsymbol{\vartheta}) = \frac{1}{\sigma^2} \boldsymbol{\mathcal{Q}}^H \boldsymbol{\mathcal{Q}}$  then  $\text{rank}(\mathcal{I}(\boldsymbol{\vartheta})) = \text{rank}(\boldsymbol{\mathcal{Q}})$ . Since  $\boldsymbol{\mathcal{Q}}$  is a  $MN \times \sum_{k=1}^K (N + L_k + M(L_k + 1))$ , then  $\text{rank}(\mathcal{I}(\boldsymbol{\vartheta})) \leq \min \left\{ MN, \sum_{k=1}^K (N + L_k + M(L_k + 1)) \right\}$

and therefore

$$\begin{aligned} \text{nullity}(\mathcal{I}(\boldsymbol{\vartheta})) &\geq \text{colsize}(\mathcal{I}(\boldsymbol{\vartheta})) - \text{rank}(\mathcal{I}(\boldsymbol{\vartheta})) \\ &\geq \max \left\{ 0, \sum_{k=1}^K (N + L_k + M(L_k + 1)) - MN \right\}. \end{aligned}$$

□

The number of columns of  $\mathcal{Q}$  (or  $\mathcal{I}(\boldsymbol{\vartheta})$ ) corresponds to the number of unknown parameters. Likewise, the number of rows of  $\mathcal{Q}$  corresponds to the number of equations (observations) in the model (4.2). In most scenarios, having more equations (rows) than unknowns (columns) is a necessary requirement for the unknowns to be solvable. Thus, the degrees of freedom of the ambiguity space is at least the number of unknowns (parameters) less the number of equations. However, if the equations are linearly dependent (redundant), as is the case in the convolutive mixture model, then the degrees of freedom is potentially greater. The following corollary combines theorems 4.7 and 4.9.

**Corollary 4.10.**

$$\text{nullity}(\mathcal{I}(\boldsymbol{\vartheta})) \geq \max \left\{ \sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+, \sum_{k=1}^K (N + L_k + M(L_k + 1)) - MN \right\}.$$

The NLB depends only on the channel orders and number of sources, whereas the “equations vs. unknowns” nullity bound depends on the channel orders, the sample size per source, the number of subchannels per source, and the number of sources. Hence, it is of interest to determine under what conditions on the sample size  $N$  and the number of channels  $M$  per source will this second bound be of no consequence. Control over the dimensions  $K$  (the number of transmitters),  $M$

(the number of receivers),  $N$  (the number of time snapshots or transmission length) is possible in the design of many communications models. The next two results determine conditions on the dimensions of the model which allow the NLB to be the minimum possible degrees of freedom. The first condition requires more subchannels per source than sources (or more receivers than transmitters in a communications context).

**Theorem 4.11.** The CFIM  $\mathcal{I}(\boldsymbol{\vartheta})$  can attain the CFIM nullity lower bound only if  $M > K$ .

*Proof.* If  $M \leq K$  then  $\mathbf{H}_M$ , which is a  $MN \times KN + \sum_{k=1}^K L_k$  matrix, cannot be full column rank by theorem 4.2. This implies the existence of a null space of a larger dimension than the space spanned by the columns of  $\mathcal{N}$  in (4.14). Hence  $\text{nullity}(\mathcal{I}(\boldsymbol{\vartheta})) > \sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+$ .  $\square$

Under the assumption that more receivers than transmitters are in the model, then corollary 4.10 implies a minimal requirement on the transmission snapshots.

**Theorem 4.12.** Provided  $M > K$ , the CFIM  $\mathcal{I}(\boldsymbol{\vartheta})$  can attain the CFIM nullity lower bound only if

$$N \geq \frac{K+1}{M-K} \sum_{j=1}^K L_j + \sum_{j=1}^K L_j + K + \frac{1}{M-K} \left( K^2 - \sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+ \right) \quad (4.15)$$

*Proof.* We desire  $\sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+ \geq \sum_{k=1}^K (N + L_k + M(L_k + 1)) - MN$ . Solving for  $N$  shows the result.  $\square$

As the channel diversity (in terms of the number of channels per source) increases, the necessary size on the data to attain the CFIM NLB decreases. As

$M \rightarrow \infty$  then the bound on the sample size becomes simply  $N \geq K + \sum_{j=1}^K L_j$ , which is comparable to theorem 4.4(d). However, for any finite  $M$  and nontrivial channel orders, the bound is effectively  $N \geq K + \sum_{j=1}^K L_j + 1$  since  $-\sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+ \geq -K \sum_{j=1}^K (L_j + 1)$ , and hence a looser bound than (4.15) is

$$\begin{aligned} N &\geq \frac{K+1}{M-K} \sum_{j=1}^K L_j + \sum_{j=1}^K L_j + K + \frac{1}{M-K} \left( K^2 - K \sum_{j=1}^K (L_j + 1) \right) \\ &\geq \sum_{j=1}^K L_j + K + \frac{1}{M-K} \sum_{j=1}^K L_j. \end{aligned}$$

In the scenario requiring the most data samples,  $M = K + 1$ , then the system requires  $N \geq (K+1) \sum_{j=1}^K L_j + \sum_{j=1}^K L_j + K + \left( K^2 - \sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+ \right)$ . (This last term is always nonpositive.) In the SIMO scenario ( $K = 1$ ), this necessary condition becomes  $N \geq 3L + 1$  which agrees with the requirement in [31].<sup>6</sup>

Given that  $M > K$  and  $N$  satisfies (4.15), then the minimum possible nullity of  $\mathcal{I}(\boldsymbol{\vartheta})$  is the nullity lower bound,

$$\sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+,$$

i.e., it is possible to limit the ambiguity strict to the mixing of sources in the convolutive channel. If it is possible to understand the conditions under which the CFIM attains this nullity lower bound, then it is known that the null space is completely characterized by the columns of  $\boldsymbol{\mathcal{N}}$  in (4.14). The importance of this is that when

---

<sup>6</sup>For the SIMO ( $K = 1$ ) case, this nullity lower bound is exactly 1 and parameters for which the CFIM attains this bound are said to be *Fisher information identifiable* in [30]. The notion of identifiability for the SIMO case is sensible since by scaling a single parameter it is possible to obtain a bound on all the remaining parameters relative to the scaled parameter. This naturally connects with the notion of strict identifiability in section 4.1.2. This interpretation does not extend simply to the MIMO ( $K > 1$ ) case; however, it shall be shown that there does exist an inherent connection between a CFIM attaining the NLB and the notion of strict identifiability.

the null space can be parameterized, then it is possible to use theorems 3.23 and 3.24 to determine constraints that lead to regularity of the CFIM (and FIM). Using the concepts of signal excitation (modes) and channel diversity (irreducibility and column-reducedness) as defined in section 4.1.1.2, the following theorems also establish a correlation of conditions between the idea of near-regularity when the FIM attains the NLB and of near-identifiability (or strict identifiability) of section 4.1.2.

**Theorem 4.13** (CFIM NLB necessary conditions). The  $M$ -channel  $K$ -source FIR system Fisher information matrix has a nullity of exactly the NLB in (4.13) only if

- (a)  $\mathbf{H}(z)$  is irreducible and column-reduced,
- (b)  $p_{\text{total}} \geq K + \sum_{j=1}^K L_j$ ,
- (c)  $p_k \geq L_k + 2$  for  $k = 1, \dots, K$  or  $p_k \geq 1$  if  $L_k = 0$ ,
- (d)  $N \geq K + \sum_{j=1}^K L_j$ , and
- (e)  $M > K$ .

**Theorem 4.14** (CFIM NLB sufficiency conditions). The  $M$ -channel  $K$ -source FIR system FIM has a nullity of exactly the NLB in (4.13) if

- (a)  $\mathbf{H}(z)$  is irreducible and column-reduced,
- (b)  $p_{\text{total}} \geq K + (K + 1) \sum_{j=1}^K L_j$ ,
- (c)  $p_k \geq L_k + 1 + \sum_{j=1}^K L_j$  for  $k = 1, \dots, K$ ,

Table 4.2: Necessary and sufficient data sizes for the FIM to attain the NLB

# of sources $K$	1	2	3	4	5
necessary data size $N$	8	20	38	74	180
sufficient data size $N$	16	42	78	124	180

( $M = 6$  channels,  $L_k = 5$  for all sources)

(d)  $N \geq K + (K + 2) \sum_{j=1}^K L_j$ , and

(e)  $M > K$ .

With the exception of condition (e), then theorems 4.13 and 4.14 are identical to theorems 4.4 and 4.5, respectively. This is an expected result since the Fisher information for identifiable parametric components within the model should meet certain regularity conditions (see theorem 2.5). For theorem 4.13(d), the condition is weaker than the condition in theorem 4.12, but it is left weaker to agree with the necessary condition for strict identifiability in theorem 4.4(d). (It is unclear if that condition would change if dependence on channel size  $M$  is considered.) Table 4.1.3.3 show the growth in the necessary and sufficient data size  $N$  as the number of sources increases for fixed channel size and channel orders, using the stronger condition in theorem 4.12 instead of that in theorem 4.13(d). It is clear that in comparison, as the number of sources increases the number of channels, the necessary data size grows to match the sufficient data size, which lead to theorem 4.15 and corollary 4.16.

First, the special case when the CFIM attains the minimal ambiguity is considered. As stated earlier this occurs when the windows of the channels are, in some sense, minimally spread.

**Theorem 4.15.** Given sufficient channel diversity and modes in the signals, e.g., the sufficient conditions of theorem 4.14, then  $\text{nullity}(\mathcal{I}(\vartheta)) = K^2$  if and only if  $L_j \in \{L_0, L_0 + 1\}$  for all  $j = 1, \dots, K$ , for some integer  $L_0$ .

*Proof.* Under the conditions of theorem 4.14, then  $\text{nullity}(\mathcal{I}(\vartheta)) = \text{NLB}$  and from the proof of corollary 4.8,  $\text{NLB} = K^2$  if and only if  $(L_i - L_j + 1)_+ + (L_j - L_i + 1)_+ \equiv 2$  for each  $i, j = 1, \dots, K$ . Either  $(L_i - L_j + 1)_+ = (L_j - L_i + 1)_+ = 1$  or  $(L_i - L_j + 1)_+ = 2$  and  $(L_j - L_i + 1)_+ = 0$  or vice versa.  $\square$

The next corollary details the necessary and sufficient data size for the sources under the special case of equivalent channel orders.

**Corollary 4.16.** Given sufficient signal modes and channel diversity, e.g., as in theorem 4.14, if  $L_k = L$  for each  $k$  and  $M = K + 1$ , then the CFIM attains the NLB if and only if

$$N \geq (K + 2)KL + K.$$

#### 4.1.4 Constraints for the convolutive mixture model

In the previous section, necessary and sufficient conditions on characteristics of the signal source and channel properties were detailed for the CFIM to have a null space with minimal dimensions equaling the nullity lower bound  $\sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+$ . In this section, pathways to regularity in the CFIM (and hence the FIM) as well as several typical constraint sets are considered.



#### 4.1.4.1 Pathways to regularity

Requiring constraints to attain regularity in the Fisher information is unnecessary to use this CCRB theory (theorem 3.17), but often without constraints the identifiable parameters are not in the desired framework to be of use to the practitioner. In the communications context, constraints on a model assumed a priori are simply a common method used to maintain a particular parametric structure in the model.

Guided by theorems 3.23, 3.24, 3.25, 3.27, and 3.29, then an objective of communications system design is to discover constraints for which identifiability of the parameter, as they are defined, is achieved, i.e., what properties can be imposed on the signal or channel to guarantee parametric identifiability. For the convolutive mixture model in (4.2), this involves an examination of the Fisher information in (4.10) with the CFIM in (4.12). In some sense, this has already been done. In section 4.1.3.3, conditions under which the CFIM has a null space spanned by the columns of the matrix  $\mathcal{N}$ , defined in (4.14), are derived.

If  $\mathbf{I}(\boldsymbol{\theta})$  is a Fisher information with Cholesky factorization  $\mathbf{L}(\boldsymbol{\theta})\mathbf{L}^T(\boldsymbol{\theta})$  where  $\mathbf{V}(\boldsymbol{\theta})$  is an orthogonal complement to  $\mathbf{L}^T(\boldsymbol{\theta})$  [20, p.194]. Then theorem 3.23 implies for any constraint to achieve local identifiability, its Jacobian must satisfy  $\mathbf{F}(\boldsymbol{\theta}) = \mathbf{A}\mathbf{L}^T(\boldsymbol{\theta}) + \mathbf{B}\mathbf{V}^T(\boldsymbol{\theta})$  where  $\mathbf{B}\mathbf{V}^T$  has full row rank. Similarly, theorem 3.24 implies the constraints must have an orthonormal complement  $\mathbf{U}(\boldsymbol{\theta})$  satisfying  $\mathbf{U}(\boldsymbol{\theta}) = \mathbf{L}(\boldsymbol{\theta})\mathbf{C} + \mathbf{V}(\boldsymbol{\theta})\mathbf{D}$  where  $\mathbf{L}(\boldsymbol{\theta})\mathbf{C}$  if full column rank.

#### 4.1.4.2 Norm channel + real-valued source constraint

The norm channel constraint

$$\|\mathbf{h}^{(1)}\|^2 = 1 \quad (4.16)$$

is a common scaling “trick” used in SIMO ( $K = 1$  source) models to obtain channel estimates under second-order statistics assumptions. Combined with the (rotational) constraint of restricting the source elements to be real-valued, i.e.,

$$\text{Im}(s^{(1)}(n)) = 0 \quad (4.17)$$

for  $n = -L_1, -L_1 + 1, \dots, N - 1$ , it is clear that the multiplicative ambiguity which is the basis of  $\mathcal{N}$  is eliminated. For the parameter vector  $\boldsymbol{\theta} = [\text{Re}(\boldsymbol{\vartheta}^T), \text{Im}(\boldsymbol{\vartheta}^T)]^T$  with  $\boldsymbol{\vartheta}$  defined in (4.11), the constraints are essentially separated as

$$\mathbf{F}(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{F}_{\text{Re}(\mathbf{h})} & \mathbf{0} & \mathbf{F}_{\text{Im}(\mathbf{h})} & \mathbf{0} \\ \mathbf{0} & \mathbf{F}_{\text{Re}(\mathbf{s})} & \mathbf{0} & \mathbf{F}_{\text{Im}(\mathbf{s})} \end{bmatrix}$$

where  $\mathbf{h} = \mathbf{h}^{(1)}$  and  $\mathbf{s} = \mathbf{s}^{(1)}$ . An orthonormal complement, or a basis for the null space, of  $[\mathbf{F}_{\text{Re}(\mathbf{s})}, \mathbf{F}_{\text{Im}(\mathbf{s})}]$  is simply  $\begin{bmatrix} \mathbf{I}_{(N+L_1) \times (N+L_1)} \\ \mathbf{0} \end{bmatrix}$ , but an analytic formula for the null space relating to the channel components is not so simple and requires numerical programming for arbitrary sizes  $L_1$  and  $N$ .

**Example 4.17.** As a particular example of this constraint, consider the  $M = 2$  SIMO convolutive channel model with  $L = 3$  taps, where the channel is predefined by

$$\begin{bmatrix} \mathbf{h}_1^{(1)T} \\ \mathbf{h}_2^{(1)T} \end{bmatrix} = \begin{bmatrix} 0.3079 + j0.0698 & 0.1657 + j0.2304 & 0.0198 - j0.3823 & 0.0929 - j0.1853 \\ -0.1841 + j0.3294 & 0.4484 - j0.1689 & 0.0156 + j0.1526 & 0.4750 - j0.0952 \end{bmatrix},$$

under a white noise assumption (with known variance  $\sigma^2$ ) and the constraints in (4.16) and (4.17). The 54 ( $N = 50$ ) signal elements are BPSK ( $\pm 1$ ) symbols randomly generated with equal probability (i.e., Bernoulli( $\frac{1}{2}$ )) and fixed for the simulation. The subspace method [52, 68] is used with a smoothing factor of 8. Additionally scoring with constraints (CSA) using (3.36) is applied on the subspace method's estimate  $\hat{\mathbf{h}}^{(1)}$  for possible improvement. The results are shown in figure 4.2, where the mean-square error (MSE) per real channel coefficient is evaluated by  $\frac{1}{16R} \sum_{r=1}^R \left\| \hat{\mathbf{h}}^{(1)} - \mathbf{h}^{(1)} \right\|^2$  over  $R = 100$  runs or trials and the SNR is given by  $10 \log_{10} \frac{1}{\sigma^2}$ . This example is also in [68] and shows how the subspace method, which had no prior theory-based performance metric, tracks with the CCRB. The additional overlaying estimation using the method of scoring demonstrates the local maximum likelihood properties in the subspace method.

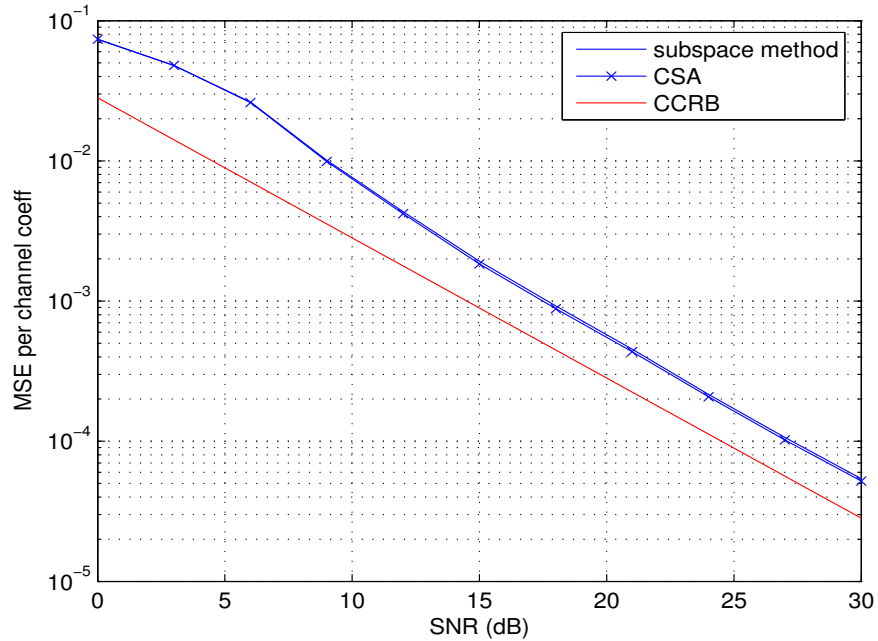


Figure 4.2: Norm-constrained channel estimation performance.

#### 4.1.4.3 Semiblind constraints: $s^{(k)}(t) = p(t)$ for $t \in \mathbb{T}$

For multiple sources, designing communications with channel constraints, such as in the previous section, is less tenable since it is not always possible to establish guarantees on functions of the channel elements. The sources, however, are often entirely designable elements by restricting the class of signals which are to be passed through the channels. The class only needs to be defined by some functional constraint.

Perhaps, the simplest constraint is knowledge of a parameter element, i.e.,  $\theta_i = a$ . This constraint produces a row vector  $\mathbf{e}_i^T$  in the Jacobian  $\mathbf{F}(\boldsymbol{\theta})$ , where  $\mathbf{e}_i$  is the unit vector with unity in the  $i$ th position and zero values in the other positions. Therefore, the corresponding orthonormal complement  $\mathbf{U}(\boldsymbol{\theta})$  for this unit vector eliminates or nulls out the  $i$  row and column of the Fisher information  $\mathbf{I}(\boldsymbol{\theta})$  while preserving the other elements. Any other nonredundant constraint in  $\mathbf{f}$  does not change this.

**Theorem 4.18.** Assume the conditions of theorem 4.14, then  $\text{nullity}(\mathcal{I}(\boldsymbol{\vartheta})) = \sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+$ . Let  $\mathcal{I}_*(\boldsymbol{\vartheta})$  be denoted by  $\mathcal{I}(\boldsymbol{\vartheta})$  with the rows  $\{i_p : p = 1, \dots, \sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+\}$  and the corresponding columns removed. Then the  $\text{nullity}(\mathcal{I}_*(\boldsymbol{\vartheta})) = 0$  if and only if  $\mathcal{N}_*$  is full column rank, where  $\mathcal{N}_*$  is the matrix formed by taking the rows  $\{i_p\}$  of  $\mathcal{N}$  in (4.14).

