

AN ABSTRACT OF THE THESIS OF

GARY JOE SEXTON for the degree DOCTOR OF PHILOSOPHY
(Name) (Degree)
in STATISTICS presented on February 21, 1975
(Major Department) (Date)

Title: MULTIVARIATE DENSITY ESTIMATION--A BAYESIAN
APPROACH Redacted for privacy

Abstract approved: _____
Dr. H. D. Brunk

Methods of estimating a multivariate density over \mathbb{R}^m are investigated from a Bayesian point of view.

The motivation for this thesis comes from the classification problem where the goal is to estimate the probability of class C_i , $i = 1, \dots, N$ given an observed vector. Initial transformations of the data are introduced so that the components of the transformed data vector \underline{U} are uncorrelated. In this case the density of \underline{U} is close to multivariate normal. Under the assumption that both the density $p(u)$ and $\log p(u)$ have expansions in terms of an orthonormal basis, estimates are derived for the coefficients in the expansions. The approximation of the density is in terms of the best L_2 approximation.

When $p(u)$ is a multivariate normal density, estimating the leading coefficients in the expansion of $p(u)$ is equivalent to

estimating the covariance matrix. Estimating the corresponding coefficients in the expansion of $\log p(u)$ is equivalent to estimating the inverse of the covariance matrix. In this case, estimates of the covariance matrix are proposed. The "goodness" of these estimates is measured by a quadratic loss function.

Multivariate Density Estimation--A
Bayesian Approach

by

Gary Joe Sexton

A THESIS

submitted to

Oregon State University

in partial fulfillment of
the requirements for the
degree of

Doctor of Philosophy

June 1975

APPROVED:

Redacted for privacy

Professor of Statistics

in charge of major

Redacted for privacy

Chairman of Department of Statistics

Redacted for privacy

Dean of Graduate School

Date thesis is presented February 21, 1975

Typed by Clover Redfern for Gary Joe Sexton

ACKNOWLEDGMENT

The author wishes to thank the many members of the Department of Statistics with whom he has had valuable course work and discussions. This would include Dr. Mark Lembersky and Dr. Charles Land.

The author wishes to express a special note of appreciation to Dr. Justus Seely and Dr. David Birkes for many stimulating and rewarding conversations during the time this thesis was being written.

Deepest appreciation is expressed to Professor H. D. Brunk who served as the author's major professor. Professor Brunk originally suggested the problem area to the author and gave guidance during its investigation. In addition to technical guidance, Professor Brunk gave abundantly of his time, encouragement, and above all, patience, without which this thesis would never have been completed.

The author is grateful for the financial support provided by the Department of Statistics through teaching assistantships and an N. I. H. Biometry traineeship.

Final thanks go to the author's wife, Mary Lou, and their children, Craig and Stacey, for their patience, understanding and moral support.

TABLE OF CONTENTS

<u>Chapter</u>	<u>Page</u>
I. INTRODUCTION	1
II. THE CLASSIFICATION PROBLEM--SPECIFICS	7
2.1 Initialization of the Data	7
2.2 Expansion of Densities	17
2.3 Orthonormal Families, Fourier Expansions and Hermite Polynomials	18
2.4 The Model Specifics	21
2.5 Details of the Expansion	25
III. DISTRIBUTION THEORY	32
3.1 The Wishart Distribution	32
3.2 Properties of the Wishart Distribution	34
3.3 The Inverted Wishart Distribution	45
3.4 Properties of the Inverted Wishart Distribution	46
IV. ESTIMATION OF PARAMETERS	50
4.1 A Natural Conjugate Family of Priors	50
4.2 An Estimator for β_r	58
4.3 An Estimator for α_r	65
4.4 Properties of Estimators	84
BIBLIOGRAPHY	101
APPENDIX	104

MULTIVARIATE DENSITY ESTIMATION--A BAYESIAN APPROACH

I. INTRODUCTION

We initially consider the following two problems of practical importance. Simply stated they are:

- (i) An investigator observes some grid, say $k \times k$ with k^2 lattice points. At each lattice point the investigator observes a number. This number could be an index number that is an index of greyness as in the data recorded by the NASA ERTS satellite and in a sense represents an aerial photograph. The investigator would like to identify the object of the photograph.
- (ii) The second problem corresponds to the situation where several measurements are taken on an individual as in the case of diagnosing a disease.

In both of the above situations we are observing a random vector and the objective is to determine to which of several classes or populations the object of measurement belongs. In the first case we might be interested in differentiating urban areas from rural areas or differentiating between grain fields and grass fields for example. In the second case we might be interested in classifying a patient as to the type of disease he has.

In both of the above examples we assume that we are dealing with several classes which are denoted by C_1, C_2, \dots, C_N . Having observed a random vector $X^T = (X_1, \dots, X_k)$, the problem of fundamental interest is to determine to which of the N classes the object represented by X^T belongs. We would also like to give odds that realistically reflect the accuracy of our classification scheme. That is, we would like to know the probability of class C_i given the observed vector X which is denoted $P(C_i|X)$ $i = 1, 2, \dots, N$. If we know $P(C_i|X)$ for each i the the "posterior odds" for class C_i as opposed to class C_j will be $P(C_i|X)/P(C_j|X)$.

In the broadest sense it is the classification problem given in (i) and (ii) that motivated this thesis.

The determination of the quantity $P(C_i|X)$ may be approximated by a scalar multiple of the quantity $P(X|C_i)$. The emphasis of this thesis is the estimation of $P(X|C_i)$ where it is assumed that $P(X|C_i)$ is a multivariate continuous density for $i = 1, \dots, N$.

Chapter II contains a discussion of certain initial transformations of the data vector that lead to convenient methods of reducing the number of components in the data vector. We also discuss the expansion of densities in terms of multidimensional Hermite polynomials, which form a complete orthonormal family on the m -dimensional Euclidean space \mathbb{R}^m . We assume that the unknown

density $p(u)$ has an expansion given by

$$p(u) = p_0(u) \sum_r \beta_r H_r(u)$$

where $p_0(u)$ is the multidimensional weight function for the Hermite polynomials $H_r(u)$. We also assume that $\log p(u)$ has a similar expansion so that

$$p(u) = p_0(u) e^{\sum_r a_r H_r(u)}.$$

The details of these expansions are given in terms of certain model specifics.

Chapter III contains a discussion of Wishart and inverted Wishart distributions including properties that are important to the estimation problem of Chapter IV.

Chapter IV contains a discussion of estimating Fourier type coefficients of the expansions developed in Chapter II. The general approach is Bayesian in nature. We initially derive an estimator for the coefficients β_r using very general considerations. We also consider the problem of estimating the coefficients a_r in the expansion of $\log p(u)$. The problem of estimating a_r when p is multivariate normal is equivalent to estimating the precision matrix in a multivariate normal density while estimating the corresponding

β_r is equivalent to estimating the covariance matrix. In this case estimators for the covariance and precision matrices are derived.

The general problem of estimating densities by way of orthogonal expansions seems to be of fairly recent origin.

The earliest work known to the author was a paper by Čencov (1962) where he considered the Fourier type expansion of a univariate density in terms of an orthonormal basis. An unbiased estimate of the Fourier coefficients and properties of the corresponding density estimate are derived in terms of an L_2 -norm.

Following the paper by Čencov, the orthogonal function approach was pursued by Schwartz (1967), Kronmal and Tarter (1968), Watson (1969) and Crain (1974). The emphasis in these papers is the estimation of univariate densities.

Schwartz addressed himself to the expansion of a univariate density in terms of Hermite functions and used an unbiased estimate of the coefficients. Under certain conditions he shows that his estimate is consistent with respect to mean squared error,

$$E[f(x) - f_n(x)]^2, \quad \text{and mean integrated squared error,}$$

$$E \int [f(x) - f_n(x)]^2 dx .$$

Kronmal and Tartar used standard Fourier series and derived trigonometric estimates along with their properties.

Watson used general orthogonal series to derive an estimate involving a weighting factor dependent upon sample size where the

weighting factor is chosen to minimize $E \int [f(x) - f_n(x)]^2 dx$.

Crain studied extensively the estimation of univariate densities and distribution functions when both the density and the log of the density have expansions in terms of orthonormal functions over compact sets. In particular, normalized Legendre polynomials were considered.

It seems to the author that very little has been done by way of estimating multivariate densities by way of orthogonal functions from a Bayesian viewpoint.

One notable paper is that of Brunk and Pierce (1974). They develop estimators for a random vector with dichotomous components by way of an orthonormal basis. They assume that $p(x)$ and $\log p(x)$ both have expansions. The coefficients in each of these expansions is related through the notion of "near independence" of the components of X . The notion of "near independence" is represented through a prior distribution on the coefficients in the expansions and resulting estimates of these coefficients are derived.

An interesting paper on the estimation of a covariance matrix in a multivariate normal distribution using a Bayesian approach is found in Stein, Efron and Morris (1972). Stein considers estimators of the form $[aS^{-1} + (b/\text{tr } S)I]^{-1}$ where S is a matrix that is a function of the data. Stein shows that his estimate strictly dominates the maximum likelihood estimate, which is in the above class, when

the "loss" function is given by

$$L(\hat{\Phi}, \Phi) = \frac{\text{tr}(\hat{\Phi}^{-1} - \Phi^{-1})S(\hat{\Phi}^{-1} - \Phi^{-1})}{k \text{tr} \Phi^{-1}}.$$

This "loss" function arises in a natural way from considerations found in a paper by Efron and Morris (1972).

We now introduce some of the notation used in this thesis.

Other notation is defined as it is introduced.

We use the symbol $:=$ to mean the quantity on the left of the symbol is defined by the quantity on the right of the symbol. We use $\langle \cdot, \cdot \rangle$ to denote inner product. The symbol $(U|X) \sim$ will mean "U given X is distributed as" while the symbol $\hat{\sim}$ will mean "is approximately distributed as." We use the generic symbol Φ to denote a covariance matrix with σ_{ij} and σ^{ij} denoting the ij elements of Φ and Φ^{-1} respectively. We use the symbol $|\cdot|$ to denote the determinant of a matrix. T as a superscript will denote the transpose. The symbols $N_m(\mu, \Phi)$ and $W_m(\nu, \Phi)$ are used to denote the m-dimensional normal and Wishart distributions with mean μ and covariance Φ in the normal case and with ν denoting the "degrees of freedom" and Φ denoting the parametric matrix in the Wishart case. The term "precision" will refer to the reciprocal of the variance in the case of a univariate distribution and the inverse of the covariance matrix in the case of a multivariate distribution.

II. THE CLASSIFICATION PROBLEM--SPECIFICS

2.1 Initialization of the Data

If $\mathbf{X}^T = (X_1, \dots, X_n)$ represents a vector of observations then we know that the posterior probabilities for the various classes satisfy the relationship

$$P(C_i | \mathbf{X}) \propto P(C_i)P(\mathbf{X} | C_i).$$

We will assume that there is sufficient information to give reliable estimates of the probabilities $P(C_i)$, $i = 1, \dots, N$. Under this assumption we will be able to concentrate our attention on $P(\mathbf{X} | C_i)$.

We also assume that population i has a covariance structure given by Σ_i , mean μ_i and the probability of class i occurring is p_i with $p_i > 0$.

The problem of estimating the mean has been studied extensively and many properties as well as the accuracy are known about the estimates. With this in mind, we will initially concentrate on the problem of estimating other structure acting as though the means were known.

Prior to discussing the estimation problem we will discuss representation of the data.

If the random vector \mathbf{X} has k components and k is large, it will be desirable in many situations to reduce the number of

components. In order to perform this reduction in some optimal manner it will be convenient to initially transform the vector \underline{X} .

Let us first suppose that the data is centered in the usual way by the transformation $\underline{Y} = \underline{X} - \mu_i$ for each of the N classes.

The formula $P_{\underline{X}}(X|C_i) = P_{\underline{Y}}(Y+\mu_i|C_i)$ will allow us to concentrate on the centered data \underline{Y} and relate the results to the original problem.

The joint covariance structure for the N classes can be expressed in terms of the marginal covariance structure by way of the following formula.

$$\begin{aligned}
 K(t, s) &:= E(\underline{Y}_s \underline{Y}_t) - E(\underline{Y}_s)E(\underline{Y}_t) \\
 &= E\{E(\underline{Y}_s \underline{Y}_t | C_i)\} - E\{E(\underline{Y}_s | C_i)\}E\{E(\underline{Y}_t | C_i)\} \\
 (2.1) \quad &= E\{\text{cov}(\underline{Y}_s, \underline{Y}_t | C_i) + E(\underline{Y}_s | C_i)E(\underline{Y}_t | C_i)\} \\
 &\quad - E\{E(\underline{Y}_s | C_i)\}E\{E(\underline{Y}_t | C_i)\} \\
 &= E\{\text{cov}(\underline{Y}_s, \underline{Y}_t | C_i)\} \\
 &= \sum_{i=1}^N E(\underline{Y}_s \underline{Y}_t | C_i).
 \end{aligned}$$

Let Φ denote the joint covariance matrix of the N classes.

We can represent \underline{Y} without error by an expansion of the form $\underline{Y} = \psi \underline{V}$ where ψ is a deterministic $k \times k$ matrix of

rank k and \underline{Y} is a random vector. From this exact expansion we may write \underline{V} as $\underline{V} = \psi^{-1} \underline{Y}$ so that \underline{V} is simply a transformation of \underline{Y} .

For the general classification problem, it will be desirable, if we are going to transform the data, to preserve the original structure as much as possible. In particular, it is especially desirable to preserve distance. With this in mind we state the following known result from linear theory.

Lemma 2.1.1. If $\psi: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a linear transformation then $\|\mathbf{x}\| = \|\psi\mathbf{x}\|$ if and only if ψ is orthogonal.

Proof: See Smith (1971). \square

In our situation the primary implication of Lemma 2.1.1 is that we shall restrict ourselves to those transformations ψ which are orthonormal. That is, those transformations ψ for which $\psi^T \psi = I$. Additionally, by restricting our attention to orthogonal transformations, the distributions are related in the following simple manner:

$$\begin{aligned} F_{\underline{V}}(\mathbf{v}) &= P(\underline{V} \leq \mathbf{v}) \\ &= P(\psi^T \underline{Y} \leq \mathbf{v}) \\ &= P(\underline{Y} \leq \psi^T \mathbf{v}) \\ &= F_{\underline{Y}}(\psi^T \mathbf{v}). \end{aligned}$$

Under this restriction we have \underline{Y} represented without error by the transformation $\underline{Y} = \psi \underline{V}$ where ψ is orthonormal and in this case $\underline{V} = \psi^T \underline{Y}$ and \underline{V} is simply an orthonormal transformation of \underline{Y} .

It should be noted that the covariance structure of \underline{V} is given by:

$$\text{cov } \underline{V} = \text{cov}(\psi^T \underline{Y}) = \psi^T (\text{cov } \underline{Y}) \psi = \psi^T \Phi \psi$$

From among the many orthonormal transformations we single out the Karhunen-Loeve transformation which will be denoted by κ . In matrix form, the columns of κ will consist of the k eigenvectors of Φ with the i th column, κ_i , being the eigenvector associated with the eigenvalue λ_i . Moreover, the eigenvalues have been ordered so that they satisfy $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$.

In using the transformation κ we now have that $\text{cov } \underline{V} = \kappa^T \Phi \kappa = D$ where D is the diagonal matrix with $D_{ii} = \lambda_i$, $i = 1, 2, \dots, k$. By rotating (or reflecting) \underline{V} in this manner we have arrived at a vector \underline{V} that has the property that the components of \underline{V} are uncorrelated and in this sense we say they are nearly independent.

We make special note of the Karhunen-Loève transformation in as much as it has certain optimal properties as we shall see. Moreover, it is the transformation that we will use.

In order to justify our choice of the Karhunen-Loeve transformation we return to the situation where ψ is an arbitrary orthonormal transformation so that $\underline{Y} = \psi \underline{V}$ and $\underline{V} = \psi^T \underline{Y}$.

In the language of pattern recognition, the components of \underline{Y} are called features or indicants representing the vector \underline{Y} . As commented earlier, if the number of features, k , is large, it will be desirable in certain situations to reduce the number of features used to represent the vector \underline{Y} . That is, suppose we only determine m of the components of \underline{Y} . Moreover, we may assume that it is the first m components that are determined.

If we replace each of the indicants \underline{Y}_i with a constant b_i , $i = m+1, \dots, k$, we will have an approximate representation of \underline{Y} . In introducing this approximation we wish to do so in a manner that will minimize the mean square error in the representation.

For notation, set $\underline{Y}_1 := (v_1, \dots, v_m)^T$, $\underline{Y}_2 := (v_{m+1}, \dots, v_k)^T$, $\underline{b} := (b_{m+1}, \dots, b_k)^T$, $\underline{Y}_b = \psi \begin{pmatrix} \underline{Y}_1 \\ \underline{b} \end{pmatrix} = \sum_{i=1}^m \psi_i \underline{Y}_i + \sum_{i=m+1}^k \psi_i b_i$ where ψ_i is the i th column of ψ . Also, suppose we partition ψ as $(\psi_1 \parallel \psi_2)$ where ψ_2 is the last $k-m$ columns of ψ .

We now have the following relationships among the partitioned matrices:

$$\underline{V} = \begin{pmatrix} \underline{V}_1 \\ \underline{V}_2 \end{pmatrix} = \begin{pmatrix} \psi_1^T \\ \psi_2^T \end{pmatrix} \underline{Y}$$

$$\text{cov } \underline{y} = \begin{pmatrix} \psi_1^T \\ \psi_2^T \end{pmatrix} \Phi(\psi_1, \psi_2) = \begin{pmatrix} \psi_1^T \Phi \psi_1 & \psi_1^T \Phi \psi_2 \\ \psi_2^T \Phi \psi_1 & \psi_2^T \Phi \psi_2 \end{pmatrix}$$

$$\underline{v}_2 = \begin{pmatrix} 0 \\ I_{k-m} \end{pmatrix} \begin{pmatrix} \underline{y}_1 \\ \underline{y}_2 \end{pmatrix} = \begin{pmatrix} 0 \\ I_{k-m} \end{pmatrix} \begin{pmatrix} \psi_1^T \\ \psi_2^T \end{pmatrix} \underline{y} = \psi_2^T \underline{y}$$

$$\text{cov } \underline{v}_2 = \psi_2^T \Phi \psi_2$$

$$\psi_2^T \psi_2 = I_{k-m}$$

The following known result gives us the proper choice of b_i , $i = m+1, \dots, k$ to minimize the mean square error in the representation.

Lemma 2.1.2. Let ψ be an arbitrary orthogonal transformation and define \underline{y} by $\underline{y} := \psi^T \underline{Y}$ so that $\underline{Y} = \psi \underline{y}$. Set $\underline{y}_b = \psi \begin{pmatrix} \underline{y}_1 \\ b \end{pmatrix}$ where b is a $(k-m) \times 1$ vector of constants so that \underline{y}_b approximates \underline{y} . The choice of b that minimizes the mean square error in the approximation is $b = E \underline{v}_2$.

Proof: The mean square error in the approximation is given by $E \|\underline{y} - \underline{y}_b\|^2$ where $\|\cdot\|$ denotes the usual Euclidean norm. If we set $\underline{\xi}_\psi := \underline{y} - \underline{y}_b$ we then have the following:

$$\begin{aligned} E\|\underline{Y}-\underline{Y}_b\|^2 &= E\|\underline{\xi}_\psi\|^2 = E\ \underline{\xi}_\psi^T \underline{\xi}_\psi = \text{tr } E\ \underline{\xi}_\psi \underline{\xi}_\psi^T \\ &= \sum_{i=m+1}^k E(\underline{Y}_i - b_i)^2. \end{aligned}$$

Each term in the sum is minimum for $b_i = E\underline{Y}_i$. \square

We should note that in our case $E\underline{Y}_i = 0$ so we set $b_i = 0$, $i = m+1, \dots, k$. We note that $\underline{\xi}_\psi = \psi \begin{pmatrix} 0 \\ \underline{V}_2 - \underline{b} \end{pmatrix}$ and in general, when $\underline{b} = E\underline{V}_2$ so that $\underline{\xi}_\psi \begin{pmatrix} 0 \\ \underline{V}_2 - E\underline{V}_2 \end{pmatrix} = \psi_2(\underline{V}_2 - E\underline{V}_2)$, we have:

$$\begin{aligned} E\|\underline{\xi}_\psi\|^2 &= \sum_{i=m+1}^k E(\underline{V}_i - E\underline{V}_i)^2 = \sum_{i=m+1}^k \text{var } \underline{V}_i \\ &= \text{tr } E\ \underline{\xi}_\psi \underline{\xi}_\psi^T \\ &= \text{tr } \psi_2 \{E(\underline{V}_2 - E\underline{V}_2)(\underline{V}_2 - E\underline{V}_2)^T\} \psi_2^T \\ &= \text{tr } \psi_2 (\text{cov } \underline{V}_2) \psi_2^T \\ &= \text{tr } \text{cov } \underline{V}_2 \\ &= \text{tr } \psi_2^T \Phi \psi_2. \end{aligned}$$

Having chosen the vector \underline{b} to minimize $E\|\underline{\xi}_\psi\|^2$ for any fixed orthogonal transformation ψ the next problem is to select an optimal transformation. The selection of an optimal transformation will rest upon the following result attributed to Poincaré.

