

A Bayesian Framework for Noise Covariance Estimation Using the Facet Model

Desikachari Nadadur, *Member, IEEE*, Robert Martin Haralick, *Fellow, IEEE*, and David Earl Gustafson

Abstract—In image processing literature, thus far, researchers have assumed the perturbation in the data to be *white* (or uncorrelated) having a covariance matrix $\sigma^2\mathbf{I}$, i.e., assumption of equal variance for all the data samples and that no correlation exists between the data samples. However, there have been very few attempts to estimate noise characteristics under the assumption that there is a correlation between data samples. In this work, we propose a new and a novel approach for the simultaneous Bayesian estimation of the *unknown* colored or correlated noise (population) covariance matrix and the hyperparameters of the covariance model using the well-known *facet model*. We also estimate the facet model coefficients. We use the facet model because of its simple, yet elegant, mathematical formulation. We use the *generalized inverted Wishart density* as the prior model for the noise covariance matrix. We place a structure on the covariance matrix using the parameters of a *correlation filter*. These hyperparameters are estimated by a *new* extension of the expectation-maximization algorithm called the *generalized constrained expectation maximization algorithm* that we developed.

Index Terms—Colored noise, constraints, correlation filter, expectation-maximization (EM) algorithm, generalized constrained expectation maximization (GCEM) algorithm, generalized inverted Wishart (GIW), hypercovariance, hyperparameters, inverted Wishart (IW), noise covariance matrix, nonlinear programming, white noise.

I. INTRODUCTION

THE field of *image processing* and *computer vision*, in layman's terms, is a science of making the computers to perceive and understand the real world scenes. Commonly, one of the first and foremost steps involved in any *computer vision* task is the preprocessing of the digital images that are created by a sensor [for example, a camera, a medical ultrasound (US) scanning system, and a magnetic resonance imaging (MRI) system]. This involves noise cleaning and enhancement. Images, so enhanced, are passed through the subsequent steps that are to be carried out to realize a computer vision task. A typical computer vision system is illustrated in Fig. 1. An *ellipse* represents an algorithm, procedure or method acting

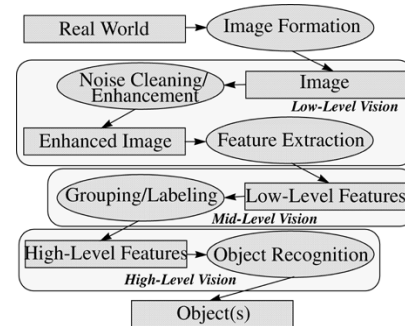


Fig. 1. Typical computer vision system.

upon a data object represented by a *rectangle*. A digital image formed by sensing the real world is preprocessed to clean up the noise and enhance the image content for subsequent processing. The enhanced image might then be processed to extract features that are of interest to us in a particular application. The features of interest could be gradient magnitude and direction, edge, ridge, and/or valley pixels, or some other information, such as histogram statistics, to name a few. This step is usually referred to as *feature extraction*. These two steps collectively are called the *low-level computer vision task*. Once these low-level features are detected, the next step could be to group them together and label them into more reasonable entities. This process could involve, for example, connecting detected edge pixels to form longer edge segments, by using some criteria. This step is referred to as *grouping and labeling* or *transformation to higher level entities*. This step is sometimes called *mid-level computer vision task*. Finally, these features are input to the *high-level computer vision task* of *object recognition* to detect object(s) of interest. As one can clearly see, that output of one step forms the input to the next. Consequently, any errors occurring in one step get propagated to the subsequent step and the output of the high-level vision task may not produce desired results. In recent years, increased attention has been paid to the development of algorithms for performing mid- and high-level vision tasks. Researchers are concentrating less on the low-level task of *noise estimation and enhancement*. For example, recent works on boundary detection and three dimensional reconstruction of the organs of the body using deformable models and templates [1]–[6] concentrated on finding the object of interest without heed to the accuracy of the results under the chosen noise model. They perform Gaussian smoothing of images assuming white noise with an assumed variance (covariance matrix $\sigma^2\mathbf{I}$). That is, an assumption of equal variance for all the data samples with no correlation among them. They did not even estimate

Manuscript received May 6, 2002; revised September 21, 2004. This work was supported by Siemens Medical Solutions USA, Inc., Ultrasound Division. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Patrick Boutheymy.

D. Nadadur is with the Advanced Imaging Applications, Developing Competency Department, Siemens Medical Solutions USA, Inc., Ultrasound Division, Issaquah, WA 98029 USA (e-mail: desikachari.nadadur@siemens.com; dnadadur@comcast.net).

R. M. Haralick is with the Department of Computer Science, Graduate Center, City University of New York, NY 10016 USA (e-mail: haralick@gc.cuny.edu).

D. E. Gustafson is with the Developing Competency Department, Siemens Medical Solutions USA, Inc., Ultrasound Division, Issaquah, WA 98029 USA (e-mail: david.gustafson@siemens.com).

Digital Object Identifier 10.1109/TIP.2005.854480

this variance from data. Others went on to estimate the noise variance σ^2 by techniques such as least squares fitting of the observed image data. Use of white noise assumption does not reflect the true correlation between samples that are observed in the real world. Therefore, this modeling error carries on to the higher level computer vision tasks such as feature extraction, perceptual grouping, and final scene reconstruction, thereby producing suboptimal results. We believe that, by properly modeling the correlation of the samples, it is possible to reduce the modeling errors, which then propagate through the various steps of computer vision algorithms to produce optimal results. Our belief is that a computer vision system without a proper modeling of the noise and the estimation of parameters of these models at each stage will not produce optimal and stable results.

This paper concentrates on the above mentioned *low-level vision* tasks of *noise covariance estimation* and *feature extraction*.

A. Background and Motivation

As stated earlier, researchers have focused very little energy, in this very important field of low-level image processing to estimate noise characteristics under the assumption that there is correlation between data samples. Most of the recent works in the covariance matrix estimation could only be found in statistical literature, for example, non-Bayesian approaches by Stein [7], [8], Lin and Perlman [9], and empirical Bayes procedure by Haff [10], fully Bayesian approach by Dickey *et al.* [11], the path models of Wright [12], the LISREL models of Jöreskog [13] and the factor analysis models of Spearman [14], mixture models of Hoffbeck [15], Toeplitz models of Cadre [16], matrix logarithm models of Leonard and Hsu [17], and hierarchical nonconjugate prior models of Daniels [18]. Brown *et al.* [19] propose a new class of priors for the covariance matrix, called the *generalized inverted Wishart (GIW)* prior, which is based on the Bartlett decomposition and apply it to the estimation of the covariance matrix in the environmental monitoring problem. We used this prior model as it naturally applies to the problem that we are trying to solve. Therefore, a lack of proper (general) framework for noise covariance matrix estimation in the field of image processing and computer vision, motivated us to undertake this research work.

B. Organization of the Paper

This paper is organized as follows. In Section II, we introduce the facet model and the additive Gaussian noise model that we will be using in the formulation of the problem of simultaneous estimation of the *noise covariance matrix* and the *facet model coefficients* in a Bayesian framework. We make use of the *GIW* density as the prior probability function for the unknown noise covariance matrix. We derive an objective function for the maximum *a posteriori* (MAP) estimation of the facet model coefficients and the noise covariance matrix. This section formally states the problem that needs to be solved. Section III discusses the correlation structure that we place on the hypercovariance matrix describing the *GIW* distribution. We design a *generalized constrained expectation maximization* (GCEM) algorithm for the estimation of the hyperparameters used to define this structure. In Section IV, we implement the *constrained*

maximization (M-step) of the GCEM algorithm using *sequential unconstrained minimization technique* (SUMT). We convert the constrained problem into a sequence of unconstrained problems using *barrier functions*. Section V discusses the evaluation protocol used to assess the performance of the algorithm and validate the results produced. We do this via the statistical hypothesis testing. Section VI discusses briefly an application that we developed to take advantage of the noise covariance matrix estimation procedure. In Section VII, we summarize the results reported in this paper.

II. PROBLEM FORMULATION

A. Introduction

In general, the solution to an image processing problem starts by specifying a model for the underlying true noise-free image, a noise or perturbation model that describes how the underlying true noise-free data might have been contaminated, and designing an algorithm to use these models to solve the problem at hand more accurately. By having such a statistical model, one can produce error estimates for the results produced so that the performance of the algorithm can be evaluated. In this section, we formulate the problem of noise covariance matrix estimation by using the *facet model* to describe the underlying noise free image and use an *additive* Gaussian noise model. The facet model was first introduced by Haralick [20]–[22]. This has been extensively used to extract features, such as edges, ridges, and corners of objects in images. See [23]–[30] for examples.

B. Facet Model

The principle of the facet model states that the image can be thought of as an underlying continuum or piecewise continuous gray level intensity surface. The observed digital image is a noisy, discretized sampling of a distorted version of this surface.

For the n th image neighborhood, we can write the facet model [22] representing the ideal noise free signal energy as

$$\mathbf{s}_n = \mathbf{B}\boldsymbol{\alpha}_n \quad (1)$$

where \mathbf{s}_n represents the $n = 1, \dots, N$ K -dimensional noiseless vectors from the signal space and \mathbf{B} is a matrix whose columns represent the discrete orthonormal polynomial (DOP) basis of the space which is modeled to contain the signal energy and $\boldsymbol{\alpha}_n$ are $n = 1, \dots, N$ M -dimensional vectors of coefficients of the facet model, and $K = (2R + 1) \times (2R + 1)$ where R is the half-width of the discrete support of the neighborhood. \mathbf{B} is obtained as discussed in [22], [31], and [32].

C. Noise Model

Let $\mathbf{x}_n; n = 1, \dots, N$ be the N K -dimensional independent samples of the noisy observed signal. Then, the facet model represents this noisy signal as

$$\mathbf{x}_n = \mathbf{B}\boldsymbol{\alpha}_n + \boldsymbol{\eta}_n \quad (2)$$

where $\boldsymbol{\eta}_n \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\eta)$ are the independent identically distributed Gaussian random variates. Arranging these N samples

together into columns of a matrix \mathbf{X} and using the notation given in [33] for matrix variate normal distribution, we can write

$$\mathbf{X} = \mathbf{B}\mathbf{\Lambda} + \boldsymbol{\eta} \quad (3)$$

where

$$\begin{aligned} \mathbf{X} &\in \mathbb{R}^{K \times N} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) \\ \mathbf{\Lambda} &\in \mathbb{R}^{M \times N} = (\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \dots, \boldsymbol{\alpha}_N) \\ \boldsymbol{\eta} &\in \mathbb{R}^{K \times N} = (\boldsymbol{\eta}_1, \boldsymbol{\eta}_2, \dots, \boldsymbol{\eta}_N). \end{aligned} \quad (4)$$

In this notation, we write the matrix normal distribution as $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\eta \otimes \mathbf{I}_N)$, where \otimes is the *Kronecker product* of matrices. Hence, $\mathbf{X} \sim \mathcal{N}(\mathbf{B}\mathbf{\Lambda}, \boldsymbol{\Sigma}_\eta \otimes \mathbf{I}_N)$. This distribution denotes that the columns of \mathbf{X} are independent and that the covariance matrix within each column is $\boldsymbol{\Sigma}_\eta \in \mathbb{R}^{K \times K}$. For some details on matrix normal distribution, see Appendix I [32]. Our goal is to estimate the polynomial coefficients $\boldsymbol{\alpha}_n$ and the common covariance matrix $\boldsymbol{\Sigma}_\eta$.

D. Estimation Problem

Given \mathbf{X} , we wish to estimate $\mathbf{\Lambda}$ and the common covariance matrix $\boldsymbol{\Sigma}_\eta$. We formulate the problem solution in a Bayesian framework to estimate $\hat{\mathbf{\Lambda}}$ of $\mathbf{\Lambda}$ and $\hat{\boldsymbol{\Sigma}}_\eta$ of $\boldsymbol{\Sigma}_\eta$ to *maximize* the posterior probability function

$$p(\mathbf{\Lambda}, \boldsymbol{\Sigma}_\eta | \mathbf{X}). \quad (5)$$

The values of $\mathbf{\Lambda}$ and $\boldsymbol{\Sigma}_\eta$ that maximize $p(\mathbf{\Lambda}, \boldsymbol{\Sigma}_\eta | \mathbf{X})$ will also maximize the joint probability

$$p(\mathbf{\Lambda}, \boldsymbol{\Sigma}_\eta, \mathbf{X}) = p(\mathbf{X} | \mathbf{\Lambda}, \boldsymbol{\Sigma}_\eta) p(\mathbf{\Lambda}, \boldsymbol{\Sigma}_\eta). \quad (6)$$

We model $\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_N, \boldsymbol{\Sigma}_\eta$ as being independent *a priori* and we take the prior probability for each $\boldsymbol{\alpha}_n$ to be uniform and we denote the prior probability for $\boldsymbol{\Sigma}_\eta$ as $p(\boldsymbol{\Sigma}_\eta)$ (see Section II-H for a discussion of the prior probability function used). For detailed derivations of all the expressions in this section, consult [32].

Under these assumptions, we have

$$p(\mathbf{\Lambda}, \boldsymbol{\Sigma}_\eta, \mathbf{X}) = p(\mathbf{X} | \mathbf{\Lambda}, \boldsymbol{\Sigma}_\eta) p(\mathbf{\Lambda}, \boldsymbol{\Sigma}_\eta). \quad (7)$$

We assume the probability function for the likelihood function as

$$p(\mathbf{X} | \mathbf{\Lambda}, \boldsymbol{\Sigma}_\eta) = \frac{\exp \left\{ -\frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \mathbf{B}\boldsymbol{\alpha}_i)' \boldsymbol{\Sigma}_\eta^{-1} (\mathbf{x}_i - \mathbf{B}\boldsymbol{\alpha}_i) \right\}}{(2\pi)^{\frac{NK}{2}} \det \boldsymbol{\Sigma}_\eta^{\frac{N}{2}}}. \quad (8)$$

Using the stated assumptions and likelihood function and neglecting the constant terms, we write the maximization problem as a minimization of the negative logarithm of the joint probability function as

$$\begin{aligned} \epsilon^2(\mathbf{\Lambda}, \boldsymbol{\Sigma}_\eta) &= -\frac{2}{N} \log p(\boldsymbol{\Sigma}_\eta) + \log \det \boldsymbol{\Sigma}_\eta \\ &+ \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mathbf{B}\boldsymbol{\alpha}_n)' \boldsymbol{\Sigma}_\eta^{-1} (\mathbf{x}_n - \mathbf{B}\boldsymbol{\alpha}_n). \end{aligned} \quad (9)$$

1) *Estimation of the Facet Model Coefficients:* This section describes the minimization of ϵ^2 with respect to $\boldsymbol{\alpha}_k$. Taking partial derivatives of ϵ^2 with respect to $\boldsymbol{\alpha}_k$, equating it to zero and solving the resulting equation, we get the estimate $\hat{\boldsymbol{\alpha}}_k$ of $\boldsymbol{\alpha}_k$ as

$$\hat{\boldsymbol{\alpha}}_k = \left(\mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1} \mathbf{B} \right)^{-1} \mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1} \mathbf{x}_k. \quad (10)$$

2) *Estimation of the Common Noise Covariance Matrix:* In this section, we carry out the minimization of ϵ^2 with respect to $\boldsymbol{\Sigma}_\eta$. Eliminating $\boldsymbol{\alpha}_k$ using $\hat{\boldsymbol{\alpha}}_k$ from the objective function ϵ^2 , we get

$$\begin{aligned} \epsilon^2(\boldsymbol{\Sigma}_\eta) &= -\frac{2}{N} \log p(\boldsymbol{\Sigma}_\eta) + \log \det \boldsymbol{\Sigma}_\eta \\ &+ \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n' \left(\boldsymbol{\Sigma}_\eta^{-1} - \boldsymbol{\Sigma}_\eta^{-1} \mathbf{B} \left(\mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1} \mathbf{B} \right)^{-1} \mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1} \right) \mathbf{x}_n. \end{aligned} \quad (11)$$

Notice that the operator $(\boldsymbol{\Sigma}_\eta^{-1} - \boldsymbol{\Sigma}_\eta^{-1} \mathbf{B} (\mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1} \mathbf{B})^{-1} \mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1})$ is one whose columns and rows are orthogonal to the columns of \mathbf{B} . That is, $(\boldsymbol{\Sigma}_\eta^{-1} - \boldsymbol{\Sigma}_\eta^{-1} \mathbf{B} (\mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1} \mathbf{B})^{-1} \mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1}) \mathbf{B} = 0$ and $\mathbf{B}' (\boldsymbol{\Sigma}_\eta^{-1} - \boldsymbol{\Sigma}_\eta^{-1} \mathbf{B} (\mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1} \mathbf{B})^{-1} \mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1}) = 0$.