*Proof.* Let  $\mathbf{F}$  be a matrix that when multiplied by the matrix  $\mathcal{N}$  selects the rows  $\{i_p\}$ , i.e.,  $\mathbf{F}$  is the matrix consisting of the row vectors  $\{\mathbf{e}_i^T : i \in \{i_p\}\}$ . Then  $\mathbf{F}\mathcal{N} = \mathcal{N}_*$ , a square full rank matrix. Therefore  $\mathbf{F} = \mathbf{A}\mathbf{I}(\boldsymbol{\theta}) + \mathbf{B}\mathcal{N}^T$  for some

$\mathbf{A}$  and some  $\mathbf{B}$  where  $\mathbf{B}\mathbf{N}^T$  is full row rank. As discussed in section 4.1.4.1, this identifies the parameters and hence  $\text{nullity}(\mathcal{I}_*(\boldsymbol{\theta})) = 0$  by theorem 3.24.  $\square$

This theorem can be used to specify either channel or source parameters to achieve identifiability and FIM regularity. For example, if only source signal parameters are specified, then it is necessary and sufficient to specify  $\sum_{j=1}^K (L_i - L_j + 1)_+$  parameters of  $\mathbf{s}^{(i)}$  for each  $i = 1, \dots, K$ , under the conditions of theorem 4.18.

#### 4.1.4.4 Unit Modulus constraint + Semiblind constraint

The unit or constant modulus constraint is a particularly powerful and reasonable assumption. All the source elements are assumed to have unit modulus, i.e.,

$$|s^{(k)}(n)|^2 = 1, \quad (4.18)$$

for every  $n = -L_k, \dots, N-1$  and for each  $k = 1, \dots, K$ . This constraint is useful for modeling  $P$ -ary phase-shift keying (PSK) in communications models with unknown  $P$ , where the signals are assumed to be derived from a finite constellation of  $P$  equispaced points on the unit circle. While in practice, this assumption is often viewed as a single constraint, it is ultimately  $\sum_{i=1}^K (N + L_i)$  constraints. Nevertheless, despite this, it is insufficient to identify the parameters by itself. For the parameter vector in (4.11), the constraint has a gradient

$$\mathbf{F}(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{0} & 2\text{Re}(\mathbf{S}_d^{(1)}) & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & 2\text{Im}(\mathbf{S}_d^{(1)}) & \mathbf{0} & \mathbf{0} & \cdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & 2\text{Re}(\mathbf{S}_d^{(2)}) & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & 2\text{Im}(\mathbf{S}_d^{(2)}) & \cdots \\ & & & & \ddots & & & & & \vdots \end{bmatrix},$$

which has an orthonormal complement

$$\mathbf{U}(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{0} & \mathbf{I}_{M(L_1+1)} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots \\ -\text{Im}(\mathbf{S}_d^{(1)}) & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{I}_{M(L_2+1)} & \mathbf{0} & \dots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & -\text{Im}(\mathbf{S}_d^{(2)}) & \mathbf{0} & \mathbf{0} & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_{M(L_1+1)} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots \\ \text{Re}(\mathbf{S}_d^{(1)}) & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \text{Re}(\mathbf{S}_d^{(2)}) & \mathbf{0} & \mathbf{I}_{M(L_2+1)} & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

satisfying (3.6) and where  $\mathbf{S}_d^{(i)} = \text{diag}(\mathbf{s}^{(i)})$ . Using this complement with (4.10) and

(4.12), we have that the  $i, j$  subblock of  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  has the structure  $\mathbf{C}(i, j) =$

$$\begin{bmatrix} \text{Re}[\mathbf{S}_d^{(i)H} \mathbf{H}_M^{(i)H} \mathbf{H}_M^{(j)} \mathbf{S}_d^{(j)}] & \text{Im}[\mathbf{S}_d^{(i)H} \mathbf{H}_M^{(i)H} (\mathbf{I}_m \otimes \mathbf{S}^{(j)})] & \text{Re}[\mathbf{S}_d^{(i)H} \mathbf{H}_M^{(i)H} (\mathbf{I}_m \otimes \mathbf{S}^{(j)})] \\ -\text{Im}[(\mathbf{I}_m \otimes \mathbf{S}^{(i)H}) \mathbf{H}_M^{(j)} \mathbf{S}_d^{(j)}] & \text{Re}[\mathbf{I}_m \otimes \mathbf{S}^{(i)H} \mathbf{S}^{(j)}] & \text{Im}[\mathbf{I}_m \otimes \mathbf{S}^{(i)H} \mathbf{S}^{(j)}] \\ \text{Re}[\mathbf{H}_M^{(i)H} \mathbf{H}_M^{(j)} \mathbf{S}_d^{(j)}] & \text{Im}[\mathbf{H}_M^{(i)H} (\mathbf{I}_m \otimes \mathbf{S}^{(j)})] & \text{Re}[\mathbf{H}_M^{(i)H} (\mathbf{I}_m \otimes \mathbf{S}^{(j)})] \end{bmatrix}.$$

Lack of identifiability is noted since

$$\mathbf{C}(k, k) \cdot \begin{bmatrix} \mathbf{1}_{N+L_k} \\ \text{Im}(\mathbf{h}^{(k)*}) \\ \text{Re}(\mathbf{h}^{(k)*}) \end{bmatrix} = \mathbf{0}$$

for each  $k = 1, \dots, K$ . Careful examination determines these are the only vectors

in the null space of  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$ , resulting in the following theorem.

**Theorem 4.19.** Assume  $\text{nullity}(\mathbf{I}(\boldsymbol{\theta})) = 2 \sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+$ . If all sources are assumed to be unit modulus and one complex-valued parameter for each source is assumed known, then  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  will be regular (and the model locally identifiable).

Intuitively, the unit modulus constraint eliminates the convolutive mixture and intra-source multiplicative amplitude ambiguity leaving only an intra-source multiplicative phase ambiguity.

**Example 4.20.** Consider a  $K = 2$  source- $M = 2$  sensor linear instantaneous mixing model ( $L_k = 0$  for all sources), with  $N = 30$  samples per source, under a  $R = 3$ -ray multipath subchannels, i.e., the spatial signature of the  $k$ th source is expressed as a weighted sum of steering vectors, i.e.,  $\mathbf{h}_k = \sum_{r=1}^R \beta_{kr} \mathbf{a}(\psi_{kr})$  where  $\beta_{kr}$  and  $\psi_{kr}$  are the complex-valued amplitude and the real-valued AOA of the  $r$ th multipath of the  $k$ th source (see Figure 4.3). The AOAs and corresponding amplitudes are  $\{\psi-1, \psi, \psi+4\}$

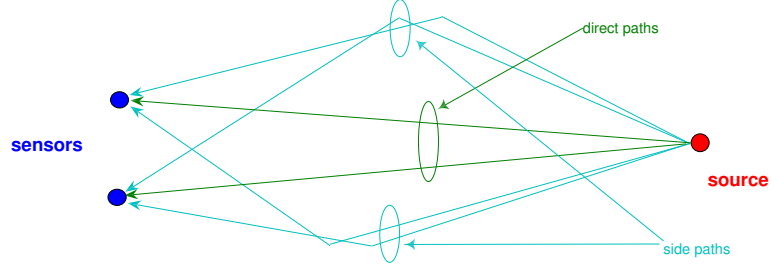


Figure 4.3: Example of multipath channel.

and  $\{\sqrt{0.2}\angle(-\frac{\pi}{6}), \sqrt{0.5}, \sqrt{0.15}\angle(-\frac{\pi i}{5})\}$  for source  $\mathbf{s}^{(1)}$  and  $\{\psi + \Delta\psi - 5, \psi + \Delta\psi, \psi + \Delta\psi + 6\}$  and  $\{\sqrt{0.15}\angle(-\frac{\pi i}{5}), \sqrt{0.6}, \sqrt{0.25}\angle(\frac{\pi}{3})\}$  for source  $\mathbf{s}^{(2)}$ . (See section 4.2 for a greater description on steering vectors.) The source elements come from an  $8PSK$  alphabet with signal powers  $\text{SNR}(\mathbf{s}^{(1)}) = 20\text{dB}$  and  $\text{SNR}(\mathbf{s}^{(2)}) = 15\text{dB}$ . The channel elements are normalized so that  $\text{SNR}(\mathbf{s}^{(k)}) = \frac{\|\mathbf{h}^{(k)}\|^2}{M\sigma^2}$  with  $\sigma^2 = 1$ . The constraints assumed are unit modulus ( $8PSK$ ) as well as knowledge of the first  $T = 2$  symbols for each source, more than sufficient for identifiability and FIM regularity (theorem 4.19).

An initial estimate is obtained using the zero-forcing (bias reducing) variant of the algebraic constant modulus algorithm (ZF-ACMA) [71]. This algorithm is a useful tool for estimation of constant modulus source parameters in short data length

experiments (only  $N \geq K^2$  or  $2K$  required), but has no means of incorporating the training side information. The ZF-ACMA estimate is projected onto  $\Theta_f$  to establish an initialization for the method of scoring with constraints. The step size rule chose  $\alpha^{(k)} = 2^{-m}$  for the least positive integer  $m$  satisfying (3.37). The average MSE (per real parameter coefficient) at each iteration over 5000 trails is compared with the CCRB for each source in figure 4.4.

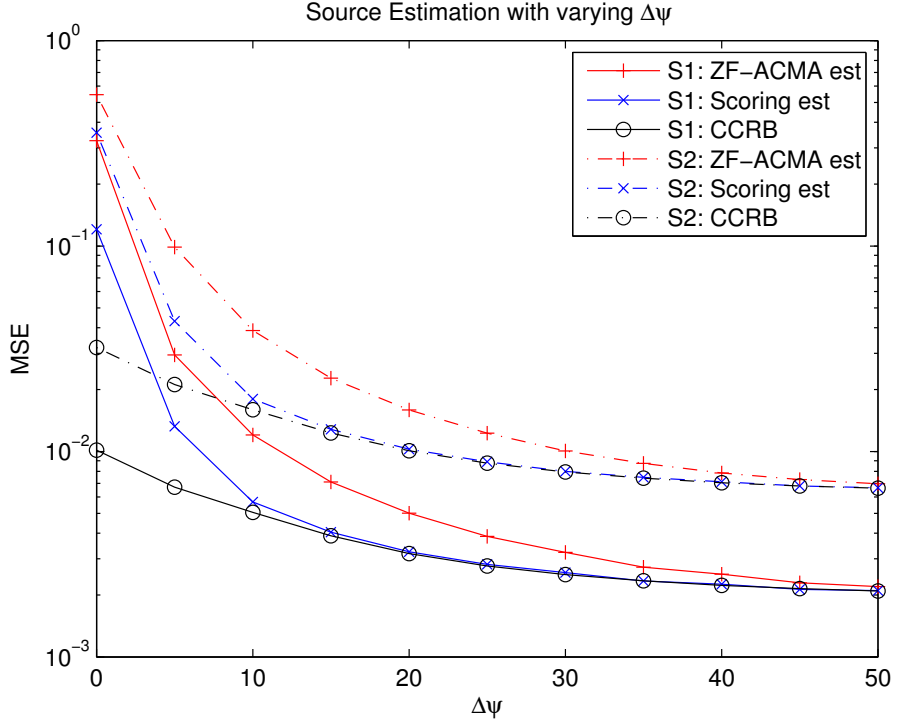


Figure 4.4: Source estimation with varying  $\Delta\psi$ .

The mean-square error improvement by utilizing the complete side information in scoring maintains efficiency with respect to the constrained CRB for moderately separated angle of arrivals compared with ZF-ACMA. As should be expected, the estimation performance degrades as the primary source angles overlap, but even in



the worst case scenario with  $\Delta\psi = 0^\circ$  of separation, the approximate corresponding 8PSK bit-error-rate (BER) for ZF-ACMA and scoring with constraints is .2914 and .0254, meaning the estimation schemes result in bit decision errors roughly 29% and 3% of the time, respectively. The departure of the estimation performance from the CCRB as  $\Delta\psi$  approaches  $0^\circ$  is possibly due to a loss of unbiasedness in the estimation.

## 4.2 Calibrated Array Model

The narrowband (calibrated) array model may be written as

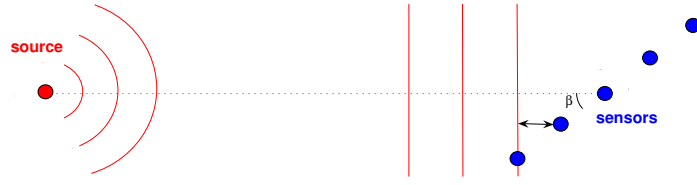
$$y_m(n) = \sum_{k=1}^K a_m(\omega_k) \gamma_k s^{(k)}(n) + w_m(n) \quad (4.19)$$

for  $n = 1, \dots, N$  and  $m = 1, \dots, M$ , where  $s^{(k)}(n)$  is the value of the  $k$ th input source at time index  $n$ ,  $\gamma_k$  is the complex-valued channel gain for the  $k$ th input,  $a_m(\omega_k)$  is the  $m$ th sensor response to the  $k$ th input source,  $\omega_k$  is the angle-of-arrival (AOA) of the  $k$ th source, and  $w_m(n)$  is the noise, modeled as zero-mean circular Gaussian with variance  $\sigma^2$  iid in both time and space. In vector-matrix notation, the model for each time slot can be written as

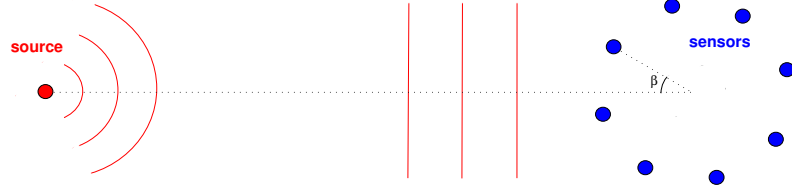
$$\mathbf{y}(n) = \sum_{k=1}^K \mathbf{a}(\omega_k) \gamma_k s^{(k)}(n) + \mathbf{w}(n) = \mathbf{A}(\boldsymbol{\omega}) \boldsymbol{\Gamma} \mathbf{s}(n) + \mathbf{w}(n) \quad (4.20)$$

where the input is given by  $\mathbf{s}^T(n) = [s^{(1)}(n), \dots, s^{(K)}(n)]$ , the channel gain matrix is  $\boldsymbol{\Gamma} = \text{diag}(\gamma_1, \dots, \gamma_K)$ , and the response matrix is given by

$$\mathbf{A}(\boldsymbol{\omega}) = [\mathbf{a}(\omega_1) \quad \mathbf{a}(\omega_2) \quad \cdots \quad \mathbf{a}(\omega_K)] \quad (4.21)$$



(a) uniform linear array



(b) uniform circular array

Figure 4.5: Calibrated array model geometries for (a) uniform linear and (b) circular arrays.

with the array response vectors  $\mathbf{a}(\omega_k)$  depending on the physical design of the array and the AOA. Calibration posits this known array geometry. The examples considered in this section consist of the uniform linear array (ULA) and uniform circular array (UCA), as shown in Figure 4.5. For example, for the ULA the response vectors typically have a Vandermonde vector structure with base  $e^{-i\pi \sin(\omega_k)}$ .

## 4.2.1 The Fisher information of the calibrated array model

### 4.2.1.1 Indirect derivation of the FIM

As was mentioned in section 4.1, the calibrated array model can be a special case of the convolutive mixture model using constraints. The method for transform-

ing the model is to keep the parameters for the instantaneous mixing model and to the parameter vector add extra parameters for the calibrated array and channel gain. When evaluating the FIM, the elements in the rows and columns corresponding to these extra parameters are zero. The constraint then represents the model reparameterization, e.g., for this case, we choose the constraints which define  $\mathbf{H} = \mathbf{A}(\boldsymbol{\omega})\Gamma$ , element by element. The resulting CCRB submatrix corresponding to the extra parameters will be equivalent to the CRB of those parameters. This procedure is made clear in the following example.

**Example 4.21.** Assume  $x \sim \mathcal{N}(ab, 1)$ . The FIM for  $\boldsymbol{\theta}^T = [a, b]$  is  $\mathbf{I}(\boldsymbol{\theta}) = \begin{bmatrix} b^2 & ab \\ ab & a^2 \end{bmatrix}$ , which is singular. Suppose we wish to reparameterize the model replacing  $ab$  with  $c$ . This is equivalent to the constraint  $f(\boldsymbol{\theta}^*) = ab - c$  for the expanded parameter vector  $\boldsymbol{\theta}^{*T} = [a, b, c]$ . For this orthonormal complement

$$\mathbf{U}(\boldsymbol{\theta}^*) = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ b & a \end{bmatrix}$$

of the Jacobian of the constraints, then  $\mathbf{U}^T(\boldsymbol{\theta}^*)\mathbf{I}(\boldsymbol{\theta}^*)\mathbf{U}(\boldsymbol{\theta}^*) = \mathbf{I}(\boldsymbol{\theta})$ , which is still singular. Hence  $\text{CCRB}(\boldsymbol{\theta}^*) = \mathbf{U}(\boldsymbol{\theta}^*)\mathbf{I}^\dagger(\boldsymbol{\theta})\mathbf{U}^T(\boldsymbol{\theta}^*)$  where the pseudoinverse is  $\mathbf{I}^\dagger(\boldsymbol{\theta}) = \frac{1}{(a^2+b^2)^2}\mathbf{I}(\boldsymbol{\theta})$ , or

$$\text{CCRB}(\boldsymbol{\theta}^*) = \begin{bmatrix} b^2 & ab & b^3 + a^2b \\ ab & a^2 & ab^2 + a^3 \\ b^3 + a^2b & ab^2 + a^3 & b^4 + 2a^2b^2 + a^4 \end{bmatrix}.$$

The component corresponding to  $\text{CCRB}(c) = \frac{1}{(a^2+b^2)^2}(b^4 + 2a^2b^2 + a^4) = 1$ , which agrees with the value of  $\text{CRB}(c)$  for the model  $x \sim \mathcal{N}(c, 1)$ .

Note that the original model need not be identifiable, nor does the replacement model. Additional constraints under either model can also be included in the

constraint function. The example verifies that even for the more difficult case it is possible to find the CRB and CCRB for the calibrated array model indirectly from the instantaneous mixing model using constraints. However, in the interests of clarity, in the next section the FIM and CCRB will be derived for the calibrated array model directly.