For notation, if Q is any matrix we will denote the i th ordered eigenvalue of Q by $\lambda_i(Q)$ with $\lambda_i(Q) \geq \lambda_{i+1}(Q)$.

Lemma 2.1.3. Let A be a $k \times k$ symmetric matrix with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$ its eigenvalues. Let B be a $k \times q$ matrix such that $B^T B = I_q$. It then follows that

$$\lambda_i(B^T A B) \leq \lambda_i(A), \quad i = 1, \dots, q,$$

and

$$\lambda_{q-j}(B^T A B) \geq \lambda_{k-j}(A), \quad j = 0, 1, \dots, q-1$$

Proof: See Rao (1965). \square

The following known result will follow immediately from Lemma 2.1.3.

Under the hypothesis of Lemma 2.1.2, if we set $b_i = EV_i$ for $i = m+1, \dots, k$ then the problem

$$\underset{\psi}{\text{minimize}} E \|\xi_{\psi}\|^2$$

has a solution given by $\psi = \kappa$ where κ is the Karhunen-Loève transformation as described on page 10.

Proof: As noted above, for any orthogonal transformation ψ ,

$$E \|\underline{\xi}_\psi\|^2 = \text{tr } \psi_2^T \Phi \psi_2.$$

It, therefore, suffices to show that

$$\text{tr } \psi_2^T \Phi \psi_2 \geq \text{tr } \kappa_2^T \Phi \kappa_2.$$

Note that ψ_2 is $k \times (k-m)$ and $\psi_2^T \psi_2 = I_{k-m}$ so that we may invoke Lemma 2.1.3 to conclude that

$$\lambda_{k-m-j}(\psi_2^T \Phi \psi_2) \geq \lambda_{k-j}(\Phi) \geq 0 \quad j = 0, 1, \dots, k-m-1.$$

Therefore,

$$\sum_{j=0}^{k-m-1} \lambda_{k-m-j}(\psi_2^T \Phi \psi_2) \geq \sum_{j=0}^{k-m-1} \lambda_{k-j}(\Phi).$$

Now, $\lambda_{k-j}(\Phi) = \lambda_{k-j}$ so that

$$\sum_{j=0}^{k-m-1} \lambda_{k-j}(\Phi) = \sum_{j=0}^{k-m-1} \lambda_{k-j} = \sum_{i=m+1}^k \lambda_i = \text{tr } \kappa_2^T \Phi \kappa_2.$$

That is, the sum of the eigenvalues, hence the trace, of $\psi_2^T \Phi \psi_2$ is bounded below by $\text{tr } \kappa_2^T \Phi \kappa_2$ and the result is established. \square

Before proceeding we present here a brief summary of the initialization process. The primary function of the initialization process is to obtain a convenient coordinate system with which to

represent the data. We have chosen to use the Karhunen-Loève coordinate system. To use this system we require the knowledge of Φ , the overall covariance structure of the N populations. There are two basic viewpoints that we have considered to obtain an estimate of Φ .

In the first case we assume that we have some reliable estimate of p_i for $i = 1, 2, \dots, N$. We also have a training set selected at random from each of the N populations. From each of these training sets we estimate μ_i in the event μ_i is unknown. We also estimate Φ_i from the sample training set corresponding to population i . We then obtain an estimate of Φ using (2.1).

In the second case we consider one sample from the combined training sets. After obtaining the overall sample we are able to identify to which population each observation belongs. We now estimate Φ from the overall sample. Moreover, we are able to directly obtain estimates of p_i by taking p_i to be the proportion of times that we observe an observation from population i .

Inasmuch as the second case allows us to perform all estimations from the one overall sample we have focused on the second viewpoint in this thesis.

Since the fundamental purpose of our estimate of Φ is to establish an initial coordinate system, the accuracy of this estimate does not seem to be critical. This was found to be the case in an

analogous situation by Martin and Bradley (1972).

With this in mind we are content to use the usual estimate of Φ . That is, we set $\overline{x_i x_j} := \frac{1}{n} \sum_{t=1}^n x_i^t x_j^t$, $\overline{x_i} := \frac{1}{n} \sum_{t=1}^n x_i^t$ and estimate σ_i by $\overline{x_i x_i} - \overline{x_i}^2$. In matrix form,

$$\hat{\Phi} = \frac{1}{n} \sum_{k=1}^n k_X k_X^T - \overline{X} \overline{X}^T.$$

Given that we have an estimate of Φ we find the Karhunen-Loève transformation κ whose columns will be the eigenvectors of Φ arranged in decreasing order by the magnitudes of the associated eigenvalues.

Now suppose we observe a random vector \underline{X} from the population at large. For fixed i , we center this observation and set $\underline{Y} = \underline{X} - \mu_i$. We now transform \underline{Y} by κ^T to obtain a transformed observation \underline{V} where $\underline{V} = \kappa^T \underline{Y}$. We now have an observation \underline{V} and the components \underline{V} are uncorrelated.

2.2 Expansion of Densities

From this point on, we shall assume that all of the previous initializations have been performed. We will now denote \underline{Y}_b as simply \underline{Y} and assume \underline{Y} is an $m \times 1$ random vector, κ is the $k \times m$ orthonormal transformation consisting of the eigenvectors

of \mathbb{R}^k . $\underline{Y} = \kappa^T \underline{X}$ so that the feature vector \underline{Y} has its components uncorrelated.

We shall concentrate on estimating the density $p(V|C_i)$ for fixed i which we will denote simply by $p(v)$. The estimation procedure will be to represent $p(v)$ by an expansion of Fourier type in terms of a complete orthonormal family and estimating the coefficients in the expansion.

We shall outline the general procedure in the univariate case and then proceed to the details of our model.

2.3 Orthonormal Families, Fourier Expansions and Hermite Polynomials

In general, let us assume that P is a probability measure on the Borel subsets \mathcal{B}^1 of \mathbb{R}^1 and P is absolutely continuous with respect to some measure ℓ so that its density is given by

$p = \frac{dP}{d\ell}$. Let P_0 be any measure on $(\mathbb{R}^1, \mathcal{B}^1)$ that is absolutely

continuous with respect to the measure ℓ , so that its density is

$$P_0 = \frac{dP_0}{d\ell}.$$

Consider the Hilbert space, $L_2(\mathbb{R}^1, \mathcal{B}^1, P_0)$ with the usual inner product and norm given by $\langle f, g \rangle := \int fg dP_0$ and

$\|f\|^2 = \int |f|^2 dP_0$. Let $\{T_s : s \in S\}$ be a family of functions from

\mathbb{R}^1 to \mathbb{R}^1 that is orthonormal with respect to P_0 so that

$\langle T_s, T_r \rangle = \delta_{r,s}$ where $\delta_{r,s}$ is the Kronecker delta. Moreover,

if it happens that $\{T_s : s \in S\}$ is complete, that is, $\langle f, T_s \rangle = 0$ for every $s \in S$ implies $f = 0$, then any function $f \in L_2(P_0)$ has an expansion of the form

$$f = \sum_s \langle f, T_s \rangle T_s$$

If we restrict our attention to densities p such that $\frac{p}{P_0}$ and $\log \frac{p}{P_0}$ are members of $L_2(P_0)$ then we have representations for p given by

$$p(v) = p_0(v) \sum_s \beta_s T_s(v)$$

and

$$p(v) = p_0(v) e^{\sum_s \alpha_s T_s(v)}$$

Regarding the terms β_s and α_s as parameters, the problem of estimating p is one of estimating the parameters in at least one of the above expansions.

A class of functions that we will find useful is the class of Hermite polynomials. Let h be the weighting function on \mathbb{R}^1 given by $h(v) = e^{-v^2/2}$ and set $g(v) = \frac{1}{\sqrt{2\pi}} h(v)$ so that $g(v)$ may be regarded as the density of a standardized normal variate.

The Hermite polynomials are defined in terms of the above functions. For $n = 0$, set $H_0(v) := 1$. For $n = 1, 2, 3, \dots$, set $H_n(v) := \frac{(-1)^n}{\sqrt{n!}} e^{v^2/2} h^{(n)}(v)$ where $h^{(n)}(v) := \frac{d^n}{dv^n} h(v)$. $H_n(v)$ is

called the n th order normalized Hermite polynomial. We should note here that we are considering the space $L_2(\mathbb{R}^1, \mathcal{B}^1, P_0)$ where P_0 is determined by $g = \frac{dP_0}{d\ell}$ and ℓ is Lebesgue measure.

Let m and n be integers and assume $m \leq n$. We then have:

$$\begin{aligned} \langle H_m, H_n \rangle &= \int H_m(v) H_n(v) g(v) dv = \frac{1}{\sqrt{2\pi}} \int e^{-v^2/2} H_m(v) H_n(v) dv \\ &= \frac{(-1)^n}{\sqrt{n!}} \int H_m(v) g^{(n)}(v) dv \end{aligned}$$

Integrating by parts m times gives:

$$\langle H_m, H_n \rangle = \frac{(-1)^{n+m}}{\sqrt{n!}} \int H_m^{(m)}(v) g^{(n-m)}(v) dv$$

If $m < n$, then $n-m > 0$ and one more integration by parts will give $\langle H_m, H_n \rangle = 0$ since H_m is a polynomial of degree m and we will be taking the $(m+1)$ st derivative.

If $m = n$ then $m-n = 0$ and we have

$$\begin{aligned} \langle H_m, H_n \rangle &= \int [H_n(v)]^2 g(v) dv \\ &= \frac{1}{\sqrt{n!}} \int H_n^{(n)}(v) g(v) dv \end{aligned}$$

$H_n(v)$ is a polynomial of degree n with leading coefficient being $\frac{1}{\sqrt{n!}}$. In this case $H_n^{(n)} = \frac{n!}{\sqrt{n!}} = \sqrt{n!}$ and we have

$$\langle H_m, H_n \rangle = \int g(v) dv = 1.$$

From this argument we can see that $\langle H_m, H_n \rangle = \delta_{m,n}$ for $m, n = 0, 1, 2, \dots$ and the normalized Hermite polynomials are orthonormal on \mathbb{R}^1 with respect to the density $g(v) = \frac{1}{\sqrt{2\pi}} e^{-v^2/2}$. It is also true that $\{H_n : n = 0, 1, \dots\}$ is complete with respect to g (see Burrill (1972)).

The first few normalized polynomials are:

$$\begin{aligned} H_0(v) &= 1 & H_3(v) &= \frac{1}{\sqrt{6}} (v^3 - 3v) \\ H_1(v) &= v & H_4(v) &= \frac{1}{\sqrt{24}} (v^4 - 6v^2 + 3) \\ H_2(v) &= \frac{1}{\sqrt{2}} (v^2 - 1) & H_5(v) &= \frac{1}{\sqrt{120}} (v^5 - 10v^3 + 5v) \end{aligned}$$

2.4 The Model Specifics

In constructing this model we assume that we are dealing with a random vector \underline{Y} defined on \mathbb{R}^m and that P is some probability measure on $(\mathbb{R}^m, \mathcal{B}^m)$ associated with \underline{Y} . Additionally, we assume that P is absolutely continuous with respect to some measure ℓ on $(\mathbb{R}^m, \mathcal{B}^m)$ so that the joint density function

is given by $p = \frac{dP}{d\ell}$.

Consider the normalized Hermite polynomials described previously which form a complete orthonormal family on \mathbb{R}^1 with respect to the density $\frac{1}{\sqrt{2\pi}} e^{-v^2/2}$ when ℓ is Lebesgue measure. The following known theorem will allow us to find a suitable orthonormal family that is complete on \mathbb{R}^m .

Lemma 2.3.1. Suppose $\{X_m\}$ and $\{Y_n\}$ are complete orthonormal families with respect to probability measures μ and ν . It then follows that the product set $\{X_m Y_n\}$ is a complete orthonormal family on the product space with respect to the product measure $\mu \times \nu$.

Proof: See Lancaster (1969) or Zygmund (1959). \square

By applying this theorem and letting $\{X_m\} = \{Y_m\} = \{H_m\}$ we have, extending by induction, that all finite products consisting of m terms of Hermite polynomials forms a complete orthonormal system on \mathbb{R}^m with respect to the measure P_0 whose density is given by

$$P_0(v) = \frac{dP_0}{d(\ell_1 \times \dots \times \ell_n)} = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\lambda_i}} e^{-v_i^2/2\lambda_i}.$$

Let $B = \{r = (r_1, \dots, r_m) : r_i = 0, 1, 2, \dots, i = 1, \dots, m\}$. For $r \in B$, a non zero entry in the j th position will indicate the

presence of a Hermite polynomial whose argument is the j th feature and the value of r_j is the degree of the polynomial. With this indexing scheme we can denote the complete orthonormal family consisting of m -products of Hermite polynomials. Writing $u_i = v_i / \sqrt{\lambda_i}$ we have that $\{H_r(u) : r \in B\}$ is a complete orthonormal family on $(\mathbb{R}^m, \mathcal{B}^m, P_0)$ where P_0 has a density with respect to the measure ℓ on \mathbb{R}^m given by

$$p_0(u) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}} e^{-u_i^2/2}$$

We assume that $\frac{p(u)}{p_0(u)}$ and $\log \frac{p(u)}{p_0(u)}$ have expansions given

by:

$$(2.2) \quad p(u) = p_0(u) \sum_r \beta_r H_r(u)$$

$$(2.3) \quad p(u) = p_0(u) e^{\sum_r \alpha_r H_r(u)}$$

A statistician, when faced with a real problem, will formulate a model that he believes most accurately reflects the true situation. In attempting the analysis he may encounter many mathematical and/or computational difficulties that are not easily, if ever, overcome. He may reject his model on the grounds that it is of little practical value if he cannot carry out the analysis. He then may reformulate the problem from an entirely different viewpoint perhaps

encountering the same or similar problems. Alternatively, he may hold with his initial viewpoint but adopt simplifying assumptions in the model to allow for the analysis. In either case, his final model is a formulation of the problem that most accurately reflects the true situation subject to the condition that the analysis is possible.

When faced with the problem of estimating multivariate densities two simplifying assumptions come immediately to mind.

(i) The components of the vector \underline{U} are independent.

(ii) The vector \underline{U} is multivariate normal.

Both of these assumptions are implicit in the use of P_0 as a distribution for \underline{U} and we will later see that assumption (ii) may be reflected in (2.3) by setting certain coefficients equal to zero.

Before proceeding to the details of the expansion, a final comment is in order concerning expansions (2.2) and (2.3). Both expansions involve expressions of the form $\sum_{\mathbf{r}} \beta_{\mathbf{r}} H_{\mathbf{r}}(\mathbf{u})$ where \mathbf{r} is an m -component vector. There are several ways to interpret this sum.

Perhaps the most natural way is to write

$$\sum_{\mathbf{r}} \beta_{\mathbf{r}} H_{\mathbf{r}}(\mathbf{u}) := \sum_{r_1=0}^{\infty} \cdots \sum_{r_m=0}^{\infty} \beta_{(r_1, \dots, r_m)} H_{(r_1, \dots, r_m)}(\mathbf{u}).$$

For our purposes it will be more convenient to interpret the

sum in the following manner. Set

$$n(\mathbf{r}) := \sum_{i=1}^m r_i$$

for $\mathbf{r} \in B$ and we will call $n(\mathbf{r})$ the rank of \mathbf{r} . We will also say that the coefficients of $H_{\mathbf{r}}(u)$ in (2.1) and (2.2) have rank $n(\mathbf{r})$.

We now think of our countable collection of basis functions as being grouped by blocks according to increasing rank with an arbitrary ordering within blocks. That is, block zero contains $H_{\mathbf{r}}(u)$ for which $n(\mathbf{r}) = 0$, block one contains $H_{\mathbf{r}}(u)$ for which $n(\mathbf{r}) = 1$ with arbitrary ordering within the block and so on. We now have the functions $H_{\mathbf{r}}(u)$ ordered by rank but arbitrary within rank with the result that the coefficients will be ordered similarly.

We now write:

$$\sum_{\mathbf{r}} \beta_{\mathbf{r}} H_{\mathbf{r}}(u) := \sum_{\nu=0}^{\infty} \sum_{\mathbf{r}:n(\mathbf{r})=\nu} \beta_{\mathbf{r}} H_{\mathbf{r}}(u)$$

2.5 Details of the Expansion

We develop our expansion in a manner analogous to that in Brunk and Pierce (1974).

Suppose we denote $(0, 0, \dots, 0) \in B$ as simple 0 and then

from (2.3),

$$\begin{aligned} 1 &= \int p(u) du \\ &= \int p_0(u) e^{\sum_{r \neq 0} a_r H_r(u)} du \end{aligned}$$

from which we can write

$$(2.4) \quad e^{-a_0} = \int p_0(u) e^{\sum_{r \neq 0} a_r H_r(u)} du.$$

Let α denote a vector whose components are the coefficients a_s for which $s \neq 0$. Where α is specified then a_0 is determined by the requirement given in (2.4). Setting $A(\alpha) = -a_0$ we have defined a function $A(\cdot)$ whose domain is a set of vectors. We can now rewrite (2.3) as

$$(2.5) \quad p(u) = p_0(u) e^{\sum_{r \neq 0} a_r H_r(u) - A(\alpha)}.$$

Recalling that $H_0(u) = 1$ and that

$$\begin{aligned} E[H_r(u)H_s(u) | p_0] &= \langle H_r(u), H_s(u) \rangle \\ &= \delta_{r,s} \end{aligned}$$

and using the dominated convergence theorem in conjunction with (2.4)

we can write

$$\begin{aligned}
\frac{\partial}{\partial a_s} e^{A(a)} &= \frac{\partial A(a)}{\partial a_s} e^{A(a)} \\
&= \frac{\partial}{\partial a_s} \int p_0(u) e^{\sum_{r \neq 0} a_r H_r(u)} du \\
&= \int \frac{\partial}{\partial a_s} p_0(u) e^{\sum_{r \neq 0} a_r H_r(u)} du \\
&= \int p_0(u) e^{\sum_{r \neq 0} a_r H_r(u)} H_s(u) du
\end{aligned}$$

so that

$$(2.6) \quad \frac{\partial}{\partial a_s} A(a) = \int p_0(u) e^{\sum_{r \neq 0} a_r H_r(u) - A(a)} H_s(u) du .$$

Using (2.3),

$$\begin{aligned}
\frac{\partial}{\partial a_s} A(a) &= \int H_s(u) p(u) du \\
(2.7) \quad &= E[H_s(u) | p] , \quad \text{provided } s \neq 0 .
\end{aligned}$$

By writing

$$\int H_s(u) p(u) du = \int H_s(u) \frac{p(u)}{p_0(u)} p_0(u) du$$

and using (2.2), we have

$$\begin{aligned}
(2.8) \quad \frac{\partial A(a)}{\partial a_s} &= \int H_s(u) \sum_r \beta_r H_r(u) p_0(u) du \\
&= \sum_r \beta_r \int H_s(u) H_r(u) p_0(u) du =
\end{aligned}$$

$$\begin{aligned}
&= \sum_r \beta_r E[H_s(u)H_r(u)|p_0(u)] \\
&= \sum_r \beta_r \delta_{r,s} \\
&= \beta_s \quad \text{provided } s \neq 0.
\end{aligned}$$

Similarly, finding the second partials of $A(a)$ gives

$$\begin{aligned}
(2.9) \quad e^{A(a)} \left[\frac{\partial^2 A(a)}{\partial a_s \partial a_t} + \frac{\partial A(a)}{\partial a_s} \frac{\partial A(a)}{\partial a_t} \right] \\
= \int p_0(u) e^{\sum_{r \neq 0} a_r H_r(u)} H_s(u) H_t(u) du
\end{aligned}$$

so that

$$(2.10) \quad \frac{\partial^2 A(a)}{\partial a_s \partial a_t} + \frac{\partial A(a)}{\partial a_s} \frac{\partial A(a)}{\partial a_t} = E[H_s(u)H_t(u)|p]$$

and

$$(2.11) \quad \frac{\partial^2 A(a)}{\partial a_s \partial a_t} = \text{cov}[H_s(u), H_t(u)|p].$$

From (2.3) and (2.5),

a. If $a = 0$ then $A(a) = A(0) = 0$

b. If $s \neq 0$ then

$$\begin{aligned}
(2.12) \quad \frac{\partial A(0)}{\partial a_s} &= \int p_0(u) H_s(u) du = \int p_0(u) H_s(u) H_0(u) du \\
&= E[H_s(u)H_0(u)|p_0] \\
&= \delta_{s,0} \\
&= 0
\end{aligned}$$

c. If $s \neq 0$ then

$$\begin{aligned} \frac{\partial^2 A(0)}{\partial \alpha_s^2} &= E[H_s(u)H_s(u)|p_0] \\ &= 1 \end{aligned}$$

If $t \neq 0$ and $t \neq s$ then

$$\begin{aligned} \frac{\partial^2 A(0)}{\partial \alpha_s \partial \alpha_t} &= E[H_s(u)H_t(u)|p_0] \\ &= 0 \end{aligned}$$

That is,

$$(2.13) \quad \frac{\partial^2 A(0)}{\partial \alpha_s \partial \alpha_t} = \delta_{s,t} \quad \text{if } s \neq 0, t \neq 0.$$

As an aside we note here that by (2.7) and (2.8) the coefficients of rank one in (2.2) all vanish. That is, if $n(r) = 1$ then

$H_r(u) = H_1(u_i)$ for some i and

$$\begin{aligned} \beta_r &= E[H_r(u)|p] \\ &= E[H_1(u_i)|p] \\ &= E[u_i|p] \\ &= 0. \end{aligned}$$

Recall that in the general classification problem Φ represents the covariance structure of the mixture of all classes. In fact, Φ will be taken to be the usual estimate of the mixture covariance

in case the individual covariance structures are not known. Recall also that U is a transformation of our original vector so that the components are uncorrelated.