This means that the only part of \mathbf{x}_n that makes a difference in the quadratic term is that part that is in the complement space to the column space of \mathbf{B} . This part can be written as $(\mathbf{I} - \mathbf{B}\mathbf{B}')\mathbf{x}_n$, since $\mathbf{B}\mathbf{B}'$ is the orthogonal projection operator to the space spanned by the columns of \mathbf{B} . Let $\mathbf{y}_n = (\mathbf{I} - \mathbf{B}\mathbf{B}')\mathbf{x}_n$, $\boldsymbol{\Sigma}_{yy} = (1/N) \sum_{n=1}^N \mathbf{y}_n \mathbf{y}_n'$ and removing terms independent of the summation index, n out of the summation, we get

$$\begin{aligned} \epsilon^2(\boldsymbol{\Sigma}_\eta) &= -\frac{2}{N} \log p(\boldsymbol{\Sigma}_\eta) + \log \det \boldsymbol{\Sigma}_\eta \\ &+ \text{tr} \left(\left(\boldsymbol{\Sigma}_\eta^{-1} - \boldsymbol{\Sigma}_\eta^{-1} \mathbf{B} \left(\mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1} \mathbf{B} \right)^{-1} \mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1} \right) \boldsymbol{\Sigma}_{yy} \right). \end{aligned} \quad (12)$$

E. When is This Approach an Improvement Over the White Noise Assumption?

Theorem 1 (Case of Uncorrelated Signal and Noise Spaces): When the noise covariance matrix, expressed in the coordinates of the signal space and its complement space is block diagonal with no coupling (i.e., no correlation) between the two spaces, then the same result is obtained for $\boldsymbol{\alpha}_k$ as in the white noise case.

Proof: Under the white noise assumption, we have $\boldsymbol{\Sigma}_\eta = \sigma^2 \mathbf{I}$. Therefore, $\boldsymbol{\Sigma}_\eta^{-1} = (1/\sigma^2) \mathbf{I}$. Now, the DOP basis coefficients $\boldsymbol{\alpha}_k$ become

$$\hat{\boldsymbol{\alpha}}_k = \left(\mathbf{B}' \frac{1}{\sigma^2} \mathbf{I} \mathbf{B} \right)^{-1} \mathbf{B}' \frac{1}{\sigma^2} \mathbf{I} \mathbf{x}_k = \mathbf{B}' \mathbf{x}_k. \quad (13)$$

Let \mathbf{C} be a matrix whose columns form an orthonormal basis for the space complement to that spanned by the columns of \mathbf{B} . In this case, the matrix $(\mathbf{B} \ \mathbf{C})$ is an orthonormal basis with inverse $(\mathbf{B} \ \mathbf{C})'$. Now

$$\boldsymbol{\Sigma}_\eta = (\mathbf{B} \ \mathbf{C})(\mathbf{B} \ \mathbf{C})' \boldsymbol{\Sigma}_\eta (\mathbf{B} \ \mathbf{C})(\mathbf{B} \ \mathbf{C})'. \quad (14)$$

The matrix $(\mathbf{B} \ \mathbf{C})' \boldsymbol{\Sigma}_\eta (\mathbf{B} \ \mathbf{C})$ is the noise covariance matrix expressed in the coordinate system of the basis formed by the columns of $(\mathbf{B} \ \mathbf{C})$.

Let us examine the case when this matrix is block diagonal, i.e.,

$$(\mathbf{B} \ \mathbf{C})' \boldsymbol{\Sigma}_\eta (\mathbf{B} \ \mathbf{C}) = \begin{pmatrix} \boldsymbol{\Sigma}_{\mathbf{B}\mathbf{B}} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{\mathbf{C}\mathbf{C}} \end{pmatrix}. \quad (15)$$

Using $\boldsymbol{\Sigma}_\eta^{-1}$ computed using the above equation, $\hat{\boldsymbol{\alpha}}_k$ becomes

$$\hat{\boldsymbol{\alpha}}_k = \mathbf{B}' \mathbf{x}_k. \quad (16)$$

We understand the above result in the following manner: Because of the relationship $\mathbf{x}_k = \mathbf{B}\boldsymbol{\alpha}_k + \boldsymbol{\eta}_k$, when the least squares estimate $\hat{\boldsymbol{\alpha}}_k$ of $\boldsymbol{\alpha}_k$ is produced, all the energy of the noise in the space spanned by the columns of \mathbf{B} goes into the computed $\hat{\boldsymbol{\alpha}}_k$. The noise in the space (spanned by the columns of \mathbf{C}) complement to the space spanned by the columns of \mathbf{B} will only influence the estimate $\hat{\boldsymbol{\alpha}}_k$ of $\boldsymbol{\alpha}_k$ to the extent that there is coupling between the signal space and the signal complement space (noise space). It is also clear that we only observe the signal from the space spanned by the columns of \mathbf{B} . We have no *direct* observations from the space spanned by the columns of \mathbf{C} , but only through the observations from the space spanned by the columns of \mathbf{B} .

F. General Case

In this section, we discuss the general case when there may be correlation between the space spanned by the columns of \mathbf{B} and that spanned by the columns of \mathbf{C} .

Theorem 2 (Case of Correlated Signal and Noise Spaces): Let $\boldsymbol{\Sigma}_{\mathbf{B}\mathbf{B}}$ be the component of the noise covariance matrix in the space spanned by the columns of \mathbf{B} , $\boldsymbol{\Sigma}_{\mathbf{C}\mathbf{C}}$ be the component of the noise covariance matrix in the space spanned by the columns of \mathbf{C} , and $\boldsymbol{\Sigma}_{\mathbf{B}\mathbf{C}}$ and $\boldsymbol{\Sigma}_{\mathbf{C}\mathbf{B}}$ be the cross-covariance matrices of the spaces spanned by \mathbf{B} and \mathbf{C} , such that $\boldsymbol{\Sigma}_{\mathbf{B}\mathbf{C}} = \boldsymbol{\Sigma}'_{\mathbf{C}\mathbf{B}}$. Then

$$\hat{\boldsymbol{\alpha}}_k = \left(\mathbf{B}' - \boldsymbol{\Sigma}_{\mathbf{B}\mathbf{C}} \boldsymbol{\Sigma}_{\mathbf{C}\mathbf{C}}^{-1} \mathbf{C}' \right) \mathbf{x}_k. \quad (17)$$

Proof: Consider the equation

$$\hat{\boldsymbol{\alpha}}_k = \left(\mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1} \mathbf{B} \right)^{-1} \mathbf{B}' \boldsymbol{\Sigma}_\eta^{-1} \mathbf{x}_k \quad (18)$$

where the noise covariance matrix $\boldsymbol{\Sigma}_\eta$ is represented in the basis of $(\mathbf{B} \ \mathbf{C})$. Let $\boldsymbol{\Sigma}_\eta$ represented in the basis of $(\mathbf{B} \ \mathbf{C})$ be partitioned as

$$\begin{aligned} \boldsymbol{\Sigma}_\eta^* \in \mathbb{R}^{K \times K} &= (\mathbf{B} \ \mathbf{C})' \boldsymbol{\Sigma}_\eta (\mathbf{B} \ \mathbf{C}) \\ &= \begin{pmatrix} \boldsymbol{\Sigma}_{\mathbf{B}\mathbf{B}} \in \mathbb{R}^{M \times M} & \boldsymbol{\Sigma}_{\mathbf{B}\mathbf{C}} \in \mathbb{R}^{M \times Q} \\ \boldsymbol{\Sigma}_{\mathbf{C}\mathbf{B}} \in \mathbb{R}^{Q \times M} & \boldsymbol{\Sigma}_{\mathbf{C}\mathbf{C}} \in \mathbb{R}^{Q \times Q} \end{pmatrix}. \end{aligned} \quad (19)$$

Using $\boldsymbol{\Sigma}_\eta^{-1}$ computed using the above partitioned structure, we get the estimate $\hat{\boldsymbol{\alpha}}_k$ of $\boldsymbol{\alpha}_k$ as

$$\hat{\boldsymbol{\alpha}}_k = \left(\mathbf{B}' - \boldsymbol{\Sigma}_{\mathbf{B}\mathbf{C}} \boldsymbol{\Sigma}_{\mathbf{C}\mathbf{C}}^{-1} \mathbf{C}' \right) \mathbf{x}_k. \quad (20)$$

Let

$$\tilde{\boldsymbol{\Sigma}}_{\mathbf{C}\mathbf{C}} = \frac{1}{N} \sum_{n=1}^N \mathbf{C}' \mathbf{x}_n \mathbf{x}_n' \mathbf{C}. \quad (21)$$

Using these results, the objective function in (9) becomes

$$\epsilon^2(\boldsymbol{\Sigma}_\eta) = -\frac{2}{N} \log p(\boldsymbol{\Sigma}_\eta) + \log \det \boldsymbol{\Sigma}_\eta + \text{tr} \left(\boldsymbol{\Sigma}_{\mathbf{C}\mathbf{C}}^{-1} \tilde{\boldsymbol{\Sigma}}_{\mathbf{C}\mathbf{C}} \right). \quad (22)$$

From the above equations, we can write $\log \det \boldsymbol{\Sigma}_\eta$ as

$$\log \det \boldsymbol{\Sigma}_\eta = \log \det \boldsymbol{\Sigma}_{\mathbf{C}\mathbf{C}} + \log \det \left(\boldsymbol{\Sigma}_{\mathbf{B}\mathbf{B}} - \boldsymbol{\Sigma}_{\mathbf{B}\mathbf{C}} \boldsymbol{\Sigma}_{\mathbf{C}\mathbf{C}}^{-1} \boldsymbol{\Sigma}_{\mathbf{C}\mathbf{B}} \right) \quad (23)$$

where $\det(\mathbf{B} \ \mathbf{C})(\mathbf{B} \ \mathbf{C})' = 1$. We can replace $p(\boldsymbol{\Sigma}_\eta)$ by $p(\boldsymbol{\Sigma}_\eta^*)$ since $\det(\mathbf{B} \ \mathbf{C})'(\mathbf{B} \ \mathbf{C})$, the determinant of the *Jacobian* of the transformation is unity. Therefore, the objective function can be written as

$$\begin{aligned} \epsilon^2(\boldsymbol{\Sigma}_\eta^*) &= -\frac{2}{N} \log p(\boldsymbol{\Sigma}_\eta^*) + \log \det \boldsymbol{\Sigma}_{\mathbf{C}\mathbf{C}} \\ &+ \log \det \left(\boldsymbol{\Sigma}_{\mathbf{B}\mathbf{B}} - \boldsymbol{\Sigma}_{\mathbf{B}\mathbf{C}} \boldsymbol{\Sigma}_{\mathbf{C}\mathbf{C}}^{-1} \boldsymbol{\Sigma}_{\mathbf{C}\mathbf{B}} \right) + \text{tr} \left(\boldsymbol{\Sigma}_{\mathbf{C}\mathbf{C}}^{-1} \tilde{\boldsymbol{\Sigma}}_{\mathbf{C}\mathbf{C}} \right). \end{aligned} \quad (24)$$

In Section II-H, we discuss the probability density function used to describe our *a priori* knowledge about $\boldsymbol{\Sigma}_\eta^*$.

G. Noise Model Revisited

The random vectors in the $(\mathbf{B} \ \mathbf{C})$ space can be written as

$$\begin{aligned} \mathbf{y}_n \in \mathbb{R}^{K \times 1} &= (\mathbf{B} \ \mathbf{C})' \mathbf{x}_n = \begin{pmatrix} \mathbf{y}_{\mathbf{B}_n} \in \mathbb{R}^{M \times 1} \\ \mathbf{y}_{\mathbf{C}_n} \in \mathbb{R}^{Q \times 1} \end{pmatrix} \\ &= \begin{pmatrix} \boldsymbol{\alpha}_n + \mathbf{B}' \boldsymbol{\eta}_n \\ \mathbf{0} + \mathbf{C}' \boldsymbol{\eta}_n \end{pmatrix}. \end{aligned} \quad (25)$$

In matrix form, data \mathbf{X} in the $(\mathbf{B} \ \mathbf{C})$ space be represented by \mathbf{Y} and given by

$$\mathbf{Y} \in \mathbb{R}^{K \times N} = (\mathbf{B} \ \mathbf{C})' \mathbf{X} = \begin{pmatrix} \mathbf{Y}_{\mathbf{B}} \in \mathbb{R}^{M \times N} \\ \mathbf{Y}_{\mathbf{C}} \in \mathbb{R}^{Q \times N} \end{pmatrix}. \quad (26)$$

Under this updated noise model, the estimate $\hat{\boldsymbol{\Lambda}}$ of $\boldsymbol{\Lambda}$ is

$$\hat{\boldsymbol{\Lambda}} = \mathbf{Y}_{\mathbf{B}} - \boldsymbol{\Sigma}_{\mathbf{B}\mathbf{C}} \boldsymbol{\Sigma}_{\mathbf{C}\mathbf{C}}^{-1} \mathbf{Y}_{\mathbf{C}}. \quad (27)$$

H. Choice of Prior Probability Function for the Noise Covariance Matrix

Brown *et al.* [19] proposed what they call the *GIW* distribution, which is based on the partitioned structure of the covariance matrix, and derived using the 1-1 Bartlett [34] matrix decomposition. For some details on the inverted Wishart (*IW*) and *GIW* distributions, see Appendices II and III.

Let $\boldsymbol{\Sigma}_\eta \sim \mathcal{IW}(\delta, \boldsymbol{\Psi}_\eta)$. Then, $\boldsymbol{\Sigma}_\eta^* \sim \mathcal{IW}(\delta, \boldsymbol{\Psi}_\eta^*)$. If $\boldsymbol{\Sigma}_\eta^* \sim \mathcal{IW}(\delta, \boldsymbol{\Psi}_\eta^*)$, then this exactly corresponds to the *GIW* distribution with parameter δ and partitioned $\boldsymbol{\Psi}_\eta^*$ as given below.

The data matrix \mathbf{Y} in (26) represents the partitioning of data conformable with that of Σ_{η}^* . Also, let the matrix Ψ_{η}^* be partitioned, conformably with Σ_{η}^* , as

$$\Psi_{\eta}^* = (\mathbf{B} \ \mathbf{C})' \Psi_{\eta} (\mathbf{B} \ \mathbf{C}) = \begin{pmatrix} \Psi_{BB} & \Psi_{BC} \\ \Psi_{CB} & \Psi_{CC} \end{pmatrix}. \quad (28)$$

The \mathcal{GTW} distribution for our problem is then described as

$$\Sigma_{\eta}^* \in \mathbb{R}^{K \times K} \sim \mathcal{GTW}(\delta, \Psi_{\eta}^*) \quad (29)$$

implies that

- 1) $\Sigma_{CC} \in \mathbb{R}^{Q \times Q} \sim \mathcal{IW}(\delta, \Psi_{CC})$, is distributed independently of $\boldsymbol{\tau}$ and ζ ;
- 2) $\zeta \in \mathbb{R}^{M \times M} \sim \mathcal{IW}(\delta + Q, \zeta_0)$;
- 3) $\boldsymbol{\tau} \in \mathbb{R}^{Q \times M} | \zeta \sim \mathcal{N}(\boldsymbol{\tau}_0, \Psi_{CC}^{-1} \otimes \zeta)$;

where $\boldsymbol{\tau} = \Sigma_{CC}^{-1} \Sigma_{CB}$; $\zeta = \Sigma_{BB} - \Sigma_{BC} \Sigma_{CC}^{-1} \Sigma_{CB} = \Sigma_{BB} - \boldsymbol{\tau}' \Sigma_{CC} \boldsymbol{\tau}$; $\boldsymbol{\tau}_0 = \Psi_{CC}^{-1} \Psi_{CB}$; $\zeta_0 = \Psi_{BB} - \Psi_{BC} \Psi_{CC}^{-1} \Psi_{CB}$; $K = M + Q$.