#### 4.2.1.2 Direct derivation of the FIM

The calibrated model in (4.19) has a likelihood given by

$$p(\mathbf{y}(1), \dots, \mathbf{y}(N); \boldsymbol{\theta}) = \frac{1}{(\pi\sigma^2)^{MN}} \exp \left\{ -\frac{1}{\sigma^2} \sum_{n=1}^N (\mathbf{y}(n) - \mathbf{A}(\boldsymbol{\theta})\boldsymbol{\Gamma}\mathbf{s}(n))^H (\mathbf{y}(n) - \mathbf{A}(\boldsymbol{\theta})\boldsymbol{\Gamma}\mathbf{s}(n)) \right\}. \quad (4.22)$$

For clarity, since there exists a mixture of complex and real parameters requiring estimation, then using the parameter vector<sup>7</sup>

$$\boldsymbol{\theta} = \begin{bmatrix} \text{Re}(\mathbf{s}(1)) \\ \text{Im}(\mathbf{s}(1)) \\ \vdots \\ \text{Re}(\mathbf{s}(N)) \\ \text{Im}(\mathbf{s}(N)) \\ \boldsymbol{\gamma} \\ \boldsymbol{\omega} \end{bmatrix}, \quad (4.23)$$

with  $\boldsymbol{\gamma}^T = [\gamma_1, \dots, \gamma_K]$  and  $\boldsymbol{\omega}^T = [\omega_1, \dots, \omega_K]$ , the Fisher information matrix is given by

$$\mathbf{I}(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{M} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{M}_1 & \mathbf{L}_1 \\ \mathbf{0} & \mathbf{M} & & \mathbf{0} & \mathbf{M}_2 & \mathbf{L}_2 \\ \vdots & & \ddots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{M} & \mathbf{M}_N & \mathbf{L}_N \\ \mathbf{M}_1^T & \mathbf{M}_2^T & \cdots & \mathbf{M}_N^T & \mathbf{K}_\gamma & \mathbf{L} \\ \mathbf{L}_1^T & \mathbf{L}_2^T & \cdots & \mathbf{L}_N^T & \mathbf{L}^T & \mathbf{K}_\omega \end{bmatrix} \quad (4.24)$$

---

<sup>7</sup>It is also possible to include the noise variance parameter  $\sigma^2$  in the parameter vector. However, this parameter decouples from the other parameters resulting in an optimistic CRB of  $\frac{\sigma^4}{MN}$  [16, 59], and so is uninteresting to the results herein.

where

$$\mathbf{M} = \begin{bmatrix} \text{Re}(\mathcal{M}) & -\text{Im}(\mathcal{M}) \\ \text{Im}(\mathcal{M}) & \text{Re}(\mathcal{M}) \end{bmatrix}, \quad \mathcal{M} = \frac{2}{\sigma^2} \mathbf{\Gamma}^H \mathbf{A}^H(\omega) \mathbf{A}(\omega) \mathbf{\Gamma}, \quad (4.25)$$

$$\mathbf{M}_n = \begin{bmatrix} \text{Re}(\mathcal{M}_n) & -\text{Im}(\mathcal{M}_n) \\ \text{Im}(\mathcal{M}_n) & \text{Re}(\mathcal{M}_n) \end{bmatrix}, \quad \mathcal{M}_n = \frac{2}{\sigma^2} \mathbf{\Phi}^H \mathbf{A}^H(\omega) \mathbf{A}(\omega) \mathbf{S}(n), \quad (4.26)$$

$$\mathbf{L}_n = \begin{bmatrix} \text{Re}(\mathcal{L}_n) \\ \text{Im}(\mathcal{L}_n) \end{bmatrix}, \quad \mathcal{L}_n = \frac{2}{\sigma^2} \mathbf{\Gamma}^H \mathbf{A}^H(\omega) \mathbf{D}(\omega) \mathbf{\Gamma} \mathbf{S}(n), \quad (4.27)$$

$$\mathbf{L} = \begin{bmatrix} \text{Re}(\mathcal{L}) \\ \text{Im}(\mathcal{L}) \end{bmatrix}, \quad \mathcal{L} = \frac{2}{\sigma^2} \sum_{n=1}^N \mathbf{S}^H(n) \mathbf{A}^H(\omega) \mathbf{D}(\omega) \mathbf{\Gamma} \mathbf{S}(n), \quad (4.28)$$

$$\mathbf{K}_\gamma = \begin{bmatrix} \text{Re}(\mathcal{K}_\gamma) & -\text{Im}(\mathcal{K}_\gamma) \\ \text{Im}(\mathcal{K}_\gamma) & \text{Re}(\mathcal{K}_\gamma) \end{bmatrix}, \quad \mathcal{K}_\gamma = \frac{2}{\sigma^2} \sum_{n=1}^N \mathbf{S}^H(n) \mathbf{A}^H(\omega) \mathbf{A}(\omega) \mathbf{S}(n), \quad (4.29)$$

$$\text{and } \mathbf{K}_\omega = \frac{2}{\sigma^2} \sum_{n=1}^N \text{Re}(\mathbf{S}^H(n) \mathbf{\Gamma}^H \mathbf{D}^H(\omega) \mathbf{D}(\omega) \mathbf{\Gamma} \mathbf{S}(n)). \quad (4.30)$$

In the equations above, we redefine  $\mathbf{S}(n) = \text{diag}(\mathbf{s}^{(1)}(n), \dots, \mathbf{s}^{(K)}(n))$  and

$$\mathbf{D}(\omega) = \begin{bmatrix} \frac{\partial \mathbf{a}(\omega_1)}{\partial \omega_1} & \frac{\partial \mathbf{a}(\omega_2)}{\partial \omega_2} & \dots & \frac{\partial \mathbf{a}(\omega_K)}{\partial \omega_K} \end{bmatrix}.$$

#### 4.2.1.3 Properties of the FIM

The model in (4.19) admits an ambiguity as  $\gamma_k s^{(k)}(n)$  is indistinguishable from  $(c_k \gamma_k) \left( \frac{s^{(k)}(n)}{c_k} \right)$  for any nonzero  $c_k \in \mathbb{C}$ . From Section 2.2, then it should be expected that the FIM in (4.24) is singular. This is indeed the case, e.g., note that

$$\begin{bmatrix} \mathcal{M} & \mathbf{0} & \dots & \mathbf{0} & \mathcal{M}_1 \\ \mathbf{0} & \mathcal{M} & & \mathbf{0} & \mathcal{M}_2 \\ \vdots & & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathcal{M} & \mathcal{M}_N \\ \mathcal{M}_1^T & \mathcal{M}_2^T & \dots & \mathcal{M}_N^T & \mathcal{K}_\gamma \\ \mathcal{L}_1^T & \mathcal{L}_2^T & \dots & \mathcal{L}_N^T & \mathcal{L}^T \end{bmatrix} \cdot \begin{bmatrix} \mathbf{S}(1) \\ \mathbf{S}(2) \\ \vdots \\ \mathbf{S}(N) \\ -\mathbf{\Gamma} \end{bmatrix} = \mathbf{0}.$$

And since the FIM consists of real and imaginary parts of this matrix, it too has a null space, namely, the columns of

$$\mathbf{N}(\boldsymbol{\theta}) = \begin{bmatrix} \text{Re}(\mathbf{S}(1)) & -\text{Im}(\mathbf{S}(1)) \\ \text{Im}(\mathbf{S}(1)) & \text{Re}(\mathbf{S}(1)) \\ \vdots & \vdots \\ \text{Re}(\mathbf{S}(N)) & -\text{Im}(\mathbf{S}(N)) \\ \text{Im}(\mathbf{S}(N)) & \text{Re}(\mathbf{S}(N)) \\ -\text{Re}(\boldsymbol{\Gamma}) & \text{Im}(\boldsymbol{\Gamma}) \\ -\text{Im}(\boldsymbol{\Gamma}) & -\text{Re}(\boldsymbol{\Gamma}) \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$$

are a basis for the null space (or at least a null subspace) of  $\mathbf{I}(\boldsymbol{\theta})$ . So the nullity, or dimension of the null space, of  $\mathbf{I}(\boldsymbol{\theta})$  is at least  $2K$ . In fact, it is exactly  $2K$  (provided  $N \geq \frac{K}{M-K}$  for reasons similar to that given in theorem 4.12).

## 4.2.2 Constraints for the calibrated array model

### 4.2.2.1 Constraints on the complex-valued channel gain: $\boldsymbol{\Gamma} = \mathbf{I}_{K \times K}$

One approach to eliminating the ambiguity between the complex-valued gain and the source input is to incorporate this gain into the signal. Instead of remodeling the mean of (4.19) to be

$$\boldsymbol{\mu}(n, \boldsymbol{\theta}) = \sum_{k=1}^K \mathbf{a}(\omega_k) s^{(k)}(n) \quad (4.31)$$

by eliminating the unknown gain, it is equivalent to impose the constraints  $\gamma_k = 1$  for  $k = 1, \dots, K$ . The model in (4.31) is a model presented in a paper on “direction finding with narrow-band sensor arrays” by Stoica and Nehorai [67]. Hence, if the theory of Chapter 3 is to be trusted, then results found by imposing proper constraints that  $\boldsymbol{\Gamma} = \mathbf{I}_{K \times K}$  (or  $\boldsymbol{\gamma} = \mathbf{1}_K$ ) should be equivalent to the results of Stoica

and Nehorai. The  $K$  constraints of a complex-valued parameter  $\boldsymbol{\gamma}$  are  $2K$  constraints of the corresponding real-valued parameters, i.e., the vector  $\boldsymbol{f}$  of constraints can be defined as

$$\begin{aligned} f_k(\boldsymbol{\theta}) &= \operatorname{Re}(\gamma_k) - 1 = 0 \\ f_{K+k}(\boldsymbol{\theta}) &= \operatorname{Im}(\gamma_k) = 0 \end{aligned}$$

for  $k = 1, \dots, K$ . For the parameter vector in (4.23), the Jacobian of these constraints is given by

$$\boldsymbol{F}(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{0}_{2K \times 2KN} & \boldsymbol{I}_{2K \times 2K} & \mathbf{0}_{2K \times K} \end{bmatrix}.$$

Since  $\boldsymbol{F}(\boldsymbol{\theta})\boldsymbol{N}(\boldsymbol{\theta}) = \boldsymbol{I}_{2K \times 2K}$ , then by theorem 3.23, this constraint is sufficient to (locally) identify the parameters and by theorem 3.24 the matrix  $\boldsymbol{U}^T(\boldsymbol{\theta})\boldsymbol{I}(\boldsymbol{\theta})\boldsymbol{U}(\boldsymbol{\theta})$  will be regular for any matrix  $\boldsymbol{U}(\boldsymbol{\theta})$  satisfying (3.6). An orthonormal basis for the null space of the Jacobian would be the columns of

$$\boldsymbol{U}(\boldsymbol{\theta}) = \begin{bmatrix} \boldsymbol{I}_{2KN \times 2KN} & \mathbf{0}_{2KN \times K} \\ \mathbf{0}_{2K \times 2KN} & \mathbf{0}_{2K \times K} \\ \mathbf{0}_{K \times 2KN} & \boldsymbol{I}_{K \times K} \end{bmatrix},$$

and this generates a reduced ‘‘FIM’’

$$\boldsymbol{U}^T(\boldsymbol{\theta})\boldsymbol{I}(\boldsymbol{\theta})\boldsymbol{U}(\boldsymbol{\theta}) = \begin{bmatrix} \boldsymbol{M} & \mathbf{0} & \cdots & \mathbf{0} & \boldsymbol{L}_1 \\ \mathbf{0} & \boldsymbol{M} & & \mathbf{0} & \boldsymbol{L}_2 \\ \vdots & & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \boldsymbol{M} & \boldsymbol{L}_N \\ \boldsymbol{L}_1^T & \boldsymbol{L}_2^T & \cdots & \boldsymbol{L}_N^T & \boldsymbol{K}_\omega \end{bmatrix},$$

which is equivalent to the Fisher information in [67, equation (E.9)]. Since the  $\boldsymbol{\gamma}$  parameters are known and therefore it is unnecessary to understand the performance potential (or CCRB) of a known parameter, it is only of interest to have a bound

on the performance of estimators of the transformation

$$\boldsymbol{\alpha} = \mathbf{k}(\boldsymbol{\theta}) = \begin{bmatrix} \text{Re}(\mathbf{s}(1)) \\ \text{Im}(\mathbf{s}(1)) \\ \vdots \\ \text{Re}(\mathbf{s}(N)) \\ \text{Im}(\mathbf{s}(N)) \\ \boldsymbol{\omega} \end{bmatrix},$$

which has the Jacobian  $\frac{\partial}{\partial \boldsymbol{\theta}^T} \boldsymbol{\alpha} = \mathbf{U}^T(\boldsymbol{\theta})$ . Hence,  $\text{CCRB}(\boldsymbol{\alpha}) = (\mathbf{U}^T(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}))^{-1}$  is the same as the CRB found by Stoica and Nehorai. This equivalence serves as further validation of the CCRB approach.

#### 4.2.2.2 Semiblind constraints: $\mathbf{s}(t) = \mathbf{p}(t)$ for $t \in \mathbb{T}$

An alternative approach to eliminating the ambiguity between the source and coefficient is to have prior knowledge of some of the source signals. These known elements are often referred to as *training* or *pilot symbols* in communications. Knowledge of any  $k$ th source element  $s^{(k)}(t) = p^{(k)}(t)$  at any time sample  $t$  resolves the ambiguity between  $\gamma_k s^{(k)}(n)$  for all  $n$  since  $\gamma_k$  is solvable in the observation corresponding to time sample  $t$  and can thus be used to solve the unknown source values when  $n \notin \mathbb{T}$ . This model can be written as

$$\boldsymbol{\mu}(n, \boldsymbol{\theta}) = \sum_{k=1}^K \mathbf{a}(\omega_k) (\gamma_k p^{(k)}(n) \delta_{n \in \mathbb{T}} + \gamma_k s^{(k)}(n) \delta_{n \notin \mathbb{T}}) \quad (4.32)$$

where  $\delta_{\text{statement}} = 1$  when the statement is true and  $= 0$  when the statement is false. This model is equivalent to the model designed by Kozick and Sadler [39, equation (9)] except with  $\gamma_k s^{(k)}(n)$  being simply  $s^{(k)}(n)$ , a distinction that still allows for a match of results of the CCRB of a properly chosen transformation. The model is also equivalent to the model designed by Li and Compton [41, equation (2)] when



$\mathbb{T} = \{1, 2, \dots, N\}$ . Equivalence of the reparameterized CRBs in [39, 41] with the CCRB can be found in [59]. As in the previous example, the CCRB on the unknown parameters is the inverse of the (unconstrained) FIM after the elimination of its rows and columns corresponding to the known or specified parameters.

#### 4.2.2.3 Finite alphabet constraint: $s^{(k)}(n) \in \mathbb{S}$

This constraint derives from the assumption that the source elements exists in a (finite) discrete set  $\mathbb{S}$ . In communications models, this corresponds to digital modulation designs such as pulse amplitude modulation (PAM), quadrature amplitude modulation (QAM), phase-shift keying (PSK), etc. As such, the model can also be the same as any of the previous models in (4.19), (4.31), or (4.32), but the defining characteristic is that the source samples only exist on a discrete set. There do not exist any CRB-type bounds for this model in the literature due to the problem of differentiability with respect to the parameters.<sup>8</sup> It is certainly possible to constrain any single real-valued parameter to a discrete set by creating a polynomial (or sine function) whose zeros match the set values. But what information can be gained from a constraint formulation of this discrete-alphabet model? This is perhaps best answered with the following example.

**Example 4.22.** Reconsider the model in example 2.2,  $y \sim \mathcal{CN}(\vartheta, \sigma^2)$ . Suppose the

---

<sup>8</sup>A Chapman-Robbins or Barankin-type bound, which does not require differentiability with respect to the parameters would be possible but even this approach does not seem to exist in the literature. Many communications engineers also discount the importance of a mean-square error bound for a digital signal and rely an alternative performance criteria, such as bit-error rate (BER).

parameter is required to satisfy the constraints

$$\begin{aligned} f_1(\vartheta) &= (\operatorname{Re}(\vartheta))^2 - \frac{1}{2} = 0 \\ f_2(\vartheta) &= (\operatorname{Im}(\vartheta))^2 - \frac{1}{2} = 0. \end{aligned}$$

This constraint restricts the real and imaginary part of  $\vartheta$  to reside in the discrete set  $\{\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2}\}$ , i.e.

$$\vartheta \in \mathbb{S} = \left\{ \frac{\sqrt{2}}{2} + i\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2} - i\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2} + i\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2} - i\frac{\sqrt{2}}{2} \right\}.$$

(In the communications vernacular, this is quadrature PSK or 4-QAM.) The Jacobian for this constraint is

$$\mathbf{F}(\boldsymbol{\theta}) = \begin{bmatrix} 2\operatorname{Re}(\vartheta) & 0 \\ 0 & 2\operatorname{Im}(\vartheta) \end{bmatrix},$$

which is full column rank for the possible set of values for  $\vartheta$ . Hence, to satisfy (3.6),  $\mathbf{U}(\boldsymbol{\theta}) = [\ ]$  is a  $2 \times 0$  null matrix.

This is analogous to the determinate case of example 3.14, therefore knowledge that a parameter exists in a discrete set is equivalent to complete knowledge of the parameter value in terms of mean-square error performance potential-the result being a Cramér-Rao bound of zero. This does not mean that the mean-square error will be zero (the decision of which set value the parameter actually is can be wrong given sufficient power in the noise), but it does mean that the mean-square error bound is trivial and not particularly helpful. This result is not surprising considering the degrees of freedom of the parameter that are restricted from such constraints. The restriction of a real-valued parameter to satisfying a single root equation eliminates its single degree of freedom.

#### 4.2.2.4 Unit modulus constraints: $|\mathbf{s}(n)| = 1$ for all $n$

Since knowledge of parameters from a discrete (finite) alphabet results in a trivial bound, then to obtain a relevant and useful measure of performance potential a relaxation of the side information is necessary. One such example of this approach is using a constant or unit modulus constraint on the source elements as in section 4.1.4.4 for the convolutive mixture model. This approach remodels the mean as

$$\boldsymbol{\mu}(n, \boldsymbol{\theta}) = \sum_{k=1}^K \mathbf{a}(\omega_k) \gamma_k e^{j\phi^{(k)}(n)}. \quad (4.33)$$

This is the constraint considered in example 3.4 applied to the communications context. Therefore, imposing the constraints

$$f_{(n-1)K+k}(\boldsymbol{\theta}) = |s^{(k)}(n)|^2 - 1 = 0,$$

for  $k = 1, \dots, K, n = 1, \dots, N$ , is an alternative approach than rederiving the Fisher information for the model in (4.33). The Jacobian for this constraint has the form

$$\mathbf{F}(\boldsymbol{\theta}) = \begin{bmatrix} 2\text{Re}(\mathbf{S}(1)) & 2\text{Im}(\mathbf{S}(1)) & & & \mathbf{0}_{K \times 3K} \\ & & \ddots & & \vdots \\ & & & 2\text{Re}(\mathbf{S}(N)) & 2\text{Im}(\mathbf{S}(N)) & \mathbf{0}_{K \times 3K} \end{bmatrix}, \quad (4.34)$$

which has a null space generated by the columns of the matrix

$$\mathbf{U}(\boldsymbol{\theta}) = \begin{bmatrix} -\text{Im}(\mathbf{S}(1)) \\ \text{Re}(\mathbf{S}(1)) \\ & \ddots & \\ & & -\text{Im}(\mathbf{S}(N)) \\ & & & \text{Re}(\mathbf{S}(N)) \\ & & & & \mathbf{I}_{3K \times 3K} \end{bmatrix}. \quad (4.35)$$

From this we can check the (local) identifiability of this model under the (unit) constant modulus constraint using theorem 3.24. The matrix  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  is

$$\begin{bmatrix} \text{Re}(\mathbf{S}^*(1)\mathbf{M}\mathbf{S}(1)) & & \text{Im}(\mathbf{S}^*(1)\mathbf{M}_1) & \text{Re}(\mathbf{S}^*(1)\mathbf{M}_1) & \text{Im}(\mathbf{S}^*(1)\mathbf{L}_1) \\ & \ddots & \vdots & \vdots & \vdots \\ & & \text{Re}(\mathbf{S}^*(N)\mathbf{M}\mathbf{S}(N)) & \text{Im}(\mathbf{S}^*(N)\mathbf{M}_N) & \text{Re}(\mathbf{S}^*(N)\mathbf{L}_N) \\ -\text{Im}(\mathbf{M}_1^H \mathbf{S}(1)) & \cdots & -\text{Im}(\mathbf{M}_N^H \mathbf{S}(N)) & \text{Re}(\mathbf{K}_\gamma) & -\text{Im}(\mathbf{K}_\gamma) & \text{Re}(\mathbf{L}) \\ \text{Re}(\mathbf{M}_1^H \mathbf{S}(1)) & \cdots & \text{Re}(\mathbf{M}_N^H \mathbf{S}(N)) & \text{Im}(\mathbf{K}_\gamma) & \text{Re}(\mathbf{K}_\gamma) & \text{Im}(\mathbf{L}) \\ -\text{Im}(\mathbf{L}_1^H \mathbf{S}(1)) & \cdots & -\text{Im}(\mathbf{L}_N^H \mathbf{S}(N)) & \text{Re}(\mathbf{L}^H) & -\text{Im}(\mathbf{L}^H) & \mathbf{K}_\omega \end{bmatrix}$$

Since

$$\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) \cdot \begin{bmatrix} \mathbf{I}_{K \times K} \\ \vdots \\ \mathbf{I}_{K \times K} \\ -\text{Im}(\boldsymbol{\Gamma}) \\ \text{Re}(\boldsymbol{\Gamma}) \\ \mathbf{0}_{K \times K} \end{bmatrix} = \mathbf{0},$$

then  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  is singular and by theorem 3.24 the (unit) constant modulus constraints are not sufficient for identifiability. This was also the case for the convolutive mixture model in section 4.1.4.4. Furthermore, reviewing the model in (4.33), this result has more reason to be expected. The original identifiability issue in the calibrated model (4.19) is the multiplicative ambiguity between the sources and the channel gain. While the constant modulus constraint resolves any amplitude ambiguity of the channel gain, i.e.,  $|\gamma_k| = |\gamma_k s^{(k)}(n)|$  for any  $n$ , there still exist a phase rotation ambiguity. It is clear from both the model and from theorem 3.23, that for each source  $k$  knowledge of an element for either the real or imaginary part of either the channel gain or a source sample will be sufficient for identifiability and a regular  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$ . Of course, given this constant unit modulus constraint, knowledge of the real (imaginary) part of any source sample is equivalent to restricting the imaginary (real) part of the source sample to a finite discrete alphabet, which as discussed in section 4.2.2.3 is the same as knowledge of the parameter in regards to the CCRB performance potential but not in regards to the estimation. Hence for estimation performance it is a necessity to constrain the real and/or imaginary part of the source sample.

#### 4.2.2.5 Unit modulus constraint; real-valued channel gain: $\text{Im}(\gamma_k) = 0$ for all $k$

This model is that given in (4.33) except with  $\gamma_k$  being real-valued. As seen in sections 4.2.2.2 and 4.2.2.3, the constraint is knowledge of the imaginary parts of  $\boldsymbol{\gamma}$  and the effect of the corresponding  $\mathbf{U}(\boldsymbol{\theta})$  matrix is the elimination of the columns and rows corresponding to the  $\text{Im}(\gamma_k)$  parameters. This model is equivalent to that of Leshem and van der Veen [40]. Verification that the CCRB is equivalent can be found in [59]. The bound is essentially the inverse of the *reduced-parameter-space Fisher information*  $\mathbf{U}^T(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta})$  from section 4.2.2.4 after the elimination of the rows and columns corresponding to the channel gain parameters.