If we were initially to assume that Y is $N(0, \Sigma)$ then $P_0(u)$ would be the distribution resulting from the application of κ to Y . That is, $N(0, I)$ would be an initial estimate of the distribution of U .

With this in mind we will be interested in the situation where the components of U are nearly independent. If we interpret the parameters a_r as interactions and assume that interactions of high order are relatively small then we may interpret near independence in terms of a_r of rank larger than one being near zero.

In this context we will consider the Maclaurin expansion of the function $A(a)$.

$$\begin{aligned} A(a) &= A(0) + \sum_{r \neq 0} a_r^2 + \dots \\ &\doteq \sum_{r \neq 0} a_r^2 + \dots \end{aligned}$$

From (2.7) and (2.8) we have $\beta_s = \frac{\partial A(a)}{\partial a_s} = E[H_s(u) | p]$ and β is a function of a through p . We shall denote this fact by writing β as $\beta(a)$. Moreover,

$$\beta_0(a) = E[H_0(u) | p] = 1$$

and

$$\beta_s(0) = \frac{\partial A(0)}{\partial a_s} = 0 \quad \text{by (2.12).}$$

These two facts give $\beta_s(0) = \delta_{s,0}$.

On the other hand, from (2.8) we have

$$\frac{\partial \beta_s(a)}{\partial a_t} = \frac{\partial^2 A(a)}{\partial a_s \partial a_t} = \text{cov}[H_s(u), H_t(u) | p]$$

and

$$\frac{\partial \beta_s(0)}{\partial a_t} = \delta_{s,t} \quad \text{for } s \neq 0, t \neq 0.$$

If we consider the expansion of $\beta_s(a)$ about zero,

$$\begin{aligned} \beta_s(a) &= \beta_s(0) + a_s + \text{higher order terms} \\ &= a_s + \text{higher order terms} \end{aligned}$$

From this we can infer that when the a 's are near zero then

$$\beta_s \doteq a_s \quad \text{if } s \neq 0.$$

In view of the two basic equations (2.2) and (2.3), the estimation of the density p involves the estimation of the parameters β_r in (2.2) and/or the estimation of the parameters a_r in (2.3).

In either case our aim is to derive some suitable estimate of p .

III. DISTRIBUTION THEORY

In this chapter we will give some results in distribution theory that will be useful in the parameter estimation problem of Chapter IV. For our purposes, the distributions of most interest are distributions from families of Wishart and Inverted Wishart distributions.

Generally, the results of this chapter may be found in various texts such as Anderson (1958) and Press (1972). The proofs supplied are believed by the author to be his own and are presented here for the sake of completeness.

Selected results concerning characteristic functions and the differentiation of matrices may be found in the Appendix.

3.1 The Wishart Distribution

Historically, the derivation of the Wishart distribution originates from observing a sample of size ν from the distribution of an m -component random vector \underline{z} that is $N_m(0, \Phi)$ and considering the random matrix

$$\underline{Y} := \sum_{\alpha=1}^{\nu} \underline{z}_{\alpha} \underline{z}_{\alpha}^T.$$

The random matrix \underline{Y} is symmetric and may therefore be considered as a collection of $\frac{m(m+1)}{2}$ distinct random variables

V_{ij} , $i, j = 1, \dots, m$, $i \leq j$. In the event that $\nu > m-1$ and Φ is non singular, the joint distribution of the random variables \underline{V}_{ij} , $i, j = 1, \dots, m$, $i \leq j$ is absolutely continuous with respect to Lebesgue measure on $\mathbb{R}^{m(m+1)/2}$ and therefore has a density.

The value of the density function is given by

$$(3.1) \quad f(V|\nu, \Phi) := C |\Phi|^{-\nu/2} |V|^{-\nu-m-1/2} e^{-(1/2)\text{tr}(\Phi^{-1}V)} \mathcal{I}_{\mathcal{V}}(V)$$

where $\text{tr}(\cdot)$ denotes, as usual, the trace operator,

$$\Gamma_m\left(\frac{\nu}{2}\right) := \prod_{i=1}^m \Gamma\left(\frac{\nu+1-i}{2}\right)$$

$$C := \{2^{\nu m/2} \Gamma^{m(m-1)/4} \Gamma_m\left(\frac{\nu}{2}\right)\}^{-1}$$

and $\mathcal{I}_{\mathcal{V}}(\cdot)$ represents the indicator function of the class of positive definite matrices \mathcal{V} .

In case \underline{V} has a density given as above, we say that \underline{V} has a non-singular Wishart distribution with ν degrees of freedom and parametric matrix Φ . We will denote this fact by simply writing $\underline{V} \sim W(\nu, \Phi)$. In many instances it is convenient to set $\tau := \Phi^{-1}$ and in this case we call τ the precision matrix of the variate \underline{V} .

For a complete discussion and derivation of the above density one can see for example Anderson (1958).

In specifying a set of values v_{ij} for the random variables \underline{y}_{ij} , we may regard these values as a vector in $\mathbb{R}^{m(m+1)/2}$. Equivalently, we may think of these values as characterizing an $m \times m$ symmetric matrix V . In either case we will refer to (3.1) as the density function.

At this point we note two additional properties of the Wishart distribution that are apparent from the definition given in (3.1). The first is that when $\nu > m-1$ and when Φ is non singular then the random matrix \underline{y} with density (3.1) is positive definite with probability one. The second is that while we initially interpret the random matrix \underline{y} as a finite sum, that is, ν is an integer, the function given by (3.1) will be a density for non integer ν , provided that $\nu > m-1$ and Φ is positive definite.

3.2 Properties of the Wishart Distribution

For notation in all that follows we let $\hat{\nu}$ denote $\nu - 1$.

Lemma 3.2.1. Suppose the random matrix \underline{y} has a Wishart distribution with ν degrees of freedom and parametric matrix Φ , then the joint characteristic function of the distinct random variables \underline{y}_{ij} , $i, j = 1, 2, \dots, m$, $i \leq j$ evaluated at the random vector t_{ij} , $1 \leq i \leq j \leq m$ is given by the formula

$$(3.2) \quad \phi_{\underline{Y}}(t) := |\Phi^{-1}|^{\nu/2} |\Phi^{-1} - it|^{-\nu/2}$$

where

$$t = \begin{bmatrix} 2t_{11} & t_{12} & \cdots & t_{1m} \\ t_{21} & 2t_{22} & \cdots & t_{2m} \\ \vdots & \vdots & & \vdots \\ t_{m1} & t_{m2} & \cdots & 2t_{mm} \end{bmatrix}.$$

Proof: See Kshirsagar (1972). \square

We note here that, consistent with the comment following example A.2.2, the characteristic function given in (3.2) can be viewed as the characteristic function of the random matrix \underline{Y} evaluated at a symmetric matrix t as in (3.2) with $\langle t, \underline{Y} \rangle = \frac{1}{2} \text{tr}(t\underline{Y})$. If we view \underline{Y} as a random matrix and use the inner product $\text{tr}(t\underline{Y})$ as in example A.2.2, we then have a characteristic function evaluated at any symmetric matrix t given by

$$(3.3) \quad \phi_{\underline{Y}}(t) = |\Phi^{-1}|^{\nu/2} |\Phi^{-1} - 2it|^{-\nu/2}.$$

For this viewpoint one can see Anderson (1958).

We should comment here that in interpreting the characteristic function of a Wishart distribution, (3.2) is the joint characteristic function of the random variables $\underline{Y}_{11}, \underline{Y}_{12}, \dots, \underline{Y}_{1m}, \underline{Y}_{22}, \underline{Y}_{23}, \dots, \underline{Y}_{2m}, \dots, \underline{Y}_{mm}$ while (3.3) is the joint characteristic function of the

variables $\underline{Y}_{11}, \underline{Y}_{22}, \dots, \underline{Y}_{33}, 2\underline{Y}_{12}, 2\underline{Y}_{13}, \dots, 2\underline{Y}_{m-1, m}$. In this thesis we have adopted the viewpoint consistent with (3.2) and use (3.2) as the characteristic function for a Wishart variate with ν degrees of freedom and precision Φ .

It should also be noted here that the proofs of this lemma that were cited initiate their discussion by assuming a sample of size ν from a $N_m(0, \Phi)$ distribution. An examination of the proofs reveals that the lemma is true when $\nu \leq m-1$ and integer even though the density exists only when $\nu > m-1$. In case $\nu < m-1$ we say that we have a singular Wishart distribution with ν degrees of freedom. We also note that Lukacs (1969) has shown that if $\nu > m-1$ and not necessarily integer then \underline{Y} has a moment generating function of the form $\frac{|\Phi^{-1}|^{\nu/2}}{|\Phi^{-1}-t|^{\nu/2}}$ and hence a characteristic function of the form given by the lemma. Olkin and Rubin (1962) have shown that $\frac{|\Phi^{-1}|^{\nu/2}}{|\Phi^{-1}-t|^{\nu/2}}$ is not the moment generating function of a probability distribution if $\nu < m-1$ and non integer.

Lemma 3.2.2. Suppose we have the transformation $\underline{Y} = \underline{C}\underline{A}\underline{C}^T$ where \underline{C} is some fixed non singular matrix. If $\underline{Y} \sim W_m(\nu, \Phi)$ then $\underline{A} \sim W_m(\nu, \underline{C}^{-1}\Phi(\underline{C}^T)^{-1})$.

Proof: If $\underline{Y} = \underline{C}\underline{A}\underline{C}^T$ then $\underline{A} = \underline{C}^{-1}\underline{Y}(\underline{C}^T)^{-1}$. Now apply Lemma A.2.5 from the Appendix to conclude that the characteristic function of \underline{A} is given by

$$\begin{aligned}
\phi_{\underline{A}}(t) &:= \phi_{\underline{Y}}[((C^T)^{-1}tC^{-1})^T] \\
&= \phi_{\underline{Y}}[(C^{-1})^T t C^{-1}] \\
&= \frac{|\underline{\Phi}^{-1}|^{\nu/2}}{|\underline{\Phi}^{-1} - i(C^{-1})^T t C^{-1}|^{\nu/2}} \\
&= \frac{|\underline{\Phi}^{-1}|^{\nu/2}}{|(C^{-1})^T (C^T \underline{\Phi}^{-1} C - i t) C^{-1}|^{\nu/2}} \\
&= \frac{|C^T \underline{\Phi}^{-1} C|^{\nu/2}}{|C^T \underline{\Phi}^{-1} C - i t|^{\nu/2}}.
\end{aligned}$$

This is the characteristic function of a Wishart variate with parameters ν and $(C^T \underline{\Phi}^{-1} C)^{-1} = C^{-1} \underline{\Phi} (C^T)^{-1}$. \square

This lemma may be used in the following way to derive conclusions concerning a Wishart random matrix. Suppose we are interested in some property \mathcal{P} regarding a Wishart variate $\underline{Y} \sim W_m(\nu, \underline{\Phi})$. In general we may verify \mathcal{P} for a variate $\underline{A} \sim W(\nu, I)$. Next set $\underline{\Phi} = CC^T$ where C is $m \times m$ and non-singular and $\underline{A} := C^{-1} \underline{Y} (C^T)^{-1}$ so that $\underline{Y} = C \underline{A} C^T$. Now by Lemma 3.2.2, $\underline{Y} \sim W(\nu, \underline{\Phi})$ implies $\underline{A} \sim W(\nu, I)$. In many instances, having verified \mathcal{P} for the variate \underline{A} we will be able to invoke simple properties of transformations to conclude the validity of \mathcal{P} for \underline{Y} . We will use this result in the proofs of several following lemmas.

We first prove several preliminary lemmas.

Lemma 3.2.3. If Φ and t are symmetric matrices with t as in (3.2) then $\frac{\partial |\bar{A}|}{\partial t_{ij}} = -2i |\bar{A}_{ij}|$ where $\bar{A} := \Phi^{-1} - it$ and $|\bar{A}_{ij}|$ is the ij cofactor of \bar{A} .

Proof: By A.1.1 (iii) of the Appendix,

$$\frac{\partial |\bar{A}|}{\partial t_{ij}} = \text{tr} \left[\bar{A}^{\#} \frac{\partial \bar{A}^{-T}}{\partial t_{ij}} \right].$$

By A.1.1 (i) and (iv),

$$\begin{aligned} \frac{\partial \bar{A}^{-T}}{\partial t_{ij}} &= \frac{\partial (\Phi^{-1} - it)}{\partial t_{ij}} \\ &= -i \frac{\partial t}{\partial t_{ij}} \\ &= \begin{cases} -i \Delta_{ij}^* & \text{if } i \neq j \\ -2i \Delta_{ii}^* & \text{if } i = j \end{cases} \end{aligned}$$

where Δ_{ij}^* is a matrix of zeroes except for ones as the ij and ji elements. Therefore,

$$\frac{\partial |\bar{A}|}{\partial t_{ij}} = \begin{cases} -i \text{tr}[\bar{A}^{\#} \Delta_{ij}^*] & \text{if } i \neq j \\ -2i \text{tr}[\bar{A}^{\#} \Delta_{ii}^*] & \text{if } i = j \end{cases}.$$

If $i \neq j$ then $\bar{A}^{\#} \Delta_{ij}^*$ will have as its i -th column the j -th column of $\bar{A}^{\#}$ and its j -th column will be the i -th column of

$\bar{A}^\#$ with all other entries being zero. Therefore,

$$\text{tr}[\bar{A}^\# \Delta_{ij}^*] = |\bar{A}_{ij}| + |\bar{A}_{ji}|$$

and

$$\frac{\partial |\bar{A}|}{\partial t_{ij}} = -i [|\bar{A}_{ij}| + |\bar{A}_{ji}|].$$

Since \bar{A} is symmetric, $|\bar{A}_{ij}| = |\bar{A}_{ji}|$ so that

$$\frac{\partial |\bar{A}|}{\partial t_{ij}} = -2i |\bar{A}_{ij}|.$$

If $i = j$ then $\bar{A}^\# \Delta_{ii}^*$ will have as its i -th column the i -th column of $\bar{A}^\#$ and

$$\text{tr}[\bar{A}^\# \Delta_{ii}^*] = |\bar{A}_{ii}|.$$

Therefore,

$$\frac{\partial |\bar{A}|}{\partial t_{ij}} = -2i |\bar{A}_{ii}|. \quad \square$$

Lemma 3.2.4. If \bar{A} is as in the previous lemma then

$$\frac{\partial |\bar{A}_{ij}|}{\partial t_{kl}} = |\bar{A}| (a^{-ik} a^{-jl} + a^{-il} a^{-jk}) - 2i \frac{|\bar{A}_{ij}|}{|\bar{A}|} |\bar{A}_{kl}|$$

where a^{gh} denotes the gh element of \bar{A}^{-1} .

Proof: We first note that for any non singular matrix \bar{M} ,

$$\frac{|\bar{M}_{ij}|}{|\bar{M}|} = m^{ji}.$$

Now,

$$\frac{\partial |\bar{A}_{ij}|}{\partial t_{kl}} = \frac{\partial}{\partial t_{kl}} |\bar{A}| \frac{|\bar{A}_{ij}|}{|\bar{A}|}.$$

By the chain rule,

$$\frac{\partial |\bar{A}_{ij}|}{\partial t_{kl}} = |\bar{A}| \frac{\partial}{\partial t_{kl}} \frac{|\bar{A}_{ij}|}{|\bar{A}|} + \frac{|\bar{A}_{ij}|}{|\bar{A}|} \frac{\partial}{\partial t_{kl}} |\bar{A}|.$$

From Lemma 3.2.3,

$$\frac{\partial |\bar{A}_{ij}|}{\partial t_{kl}} = |\bar{A}| \frac{\partial}{\partial t_{kl}} \frac{|\bar{A}_{ij}|}{|\bar{A}|} - 2i \frac{|\bar{A}_{ij}|}{|\bar{A}|} |\bar{A}_{kl}|.$$

From the initial remark,

$$\begin{aligned} \frac{\partial}{\partial t_{kl}} \frac{|\bar{A}_{ij}|}{|\bar{A}|} &= \frac{\partial}{\partial t_{kl}} \bar{a}^{ji} \\ &= \sum_{r,s} \frac{\partial \bar{a}^{ji}}{\partial \bar{a}_{rs}} \frac{\partial \bar{a}_{rs}}{\partial t_{kl}} \\ &= \frac{\partial \bar{a}^{ji}}{\partial \bar{a}_{kl}} \frac{\partial \bar{a}_{kl}}{\partial t_{kl}} \\ &= -i \frac{\partial \bar{a}^{ji}}{\partial \bar{a}_{kl}}. \end{aligned}$$

The expression $\frac{\partial \bar{a}^{-ji}}{\partial \bar{a}_{kl}}$ is simply the ji element of the matrix

$$\frac{\partial \bar{A}^{-1}}{\partial \bar{a}_{kl}} = -\bar{A}^{-1} \Delta_{kl}^* \bar{A}^{-1}$$

where as before Δ_{kl}^* is a matrix of zeroes except for ones at the kl and lk elements. A simple multiplication gives

$$(-\bar{A}^{-1} \Delta_{kl}^* \bar{A}^{-1})_{ji} = -(a^{-ik} a^{-jl} + a^{-kl} a^{-jk}).$$

Therefore,

$$\frac{\partial}{\partial t_{kl}} \frac{|\bar{A}_{ij}|}{|\bar{A}|} = (a^{-ik} a^{-jl} + a^{-il} a^{-jk})$$

and the result follows. \square

We are now ready to find some of the moments of the Wishart distribution.

Lemma 3.2.5. If $\underline{y} \in W(\nu, \Phi)$ then $E\underline{y} = \nu\Phi$.

Proof: If we are restricting ν to be integer then the result follows immediately from the characterization

$$\underline{y} = \sum_{\alpha=1}^{\nu} \underline{z}_{\alpha} \underline{z}_{\alpha}^T.$$

The result is not quite so apparent if we allow ν to be non integer.