I. Updated Objective Function

In the following, we make the dependence of Σ_{η}^* on δ and Ψ_{η}^* explicit in notation. From the preceding analysis, we write $p(\Sigma_{\eta}^*)$ as

$$p(\Sigma_{\eta}^*) = p(\Sigma_{\eta}^* | \delta, \Psi_{\eta}^*) \\ = p(\Sigma_{CC} | \delta, \Psi_{CC}) p(\boldsymbol{\tau} | \boldsymbol{\tau}_0, \Psi_{CC}, \zeta) p(\zeta | \delta, \zeta_0). \quad (30)$$

Using the appropriate density functions for matrix normal and \mathcal{IW} distributions (given in Appendices I and II), we have the objective function to be minimized as

$$\begin{aligned} \epsilon^2(\Sigma_{CC}, \boldsymbol{\tau}, \zeta, \mathcal{H}) \\ = \frac{\delta + 2Q + N}{2} \log \det \Sigma_{CC} \\ + \frac{\delta + 2(Q + M) + N}{2} \log \det \zeta \\ + \frac{1}{2} \text{tr} \left([\Psi_{CC} + N \tilde{\Sigma}_{CC}] \Sigma_{CC}^{-1} \right) + \frac{1}{2} \text{tr}(\zeta_0 \zeta^{-1}) \\ + \frac{1}{2} \text{tr} \left([\Psi_{CC}(\boldsymbol{\tau} - \boldsymbol{\tau}_0)] [(\boldsymbol{\tau} - \boldsymbol{\tau}_0) \zeta^{-1}]' \right) \end{aligned} \quad (31)$$

where \mathcal{H} represents the hyperparameters to be elaborated later. Therefore, our goal is to estimate Σ_{CC} , ζ and $\boldsymbol{\tau}$, by minimizing the objective function given in (31) given the hyperparameters, \mathcal{H} . Once Σ_{CC} , ζ and $\boldsymbol{\tau}$ are estimated, we can then construct Σ_{η}^* and Ψ_{η}^* by performing an inverse of the Bartlett decomposition.

1) *Estimating Σ_{CC}* : To estimate Σ_{CC} , minimize the following function:

$$\begin{aligned} \epsilon^2(\Sigma_{CC}) = \frac{\delta + 2Q + N}{2} \log \det \Sigma_{CC} \\ + \frac{1}{2} \text{tr} \left([\Psi_{CC} + N \tilde{\Sigma}_{CC}] \Sigma_{CC}^{-1} \right). \end{aligned} \quad (32)$$

Using Theorem 6, we get the estimate as

$$\hat{\Sigma}_{CC} = \frac{1}{\delta + 2Q + N} (\Psi_{CC} + N \tilde{\Sigma}_{CC}). \quad (33)$$

2) *Estimating $\boldsymbol{\tau}$* : By minimizing the following function:

$$\epsilon^2(\boldsymbol{\tau}) = \frac{1}{2} \text{tr} \left([\Psi_{CC}(\boldsymbol{\tau} - \boldsymbol{\tau}_0)] [(\boldsymbol{\tau} - \boldsymbol{\tau}_0) \zeta^{-1}]' \right) \quad (34)$$

the estimate $\hat{\boldsymbol{\tau}}$ of $\boldsymbol{\tau}$ is obtained as

$$\hat{\boldsymbol{\tau}} = \boldsymbol{\tau}_0. \quad (35)$$

3) *Estimating ζ* : An estimate of ζ is obtained by minimizing the following function:

$$\begin{aligned} \epsilon^2(\zeta) = \frac{\delta + 2(Q + M) + N}{2} \log \det \zeta \\ + \frac{1}{2} \text{tr} \left([\zeta_0 + (\boldsymbol{\tau} - \boldsymbol{\tau}_0)' \Psi_{CC}(\boldsymbol{\tau} - \boldsymbol{\tau}_0)] \zeta^{-1} \right). \end{aligned} \quad (36)$$

Using Theorem 6, and $\boldsymbol{\tau} = \hat{\boldsymbol{\tau}}$, we get

$$\hat{\zeta} = \frac{1}{\delta + 2(Q + M) + N} \zeta_0. \quad (37)$$

The above expressions are derived assuming that the hyperparameters are known. However, we do not know these in advance, and, hence, we need to employ an estimation procedure. Also, the optimization becomes tricky if we do not know how the space (signal space) spanned by the columns of \mathbf{B} is related to the space (noise space) spanned by the columns of \mathbf{C} . This is because there exist several choices for the matrices Σ_{BB} , $\Sigma_{BC} = \Sigma_{CB}'$ and Σ_{CC} , so that one can obtain a covariance matrix that satisfies the positive-definiteness constraint. However, this may not be the optimal solution. Therefore, it is best to have a model for the correlation between the spaces spanned by the columns of \mathbf{B} and \mathbf{C} . This problem is best solved if we assume that the observations \mathbf{y}_n were obtained by passing *white noise* through a filter h defined by a functional form with free parameters. Then, the entries in the covariance matrix become a function of these free parameters. It is better to place this structure on the hypercovariance matrix, Ψ_{η}^* rather than on Σ_{η}^* . This is because if a structure is placed at the primary level, then any deviation from the assumed structure during the estimation procedure will result in serious errors [11], while placing the structure on a secondary level induces more flexibility and reduces the total number of parameters that need to be estimated. This is discussed in the next section along with the hyperparameter estimation.

III. CORRELATION STRUCTURE AND HYPERPARAMETER ESTIMATION

A. Introduction

In this section, we discuss an approach to estimating the *hypercovariance matrix* and other *hyperparameters* by specifying a correlation structure on the hypercovariance matrix Ψ_{η}^* (this is the hypercovariance matrix Ψ_{η} expressed in $(\mathbf{B} \ \mathbf{C})$ space). By looking at the noise model in (25), it is clear that the observations \mathbf{X} provide only samples from the signal space plus the noise proportional to the correlation between the spaces spanned by the columns of the matrices \mathbf{B} and \mathbf{C} . In other words, we do not have direct observations of \mathbf{Y}_C , but only through \mathbf{X} . Therefore, we need to, somehow, estimate the covariance matrix Σ_{CC}

in the space spanned by the columns of the matrix \mathbf{C} (noise space), from the observations \mathbf{X} . This estimation procedure becomes tricky if we do not have a model for the *correlation* between the space spanned by the columns of \mathbf{B} and that spanned by the columns of \mathbf{C} , as there exist several choices for the matrices $\Sigma_{\mathbf{B}\mathbf{B}}$, $\Sigma_{\mathbf{B}\mathbf{C}} = \Sigma'_{\mathbf{C}\mathbf{B}}$, and $\Sigma_{\mathbf{C}\mathbf{C}}$ so that one can satisfy the *positive definiteness constraint* on the covariance matrix, Σ_{η}^* . This may not be an optimal choice.

B. Correlation Structure

Our problem is best solved if we assume that the observations \mathbf{Y} were obtained by passing *white noise* through a filter h defined by a functional form with free parameters. Associated with any covariance matrix Σ_{η}^* is a filter h . If *white noise* is put through the filter h , the resulting colored noise signal will have covariance matrix directly computable in closed form depending on the values of the components of h . Now, suppose that we consider a class of parametric filters, where the number of parameters of the filter is substantially less than the dimension of η . In this case, Σ_{η}^* is a matrix each of whose entries is a function of the parameters of the filter h . As discussed earlier, any requisite structure is placed on the hypercovariance matrix Ψ_{η}^* , instead of on the matrix Σ_{η}^* , as it offers flexibility in estimating Σ_{η}^* , by allowing us to waver from the assumed structure secondary level when estimating the primary level variables.

To model the correlation, we assume our observations to be coming from a S th order *moving average process* (MA), indexed by t , and defined by

$$y(t) = \sum_{i=0}^S a_i \omega(t-i) \quad (38)$$

where $\omega(t)$ is the *white noise* process with $\mathbf{E}[\omega(t)] = 0$ and $\mathbf{V}[\omega(t)] = \sigma_{\omega}^2$.

In the current problem at hand, $t = 1, \dots, K$ and set $S = K - 1$ so that we make all the samples $y(t)$ correlated with one another. We use ρ to denote the correlation term and γ to denote the covariance term.

The condition on the covariances implies that the covariance (and, hence, the correlation) matrix of the vector (y_1, y_2, \dots, y_K) has the form

$$\Psi_{\eta}^* = \begin{pmatrix} \gamma_0 & \gamma_1 & \gamma_2 & \dots & \gamma_{(K-1)} \\ \gamma_1 & \gamma_0 & \gamma_1 & \dots & \gamma_{(K-2)} \\ \gamma_2 & \gamma_1 & \gamma_0 & \dots & \gamma_{(K-3)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \gamma_{(K-1)} & \gamma_{(K-2)} & \gamma_{(K-3)} & \dots & \gamma_0 \end{pmatrix} \quad (39)$$

wherein the generic element in the (i, j) th position is $\gamma_{|i-j|} = C(y_i, y_j)$, the covariance of y_i and y_j .

We have

$$\gamma_{\lambda} = \sigma_{\omega}^2 \sum_j a_j a_{j+\lambda} \quad (40)$$

where $\lambda = 0, \dots, K - 1$ and $\lambda + j \leq K - 1$.

The correlation coefficient ρ_{λ} is computed as

$$\rho_{\lambda} = \frac{\gamma_{\lambda}}{\gamma_0} \quad (41)$$

and is elaborated into

$$\rho_{\lambda} = \begin{cases} 1, & \lambda = 0 \\ \frac{\sum_j a_j a_{j+\lambda}}{\sum_i a_i^2}, & \lambda = 1, \dots, (K-1), \text{ and } \lambda + j \leq K-1. \\ 0, & \lambda > K-1 \end{cases} \quad (42)$$

In general, a process like this is normalized by setting either $a_0 = 1$ or $\sigma_{\omega}^2 = 1$. We make the choice of $\sigma_{\omega}^2 = 1$ and estimate the remaining parameters. We call $\mathcal{H} = (\delta, a_0, a_1, \dots, a_{K-1})$ as our set of *hyperparameters* that are to be estimated. Now, the problem is to estimate the hyperparameters \mathcal{H} given the data \mathbf{Y} under the constraints that

$$\left. \begin{aligned} \delta &> 0 \\ |\rho_{\lambda}| &\leq 1, \quad \text{for } \lambda = 1, \dots, K-1 \\ \Psi_{\eta}^* &\succ 0 \\ \Psi_{\mathbf{C}\mathbf{C}} &\succ 0 \\ \zeta_0 &\succ 0 \end{aligned} \right\}. \quad (43)$$

The first two equations give the element level constraints on the hypercovariance matrix and the rest of the equations provide the higher or matrix-level constraints.

In the following sections, we derive an expectation maximization algorithm for estimating the hyperparameters \mathcal{H} . In this process, we derive the much needed joint density of data $(\mathbf{Y}'_{\mathbf{B}}, \mathbf{Y}'_{\mathbf{C}})'$, and $\Sigma_{\mathbf{C}\mathbf{C}}$, τ , ζ , and the joint posterior density of $\Sigma_{\mathbf{C}\mathbf{C}}$, τ , ζ given data $(\mathbf{Y}'_{\mathbf{B}}, \mathbf{Y}'_{\mathbf{C}})'$. This also serves to illustrate that the objective function derived in earlier sections by starting with data \mathbf{X} and then transforming the objective function to the space spanned by the columns of $(\mathbf{B} \ \mathbf{C})$ is exactly the same as if we derived the objective function by first transforming the data \mathbf{X} to the space spanned by the columns of $(\mathbf{B} \ \mathbf{C})$ as \mathbf{Y} and then worked directly in this space. Directly working in the space spanned by the columns of $(\mathbf{B} \ \mathbf{C})$ leads to the result in a more straight forward manner.

C. Posterior Density of $\Sigma_{\mathbf{C}\mathbf{C}}$, τ and ζ

In this section, we derive an expression for the *joint posterior density* of $\Sigma_{\mathbf{C}\mathbf{C}}$, τ , and ζ given the data $\mathbf{Y}_{\mathbf{B}}$ and $\mathbf{Y}_{\mathbf{C}}$. In this treatment, we approximate Λ by its estimate $\hat{\Lambda}$ given by [rewritten from (27)]

$$\hat{\Lambda} = \mathbf{Y}_{\mathbf{B}} - \tau' \mathbf{Y}_{\mathbf{C}}. \quad (44)$$

Theorem 3 (Joint Posterior Density of $\Sigma_{\mathbf{C}\mathbf{C}}$, τ and ζ): Given

$$\mathbf{Y}_{\mathbf{B}} | \mathbf{Y}_{\mathbf{C}} \sim \mathcal{N}(\Lambda + \tau' \mathbf{Y}_{\mathbf{C}}, \zeta \otimes \mathbf{I}_N) \quad (\text{see Theorem 7}) \quad (45)$$

$$\mathbf{Y}_{\mathbf{C}} \sim \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{C}\mathbf{C}} \otimes \mathbf{I}_N) \quad (46)$$

$$\Sigma_{\mathbf{C}\mathbf{C}} \sim \mathcal{IW}(\delta, \Psi_{\mathbf{C}\mathbf{C}}) \quad (47)$$

$$\zeta \sim \mathcal{IW}(\delta + Q, \zeta_0) \quad (48)$$

$$\tau | \zeta \sim \mathcal{N}(\tau_0, \Psi_{\mathbf{C}\mathbf{C}}^{-1} \otimes \zeta) \quad (49)$$

Λ is assumed to be a *a priori* uniform and it is a *a priori* independent of $\Sigma_{\mathbf{C}\mathbf{C}}$, τ , ζ . (50)

Then, the joint posterior density of $(\Sigma_{\mathbf{C}\mathbf{C}}, \tau, \zeta)$ is given by

$$p(\Sigma_{\mathbf{C}\mathbf{C}}, \tau, \zeta | \Lambda = \hat{\Lambda}, \mathbf{Y}_{\mathbf{B}}, \mathbf{Y}_{\mathbf{C}}) = p(\Sigma_{\mathbf{C}\mathbf{C}} | \mathbf{Y}_{\mathbf{B}}, \mathbf{Y}_{\mathbf{C}}, \delta, \Psi_{\mathbf{C}\mathbf{C}}) \times p(\tau | \mathbf{Y}_{\mathbf{B}}, \mathbf{Y}_{\mathbf{C}}, \tau_0, \Psi_{\mathbf{C}\mathbf{C}}, \zeta) \times p(\zeta | \mathbf{Y}_{\mathbf{B}}, \mathbf{Y}_{\mathbf{C}}, \delta, \zeta_0) \quad (51)$$

where

$$\Sigma_{CC}|Y_B, Y_C, \delta, \Psi_{CC} \sim \mathcal{IW}(\delta + N, \Psi_{CC} + Y_C Y_C') \quad (52)$$

$$\tau|Y_B, Y_C, \tau_0, \Psi_{CC}, \zeta \sim \mathcal{N}(\tau_0, \Psi_{CC}^{-1} \otimes \zeta) \quad (53)$$

$$\zeta|Y_B, Y_C, \delta, \zeta_0 \sim \mathcal{IW}(\delta + Q + N, \zeta_0). \quad (54)$$

Proof: For the proof, see Appendix VI ■

D. Hyperparameter Estimation via the Expectation Maximization Algorithm

In the previous section, we carried out the analysis under the assumption that the hyperparameters \mathcal{H} of the $\mathcal{G}\mathcal{I}\mathcal{W}$ distribution were known. This works well if we know the values of \mathcal{H} *a priori*. However, in reality, we do not know these values. One must resort to some form of estimation procedure. One of the estimation procedures called the *EM* algorithm was proposed by Dempster [35] for this purpose. The EM algorithm works well when the probability density function for the incomplete data (the covariance matrix, in our case) cannot be computed easily. In that case, the conditional expectation of the logarithm of the probability density function of the incomplete data given the observables and the current point estimate of the parameters, is maximized to obtain the estimate of the parameters for the next iteration. Recall that we placed a structure on Ψ_{η}^* in Section III-B based on the filter coefficients a_i ; $i = 0, \dots, K-1$. This translates into Ψ_{CC} , τ_0 , and ζ_0 being functions of these parameters. In this section, we design an algorithm using EM theory for estimating the hyperparameters, $\mathcal{H} = \{\delta, a_0, a_1, \dots, a_{K-1}\}$ simultaneously, from the observed data, under the constraints imposed in (43). Additional examples of the usage of EM theory can be found in [36]–[38].