#### 4.2.2.6 Semi-blind and unit modulus constraint

This model merges the models in (4.33) and (4.32). Without loss of generality, assume the source elements are known for the first  $T$  time slots for each source. Then the constraint Jacobian is  $\mathbf{F}(\boldsymbol{\theta}) =$

$$\begin{bmatrix} \mathbf{I}_{2TK \times 2TK} & & & & \mathbf{0}_{2TK \times 3K} \\ & 2\text{Re}(\mathbf{S}(T+1)) & 2\text{Im}(\mathbf{S}(T+1)) & & \mathbf{0}_{K \times 3K} \\ & & \ddots & & \vdots \\ & & & 2\text{Re}(\mathbf{S}(N)) & 2\text{Im}(\mathbf{S}(N)) & \mathbf{0}_{K \times 3K} \end{bmatrix}, \quad (4.36)$$

which has an orthonormal complement

$$\mathbf{U}(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{0}_{2TK \times K} \\ -\text{Im}(\mathbf{S}(1)) \\ \text{Re}(\mathbf{S}(1)) \\ \vdots \\ -\text{Im}(\mathbf{S}(N)) \\ \text{Re}(\mathbf{S}(N)) \\ \mathbf{I}_{3K \times 3K} \end{bmatrix}. \quad (4.37)$$

The reduced FIM is then

$$U(\boldsymbol{\theta})\mathbf{I}(\boldsymbol{\theta})U(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{G}(T+1) & & & & \mathbf{C}(T+1) & \mathbf{B}(T+1) \\ & \mathbf{G}(T+2) & & & \mathbf{C}(T+2) & \mathbf{B}(T+2) \\ & & \ddots & \vdots & \vdots & \\ & & & \mathbf{G}(N) & \mathbf{C}(N) & \mathbf{B}(N) \\ \mathbf{C}^T(T+1) & \mathbf{C}^T(T+2) & \cdots & \mathbf{C}^T(N) & \mathbf{K}_\gamma & \mathbf{L} \\ \mathbf{B}^T(T+1) & \mathbf{B}^T(T+2) & \cdots & \mathbf{B}^T(N) & \mathbf{L}^T & \mathbf{K}_\omega \end{bmatrix} \quad (4.38)$$

where  $\mathbf{G}(t) = \text{Re}[\mathbf{S}^*(t)\mathbf{M}\mathbf{S}(t)]$  and  $\mathbf{M}$  is defined as in (4.25), where  $\mathbf{C}(t) = [\text{Im}[\mathbf{S}^*(t)\mathbf{M}_t] \quad \text{Re}[\mathbf{S}^*(t)\mathbf{M}_t]]$  and  $\mathbf{M}_n$  is defined in (4.26),  $\mathbf{B}(t) = \text{Im}[\mathbf{S}^*(t)\mathbf{L}_t]$  and  $\mathbf{L}_t$  is defined in (4.27). To analytically invert this matrix, the Schur complement,

$$\Phi = \begin{bmatrix} \mathbf{K}_\gamma - \sum_{t=T+1}^N \mathbf{C}^T(t)\mathbf{G}^{-1}(t)\mathbf{C}(t) & \mathbf{L} - \sum_{t=T+1}^N \mathbf{C}^T(t)\mathbf{G}^{-1}(t)\mathbf{B}(t) \\ \mathbf{L}^T - \sum_{t=T+1}^N \mathbf{B}^T(t)\mathbf{G}^{-1}(t)\mathbf{C}(t) & \mathbf{K}_\omega - \sum_{t=T+1}^N \mathbf{B}^T(t)\mathbf{G}^{-1}(t)\mathbf{B}(t) \end{bmatrix}, \quad (4.39)$$

is useful. If  $\Phi$  is partitioned into corresponding subblocks  $\begin{bmatrix} \Phi_{\mathbf{K}_\gamma} & \Phi_{\mathbf{L}} \\ \Phi_{\mathbf{L}^T} & \Phi_{\mathbf{K}_\omega} \end{bmatrix}$ , then the CCRB subblocks for the unknown elements are given by

$$\begin{aligned} \text{CCRB}(\boldsymbol{\omega}) &= \left[ \Phi_{\mathbf{K}_\omega} - \Phi_{\mathbf{L}^T} \Phi_{\mathbf{K}_\gamma}^{-1} \Phi_{\mathbf{L}} \right]^{-1} \\ \text{CCRB}(\boldsymbol{\gamma}) &= \left[ \Phi_{\mathbf{K}_\gamma} - \Phi_{\mathbf{L}} \Phi_{\mathbf{K}_\omega}^{-1} \Phi_{\mathbf{L}^T} \right]^{-1} \\ \text{CCRB}\left(\begin{bmatrix} \text{Re}(\mathbf{s}(t)) \\ \text{Im}(\mathbf{s}(t)) \end{bmatrix}\right) &= \begin{bmatrix} -\text{Im}(\mathbf{S}(t)) \\ \text{Re}(\mathbf{S}(t)) \end{bmatrix} \mathbf{X}(t) \begin{bmatrix} -\text{Im}(\mathbf{S}(t)) & \text{Re}(\mathbf{S}(t)) \end{bmatrix} \end{aligned}$$

where  $\mathbf{X}(t) = \mathbf{G}^{-1}(t) + \mathbf{G}^{-1}(t) [\mathbf{C}(t) \quad \mathbf{B}(t)] \Phi^{-1} \begin{bmatrix} \mathbf{C}^T(t) \\ \mathbf{B}^T(t) \end{bmatrix} \mathbf{G}^{-1}(t)$ .

**Example 4.23.** A distinct advantage of the CCRB, as implemented in this treatise, is the ability to compare the performance potential of a number of different models seamlessly. Suppose we consider  $M = 5$  omni-directional sensors with a beamwidth

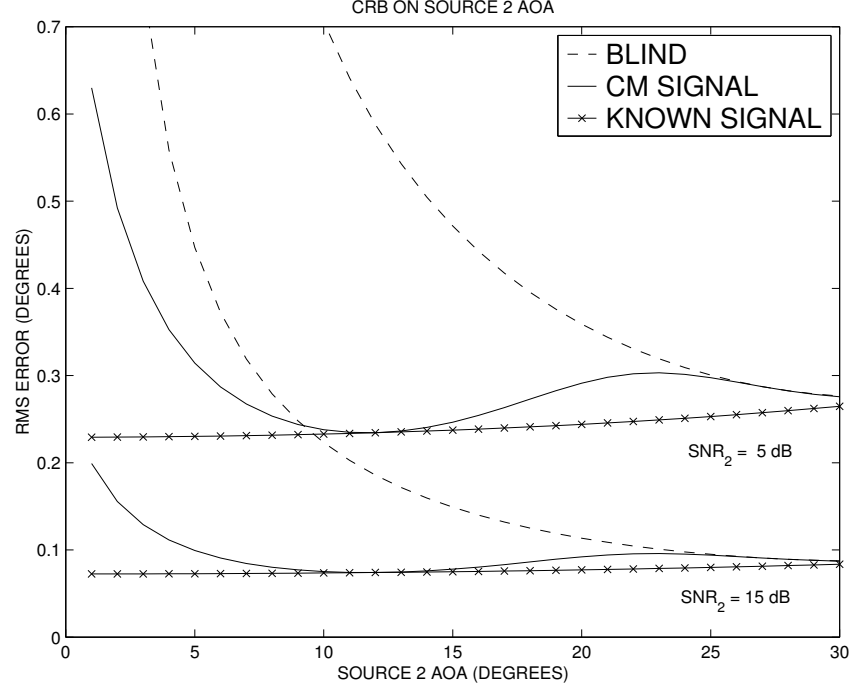


Figure 4.6: CCRBs on AOA for blind, constant modulus, and known signal models.

of  $\approx 23.6^\circ$  in a uniform linear array receiving  $K = 2$  source signals over  $N = 100$  time samples. Then, it is possible to compare various communications design scenarios, e.g.,

- (a) the “blind” case<sup>9</sup>:  $s_k(1)$  known for  $k = 1, 2$ ,
- (b) the unit modulus case:  $|s_k(t)|^2 = 1$  for  $k = 1, 2$  and  $t = 1, \dots, 100$ ,
- (c) the semiblind case:  $s_k(t)$  known for  $k = 1, 2$  and  $t = 1, \dots, T = 20$ ,
- (d) the unit modulus and semiblind case:  $s_k(t)$  known for  $k = 1, 2$  and  $t = 1, \dots, T = 20$ , and  $|s_k(t)|^2 = 1$  for  $k = 1, 2$  and  $t = T + 1, \dots, 100$ , and
- (e) the known signal case:  $s_k(t)$  known for  $k = 1, 2$  and  $t = 1, \dots, 100$ .

<sup>9</sup>This is not a truly blind scenario, but is often referred to as blind in the literature.

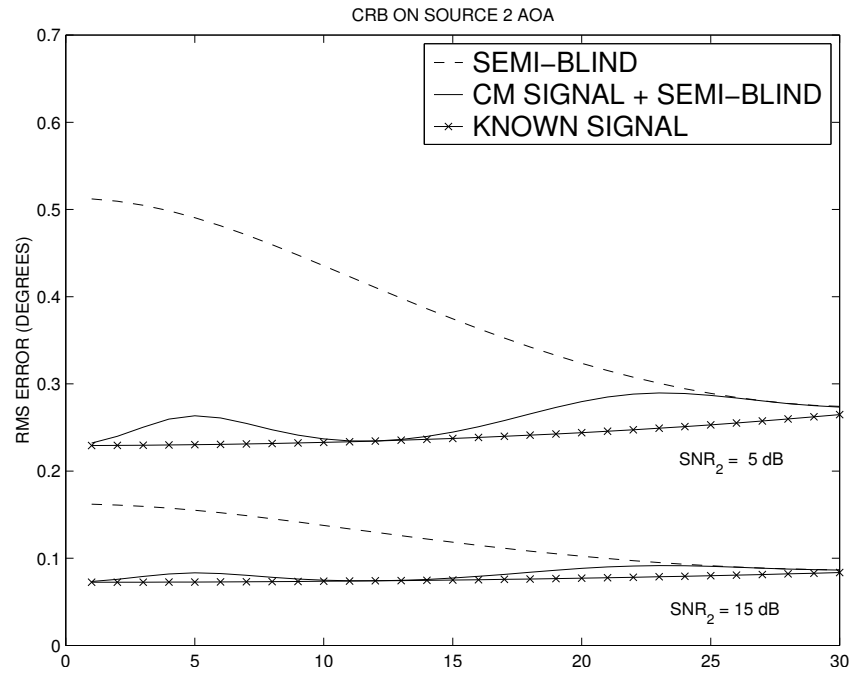


Figure 4.7: CCRBs on AOA for semiblind, constant modulus + semiblind, and known signal models.



Figure 4.6 displays a comparison between blind, unit modulus constraints, and known signal model designs of the CCRBs on angle-of-arrival (AOA) estimation for the second source signal over varying directions and different signal-to-noise ratios (SNRs) with the first source signal arriving at  $0^\circ$ . Figure 4.7 displays a comparison of CCRBs on AOA estimation between semiblind constraints, unit modulus with semiblind constraints, and known signal model designs. And finally, figure 4.8 displays CCRBs on signal phase estimation between blind, unit modulus constraints, semiblind constraints, and a mixture of unit modulus and semiblind constraints. The known signal model is the best case scenario for AOA estimation potential and is a useful guideline for more desirable scenarios where information (data or unknown parameters) is included in the transmission. Figures 4.6 and 4.7 demonstrate the characteristic loss of performance when the sources' AOAs differ by roughly the beamwidth. The value of the unit modulus constraint is evident when the sources' AOAs are closely spaced as the CCRB performance potential approximates that of the known signal model. In figures 4.7 and 4.8, the estimation potential actually improves for closely space sources with the semiblind and unit modulus constraint mixture.

### 4.3 Discussion

This chapter includes extensions of only a brief sampling of my prior research [59, 39, 48, 51, 39, 48, 59, 49] as it relates to the practical application of the CCRB. In this chapter, the convolutive mixture model and the calibrated array model were

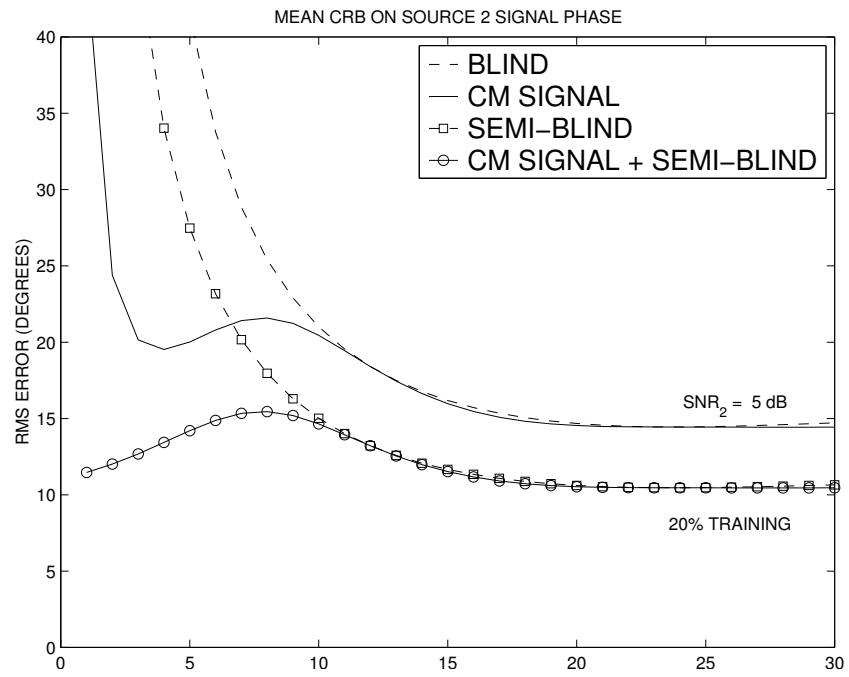


Figure 4.8: CCRBs on signal phase for blind, constant modulus, semibind, and constant modulus + semibind models.

treated as base models for which the Fisher information was derived a single time and then a series of variations on the models in the form of differentiable parametric constraints were considered. This approach presents a simple procedure to compare and contrast a large class of constraints, essentially different models, in an efficient manner to determine the value of particular formulation in terms of performance potential as measured in the CCRB.

## Appendix A

### Appendices

#### A.1 A proof of the CCRB using the Chapman-Robbins version of the Barankin bound

Gorman and Hero developed a CCRB using the multiparameter version of the Hammersley-Chapman-Robbins bound (HCRB) [16, 26]. However, the result produced a variant form of the CCRB

$$\mathbf{I}^{-1}(\boldsymbol{\theta}) - \mathbf{I}^{-1}(\boldsymbol{\theta}) \mathbf{F}^T(\boldsymbol{\theta}) (\mathbf{F}(\boldsymbol{\theta}) \mathbf{I}^{-1}(\boldsymbol{\theta}) \mathbf{F}^T(\boldsymbol{\theta}))^{-1} \mathbf{F}(\boldsymbol{\theta}) \mathbf{I}^{-1}(\boldsymbol{\theta}),$$

which requires a nonsingular FIM. What follows is a shorter variation of their approach that does not assume a nonsingular FIM, starting with a brief description of the HCRB.

Rather than relying on the Fisher score, which is the derivative of the log-likelihood, i.e.,

$$\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \log p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{p(\mathbf{x}; \boldsymbol{\theta})} \frac{\partial p(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

the regularity conditions requiring a differentiable likelihood can be relaxed by considering finite differences, i.e., for each  $i = 1, \dots, m$ ,

$$\frac{1}{p(\mathbf{x}; \boldsymbol{\theta})} \frac{p(\mathbf{x}; \boldsymbol{\theta} + \epsilon_i \mathbf{e}_i) - p(\mathbf{x}; \boldsymbol{\theta})}{\epsilon_i}$$

where the  $\mathbf{e}_i$  are canonical unit vectors. If the likelihood is differentiable then the limit as each  $\epsilon_i \rightarrow 0$  is the Fisher score. Of course, the finite differences need not

be with respect to the canonical axis and the number of finite differences need not be the same as the dimension of the parameters. This is the generalization of the CRB known as the Hammersley-Chapman-Robbins (HCR) version of the Barankin bound.

**Theorem A.1.** If  $\mathbf{t}(\mathbf{x})$  be an unbiased estimate of  $\mathbf{h}(\boldsymbol{\theta}) \in \mathbf{h}(\mathbb{R}^m) \subset \mathbb{R}^t$ , and  $\boldsymbol{\psi} = [\boldsymbol{\psi}^{(1)}, \dots, \boldsymbol{\psi}^{(p)}]$  is a matrix whose columns are test points in  $\mathbb{R}^m$ , all distinct from each other as well as from  $\boldsymbol{\theta}$ , then the variance of  $\mathbf{t}(\mathbf{x})$  is bounded below by the inequality

$$\text{Var}(\mathbf{t}(\mathbf{x})) \geq \sup_{\boldsymbol{\psi}^{(1)}, \dots, \boldsymbol{\psi}^{(p)}, p} \boldsymbol{\Delta}(\boldsymbol{\theta}, \boldsymbol{\psi}) \boldsymbol{\Upsilon}^{-1}(\boldsymbol{\theta}, \boldsymbol{\psi}) \boldsymbol{\Delta}^T(\boldsymbol{\theta}, \boldsymbol{\psi})$$

where  $\boldsymbol{\Delta}(\boldsymbol{\theta}, \boldsymbol{\psi})$  is called a *translation matrix* defined by

$$\boldsymbol{\Delta}(\boldsymbol{\theta}, \boldsymbol{\psi}) = [\mathbf{h}(\boldsymbol{\psi}^{(1)}) - \mathbf{h}(\boldsymbol{\theta}), \dots, \mathbf{h}(\boldsymbol{\psi}^{(p)}) - \mathbf{h}(\boldsymbol{\theta})]$$

and  $\boldsymbol{\Upsilon}(\boldsymbol{\theta}, \boldsymbol{\psi})$  is called an *HCR information matrix* defined by

$$\Upsilon_{ij}(\boldsymbol{\theta}, \boldsymbol{\psi}) = E_{\boldsymbol{\theta}} \frac{p(\mathbf{x}; \boldsymbol{\psi}^{(i)}) - p(\mathbf{x}; \boldsymbol{\theta})}{p(\mathbf{x}; \boldsymbol{\theta})} \frac{p(\mathbf{x}; \boldsymbol{\psi}^{(j)}) - p(\mathbf{x}; \boldsymbol{\theta})}{p(\mathbf{x}; \boldsymbol{\theta})}.$$

This result encompasses the CRB result when the test points satisfy  $\boldsymbol{\psi}^{(i)} = \boldsymbol{\theta} + \epsilon_i \mathbf{e}_i$  and the  $\epsilon_i \rightarrow 0$  for  $i = 1, \dots, m = p$ , i.e., as a properly chosen set of test points approach the parameter. If the set of vectors  $\boldsymbol{\psi}^{(1)} - \boldsymbol{\theta}, \dots, \boldsymbol{\psi}^{(m)} - \boldsymbol{\theta}$  span an  $m$ -dimensional space, then the limit as the finite differences approach derivatives will still obtain the CRB.