Let us initially suppose that $\Phi = I$. Regarding the random variables \underline{Y}_{ij} , $1 \leq i \leq j \leq m$ as a vector in $\mathbb{R}^{m(m+1)/2}$ we have from Lemma A.2.3,

$$E \underline{Y}_{ij} = \frac{1}{i} \left. \frac{\partial \phi_{\underline{Y}}(t)}{\partial t_{ij}} \right|_{t=0}$$

Now,

$$\begin{aligned} \left. \frac{\partial \phi_{\underline{Y}}(t)}{\partial t_{ij}} \right|_{t=0} &= \left. \frac{\partial}{\partial t_{ij}} |I|^{\nu/2} |I - it|^{-\nu/2} \right|_{t=0} \\ &= \left. \left\{ -\frac{\nu}{2} |I - it|^{-\nu/2-1} \frac{\partial}{\partial t_{ij}} |I - it| \right\} \right|_{t=0}. \end{aligned}$$

By Lemma 3.2.3 with $\Phi = I$ we have

$$\begin{aligned} \left. \frac{\partial \phi_{\underline{Y}}(t)}{\partial t_{ij}} \right|_{t=0} &= \left. \left\{ i \nu |I - it|^{-\nu/2-1} |(I - it)_{ij}| \right\} \right|_{t=0} \\ &= i \nu \delta_{ij} \end{aligned}$$

and

$$E \underline{Y}_{ij} = \nu \delta_{ij}.$$

Therefore,

$$E \underline{Y} = \nu I.$$

For the general case, if $\underline{Y} \sim W(\nu, \Phi)$ set $\Phi = CC^T$ where C is $m \times m$ and non singular and set $\underline{A} := C^{-1} \underline{Y} (C^T)^{-1}$ so that

$\underline{Y} = \underline{C}\underline{A}\underline{C}^T$. By Lemma 3.2.2, $\underline{A} \sim W(\nu, I)$. We now have

$$\begin{aligned} E\underline{Y} &= E\underline{C}\underline{A}\underline{C}^T \\ &= \underline{C}(E\underline{A})\underline{C}^T \\ &= \underline{C}(\nu I)\underline{C}^T \\ &= \nu \underline{\Phi}. \quad \square \end{aligned}$$

Additional useful moments are given in the following lemma.

Lemma 3.2.6. If $\underline{Y} \in W_m(\nu, \underline{\Phi})$ then

- i) $E(\underline{Y}_{ij}) = \nu \sigma_{ij}$
- ii) $E(\underline{Y}_{ij}\underline{Y}_{kl}) = \nu^2 \sigma_{ij} \sigma_{kl} + \nu \sigma_{ik} \sigma_{jl} + \nu \sigma_{il} \sigma_{jk}$
- iii) $\text{cov}(\underline{Y}_{ij}, \underline{Y}_{kl}) = \nu(\sigma_{ik} \sigma_{jl} + \sigma_{il} \sigma_{jk})$
- iv) $\text{var}(\underline{Y}_{ij}) = \nu(\sigma_{ii} \sigma_{jj} + \sigma_{ij}^2)$.

Proof: Result (i) was proved in Lemma 3.2.5. (iii) follows from (i) and (ii) in the usual way by writing

$$\text{cov}(\underline{Y}_{ij}, \underline{Y}_{kl}) = E(\underline{Y}_{ij}\underline{Y}_{kl}) - (E\underline{Y}_{ij})(E\underline{Y}_{kl})$$

and (iv) follows from (iii) by the relationship

$$\text{var}(\underline{Y}_{ij}) = \text{cov}(\underline{Y}_{ij}, \underline{Y}_{ij}).$$

All that needs to be proved is (ii) and the proof of this is a direct computation using the characteristic function and results on matrix

differentiation found in the Appendix. We can write

$$(3.4) \quad \mathbb{E} \tilde{Y}_{ij} \tilde{Y}_{kl} = \frac{1}{i^2} \frac{\partial^2 \phi_Y(t)}{\partial t_{kl} \partial t_{ij}} \Big|_{t=0}$$

For notation, as before let $\bar{A} := \Phi^{-1} - it$ and $A := \Phi^{-1}$. Using Lemma 3.2.3 we can write

$$\frac{\partial \phi_Y(t)}{\partial t_{ij}} = i\nu |\Phi^{-1}|^{\nu/2} |\Phi^{-1} - it|^{-\nu/2-1} |\bar{A}_{ij}|.$$

Taking second partials with respect to t_{kl} gives

$$\begin{aligned} \frac{\partial^2 \phi_Y(t)}{\partial t_{kl} \partial t_{ij}} &= \frac{\partial}{\partial t_{kl}} \left\{ i\nu |\Phi^{-1}|^{\nu/2} |\Phi^{-1} - it|^{-\nu/2-1} |\bar{A}_{ij}| \right\} \\ &= i\nu |\Phi^{-1}|^{\nu/2} \left\{ -\left(\frac{\nu}{2}+1\right) |\bar{A}_{ij}| |\bar{A}|^{-\nu/2-2} \frac{\partial}{\partial t_{kl}} |\bar{A}| \right. \\ &\quad \left. + |\bar{A}|^{-\nu/2-1} \frac{\partial}{\partial t_{kl}} |\bar{A}_{ij}| \right\} \\ &= i\nu |\Phi^{-1}|^{\nu/2} \left\{ 2i \left(\frac{\nu}{2}+1\right) |\bar{A}_{ij}| |\bar{A}|^{-\nu/2-2} |\bar{A}_{kl}| \right. \\ &\quad \left. + |\bar{A}|^{-\nu/2-1} \frac{\partial}{\partial t_{kl}} |\bar{A}_{ij}| \right\} \end{aligned}$$

Now,

$$\frac{\partial^2 \phi_Y(t)}{\partial t_{kl} \partial t_{ij}} \Big|_{t=0} = i^2 (\nu^2 + 2\nu) \left\{ \frac{|\bar{A}_{ij}| |\bar{A}_{kl}|}{|\Phi^{-1}| |\Phi^{-1}|} \right\} \Big|_{t=0} + i\nu |\Phi| \left\{ \frac{\partial}{\partial t_{kl}} |\bar{A}_{ij}| \right\} \Big|_{t=0}.$$

When $t = 0$, $|\bar{A}_{ij}|$ is the (i, j) cofactor of Φ^{-1} so that

$$\left. \frac{\partial^2 \phi_{\underline{Y}}(t)}{\partial t_{kl} \partial t_{ij}} \right|_{t=0} = i^2 (\nu^2 + 2\nu) \frac{|\Phi_{ij}^{-1}| |\Phi_{kl}^{-1}|}{|\Phi^{-1}| |\Phi^{-1}|} + i\nu |\Phi| \left\{ \frac{\partial}{\partial t_{kl}} |\bar{A}_{ij}| \right\} \Big|_{t=0}.$$

Noting that $\frac{|\Phi_{ij}^{-1}|}{|\Phi^{-1}|}$ is the (i, j) element of Φ and using (3.4) we have

$$EV_{ij} V_{kl} = (\nu^2 + 2\nu) \sigma_{ij} \sigma_{kl} + \frac{1}{2} \nu |\Phi| \left\{ \frac{\partial}{\partial t_{kl}} |\bar{A}_{ij}| \right\}_{t=0}$$

By Lemma 3.2.4

$$\frac{\partial}{\partial t_{kl}} \frac{\partial}{\partial t_{kl}} \left\{ |\bar{A}_{ij}| \right\}_{t=0} = i |\Phi^{-1}| (\sigma_{ik} \sigma_{jl} + \sigma_{il} \sigma_{jk}) - 2i \frac{|\Phi_{ij}^{-1}|}{|\Phi^{-1}|} |\Phi_{kl}^{-1}|$$

so that

$$EV_{ij} V_{kl} = \nu^2 \sigma_{ji} \sigma_{lk} + 2\nu \sigma_{ij} \sigma_{kl} - 2\nu \frac{|\Phi_{ij}^{-1}|}{|\Phi^{-1}|} \frac{|\Phi_{kl}^{-1}|}{|\Phi^{-1}|} + \nu (\sigma_{ik} \sigma_{jl} + \sigma_{il} \sigma_{jk})$$

and the result is established. \square

3.3 The Inverted Wishart Distribution

Another distribution that will be of use is the distribution of a random matrix \underline{U} where $\underline{U} = \underline{V}^{-1}$ and $\underline{V} \sim W(\nu, \Phi)$.

Definition 3.3.1. An $m \times m$ random matrix \underline{U} is said to have an inverted Wishart distribution with parameters S and t

where S is a positive definite matrix and t is a real number in case \underline{U} has a density given by

$$f(\underline{U} | S, t) = \begin{cases} \frac{|\underline{U}|^{-t/2} |S|^{(t-m-1)/2} e^{-(1/2)\text{tr}(\underline{U}^{-1}S)}}{2^{(t-m-1)(m/2)} \Pi_m^{m(m-1)/2} \Gamma_m\left(\frac{t}{2}\right)} & \text{if } \underline{U} > 0 \\ & 2m < t \\ 0 & \text{otherwise} \end{cases}$$

When \underline{U} has an inverted Wishart distribution we will write $\underline{U} \sim W_m^{-1}(t, S)$.

The inverse Wishart that will be of interest to us arises when we have $\underline{Y} \sim W_m(\nu, \Phi)$ and set $\underline{U} = \underline{Y}^{-1}$. A direct computation shows that in this case $\underline{U} \sim W_m^{-1}(\nu+m+1, \Phi^{-1})$. A discussion of this density may be found in Press (1972).

3.4 Properties of the Inverted Wishart Distribution

Let us first recall that if a random variable \underline{X} is distributed as χ^2 with K degrees of freedom then $E\left[\frac{1}{\underline{X}}\right] = \frac{1}{K-2}$. Also, we should recall the following long known lemma.

Lemma 3.4.1. If $\underline{Y} \sim W_m(\nu, \Phi)$ and L is a fixed vector then

$$\frac{L^T \Phi^{-1} L}{L^T \underline{Y}^{-1} L} \sim \chi^2_{[\nu-(m-1)]}$$

Proof: See Rao (1965). \square

We now can prove the following useful result.

Lemma 3.4.2. If $\underline{Y} \sim W(\nu, \Phi)$ then

$$E\underline{Y}^{-1} = \frac{1}{\nu - m - 1} \Phi^{-1}$$

Proof: Let us first assume $\Phi = I$ so that $\underline{Y} \sim W_m(\nu, I)$ and let $L := e_i$ where e_i is a standard basis vector with a one in the i th component. By Lemma 3.4.1.

$$\frac{1}{\underline{U}_{ii}} \sim \chi^2_{[\nu - (m-1)]}$$

where $\underline{U} = \underline{Y}^{-1}$. It now follows that each diagonal element of \underline{U} has the same expectation. That is,

$$\begin{aligned} E\underline{U}_{ii} &= \frac{1}{\nu - (m-1) - 2} & i = 1, 2, \dots, m \\ &= \frac{1}{\nu - m - 1} & i = 1, 2, \dots, m. \end{aligned}$$

Now, let $e_{ij} := e_i + e_j$ and set $L := e_{ij}$. Again, by Lemma 3.4.1,

$$\frac{L^T \underline{U} L}{L^T \underline{U} L} \sim \chi^2_{[\nu - (m-1)]}$$

That is,

$$\frac{2}{\underline{U}_{ii} + \underline{U}_{jj} + 2\underline{U}_{ij}} \sim \chi^2_{(\nu - m + 1)} \quad \text{for } i, j = 1, 2, \dots, m, i \neq j.$$

Again taking the expectation of the reciprocal of a χ^2 variate we have

$$\frac{1}{2} E(U_{\sim ii} + U_{\sim jj} + 2U_{\sim ij}) = \frac{1}{\nu - m - 1}$$

or

$$EU_{\sim ii} + EU_{\sim jj} + 2EU_{\sim ij} = \frac{2}{\nu - m - 1}$$

so that

$$\frac{1}{\nu - m - 1} + \frac{1}{\nu - m - 1} + 2EU_{\sim ij} = \frac{2}{\nu - m - 1}$$

and

$$EU_{\sim ij} = 0$$

To summarize, when $\Phi = I$.

$$\begin{aligned} EV^{-1} &= EU \\ &= (EU_{\sim ij})_{m \times m} \\ &= \frac{1}{\nu - m - 1} I \end{aligned}$$

For the general case the result follows by setting $\Phi = CC^T$, C non singular, $\underline{A} = C^{-1}\underline{V}(C^T)^{-1}$ so that $\underline{V} = C\underline{A}C^T$, and applying Lemma 3.2.2 to conclude that $\underline{A} \sim W_m(\nu, I)$. Now,

$$\begin{aligned} \underline{V}^{-1} &= (C\underline{A}C^T)^{-1} \\ &= (C^T)^{-1}\underline{A}^{-1}C^{-1} \end{aligned}$$

and

$$\begin{aligned}
E\tilde{Y}^{-1} &= (C^T)^{-1} E(\tilde{A}^{-1}) C^{-1} \\
&= (C^T)^{-1} \left(\frac{1}{\nu-m-1} I \right) C^{-1} \\
&= \frac{1}{\nu-m-1} (CC^T)^{-1} \\
&= \frac{1}{\nu-m-1} \Phi^{-1}. \quad \square
\end{aligned}$$

Other moments for an inverted Wishart distribution are given in the following lemma.

Lemma 3.4.3. If $\tilde{Y} \sim W_m(\nu, \Phi)$ and $\tilde{U} = \tilde{Y}^{-1}$ so that $\tilde{U} \sim W_m^{-1}(\nu+m+1, \Phi^{-1})$ then

$$\text{i) } \text{cov}(\tilde{U}_{ij}, \tilde{U}_{kl}) = \frac{\left[\frac{2\sigma^{ij}\sigma^{kl}}{\nu-m-1} + \sigma^{ik}\sigma^{jl} + \sigma^{il}\sigma^{kj} \right]}{(\nu-m)(\nu-m-1)(\nu-m-3)}$$

$$\text{ii) } \text{var}(\tilde{U}_{ii}) = \frac{2(\sigma^{ii})^2}{(\nu-m-1)^2(\nu-m-3)}$$

$$\text{iii) } \text{var}(\tilde{U}_{ij}) = \frac{\sigma^{ii}\sigma^{jj} + \frac{\nu-m+1}{\nu-m-1}(\sigma^{ij})^2}{(\nu-m)(\nu-m-1)(\nu-m-3)}$$

Proof: See Press (1972). \square

IV. ESTIMATION OF PARAMETERS

4.1 A Natural Conjugate Family of Priors

Our model assumes, after suitable transformations, that we are observing a random vector U defined on $(\mathbb{R}^m, \mathcal{B}^m)$ having a multivariate density given by $p(u)$ with $EU_i = 0$. We also assume that the density $p(u)$ has expansions given by (2.2) and (2.3) both of which involve parameters to be estimated. In most applications we will be interested in truncating these series at some point. In examining expansion (2.3), we see that by truncating the series at some point we have a representation of $p(u)$ by a distribution that is a member of an exponential family.

With this in mind we shall initiate our discussion of the estimation problem with a general comment concerning exponential families of distributions and certain related priors.

Suppose we have an m -dimensional random vector \underline{U} whose distribution is in some exponential family. As an example let us assume that \underline{U} is multivariate normal. Furthermore, let us assume that the mean is known to be zero and we are interested in the problem of estimating the precision matrix R . In such a problem, one might be interested in using a Bayesian approach by introducing a prior distribution for the parameters to be estimated.

In general, a useful family of priors when considering a Bayesian approach to an estimation problem is a family of priors that is a "natural conjugate" family with respect to the likelihood function. We mean by this that both the prior and posterior distributions of the parameters involved are in the same family of distributions. Having the prior and the posterior distribution both in the same family has the advantage of mathematical tractability. A natural conjugate family also allows us, in many instances, to determine in an intuitive way the relative contributions of the prior distribution and the sample data to the posterior knowledge of the parameter. For example, suppose we consider a random sample of size n from a normal distribution with unknown mean μ_1 and known variance σ_1^2 . It is well known that the family of normal distributions is a natural conjugate family for this problem and that the posterior distribution of μ_1 is normal with mean and variance given by

$$\mu_2 = \frac{\frac{\bar{nx}}{\sigma_1^2} + \frac{\mu_0}{\sigma_0^2}}{\frac{n}{\sigma_1^2} + \frac{1}{\sigma_0^2}}$$

and

$$\sigma_2^2 = \frac{1}{\left(\frac{n}{\sigma_1^2} + \frac{1}{\sigma_0^2}\right)}$$

where μ_0 and σ_0^2 are the mean and variance of the prior distribution for μ_1 . From these relationships we may interpret the prior distribution for μ_1 as being equivalent to a prior sample of size $n_0 = \sigma_1^2 / \sigma_0^2$ from the normal distribution with unknown mean μ_1 and known variance σ_1^2 starting with a flat prior. That is, we suppose we start with no observations and an improper uniform prior distribution for μ_1 . We then suppose we have a sample of size n_0 and sample mean μ_0 from the given distribution. This "equivalent" sample produces a posterior distribution that is normal with mean μ_0 and variance σ_0^2 . When this equivalent sample is combined with the actual sample the composite sample yields the posterior distribution with mean μ_2 and variance σ_2^2 . For a more complete discussion of natural conjugate families the reader is referred to Raiffa and Schlaifer (1961) or Barnett (1973).

Suppose we now consider an exponential family of distributions determined by the vector θ of parameters. The density may be written as

$$f(u|\theta) = e^{\theta^T W(u) - \psi(\theta)}$$

where $W(u)$ is a vector dependent upon U and $\psi(\theta)$ is a normalizing constant.

The likelihood function for a sample of size n from the above distribution may be written as

$$\begin{aligned}
 (4.1) \quad f(1_u, 2_u, \dots, n_u | \theta) &= e^{\sum_{t=1}^n \theta^T W(t_u) - n\psi(\theta)} \\
 &= e^{n\theta^T \bar{w} - n\psi(\theta)}
 \end{aligned}$$

where

$$\bar{w} = \frac{1}{n} \sum_{t=1}^n W(t_u)$$

Now suppose that we have a prior distribution for θ that has a density given by

$$(4.2) \quad g(\theta) \propto e^{b^*[\theta^T a^* - \psi(\theta)]}$$

where b^* is some specified constant and a^* is some prescribed vector.

The posterior distribution of θ given the observations is given by

$$\begin{aligned}
 (4.3) \quad g(\theta | 1_u, 2_u, \dots, n_u) &\propto e^{\theta^T [n\bar{w} + b^*a^*] - (n+b^*)\psi(\theta)} \\
 &= e^{(n+b^*)[\theta^T (\frac{n\bar{w} + b^*a^*}{n+b^*}) - \psi(\theta)]} \\
 &= e^{b[\theta^T a - \psi(\theta)]}
 \end{aligned}$$

where b is the constant $n+b^*$ and a is the vector $\frac{n\bar{w} + b^*a^*}{n+b^*}$.

From Equations (4.2) and (4.3) we see that the family of distributions given by (4.2) forms a natural conjugate family of priors for the

exponential family of distributions.

One very useful result concerning families which are natural conjugate families for exponential distributions is the following result.

Lemma 4.1.1. Suppose the density of the random vector \underline{U} is given by $f(u|\theta) = e^{\theta^T W(u) - \psi(\theta)}$. If a prior distribution for θ is given by (4.2) so that the corresponding posterior distribution is given by (4.3) and we set $\psi_i(\theta) := \frac{\partial \psi(\theta)}{\partial \theta_i}$, then under mild regularity conditions $\psi_i(\theta) = E[W_i(u)|\theta]$ and $a_i = E\psi_i(\theta)$ ($a_i^* = E\psi_i(\theta)$) where the expectation is with respect to the posterior (prior) distribution of θ . Moreover, a_i (a_i^*) is the modal value of the posterior (prior) in the sense that if $\bar{\theta}$ is defined by $\psi_i(\bar{\theta}) = a_i$ then $\bar{\theta}$ is the mode of the posterior.