1) *Background of EM Theory:* Let \mathbf{y} be the vector of hidden parameters. The vector \mathbf{x} is called the *incomplete data*. The $\mathbf{z} = (\mathbf{x}, \mathbf{y})$ is called the *complete data*. Then, the joint density function or the *complete data density function* is

$$p(\mathbf{z}|\boldsymbol{\theta}) = p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}) = p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})p(\mathbf{x}|\boldsymbol{\theta}). \quad (55)$$

This defines the *complete data likelihood*. This function is a random variable since the missing information \mathbf{y} is unknown, random and presumably governed by an underlying distribution. Therefore, this function can be thought of as a function of \mathbf{y} given that $\boldsymbol{\theta}$ and \mathbf{x} are constants.

The EM algorithm then finds the expected value of the complete data log-likelihood $\log p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})$ with respect to the unknown data \mathbf{y} given the observed data \mathbf{x} and the current parameter estimate $\boldsymbol{\theta}^{(t)}$. Therefore, let us define

$$\begin{aligned} \Omega(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) &= \mathbf{E} \left[\log p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}) | \mathbf{x}, \boldsymbol{\theta}^{(t)} \right] \\ &= \int_{\mathbf{y}} \log p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}) p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}^{(t)}) d\mathbf{y} \end{aligned} \quad (56)$$

where $\boldsymbol{\theta}^{(t)}$ is the current parameter estimate that is used to evaluate the expectation and $\boldsymbol{\theta}$ are the parameters that we find to

maximize Ω . The function $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}^{(t)})$ is the marginal distribution of the unobserved data and is dependent on both the observations \mathbf{x} and the current parameter estimates. Convergence properties of EM algorithm are well understood [35], [39]. Algorithm III.1 describes the steps of the basic EM algorithm.

Algorithm III.1: Basic EM Algorithm

Once $\Omega(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$ is defined the EM algorithm has two steps, the *evaluation of expectation step* or the *E-step* and the *maximization of expectation step* or the *M-step*. These two steps of the algorithm are described below.

E-step Given the data \mathbf{x} and the parameter estimates $\boldsymbol{\theta}^{(t)}$ at iteration t , evaluate the expectation

$$\Omega(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) = \mathbf{E} \left[\log p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}) | \mathbf{x}, \boldsymbol{\theta}^{(t)} \right]. \quad (57)$$

Go to M-step.

M-step Given the expectation computed in E-step, optimize $\Omega(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$ to find $\boldsymbol{\theta}^{(t+1)}$ as

$$\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} \Omega(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}). \quad (58)$$

Set $t = t + 1$ and go to E-step.

The E-step and M-step are carried out to convergence.

There is another variation of the algorithm called the *generalized EM algorithm* (GEM) in which, instead of maximizing $\Omega(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$, we just find some $\boldsymbol{\theta}^{(t+1)}$ such that

$$\Omega(\boldsymbol{\theta}^{(t+1)}|\boldsymbol{\theta}^{(t)}) \geq \Omega(\boldsymbol{\theta}^{(t)}|\boldsymbol{\theta}^{(t)}) \quad (59)$$

in the M-step. This because the EM algorithm in its original form is guaranteed to increase the likelihood. We will use the GEM algorithm in the hyperparameter estimation to be discussed in the next section. Recall that the estimates of the hyperparameters must satisfy the constraints given in (43). Therefore, in the next section, we modify the GEM algorithm to incorporate these constraints during the M-step and we christen this approach GCEM algorithm.

2) *Hyperparameter Estimation Using the GCEM Algorithm:* In this section, we derive $\Omega(\mathcal{H}|\mathcal{H}^{(t)})$ for estimating the hyperparameters defined by $\mathcal{H} = \{\delta, a_0, a_1, \dots, a_{K-1}\}$. Since \mathcal{H} are the hyperparameters defining Σ_{η}^* partitioned as Σ_{CC} , τ , and ζ , we conveniently choose these as the hidden variables. We remove the facet model coefficients Λ by conditioning as $\Lambda = \hat{\Lambda}$. This is valid as we are assuming Λ to be *a priori* uniform, and, therefore, there are no hyperparameters of Λ to be estimated. The data matrix is $\mathbf{Y} = (\mathbf{Y}'_B, \mathbf{Y}'_C)'$.

Note that Ψ_{CC} , τ_0 and ζ_0 are functions of $(a_0, a_1, \dots, a_{K-1})$. The objective function $\Omega(\mathcal{H}|\mathcal{H}^{(t)})$,

the terms $\log p(\mathbf{Y}_B | \mathbf{Y}_C, \boldsymbol{\tau}, \boldsymbol{\zeta})$ and $\log p(\mathbf{Y}_C | \boldsymbol{\Sigma}_{CC})$ are independent of \mathcal{H} ; therefore, we can write it as

$$\begin{aligned} \Omega^* (\mathcal{H} | \mathcal{H}^{(t)}) &= \mathbf{E} \left[\log p(\boldsymbol{\Sigma}_{CC} | \mathcal{H}) | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] \\ &+ \mathbf{E} \left[\log p(\boldsymbol{\tau} | \boldsymbol{\zeta}, \mathcal{H}) | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] \\ &+ \mathbf{E} \left[\log p(\boldsymbol{\zeta} | \mathcal{H}) | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right]. \end{aligned} \quad (60)$$

The above equation shows that we take expectations of the logarithms of the prior probability density functions of $(\boldsymbol{\Sigma}_{CC}, \boldsymbol{\tau}, \boldsymbol{\zeta})$ with respect to the posterior density of $(\boldsymbol{\Sigma}_{CC}, \boldsymbol{\tau}, \boldsymbol{\zeta})$ (computed using Theorem 3), respectively. The expectations are computed using the properties of matrix normal, Wishart and \mathcal{IW} distributions given in Appendices I and II. The objective function to be maximized is

$$\begin{aligned} \Omega^* (\mathcal{H} | \mathcal{H}^{(t)}) &= -\frac{\delta(M+Q)}{2} \log 2 - \sum_{i=1}^Q \log \Gamma \left[\frac{\delta+Q-i}{2} \right] \\ &- \sum_{i=1}^M \log \Gamma \left[\frac{\delta+Q+N+M-i}{2} \right] \\ &+ \left(\frac{\delta+Q+M-1}{2} \right) \log \det \boldsymbol{\Psi}_{CC} \\ &- \left(\frac{\delta+2Q}{2} \right) \mathbf{E} \left[\log \det \boldsymbol{\Sigma}_{CC} | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] \\ &- \frac{1}{2} \text{tr} \left(\mathbf{E} \left[\boldsymbol{\Sigma}_{CC}^{-1} | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] \boldsymbol{\Psi}_{CC} \right) \\ &- \left(\frac{\delta+2(M+Q)+N}{2} \right) \mathbf{E} \left[\log \det \boldsymbol{\zeta} | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] \\ &- \frac{1}{2} \text{tr} \left(\mathbf{E} \left[\boldsymbol{\tau} \boldsymbol{\zeta}^{-1} \boldsymbol{\tau}' | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] \boldsymbol{\Psi}_{CC} \right) \\ &+ \frac{1}{2} \text{tr} \left(\mathbf{E} \left[\boldsymbol{\tau} \boldsymbol{\zeta}^{-1} | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] \boldsymbol{\tau}'_0 \boldsymbol{\Psi}_{CC} \right) \\ &+ \frac{1}{2} \text{tr} \left(\boldsymbol{\tau}_0 \mathbf{E} \left[\boldsymbol{\zeta}^{-1} \boldsymbol{\tau}' | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] \boldsymbol{\Psi}_{CC} \right) \\ &- \frac{1}{2} \text{tr} \left(\boldsymbol{\tau}_0 \mathbf{E} \left[\boldsymbol{\zeta}^{-1} | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] \boldsymbol{\tau}'_0 \boldsymbol{\Psi}_{CC} \right) \\ &+ \left(\frac{\delta+Q+N+M-1}{2} \right) \log \det \boldsymbol{\zeta}_0 \\ &- \frac{1}{2} \text{tr} \left(\mathbf{E} \left[\boldsymbol{\zeta}^{-1} | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] \right) \end{aligned} \quad (61)$$

where

$$\begin{aligned} \mathbf{E} \left[\log \det \boldsymbol{\Sigma}_{CC} | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] &= -Q \log 2 - \sum_{i=1}^Q \mathcal{D} \left(\frac{\delta^{(t)} + N + Q - i}{2} \right) \\ &+ \log \det \left(\boldsymbol{\Psi}_{CC}^{(t)} + \mathbf{Y}_C \mathbf{Y}_C' \right) \end{aligned} \quad (62)$$

$$\begin{aligned} \mathbf{E} \left[\boldsymbol{\Sigma}_{CC}^{-1} | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] &= \left(\delta^{(t)} + N + Q - 1 \right) \left(\boldsymbol{\Psi}_{CC}^{(t)} + \mathbf{Y}_C \mathbf{Y}_C' \right)^{-1} \end{aligned} \quad (63)$$

$$\begin{aligned} \mathbf{E} \left[\log \det \boldsymbol{\zeta} | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] &= -M \log 2 - \sum_{i=1}^M \mathcal{D} \left(\frac{\delta^{(t)} + Q + N + M - i}{2} \right) \\ &+ \log \det \boldsymbol{\zeta}_0^{(t)} \end{aligned} \quad (64)$$

$$\begin{aligned} \mathbf{E} \left[\boldsymbol{\zeta}^{-1} | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] &= \left(\delta^{(t)} + Q + N + M - 1 \right) \left(\boldsymbol{\zeta}_0^{(t)} \right)^{-1} \end{aligned} \quad (65)$$

$$\begin{aligned} \mathbf{E} \left[\boldsymbol{\tau} \boldsymbol{\zeta}^{-1} \boldsymbol{\tau}' | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] &= \left(\delta^{(t)} + Q + N + M - 1 \right) \boldsymbol{\tau}_0^{(t)} \left(\boldsymbol{\zeta}_0^{(t)} \right)^{-1} \boldsymbol{\tau}_0^{(t)'} \\ &+ M \left(\boldsymbol{\Psi}_{CC}^{(t)} \right)^{-1} \end{aligned} \quad (66)$$

$$\begin{aligned} \mathbf{E} \left[\boldsymbol{\tau} \boldsymbol{\zeta}^{-1} | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)} \right] &= \left(\delta^{(t)} + Q + N + M - 1 \right) \boldsymbol{\tau}_0^{(t)} \left(\boldsymbol{\zeta}_0^{(t)} \right)^{-1}. \end{aligned} \quad (67)$$

Note that $\mathbf{E}[\boldsymbol{\zeta}^{-1} \boldsymbol{\tau}' | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)}] = \mathbf{E}[\boldsymbol{\tau} \boldsymbol{\zeta}^{-1} | \mathbf{Y}_B, \mathbf{Y}_C, \mathcal{H}^{(t)}]'$.

In the GCEM algorithm, the M-step involves the constraints that are to be satisfied by the hyperparameters. Therefore, we implement generalized M-step by converting the problem to an unconstrained one using SUMT to compute some $\mathcal{H}^{(t+1)}$ which satisfies the condition (59) with $\boldsymbol{\Omega}$ replaced by $\boldsymbol{\Omega}^*$. This optimization procedure is discussed in the next section.

IV. NONLINEAR PROGRAMMING

A. Introduction

In this section, we discuss the nonlinear programming method incorporated in solving the maximization (minimization) problem for estimating the hyperparameters. Further details can be obtained from the excellent reference text [40].

B. Implementation of Generalized Constrained M-Step Using SUMT

We use the *interior point* or *barrier function* methods pioneered by Fiacco and McCormick [41] in solving the problem at hand. Since these methods involve minimization of a sequence of unconstrained problems, they are referred to as SUMTs. In the following sections, we develop the theory necessary for implementing the nonlinear programming algorithm using barrier functions to estimate the hyperparameters $\mathcal{H}^{(t+1)}$ in the *generalized constrained M-step*.

1) *BFGS Method*: Once the constrained optimization problem is transformed into a sequence of unconstrained problems using barrier functions, we can use any unconstrained optimization algorithm to solve this sequence of problems. We use the Broyden–Fletcher–Goldfarb–Shanno (BFGS) algorithm [42] for multidimensional minimization by replacing the standard backtracking line search algorithm by that suggested by Moré and Thuente [43], [44] which uses the gradient information and guarantees sufficient function decrease at any chosen step size.

2) *Barrier Function*: Barrier function used to implement the constraints given in (43) is a logarithmic barrier function defined as

$$\Theta(\mathcal{H}) = \begin{cases} \log \delta^{-1} + \sum_{\lambda=1}^{K-1} \log (1 - \rho_{\lambda}^2)^{-1} \\ + \log \det \Psi_{\eta}^{*-1} \\ + \log \det \Psi_{CC}^{-1} + \log \det \zeta_0^{-1} \\ + \infty, \end{cases} \begin{cases} \text{if Condition - 1} \\ \text{otherwise} \end{cases} \quad (68)$$

where Condition-1 is

$$\begin{aligned} \delta &> 0 \\ |\rho_{\lambda}| &\leq 1 \\ \text{for } \lambda &= 1, \dots, K-1 \\ \Psi_{\eta}^* &> 0 \\ \Psi_{CC} &> 0 \\ \zeta_0 &> 0. \end{aligned} \quad (69)$$

Since the maximization of $\Omega^*(\mathcal{H}|\mathcal{H}^{(t)})$ is the same as the minimization of $-\Omega^*(\mathcal{H}|\mathcal{H}^{(t)})$ (negative of the original objective function), we can write the GCEM algorithm as given in Algorithm IV.1.

Algorithm IV.1: GCEM Algorithm for the Estimation of \mathcal{H} With M-Step Using SUMT

Let $\varepsilon > 0$ be a termination scalar, and choose a point $\mathcal{H}^{(1)} \in \chi$ with $\mathbf{g}(\mathbf{x}^{(1)}) < \mathbf{0}$. Let $\mu^{(1)} > 0$, $\theta \in (0,1)$, let $t = 1$, and go to E-step.

E-step Given the data $(\mathbf{Y}_B, \mathbf{Y}_C)$ and the parameter estimates $\mathcal{H}^{(t)}$ at iteration t , evaluate the expectation

$$\Omega^*(\mathcal{H}|\mathcal{H}^{(t)}). \quad (70)$$

Go to M-step.

M-step Given the expectation computed in E-step, minimize using SUMT, the auxiliary function

$$-\Omega^*(\mathcal{H}|\mathcal{H}^{(t)}) + \mu^{(t)}\Theta(\mathcal{H}) \quad \mu > 0 \quad (71)$$

to find some $\mathcal{H}^{(t+1)}$ such that

$$\Omega^*(\mathcal{H}^{(t+1)}|\mathcal{H}^{(t)}) \geq \Omega^*(\mathcal{H}^{(t)}|\mathcal{H}^{(t)}). \quad (72)$$

If $\mu^{(t)}\Theta(\mathcal{H}^{(t+1)}) < \varepsilon$, stop, achieved convergence. Otherwise, set $\mu^{(t+1)} = \theta\mu^{(t)}$, $t = t+1$ and go to E-step.