Since  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$  it makes sense to restrict the test points to also satisfy the constraints,  $\mathbf{f}(\boldsymbol{\psi}^{(i)}) = \mathbf{0}$ , and examine the limit of the HCRB as the finite differences

approach derivatives. The Taylor series approximation of  $\mathbf{f}(\boldsymbol{\psi}^{(i)})$  about  $\boldsymbol{\theta}$  is

$$\mathbf{f}(\boldsymbol{\psi}^{(i)}) = \mathbf{f}(\boldsymbol{\theta})\mathbf{F}(\boldsymbol{\theta}) (\boldsymbol{\psi}^{(i)} - \boldsymbol{\theta}) + o(\|\boldsymbol{\psi}^{(i)} - \boldsymbol{\theta}\|).$$

Since  $\mathbf{f}(\boldsymbol{\psi}^{(i)}) = \mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$ , then  $\boldsymbol{\psi}^{(i)} - \boldsymbol{\theta}$  almost entirely resides in  $\text{null}(\mathbf{F}(\boldsymbol{\theta}))$ . So without loss of generality we can allow

$$\boldsymbol{\psi}^{(i)} = \boldsymbol{\theta} + \delta_i \mathbf{u}_i(\boldsymbol{\theta}) + o(\|\boldsymbol{\psi}^{(i)} - \boldsymbol{\theta}\|)$$

where  $\mathbf{u}_i(\boldsymbol{\theta})$  is the  $i$ th column of the matrix  $\mathbf{U}(\boldsymbol{\theta})$  satisfying (3.6). Then

$$\begin{aligned} \frac{p(\mathbf{x}; \boldsymbol{\psi}^{(i)}) - p(\mathbf{x}; \boldsymbol{\theta})}{\delta_i p(\mathbf{x}; \boldsymbol{\theta})} &= \frac{p(\mathbf{x}; \boldsymbol{\theta} + \delta_i \mathbf{u}_i(\boldsymbol{\theta}) + o(\|\boldsymbol{\psi}^{(i)} - \boldsymbol{\theta}\|)) - p(\mathbf{x}; \boldsymbol{\theta})}{\delta_i p(\mathbf{x}; \boldsymbol{\theta})} \\ &\rightarrow \mathbf{u}_i^T(\boldsymbol{\theta}) \frac{\partial p(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \frac{1}{p(\mathbf{x}; \boldsymbol{\theta})} = \mathbf{u}_i^T(\boldsymbol{\theta}) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}) \end{aligned}$$

and

$$\frac{\boldsymbol{\psi}^{(i)} - \boldsymbol{\theta}}{\delta_i} = \mathbf{u}_i(\boldsymbol{\theta}) + \frac{1}{\delta_i} o(\|\boldsymbol{\psi}^{(i)} - \boldsymbol{\theta}\|) \rightarrow \mathbf{u}_i(\boldsymbol{\theta})$$

as  $\delta_i \rightarrow 0$ . This is true for any  $i$ , so if the test points are chosen such that each column of  $\mathbf{U}(\boldsymbol{\theta})$  is used, this gives us the CCRB

$$\mathbf{U}(\boldsymbol{\theta}) (\mathbf{U}^T(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\theta}) \mathbf{U}(\boldsymbol{\theta}))^{-1} \mathbf{U}^T(\boldsymbol{\theta})$$

as the limit of a constrained HCRB when the finite differences become derivatives.

## A.2 A proof of the CCRB using the method of implicit differentiation

Suppose  $\boldsymbol{\theta}$  is restricted to the zeros of  $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^k$  with Jacobian  $\mathbf{F}(\boldsymbol{\theta}) = \frac{\partial \mathbf{f}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^T}$  having rank  $k$  whenever  $\mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}$ . The method of implicit differentiation assumes that parameters that would be eliminated under a reparameterization can

be written in terms of the remaining parameters, since the conditions satisfy the implicit function theorem. The actual function generally remains unknown, but its derivative is calculated by first taking partial derivatives of constraint function and using linear algebra to solve for the partial of the eliminated parameters in terms of the remaining parameters. This approach was also used by Marzetta [47, proof of theorem 1] to prove the regularity conditions given in (3.20).

The parameter vector may be separated as  $\boldsymbol{\theta} = \begin{bmatrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \end{bmatrix}$  and the constraint  $\mathbf{f}$  may be rewritten as  $\mathbf{f}^* : \mathbb{R}^{m-k} \times \mathbb{R}^k \rightarrow \mathbb{R}^k$  via the mapping  $\mathbf{f}^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \mathbf{f}\left(\begin{bmatrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \end{bmatrix}\right)$ . Then the Jacobian of  $\mathbf{f}$  can be represented as

$$\mathbf{F}(\boldsymbol{\theta}) = [\mathbf{f}_{\boldsymbol{\theta}_1}^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \quad \mathbf{f}_{\boldsymbol{\theta}_2}^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)]$$

where  $\mathbf{f}_{\boldsymbol{\theta}_i}^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \frac{\partial}{\partial \boldsymbol{\theta}_i^T} \mathbf{f}^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$  for each  $i = 1, 2$ .

Without loss of generality, assume  $\boldsymbol{\theta}_2 \in \mathbb{R}^k$  is a function of  $\boldsymbol{\theta}_1 \in \mathbb{R}^{m-k}$ , i.e.,  $\boldsymbol{\theta}_2 = \boldsymbol{\theta}_2(\boldsymbol{\theta}_1)$  is an implicit function. Therefore,  $\mathbf{f}^*$  is implicitly only a parameter of  $\boldsymbol{\theta}_1$  and

$$\frac{\partial}{\partial \boldsymbol{\theta}_1^T} \mathbf{f}^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2(\boldsymbol{\theta}_1)) = \mathbf{f}_{\boldsymbol{\theta}_1}^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2(\boldsymbol{\theta}_1)) + \mathbf{f}_{\boldsymbol{\theta}_2}^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2(\boldsymbol{\theta}_1)) \frac{\partial \boldsymbol{\theta}_2(\boldsymbol{\theta}_1)}{\partial \boldsymbol{\theta}_1^T}.$$

If this derivative is only taken where  $\mathbf{f}^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \mathbf{0}$ , then in matrix form,

$$[\mathbf{f}_{\boldsymbol{\theta}_1}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \quad \mathbf{f}_{\boldsymbol{\theta}_2}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)] \begin{bmatrix} \mathbf{I}_{m-k, m-k} \\ \frac{\partial \boldsymbol{\theta}_2(\boldsymbol{\theta}_1)}{\partial \boldsymbol{\theta}_1} \end{bmatrix} = \mathbf{0}.$$

The first matrix above is  $\mathbf{F}(\boldsymbol{\theta})$ ; the second matrix above consists of  $m - k$  linearly independent columns which exist in the null space of the row vectors of  $\mathbf{F}(\boldsymbol{\theta})$ , i.e., the second matrix is merely some transformation of some matrix  $\mathbf{U}(\boldsymbol{\theta})$  defined as in (3.6).

### A.3 Alternative proof of asymptotic normality

Crowder [18] proved the following theorem.

**Theorem A.2** (Crowder). If

1. there is a consistent solution  $(\dot{\boldsymbol{\theta}}_n, \dot{\boldsymbol{\lambda}}_n)$  of the likelihood equations,
2.  $\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}_0) \stackrel{d}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}(\boldsymbol{\theta}_0))$ ,
3.  $\mathbf{D}^{-1}(\dot{\boldsymbol{\theta}}_n) \left( \mathbf{I}(\boldsymbol{\theta}_0) + \frac{\partial}{\partial \boldsymbol{\theta}^T} \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}_0, \dot{\boldsymbol{\theta}}_n) \right) \xrightarrow{p} \mathbf{0}$ ,
4.  $\mathbf{Q}(\dot{\boldsymbol{\theta}}_n) - \mathbf{Q}(\boldsymbol{\theta}_0) \xrightarrow{p} \mathbf{0}$ , and
5.  $\det \mathbf{Q}(\dot{\boldsymbol{\theta}}_n) \leq K < \infty$
6.  $\det \mathbf{D}^{-1}(\dot{\boldsymbol{\theta}}_n) \mathbf{F}^T(\dot{\boldsymbol{\theta}}_n) (\mathbf{F}(\boldsymbol{\theta}, \dot{\boldsymbol{\theta}}_n) \mathbf{D}^{-1}(\dot{\boldsymbol{\theta}}_n) \mathbf{F}^T(\dot{\boldsymbol{\theta}}_n))^{-1} \leq K < \infty$

where  $\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}_0, \boldsymbol{\theta})$  is a matrix in the form of the Fisher score  $\mathbf{s}(\mathbf{x}; \boldsymbol{\theta})$  but with each row evaluated at possibly different points on the line between  $\boldsymbol{\theta}_0$  and  $\dot{\boldsymbol{\theta}}_n$ , and similarly for  $\mathbf{F}(\boldsymbol{\theta}, \dot{\boldsymbol{\theta}}_n)$ , and  $\mathbf{Q}(\boldsymbol{\theta}) = \mathbf{D}^{-1}(\boldsymbol{\theta}) \mathbf{F}^T(\boldsymbol{\theta}) (\mathbf{F}(\boldsymbol{\theta}_0, \boldsymbol{\theta}) \mathbf{D}^{-1}(\boldsymbol{\theta}) \mathbf{F}^T(\boldsymbol{\theta}))^{-1}$ . Then

$$\sqrt{n} (\dot{\boldsymbol{\theta}} - \boldsymbol{\theta}) \stackrel{d}{\rightarrow} \mathcal{N} \left( \mathbf{0}, \mathbf{D}^{-1}(\boldsymbol{\theta}) - \mathbf{D}^{-1}(\boldsymbol{\theta}) \mathbf{F}^T(\boldsymbol{\theta}) (\mathbf{F}(\boldsymbol{\theta}) \mathbf{D}^{-1}(\boldsymbol{\theta}) \mathbf{F}^T(\boldsymbol{\theta}))^{-1} \mathbf{F}(\boldsymbol{\theta}) \mathbf{D}^{-1}(\boldsymbol{\theta}) \right)$$

where  $\mathbf{D}(\boldsymbol{\theta}) = \mathbf{I}(\boldsymbol{\theta}) + \mathbf{F}^T(\boldsymbol{\theta}) \mathbf{K} \mathbf{F}(\boldsymbol{\theta})$  for an arbitrary positive semi-definite matrix  $\mathbf{K}$ .

Crowder's asymptotic normality theorem shows that variance of the CMLE satisfies

$$\text{Var}(\sqrt{n} \hat{\boldsymbol{\theta}}_n) \xrightarrow{d} \mathbf{D}^{-1}(\boldsymbol{\theta}_0) - \mathbf{D}^{-1}(\boldsymbol{\theta}_0) \mathbf{F}^T(\boldsymbol{\theta}_0) (\mathbf{F}(\boldsymbol{\theta}_0) \mathbf{D}^{-1}(\boldsymbol{\theta}_0) \mathbf{F}^T(\boldsymbol{\theta}_0))^{-1} \mathbf{F}(\boldsymbol{\theta}_0) \mathbf{D}^{-1}(\boldsymbol{\theta}_0)$$



as  $n \rightarrow \infty$ , where  $\mathbf{D}(\boldsymbol{\theta}) = \mathbf{I}(\boldsymbol{\theta}) + \mathbf{F}^T(\boldsymbol{\theta})\mathbf{K}\mathbf{F}(\boldsymbol{\theta})$ . This asymptotic variance has the exact structure as the Marzetta form of the CCRB when  $\mathbf{K} = \mathbf{0}$  and the FIM is full rank. Applying the algebraic identity of Lemma 3.8, then

$$\text{Var}(\sqrt{n}\hat{\boldsymbol{\theta}}_n) \xrightarrow{d} \mathbf{U}(\boldsymbol{\theta}_0) \left( \mathbf{U}^T(\boldsymbol{\theta}_0) \mathbf{D}(\boldsymbol{\theta}_0) \mathbf{U}(\boldsymbol{\theta}_0) \right)^{-1} \mathbf{U}^T(\boldsymbol{\theta}_0).$$

It only remains to note

$$\mathbf{U}^T(\boldsymbol{\theta}_0) \mathbf{D}(\boldsymbol{\theta}_0) \mathbf{U}(\boldsymbol{\theta}_0) = \mathbf{U}^T(\boldsymbol{\theta}_0) \mathbf{I}(\boldsymbol{\theta}_0) \mathbf{U}(\boldsymbol{\theta}_0)$$

since  $\mathbf{F}(\boldsymbol{\theta})\mathbf{U}(\boldsymbol{\theta}) = \mathbf{0}$ .

## Appendix B

### Proofs of Convergence Properties of Constrained Scoring

**Theorem** (Theorem 3.33). If for any iterate  $\dot{\boldsymbol{\theta}}^{(k)} \in \Theta_f$  there does not exist an  $\alpha^{(k)} > 0$  that satisfies (3.37), then  $\dot{\boldsymbol{\theta}}^{(k)}$  is a stationary point.

*Proof.* Let  $\dot{\boldsymbol{\theta}} \in \Theta_f$  and define  $\ddot{\boldsymbol{\theta}}(\alpha) = \boldsymbol{\pi} [\dot{\boldsymbol{\theta}} + \alpha \text{CCRB}(\dot{\boldsymbol{\theta}}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}})]$ . By a property of the natural projection of convex sets,

$$\left\| \ddot{\boldsymbol{\theta}}(\alpha) - \dot{\boldsymbol{\theta}} \right\|_{I(\dot{\boldsymbol{\theta}})} \leq \alpha \left\| \text{CCRB}(\dot{\boldsymbol{\theta}}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}) \right\|_{I(\dot{\boldsymbol{\theta}})}.$$

Hence it is sufficient to show there exists an  $\alpha > 0$  such that

$$\frac{(\log p(\mathbf{x}; \ddot{\boldsymbol{\theta}}(\alpha)) - \log p(\mathbf{x}; \dot{\boldsymbol{\theta}}))}{\alpha} \geq \kappa \mathbf{s}^T(\mathbf{x}; \dot{\boldsymbol{\theta}}) \text{CCRB}(\dot{\boldsymbol{\theta}}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}).$$

To show this by contradiction, assume not and take the limit as  $\alpha \rightarrow 0$ . Then

$$\begin{aligned} \mathbf{s}^T(\mathbf{x}; \dot{\boldsymbol{\theta}}) \text{CCRB}(\dot{\boldsymbol{\theta}}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}) &\leq \kappa \mathbf{s}^T(\mathbf{x}; \dot{\boldsymbol{\theta}}) \text{CCRB}(\dot{\boldsymbol{\theta}}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}) \\ 0 &\leq (\kappa - 1) \mathbf{s}^T(\mathbf{x}; \dot{\boldsymbol{\theta}}) \text{CCRB}(\dot{\boldsymbol{\theta}}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}). \end{aligned}$$

This inequality implies  $\text{CCRB}(\dot{\boldsymbol{\theta}}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}) = \mathbf{0}$  and  $\mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}) \in \text{span}(\mathbf{F}^T(\dot{\boldsymbol{\theta}}))$  since  $\kappa < 1$ , so  $\dot{\boldsymbol{\theta}}$  satisfies the stationarity condition (3.32).  $\square$

**Theorem** (Theorem 3.34). The sequence  $\{p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})\}$  is a monotone increasing sequence. Furthermore, if  $p(\mathbf{x}; \cdot)$  is bounded above, then  $\{p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})\}$  converges.

*Proof.* Since  $\kappa \geq 0$  and  $\|\dot{\boldsymbol{\theta}}^{(k+1)} - \dot{\boldsymbol{\theta}}^{(k)}\|_{I(\dot{\boldsymbol{\theta}}^{(k)})}^2 \geq 0$ , then by the rule in (3.37), the value of the likelihood function can only increase after each iteration. The second

statement is a consequence of the monotone convergence principle, i.e., a bounded monotone sequence converges [38, p. 44, theorem 2-6].  $\square$

**Theorem** (Theorem 3.35). If the likelihood  $p(\mathbf{x}; \cdot)$  is bounded above, then the sequence

$$\{\log p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k+1)}) - \log p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})\}$$

vanishes.

*Proof.* Since  $\{p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})\}$  converges, then  $\frac{p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k+1)})}{p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})} \rightarrow 1$ .<sup>1</sup> And since  $\log(\cdot)$  is continuous, then  $\log \frac{p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k+1)})}{p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})} \rightarrow 0$ .  $\square$

**Theorem** (Theorem 3.36). If the likelihood  $p(\mathbf{x}; \cdot)$  is bounded above, then the sequence

$$\{\|\dot{\boldsymbol{\theta}}^{(k+1)} - \dot{\boldsymbol{\theta}}^{(k)}\|_{I(\dot{\boldsymbol{\theta}}^{(k)})}\}$$

vanishes as  $k \rightarrow \infty$ .

*Proof.* Again, by the rule in (3.37), the sequence  $\{\|\dot{\boldsymbol{\theta}}^{(k+1)} - \dot{\boldsymbol{\theta}}^{(k)}\|_{I(\dot{\boldsymbol{\theta}}^{(k)})}^2\}$  is bounded above by the product of a bounded sequence  $\{\alpha^{(k)}\}$  and a vanishing sequence  $\{\log p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k+1)}) - \log p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})\}$ , and clearly each element of the sequence is non-negative. Hence, by the squeezing theorem<sup>2</sup>, the sequence vanishes as  $k \rightarrow \infty$ .  $\square$

**Theorem** (Theorem 3.37). If  $\Theta_{\dot{\boldsymbol{\theta}}(1)}$  is compact and convex, then limit points of the sequence  $\{\dot{\boldsymbol{\theta}}^{(k)}\}$  are also stationary points.

---

<sup>1</sup>If  $a_k \rightarrow a$ ,  $b_k \rightarrow b$ , and  $b_k \neq 0$  for any  $k$ , then  $\frac{a_k}{b_k} \rightarrow \frac{a}{b}$  [38, p. 41, theorem 2-4(d)].

<sup>2</sup>If  $0 \leq a_k \leq b_k$  and  $b_k \rightarrow 0$ , then  $a_k \rightarrow 0$  [38, p. 43, theorem 2-5(c)].

*Proof.* Let  $\boldsymbol{\theta}^*$  be a limit point of the sequence  $\{\dot{\boldsymbol{\theta}}^{(k)}\}$ . Then by virtue of Bolzano-Weierstrass [38, p. 54, theorem 2-14], there exists a convergent subsequence  $\{\dot{\boldsymbol{\theta}}^{(k_{i_j})}\}$  that converges to  $\boldsymbol{\theta}^*$ . Additionally, since  $\{\alpha^{(k_i)}\}$  is a bounded sequence, it contains a convergent subsequence  $\{\alpha^{(k_{i_j})}\}$  with a limit point we shall denote  $\alpha^*$ . It is still true that  $\dot{\boldsymbol{\theta}}^{(k_{i_j})} \rightarrow \boldsymbol{\theta}^*$  [38, p.49, theorem 2-10]. Then we can bound the norm-distance between  $\boldsymbol{\pi} [\boldsymbol{\theta}^* + \alpha^* \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*)]$  and  $\boldsymbol{\theta}^*$  using the triangle inequality, e.g.,

$$\begin{aligned}
& \left\| \boldsymbol{\pi} [\boldsymbol{\theta}^* + \alpha^* \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*)] - \boldsymbol{\theta}^* \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} \\
& \leq \left\| \boldsymbol{\pi} [\boldsymbol{\theta}^* + \alpha^* \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*)] - \boldsymbol{\pi} [\dot{\boldsymbol{\theta}}^{(k_{i_j})} + \alpha^{(k_{i_j})} \text{CCRB}(\dot{\boldsymbol{\theta}}^{(k_{i_j})}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k_{i_j})})] \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} \\
& \quad + \left\| \boldsymbol{\pi} [\dot{\boldsymbol{\theta}}^{(k_{i_j})} + \alpha^{(k_{i_j})} \text{CCRB}(\dot{\boldsymbol{\theta}}^{(k_{i_j})}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k_{i_j})})] - \dot{\boldsymbol{\theta}}^{(k_{i_j})} \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} + \left\| \dot{\boldsymbol{\theta}}^{(k_{i_j})} - \boldsymbol{\theta}^* \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} \\
& \leq \left\| \boldsymbol{\theta}^* + \alpha^* \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*) - \dot{\boldsymbol{\theta}}^{(k_{i_j})} - \alpha^{(k_{i_j})} \text{CCRB}(\dot{\boldsymbol{\theta}}^{(k_{i_j})}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k_{i_j})}) \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} \\
& \quad + \left\| \dot{\boldsymbol{\theta}}^{(k_{i_j}+1)} - \dot{\boldsymbol{\theta}}^{(k_{i_j})} \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} + \left\| \dot{\boldsymbol{\theta}}^{(k_{i_j})} - \boldsymbol{\theta}^* \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)}.
\end{aligned}$$

The second inequality is a result of a property of projections on convex sets for the first term and the definition of our method of scoring with constraints for the second term. Note this second term will vanish by theorem 3.36 as  $k_{i_j} \rightarrow \infty$ . Also, the third term will vanish as  $k_{i_j} \rightarrow \infty$  since  $\boldsymbol{\theta}^*$  is the limit of the sequence  $\{\dot{\boldsymbol{\theta}}^{(k_{i_j})}\}$ . The first term is bounded, using the triangle inequality again, as in

$$\begin{aligned}
& \left\| \boldsymbol{\theta}^* + \alpha^* \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*) - \dot{\boldsymbol{\theta}}^{(k_{i_j})} - \alpha^{(k_{i_j})} \text{CCRB}(\dot{\boldsymbol{\theta}}^{(k_{i_j})}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k_{i_j})}) \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} \\
& \leq \left\| \dot{\boldsymbol{\theta}}^{(k_{i_j})} - \boldsymbol{\theta}^* \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} + \left\| \alpha^* \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*) - \alpha^{(k_{i_j})} \text{CCRB}(\dot{\boldsymbol{\theta}}^{(k_{i_j})}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k_{i_j})}) \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)}.
\end{aligned}$$

Again, this first term will vanish as  $k_{i_j} \rightarrow \infty$ . This last term, using the triangle

inequality, satisfies

$$\begin{aligned}
& \left\| \alpha^* \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*) - \alpha^{(k_{i_j})} \text{CCRB}(\dot{\boldsymbol{\theta}}^{(k_{i_j})}) \mathbf{s}(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k_{i_j})}) \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} \\
& \leq \left\| \alpha^* \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*) - \alpha^{(k_{i_j})} \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*) \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} \\
& \quad + \left\| \alpha^{(k_{i_j})} \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*) - \alpha^{(k_{i_j})} \text{CCRB}(\boldsymbol{\theta}^{(k_{i_j})}) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^{(k_{i_j})}) \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} \\
& \leq \left\| \alpha^{(k_{i_j})} - \alpha^* \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} \left\| \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*) \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} \\
& \quad + \left\| \alpha^{(k_{i_j})} \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)} \left\| \text{CCRB}(\boldsymbol{\theta}^{(k_{i_j})}) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^{(k_{i_j})}) - \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*) \right\|_{\mathbf{I}(\boldsymbol{\theta}^*)}.
\end{aligned}$$

The second inequality used the distributive property of norms [19, p.170, theorem 6.9.2]. By compactness,  $\{\|\text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*)\|_{\mathbf{I}(\boldsymbol{\theta}^*)}\}$  is bounded. So the first term above will vanish as  $k_{i_j} \rightarrow \infty$  since  $\alpha^*$  is the limit of the sequence  $\{\alpha^{(k_{i_j})}\}$ . Similarly,  $\{\|\alpha^{(k_{i_j})}\|_{\mathbf{I}(\boldsymbol{\theta}^*)}\}$  is a bounded sequence, and so the second term above will vanish as  $k_{i_j} \rightarrow \infty$  since  $\text{CCRB}(\boldsymbol{\theta}^{(k_{i_j})}) \rightarrow \text{CCRB}(\boldsymbol{\theta}^*)$  and  $\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^{(k_{i_j})}) \rightarrow \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*)$  by continuity [38, p.78, corollary 4-2]. Therefore,

$$\pi [\boldsymbol{\theta}^* + \alpha^* \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*)] = \boldsymbol{\theta}^*$$

and one of the following holds:

- (a)  $\alpha^* = 0$ ,
- (b)  $\text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*) = \mathbf{0}$ , or
- (c) the step projection  $\boldsymbol{\theta}^* + \alpha^* \text{CCRB}(\boldsymbol{\theta}^*) \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*)$  is perpendicular to  $\Theta_f$  at  $\boldsymbol{\theta}^*$ .