Proof: The first statement in the lemma may be found in Barndorff-Nielsen (1973). The result $a_i = E\psi_i(\theta)$ is due to Morgan (1969). The regularity conditions given by Morgan are given in the Appendix. \square

In the event that \underline{U} is an m -dimensional random vector that has a normal distribution with mean zero and precision R , we may write the density as

$$\begin{aligned}
(4.4) \quad f(u|R) &= (2\pi)^{-m/2} |R|^{1/2} e^{-\frac{1}{2} u^T R u} \\
&= (2\pi)^{-m/2} |R|^{1/2} e^{-\frac{1}{2} \sum_{ij} r_{ij} u_i u_j} \\
&= (2\pi)^{-m/2} e^{\sum_{i=1}^m r_{ii} w_{ii}(u) + \sum_{i < j} r_{ij} w_{ij}(u) - \psi(R)}
\end{aligned}$$

where

$$\begin{aligned}
w_{ii}(u) &:= -\frac{1}{2} u_i^2 \\
w_{ij}(u) &:= -u_i u_j \quad \text{for } i < j \\
\psi(R) &:= -\frac{1}{2} \log |R| .
\end{aligned}$$

The likelihood function for a sample of size n is given by

$$(4.5) \quad f(u_1, u_2, \dots, u_n | R) \propto e^{n \sum_{i=1}^m r_{ii} \bar{w}_{ii} + n \sum_{i < j} r_{ij} \bar{w}_{ij} - n \psi(R)}$$

where $\bar{w}_{ij} = \frac{1}{n} \sum_{t=1}^n w_{ij}(u_t)$. Now suppose that R is given a prior distribution that is Wishart with ν degrees of freedom and precision matrix τ . If $\nu > m-1$ the density of R is given by

$$\begin{aligned}
f(R|\nu, \tau) &\propto |R|^{\frac{\nu-m-1}{2}} e^{-\frac{1}{2} \text{tr}(\tau R)} \\
&= e^{-\frac{1}{2} \text{tr}(\tau R) + \frac{\nu-m-1}{2} \log |R|}
\end{aligned}$$

where we are assuming R and τ are positive definite. Since τ and R are symmetric we have

$$\begin{aligned} \text{tr}(\tau R) &= \sum_{i=1}^m \sum_{j=1}^m \tau_{ij} r_{ij} \\ &= \sum_{i=1}^m \tau_{ii} r_{ii} + 2 \sum_{i < j} \tau_{ij} r_{ij} . \end{aligned}$$

Therefore,

$$\begin{aligned} (4.6) \quad f(R | \nu, \tau) &\propto e^{-\frac{1}{2}[\sum_{i=1}^m \tau_{ii} r_{ii} + 2\sum_{i < j} \tau_{ij} r_{ij}]} + \frac{\nu - m - 1}{2} \log |R| \\ &= e^{(\nu - m - 1) \left[-\frac{1}{2(\nu - m - 1)} \sum_{i=1}^m \tau_{ii} r_{ii} - \sum_{i < j} \frac{\tau_{ij} r_{ij}}{(\nu - m - 1)} + \frac{1}{2} \log |R| \right]} \\ &= e^{b^* [\sum_{i=1}^m r_{ii} a_{ii}^* + \sum_{i < j} r_{ij} a_{ij}^* - \psi(R)]} \end{aligned}$$

where

$$b^* := \nu - m - 1$$

$$a_{ii}^* := -\frac{\tau_{ii}}{2(\nu - m - 1)}$$

$$a_{ij}^* := -\frac{\tau_{ij}}{\nu - m - 1} .$$

The posterior distribution of R given the observations is given by

$$\begin{aligned}
(4.7) \quad & f(R | \mathbf{l}_u, \dots, \mathbf{n}_u, \nu, \tau) \\
& \propto e^{\sum_{i=1}^m r_{ii} [n\bar{w}_{ii} + b^* a_{ii}^*] + \sum_{i < j} r_{ij} [n\bar{w}_{ij} + b^* a_{ij}^*] - (n+b^*)\psi(R)} \\
& = e^{(n+b^*) \left\{ \sum_{i=1}^m r_{ii} \left[\frac{n\bar{w}_{ii} + b^* a_{ii}^*}{n+b^*} \right] + \sum_{i < j} r_{ij} \left[\frac{n\bar{w}_{ij} + b^* a_{ij}^*}{n+b^*} \right] - \psi(R) \right\}} \\
& = e^{b \left[\sum_{i=1}^m r_{ii} a_{ii} + \sum_{i < j} r_{ij} a_{ij} - \psi(R) \right]}
\end{aligned}$$

where

$$\begin{aligned}
b & := n+b^* \\
a_{ii} & := \frac{n\bar{w}_{ii} + b^* a_{ii}^*}{n+b^*} \\
a_{ij} & := \frac{n\bar{w}_{ij} + b^* a_{ij}^*}{n+b^*}.
\end{aligned}$$

From (4.6) and (4.7) we note that the prior and posterior distributions for the precision matrix R are members of the same family. Therefore, we conclude that the family of Wishart distributions is a natural conjugate family for the precision matrix in a multivariate normal family with zero mean.

We note in passing that from this result and the relationship between a Wishart and an inverted Wishart distribution the family of inverted Wishart distributions is a natural conjugate family for the covariance matrix in the same problem.

4.2 An Estimator for β_r

In this section we concentrate on deriving an estimator $\hat{\beta}_r$ for the coefficient β_r in the expansion given by (2.2). We assume that we have a sample training set $\{^t\mathbf{U}, t=1, 2, \dots, n\}$ from one of our classes where \mathbf{U} is a vector of m components. We recall that $H_r(^t\mathbf{u})$ denotes the m -dimensional basis function evaluated at the observation $^t\mathbf{u}$. Set

$$\bar{H}_r := \frac{1}{n} \sum_{t=1}^n H_r(^t\mathbf{u})$$

so that if n is large, \bar{H}_r is approximately multivariate normal.

We now consider the likelihood function for α . (Recall that α was defined on page 26.) We also recall that if α_r is near zero for all r then $\beta_r \doteq \alpha_r$. The likelihood function is given by

$$(4.8) \quad L(\alpha) = \left\{ \prod_{t=1}^n p_0(^t\mathbf{u}) \right\} e^{n[\sum_{r \neq 0} \alpha_r \bar{H}_r - A(\alpha)]}$$

From the relationships

$$E[H_r(\mathbf{U}) | p] = \beta_r$$

and

$$\begin{aligned}\text{var}[H_r(U)|p] &= \text{cov}[H_r(U), H_r(U)|p] \\ &= \frac{\partial^2 A(\alpha)}{\partial \alpha_r^2}\end{aligned}$$

and the fact that when the α 's are near zero,

$$\frac{\partial^2 A(\alpha)}{\partial \alpha_r^2} \doteq 1$$

we can conclude that

$$(4.9) \quad E[\bar{H}_r | p] = \beta_r$$

and

$$(4.10) \quad \begin{aligned}\text{var}[\bar{H}_r | p] &\doteq \frac{1}{n} \sum_{t=1}^n \text{var}[H_r(t\tilde{U}) | p] \\ &\doteq \frac{1}{n}.\end{aligned}$$

From (4.8) we note that the statistics \bar{H}_r , $r \neq 0$, are sufficient for α .

In order to derive a possible estimator we can, as Brunk and Pierce (1974) point out in an analogous situation, represent the notion of near independence among indicants by way of a prior distribution on the α_r 's. As Brunk and Pierce (1974) also point out such a prior might be of the form

$$(4.11) \quad \pi(\alpha) \propto e^{-\frac{1}{2} \sum_{r \neq 0} c_r \alpha_r^2}$$

where c_r is the precision associated with α_r . That is, we assume a prior whereby the α_r are approximately independent and normally distributed with mean 0 and precision c_r where c_r increases with rank. As a possible estimator, suppose we consider the mode of the posterior distribution of α for such a prior.

The posterior mode will maximize $\log L(\alpha)\pi(\alpha)$ as well as $L(\alpha)\pi(\alpha)$. That is, it will maximize

$$n \left\{ \sum_{r \neq 0} \alpha_r \bar{H}_r - A(\alpha) \right\} - \frac{1}{2} \sum_{r \neq 0} c_r \alpha_r^2.$$

In the usual way, taking partials and using (2.8) we have

$$(4.12) \quad \frac{\partial}{\partial \alpha_r} \left[n \left\{ \sum_{r \neq 0} \alpha_r \bar{H}_r - A(\alpha) \right\} - \frac{1}{2} \sum_{r \neq 0} c_r \alpha_r^2 \right] = n[\bar{H}_r - \beta_r] - c_r \alpha_r = 0.$$

When $\beta_r \doteq \alpha_r$ then (4.12) reduces to

$$(4.13) \quad \hat{\beta}_r = \frac{n\bar{H}_r}{n+c_r}.$$

We first note that the estimator (4.13) of β_r is a function of the

sufficient statistic \bar{H}_r and the precision c_r of a_r . We next note that from (4.9),

$$\begin{aligned} E[\hat{\beta}_r | p] &= \frac{n}{n+c_r} E[\bar{H}_r | p] \\ &= \frac{n}{n+c_r} \beta_r \\ &= \left(1 - \frac{c_r}{n+c_r}\right) \beta_r . \end{aligned}$$

That is $\hat{\beta}_r$ is a biased estimator of β_r with negative bias.

Finally we note that the expected mean square error of the estimate

$\hat{\beta}_r$ is given by

$$E[(\hat{\beta}_r - \beta_r)^2 | p] = \text{var}[\hat{\beta}_r | p] + \beta_r^2 \left(\frac{n}{n+c_r} - 1\right)^2$$

and

$$\begin{aligned} \text{var}[\hat{\beta}_r | p] &= \text{var}\left[\left(1 - \frac{c_r}{n+c_r}\right) \bar{H}_r | p\right] \\ &= \left(1 - \frac{c_r}{n+c_r}\right)^2 \text{var}[\bar{H}_r | p] . \end{aligned}$$

Now, $\beta_r^2 \left(\frac{n}{n+c_r} - 1\right)^2 \rightarrow 0$ as $n \rightarrow \infty$ and since $\text{var}[\bar{H}_r | p] < \infty$, $\text{var}[\hat{\beta}_r | p] \rightarrow 0$ as $n \rightarrow \infty$ so that

$$E[(\hat{\beta}_r - \beta_r)^2 | p] \rightarrow 0 \text{ as } n \rightarrow \infty$$

and $\hat{\beta}_r$ is a consistent estimator for β_r .

If one is interested in using a prior such as (4.11) then it is necessary to specify the precisions c_r . Brunk and Pierce (1974) have indicated that one possible way of specifying precisions c_r for such a prior is the following. The investigator might feel that all a_r of the same rank have the same precision and precision increases with rank. For example,

$$(4.14) \quad c_r = \pi \lambda^{n(r)-1}.$$

This condition is equivalent to the condition that if $n(s) = n(r)+1$ then for $\epsilon > 0$ and $\lambda > 1$, $\Pr\{|a_s| < \epsilon\} = \Pr\{|a_r| < \epsilon\sqrt{\lambda}\}$. This statement in turn is equivalent to the condition $\Pr\{|\sqrt{c_s} a_s| < \epsilon\sqrt{c_s}\} = \Pr\{|\sqrt{c_r} a_r| < \epsilon\sqrt{\lambda c_r}\}$. Now, $\sqrt{c_s} a_s$ and $\sqrt{c_r} a_r$ are standardized normal variables so that this last statement is equivalent to asserting that $\epsilon\sqrt{c_s} = \epsilon\sqrt{\lambda c_r}$ which finally reduces to $c_s = \lambda c_r$.

With a prior as in (4.11) and the precisions as in (4.14), the investigator has reduced the problem of specifying a prior with many parameters to a problem of specifying a two parameter prior. That is, he must specify λ and π . If the investigator feels that he can do an adequate job of determining λ and π he might well consider using (4.13) with no further ado.

Suppose we now consider the situation where the determination

of λ and π is a difficult or costly process. That is, the investigator does not feel competent enough to specify λ and π . We consider the problem from a different viewpoint with the idea that we will let the data aid us in specifying the prior precisions.

Let us consider β_r as a parameter and \bar{H}_r as an observation. When n is large and when β_r and α_r are near zero then, as before, the Central Limit Theorem gives us the validity of (4.9) and (4.10) so that $(\bar{H}_r | \beta_r) \sim N(\beta_r, \frac{1}{n})$. Moreover, $\text{cov}[\bar{H}_r, \bar{H}_s | p]$ is approximately zero so that we may consider the collection of means $\bar{H}_r, r \in B$ as a collection of approximately independent random variables. Now consider a prior on the parameters β_r whereby they are independent normal variables with zero means. Moreover, assume that all β_r with the same rank have the same precision. Denote by φ_k the common variance of all β_r with $n(r) = k$. If the φ_k are all specified then $\beta_r | \bar{H}_r$ is distributed as a normal variate with mean $n\bar{H}_r / (n + (1/\varphi_k))$ and precision $n + \frac{1}{\varphi_k}$. Additionally, the marginal distribution of \bar{H}_r is normal with mean 0 and variance $\varphi_k + \frac{1}{n}$. In this case we have a "natural estimate" for β_r given by (4.13) with $\frac{1}{\varphi_k} = c_r$ whenever $n(r) = k$. We also have a "natural estimate" for $\frac{1}{n} + \varphi_k$ which is to let

$$(4.15) \quad \hat{D} = \frac{1}{\#\{r : n(r)=k\}} \sum_{r : n(r)=k} \bar{H}_r^2$$

where $\#\{r: n(r)=k\}$ denotes the number of indices r having rank k . With this estimate for $\frac{1}{n} + \varphi_k$, we let

$$(4.16) \quad \hat{\varphi}_k = \hat{D} - \frac{1}{n} .$$

At this point we have derived using Bayesian ideas and some rather general considerations our estimator for β_r under the assumption that the parameters β_r satisfy the condition that all coefficients of the same rank have the same precision. We see also that, all other considerations aside, the estimator given by (4.13) or (4.14) and (4.16) is extremely simple. We shall also see in Section 4.3 that a similar estimator may be derived without the Central Limit Theorem using a method of moments.

We have concentrated in this section on the estimation of the coefficients in (2.2). In principle, $\beta_r \doteq \alpha_r$ when α is near zero so that we can use

$$(4.17) \quad \hat{\alpha}_r = \frac{\overline{nH}_r}{n + \frac{1}{\varphi_k}}$$

as an estimate of α_s in (2.3). However, taking as an estimate of α_r the same estimate as for β_r could lead to a poor estimate for α_s . That is, if the α 's are "too far" from zero then it might not be true that $\alpha_r \doteq \beta_r$. With this in mind we shall next take up the

problem of estimating α_r independently of β_r .

4.3 An Estimator for α_r

Before discussing an estimator for α_r we shall briefly consider α_r and β_r for certain specific ranks in order to determine what these coefficients are measuring.

We note here that if $n(r) = 1$ then $\beta_r = 0$. Suppose we now consider for a moment only those r 's with rank two.

For notation we shall let i_i denote an entry of two in the i th component of r and i_j denote an entry of one in each of the i and j components of r . All other components will be zero since $n(r) = 2$.

If we consider our estimator for β_r given in (4.13) we see that it is a simple linear function of \bar{H}_r . If we restrict ourselves to those r 's with rank two we then have

$$\bar{H}_r = \frac{1}{n} \sum_{t=1}^n H_r(t_u)$$

with

$$\begin{aligned} (4.18) \quad \bar{H}_{i_i}(t_u) &= H_2(t_{u_i}) \\ &= \frac{1}{\sqrt{2}} (t_{u_i}^2 - 1) \end{aligned}$$

and

$$\begin{aligned}
 H_{ij}(t_u) &= H_i(t_{u_i})H_j(t_{u_j}) \\
 &= t_{u_i} t_{u_j}.
 \end{aligned}$$

The relationships in (4.18) may be rewritten as

$$(4.19) \quad \bar{H}_{ii} = \frac{1}{\sqrt{2}} (\sigma_i^2 - 1)$$

and

$$\bar{H}_{ij} = \sigma_{ij}$$

where

$$\sigma_i^2 := \frac{1}{n} \sum_{t=1}^n t_{u_i}^2$$

and

$$\sigma_{ij} := \frac{1}{n} \sum_{t=1}^n t_{u_i} t_{u_j}.$$

From these relationships we see that when the rank of r is two, $\hat{\beta}_r$ is a simple linear function of the sample covariances.

Suppose we now consider (2.3) and introduce a prior on the α 's so that all α 's of rank larger than two are zero. We note that in this case if $n(r) = 1$ then $\alpha_r = 0$. We now write (2.3) as

$$\begin{aligned}
 (4.20) \quad p(u) &= p_0(u) e^{\sum_{r \neq 0} \alpha_r H_r(u) - A(\alpha)} \\
 &= (2\pi)^{-\frac{m}{2}} e^{-\frac{1}{2} \sum_{i=1}^m u_i^2} e^{\sum_{r: n(r)=2} \alpha_r H_r(u) - A(\alpha)}.
 \end{aligned}$$

The exponent in (4.20) may be written as

$$\begin{aligned}
 & -\frac{1}{2} \sum_{i=1}^m u_i^2 + \sum_{n(r)=2} a_r H_r(u) \\
 & = \frac{1}{2} \sum_{i=1}^m u_i^2 + \frac{\sqrt{2}}{2} \sum_{i=1}^m a_{ii} (u_i^2 - 1) + a_{12} u_1 u_2 + \dots + a_{1m} u_1 u_m \\
 & \quad + a_{23} u_2 u_3 + \dots \\
 & = -\frac{1}{2} \sum_{i=1}^m (1 - \sqrt{2} a_{ii}) u_i^2 - \sum_{i=1}^m \frac{a_{ii}}{\sqrt{2}} + \sum_{i=1}^m \sum_{j=i+1}^m a_{ij} u_i u_j \\
 & = -\frac{1}{2} \sum_{i=1}^m (1 - \sqrt{2} a_{ii}) u_i^2 - 2 \sum_{i=1}^m \sum_{j=i+1}^m a_{ij} u_i u_j - \sum_{i=1}^m \frac{a_{ii}}{\sqrt{2}} .
 \end{aligned}$$

Now the quadratic form within the braces is easily recognized and

(4.20) may be written as

$$p(u) = (2\pi)^{-m/2} e^{-\frac{1}{2} u^T \Phi^{-1} u - \sum_{i=1}^m \frac{a_{ii}}{\sqrt{2}}} - A(a)$$

where Φ^{-1} is the symmetric matrix

$$\begin{bmatrix}
 1 - \sqrt{2} a_{11} & -a_{12} & \dots & -a_{1m} \\
 & 1 - \sqrt{2} a_{22} & & -a_{2m} \\
 & & \ddots & \vdots \\
 & & & 1 - \sqrt{2} a_{mm}
 \end{bmatrix}$$

$A(\alpha)$ is the normalizing constant and when

$$A(\alpha) = -\frac{1}{2} \left\{ \sum_{i=1}^m \sqrt{2} \alpha_{ii} + \ln[|\Phi^{-1}|] \right\}$$

we may write (4.20) as

$$(4.21) \quad p(u) = \frac{|\Phi|^{-1/2}}{(2\pi)^{m/2}} e^{-\frac{1}{2} u^T \Phi^{-1} u}$$

We may summarize here that when the rank of r is two, β_r is a simple linear function of the sample covariances. It is also true that from the relationship $\beta_r = E[H_r(\underline{u})|p]$ we can conclude that when $n(r) = 2$, β_r is a simple linear function of the covariance matrix Φ . If in addition $n(r) > 2$ implies $\alpha_r = 0$ then α_r with $n(r) = 2$ is a simple linear function of the precision matrix Φ^{-1} . That is, to estimate β_r when $n(r) = 2$ is equivalent to estimating Φ and to estimate α_r when $n(r) = 2$ and when $n(r) > 2$ implies $\alpha_r = 0$ is equivalent to estimating Φ^{-1} . We also note here that the condition $n(r) > 2$ implies $\alpha_r = 0$ tacitly assumes that $p(u)$ is an m -dimensional normal density.

Estimation of the Precision Matrix of a Multivariate Normal Distribution. As in Section 4.2 let us suppose that we have the random observations ${}^t\tilde{U}$, $t = 1, \dots, n$ from one of our classes.

Furthermore we make the assumption that the distribution of \underline{t}_U is multivariate normal with a non singular covariance matrix denoted by Φ_i . We will henceforth suppress the i and simply write Φ .