C. Final Algorithm for Hyperparameter Estimation

In this section, we summarize the algorithmic steps required to carry out the estimation of the hyperparameters, \mathcal{H} . Since the

GCEM algorithm converges to a local maximum (and global maximum if only one maximum is present) when the objective function is extremely nonlinear. One way to solve this problem is to perform the optimization by restarting several times from randomly selected initial values for the parameters being estimated and finding the (possibly, local) maximum in each case. Then, the parameter set that gives the maximum among the many estimated local maxima is said to be the *global maximum*. This is detailed in Algorithm IV.2. The next section discusses the performance evaluation of the hyperparameter estimation and the covariance matrix estimation algorithms and discuss the results.

Algorithm IV.2: Overall Algorithm for Hyperparameter Estimation

```

1: for  $i = 0$  to  $i = \text{MAXRAND}$  in increments of 1 do
2:  $\mathcal{H}^{(1)} \leftarrow$  Random Initialization that is strictly feasible.
3: Initialize  $t = 1$ .
4: Given  $\mathcal{H}^{(t)}$ , apply Algorithm IV.1 to estimate the result  $\mathcal{H}^{(f)}$  via GCEM.
5: Save the function value  $\Omega^*(\mathcal{H}^{(f)}|\mathcal{H}^{(t)})$ , and the parameters  $\mathcal{H}^{(f)}$ .
6: Go to step 1.
7: end for
8:  $\mathcal{H} = \arg \max_{\mathcal{H}^{(f)}} \Omega^*(\mathcal{H}^{(f)}|\mathcal{H}^{(t)})$  is the global maximum of  $\Omega^*(\mathcal{H}|\mathcal{H}^{(t)})$ .

```

V. PERFORMANCE EVALUATION OF HYPERPARAMETER ESTIMATION

In this section, we discuss the protocol used for performing experiments and validating the estimates computed by the GCEM algorithm discussed in the previous sections. We do this indirectly by performing hypothesis testing of the estimated facet model coefficients which are the functions of the estimated population covariance matrix which is in turn a function of the estimated hyperparameters. To do this, we first need a way of computing the covariance matrix of the facet model coefficients. This is described in the next section.

A. Uncertainty in the Estimated Facet Model Coefficients

In this section, we express the uncertainty in the estimated facet model coefficients in terms of the input perturbation described in terms of the noise covariance matrix Σ_{η} . From (20), it is clear that $\hat{\alpha}_n$ is distributed *a posteriori*, with mean α_n (the true unperturbed value) and the covariance matrix $\Sigma_{\hat{\alpha}}$. This is illustrated in this section.

Theorem 4 (Mean and Covariance of the Estimate $\hat{\Lambda}$ of Λ): Given that an estimate $\hat{\Lambda}$ of Λ is given by the (20), the mean and covariance matrix of the estimate $\hat{\Lambda}$, respectively, are

$$\mu_{\hat{\Lambda}} = \Lambda; \quad \Sigma_{\hat{\alpha}} = \zeta. \quad (73)$$

Proof: For the proof, see Appendix VII. ■

B. Verification Through Hypothesis Testing

The approach we take in testing that the software algorithm is producing the right answers is to test the statistical properties of the answers. In other words, we can statistically test whether the statistical properties of the answers obtained are similar the statistical properties we expect. These expectations involve whether the mean of the computed estimates is sufficiently close to the population mean and whether the estimated covariance matrix of the estimates is sufficiently close to the population covariance matrix or both. In our tests, we test to see if the mean of the estimates of the facet model coefficients is close to the assumed mean of the population from which the facet model coefficients arose. We assume that we know the covariance matrix of the facet model coefficients. Using the closed form expression, given by (73), we can compute the expected covariance matrix of the facet model coefficients using the assumed true noise covariance matrix of the data. In doing this, we perform several *monte carlo* type experiments and test our hypothesis about the outcome. For more details on hypothesis testing one may consult [45] and [46].

1) *Hypothesis Testing*: We want to test the hypothesis that $H_0 : \alpha = \alpha_0$ called the *null hypothesis* against the *alternative hypothesis*, $H_A : \alpha \neq \alpha_0$ with the covariance matrix Σ_α is assumed known and fixed at Σ_{α_1} . To test this hypothesis a general method is as follows: a significance level, θ is selected. When the test is run, a test statistic, say ϕ , is computed. This test statistic is typically chosen so that in the case that the hypothesis is true, the test statistic will tend to have its values, distributed around zero, in accordance with a known distribution. If the test statistic has a value, say, higher than a given ϕ_0 , we reject the hypothesis that the computed estimated has statistically behaved as we expected. If we do not reject this hypothesis, it is, in effect, accepting the *null hypothesis*. The value ϕ_0 is chosen so that the probability that we reject the hypothesis, given that the hypothesis is true is less than the significance level θ .

We set up an experiment in which we know what the correct answer α_0 for the no noise ideal case would be. We generate the perturbed data by adding a normally distributed vector from a population having a zero mean and a *given* covariance matrix Σ_η , according to our noise model $x_i = B\alpha_0 + \eta_i$. By using (73), we derive the covariance matrix Σ_{α_1} of the estimates of α from Σ_η . Using the perturbed data, we run our noise covariance matrix estimation procedure and compute the estimates of α . If we repeat this experiment many times by just changing the perturbed realizations and leaving everything else the same, the experiment produces estimates $\alpha_1, \dots, \alpha_N$ that will come from a normal population having mean α_0 , the correct answer for the ideal no noise case, and the covariance matrix Σ_{α_1} . Now, we want to test the hypothesis that the observations $\alpha_1, \dots, \alpha_N$ come from a normal distribution with mean α_0 with known covariance matrix Σ_{α_1} .

Define the test statistic, $\phi = N(\bar{\alpha} - \alpha_0)' \Sigma_{\alpha_1}^{-1} (\bar{\alpha} - \alpha_0)$ where $\bar{\alpha} = (1/N) \sum_{i=1}^N \hat{\alpha}_i$. Distribution ϕ under *null hypothesis* is Chi-squared [45] given by $\phi \sim \chi_M^2$ where M is the dimension of the vector α_i . The distribution under the alternative hypothesis is a noncentral Chi-squared and is given by $\phi \sim \chi_{M,d}^2$ where $d = N(\alpha - \alpha_0)' \Sigma_{\alpha_1}^{-1} (\alpha - \alpha_0)$ is the noncentrality parameter.

2) *Experiments*: We performed the experiment several times for different known values of α_0 , and different values of noise covariance matrices Σ_η . We chose the significance level $\theta = 0.05$. We choose an α_0 and a Σ_η , and we generated $N = 2000$ perturbed realizations x_i for each run of the algorithm. We performed 250 such runs with this combination of α_0 and Σ_η . We then repeated the experiments for different combinations of α_0 and Σ_η and at each run of the algorithm we performed the hypothesis testing. For noise estimation we chose a neighborhood size of 5×5 , resulting in an observation vector size $K = 25$. The covariance matrix Σ_η is of size 25×25 . The DOP basis in this neighborhood has a total of 25 coefficients. We use a fourth order DOP to represent the signal space. That is, the first 15 coefficients are used to represent the signal space. These 15 coefficients constitute the vector α . The remaining ten coefficients are used to represent the noise space. We do not write down the true noise covariance matrix Σ_η here, due to space constraints. Known noise covariance matrix Σ_η is obtained as follows: first, a hyperparameter vector is generated at random, so that each component is uniformly distributed. For the test described here, we generated δ to be uniformly distributed in $[0,1000]$ and let it be denoted by δ_1 . The remaining 25 parameters (correlation filter coefficients) are generated so that each component is uniformly distributed in $[-10,10]$. Given the hyperparameter vector, we create the hypercovariance matrix Ψ_η^* which is premultiplied by $(B C)$ and post-multiplied by its transpose to get Ψ_η . We derive a random sample from an \mathcal{IW} distribution with number of degrees of freedom parameter δ_1 and scale matrix Ψ_η . This random sample constitutes our true covariance matrix Σ_η . This is valid because the \mathcal{GIW} distribution in the form that we used exactly corresponds to \mathcal{IW} distribution. The known mean vector α_0 is chosen at random so that each of its components is sampled uniformly from the interval $[-100,100]$. We generate noise sample vectors η_i from a Gaussian distribution with zero mean and covariance matrix Σ_η and add it to $B\alpha_0$ to generate perturbed input vectors x_i . Results from the first 50 sample experiments are discussed here for brevity. Fig. 2(a) shows the p-values plotted against the experiment number. The horizontal line shows the significance level $\theta = 0.05$ chosen for the experiment. All the points that fell below this line indicate a failure of the test. Fig. 2(b) shows the test statistic ϕ plotted against the experiment number.

VI. APPLICATIONS

As an application of the noise covariance estimation procedure, we designed a new ridge operator that is based on the facet model. The ridge operator uses the integral of the second directional derivative of the facet model in estimating the optimal ridge direction. Once the ridge direction is found it is easy to label a pixel to be a ridge or not, using the first and second directional derivatives. The ridge operator is abbreviated as ISDDRO. The ISDDRO operator is designed to work both under additive colored noise as well as white noise assumptions. Under colored noise assumption, we used the noise covariance matrix estimation procedure discussed in this paper to estimate the facet model coefficients which are then used in ridge detection.

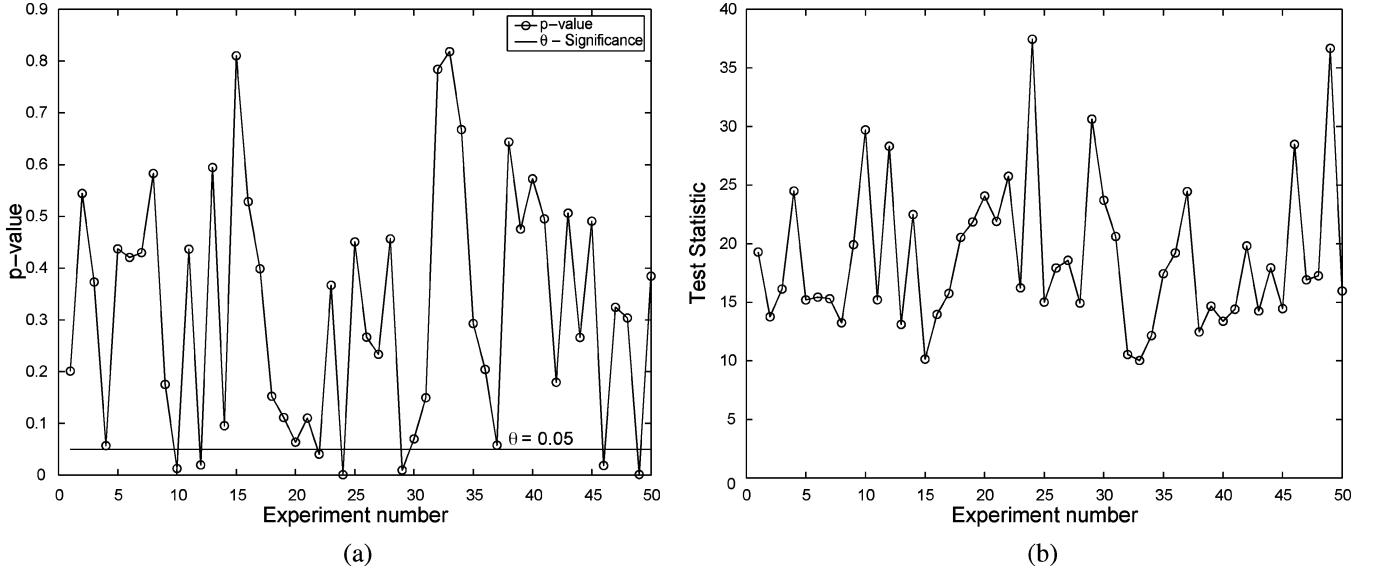


Fig. 2. Hypothesis testing. (a) Plot of p-value. (b) Plot of test statistic, ϕ .

The results of this application are reported in other papers [32], [47]. We evaluated the ridge detectors based on the *ridge orientation estimate*. The *mean orientation bias* which is the difference between the sample mean of the estimates and true orientation, and the *orientation standard deviation* which is the deviation of the estimates from their sample mean are chosen as the objective measures of performance. To summarize the results from the said references, using the chosen performance measures, ISDDRO operator that uses the colored noise assumption (ISDDRO-CN) displays superior noise sensitivity characteristics as compared to the same operator under white noise assumption (ISDDRO-WN) and both versions out performed the most competing operator MLSEC [48]. In fact, ISDDRO under colored noise assumption performs very well under high noise conditions where the other two operators failed.

VII. CONCLUSION

In this paper, we derived the theory and formulated the problem of noise covariance estimation as problem of minimizing an *objective* or *criterion* function in a Bayesian framework. The problem formulation was such that we decomposed the noise covariance matrix into its components in the space spanned by the columns of the signal space basis matrix \mathbf{B} and its complement space spanned by the columns of the matrix \mathbf{C} orthonormal to \mathbf{B} . This facilitated theorizing and proving that the same result for α_n as in the case of white noise is obtained when the signal space and its complement space have no correlation between them, which implied that all energy of the noise in the space spanned by the columns of \mathbf{B} goes into the computed α_n , but the noise in the space spanned by the columns of \mathbf{C} will only influence the estimate of α_n to the extent that there is coupling between the signal space and its complement space. We also derived an expression for α_n for the general case of correlation between the two spaces. To reduce the dimensionality and to ease the optimization we introduced a correlation model for the two spaces that was necessary for the estimation of the population covariance matrix

and derived the constraints which the hyperparameters must obey. We also derived a GCEM algorithm for the estimation of the hyperparameters. We put forward a new way of implementing the constrained M-step of the GCEM algorithm using SUMT. We evaluated the performance of our noise covariance matrix estimation algorithm using statistical hypothesis testing and concluded that the algorithm performs as expected. We also demonstrated, by designing a new ridge detector, that we can design feature extractors with better noise sensitivity characteristics if we assume additive colored noise perturbation as compared to the additive white noise assumption.

APPENDIX I MATRIX NORMAL DISTRIBUTION

In this section, we give only the major results about matrix normal distribution. Further details can be found in [33], [45].

Proposition 1 (Matrix Normal Distribution): A $M \times N$ random matrix \mathbf{X} is said to have the matrix normal distribution $\mathcal{N}(\mathbf{X}_0, \mathbf{\Lambda} \otimes \mathbf{\Upsilon})$, if it has the joint density

$$p(\mathbf{X}) = \left(\frac{1}{(2\pi)^{\frac{NM}{2}} \det \mathbf{\Lambda}^{\frac{N}{2}} \det \mathbf{\Upsilon}^{\frac{M}{2}}} \right) \times \text{etr} \left\{ -\frac{1}{2} \left[\mathbf{\Lambda}^{-1} (\mathbf{X} - \mathbf{X}_0) \right] \left[(\mathbf{X} - \mathbf{X}_0) \mathbf{\Upsilon}^{-1} \right]' \right\} \quad (74)$$

$\mathbf{\Lambda}$ is the covariance matrix within each column, $\mathbf{\Upsilon}$ is the covariance matrix between columns and \mathbf{X}_0 is the mean matrix.

Proposition 2 (Matrix Normal Distrib., With i.i.d. Columns): Let the normal random matrix \mathbf{X} be given in its component form as, $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$, where \mathbf{x}_i is a column vector. Let a_{ij} be an element of $\mathbf{\Lambda} \in \mathbb{R}^{M \times M}$ and b_{ij} be an element of $\mathbf{\Upsilon} \in \mathbb{R}^{N \times N}$. If \mathbf{x}_i are independent, then $\mathbf{\Upsilon} = \mathbf{I}_N$, where $\mathbf{I}_N \in \mathbb{R}^{N \times N}$ is an identity matrix. Therefore, $\mathbf{X} \sim \mathcal{N}(\mathbf{X}_0, \mathbf{\Lambda} \otimes \mathbf{I}_N)$ and the density function is given by

$$p(\mathbf{X}) = \left(\frac{1}{(2\pi)^{\frac{NM}{2}} \det \mathbf{\Lambda}^{\frac{N}{2}}} \right) \times \text{etr} \left\{ -\frac{1}{2} \left[\mathbf{\Lambda}^{-1} (\mathbf{X} - \mathbf{X}_0) \right] \left[\mathbf{X} - \mathbf{X}_0 \right]' \right\} \quad (75)$$

which is equivalent to

$$p(\mathbf{X}) = \left(\frac{1}{(2\pi)^{\frac{NM}{2}} \det \mathbf{\Lambda}^{\frac{N}{2}}} \right) \times \exp \left\{ -\frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \mathbf{x}_{0i})' \mathbf{\Lambda}^{-1} (\mathbf{x}_i - \mathbf{x}_{0i}) \right\} \quad (76)$$

where \mathbf{x}_{0i} are the column vectors of the mean matrix \mathbf{X}_0 .