This last case (c) is impossible since the step is directed by linear combinations of the vectors of  $\mathbf{U}(\boldsymbol{\theta}^*)$  which are tangent to the constraint space at  $\boldsymbol{\theta}^*$ . This first case (a) implies stationarity by applying continuity on the step size rule condition (3.37)

and then theorem 3.33. And (b) implies  $\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}^*)$  is some linear combination of the columns of  $\mathbf{F}^T(\boldsymbol{\theta}^*)$ , i.e.,  $\boldsymbol{\theta}^*$  satisfies (3.32). Therefore,  $\boldsymbol{\theta}^*$  is a stationary point.  $\square$

**Theorem** (Theorem 3.38). If  $\Theta_{\dot{\boldsymbol{\theta}}^{(1)}}$  is compact for all sequences in a closed set of  $\Theta_f$  and if there is a unique limit point  $\boldsymbol{\theta}^*$  for all such sequences then  $\lim_{k \rightarrow \infty} \dot{\boldsymbol{\theta}}^{(k)} = \boldsymbol{\theta}^*$  for every sequence  $\{\dot{\boldsymbol{\theta}}^{(k)}\}$ . Also,  $\boldsymbol{\theta}^*$  is the maximum of  $p(\mathbf{x}; \cdot)$ .

*Proof.* Let  $\dot{\boldsymbol{\theta}}^{(2)}$  be any point in the compact set  $\Theta_{\dot{\boldsymbol{\theta}}^{(1)}}$ . Since  $\{\dot{\boldsymbol{\theta}}^{(k)}\}$  resides in a compact set it has a limit point (Bolzano-Weierstrass [38, p. 52, theorem 2-12]), which must be unique and therefore  $\lim_{k \rightarrow \infty} \dot{\boldsymbol{\theta}}^{(k)} = \boldsymbol{\theta}^*$ . By theorem 3.34,  $p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k+d)}) \geq p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)})$ , and by continuity  $p(\mathbf{x}; \dot{\boldsymbol{\theta}}^{(k)}) \rightarrow p(\mathbf{x}; \boldsymbol{\theta}^*)$ . Hence,  $p(\mathbf{x}; \boldsymbol{\theta}^*) \geq p(\mathbf{x}; \boldsymbol{\theta}')$  for every  $\boldsymbol{\theta}' \in \Theta_{\dot{\boldsymbol{\theta}}^{(1)}}$ .  $\square$

## Appendix C

### Proofs of Theorems in Chapter 4

**Theorem** (Theorem 4.7). The CFIM is singular and the dimension of its null space is lower bounded as

$$\text{nullity}(\mathcal{I}(\boldsymbol{\vartheta})) \geq \sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+, \quad (\text{C.1})$$

where  $(a)_+ = a$  for  $a \geq 0$  and  $(a)_+ = 0$  for  $a < 0$ . This limit quantity is the *nullity lower bound (NLB)*.

*Proof.* This proof is by construction. We shall develop a null subspace of the submatrix  $[\mathbf{Q}_i \ \mathbf{Q}_j]$  of  $\mathcal{Q}$ . In particular, consider the submatrix of consisting of the  $i$ th source elements and  $j$  channel elements corresponding to the  $m$ th channel, i.e.,  $[\mathbf{S}^{(i)} \ \mathbf{H}_{(m)}^{(j)}]$ . There exists three case to consider: (1)  $L_i = L_j$ , (2)  $L_i > L_j$ , and (3)  $L_i < L_j$ . But first, for use in the proof, define the  $M(L_j + 1) \times (L_j - L_i + 1)_+$  matrix

$$\boldsymbol{\mathcal{H}}_{(j)}^{(i)} = \begin{bmatrix} \boldsymbol{\mathcal{H}}_{(j)1}^{(i)} \\ \boldsymbol{\mathcal{H}}_{(j)2}^{(i)} \\ \vdots \\ \boldsymbol{\mathcal{H}}_{(j)M}^{(i)} \end{bmatrix} \quad (\text{C.2})$$

where  $\mathcal{H}_{(j)m}^{(i)}$  is the  $(L_j + 1) \times (L_j - L_i + 1)_+$  matrix

$$\mathcal{H}_{(j)m}^{(i)} = \begin{bmatrix} h_m^{(i)}(0) & 0 & \cdots & 0 \\ h_m^{(i)}(1) & h_m^{(i)}(0) & \cdots & 0 \\ \vdots & h_m^{(i)}(1) & \ddots & 0 \\ h_m^{(i)}(L_i) & \vdots & \cdots & h_m^{(i)}(0) \\ 0 & h_m^{(i)}(L_i) & \cdots & h_m^{(i)}(1) \\ \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & h_m^{(i)}(L_i) \end{bmatrix}. \quad (\text{C.3})$$

Then  $\mathcal{H}_{(j)m}^{(i)} = \mathbf{h}_m^{(i)}$  and  $\mathcal{H}_{(j)}^{(i)} = \mathbf{h}^{(i)}$  from the secondary vector-matrix model in (4.6)

if  $L_i = L_j$ . And also define the  $(N + L_j) \times (L_i - L_j + 1)_+$  matrix

$$\mathcal{S}_{(j)}^{(i)} = \begin{bmatrix} s^{(i)}(-L_j) & s^{(i)}(-L_j - 1) & \cdots & s^{(i)}(-L_i) \\ s^{(i)}(-L_j + 1) & s^{(i)}(-L_j) & \cdots & s^{(i)}(-L_i + 1) \\ \vdots & \vdots & \ddots & \vdots \\ s^{(i)}(N - 1) & s^{(i)}(N - 2) & \cdots & s^{(i)}(N - L_i + L_j - 1) \end{bmatrix}. \quad (\text{C.4})$$

Then  $\mathcal{S}_{(j)}^{(i)} = \mathbf{s}^{(i)}$  if  $L_i = L_j$ . Now consider the cases:

(1) From the dual interpretation of the model in (4.3) and (4.6), it is clear that

$$\mathbf{H}_{(m)}^{(j)} \mathbf{s}^{(i)} = \mathcal{S}^{(i)} \mathbf{h}_m^{(j)}; \text{ hence,}$$

$$\begin{bmatrix} \mathbf{I}_M \otimes \mathcal{S}^{(i)} & \mathbf{H}_M^{(j)} \end{bmatrix} \begin{bmatrix} \mathcal{H}_{(i)m}^{(j)} \\ -\mathcal{S}_{(j)}^{(i)} \end{bmatrix} = \mathbf{0}.$$

Also,  $(L_i - L_j + 1)_+ = \dim(\begin{bmatrix} \mathbf{h}^{(j)} \\ -\mathbf{s}^{(i)} \end{bmatrix}) = 1$  unless  $\mathbf{h}^{(j)} = \mathbf{0}$  and  $\mathbf{s}^{(i)} = \mathbf{0}$ , in which case,  $\text{nullity}(\begin{bmatrix} \mathbf{I}_M \otimes \mathcal{S}^{(i)} & \mathbf{H}_M^{(j)} \end{bmatrix}) > (L_i - L_j + 1)_+$ .

(2)  $\mathcal{H}_{(i)m}^{(j)}$  in (C.3) has rank  $(L_i - L_j + 1)$  unless  $\mathbf{h}_m^{(j)} = \mathbf{0}$ , which is impossible by definition. Likewise,  $\mathcal{S}_{(j)}^{(i)}$  is full column rank unless  $\mathbf{s}^{(i)}$  has fewer than  $(L_i - L_j + 1)$  modes or  $N < L_i - 2L_j + 1$  (see theorem 4.3). Therefore, since

$$\begin{bmatrix} \mathbf{I}_M \otimes \mathcal{S}^{(i)} & \mathbf{H}_M^{(j)} \end{bmatrix} \begin{bmatrix} \mathcal{H}_{(i)m}^{(j)} \\ -\mathcal{S}_{(j)}^{(i)} \end{bmatrix} = \mathbf{0}$$

then  $\text{nullity}(\begin{bmatrix} \mathbf{I}_M \otimes \mathcal{S}^{(i)} & \mathbf{H}_M^{(j)} \end{bmatrix}) \geq (L_i - L_j + 1)_+$ .



(3) Finally,  $\text{nullity}\left(\begin{bmatrix} \mathbf{I}_M \otimes \mathbf{S}^{(i)} & \mathbf{H}_M^{(j)} \end{bmatrix}\right) \geq 0 \geq (L_i - L_j + 1).$

□

**Theorem** (CFIM NLB necessary conditions, Theorem 4.13). The  $M$ -channel  $K$ -source FIR system Fisher information matrix has a nullity of exactly the NLB in (4.13) only if

(a)  $\mathbf{H}(z)$  is irreducible and column-reduced,

(b)  $p_{\text{total}} \geq K + \sum_{j=1}^K L_j,$

(c)  $p_k \geq L_k + 2$  for  $k = 1, \dots, K$  or  $p_k \geq 1$  if  $L_k = 0,$

(d)  $N \geq K + \sum_{j=1}^K L_j,$  and

(e)  $M > K.$

*Proof.* If any of these conditions fail, then the null space of  $\mathcal{I}(\boldsymbol{\vartheta})$  is greater than the NLB.

(a) If  $\mathbf{H}(z)$  is reducible or not column-reduced, then by theorem 4.2  $\mathbf{H}_M$  is not full column rank. Let  $\mathbf{v} \in \text{null}(\mathbf{H}_M)$  and partition  $\mathbf{v}$  as

$$\mathbf{v} = \begin{bmatrix} \mathbf{v}^{(1)} \\ \mathbf{v}^{(2)} \\ \vdots \\ \mathbf{v}^{(K)} \end{bmatrix}$$

where  $\mathbf{v}^{(k)}$  is a  $N + L_k$  length vector. (It is assumed that  $\mathbf{v}$  is independent of  $\mathbf{s}$  otherwise  $E_{\boldsymbol{\vartheta}}\mathbf{y} = \mathbf{0}$  and the model itself is not identifiable.) Then

$$\mathbf{v}^* = \begin{bmatrix} \mathbf{0}^{(1)} \\ \mathbf{v}^{(1)} \\ \vdots \\ \mathbf{0}^{(K)} \\ \mathbf{v}^{(K)} \end{bmatrix} \in \text{null}(\mathcal{I}(\boldsymbol{\vartheta}))$$

where  $\mathbf{0}^{(k)}$  is a  $M(L_k + 1)$  length zero vector. Assume  $\mathbf{v}^*$  is a linear combination of the columns of  $\mathcal{N}$  in (4.14). Then for some  $k$  the columns of  $\begin{bmatrix} \mathcal{H}_{(k)}^{(1)} & \mathcal{H}_{(k)}^{(2)} & \cdots & \mathcal{H}_{(k)}^{(K)} \end{bmatrix}$  must be linearly dependent. (Otherwise without this assumption, then  $\mathbf{v}^* \notin \text{span}(\mathcal{N})$ . This corresponds to the  $\mathbf{0}^{(k)}$  subvector in  $\mathbf{v}^*$ .) Let

$$\mathbf{u}_{(k)} = \begin{bmatrix} \mathbf{u}_{(k)}^{(1)} \\ \mathbf{u}_{(k)}^{(2)} \\ \vdots \\ \mathbf{u}_{(k)}^{(K)} \end{bmatrix} \in \text{null}(\begin{bmatrix} \mathcal{H}_{(k)}^{(1)} & \mathcal{H}_{(k)}^{(2)} & \cdots & \mathcal{H}_{(k)}^{(K)} \end{bmatrix})$$

with  $\mathbf{u}_{(k)}^{(i)}$  a  $(L_k - L_i + 1)_+$  length vector. Then

$$\mathbf{U}_{(k)} = \begin{bmatrix} \mathbf{U}_{(k)}^{(1)} \\ \mathbf{U}_{(k)}^{(2)} \\ \vdots \\ \mathbf{U}_{(k)}^{(K)} \end{bmatrix} \in \text{null}(\mathbf{H}_M)$$

where

$$\mathbf{U}_{(k)}^{(i)} = \begin{bmatrix} \mathbf{u}_{(k)}^{(i)T} & 0 & \cdots & 0 \\ 0 & \mathbf{u}_{(k)}^{(i)T} & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & \cdots & & \mathbf{u}_{(k)}^{(i)T} \end{bmatrix}_{N+L_i \times N+L_k}.$$

(If  $L_i > L_k$  then  $\mathbf{U}_{(k)}^{(i)}$  is a null matrix. Also,  $\mathbf{U}_{(k)}^{(k)} = u_{(k)} \mathbf{I}_{N+L_k \times N+L_k}$ .) The matrix  $\mathbf{U}_{(k)}$  can be arranged to create  $N + L_k$  linearly independent columns in the null space of  $\mathcal{I}(\boldsymbol{\vartheta})$ , where the submatrices  $\mathbf{U}_{(k)}^{(i)}$  correspond to the rows of  $\mathcal{N}$  containing  $\begin{bmatrix} -\mathcal{S}_{(i)}^{(1)} & -\mathcal{S}_{(i)}^{(2)} & \cdots & -\mathcal{S}_{(i)}^{(K)} \end{bmatrix}$ , which has rank at most  $\sum_{j=1}^K (L_j - L_i + 1)_+ \leq \sum_{j=1}^K (L_j + 1)$ . For the columns of  $\mathcal{N}$  to be a basis of the null space of  $\mathcal{I}(\boldsymbol{\vartheta})$ , it is needed then that  $\sum_{j=1}^K (L_j + 1) \geq \sum_{j=1}^K (L_j - L_k + 1)_+ \geq N + L_k \geq N + L_i$ ,

which implies at most  $N = K + \sum_{j=1}^K L_j$  contradicting theorem 4.12 unless all the channel orders are zero. In this latter case, then note that

$$\mathcal{N} \cdot \left( \mathbf{I}_{K \times K} \otimes \begin{bmatrix} u_{(k)}^{(1)} \\ u_{(k)}^{(2)} \\ \vdots \\ u_{(k)}^{(K)} \end{bmatrix} \right)$$

will not be full rank, and hence neither can  $\mathcal{N}$  be so.

(b) If  $p_{\text{total}} < K + \sum_{k=1}^K L_k$ , then there exists  $\mathbf{v} \in \text{null}(\mathbf{S})$  where the matrix  $\mathbf{S} = [\mathbf{S}^{(1)} \ \mathbf{S}^{(2)} \ \dots \ \mathbf{S}^{(K)}]$ . If  $\mathbf{v}$  is partitioned as

$$\mathbf{v} = \begin{bmatrix} \mathbf{v}^{(1)} \\ \mathbf{v}^{(2)} \\ \vdots \\ \mathbf{v}^{(K)} \end{bmatrix}$$

where  $\mathbf{v}^{(k)}$  is a  $L_k + 1$  length vector, then

$$\mathbf{v}^* = \begin{bmatrix} \mathbf{1}_{M \times 1} \otimes \mathbf{v}^{(1)} \\ \mathbf{0}^{(1)} \\ \vdots \\ \mathbf{1}_{M \times 1} \otimes \mathbf{v}^{(K)} \\ \mathbf{0}^{(K)} \end{bmatrix} \in \text{null}(\mathcal{I}(\boldsymbol{\vartheta}))$$

with  $\mathbf{0}^{(j)}$  is an  $N + L_j$  length zero vector. If  $\mathbf{v}^* \in \text{span}(\mathcal{N})$ , then for some  $k$  we have  $\text{rank}([\mathbf{S}_{(k)}^{(1)} \ \mathbf{S}_{(k)}^{(2)} \ \dots \ \mathbf{S}_{(k)}^{(K)}]) < \sum_{j=1}^K (L_j - L_k + 1)_+$ . Then by a construction similar to part (a),  $\text{nullity}(\mathbf{S}) \geq L_k + 1$  and this contribution to

the null space of  $\mathcal{I}(\boldsymbol{\vartheta})$  has rank at least  $M(L_k + 1)$ , with submatrices that coincide with  $[\mathcal{H}_{(k)}^{(1)} \ \mathcal{H}_{(k)}^{(2)} \ \dots \ \mathcal{H}_{(k)}^{(K)}]$ , which has rank at most  $\sum_{j=1}^K (L_k -$

$$L_j + 1)_+ \leq \sum_{j=1}^K (L_k + 1) = K(L_k + 1) < M(L_k + 1) \text{ for any } L_k \text{ since } M > K.$$

Therefore,  $\text{nullity}(\mathcal{I}(\boldsymbol{\vartheta})) > \sum_{i=1}^K \sum_{j=1}^K (L_i - L_j + 1)_+$ .

(c) If  $p_k < L_k + 1$  then by [31, lemma 1],  $\mathbf{S}^{(k)}$  and, consequently  $\mathbf{S}$ , has a null space and the argument in (b) applies. So assume  $p_k = L_k + 1$ . If  $N < L_k + 1$  (and  $L_k \neq 0$ ) then  $\mathbf{S}^{(k)}$  has a null space [76, lemma 1], so assume  $N \geq L_k + 1$ . From [31, the proof of theorem 1], it is possible to construct a  $\mathbf{v}$  independent from  $\mathbf{s}^{(k)}$  such that  $\text{span}(\mathbf{V}) = \text{span}(\mathbf{S}^{(k)})$  for  $\mathbf{V}$  defined similarly as  $\mathbf{S}^{(k)}$ . So for any  $\mathbf{h}^{(k)}$  there exists a  $\mathbf{h}^*$  such that  $\mathbf{H}_M^{(k)} \mathbf{v} = (\mathbf{I}_{M \times M} \otimes \mathbf{V}) \mathbf{h}^{(k)} = (\mathbf{I}_{M \times M} \otimes \mathbf{S}^{(k)}) \mathbf{h}^*$ . Therefore both  $\begin{bmatrix} -\mathbf{v} \\ \mathbf{h}^* \end{bmatrix}$  and  $\begin{bmatrix} -\mathbf{s}^{(k)} \\ \mathbf{h}^{(k)} \end{bmatrix}$  reside in  $\text{null}\left(\begin{bmatrix} \mathbf{H}_M^{(k)} & (\mathbf{I}_{M \times M} \otimes \mathbf{S}^{(k)}) \end{bmatrix}\right)$ , which increases the nullity lower bound by at least one.

(d) This is a looser bound and hence required by theorem 4.12.

□

**Theorem** (CFIM NLB sufficiency conditions, theorem 4.14). The  $M$ -channel  $K$ -source FIR system FIM has a nullity of exactly the NLB in (4.13) if

(a)  $\mathbf{H}(z)$  is irreducible and column-reduced,

$$(b) \ p_{\text{total}} \geq K + (K + 1) \sum_{j=1}^K L_j,$$

$$(c) \ p_k \geq L_k + 1 + \sum_{j=1}^K L_j \text{ for } k = 1, \dots, K,$$

$$(d) \ N \geq K + (K + 2) \sum_{j=1}^K L_j, \text{ and}$$

(e)  $M > K$ .