The joint density of the n random vectors is then given by

$$(4.22) \quad (\underline{t}_U, \dots, \underline{t}_U | R) \sim (2\pi)^{-mn/2} |R|^{n/2} e^{-\frac{1}{2} \sum_{t=1}^n \underline{t}_U^T R \underline{t}_U}$$

where $R := \Phi^{-1}$. Let us also suppose that the precision matrix R is given a prior distribution that is a Wishart distribution with ν degrees of freedom and precision τ so that we have the density

$$(4.23) \quad (R | \nu, \tau) \sim c |\tau|^{m/2} |R|^{\nu-m-1/2} e^{-\frac{1}{2} \text{tr}(\tau R)}$$

where

$$c := \left\{ 2^{\nu m/2} \pi^{m(m-1)/4} \prod_{i=1}^m \tau \left(\frac{\nu+1-i}{2} \right) \right\}^{-1}$$

and

$$\nu > m-1.$$

In order to estimate the α 's our aim is to find estimates of the α 's in terms of the sample random vectors $\underline{t}_U, \dots, \underline{t}_U$. Inasmuch as estimating α_r for $n(r) = 2$ when $n(s) > 2$ implies $\alpha_s = 0$ is equivalent to estimating R we focus our attention on R . We have assumed a prior on \underline{R} that is $W(\nu, \tau^{-1})$ so that we have from (4.6) and (4.7) the result that the posterior distribution of \underline{R}

is also Wishart. In fact,

$$(\underline{R} | {}^1U, \dots, {}^nU, \nu, \tau) \sim W(\nu+n, \tau^{*-1})$$

where

$$(4.24) \quad \tau^* = \tau + \mathcal{U}\mathcal{U}^T$$

and

$$\mathcal{U}\mathcal{U}^T := \sum_{t=1}^n t_u t_u^T.$$

Since $t_u^T R t_u$ is a scalar for each $t = 1, \dots, n$, we have, using properties of the trace operator

$$(4.25) \quad \begin{aligned} \text{tr} \sum_{t=1}^n t_u^T R t_u &= \sum_{t=1}^n \text{tr}[t_u^T R t_u] \\ &= \sum_{t=1}^n \text{tr}[t_u t_u^T R] \\ &= \text{tr} \left[\left(\sum_{t=1}^n t_u t_u^T \right) R \right] \\ &= \text{tr}[\mathcal{U}\mathcal{U}^T R]. \end{aligned}$$

With this relationship we note that (4.22) depends upon the observations only through $\mathcal{U}\mathcal{U}^T$ so that $\mathcal{U}\mathcal{U}^T$ is a sufficient statistic for R .

If ν and τ are specified then the posterior distribution of R will give us an estimate of R based upon the sufficient statistic $u u^T$. An estimate that comes immediately to mind is the mean of the posterior. That is, by Lemma 3.2.7 we have as an estimate $(\nu+n)\tau^{*-1}$. From the relationship $R = \Phi^{-1}$ and Lemma 3.4.2 we have as another candidate

$$(4.26) \quad \hat{R} = \{E[R^{-1} | u, \dots, u, \nu, \tau]\}^{-1} \\ = (\nu+n-m-1)\tau^{*-1}.$$

We comment here that the estimate given in (4.26) has certain intuitive appeal when τ is of a certain form which will be discussed later.

If ν and τ are specified by the investigator then we have an estimate of R available. We next address ourselves to the problem of specifying ν and τ .

We may view $p_0(u)$ as initially approximating the density $p(u)$ after the Karhunen-Loève transformation. From this viewpoint we initially estimate R^{-1} by I . Under the additional assumption that $(R | \nu, \tau) \sim W(\nu, \tau^{-1})$ we have by Lemma 3.4.2

$$E[R^{-1} | \nu, \tau] = \frac{1}{\nu-m-1} \tau.$$

Combining these two ideas we set

$$\frac{1}{\nu-m-1} \tau = I$$

so that

$$\tau = (\nu-m-1)I.$$

We have used these interpretations to derive τ as a function of ν and thereby reduce the problem of specifying ν and τ to a problem of specifying only ν .

Putting aside for a moment any previous estimators for R but retaining the assumption that $\tau = (\nu-m-1)I$, we could use a procedure similar to that suggested by Lindley and Smith (1971). That is, we could think of ν as a hyperparameter and give ν a prior distribution. We then would consider the joint posterior distribution of \underline{R} and ν and integrate with respect to ν to obtain the marginal posterior distribution of \underline{R} . We then would consider some function of the marginal posterior moments as an estimate of R .

We choose to pursue a different option. That is, we choose to let the data furnish direct estimates of ν .

At this point let us recall the discussion involving equations (4.4) through (4.7). In the case that $\tau = (\nu-m-1)I$ we have the constants in (4.6) given by

$$b^* = v - m - 1$$

$$a_{ii}^* = -\frac{1}{2}$$

$$a_{ij}^* = 0.$$

The constants in (4.7) are given by

$$b = v + n - m - 1$$

$$a_{ii} = \frac{\overline{nw}_{ii} + b^* a_{ii}^*}{n + b^*}$$

$$a_{ij} = \frac{\overline{nw}_{ij}}{n + b^*}.$$

If we now apply Lemma 4.1.1 we have

$$\begin{aligned} (4.27) \quad a_{ii} &= E\{E[\tilde{W}_{ii} | R] | {}^1U, \dots, {}^nU\} \\ &= -\frac{1}{2} E\{E[u_i^2 | R] | {}^1U, \dots, {}^nU\} \\ &= -\frac{1}{2} E\{\sigma_i^2 | {}^1U, \dots, {}^nU\} \end{aligned}$$

and

$$\begin{aligned} (4.28) \quad a_{ij} &= E\{E[\tilde{W}_{ij} | R] | {}^1U, \dots, {}^nU\} \\ &= -E\{E[u_i u_j | R] | {}^1U, \dots, {}^nU\} \\ &= -E\{\sigma_{ij} | {}^1U, \dots, {}^nU\} \end{aligned}$$

What we have here is that when $\tau = (\nu - m - 1)I$ and the posterior distribution is written as in (4.7) the a_{ij} terms in the posterior distribution give us the Bayes estimates for β_r when $n(r) = 2$.

With this development we shall digress for a moment and compare this estimate with the one derived in Section 4.2.

Recall from Section 4.2 that when $n(r) = 2$, the estimator takes the form

$$\hat{\beta}_{ij} = \frac{n\bar{H}_{ij}}{n+c_k}$$

where c_k is the precision common to all β_r where $n(r) = 2$ and

$$\bar{H}_{ij} = \frac{1}{n} \sum_{t=1}^n H_{ij}(t_{u_i}).$$

Also,

$$H_{ij}(u) = \begin{cases} \frac{1}{\sqrt{2}} (u_i^2 - 1) & \text{if } i = j \\ u_i u_j & \text{if } i \neq j \end{cases}.$$

Therefore,

$$\bar{H}_{ii} = \frac{1}{n\sqrt{2}} \sum_{t=1}^n (t_{u_i}^2 - 1) \quad i = j$$

$$\bar{H}_{ij} = \frac{1}{n} \sum_{t=1}^n t_{u_i} t_{u_j} \quad i \neq j.$$

From the relationships in (4.4) and this last relationship we can write

$$\bar{w}_{ii} = \frac{1}{n} \sum_{t=1}^n \frac{1}{2} t_{u_i}^2 \quad i = j$$

$$\bar{w}_{ij} = \frac{1}{n} \sum_{t=1}^n -t_{u_i} t_{u_j} \quad i \neq j$$

$$= -\frac{1}{n} \sum_{t=1}^n t_{u_i} t_{u_j}$$

$$= -\bar{H}_{ij}.$$

We can now write $\hat{\beta}_{ij}$ as

$$\begin{aligned} (4.29) \quad \hat{\beta}_{ij} &= \frac{\bar{nH}_{ij}}{n+c_k} \\ &= -\frac{\bar{nw}_{ij}}{n+c_k} \\ &= \left(\frac{n+b^*}{n+c_k}\right)(-a_{ij}) \quad \text{if } i \neq j. \end{aligned}$$

If $b^* = c_k$, we then have

$$\hat{\beta}_{ij} = -a_{ij} \quad \text{if } i \neq j.$$

That is, if $b^* = c_k$ then the estimate $\hat{\beta}_{ij}$ and the Bayes estimate are one and the same. Note, if we replace c_k by \hat{c}_k where $\hat{c}_k = \frac{1}{\hat{\phi}_k}$ in (4.29) then $\hat{\beta}_{ij}$ and a_{ij} may still both be considered as estimators of σ_{ij} when $i \neq j$. If $\hat{c}_k = b^*$ then they are the same estimators.

$$(4.30) \quad \hat{\beta}_{ij} = \frac{\overline{nH_{ij}}}{n+c_k}$$

$$= \frac{-\frac{2}{\sqrt{2}} \overline{w_{ij}} - \frac{n}{2}}{n+c_k}.$$

If $b^* = c_k$ then

$$\hat{\beta}_{ii} = \frac{-\frac{2}{\sqrt{2}} [(n+b^*)a_{ii} - b^*a_{ii}^*] - \frac{n}{2}}{n+b^*}$$

$$= -\frac{1}{\sqrt{2}} (2a_{ii} + 1).$$

We note that from (4.27), a_{ii} is an estimate of $-\frac{1}{2}\sigma_i^2$ so that

$$\hat{\sigma}_i^2 = 2a_{ii}$$

and we may write

$$\hat{\beta}_{ii} = \frac{1}{\sqrt{2}} (\sigma_i^2 - 1).$$

If we replace c_k by \hat{c}_k in (4.30) then

$$\hat{\beta}_{ii} = \frac{\frac{1}{\sqrt{2}} (\hat{\sigma}_i^2 - 1)^2 - \frac{1}{\sqrt{2}} n}{(\hat{\sigma}_i^2 - 1)}$$

$$\doteq \frac{1}{\sqrt{2}} (\hat{\sigma}_i^2 - 1).$$

Once again, with the assumption that $b^* = c_k$ or $b^* = \hat{c}_k$ and the fact that $2a_{ii} = -\hat{\sigma}_i^2$ we get the same estimate as in Section 4.2.

At this point we note that in Section 4.2 we derived an estimator for β_r under some very general assumptions. Under the assumption of a prior on \underline{R} that is $W(\nu, \tau^{-1})$ with τ being equal to $(\nu - m - 1)I$ and considering the posterior mean of \underline{R}^{-1} as an estimate of the covariance structure we have just observed that our original estimate is appropriate in this special case with b^* equal to c_k or \hat{c}_k .

From these considerations we have

$$\hat{c}_k = \frac{1}{\hat{\varphi}_k}$$

and

$$\hat{c}_k = b^*$$

$$= \nu - m - 1.$$

We set

$$(4.31) \quad \hat{\nu} = m + 1 + \frac{1}{\hat{\varphi}_k}$$

$$\hat{\varphi}_k + \frac{1}{n} = \left\{ \frac{2}{m(m+1)} \sum_{r: n(r)=2} \overline{H}_r^2 - \frac{1}{n} \right\}^{-1}$$

and take

$$\hat{\nu} = m+1 + \left\{ \frac{2}{m(m+1)} \sum_{r:n(r)=2} \left(\bar{H}_r^2 - \frac{1}{n} \right) \right\}^{-1}$$

to be an estimate of ν .

We now have an estimate for ν and have expressed τ as a function of ν . That is, we have a specification for ν and τ as was required in the estimation process for a_r when $n(r) = 2$.

We now consider one final aspect of the problem of determining ν .

We first note that again by (4.22) and (4.25), $u u^T$ is a sufficient statistic for ν . Conceptually we consider finding the joint distribution of $u u^T$ and R given ν and τ then integrate with respect to R . This will leave us with the marginal density of $u u^T$ in terms of ν and τ . We then consider the first few moments of this distribution and choose some function of the lower order moments as an estimate of ν .

In point of fact, we can find some of the lower order moments without specifically finding the marginal distribution.

Lemma 4.3.1. If $\tau = (\nu - m - 1)I$ then $E[u u^T | \nu, \tau] = nI$.

Proof: We first note that $(u u^T | R, \nu, \tau) \sim W(n, R^{-1})$ and by Lemma 3.2.5

$$E[\underline{u}\underline{u}^T | R, \nu, \tau] = nR^{-1}.$$

\underline{R} being Wishart implies by Lemma 3.4.2,

$$E[n\underline{R}^{-1}] = \frac{n}{\nu-m-1} \tau.$$

The result now follows from the formula

$$E[\underline{u}\underline{u}^T | \nu, \tau] = E\{E[\underline{u}\underline{u}^T | R, \nu, \tau] | R\}. \quad \square$$

If we set $\underline{S} := \underline{u}\underline{u}^T$ so that $S_{ij} = n \overline{u_i u_j}$ for $i, j = 1, \dots, m$ then Lemma 4.3.1 gives us $E[S_{ij} | \nu, \tau] = n\delta_{ij}$.

Suppose we now consider functions of the second moments of S_{ij} .

Lemma 4.3.2. If $\tau = (\nu-m-1)I$ then

$$\text{var}(\overline{u_i u_j} | \nu, \tau) = \begin{cases} \frac{2}{n} + \frac{2(\frac{2}{n} + 1)}{\nu-m-3} & \text{if } i = j \\ \frac{(\nu-m-1)(\nu+n-m-1)}{n(\nu-m)(\nu-m-3)} & \text{if } i \neq j \end{cases}$$

Proof: The proof of this lemma rests on the conditional variance formula

$$\text{var}(S_{ij} | \nu, \tau) = E[\text{var}(S_{ij} | R, \nu, \tau)] + \text{var}[E(S_{ij} | R, \nu, \tau)].$$

Since $(\underline{R} | \nu, \tau) \sim W(\nu, \tau^*^{-1})$, \underline{R}^{-1} follows an inverted Wishart distribution with $\nu+m+1$ degrees of freedom and parametric matrix

τ . With $\Phi = R^{-1}$ we have by Lemmas 3.4.2 and 3.4.3 the following.

$$(4.32) \quad \begin{aligned} E(\tilde{\sigma}_{ij}) &= \frac{1}{\nu-m-1} \tau_{ij} \\ &= \delta_{ij} \end{aligned}$$

$$(4.33) \quad \begin{aligned} E(\tilde{\sigma}_{ij}^2) &= \text{var } \tilde{\sigma}_{ij} + [E\tilde{\sigma}_{ij}]^2 \\ &= \text{var } \tilde{\sigma}_{ij} + \delta_{ij}. \end{aligned}$$

$$(4.34) \quad \text{var } \tilde{\sigma}_{ij} = \begin{cases} \frac{2}{\nu-m-3} & \text{if } i = j \\ \frac{(\nu-m-1)}{(\nu-m)(\nu-m-3)} & \text{if } i \neq j, \end{cases}$$

Combining (4.33) and (4.34),

$$(4.35) \quad E(\tilde{\sigma}_{ij}^2) = \begin{cases} \frac{\nu-m-1}{\nu-m-3} & \text{if } i = j \\ \frac{\nu-m-1}{(\nu-m)(\nu-m-3)} & \text{if } i \neq j. \end{cases}$$

Again by Lemma 3.4.3

$$(4.36) \quad \text{cov}(\tilde{\sigma}_{ii}, \tilde{\sigma}_{jj}) = \begin{cases} \frac{2}{\nu-m-3} & \text{if } i = j \\ \frac{2}{(\nu-m)(\nu-m-3)} & \text{if } i \neq j. \end{cases}$$

From the relationships

$$E\tilde{\sigma}_{ii}\tilde{\sigma}_{jj} = \text{cov}(\tilde{\sigma}_{ii}, \tilde{\sigma}_{jj}) + E\tilde{\sigma}_{ii}E\tilde{\sigma}_{jj},$$

(4.32) and (4.36) we have

$$(4.36) \quad E\sigma_{\tilde{ii}}\sigma_{\tilde{jj}} = \begin{cases} \frac{\nu-m-1}{\nu-m-3} & \text{if } i = j \\ \frac{2}{(\nu-m)(\nu-m-3)} + 1 & \text{if } i \neq j. \end{cases}$$

By noting that $(\tilde{S}|R, \nu, \tau) \sim W(n, R^{-1})$ and using Lemma 3.2.6 we have

$$(4.38) \quad E(S_{\tilde{ij}}|R, \nu, \tau) = n\sigma_{ij}$$

$$\text{var}(S_{\tilde{ij}}|R, \nu, \tau) = n(\sigma_{ij}^2 + \sigma_{ii}\sigma_{jj})$$

Therefore,

$$\text{var}(S_{\tilde{ij}}|R, \nu, \tau) = nE\sigma_{\tilde{ij}}^2 + nE\sigma_{\tilde{ii}}\sigma_{\tilde{jj}}$$

Using (4.32)-(4.38) we have

$$(4.39) \quad \text{var}(S_{\tilde{ij}}|\nu, \tau) = \begin{cases} \frac{2n(\nu-m-1)}{\nu-m-3} + \frac{2n^2}{\nu-m-3} & \text{if } i = j \\ \frac{n(\nu-m+1)}{(\nu-m)(\nu-m-3)} + \frac{n^2(\nu-m-1)}{(\nu-m)(\nu-m-3)} + n & \text{if } i \neq j. \end{cases}$$

Noting that

$$\text{var}(\overline{u_i u_j} | \nu, \tau) = \text{var}\left(\frac{1}{n} S_{\tilde{ij}} | \nu, \tau\right),$$

the result follows. \square

We note that

$$\bar{H}_{ij} = \begin{cases} \frac{1}{\sqrt{2}} \overline{u_i^2} - \frac{1}{\sqrt{2}} & \text{if } i = j \\ \overline{u_i u_j} & \text{if } i \neq j \end{cases}$$

so that

$$\text{var}(\bar{H}_{ij}) = \begin{cases} \frac{1}{2} \text{var}(\overline{u_i^2}) & \text{if } i = j \\ \text{var}(\overline{u_i u_j}) & \text{if } i \neq j. \end{cases}$$

Moreover it follows from Lemma 4.3.1 that $E(\bar{H}_{ij}) = 0$ for all i and j .

From Lemma 4.3.2 we can estimate ν by a method of moments. For example we might choose the case $i = j$ when

$$\begin{aligned} \text{var}(\bar{H}_{ii}) &= \frac{1}{2} \text{var} \overline{u_i^2} \\ &= \frac{1}{n} + \frac{\frac{2}{n} + 1}{\nu - m - 3}. \end{aligned}$$

Since $E(\bar{H}_{ii}) = 0$ we have a "natural" estimate for $\text{var}(\bar{H}_{ii})$ given by

$$\frac{1}{m} \sum_{i=1}^m \bar{H}_{ii}^2.$$

We set

$$\frac{\frac{2}{n} + 1}{\hat{\nu} - m - 3} = \frac{1}{m} \sum_{i=1}^m \bar{H}_{ii}^2 - \frac{1}{n}$$

so that

$$\hat{\nu} = m+3 + \left(\frac{2}{n} + 1\right) \left\{ \frac{1}{m} \sum_{i=1}^m \bar{H}_{ii}^2 - \frac{1}{n} \right\}^{-1}$$

and take $\hat{\nu}$ to be an estimate of ν .

Similarly, we could choose the case $i \neq j$ to derive an estimate of ν . The case $i \neq j$ is not as simple as the case $i = j$ inasmuch as

$$\text{var}(\bar{H}_{ii}) = \frac{(\nu-m-1)(\nu+n-m-1)}{n(\nu-m)(\nu-m-3)}.$$

We note here that we get results that are quite similar to the estimates derived using the assumptions of Section 4.2. That is, the estimate of ν given through (4.40) are quite similar to the estimate of ν given by (4.31). The estimate given by (4.31) has one very attractive feature--ease of computation.

In this section we have concentrated on deriving an estimator for the α_r when $n(r) = 2$ and $n(s) > 2$ implies $\alpha_s = 0$. We have noted in the process that estimating the α 's of rank two is equivalent to estimating the precision matrix R and estimating the β 's of rank two is equivalent to estimating the covariance structure Φ .

4.4 Properties of Estimators

In Section 4.2 we developed the formula for $\hat{\beta}_r$ given by (4.13) from rather general considerations. We also noted two properties of $\hat{\beta}_r$:

- (i) $\hat{\beta}_r$ is biased for β_r
- (ii) $\hat{\beta}_r$ is consistent for β_r .

In this section we further investigate properties of estimates based on (4.13) as well as estimates based on the formula given by (4.26).

We comment here that the formula

$$\hat{\beta}_r = \frac{\overline{nH}_r}{n+c_r}$$

gives a way of estimating β_r for each r in the expansion

$$\frac{p(u)}{p_0(u)} = \sum_r \beta_r H_r(u).$$

However, the fundamental estimation problem is not to estimate β_r for each r but rather to estimate $p(u)$.

In most practical situations we are not interested in estimating the entire series $\sum_r \beta_r H_r(u)$ but rather in considering a convenient truncation of this series and using this truncation to estimate $p(u)$.