Proposition 3 (Properties): The proposition states some properties of matrix normal distribution.

- $\mathbf{X} \sim \mathcal{N}(\mathbf{X}_0, \mathbf{\Lambda} \otimes \mathbf{\Upsilon})$ iff $\mathbf{X}' \sim \mathcal{N}(\mathbf{X}'_0, \mathbf{\Upsilon} \otimes \mathbf{\Lambda})$.
- $\mathbf{E}[\mathbf{X}] = \mathbf{X}_0$, $\text{var}[\text{vec}(\mathbf{X})] = \mathbf{\Lambda} \otimes \mathbf{\Upsilon}$, and $\text{var}[\text{vec}(\mathbf{X}')] = \mathbf{\Upsilon} \otimes \mathbf{\Lambda}$.
- $\text{cov}(\mathbf{x}_i, \mathbf{x}_j) = a_{ij} \mathbf{\Upsilon}$, and $\text{cov}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = b_{ij} \mathbf{\Lambda}$.
- For any matrices $\mathbf{C} \in \mathbb{R}^{K \times M}$ and $\mathbf{D} \in \mathbb{R}^{N \times J}$, $\mathbf{C}\mathbf{X}\mathbf{D} = \mathcal{N}(\mathbf{C}\mathbf{X}_0\mathbf{D}, \mathbf{C}\mathbf{\Lambda}\mathbf{C}' \otimes \mathbf{D}'\mathbf{\Upsilon}\mathbf{D})$.

where $\text{var}(\mathbf{x})$ is the variance of \mathbf{x} , and $\text{vec}(\mathbf{X})$ is the vector resulting by stacking up one column of the matrix \mathbf{X} on top of another.

APPENDIX II

WISHART AND INVERTED WISHART DISTRIBUTIONS

In this section, we present a summary of results on *Wishart* and *IW* distributions. Details can be found in [33], [36], [45], and the classic paper by J. Wishart [49]–[51].

Proposition 4 (Wishart Distribution): A $K \times K$ random matrix $\mathbf{\Sigma}$ has a *Wishart distribution* $\mathbf{\Sigma} \sim \mathcal{W}(\nu, \mathbf{\Psi})$, if its joint density has the form

$$p(\mathbf{\Sigma}) = \frac{1}{2^{\frac{\nu K}{2}} \Gamma_K(\frac{\nu}{2})} \det \mathbf{\Psi}^{-\frac{\nu}{2}} \det \mathbf{\Sigma}^{\frac{\nu-K-1}{2}} \times \text{etr} \left\{ -\frac{1}{2} \mathbf{\Psi}^{-1} \mathbf{\Sigma} \right\} \quad (77)$$

where ν is the number of degrees of freedom, $\Gamma_K(t)$ is a multidimensional Gamma function given by, $\Gamma_K(t) = \pi^{K(K-1)/4} \prod_{i=1}^K \Gamma[t - ((i-1)/2)]$.

Proposition 5 (IW Distribution): A $K \times K$ matrix $\mathbf{\Sigma}$ has an *IW distribution* $\mathbf{\Sigma} \sim \mathcal{IW}(\delta, \mathbf{\Psi})$, if it has the density function of the following form:

$$p(\mathbf{\Sigma}) = \frac{1}{2^{\frac{(\delta+K-1)K}{2}} \Gamma_K(\frac{\delta+K-1}{2})} \det \mathbf{\Psi}^{\frac{\delta+K-1}{2}} \times \det \mathbf{\Sigma}^{-\left(\frac{\delta+2K}{2}\right)} \text{etr} \left\{ -\frac{1}{2} (\mathbf{\Psi} \mathbf{\Sigma}^{-1}) \right\} \quad (78)$$

where δ is the parameter of the distribution which is chosen because it is invariant under change of dimension [11], [33], and $\nu = \delta + K - 1$ and $\Gamma_K(t)$ is a multidimensional gamma function as defined in Proposition 4.

Proposition 6 (Properties): Some properties of the *Wishart* and *IW* distributions are given below.

- $\mathbf{\Sigma} \sim \mathcal{IW}(\delta, \mathbf{\Psi})$ iff $\mathbf{\Sigma}^{-1} \sim \mathcal{W}(\nu, \mathbf{\Psi}^{-1})$.
- If $\mathbf{\Sigma} \sim \mathcal{W}(\nu, \mathbf{\Psi})$, then $\mathbf{E}[\mathbf{\Sigma}] = \nu \mathbf{\Psi}$, and $\mathbf{E}[\mathbf{\Sigma}^{-1}] = (1/(\delta-2)) \mathbf{\Psi}^{-1}$ provided $(\delta-2) > 0$.
- If $\mathbf{\Sigma} \sim \mathcal{IW}(\delta, \mathbf{\Psi})$, then $\mathbf{E}[\mathbf{\Sigma}] = (1/(\delta-2)) \mathbf{\Psi}$, and $\mathbf{E}[\mathbf{\Sigma}^{-1}] = \nu \mathbf{\Psi}^{-1}$.

- If $\mathbf{\Sigma} \sim \mathcal{IW}(\delta, \mathbf{\Psi})$, then, $\mathbf{E}[\log \det \mathbf{\Sigma}] = -K \log 2 - \sum_{i=1}^K \mathcal{D}[(\delta + K - i)/2] + \log \det \mathbf{\Psi}$, where $\mathcal{D}(x) = (d/dx) \log \Gamma(x)$ is the *digamma* function.

The following proposition states that the \mathcal{IW} distribution is closed under linear transformations.

Proposition 7 (Linear Transformation of IW): Let $\mathbf{\Sigma} \sim \mathcal{IW}(\delta, \mathbf{\Psi})$. Let \mathbf{T} be a given transformation matrix. Then, $\mathbf{T}'\mathbf{\Sigma}\mathbf{T} \sim \mathcal{IW}(\delta, \mathbf{T}'\mathbf{\Psi}\mathbf{T})$.

Theorem 5 (Orthonormal Transformation of $p(\mathbf{\Sigma})$): Let $p(\mathbf{\Sigma})$ denote a matrix variate distribution. Let \mathbf{T} be an orthonormal transformation matrix. Then, $p(\mathbf{\Sigma}) = p(\mathbf{T}'\mathbf{\Sigma}\mathbf{T})$.

Proof: The proof of this theorem is self evident as the determinant of the Jacobian matrix $\det \mathbf{J} = \det \mathbf{T}'\mathbf{T} = 1$, since \mathbf{T} is an orthonormal matrix with unity determinant. ■

APPENDIX III

GENERALIZED INVERTED WISHART DISTRIBUTION

In this section, we discuss the \mathcal{GIW} distribution introduced by Brown *et al.* [19].

Let $\mathbf{X} \in \mathbb{R}^{K \times N}$ be a data matrix containing N , K -dimensional independent observation vectors arranged in columns. Then, according to Section I, $\mathbf{X} \sim \mathcal{N}(\mathbf{X}_0, \mathbf{\Sigma} \otimes \mathbf{I}_N)$.

Let the matrix $\mathbf{\Sigma}$ be distributed as $\mathcal{IW}(\delta, \mathbf{\Psi})$ where δ is as defined in Section II, and $\mathbf{\Psi}$ is the scale matrix.

Let $\mathbf{\Sigma}$ be partitioned naturally as

$$\mathbf{\Sigma} = \begin{pmatrix} \mathbf{\Sigma}_{11} & \mathbf{\Sigma}_{12} \\ \mathbf{\Sigma}_{21} & \mathbf{\Sigma}_{22} \end{pmatrix}. \quad (79)$$

Conformably, let the data matrix \mathbf{X} be partitioned as $\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}$. Let $\mathbf{\Psi}$ be partitioned conformably as

$$\mathbf{\Psi} = \begin{pmatrix} \mathbf{\Psi}_{11} & \mathbf{\Psi}_{12} \\ \mathbf{\Psi}_{21} & \mathbf{\Psi}_{22} \end{pmatrix}. \quad (80)$$

Using Bartlett [34] decomposition, we have

$$\mathbf{\Sigma} = \mathbf{T}\mathbf{\Omega}\mathbf{T}' \quad (81)$$

where $\mathbf{\Omega} = \begin{pmatrix} \zeta & \mathbf{0} \\ \mathbf{0} & \mathbf{\Sigma}_{22} \end{pmatrix}$ and $\mathbf{T} = \begin{pmatrix} \mathbf{I} & \boldsymbol{\tau}' \\ \mathbf{0} & \mathbf{I} \end{pmatrix}$. Hence

$$\mathbf{\Sigma} = \begin{pmatrix} \zeta + \boldsymbol{\tau}' \mathbf{\Sigma}_{22} \boldsymbol{\tau} & \boldsymbol{\tau}' \mathbf{\Sigma}_{22} \\ \mathbf{\Sigma}_{22} \boldsymbol{\tau} & \mathbf{\Sigma}_{22} \end{pmatrix} \quad (82)$$

where $\mathbf{\Sigma}_{22} : Q \times Q$, $\zeta : M \times M$ and $\boldsymbol{\tau} : Q \times M$ and $\boldsymbol{\tau} = \mathbf{\Sigma}_{22}^{-1} \mathbf{\Sigma}_{21}$; $\zeta = \mathbf{\Sigma}_{11} - \mathbf{\Sigma}_{12} \mathbf{\Sigma}_{22}^{-1} \mathbf{\Sigma}_{21}$; $\boldsymbol{\tau}_0 = \mathbf{\Psi}_{22}^{-1} \mathbf{\Psi}_{21}$; $\zeta_0 = \mathbf{\Psi}_{11} - \mathbf{\Psi}_{12} \mathbf{\Psi}_{22}^{-1} \mathbf{\Psi}_{21}$; $K = M + Q$.

Proposition 8 (GIW Distribution): Let $\mathbf{\Sigma} \sim \mathcal{IW}(\delta, \mathbf{\Psi})$. If $\mathbf{\Sigma}$ and $\mathbf{\Psi}$ have the conformable partitioning given in (79) and (80), respectively, then in terms of its block partitions $\mathbf{\Sigma}$ has a \mathcal{GIW} model [19] given as

- 1) $\mathbf{\Sigma}_{22} \sim \mathcal{IW}(\delta, \mathbf{\Psi}_{22})$, is distributed independently of $\boldsymbol{\tau}$ and ζ ;
- 2) $\zeta \sim \mathcal{IW}(\delta + Q, \zeta_0)$;
- 3) $\boldsymbol{\tau} | \zeta \sim \mathcal{N}(\boldsymbol{\tau}_0, \mathbf{\Psi}_{22}^{-1} \otimes \zeta)$.

It is clear that $\boldsymbol{\tau}$ is the *slope* of the best linear predictor of the observables \mathbf{X}_1 based on the observables \mathbf{X}_2 and ζ is the residual covariance matrix of the resulting prediction errors.

Under the above assumptions Σ is said to have the $\mathcal{G}\mathcal{I}\mathcal{W}$ distribution with parameter δ and the scale matrix Ψ . This is denoted as $\Sigma \sim \mathcal{G}\mathcal{I}\mathcal{W}(\delta, \Psi)$. Under this notation, it is assumed that Σ and Ψ have the aforementioned partitioned structure. The $\mathcal{G}\mathcal{I}\mathcal{W}$ in the above form exactly corresponds to the standard $\mathcal{I}\mathcal{W}(\delta, \Psi)$ on the unpartitioned Σ . We still use the term $\mathcal{G}\mathcal{I}\mathcal{W}$ to refer to this special case of $\mathcal{G}\mathcal{I}\mathcal{W}$ just to differentiate this hierarchical structure from the standard $\mathcal{I}\mathcal{W}$.

APPENDIX IV OPTIMIZATION OF SCALAR FUNCTIONS OF MATRICES

In this section, we state and prove [45], the following theorem, which we use in this paper.

Theorem 6 (Maximum of a Scalar Function of a Matrix): If $\mathbf{D} \in \mathbb{R}^{M \times M}$ is positive definite and C is a constant, then the maximum of

$$f(\mathbf{G}) = -C \log \det \mathbf{G} - \text{tr}(\mathbf{G}^{-1} \mathbf{D}) \quad (83)$$

with respect to positive definite matrices \mathbf{G} exists, and occurs at

$$\mathbf{G} = \frac{1}{C} \mathbf{D} \quad (84)$$

and has the value of $f((1/C)\mathbf{D}) = MC \log C - C \log \det \mathbf{D} - MC$.

Proof: Let $\mathbf{D} = \mathbf{E}\mathbf{E}'$ and $\mathbf{E}'\mathbf{G}^{-1}\mathbf{E} = \mathbf{H}$. Then, $\mathbf{G} = \mathbf{E}\mathbf{H}^{-1}\mathbf{E}'$, and $\det \mathbf{G} = \det \mathbf{E} \det \mathbf{H}^{-1} \det \mathbf{E}' = \det \mathbf{H}^{-1} \det(\mathbf{E}\mathbf{E}') = (\det \mathbf{D} / \det \mathbf{H})$, and $\text{tr}(\mathbf{G}^{-1} \mathbf{D}) = \text{tr}(\mathbf{G}^{-1} \mathbf{E}\mathbf{E}') = \text{tr}(\mathbf{E}'\mathbf{G}^{-1}\mathbf{E}) = \text{tr} \mathbf{H}$. Then, the function to be maximized (with respect to the positive definite \mathbf{H}) is

$$f = -C \log \det \mathbf{D} + C \log \det \mathbf{H} - \text{tr} \mathbf{H}. \quad (85)$$

Let $\mathbf{H} = \mathbf{T}\mathbf{T}'$, where \mathbf{T} is lower triangular. Then, the maximum of

$$\begin{aligned} f &= -C \log \det \mathbf{D} + C \log(\det \mathbf{T})^2 - \text{tr}(\mathbf{T}\mathbf{T}') \\ &= -C \log \det \mathbf{D} + \sum_{i=1}^M (C \log t_{ii}^2 - t_{ii}^2) - \sum_{i>j} t_{ij}^2 \end{aligned} \quad (86)$$

occurs at $t_{ii}^2 = C$, $t_{ij} = 0$, $i \neq j$, that is, at $\mathbf{H} = C\mathbf{I}$. Then, $\mathbf{G} = (1/C)\mathbf{E}\mathbf{E}' = (1/C)\mathbf{D}$. ■

APPENDIX V CONDITIONAL DENSITY OF \mathbf{Y}_B GIVEN \mathbf{Y}_C

In this section, we derive the conditional density of \mathbf{Y}_B given \mathbf{Y}_C . This is necessary for computing the joint posterior density of Λ and Σ_{CC} , τ , and ζ given the data $(\mathbf{Y}'_B, \mathbf{Y}'_C)$.