This result is conceptually easier to prove from yet another alternative matrix model from the ones in (4.3) or (4.6) using  $\mathbf{S}(n)$  as defined in (4.9). Then define

the  $(L_k + 1 + n) \times (n + 1)$  matrix

$$\mathbf{H}_{(m)}^{(k)}(n) = \begin{bmatrix} h_m^{(k)}(0) & & \\ \vdots & \ddots & h_m^{(k)}(0) \\ h_m^{(k)}(L_k) & \ddots & \vdots \\ & & h_m^{(k)}(L_k) \end{bmatrix},$$

which is the impulse response for the  $i$ th subchannel of the  $k$ th source. Then define

$$\mathbf{H}^{(k)}(n) = [\mathbf{H}_{(1)}^{(k)}(n), \dots, \mathbf{H}_{(M)}^{(k)}(n)] \text{ and the } (K(n+1) + \sum_{k=1}^K L_k) \times M(n+1) \text{ matrix}$$

$$\mathbf{H}(n) = \begin{bmatrix} \mathbf{H}^{(1)}(n) \\ \vdots \\ \mathbf{H}^{(K)}(n) \end{bmatrix}.$$

The observations for the  $m$ th channels (receiver) can be collected into the  $(N - n) \times (n + 1)$  matrix

$$\mathbf{Y}_{(m)}(n) = \begin{bmatrix} y_m(n) & \cdots & y_m(0) \\ \vdots & & \vdots \\ y_m(N-1) & \cdots & y_m(N-1-n) \end{bmatrix}$$

with  $\mathbf{Y}(n) = [\mathbf{Y}_{(1)}(n), \dots, \mathbf{Y}_{(M)}(n)]$ . The noise matrix  $\mathbf{W}(n)$  is defined similarly.

Then the alternative model is

$$\mathbf{Y}(n) = \mathbf{S}(n)\mathbf{H}(n) + \mathbf{W}(n). \quad (\text{C.5})$$

Before this theorem is proven, a lemma will be needed. This lemma is a generalization of a result in [52, theorem 3]. The proof was originally shown in [49, Appendix].

**Lemma C.1.** Assume  $\mathbf{H}(n)$  be full row rank and  $\mathbf{h}'^{(k)}$  be any nontrivial  $M(L_k + 1)$  length vector, and define  $\mathbf{H}'^{(k)}(n^*)$  to be the  $L_k + 1 + n^* \times M(1 + n^*)$  matrix composed from  $\mathbf{h}'^{(k)}$  as in (4.4) and (4.5). Then the following two statements are equivalent:

$$(i) \quad \text{corange}\{\mathbf{H}'^{(k)}(n^*)\} \subset \text{corange}\{\mathbf{H}^{(1)}(n^*), \dots, \mathbf{H}^{(K)}(n^*)\} = \text{corange}\{\mathbf{H}(n^*)\}.$$

$$(ii) \mathbf{h}^{(k)} \in \text{range}\{\mathcal{H}_{(k)}^{(1)}, \dots, \mathcal{H}_{(k)}^{(K)}\}.$$

We shall prove a more general version of the lemma, which is essentially an extension of the identifiability theorem of [52, theorem 3] from the SIMO to the MIMO scenario. First, define  $\mathbf{h}^{(j)}(l) = [h_{(1)}^{(j)}(l), \dots, h_{(M)}^{(j)}(l)]^T$  and the  $Mn \times (n + L_j)$  matrix

$$\mathbf{H}^{(j)}(n) = \begin{bmatrix} \mathbf{h}^{(j)}(L_j) & \mathbf{h}^{(j)}(L_j - 1) & \cdots & \mathbf{h}^{(j)}(0) \\ & \mathbf{h}^{(j)}(L_j) & \mathbf{h}^{(j)}(L_j - 1) & \cdots & \mathbf{h}^{(j)}(0) \\ & & \ddots & \ddots & \ddots \\ & & & \mathbf{h}^{(j)}(L_j) & \mathbf{h}^{(j)}(L_j - 1) & \cdots & \mathbf{h}^{(j)}(0) \end{bmatrix}$$

so that  $\mathbf{H}(N) = [\mathbf{H}^{(1)}(N), \dots, \mathbf{H}^{(K)}(N)]$  is similar to  $\mathbf{H}_M$  in the original model (4.3), except for the order of the rows. Let  $L$  satisfy  $\min_j \{L_j\} \leq L \leq \max_j \{L_j\}$ .

Also define  $\mathcal{H}_L$  as the  $M(L + 1) \times (L - L_j + 1)_+$  matrix

$$\mathcal{H}_L^{(j)} = \begin{bmatrix} \mathbf{h}^{(j)} & \mathbf{0}_{M \times 1} & \cdots & \mathbf{0}_{M(L-L_j-1) \times 1} \\ \mathbf{0}_{M \times 1} & \mathbf{h}^{(j)} & \ddots & \mathbf{0}_{M \times 1} \\ \mathbf{0}_{M(L-L_j-1) \times 1} & \mathbf{0}_{M(L-L_j-1) \times 1} & & \mathbf{h}^{(j)} \end{bmatrix}$$

where  $\mathbf{h}^{(j)} = [\mathbf{h}^{(j)}(0)^T, \dots, \mathbf{h}^{(j)}(L_j)^T]^T$ , an  $M(L_j + 1) \times 1$  matrix. (The matrix is not to scale, i.e.,  $\mathbf{0}_{M \times 1}$  does not align with  $\mathbf{h}^{(j)}$ .) Note, that  $\mathcal{H}_L^{(j)}$  is null for  $L < L_j$ .

Now, we restate the lemma as the following. Instead, we show the following are equivalent:

$$(1) \text{Range}\{\mathbf{H}'(N)\} \subset \text{Range}\{\mathbf{H}^{(1)}(N), \dots, \mathbf{H}^{(K)}(N)\} = \text{Range}\{\mathbf{H}(N)\}.$$

$$(2) \mathbf{h}' \in \text{Range}\{\mathcal{H}_L^{(1)}, \dots, \mathcal{H}_L^{(K)}\}.$$

*Proof.* Assume  $N > \sum_{j=1}^K L_j$  and  $\mathbf{H}(N - 1)$  is full-column rank. Let  $\mathbf{h}'$  be any  $M(L + 1) \times 1$  nonzero complex vector, and define  $\mathbf{H}'(N)$  to be the  $MN \times N + L$  channel matrix composed from  $\mathbf{h}'$ . We'll first assume (1), and show (2). Note that

$$\mathbf{H}^{(j)}(N) = \begin{bmatrix} \mathbf{h}^{(j)}(0) & \mathbf{p}^{(j)}(N - 1) \\ \mathbf{0}_{M(N-1) \times 1} & \mathbf{H}^{(j)}(N - 1) \end{bmatrix} = \begin{bmatrix} \mathbf{H}^{(j)}(N - 1) & \mathbf{0}_{M(N-1) \times 1} \\ \mathbf{q}^{(j)}(N - 1) & \mathbf{h}^{(j)}(L_j) \end{bmatrix}$$

where

$$\begin{aligned}\mathbf{p}^{(j)}(N-1) &= [\mathbf{h}^{(j)}(1), \dots, \mathbf{h}^{(j)}(L_j), \mathbf{0}_{M \times N-1}], \text{ and} \\ \mathbf{q}^{(j)}(N-1) &= [\mathbf{0}_{M \times N-1}, \mathbf{h}^{(j)}(0), \dots, \mathbf{h}^{(j)}(L_j-1)].\end{aligned}$$

Note statement (1) implies the first column of  $\mathbf{H}'(N)$  satisfies

$$\begin{bmatrix} \mathbf{h}'(0) \\ \mathbf{0}_{M(N-1) \times 1} \end{bmatrix} = \sum_{j=1}^K \begin{bmatrix} \mathbf{h}^{(j)}(0) & \mathbf{p}^{(j)}(N-1) \\ \mathbf{0}_{M(N-1) \times 1} & \mathbf{H}^{(j)}(N-1) \end{bmatrix} \cdot \begin{bmatrix} \alpha_0^{(j)} \\ \mathbf{a}_0^{(j)} \end{bmatrix}$$

where  $\alpha_0^{(j)}$  is a constant and  $\mathbf{a}_0^{(j)}$  is an  $N-1+L_j \times 1$  vector. Hence, we have the following two linear systems:

$$(A1) \quad \mathbf{h}'(0) = \sum_{j=1}^K \alpha_0^{(j)} \mathbf{h}^{(j)}(0) + \mathbf{p}^{(j)}(N-1) \mathbf{a}_0^{(j)}$$

$$(A2) \quad \mathbf{0}_{M(N-1) \times 1} = \sum_{j=1}^K \mathbf{H}^{(j)}(N-1) \mathbf{a}_0^{(j)}$$

But  $\mathbf{H}(N-1)$  is full-column rank, thus  $\mathbf{a}_0^{(j)} = \mathbf{0}_{N+L_j-1 \times 1}$  for all  $j$ . Therefore,

$$\mathbf{h}'(0) = \sum_{j=1}^K \alpha_0^{(j)} \mathbf{h}^{(j)}(0). \quad (C.6)$$

Likewise, the next column of  $\mathbf{H}'(N)$  is given by

$$\begin{bmatrix} \mathbf{h}'(1) \\ \mathbf{h}'(0) \\ \mathbf{0}_{M(N-2) \times 1} \end{bmatrix} = \sum_{j=1}^K \begin{bmatrix} \mathbf{h}^{(j)}(0) & \mathbf{p}^{(j)}(N-1) \\ \mathbf{0}_{M(N-1) \times 1} & \mathbf{H}^{(j)}(N-1) \end{bmatrix} \cdot \begin{bmatrix} \alpha_1^{(j)} \\ \mathbf{a}_1^{(j)} \end{bmatrix}$$

where, again,  $\alpha_1^{(j)}$  is a constant and  $\mathbf{a}_1^{(j)}$  is an  $N+L_j-1 \times 1$  vector. Now, we have the following two systems:

$$(B1) \quad \mathbf{h}'(1) = \sum_{j=1}^K \alpha_1^{(j)} \mathbf{h}^{(j)}(0) + \mathbf{p}^{(j)}(N-1) \mathbf{a}_1^{(j)}$$

$$(B2) \quad \begin{bmatrix} \mathbf{h}'(0) \\ \mathbf{0}_{M(N-2) \times 1} \end{bmatrix} = \sum_{j=1}^K \mathbf{H}^{(j)}(N-1) \mathbf{a}_1^{(j)}$$

Note (C.6) implies that  $\begin{bmatrix} \mathbf{h}'(0) \\ \mathbf{0}_{M(N-2) \times 1} \end{bmatrix} = \sum_{j=1}^K \alpha_0^{(j)} \cdot \begin{bmatrix} \mathbf{h}^{(j)}(0) \\ \mathbf{0}_{M(N-2) \times 1} \end{bmatrix}$ , i.e., it is a linear combination of column vectors of  $\mathbf{H}(N-1)$ . Since  $\mathbf{H}(N-1)$  is full-column rank, then (B2) implies that  $\mathbf{a}_1^{(j)} = \begin{bmatrix} \alpha_0^{(j)} \\ \mathbf{0}_{N+L_j-2 \times 1} \end{bmatrix}$ . Thus, evaluating (B1), we have

$$\mathbf{h}'(1) = \sum_{j=1}^K \alpha_1^{(j)} \mathbf{h}^{(j)}(0) + \alpha_0^{(j)} \mathbf{h}^{(j)}(1). \quad (\text{C.7})$$

Continuing in this fashion, we arrive at the general expression

$$\mathbf{h}'(l) = \sum_{j=1}^K \sum_{i=0}^l \alpha_{l-i}^{(j)} \mathbf{h}^{(j)}(i) \quad (\text{C.8})$$

for  $0 \leq l \leq L$ . For convenience, we define  $\mathbf{h}^{(j)}(i)$  to be null if  $i < L_j$  or  $i < 0$ . Now, consider the  $(L+2)$ nd column of  $\mathbf{H}'(N)$ ,

$$\begin{bmatrix} \mathbf{0}_{M \times 1} \\ \mathbf{h}'(L) \\ \vdots \\ \mathbf{h}'(0) \\ \mathbf{0}_{M(N-L-2) \times 1} \end{bmatrix} = \sum_{j=1}^K \begin{bmatrix} \mathbf{h}^{(j)}(0) & \mathbf{p}^{(j)}(N-1) \\ \mathbf{0}_{M(N-1) \times 1} & \mathbf{H}^{(j)}(N-1) \end{bmatrix} \cdot \begin{bmatrix} \alpha_{L+1}^{(j)} \\ \mathbf{a}_{L+1}^{(j)} \end{bmatrix}$$

It follows that

$$(\text{C1}) \quad \mathbf{0}_{M \times 1} = \sum_{j=1}^K \alpha_{L+1}^{(j)} \mathbf{h}^{(j)}(0) + \mathbf{p}^{(j)}(N-1) \mathbf{a}_{L+1}^{(j)}$$

$$(\text{C2}) \quad \begin{bmatrix} \mathbf{h}'(L) \\ \vdots \\ \mathbf{h}'(0) \\ \mathbf{0}_{M(N-L-2) \times 1} \end{bmatrix} = \sum_{j=1}^K \mathbf{H}^{(j)}(N) \mathbf{a}_{L+1}^{(j)}$$

Since  $\mathbf{H}(N-1)$  is full-column rank, equation (C.8) and (C2) imply that  $\mathbf{a}_{L+1}^{(j)} = \left[ \alpha_L^{(j)}, \dots, \alpha_0^{(j)}, \mathbf{0}_{1 \times M(N-L-2)} \right]^T$ . Applied to (C1) we have

$$\mathbf{0}_{M \times 1} = \sum_{j=1}^K \alpha_{L+1}^{(j)} \mathbf{h}^{(j)}(0) + \alpha_L^{(j)} \mathbf{h}^{(j)}(1) + \dots + \alpha_{L-L_j+1}^{(j)} \mathbf{h}^{(j)}(L_j)$$



Continuing in this fashion, we see that for  $L + 1 \leq l \leq N - 1$ ,

$$\mathbf{0}_{M \times 1} = \sum_{j=1}^K \sum_{i=0}^{L_j} \alpha_{l-i}^{(j)} \mathbf{h}^{(j)}(i) \quad (\text{C.9})$$

In particular, (C.8) and (C.9) show that the  $N$ th column of  $\mathbf{H}'(N)$  can be written as

$$\begin{bmatrix} \mathbf{0}_{M(N-L-1) \times 1} \\ \mathbf{h}'(L) \\ \vdots \\ \mathbf{h}'(0) \end{bmatrix} = \sum_{j=1}^K \mathbf{H}^{(j)}(N) \cdot \begin{bmatrix} \alpha^{(j)}(N-1) \\ \vdots \\ \alpha_0^{(j)} \\ \mathbf{0}_{L_j \times 1} \end{bmatrix} \quad (\text{C.10})$$

Similarly, statement (1) implies the last column of  $\mathbf{H}'(N)$  satisfies

$$\begin{bmatrix} \mathbf{0}_{M(N-1) \times 1} \\ \mathbf{h}'(L) \end{bmatrix} = \sum_{j=1}^K \begin{bmatrix} \mathbf{H}^{(j)}(N-1) & \mathbf{0}_{M(N-1) \times 1} \\ \mathbf{q}^{(j)}(N-1) & \mathbf{h}^{(j)}(L_j) \end{bmatrix} \cdot \begin{bmatrix} \mathbf{y}_L^{(j)} \\ \beta_L^{(j)} \end{bmatrix}$$

where  $\mathbf{y}_L^{(j)}$  is an  $N + L_j + 1 \times 1$  vector and  $\beta_L^{(j)}$  a constant. Thus we have

$$(a1) \quad \mathbf{0}_{M(N-1) \times 1} = \sum_{j=1}^K \mathbf{H}^{(j)}(N-1) \mathbf{y}_L^{(j)}$$

$$(a2) \quad \mathbf{h}'(L) = \sum_{j=1}^K \mathbf{q}^{(j)}(N-1) \mathbf{y}_L^{(j)} + \beta_L^{(j)} \mathbf{h}^{(j)}(L_j)$$

Here, since  $\mathbf{H}(N-1)$  is full rank, then (a1) implies that  $\mathbf{y}_L^{(j)} = \mathbf{0}_{N+L_j-1 \times 1}$ , hence

$$\mathbf{h}'(L) = \sum_{j=1}^K \beta_L^{(j)} \mathbf{h}^{(j)}(L_j)$$

Proceeding as before, it is clear that

$$\mathbf{h}'(L-l) = \sum_{j=1}^K \sum_{i=0}^l \beta_{L-l}^{(j)} \mathbf{h}^{(j)}(L_j - l + i)$$

for  $0 \leq l \leq L$ . Thus, the  $N$ th column of  $\mathbf{H}'(N)$  can also be expressed as

$$\begin{bmatrix} \mathbf{0}_{M(N-L-1) \times 1} \\ \mathbf{h}'(L) \\ \vdots \\ \mathbf{h}'(0) \end{bmatrix} = \sum_{j=1}^K \mathbf{H}^{(j)}(N) \cdot \begin{bmatrix} \mathbf{0}_{N+L_j-L-1 \times 1} \\ \beta_L^{(j)} \\ \vdots \\ \beta_0^{(j)} \end{bmatrix} \quad (\text{C.11})$$

But  $\mathbf{H}(N)$  is full-column rank, hence (C.10) and (C.11) imply that

$$\begin{bmatrix} \alpha^{(j)}(N-1) \\ \vdots \\ \alpha_0^{(j)} \\ \mathbf{0}_{L_j \times 1} \end{bmatrix} = \begin{bmatrix} \mathbf{0}_{N+L_j-L-1 \times 1} \\ \beta_L^{(j)} \\ \vdots \\ \beta_0^{(j)} \end{bmatrix}$$

Thus,  $\alpha_l^{(j)} = 0$  for all  $l \geq L+1$ . Now, consider the following three cases:

- (i)  $L = L_j$ , then  $\alpha_0^{(j)} = \beta_L^{(j)}$  and  $\alpha_l^{(j)}, \beta_{L-l}^{(j)} = 0$  for  $1 \leq l \leq L$ .
- (ii)  $L < L_j$ , then  $\alpha_l^{(j)}, \beta_{L-l}^{(j)} = 0$  for  $0 \leq l \leq L$ .
- (iii)  $L > L_j$ , then  $\alpha_l^{(j)} = \beta_{L_j+l}^{(j)}$  for  $0 \leq l \leq L-L_j$  and  $\alpha_l^{(j)}, \beta_{L-l}^{(j)} = 0$  for  $L-L_j+1 \leq l \leq L$ .

Define an  $(L-L_j+1)^+ \times 1$  vector  $\boldsymbol{\alpha}^{(j)} = [\alpha_0^{(j)}, \dots, \alpha_{L-L_j}^{(j)}]^T$  for each  $j$ . If  $L < L_j$ , let  $\boldsymbol{\alpha}$  be null. Then,

$$\mathbf{h}' = \sum_{j=1}^K \mathbf{H}_L^{(j)} \boldsymbol{\alpha}^{(j)}$$

or  $\mathbf{h}' \in \text{Range} \{ \mathbf{H}_L^{(1)}, \dots, \mathbf{H}_L^{(K)} \}$ . This proves statement (2).

Now, conversely assume statement (2). Then  $\mathbf{h}' = \sum_{j=1}^K \mathbf{H}_L^{(j)} \boldsymbol{\gamma}^{(j)}$  where  $\boldsymbol{\gamma}^{(j)} = [\gamma_1^{(j)}, \dots, \gamma_{(L-L_j+1)^+}^{(j)}]^T$  is an  $(L-L_j+1)^+ \times 1$  vector. Define

$$\mathbf{\Gamma}^{(j)}(N) = \begin{bmatrix} \gamma_1^{(j)} & \gamma_2^{(j)} & \cdots & \gamma_{(L-L_j+1)^+}^{(j)} & & \\ & \gamma_1^{(j)} & \gamma_2^{(j)} & \cdots & \gamma_{(L-L_j+1)^+}^{(j)} & \\ & & \ddots & \ddots & \ddots & \\ & & & \gamma_1^{(j)} & \gamma_2^{(j)} & \cdots & \gamma_{(L-L_j+1)^+}^{(j)} \end{bmatrix}_{N+L_j \times N+L}$$

Note, if  $L < L_j$ , then  $\mathbf{\Gamma}^{(j)}(N) = \mathbf{0}_{N+L_j \times N+L}$ . Then, it is simple to verify that

$$\mathbf{H}'(N) = \sum_{j=1}^K \mathbf{H}^{(j)}(N) \mathbf{\Gamma}^{(j)}(N)$$

which proves statement (1). □

Now having this lemma at hand, theorem 4.14 can be proven.