With this in mind we recall that the notation $\sum_r \beta_r H_r(u)$ represents the series $\sum_{c=0}^{\infty} \sum_{r:n(r)=c} \beta_r H_r(u)$. We now introduce the following notation: let $\underline{H}_k(u)$ and B_k be the vectors whose components are $H_r(u)$ and β_r respectively when $n(r) = k$. We also consider the vector b defined by

$$(b_0, b_1, \dots, b_\ell)^T := (B_0, B_1, \dots, B_k)^T$$

where b_0 corresponds to the first element of B_0 and b_ℓ corresponds to the last element of B_k . We also denote by $H^s(u)$ that polynomial $H_r(u)$ that corresponds to b_s for $s = 1, \dots, \ell$.

Now let us suppose we wish to approximate $p(u)$ by a truncation of the series at $n(r) = k$. That is, we wish to use only those terms of the expansion (2.2) through rank k .

We have by (2.2),

$$\begin{aligned} \frac{p(u)}{p_0(u)} &= \sum_r \beta_r H_r(u) \\ &= \sum_{c=0}^{\infty} \sum_{r:n(r)=c} \beta_r H_r(u). \end{aligned}$$

We set

$$\begin{aligned}
 (4.41) \quad \frac{p_k(u)}{p_0(u)} &:= \sum_{r: n(r) \leq k} \beta_r H_r(u) \\
 &= \sum_{c=0}^k \sum_{r: n(r)=c} \beta_r H_r(u) \\
 &= \sum_{s=0}^l b_s H^s(u).
 \end{aligned}$$

Now consider any other linear combination of the $H_r(u)$ where $n(r) \leq k$. Let us denote this arbitrary linear combination by

$$\frac{\pi_k(u)}{p_0(u)} := \sum_{s=0}^l a_s H^s(u)$$

Since $H^s(u)$ is a complete orthonormal family in $L_2(P_0)$,

Parseval's equation is valid. That is,

$$\left\| \frac{p(u)}{p_0(u)} \right\|^2 = \sum_{s=0}^{\infty} b_s^2.$$

Now, it is a well known result [e.g., Natanson (1965)] that the error as measured by the $L_2(P_0)$ norm in approximating $\frac{p(u)}{p_0(u)}$ by $\frac{\pi_k(u)}{p_0(u)}$ is given by

$$(4.42) \quad \left\| \frac{p(u)}{p_0(u)} - \frac{\pi_k(u)}{p_0(u)} \right\|^2 = \sum_{s=0}^{\infty} b_s^2 - \sum_{s=0}^{\ell} b_s^2 + \sum_{s=0}^{\ell} (a_s - b_s)^2.$$

From this result we see that the error in terms of the $L_2(P_0)$ norm is minimum when $\pi_k(u) = p_k(u)$. That is, when $a_s = b_s$, $s = 0, 1, 2, \dots, n$.

Having determined that the best $L_2(P_0)$ approximation to $\frac{p(u)}{p_0(u)}$ by a linear combination of the $H_r(u)$ for which $n(r) \leq k$ is given by (4.41) we next consider as an estimate of $\frac{p(u)}{p_0(u)}$

$$(4.42) \quad \overline{\frac{p_k(u)}{p_0(u)}} := \sum_{r: n(r) \leq k} \hat{\beta}_r H_r(u).$$

We now consider properties of the expected mean square error in using (4.42) to estimate $\frac{p(u)}{p_0(u)}$.

We have immediately the following results.

Lemma 4.4.1. For k fixed, the expected mean squared error of $\overline{\frac{p_k(u)}{p_0(u)}}$ is

$$\sum_{r: n(r) > k} \beta_r^2 + \sum_{r: n(r) \leq k} \{ \text{var}(\hat{\beta}_r | p) + \beta_r^2 \left(\frac{n}{n+c} - 1 \right)^2 \}$$

Proof: The result is immediate by setting

$$\frac{\overline{p_k(u)}}{p_0(u)} = \frac{\pi_k(u)}{p_0(u)}$$

in (4.42), taking the expectation of (4.42) with respect to $p(u)$ and recalling from Section 4.2 that

$$E[(\hat{\beta}_r - \beta_r)^2 | p] = \text{var}(\hat{\beta}_r | p) + \beta_r^2 \left(\frac{n}{n+c_r} - 1 \right)^2. \quad \square$$

Lemma 4.42. For fixed k , $\lim_{n \rightarrow \infty} E \left\| \frac{p(u)}{p_0(u)} - \frac{\overline{p_k(u)}}{p_0(u)} \right\|^2 = \sum_{n(r) > k} \beta_r^2.$

Proof: This follows immediately from Lemma 4.4.1 and the consistency of $\hat{\beta}_r$ for β_r derived in Section 4.2. \square

We now turn our attention to the estimators developed in Section 4.3 under the assumption that when $n(s) > 2$, $\alpha_s = 0$.

We recall that this assumption is equivalent to assuming $p(u)$ is a multivariate normal density. Moreover, with this assumption the problem of estimating α_r when $n(r) = 2$ is identical to estimating R , the inverse of the covariance matrix. We also recall that having assumed a prior on \underline{R} that is $W(\nu, \tau^{-1})$ with $\tau = (\nu - m - 1)I$ we found that the posterior distribution of \underline{R} was $W(\nu + n, \tau^{*-1})$ where $\tau^* = \tau + uu^T$ and $uu^T = \sum_{t=1}^n t_u t_u^T$. It is also true that the posterior distribution of $\underline{\Phi} = \underline{R}^{-1}$ is inverted Wishart with

$$E\tilde{\Phi} = \frac{1}{\nu+n-m-1} \tau^* .$$

We consider the posterior mean as an estimate of Φ so that we set

$$\hat{\Phi} := \frac{1}{\nu+n-m-1} \tau^* .$$

Equation (4.26) gives an estimate of R based on the inverse of the posterior mean of Φ . That is, (4.26) gives

$$\begin{aligned} \hat{R} &= (\nu+n-m-1)\tau^{*-1} \\ &= (\hat{\Phi})^{-1} . \end{aligned}$$

We now consider some of the properties of these estimates.

We set $S := uu^T$ and $\Phi_0 := \frac{1}{n}S$ so that Φ_0 is the maximum likelihood estimate of Φ when sampling from a normal distribution with mean 0 and covariance matrix Φ .

With $\tau = (\nu-m-1)I$ we set $a := \nu-m-1$ and have the result

$$\begin{aligned} (4.43) \quad \hat{\Phi} &= \frac{1}{a+n} [aI+S] \\ &= \frac{a}{a+n} I + \frac{1}{a+n} S \\ &= \frac{a}{a+n} I + \frac{n}{a+n} \Phi_0 . \end{aligned}$$

We may write (4.43) in component form as

$$\hat{\Phi}_{ij} = \begin{cases} \frac{n}{a+n} (\Phi_0)_{ij} & \text{if } i \neq j \\ \frac{a}{a+n} + \frac{n}{a+n} (\Phi_0)_{ij} & \text{if } i = j. \end{cases}$$

From this equation we see that for $i \neq j$, $\hat{\Phi}$ tends to be shifted from the maximum likelihood estimate towards zero while if $i = j$, $\hat{\Phi}$ tends to be shifted from the maximum likelihood estimate towards one.

From (4.43) we see that

$$(4.44) \quad \begin{aligned} E \hat{\Phi} &= \frac{a}{a+n} I + \frac{n}{a+n} E \Phi_0 \\ &= \frac{a}{a+n} I + \frac{n}{a+n} \Phi. \end{aligned}$$

From this equation it is immediate that $\hat{\Phi}$ is asymptotically unbiased.

We now consider the norm on the space of all $m \times m$ symmetric matrices given by

$$\begin{aligned} \|\Phi\|^2 &= \langle \Phi, \Phi \rangle \\ &= \text{tr } \Phi^2. \end{aligned}$$

In the classical sense, recall that the relative efficiency of an unbiased estimate $\bar{\Phi}$ relative to the unbiased estimate $\bar{\Phi}$ is defined by

$$\text{r. e.} := \frac{E \|\bar{\Phi} - \Phi\|^2}{E \|\bar{\Phi} - \Phi\|^2}$$

An estimator $\bar{\Phi}$ is called most efficient in case $\text{r. e.} \leq 1$ for all $\bar{\Phi}$. It is well known that Φ_0 is most efficient in the class of unbiased estimates. We note that in the classical sense the expectations are computed with respect to the distribution of $(U|\Phi)$.

Lemma 4.4.3. Under conditions of sampling from a multivariate normal distribution with mean 0 and covariance matrix Φ , if $\hat{\Phi}$ is given by (4.43) then

$$E \|\hat{\Phi} - \Phi\|^2 = \left(\frac{a}{a+n}\right)^2 \text{tr}[I - 2\Phi + \Phi^2] + \frac{n}{(a+n)^2} [(\text{tr } \Phi)^2 + \text{tr}(\Phi^2)]$$

where the expectation is with respect to the distribution of U given Φ .

Proof: Using the fact that $(S|\Phi, \nu, \tau) \sim W(n, \Phi)$

$$\begin{aligned} E \|\hat{\Phi} - \Phi\|^2 &= E \text{tr} \left[\frac{a}{a+n} I + \frac{1}{a+n} S - \Phi \right]^2 \\ &= \text{tr} E \left[\left(\frac{a}{a+n} \right)^2 I + \frac{2a}{(a+n)^2} S - \frac{2a}{a+n} \Phi + \frac{1}{(a+n)^2} S^2 - \frac{2}{(a+n)} S\Phi + \Phi^2 \right] \\ &= \text{tr} \left[\left(\frac{a}{a+n} \right)^2 I + \frac{2an}{(a+n)^2} \Phi - \frac{2a}{a+n} \Phi + \Phi^2 + \frac{1}{(a+n)^2} E S^2 - \frac{2}{(a+n)} E S\Phi \right] = \end{aligned}$$

$$\begin{aligned}
&= \text{tr} \left[\left(\frac{a}{a+n} \right)^2 I - 2 \left(\frac{a}{a+n} \right)^2 \Phi + \frac{a-n}{a+n} \Phi^2 \right] + \text{tr} \, E \left[\frac{1}{(a+n)^2} \tilde{S}^2 \right] \\
&= \left(\frac{a}{a+n} \right)^2 \text{tr} [I - 2\Phi + \Phi^2] - \text{tr} \left[\left(\frac{n}{a+n} \right)^2 \Phi^2 \right] + \text{tr} \, E \left[\frac{1}{(a+n)^2} \tilde{S}^2 \right] \\
&= \left(\frac{a}{a+n} \right)^2 \text{tr} [I - 2\Phi + \Phi^2] - \text{tr} \left[\left(\frac{n}{a+n} \right)^2 \Phi^2 \right] + \frac{1}{(a+n)^2} E \, \text{tr} \, \tilde{S}^2.
\end{aligned}$$

Now,

$$\begin{aligned}
E \, \text{tr} \, \tilde{S}^2 &= E \sum_i \sum_j S_{ij}^2 \\
&= \sum_i \sum_j E S_{ij}^2 \\
&= \sum_i \sum_j [\text{var} \, S_{ij} + (E S_{ij})^2].
\end{aligned}$$

From Lemmas 3.2.6 we have

$$\begin{aligned}
E \, \text{tr} \, \tilde{S}^2 &= \sum_i \sum_j [n(\sigma_{ii}\sigma_{jj} + \sigma_{ij}^2) + n^2 \sigma_{ij}^2] \\
&= n(\text{tr} \, \Phi)^2 + n \, \text{tr}(\Phi^2) + n^2 \, \text{tr}(\Phi^2).
\end{aligned}$$

Therefore,

$$E \left\| \hat{\Phi}_{\tilde{S}} - \Phi \right\|^2 = \left(\frac{a}{a+n} \right)^2 \text{tr} [I - 2\Phi + \Phi^2] + \frac{n}{(a+n)^2} [(\text{tr} \, \Phi)^2 + \text{tr}(\Phi^2)]. \quad \square$$

Lemma 4.4.4. Under the same conditions as Lemma 4.4.3

we have for $\Phi_0 = \frac{1}{n} S$,

$$E \|\hat{\Phi}_0 - \Phi\|^2 = \frac{1}{n} [(tr \Phi)^2 + tr(\Phi^2)].$$

Proof:

$$\begin{aligned} E \|\hat{\Phi}_0 - \Phi\|^2 &= E \operatorname{tr} \left(\frac{1}{n} \tilde{S} - \Phi \right)^2 \\ &= \sum_i \sum_j E \left(\frac{1}{n} \tilde{S}_{ij} - \sigma_{ij} \right)^2 \\ &= \sum_i \sum_j \operatorname{var} \left(\frac{1}{n} \tilde{S}_{ij} \right) \\ &= \sum_i \sum_j \frac{1}{n^2} \operatorname{var} \tilde{S}_{ij} \\ &= \sum_i \sum_j \frac{1}{n^2} n(\sigma_{ii} \sigma_{jj} + \sigma_{ij}^2) \\ &= \frac{1}{n} [(tr \Phi)^2 + tr(\Phi^2)]. \quad \square \end{aligned}$$

We note here that $E \|\hat{\Phi} - \Phi\|^2 < E \|\hat{\Phi}_0 - \Phi\|^2$ if and only if

$$\frac{an}{a+2n} \operatorname{tr}(I - \Phi)^2 < [(tr \Phi)^2 + tr(\Phi^2)].$$

From this relationship we have that $E \|\hat{\Phi} - \Phi\|^2 < E \|\hat{\Phi}_0 - \Phi\|^2$ when

Φ is "close" to I , the prior mean of Φ ,

Lemmas 4.4.3 and 4.4.4 allow us to write

$$\lim_{n \rightarrow \infty} \frac{E \|\hat{\Phi}_0 - \Phi\|^2}{E \|\hat{\Phi} - \Phi\|^2} = \lim_{n \rightarrow \infty} \frac{\frac{1}{n} [(\text{tr } \Phi)^2 + \text{tr}(\Phi^2)]}{\left(\frac{a}{a+n}\right)^2 \text{tr}(I - \Phi)^2 + \frac{n}{(a+n)^2} [(\text{tr } \Phi)^2 + \text{tr}(\Phi^2)]}$$

$$= 1.$$

We may conclude from this last expression that $\hat{\Phi}$ is asymptotically most efficient.

We also have by Lemma 4.4.3 the following lemma.

Lemma 4.4.5. $\hat{\Phi}$ is a consistent estimator of Φ .

Proof: By Lemma 4.4.3 we write

$$E \|\hat{\Phi} - \Phi\|^2 = \left(\frac{a}{a+n}\right)^2 \text{tr}[I - \Phi]^2 + \frac{n}{(a+n)^2} [(\text{tr } \Phi)^2 + \text{tr}(\Phi^2)]$$

from which we conclude that

$$E \|\hat{\Phi} - \Phi\|^2 \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

That is, $\hat{\Phi} \rightarrow \Phi$ in mean of order two and therefore $\hat{\Phi} \rightarrow \Phi$ in probability. \square

These results involving asymptotic results are not surprising inasmuch as we can see from (4.43) that as $n \rightarrow \infty$, $\hat{\Phi}$ is

essentially the same as $\hat{\Phi}_0$.

We have been considering the quadratic loss function for Φ given by $L(\hat{\Phi}, \Phi) = \text{tr}[\hat{\Phi} - \Phi]^2$. Equivalently, we have been considering the loss function for R given by

$$L(\hat{R}, R) = \text{tr}[\hat{R}^{-1} - R^{-1}]^2$$

which is quadratic in R^{-1} .

It is well known that under quadratic loss the Bayes rule for estimating Φ is the mean of the posterior. That is, $\hat{\Phi}$ is Bayes with respect to $\text{tr}[\hat{\Phi} - \Phi]^2$ for estimating Φ and equivalently $\hat{R} = \hat{\Phi}^{-1}$ is Bayes with respect to $\text{tr}[\hat{R}^{-1} - R^{-1}]^2$ for estimating R .

Lemma 4.4.6. Under the quadratic loss function

$L(\hat{\Phi}, \Phi) = \text{tr}[\hat{\Phi} - \Phi]^2$ the Bayes risk, R_0 in estimating Φ by $\hat{\Phi}_0$ is given by

$$R_0 = \frac{1}{n} \left[\frac{(m^2 + m)a^2 + 2am}{(a+1)(a-2)} \right]$$

Proof:

$$\begin{aligned} R_0 &= E_{\Phi} \{ E_{u|\Phi} L(\hat{\Phi}_0, \Phi) \} \\ &= E_{\Phi} \{ E_{u|\Phi} \text{tr}[\hat{\Phi}_0 - \Phi]^2 \} \end{aligned}$$

From Lemma 4.4.4,

$$\begin{aligned} E_{u|\Phi} \operatorname{tr}[\Phi - \tilde{\Phi}_0]^2 &= \frac{1}{n}[(\operatorname{tr} \Phi)^2 + \operatorname{tr}(\Phi^2)] \\ &= \sum_i \sum_j \frac{1}{n} (\sigma_{ii} \sigma_{jj} + \sigma_{ij}^2). \end{aligned}$$

Therefore,

$$\begin{aligned} R_0 &= E_{\Phi} \sum_i \sum_j \frac{1}{n} (\sigma_{ii} \sigma_{jj} + \sigma_{ij}^2) \\ &= \frac{1}{n} \sum_i \sum_j [E \sigma_{ii} \sigma_{jj} + E \sigma_{ij}^2] \\ &= \frac{1}{n} \sum_i \sum_j [\operatorname{cov}(\sigma_{ii}, \sigma_{jj}) + (E \sigma_{ii})(E \sigma_{jj})] \\ &\quad + \sum_i \sum_j [\operatorname{var} \sigma_{ij} + (E \sigma_{ij})^2]. \end{aligned}$$

Now, \tilde{R} has prior that is $W(\nu, \tau^{-1})$ so that the distribution of $\tilde{\Phi}$ is inverted Wishart with parameters $\nu+m+1$ and aI so we may invoke Lemmas 3.4.2 and 3.4.3 with $\tau_{ij} = a\delta_{ij}$. Therefore,

$$\begin{aligned} R_0 &= \frac{1}{n} \left\{ \sum_i \sum_j [\operatorname{cov}(\sigma_{ii}, \sigma_{jj})] + \sum_i \sum_j 1 + \sum_i \sum_j \operatorname{var} \sigma_{ij} + \sum_i \sum_j \delta_{ij} \right\} \\ &= \frac{1}{n} \left\{ \sum_i \sum_j \frac{2a+2a^2\delta_{ij}}{(\nu-m)a(\nu-m-3)} + m^2 + \sum_i \sum_j \frac{a^2+(\nu-m+1)a\delta_{ij}}{(\nu-m)a(\nu-m-3)} + m \right\} = \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n} \left\{ m^2 + m + \sum_i \sum_j \frac{2a + 2a^2 \delta_{ij}}{(\nu - m)(\nu - m - 3)} + m^2 + \sum_i \sum_j \frac{a + (\nu - m + 1) \delta_{ij}}{(\nu - m)(\nu - m - 3)} \right\} \\
&= \frac{1}{n} \left\{ m^2 + m + \frac{2m^2}{(\nu - m)(\nu - m - 3)} + \frac{2am}{(\nu - m)(\nu - m - 3)} + \frac{am^2}{(\nu - m)(\nu - m - 3)} \right. \\
&\quad \left. + \frac{(\nu - m + 1)m}{(\nu - m)(\nu - m - 3)} \right\} \\
&= \frac{1}{n} \left\{ m^2 + m + \frac{2m^2 + 2am + am^2 + (a + 2)m}{(a + 1)(a - 2)} \right\} \\
&= \frac{1}{n} \left\{ \frac{(m^2 + m)a^2 + 2am}{(a + 1)(a - 2)} \right\}. \quad \parallel
\end{aligned}$$

Lemma 4.4.7. Under the conditions of Lemma 4.4.6, the Bayes risk \mathcal{R} in using the Bayes rule $\hat{\Phi}$ is given by

$$\mathcal{R} = \frac{a^3 m^2 + a^2 (a + 2)m}{(a + n)^2 (a + 1)(a - 2)} + \frac{n^2}{(a + n)^2} \mathcal{R}_0.$$

Proof: From Lemma 4.4.3,

$$E_{u|\Phi} \operatorname{tr} [\hat{\Phi} - \Phi]^2 = \left(\frac{a}{a + n} \right)^2 \operatorname{tr} [I - \Phi]^2 + \frac{n}{(a + n)^2} [(\operatorname{tr} \Phi)^2 + \operatorname{tr}(\Phi^2)].$$

Therefore,

$$\mathcal{R} = \left(\frac{a}{a + n} \right)^2 E \operatorname{tr} [I - \Phi]^2 + \frac{n}{(a + n)^2} E [(\operatorname{tr} \Phi)^2 + \operatorname{tr}(\Phi^2)].$$

From Lemma 4.4.4,

$$(\operatorname{tr} \Phi)^2 + \operatorname{tr}(\Phi^2) = n E_{u|\Phi} \operatorname{tr} [\hat{\Phi}_0 - \Phi]^2.$$

Therefore,

$$\begin{aligned}
 &= \left(\frac{a}{a+n}\right)^2 E \operatorname{tr}[I-\hat{\Phi}]^2 + \left(\frac{n}{a+n}\right)^2 R_0 \\
 &= \left(\frac{a}{a+n}\right)^2 \sum_i \sum_j E(\delta_{ij} - \sigma_{ij})^2 + \left(\frac{n}{a+n}\right)^2 R_0 \\
 &= \left(\frac{a}{a+n}\right)^2 \sum_i \sum_j \operatorname{var} \sigma_{ij} + \left(\frac{n}{a+n}\right)^2 R_0 \\
 &= \frac{a^2}{(a+n)^2} \frac{am^2 + (a+2)m}{(a+1)(a-2)} + \frac{n^2}{(a+n)^2} R_0 \quad |
 \end{aligned}$$

It is the content of the next lemma that the Bayes risk in using $\hat{\Phi}$ is strictly less than the Bayes risk in using Φ_0 .