Theorem 7 (Conditional Density of \mathbf{Y}_B Given \mathbf{Y}_C): If

$$\begin{aligned} \mathbf{Y} &= \begin{pmatrix} \mathbf{Y}_B \\ \mathbf{Y}_C \end{pmatrix} \\ &\sim \mathcal{N} \left(\begin{pmatrix} \Lambda \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \Sigma_{BB} & \Sigma_{BC} \\ \Sigma_{CB} & \Sigma_{CC} \end{pmatrix} \otimes \mathbf{I}_N \right) \end{aligned} \quad (87)$$

and

$$\mathbf{Y}_B \sim \mathcal{N}(\Lambda, \Sigma_{BB} \otimes \mathbf{I}_N) \quad (88)$$

$$\mathbf{Y}_C \sim \mathcal{N}(\mathbf{0}, \Sigma_{CC} \otimes \mathbf{I}_N) \quad (89)$$

then

$$\mathbf{Y}_B | \mathbf{Y}_C \sim \mathcal{N}(\Lambda + \tau' \mathbf{Y}_C, \zeta \otimes \mathbf{I}_N) \quad (90)$$

where $\zeta = \Sigma_{BB} - \Sigma_{BC} \Sigma_{CC}^{-1} \Sigma_{CB}$ and $\tau = \Sigma_{CC}^{-1} \Sigma_{CB}$.

Proof: Let us define the following linear transformations of the observations:

$$\mathbf{Y}_B^\dagger = \mathbf{Y}_B + \mathbf{A} \mathbf{Y}_C \quad (91)$$

$$\mathbf{Y}_C^\dagger = \mathbf{Y}_C. \quad (92)$$

We need to choose the transformation matrix \mathbf{A} such that \mathbf{Y}_B^\dagger and \mathbf{Y}_C^\dagger are uncorrelated. This condition means that the cross-correlation must be zero, i.e., $\mathbf{E}[(\mathbf{Y}_B^\dagger - \mathbf{E}[\mathbf{Y}_B^\dagger])(\mathbf{Y}_C^\dagger - \mathbf{E}[\mathbf{Y}_C^\dagger])'] = 0$, where $\mathbf{E}[\mathbf{Y}_B^\dagger] = \Lambda$ and $\mathbf{E}[\mathbf{Y}_C^\dagger] = \mathbf{0}$. Substituting for \mathbf{Y}_B^\dagger and \mathbf{Y}_C^\dagger from (92) and (91), and their expectations just computed, we get $\mathbf{E}[(\mathbf{Y}_B + \mathbf{A} \mathbf{Y}_C - \Lambda) \mathbf{Y}_C'] = \mathbf{0}$ and $\mathbf{E}[(\mathbf{Y}_B - \Lambda) \mathbf{Y}_C'] + \mathbf{A} \mathbf{E}[\mathbf{Y}_C \mathbf{Y}_C'] = \mathbf{0}$. Now, using the fact that $\mathbf{E}[(\mathbf{Y}_B - \Lambda) \mathbf{Y}_C'] = \Sigma_{BC}$ and $\mathbf{E}[\mathbf{Y}_C \mathbf{Y}_C'] = \Sigma_{CC}$, we get $\Sigma_{BC} + \mathbf{A} \Sigma_{CC} = \mathbf{0}$ and $\mathbf{A} = -\Sigma_{BC} \Sigma_{CC}^{-1}$. Now, (92) and (91) can be written in matrix form as $\mathbf{Y}^\dagger = \begin{pmatrix} \mathbf{Y}_B^\dagger \\ \mathbf{Y}_C^\dagger \end{pmatrix} = \begin{pmatrix} \mathbf{I} & -\Sigma_{BC} \Sigma_{CC}^{-1} \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{Y}_B \\ \mathbf{Y}_C \end{pmatrix}$. Taking expectations of this result, we get

$$\mathbf{E}[\mathbf{Y}^\dagger] = \begin{pmatrix} \Lambda \\ \mathbf{0} \end{pmatrix}. \quad (93)$$

Covariance matrix $\Sigma_{\mathbf{Y}^\dagger}$ of \mathbf{Y}^\dagger is given by

$$\begin{aligned} \Sigma_{\mathbf{Y}^\dagger} &= \mathbf{E} \left[\left(\mathbf{Y}^\dagger - \mathbf{E}[\mathbf{Y}^\dagger] \right) \left(\mathbf{Y}^\dagger - \mathbf{E}[\mathbf{Y}^\dagger] \right)' \right] \\ &= \begin{pmatrix} \mathbf{E} \left[\left(\mathbf{Y}_B^\dagger - \mathbf{E}[\mathbf{Y}_B^\dagger] \right) \left(\mathbf{Y}_B^\dagger - \mathbf{E}[\mathbf{Y}_B^\dagger] \right)' \right] & \mathbf{0} \\ \mathbf{0} & \Sigma_{CC} \end{pmatrix}. \end{aligned} \quad (94)$$

We used the fact that \mathbf{Y}_B^\dagger and \mathbf{Y}_C^\dagger are uncorrelated. Consider the term

$$\begin{aligned} &\mathbf{E} \left[\left(\mathbf{Y}_B^\dagger - \mathbf{E}[\mathbf{Y}_B^\dagger] \right) \left(\mathbf{Y}_B^\dagger - \mathbf{E}[\mathbf{Y}_B^\dagger] \right)' \right] \\ &= \mathbf{E} \left[\left(\mathbf{Y}_B - \Sigma_{BC} \Sigma_{CC}^{-1} \mathbf{Y}_C - \Lambda \right) \right. \\ &\quad \left. \times \left(\mathbf{Y}_B - \Sigma_{BC} \Sigma_{CC}^{-1} \mathbf{Y}_C - \Lambda \right)' \right] \\ &= \mathbf{E}[(\mathbf{Y}_B - \Lambda)(\mathbf{Y}_B - \Lambda)'] \\ &\quad - \mathbf{E}[(\mathbf{Y}_B - \Lambda) \mathbf{Y}_C'] \Sigma_{CC}^{-1} \Sigma_{CB} \\ &\quad - \Sigma_{BC} \Sigma_{CC}^{-1} \mathbf{E}[\mathbf{Y}_C (\mathbf{Y}_B - \Lambda)'] \\ &\quad + \Sigma_{BC} \Sigma_{CC}^{-1} \mathbf{E}[\mathbf{Y}_C \mathbf{Y}_C'] \Sigma_{CC}^{-1} \Sigma_{CB} \\ &= \Sigma_{BB} - \Sigma_{BC} \Sigma_{CC}^{-1} \Sigma_{CB}. \end{aligned} \quad (95)$$

Therefore

$$\Sigma_{\mathbf{Y}^\dagger} = \begin{pmatrix} \zeta & \mathbf{0} \\ \mathbf{0} & \Sigma_{CC} \end{pmatrix} \quad (96)$$

where $\zeta = \Sigma_{BB} - \Sigma_{BC}\Sigma_{CC}^{-1}\Sigma_{CB}$. Since uncorrelated normally distributed random variables are also *independent*, the conditional density of \mathbf{Y}_B^\dagger given \mathbf{Y}_C (note that $\mathbf{Y}_C^\dagger = \mathbf{Y}_C$) is same as the marginal density of \mathbf{Y}_B^\dagger . Therefore, $\mathbf{Y}_B^\dagger|\mathbf{Y}_C \sim \mathcal{N}(\Lambda, \zeta)$. Given \mathbf{Y}_C , we have

$$\mathbf{Y}_B = \mathbf{Y}_B^\dagger + \Sigma_{BC}\Sigma_{CC}^{-1}\mathbf{Y}_C. \quad (97)$$

Note that \mathbf{Y}_C is a constant in the conditional density of $\mathbf{Y}_B|\mathbf{Y}_C$. We have from (97) that

$$\begin{aligned} E[\mathbf{Y}_B|\mathbf{Y}_C] &= E\left[\mathbf{Y}_B^\dagger\right] + \Sigma_{BC}\Sigma_{CC}^{-1}\mathbf{Y}_C \\ &= \Lambda + \Sigma_{BC}\Sigma_{CC}^{-1}\mathbf{Y}_C. \end{aligned} \quad (98)$$

The covariance matrix of the conditional density $\mathbf{Y}_B|\mathbf{Y}_C$ is

$$\Sigma_{\mathbf{Y}_B|\mathbf{Y}_C} = E\left[(\mathbf{Y}_B - E[\mathbf{Y}_B|\mathbf{Y}_C])(\mathbf{Y}_B - E[\mathbf{Y}_B|\mathbf{Y}_C])'\right]. \quad (99)$$

Substituting for \mathbf{Y}_B and $E[\mathbf{Y}_B|\mathbf{Y}_C]$ from (97) and (98), respectively, and using (93), (96), and (95) and simplifying, we get, $\Sigma_{\mathbf{Y}_B|\mathbf{Y}_C} = \zeta$. Therefore, we can write

$$\mathbf{Y}_B|\mathbf{Y}_C \sim \mathcal{N}\left(\Lambda + \Sigma_{BC}\Sigma_{CC}^{-1}\mathbf{Y}_C, \zeta \otimes \mathbf{I}_N\right) \quad (100)$$

where $\tau = \Sigma_{CC}^{-1}\Sigma_{CB}$. ■

APPENDIX VI PROOF OF THEOREM 3

Proof: Joint density of $(\Sigma_{CC}, \tau, \zeta)$ and the data $(\mathbf{Y}_B, \mathbf{Y}_C)$ is given by

$$\begin{aligned} &p(\mathbf{Y}_B, \mathbf{Y}_C, \Lambda, \Sigma_{CC}, \tau, \zeta, \delta, \Psi_{CC}, \tau_0, \zeta_0) \\ &= p(\mathbf{Y}_B|\mathbf{Y}_C, \Lambda, \Sigma_{CC}, \tau, \zeta, \delta, \Psi_{CC}, \tau_0, \zeta_0) \\ &\quad \times p(\mathbf{Y}_C|\Lambda, \Sigma_{CC}, \tau, \zeta, \delta, \Psi_{CC}, \tau_0, \zeta_0) \\ &\quad \times p(\Lambda)p(\Sigma_{CC}|\delta, \Psi_{CC}) \\ &\quad \times p(\tau|\zeta, \delta, \Psi_{CC}, \tau_0, \zeta_0)p(\zeta|\delta, \zeta_0) \\ &= \left(\frac{1}{(2\pi)^{\frac{NM}{2}} \det \zeta^{\frac{N}{2}}}\right) \\ &\quad \times \text{etr}\left\{-\frac{1}{2}\zeta^{-1}(\mathbf{Y}_B - \Lambda - \tau'\mathbf{Y}_C)\right. \\ &\quad \quad \left.\times (\mathbf{Y}_B - \Lambda - \tau'\mathbf{Y}_C)'\right\} \\ &\quad \times \frac{1}{(2\pi)^{\frac{NQ}{2}} \det \Sigma_{CC}^{\frac{N}{2}}} \text{etr}\left\{-\frac{1}{2}\Sigma_{CC}^{-1}\mathbf{Y}_C\mathbf{Y}_C'\right\} \\ &\quad \times \frac{\det \Psi_{CC}^{\frac{\delta+Q-1}{2}} \det \Sigma_{CC}^{\frac{\delta+2Q}{2}}}{2^{\frac{(\delta+Q-1)Q}{2}} \Gamma_Q\left(\frac{\delta+Q-1}{2}\right)} \text{etr}\left\{-\frac{1}{2}\Psi_{CC}\Sigma_{CC}^{-1}\right\} \\ &\quad \times \frac{\det \zeta_0^{\frac{\delta+Q+M-1}{2}} \det \zeta^{-\frac{\delta+Q+2M}{2}}}{2^{\frac{(\delta+Q+M-1)M}{2}} \Gamma_M\left(\frac{\delta+Q+M-1}{2}\right)} \text{etr}\left\{-\frac{1}{2}\zeta_0\zeta^{-1}\right\} \\ &\quad \times \left(\frac{1}{(2\pi)^{\frac{MQ}{2}} \det \Psi_{CC}^{-\frac{M}{2}} \det \zeta^{\frac{Q}{2}}}\right) \\ &\quad \times \text{etr}\left\{-\frac{1}{2}[\Psi_{CC}(\tau - \tau_0)][(\tau - \tau_0)\zeta^{-1}]\right\}. \end{aligned} \quad (101)$$

Now, substituting $\Lambda = \hat{\Lambda}$ the first exponential in above equation becomes unity. Now, multiplying and dividing the (101) by

$$\frac{\det(\Psi_{CC} + \mathbf{Y}_C\mathbf{Y}_C')^{\frac{\delta+N+Q-1}{2}} \det \zeta_0^{\frac{N}{2}}}{\Gamma_Q\left(\frac{\delta+N+Q-1}{2}\right) \Gamma_M\left(\frac{\delta+Q+N+M-1}{2}\right)} \quad (102)$$

and rearranging, we get

$$\begin{aligned} &= \left(\frac{\det(\Psi_{CC} + \mathbf{Y}_C\mathbf{Y}_C')^{\frac{\delta+N+Q-1}{2}} \det \Sigma_{CC}^{-\frac{\delta+N+2Q}{2}}}{2^{\frac{(\delta+N+Q-1)Q}{2}} \Gamma_Q\left(\frac{\delta+N+Q-1}{2}\right)}\right) \\ &\quad \times \text{etr}\left\{-\frac{1}{2}(\Psi_{CC} + \mathbf{Y}_C\mathbf{Y}_C')\Sigma_{CC}^{-1}\right\} \\ &\quad \times \left(\frac{\det \zeta_0^{\frac{\delta+Q+N+M-1}{2}} \det \zeta^{-\frac{\delta+Q+N+2M}{2}}}{2^{\frac{(\delta+Q+N+M-1)M}{2}} \Gamma_M\left(\frac{\delta+Q+N+M-1}{2}\right)}\right) \\ &\quad \times \text{etr}\left\{-\frac{1}{2}\zeta_0\zeta^{-1}\right\} \times \left(\frac{1}{(2\pi)^{\frac{MQ}{2}} \det \Psi_{CC}^{-\frac{M}{2}} \det \zeta^{\frac{Q}{2}}}\right) \\ &\quad \times \text{etr}\left\{-\frac{1}{2}[\Psi_{CC}(\tau - \tau_0)][(\tau - \tau_0)\zeta^{-1}]\right\} \\ &\quad \times \frac{\text{Term} - 1}{\text{Term} - 2} \end{aligned} \quad (103)$$

where

$$\begin{aligned} \text{Term} - 1 &= \det \Psi_{CC}^{\frac{\delta+Q-1}{2}} \Gamma_Q\left(\frac{\delta + N + Q - 1}{2}\right) \\ &\quad \times \Gamma_M\left(\frac{\delta + Q + N + M - 1}{2}\right) \quad (104) \\ \text{Term} - 2 &= \pi^{\frac{NQ}{2}} \Gamma_Q\left(\frac{\delta + Q - 1}{2}\right) \Gamma_M\left(\frac{\delta + Q + M - 1}{2}\right) \\ &\quad \times \det(\Psi_{CC} + \mathbf{Y}_C\mathbf{Y}_C')^{\frac{\delta+N+Q-1}{2}} \det \zeta_0^{\frac{N}{2}}. \end{aligned} \quad (105)$$

Integrating the above equation, in turn, over Σ_{CC} , τ , and ζ , the first three terms yield unity as they are proper density functions leaving only the last term as the result, given by

$$= \frac{\text{Term} - 1}{\text{Term} - 2}. \quad (106)$$

Dividing the joint density in (103) by the term in (106), we get the result for the posterior density. ■

APPENDIX VII PROOF OF THEOREM 4

Proof: Mean of the $\hat{\alpha}_n$ is given by

$$\begin{aligned} E(\hat{\alpha}_n) &= (\mathbf{B}' - \Sigma_{BC}\Sigma_{CC}^{-1}\mathbf{C}') E(x_n) \\ &= (\mathbf{B}' - \Sigma_{BC}\Sigma_{CC}^{-1}\mathbf{C}') E(\mathbf{B}\alpha_n + \eta_n) \\ &= (\mathbf{B}' - \Sigma_{BC}\Sigma_{CC}^{-1}\mathbf{C}') \mathbf{B}\alpha_n \\ &= \alpha_n \end{aligned} \quad (107)$$

since $E(\eta_n) = \mathbf{0}$, $\mathbf{B}'\mathbf{B} = \mathbf{I}$ and $\mathbf{C}'\mathbf{B} = \mathbf{0}$. Therefore, $\hat{\alpha}_n$ is an unbiased estimator of α_n . Now, consider

$$\begin{aligned} \hat{\alpha}_n &= (\mathbf{B}' - \Sigma_{BC}\Sigma_{CC}^{-1}\mathbf{C}') x_n \\ &= (\mathbf{B}' - \Sigma_{BC}\Sigma_{CC}^{-1}\mathbf{C}') (\mathbf{B}\alpha_n + \eta_n) \\ &= \alpha_n + (\mathbf{B}' - \Sigma_{BC}\Sigma_{CC}^{-1}\mathbf{C}') \eta_n. \end{aligned} \quad (108)$$