*Proof.* Let  $n^* = \sum_{j=1}^K L_j$ . If (a)-(d) are satisfied, then we have  $\text{column-rank}(\mathbf{S}(n^*)) = \text{row-rank}(\mathbf{H}(n^*)) = K + (K + 1)n^*$ . Let  $\mathbf{v} \in \text{null}(\mathcal{I}(\boldsymbol{\vartheta}))$  where  $\mathbf{v}$  is partitioned as

$$\mathbf{v} = \begin{bmatrix} \mathbf{h}'^{(1)} \\ \mathbf{s}'^{(1)} \\ \vdots \\ \mathbf{h}'^{(K)} \\ \mathbf{s}'^{(K)} \end{bmatrix}.$$

Then  $\sum_{k=1}^K (\mathbf{I}_{M \times M} \otimes \mathbf{S}^{(k)}) \mathbf{h}'^{(k)} + \mathbf{H}_M^{(k)} \mathbf{s}'^{(k)} = \mathbf{0}$ , or in the alternative model in (C.5),

we have  $\sum_{k=1}^K \mathbf{S}^{(k)}(n^*) \mathbf{H}'^{(k)}(n^*) + \sum_{k=1}^K \mathbf{S}'^{(k)}(n^*) \mathbf{H}^{(k)}(n^*) = \mathbf{0}$ , or

$$[\mathbf{S}(n^*) \mathbf{S}'(n^*)] \begin{bmatrix} \mathbf{H}'(n^*) \\ \mathbf{H}(n^*) \end{bmatrix} = \mathbf{0}.$$

Therefore  $\text{nullity}([\mathbf{S}(n^*) \mathbf{S}'(n^*)]) \geq \text{rank}(\begin{bmatrix} \mathbf{H}'(n^*) \\ \mathbf{H}(n^*) \end{bmatrix})$ . Since

$$\begin{aligned} \text{nullity}([\mathbf{S}(n^*) \mathbf{S}'(n^*)]) &= \text{columns}([\mathbf{S}(n^*) \mathbf{S}'(n^*)]) - \text{rank}([\mathbf{S}(n^*) \mathbf{S}'(n^*)]) \\ &\leq 2K + 2(K + 1)n^* - \text{rank}(\mathbf{S}(n^*)) \\ &= K + (K + 1)n^* \end{aligned}$$

and

$$\text{rank}(\begin{bmatrix} \mathbf{H}'(n^*) \\ \mathbf{H}(n^*) \end{bmatrix}) \geq \text{rank}(\mathbf{H}(n^*)) = K + (K + 1)n^*,$$

then  $\text{nullity}([\mathbf{S}(n^*) \mathbf{S}'(n^*)]) = \text{rank}(\begin{bmatrix} \mathbf{H}'(n^*) \\ \mathbf{H}(n^*) \end{bmatrix}) = K + (K + 1)n^*$ . Since  $\text{rank}(\mathbf{H}(n^*))$

$= K + (K + 1)n^*$  then there exist some matrix  $\mathbf{T}$  such that  $\mathbf{H}'(n^*) = \mathbf{T} \mathbf{H}(n^*)$ .

Thus, by the lemma, for each  $k$  then  $\mathbf{h}'^{(k)} = \sum_{j=1}^K \mathcal{H}_{(k)}^{(j)} \boldsymbol{\gamma}^{(k,j)}$  for some  $(L_k - L_j + 1)_+$

length vector. In the alternative model of (C.5),  $\mathbf{H}'^{(k)}(n^*) = \sum_{j=1}^K \boldsymbol{\Gamma}^{(k,j)} \mathbf{H}^{(j)}(n^*)$

where

$$\mathbf{\Gamma}^{(k,j)} = \begin{bmatrix} \gamma_{(L_k-L_j+1)+}^{(k,j)} & & & & \\ \vdots & \ddots & & & \\ \gamma_2^{(k,j)} & & \gamma_{(L_k-L_j+1)+}^{(k,j)} & & \\ \gamma_1^{(k,j)} & & \vdots & \gamma_{(L_k-L_j+1)+}^{(k,j)} & \\ & \ddots & \gamma_2^{(k,j)} & \vdots & \\ & & \gamma_1^{(k,j)} & \gamma_2^{(k,j)} & \\ & & 0 & \gamma_1^{(k,j)} & \end{bmatrix}_{L_k+1+n^* \times L_j+1+n^*}$$

and this defines  $\mathbf{T}$  as

$$\mathbf{T} = \begin{bmatrix} \mathbf{\Gamma}^{(1,1)} & \mathbf{\Gamma}^{(1,2)} & \dots & \mathbf{\Gamma}^{(1,K)} \\ \mathbf{\Gamma}^{(2,1)} & \mathbf{\Gamma}^{(2,2)} & \dots & \mathbf{\Gamma}^{(2,K)} \\ \vdots & \vdots & \dots & \vdots \\ \mathbf{\Gamma}^{(K,1)} & \mathbf{\Gamma}^{(K,2)} & \dots & \mathbf{\Gamma}^{(K,K)} \end{bmatrix}.$$

Then since

$$(\mathbf{S}(n^*)\mathbf{T} + \mathbf{S}'(n^*))\mathbf{H}(n^*) = [\mathbf{S}(n^*) \quad \mathbf{S}'(n^*)] \begin{bmatrix} \mathbf{TH}(n^*) \\ \mathbf{H}(n^*) \end{bmatrix} = \mathbf{0}$$

and  $\mathbf{H}(n^*)$  is full row rank, it can be seen that  $\mathbf{S}'(n^*) = -\mathbf{S}(n^*)\mathbf{T}$  or  $\mathbf{S}'^{(k)}(n^*) = \sum_{j=1}^K -\mathbf{S}^{(j)}(n^*)\mathbf{\Gamma}^{(j,k)}$ , which implies in the vector-matrix model in (4.3) that  $\mathbf{s}'^{(k)} = \sum_{j=1}^K -\mathbf{s}_{(k)}^{(j)}\gamma^{(j,k)}$ . Let

$$\boldsymbol{\gamma} = \begin{bmatrix} \gamma^{(1,1)} \\ \gamma^{(1,2)} \\ \gamma^{(2,2)} \\ \gamma^{(2,1)} \\ \vdots \\ \gamma^{(1,K)} \\ \vdots \\ \gamma^{(K,K)} \\ \vdots \\ \gamma^{(K,1)} \end{bmatrix}.$$

Then  $\mathbf{v} = \mathcal{N}\boldsymbol{\gamma}$  and  $\text{span}(\mathcal{N}) \equiv \text{null}(\mathcal{I}(\boldsymbol{\vartheta}))$ . □

## Bibliography

- [1] K. Abed-Meraim, Y. Hua, "Strict identifiability of multichannel FIR systems: further results and developments," *Proceedings of the International Conference on Telecommunications*, Melbourne, Australia, pp. 1029-1032, April 1997.
- [2] John Aitchison, "Large Sample Restricted Parametric Tests," *Journal of the Royal Statistical Society, Series B (Methodological)*, Volume 24, Number 1, pp. 234-250, 1962.
- [3] J. Aitchison, S.D. Silvey, "Maximum likelihood estimation of parameters subject to restraints," *Ann. Math. Statist.*, Volume 29, pp. 813-828, 1958.
- [4] J. Aitchison, S.D. Silvey, "Maximum-likelihood estimation procedures and associated tests of significance," *Journal of the Royal Statistical Society, Series B (Methodological)*, Volume 22, No. 1, pp. 154-171, 1960.
- [5] Joshua N. Ash, "On singular estimation problems in sensor localization systems," Ph.D. Thesis, Ohio State University, 2007.
- [6] Joshua N. Ash, Randolph L. Moses, "On the relative error and absolute positioning errors in self-localization systems," *IEEE Transactions on Signal Processing*, Volume 56, Number 11, pp. 5668-5679, November 2008.
- [7] W.J. Bangs, *Array Processing with Generalized Beamformers*, Ph.D. Thesis, Yale University, New Haven, CT, 1971.
- [8] E. W. Barankin, "Locally best unbiased estimates," *Ann. Math. Statist.*, Volume 20, Number 4, pp. 477-501, 1949.
- [9] Zvika Ben-Haim, Yonina Eldar, "On the constrained Cramér-Rao bound with a singular Fisher information matrix," *IEEE Signal Processing Letters*, Volume 16, Number 6, pp. 453-456, June 2009.
- [10] Dimitri Bertsekas, *Nonlinear Programming: 2nd Edition*, Athena Scientific, 1999.
- [11] A. Bhattacharyya, "On some analogues of the amount of information and their use in statistical estimation," *Sankhya*, Volume 8, pp. 1-14, 1946.
- [12] A.A. Borovkov, A. Moullagaliev, *Mathematical Statistics*. Boca-Raton, FL: CRC, 1998.

- [13] Roger Bowden, "The theory of parametric identification," *Econometrica*, Volume 41, Number 6, pp. 1069-1074, November 1973.
- [14] George Casella, Roger L. Berger, *Statistical Inference, 2nd Edition*, Duxbury Press, 2002.
- [15] F. Cayre, C. Fontaine, T. Furon, "Watermarking security: theory and practice," *IEEE Transactions on Signal Processing*, Volume 53, Number 10, pp. 3976-3987, October 2005.
- [16] D. G. Chapman, H. Robbins, "Minimum variance estimation without regularity assumptions," *Ann. Math. Stat.*, Volume 22, pp. 581-586, 1951.
- [17] Harald Cramér, *Mathematical Methods of Statistics*, Princeton NJ: Princeton University Press, 1946.
- [18] Martin Crowder, "On constrained maximum likelihood estimation with non-i.i.d. observations," *Ann. Inst. Statist. Math.*, Volume 36, Part A, pp.239-249, 1984.
- [19] J.N. Franklin, *Matrix Theory*. New York: Dover, 1993.
- [20] James E. Gentle, *Matrix Algebra: Theory, Computations, and Applications in Statistics*. New York: Springer, 2007.
- [21] P.E. Gill, W. Murray, and M.H. Wright, *Practical Optimization*. New York: Academic, 1981.
- [22] A.A. Goldstein, "Convex programming in Hilbert space," *Bulletin of the American Mathematics Society*, Volume 70, Number 5, pp. 709-710, 1964.
- [23] John D. Gorman, Alfred O. Hero, "Lower bounds for parametric estimation with constraints," *IEEE Transactions on Information Theory*, Volume 26, Number 6, pp. 1285-1301, November 1990.
- [24] A. Gorokhov, P. Loubaton, "Subspace based techniques for blind separation of convolutive mixtures with temporally correlated sources," *IEEE Transactions on Circuits and Systems*, Volume 44, pp. 813-820, September 1997.
- [25] R.T. Haftka, Z. Gürdal, *Elements of Structural Optimization*. Norwell, MA: Kluwer Academic, 1992.

- [26] J. M. Hammersley, "On estimating restricted parameters," *Journal of the Royal Statistical Society, Series B*, Volume 12, Number 2, pp. 192-240, 1950.
- [27] Harrie Hendriks, "A Cramér-Rao type lower bound for estimators with values in a manifold," *Journal of Multivariate Analysis*, Volume 38, pp. 245-261, 1991.
- [28] A.O. Hero, R. Piramuthu, J.A. Fessler, S.R. Titus, "Minimax emission computed tomography using high-resolution anatomical side information and B-spline models," *IEEE Transactions on Information Theory*, Volume 45, Number 3, pp. 920-938, April 1999.
- [29] Bertrand Hochwald, Arye Nehorai, "On identifiability and information-regularity in parameterized normal distributions," *Circuits Systems Signal Processing*, Volume 16, Number 1, pp. 83-89, 1997.
- [30] Yingbo Hua, "Fast maximum-likelihood for blind identification of multiple FIR channels," *IEEE Transactions on Signal Processing*, Volume 44, Number 3, pp. 661-672, March 1996.
- [31] Y. Hua, M. Wax, "Strict identifiability of multiple FIR channels driven by an unknown arbitrary sequence," *IEEE Transactions on Signal Processing*, Volume 44, pp. 756-759, March 1996.
- [32] Aditya K. Jagannatham, Bhaskar D. Rao, "Cramér-Rao lower bound for constrained complex parameters," *IEEE Signal Processing Letters*, Volume 11, Number 11, pp. 875-878, November 2004.
- [33] Mortaza Jamshidian, "On algorithms for restricted maximum likelihood estimation," *Computational Statistics & Data Analysis*, Volume 45, pp. 137-157, 2004.
- [34] Jürgen Jost, *Riemannian Geometry and Geometric Analysis*, Third Edition, Springer-Verlag, 2002.
- [35] T. Kailath, *Linear Systems*. Englewood Cliffs, NH: Prentice-Hall, 1980.
- [36] Steven M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*. Englewood Cliffs, NJ: Prentice Hall, 1993.
- [37] C.G. Khatri, "A note on a MANOVA model applied to problems in growth curve," *Annals of the Institute of Statistical Mathematics*, Volume 18, pp. 75-86, 1966.

- [38] J.R. Kirkwood, *An Introduction to Analysis*. Boston, MA: PWS Kent, 1995.
- [39] Richard J. Kozick, Brian M. Sadler, "Maximum likelihood array processing: The semi-blind case," *Proceedings of the 2nd IEEE Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, pp.70-73, 1999.
- [40] Amir Leshem, Alle-Jan van der Veen, "Direction-of-arrival estimation for constant modulus signals," *IEEE Transactions on Signal Processing*, Volume 47, Number 11, pp. 3125-2319, November 1999.
- [41] Jian Li, R.T. Compton, Jr., "Maximum likelihood angle estimation for signals with known waveforms," *IEEE Transactions on Signal Processing*, Volume 41, Number 9, pp. 2850-2863, September 1993.
- [42] Z. Liu, G.B. Giannakis, S. Barbarossa, A. Scaglione, "Transmit antennae space-time block coding for generalized OFDM in the presence of unknown multipath," *IEEE Journal on Selected Areas in Communications*, Volume 19, Number 7, pp. 1352-1364, July 2001.
- [43] H. Liu, G. Xu, L. Tong, "A deterministic approach to blind channel identification of multi-channel FIR systems," *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP'94)*, Adelaide, Australia, Volume 4, pp. 581-584, April 1994.
- [44] P. Loubaton, E. Moulines, "On blind multiuser forward link channel estimation by the subspace method: Identifiability results," *IEEE Transactions on Signal Processing*, Volume 48, pp. 2366-2376, August 2000.
- [45] D.G. Luenberger, *Introduction to Linear and Nonlinear Programming*. Reading, MA: Addison-Wesley, 1973.
- [46] J.H. Manton, "An improved least squares blind channel identification algorithm for linear and affinely precoded communication systems," *IEEE Signal Processing Letters*, Volume 9, Number 9, pp. 282-285, September 2002.
- [47] Thomas L. Marzetta, "A simple derivation of the constrained multiple parameter Cramér-Rao bound," *IEEE Transactions on Signal Processing*, Volume 41, Number 6, pp. 2247-2249, June 1993.
- [48] Terrence J. Moore, Brian M. Sadler, Richard J. Kozick, "Regularity and strict identifiability in MIMO systems," *IEEE Transactions on Signal Processing*, Volume 50, Number 8, pp. 1831-1842, August 2002.



- [49] Terrence J. Moore, Brian M. Sadler, "Sufficient conditions for regularity and strict identifiability in MIMO systems," *IEEE Transactions on Signal Processing*, Volume 52, Number 9, pp. 2650-2655, September 2004.
- [50] Terrence J. Moore, Richard J. Kozick, Brian M. Sadler, "The constrained Cramér-Rao bound from the perspective of fitting a model," *IEEE Signal Processing Letters*, Volume 14, Number 8, pp. 564-567, August 2007.
- [51] Terrence J. Moore, Brian M. Sadler, Richard J. Kozick, "Maximum-likelihood estimation, the Cramér-Rao bound, and the method of scoring with parameter constraints," *IEEE Transactions on Signal Processing*, Volume 56, Number 3, pp. 895-908, March 2008.
- [52] E. Moulines, P. Duhamel, J. Cardoso, S. Mayrargue, "Subspace methods for the blind identification of multichannel FIR filters," *IEEE Transactions on Signal Processing*, Volume 43, pp. 516-525, February 1995.
- [53] J.M. Oller, J.M. Corcuera, "Intrinsic analysis of statistical estimation," *The Annals of Statistics*, Volume 23, Number 5, pp. 1562-1581, October 1995.
- [54] M.R. Osborne, "Fisher's method of scoring," *Int. Stat. Rev.*, Volume 60, pp. 99-117, 1992.
- [55] M.R. Osborne, "Scoring with constraints," *ANZIAM Journal*, Volume 42, Number 1, pp. 9-25, July 2000.
- [56] Callyampudi Radakrishna Rao, "Information and the accuracy attainable in the estimation of statistical parameters," *Bulletin of the Calcutta Mathematical Society*, Volume 37, pp. 81-89, 1945.
- [57] Alvin C. Rencher, *Linear Models In Statistics*, John Wiley & Sons, New York, 2000.
- [58] Thomas J. Rothenberg, "Identification in parametric models," *Econometrica*, Volume 39, Number 3, pp. 577-591, May 1971.
- [59] Brian M. Sadler, Richard J. Kozick, Terrence Moore, "Bounds on bearing and symbol estimation with side information," *IEEE Transactions on Signal Processing*, Volume 49, Number 4, pp. 822-834, April 2001.
- [60] L.L. Scharf, L.T. McWhorter, "Geometry of the Cramér-Rao bound," *Conference Proceedings, IEEE Sixth SP Workshop on Statistical Signal and Array Processing*, pp. 5-8, October 1992.

- [61] George A.F. Seber, Alan J. Lee, *Linear Regression Analysis: 2nd Edition*, Wiley-Interscience, 2003.
- [62] Jun Shao, *Mathematical Statistics*. New York, NY: Springer-Verlag, 2003.
- [63] S.D. Silvey, "The Lagrange multiplier test," *Ann. Math. Statist.*, Volume 30, pp. 389-407, 1959.
- [64] D.R. Smart, *Fixed Point Theorems*, Cambridge University Press, 1980.
- [65] M. Spivak, *Calculus on Manifolds*. Reading, MA: Addison-Wesley, 1965.
- [66] Petre Stoica, Thomas L. Marzetta, "Parameter estimation problems with singular information matrices," *IEEE Transactions on Signal Processing*, Volume 49, Number 1, pp. 87-90, January 2001.
- [67] Petre Stoica, Arye Nehorai, "MUSIC, maximum likelihood, and Cramér-Rao bound," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, Volume 37, Number 5, pp. 720-741, May 1989.
- [68] Petre Stoica, Boon Chong Ng, "On the Cramér-Rao bound under parametric constraints," *IEEE Signal Processing Letters*, Volume 5, Number 7, pp. 177-179, July 1998.
- [69] Petre Stoica, Torsten Söderström, "On non-singular information matrices and local identifiability," *International Journal on Control*, Volume 36, Number 2, pp. 323-329, 1982.
- [70] A. van den Bos, "A Cramér-Rao lower bound for complex parameters," *IEEE Transactions on Signal Processing*, Volume 42, Number 10, p. 2859, October 1994.
- [71] Alle-Jan van der Veen, "Asymptotic properties of the algebraic constant modulus algorithm," *IEEE Transactions on Signal Processing*, Volume 49, Number 8, pp. 1796-1807, August 2001.
- [72] Harry L. Van Trees, *Detection, Estimation, and Modulation Theory, Part I*, John Wiley and Sons, Inc. 1968.
- [73] L.J. Waldorp, H.M. Huizenga, R.P.P.P. Grasman, "The Wald test and Cramér-Rao bound for misspecified models in electromagnetic source analysis," *IEEE Transactions on Signal Processing*, Volume 53, Number 9, pp. 3427-3435, September 2005.

- [74] João Xavier, Victor Barroso, “Intrinsic distance lower bound for unbiased estimators on Riemannian manifolds,” *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP’02)*, Volume 2, pp. 1141-1144, May 2002.
- [75] João Xavier, Victor Barroso, “Intrinsic variance lower bound (IVLB): An extension of the Cramér-Rao bound to Riemannian manifolds,” *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP’05)*, Volume 5, pp. 1033-1036, March 2005.
- [76] Guanghan Xu, Hui Liu, Lang Tong, Thomas Kailath, “A least-squares approach to blind channel identification,” *IEEE Transactions on Signal Processing*, Volume 43, Number 12, pp. 2982-2993, December 1995.
- [77] Yingwei Yao, Georgios Giannakis, “On regularity and identifiability of blind source separation under constant-modulus constraints,” *IEEE Transactions on Signal Processing*, Volume 53, Number 4, pp. 1272-1281, April 2005.
- [78] J. Ziv, M. Zakai, “Some lower bounds on signal parameter estimation,” *IEEE Transactions on Information Theory*, Volume IT-15, pp. 386-391, May 1969.