The covariance matrix $\Sigma_{\hat{\alpha}_n}$ of $\hat{\alpha}_n$ is

$$\Sigma_{\hat{\alpha}} = E(\hat{\alpha}_n - \alpha_n)(\hat{\alpha}_n - \alpha_n)'. \quad (109)$$

Substituting for $\hat{\alpha}_n$ from (108), using the fact that $E(\eta_n \eta_n') = \Sigma_{\eta}$, and simplifying, we get

$$\Sigma_{\hat{\alpha}} = (B' - \Sigma_{BC} \Sigma_{CC}^{-1} C') \Sigma_{\eta} \times (B' - \Sigma_{BC} \Sigma_{CC}^{-1} C')'. \quad (110)$$

We can express Σ_{η} in $(B \ C)$ space, by using the fact that $(B \ C)^{-1} = (B \ C)'$, we have

$$\begin{aligned} \Sigma_{\eta} &= (B \ C)(B \ C)' \Sigma_{\eta} (B \ C)(B \ C)' \\ &= (B \ C) \begin{pmatrix} \Sigma_{BB} & \Sigma_{BC} \\ \Sigma_{CB} & \Sigma_{CC} \end{pmatrix} (B \ C)'. \end{aligned} \quad (111)$$

After substituting for Σ_{η} from (111) into (110) and simplifying, we get the covariance matrix of the estimated α_n as

$$\Sigma_{\hat{\alpha}} = \Sigma_{BB} - \Sigma_{BC} \Sigma_{CC}^{-1} \Sigma_{CB} = \zeta. \quad (112)$$

Note that, in the *white* noise case [52], due to the independence of the spaces spanned by the columns of B , and the columns of C , we have

$$\Sigma_{\hat{\alpha}} = \Sigma_{BB}. \quad (113)$$

REFERENCES

- [1] V. Chalana, "Deformable models for segmentation of medical ultrasound images," Ph.D. dissertation, Dept. Bioeng., Univ. Washington, Seattle, 1996.
- [2] L. D. Cohen and I. Cohen, "Deformable models for 3-D medical images using finite elements and balloons," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 1992, pp. 592–598.
- [3] D. Terzopoulos and D. Metaxas, "Dynamic 3D models with local and global deformations: deformable superquadrics," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 13, no. 7, pp. 703–714, Jul. 1991.
- [4] T. McInerney and D. Terzopoulos, "A dynamic finite element surface model for segmentation and tracking in multidimensional medical images with application to cardiac 4D image analysis," *J. Comput. Med. Imag. Graph.*, to be published.
- [5] L. D. Cohen and I. Cohen, "Finite-element methods for active contour models and balloons for 2-D and 3-D images," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 15, no. 10, pp. 1131–1147, Oct. 1993.
- [6] M. Kass, A. Witkin, and D. Terzopoulos, "Snakes: active contour models," *Int. J. Comput. Vis.*, vol. 1, pp. 321–331, 1988.
- [7] C. Stein, "Estimation of a covariance matrix," presented at the Rietz Lecture, Annu. Meeting Institute of Mathematics and Statistics, Atlanta, GA, 1975.
- [8] —, "Lectures on the theory of estimation of many parameters," in *Proc. Scientific Seminars of the Steklov Institute, Leningrad Division, Studies in the Statistical Theory of Estimation, Part I*, I. A. Ibragimov and M. S. Nikulin, Eds., 1977, pp. 4–65.
- [9] S. P. Lin and M. D. Perlman, "An improved procedure for the estimation of a correlation matrix," in *Statistical Theory and Data Analysis*, K. Matusita, Ed. Amsterdam, The Netherlands: Elsevier, 1985, pp. 369–379.
- [10] L. R. Haff, "Empirical bayes estimation of the multivariate normal covariance matrix," *Ann. Stat.*, vol. 8, pp. 586–597, 1980.
- [11] J. M. Dickey, D. V. Lindley, and S. J. Press, "Bayesian estimation of the dispersion matrix of a multivariate normal distribution," *Commun. Stat. Theory Meth.*, vol. 14, no. 5, pp. 1019–1034, 1985.
- [12] S. Wright, "The method of path coefficients," *Ann. Math. Stat.*, vol. 5, pp. 161–215, 1934.
- [13] K. G. Jöreskog, *Structural Equation Models in the Social Sciences: Specification, Estimation and Testing*. Amsterdam, The Netherlands: North-Holland, 1977, pp. 265–286.
- [14] C. Spearman, "Proof and measurement of association between two things," *Amer. J. Psych.*, vol. 15, pp. 72–202, 1904.
- [15] J. P. Hoffbeck and D. A. Landgrebe, "Covariance matrix estimation and classification with limited training data," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 18, no. 7, pp. 763–767, Jul. 1996.
- [16] J. P. Le Cadre, "Parametric methods for spatial signal processing in the presence of unknown colored noise fields," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 37, no. 3, pp. 965–983, Jun. 1989.
- [17] T. Leonard and J. S. J. Hsu, "Bayesian inference for a covariance matrix," *Ann. Stat.*, vol. 20, no. 4, pp. 1669–1696, 1992.
- [18] M. J. Daniels and R. E. Kass, "Nonconjugate Bayesian estimation of covariance matrices and its use in hierarchical models," Dept. Stat., Carnegie Mellon Univ., Pittsburgh, PA, Tech. Rep. 659, Nov. 1998.
- [19] P. J. Brown, N. D. Le, and J. V. Zidek, *Aspects of Uncertainty: A Tribute to D. V. Lindley*. New York: Wiley, 1994, ch. Inference for a covariance matrix, pp. 77–92.
- [20] R. M. Haralick and L. T. Watson, "A facet model for image data," *Comput. Graph. Image Process.*, vol. 15, pp. 113–129, 1981.
- [21] R. M. Haralick, "A Bayesian approach to robust local facet estimation," in *Maximum-Entropy and Bayesian Spectral Analysis and Estimation Problems*, C. R. Smith and G. J. Erickson, Eds. Dordrecht, The Netherlands: Reidel, 1987, pp. 85–97.
- [22] R. M. Haralick and L. G. Shapiro, *Computer and Robot Vision*. Reading, MA: Addison-Wesley, 1992, vol. 1.
- [23] R. M. Haralick, "Edge and region analysis for digital image data," *Comput. Vis., Graph., Image Process.*, vol. 12, pp. 60–73, 1980.
- [24] —, "The digital edge," in *Proc. IEEE Computer Soc. Conf. Pattern Recognition and Image Processing*, New York, 1981, pp. 285–294.
- [25] —, "Zero-crossing of second directional derivative edge operator," in *Proc. SPIE Symp. Robotic Vision*, Washington, DC, 1982, p. 23.
- [26] —, "Ridges and valleys on digital images," *Comput. Vis., Graph., Image Process.*, vol. 22, pp. 28–38, 1983.
- [27] R. M. Haralick, L. T. Watson, and T. J. Laffey, "The topographic primal sketch," *Int. J. Robot. Res.*, vol. 2, pp. 50–72, 1983.
- [28] R. M. Haralick, "Digital step edges from zero crossing of the second directional derivatives," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. PAMI-5, no. 1, pp. 58–68, Jan. 1984.
- [29] —, "Cubic facet model edge detector and ridge-valley detector: implementation details," in *Pattern Recognition in Practice II*, E. S. Gelsema and L. N. Kanal, Eds. Amsterdam, The Netherlands: Elsevier, 1986, pp. 81–90.
- [30] O. A. Zuniga and R. M. Haralick, "Integrated directional derivative gradient operator," *IEEE Trans. Syst., Man, Cybern.*, vol. SMC-17, no. 3, pp. 508–517, May/June 1987.
- [31] G. E. Forsythe, "Generation and use of orthogonal polynomials for data-fitting with a digital computer," *J. Soc. Ind. Appl. Math.*, vol. 5, pp. 74–88, 1957.
- [32] D. Nadadur, "Noise Covariance Estimation in Low-level Computer Vision," Ph.D. dissertation, Dept. Elect. Eng., Univ. Washington, Seattle, Nov. 2001.
- [33] A. P. Dawid, "Some matrix-variate distribution theory: notational considerations and a bayesian application," *Biometrika*, vol. 68, no. 1, pp. 265–274, 1981.
- [34] M. S. Bartlett, "On the theory of statistical regression," in *Proc. Roy. Stat. Soc.*, vol. 53, Edinburgh, U.K., 1933, pp. 260–283.
- [35] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *J. Roy. Stat. Soc. B*, vol. 39, no. 1, pp. 1–38, Dec. 1976.
- [36] N. D. Le, L. Sun, and J. V. Zidek, "Bayesian spatial interpolation and backcasting using Gaussian-generalized inverted Wishart model," Univ. British Columbia, Vancouver, BC, Canada, Tech. Rep., 1999.
- [37] B. M. Golam Kibria, "Multivariate Bayesian spatial interpolation using Gaussian-generalized inverted Wishart model," Univ. British Columbia, Vancouver, BC, Canada, Tech. Rep., 2000.
- [38] C. F. Chen, "Bayesian inference for a normal dispersion matrix and its application to stochastic multiple regression analysis," *J. Roy. Stat. Soc. B*, no. 41, pp. 235–248, 1979.
- [39] G. J. McLachlan and T. Krishnan, *The EM Algorithm and Extensions*. New York: Wiley, 1997.
- [40] M. S. Bazaraa, H. D. Sherali, and C. M. Shetty, *Nonlinear Programming: Theory and Algorithms*, 2nd ed. New York: Wiley, 1993.
- [41] A. V. Fiacco and G. P. McCormick, *Nonlinear Programming: Sequential Unconstrained Minimization Techniques*. New York: Wiley, 1968.

- [42] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, *Numerical Recipes in C: The Art of Scientific Computing*, 2nd ed. Cambridge, U.K.: Cambridge Univ. Press, 1992.
- [43] J. J. Moré and D. J. Thuente, Line search algorithms with guaranteed sufficient decrease, *Math. Comput. Sci. Division, Argonne Nat. Lab.*, Preprint MCS-P330-1092, Oct. 1992.
- [44] —, “Line search algorithms with guaranteed sufficient decrease,” *ACM Trans. Math. Software*, vol. 20, no. 3, pp. 286–307, 1994.
- [45] T. W. Anderson, *An Introduction to Multivariate Statistical Analysis*, 2nd ed. New York: Wiley, 1984.
- [46] T. Kanungo and R. M. Haralick, “Multivariate hypothesis testing for Gaussian samples,” *Intell. Syst. Lab., Dept. Elect. Eng., Univ. Washington, Seattle, Tech. Rep. ISL-TR-95-05*, Oct. 5, 1995.
- [47] D. Nadadur, R. M. Haralick, and D. E. Gustafson, Integrated second directional derivative ridge operator, to be published.
- [48] A. M. López, F. Lumbreras, J. Serrat, and J. J. Villanueva, “Evaluation of methods for ridge and valley detection,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 21, no. 4, pp. 327–335, Apr. 1999.
- [49] J. Wishart, “The generalized product moment distribution in samples from a normal multivariate population,” *Biometrika*, vol. 20, no. A, pp. 32–52, 1928.
- [50] —, “Proofs of the distribution law of the second order moment statistics,” *Biometrika*, vol. 35, pp. 55–57, 1948.
- [51] J. Wishart and M. S. Bartlett, “The generalized product moment distribution in a normal system,” in *Proc. Cambridge Philosophical Soc.*, vol. 29, 1933, pp. 260–270.
- [52] D. Nadadur and R. M. Haralick, “A Bayesian framework for noise covariance estimation using the facet model,” Siemens Medical Syst., Inc., Ultrasound Group (SMS-UG) and *Intell. Syst. Lab. (ISL)*, Issaquah, WA, Tech. Rep. TR (DN) 2001-0001, 2001.



Desikachari Nadadur (M’00) was born in Cuddapah, Andhra Pradesh, India, in 1970. He received the B.Tech. degree in electronics and communications engineering from Sri Venkateswara University, Tirupati, Andhra Pradesh, in 1990, the M.S. degree in electrical engineering from Gonzaga University, Spokane, WA, in 1993, and the Ph.D. degree in electrical engineering, with a specialization in image processing and computer vision, from the University of Washington (UW), Seattle, in 2001, under the supervision of Prof. R. M. Haralick. His dissertation

was titled “Noise Covariance Estimation in Low-Level Computer Vision.”

Since 1999, he has been a Senior Scientist in the Advanced Imaging Applications Group of the Developing Competency Department, Siemens Medical Solutions, Inc., Ultrasound Division, Issaquah, WA. He was a Teaching Assistant with the Department of Electrical Engineering, Gonzaga University, Spokane, from 1991 to 1993. He was a predoctoral Research Associate II at the Intelligent Systems Laboratory, UW, from 1993 to 1997. He was a Software Engineer in the Advanced Technologies Group at Electronics for Imaging, Inc., Foster City, CA, from 1997 to 1999. He has coauthored several publications in the fields of image processing, mathematical morphology, and computer vision, and holds several patents. His current research interests include application of deformable, finite-element models, PDE and level-set methods to medical image segmentation, Bayesian and other statistical methods in image analysis, and performance characterization of vision algorithms.



Robert Martin Haralick (F’84) is a Distinguished Professor in the Department of Computer Science, Graduate Center, City University of New York. He was the Boeing Clairmont Egtvedt Professor in the Department of Electrical Engineering, University of Washington, Seattle. He is responsible for developing the grayscale co-occurrence texture analysis technique and the facet model technique for the image processing. He has worked on robust methods for photogrammetry and developed fast algorithms for solving the consistent labeling problems in

high-level computer vision. He has developed shape analysis and extraction techniques using mathematical morphology, the morphological sampling theorem, and fast recursive morphology algorithms. Together with Prof. I. Phillips, he has developed a comprehensive ground-truthed set of some 1600 document image pages (most in English) and some 200 pages in Japanese in the area of document image understanding. He has also developed algorithms for document skew angle estimation, zone delineation, and word and text line bounding box delineation. His most recent research is in the area of computer vision performance characterization and covariance propagation. He has published more than 490 papers.

Dr. Haralick was made a Fellow of IEEE for his contributions in computer vision and image processing and a Fellow of the International Association for Pattern Recognition (IAPR) for his contributions in pattern recognition, image processing, and for services to IAPR. He was also the President of the IAPR.



David Earl Gustafson received the B.S. degree in physics from Hamline University, St. Paul, MN, and the Ph.D. degree in physics from the University of Virginia, Charlottesville.

He was a Nuclear Physics Postdoctoral Fellow, Florida State University, Gainesville, and a Postdoctoral Research Fellow at the Mayo Clinic, Rochester, MN. He is, at present, a Senior Director at Siemens Medical Solutions USA, Inc., Ultrasound Division, Issaquah, WA. He has extensive experience, including medical imaging research and development

in both university academic and industrial settings. five years of experience in academia, including teaching in radiology residency and technologist programs, and 24 years in the medical imaging industry. His industry positions, which included group to department managerial roles for 20 years, have focused on conducting and directing product research and development, with a strong focus on extensive strategic advanced development carrying product features from inception to product release. His experience has also included manufacturing, marketing, and business management responsibility and management of multi-site development teams, widely diversified geographically, including extensive R & D staff consisting of engineers and scientists, most with advanced degrees, and product development marketing personnel. His development projects have covered most modern medical imaging methods. Broad imaging and image processing product development experience includes: computed tomography, digital X-ray fluoroscopy, radiography and angiography products, mobile X-ray systems, image-intensified special procedures radiology and cardiology systems, mammography, ultrasound, PACS and connectivity, CAD, and nuclear medicine products.