Lemma 4.4.8. Under the conditions of Lemmas 4.4.6 and 4.4.7,

$$R < R_0.$$

Proof: From Lemmas 4.4.6 and 4.4.7,

$$\begin{aligned}
 R_0 - R &= \left[1 - \frac{n^2}{(a+n)^2}\right] R_0 - \frac{a^2}{(a+n)^2} \left[\frac{am^2 + (a+2)m}{(a+1)(a-2)} \right] \\
 &= \frac{a^2 + 2an}{(a+n)^2} \frac{(m^2 + m)a^2 + 2am}{n(a+1)(a-2)} - \frac{a^2 [am^2 + (a+2)m]}{(a+n)^2 (a+1)(a-2)} \\
 &= \frac{1}{(a+n)^2 (a+1)(a-2)} \left\{ \left(\frac{a^2 + 2an}{n}\right) [(m^2 + m)a^2 + 2am] - a^2 [am^2 + (a+2)m] \right\} =
 \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{(a+n)^2(a+1)(a-2)} \left\{ \frac{a^4}{n} (m^2+m) + a^3 m^2 + a^3 m + \frac{2a^3 m}{n} + 2a^2 m \right\} \\
&= \frac{1}{(a+n)^2(a+1)(a-2)} \left\{ \frac{a^4}{n} (m^2+m) + a^3 (m^2+m) + 2a^2 m \left(\frac{a}{n} + 1 \right) \right\} \\
&= \frac{a^2}{(a+n)^2(a+1)(a-2)} \left\{ \frac{a^2}{n} (m^2+m) + a(m^2+m) + 2m \left(\frac{a}{n} + 1 \right) \right\}.
\end{aligned}$$

We note that for the required moments of the inverted Wishart distribution to exist we require that $\nu - m - 3 > 0$. That is, $a - 2 > 0$. With this fact we see that the right hand side of the last equation is the product of two strictly positive terms and the result is established. \square

As a final comment in this section we consider a problem that arises in the classification problem.

In certain classification problems particularly in the area of pattern recognition the problem of insufficient data occurs. By this we mean that we are trying to classify an m -component vector based on a sample of size n where n may be much smaller than m . If we are using certain linear discriminant functions it will be required that we invert our estimate of Φ . If we consider Φ_0 then we need $n \geq m$. A possible way around this problem is to use a feature extractor to select the "best" m_1 components from X where $n \geq m_1$. With the requirement that $m_1 \leq n$, it is easy to

see that the investigator might be required to eliminate components which he feels are pertinent to the problem.

An alternative would be to consider another estimator of Φ . The estimator $\hat{\Phi}$ given by (4.43) has one very nice property with respect to this problem. Since

$$\begin{aligned}\hat{\Phi} &= \frac{a}{a+n} I + \frac{1}{a+n} S \\ &= \frac{a}{a+n} I + \frac{n}{a+n} \Phi_0\end{aligned}$$

and $\frac{a}{a+n} > 0$ we have that $\hat{\Phi}$ is positive definite and therefore invertible for every $n \geq 1$.

BIBLIOGRAPHY

- Anderson, T.W. 1958. An introduction to multivariate statistical analysis. New York, Wiley. 374 p.
- Barnett, Vic. 1973. Comparative statistical inference. New York, Wiley. 287 p.
- Barndorff-Nielsen, Ole. 1973. Exponential families and conditioning. New York, Wiley.
- Bellman, Richard. 1960. Introduction to matrix analysis. New York, McGraw-Hill. 328 p.
- Brunk, H.D. and D.A. Pierce. 1974. Estimation of discrete multivariate densities for computer aided differential diagnosis of disease. *Biometrika* 61:493-501.
- Burrill, Claude. 1972. Measure, integration, and probability. New York, McGraw-Hill. 464 p.
- Crain, Bradford R. 1974. Estimation of distributions using orthogonal expansions. *The Annals of Statistics* 2:454-463.
- Efron, B. and Carl Morris. 1972. Empirical Bayes on vector observations: An extension of Stein's method. *Biometrika* 59:335-347.
- Efron, B. and Carl Morris. 1973. Stein's estimation rule and its competitors--an Empirical Bayes approach. *Journal of the American Statistical Association* 68:117-130.
- Ferguson, T.S. 1967. *Mathematical Statistics a decision theoretic approach*. New York, McGraw-Hill. 396 p.
- Fukunaga, K. 1972. *Introduction to statistical pattern recognition*. New York, Academic Press. 369 p.
- Graybill, F. 1969. *Introduction to matrices with applications in statistics*. Belmont, Wadsworth. 372 p.
- Jackson, D. 1941. *Fourier series and orthogonal polynomials*. Carus Mathematical Monograph No. 6. Menasha, Wisc. Banta. 234 p.

- Kshirsagar, A. 1972. Multivariate analysis. New York, Marcel Dekker. 534 p.
- Kronmanl, R. and Michael Tarter. 1968. The estimation of probability densities and accumulatives by Fourier series methods. Journal of the American Statistical Association 63:925-952.
- Lancaster, H.O. 1969. The Chi-Squared Distribution. New York, Wiley. 356 p.
- Lindley, D.V. and A.F.M. Smith. 1972. Bayes estimates for the linear model. Journal of the Royal Statistical Society. 34:1-41.
- Lukacs, E. and R.G. Laha. 1964. Applications of characteristic functions. London. Griffin. 202 p.
- Martin, D.C. and R.A. Bradley. 1972. Probability models, estimation and classification for multivariate dichotomous populations. Biometrics 23:203-221.
- Morgan, R.L. 1969. A class of conjugate prior distributions and optimal allocation. 1969. Unpublished Ph.D. thesis at Univ. of Missouri.
- Natanson, I.P. 1965. Constructive function theory Vol II. New York. Ungar. 174 p.
- Olkin, I. and H. Rubin. 1962. A characterization of the Wishart distribution. The Annals of Mathematical Statistics: 33. No. 4. 1272-1280.
- Press, S.J. 1972. Applied Multivariate Analysis. New York. Holt, Rinehart and Winston. 521 p.
- Raiffa, H. and R. Schlaifer. 1961. Applied statistical decision theory. Cambridge, Mass. M.I.T. Press. 356 p.
- Rao, C.R. 1965. Linear statistical inference and its applications. New York. Wiley. 522 p.
- Sansone, G. 1959. Orthogonal functions. New York. Interscience. 411 p.

- Schwartz, Stuart C. 1967. Estimation of probability density by orthogonal series. *The Annals of Mathematical Statistics* 38:1262-1265.
- Smith, K.T. 1971. *Smiths primer of modern analysis*. New York. Bogden and Quigley. 412 p.
- Stein, C., Efron, B. and C. Morris. 1972. Improving the usual estimator of a normal covariance matrix. Stanford University Dept. of Statistics, Technical Report No. 37.
- Watson, Geoffrey S. 1969. Density estimation by orthogonal series. *The Annals of Mathematical Statistics* 40:1496-1498.
- Zygmund, A. 1959. *Trigonometric series*. Vol II. Cambridge Univ. Press. 354 p.

APPENDIX

This appendix contains results scattered throughout various reference works that were found to be useful in writing this thesis.

A.1 Matrix Differentiation

The results given here are formulas for differentiating matrices and functions of matrices with respect to matrices as well as scalars. Most of these results are found in Graybill (1969) as well as various books in econometrics and nonlinear programming.

Definition A.1.1. If f is a function of the real variables x_1, x_2, \dots, x_m then f is a function of the vector $(x_1, \dots, x_m)^T$. We define the derivative of f with respect to $x = (x_1, \dots, x_m)^T$ by the formula

$$\frac{\partial f}{\partial x} := \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_m} \right)^T$$

Definition A.1.2. If f is a function of the mn independent real variables $x_{11}, x_{12}, \dots, x_{mn}$ then f is a function of the $m \times n$ matrix (x_{ij}) $i = 1, \dots, m, j = 1, \dots, n$. We define the derivative of f with respect to the matrix $X = (x_{ij})$ to be the $m \times n$ matrix with

$$\left(\frac{\partial f}{\partial X}\right)_{ij} := \frac{\partial f}{\partial x_{ij}}.$$

We should note here that x_{ij} denotes the element in the i th row and j th column of X . If there are any functional relationships between the elements they are disregarded in this definition. For example if X is $m \times m$ and symmetric then we consider the matrix Y where the components of Y are functions of the matrix X given by $y_{ij}(x) = y_{ji}(x) = x_{ij}$, $1 \leq i \leq j \leq m$. We now have by the chain rule,

$$\begin{aligned} \frac{\partial f}{\partial x_{ij}} &= \sum_{k,l} \frac{\partial f}{\partial y_{kl}} \frac{\partial y_{kl}}{\partial x_{ij}} \\ &= \frac{\partial f}{\partial y_{ij}} + \frac{\partial f}{\partial y_{ji}}. \end{aligned}$$

Definition A.1.3. If Y is an $m \times m$ matrix and y_{ij} is some function of the real variable ξ we define the derivative of Y with respect to ξ to be the $m \times m$ matrix with

$$\left(\frac{\partial Y}{\partial \xi}\right)_{ij} := \frac{\partial y_{ij}}{\partial \xi}.$$

For notation, if X is a matrix we shall denote the ij cofactor of X by $|X_{ij}|$ and the matrix of cofactors of X by $X^\#$. We now state several results regarding matrix differentiation most of which

can be found in Graybill (1969).

Result A. 1. 1.

(i) If X and Y are matrices where x_{ij} and y_{ij} are functions of the real variable ξ then

$$\frac{\partial(X+Y)}{\partial\xi} = \frac{\partial X}{\partial\xi} + \frac{\partial Y}{\partial\xi} .$$

(ii) If $X = (x_{ij})$ is an $m \times m$ matrix of independent variables then

$$\frac{\partial |X|}{\partial X} = X^\#$$

(iii) If $Y = (y_{kl})$ is an $m \times m$ matrix and each element y_{kl} is a function of several independent real variables including ξ then

$$\frac{\partial |Y|}{\partial\xi} = \text{tr}[Y^\# \frac{\partial Y^T}{\partial\xi}] .$$

(iv) If $X = (x_{ij})$ is an $m \times m$ non singular matrix of independent variables then

$$\frac{\partial X^{-1}}{\partial x_{ij}} = -X^{-1} \Delta_{ij} X^{-1}$$

where Δ_{ij} is an $m \times m$ matrix with an entry of one as the ij element and zeroes elsewhere.

We comment here that if X is a symmetric matrix then we may view it as a matrix whose elements are functions of the $\frac{m(m+1)}{2}$ independent variables x_{ij} , $1 \leq i \leq j \leq m$. With this interpretation we state the following results.

Result A.1.2. If X is an $m \times m$ symmetric matrix then for $1 \leq i \leq j \leq m$:

- (i) $\frac{\partial X}{\partial x_{ij}} = \Delta_{ij}^*$ where Δ_{ij}^* contains entries of zero as the ij and ji elements and zeroes elsewhere.
- (ii) $\frac{\partial |X|}{\partial x_{ij}} = 2|X_{ij}| - \delta_{ij}|X_{ij}|$.
- (iii) If in addition X is non singular, $\frac{\partial X^{-1}}{\partial x_{ij}} = -X^{-1} \Delta_{ij}^* X^{-1}$.

We note here that (i) follows directly from Definition A.1.3., (ii) follows direction from A.1.1, (iii) and A.1.2. (i), and (iii) follows from the relationship $X^{-1}X = I$ and the chain rule.

A.2 Characteristic Functions

Definition A.2.1. Let $(W, \mathcal{W}, \lambda)$ be a measure space and let \underline{X} be a random vector defined on W into the m -dimensional Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle)$. We define the characteristic function of \underline{X} to be the function $\phi : \mathcal{H} \rightarrow \mathbb{R}^1$ defined for each $t \in \mathcal{H}$ by

$$\begin{aligned}\phi_{\underline{X}}(t) &= \int e^{i\langle t, \underline{x} \rangle} dP_{\underline{X}}(\underline{x}) \\ &= \int e^{i\langle t, X(\omega) \rangle} d\lambda(\omega).\end{aligned}$$

We note here that for random vectors $\underline{X} \in \mathbb{R}^m$ we will use

$$\langle t, \underline{X} \rangle = \sum_{i=1}^m t_i X_i.$$

We also note that when dealing with random matrices there are two Hilbert spaces that are particularly useful.

Example A.2.1. Let $\mathcal{H}_{mn} = \{A : A \text{ is an } m \times n \text{ real matrix}\}$
 $(\mathcal{H}_{mn}, \langle \cdot, \cdot \rangle)$ is a Hilbert space of dimension mn where
 $\langle A, B \rangle = \text{tr}(A^T B)$.

Example A.2.2. Let $\mathcal{H}_{mm} = \{A : A \text{ is an } m \times m \text{ real symmetric matrix}\}$
 $(\mathcal{H}_{mm}, \langle \cdot, \cdot \rangle)$ is a Hilbert space of dimension $\frac{m(m+1)}{2}$ where $\langle A, B \rangle = \text{tr}(AB)$.

We comment here that if we have the symmetric matrices V and t where

$$t = \begin{bmatrix} 2t_{11} & t_{12} & \cdots & t_{1m} \\ t_{21} & 2t_{22} & \cdots & t_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ t_{m1} & t_{m2} & \cdots & 2t_{mm} \end{bmatrix}$$

and we use the inner product $\frac{1}{2} \text{tr}(tV)$ in Example A.2.2 then numerically this is equivalent to viewing V and t as vectors V_{ij} and t_{ij} , $1 \leq i \leq j \leq m$ in $\mathbb{R}^{m(m+1)/2}$ and using the standard inner product.

Definition A.2.2. If \tilde{X} is a random vector with characteristic function $\phi_{\tilde{X}}(t) = E_{\tilde{X}}\{e^{i\langle t, \tilde{X} \rangle}\}$ then a product moment of order k is defined to be $E[\tilde{X}_1^{J_1} \tilde{X}_2^{J_2} \dots \tilde{X}_m^{J_m}]$ provided the expectation exists and where $\sum_{i=1}^m J_i = k$, $\{e_i : i = 1, \dots, m\}$ is a basis and $\tilde{X}_i = \langle \tilde{X}, e_i \rangle$.

The next three lemmas concerning the characteristic functions of random vectors are found in Lukacs (1964).

Lemma A.2.1. Suppose g is a finite real valued, measurable function defined on the range of X . Set $\tilde{Y} = g(\tilde{X})$. We then have that

$$\begin{aligned} \phi_{\tilde{Y}}(t) &= E_{\tilde{Y}}\{e^{ity}\} \\ &= E_{\tilde{X}}\{e^{itg(x)}\} \\ &= \int_{\mathbb{R}^m} e^{itg(x)} dP_X. \end{aligned}$$

Lemma A.2.2. Let g_i , $i = 1, \dots, s$ be single valued functions and set $\tilde{Y} = (\tilde{Y}_1, \dots, \tilde{Y}_s)^T$ where $\tilde{Y}_i = g_i(\tilde{X})$. The

characteristic function of \underline{Y} is then given by

$$\begin{aligned}\phi_{\underline{Y}}(t) &= E_{\underline{Y}} e^{it^T \underline{Y}} \\ &= \int_{\mathbb{R}^m} e^{it^T g(x)} dP_X.\end{aligned}$$

We should note here that if we set $g_i(x) = x_i$ in Lemma A.2.2 we then have the result that the joint characteristic function of any s components of the random vector \underline{X} may be found by setting $t_i = 0, i \neq 1, 2, \dots, s$ in the characteristic function of \underline{X} .

Lemma A.2.3. Suppose \underline{X} is a random vector with each component of \underline{X} having finite moments up to order k . It then follows that all product moments of order k exist. Moreover,

$$\left. \frac{\partial^k \phi_{\underline{X}}(t)}{\partial t_1^{j_1} \dots \partial t_m^{j_m}} \right|_{t=0} = i^{j_1 + \dots + j_m} E[\underline{X}_1^{j_1} \underline{X}_2^{j_2} \dots \underline{X}_m^{j_m}]$$

In the special case where each component of \underline{X} has a finite first moment then we define the expectation of \underline{X} to be the vector of expectations. That is,

$$E\underline{X} = \begin{pmatrix} E\underline{X}_1 \\ \vdots \\ E\underline{X}_m \end{pmatrix}.$$

Lemma A.2.3 tells us that $E\tilde{X}_j = \left. \frac{\partial \phi_{\tilde{X}}(t)}{\partial t_j} \right|_{t=0}$ so that we may write

$$\begin{aligned} E\tilde{X} &= \frac{1}{i} \left(\left. \frac{\partial \phi_{\tilde{X}}(t)}{\partial t_i} \right|_{t=0} \right)_{m \times 1} \\ &= \frac{1}{i} \left(\left. \frac{\partial \phi_{\tilde{X}}(t)}{\partial t} \right|_{t=0} \right) \end{aligned}$$

We conclude the section on characteristic functions with the following useful lemma.

Lemma A.2.5. If \tilde{X} is a random $m \times n$ matrix, set $\tilde{Z} = A\tilde{X}B + C$ where A, B, C are fixed matrices, then

$$\phi_{\tilde{Z}}(U) = e^{i \operatorname{tr}(U^T C)} \phi_{\tilde{X}}[A^T U B^T]$$

Proof:

$$\begin{aligned} \phi_{\tilde{Z}}(U) &= E_{\tilde{Z}} e^{i \langle U, \tilde{Z} \rangle} \\ &= E_{\tilde{X}} e^{i \langle U, A\tilde{X}B + C \rangle} \\ &= E_{\tilde{X}} e^{i \operatorname{tr}[U^T (A\tilde{X}B + C)]} \\ &= e^{i \operatorname{tr} U^T C} E_{\tilde{X}} e^{i \langle (BU^T A)^T, \tilde{X} \rangle} \\ &= e^{i \operatorname{tr} U^T C} \phi_{\tilde{X}}[(BU^T A)^T] \\ &= e^{i \operatorname{tr} U^T C} \phi_{\tilde{X}}[A^T U B^T] \end{aligned}$$

A.3 Remarks Concerning Lemma 4.1.1

We recall the situation in which Lemma 4.1.1 was applied

$$\begin{aligned} f(u|R) &= (2\pi)^{-m/2} |R|^{-1/2} e^{-\frac{1}{2} u^T R u} \\ &= (2\pi)^{-m/2} e^{\sum_i r_{ii} W_{ii}(u) + \sum_{i < j} r_{ij} W_{ij}(u) - \psi(R)} \end{aligned}$$

where

$$W_{ii}(u) = -\frac{1}{2} u_i^2$$

$$W_{ij}(u) = -u_i u_j, \quad i < j$$

$$\psi(R) = -\frac{1}{2} \log |R|$$

Suppose the prior density is given by

$$\begin{aligned} f(R|\nu, \tau) &\propto |R|^{\nu-m-1/2} e^{-\frac{1}{2} \text{tr}(\tau R)} \\ &= e^{b^* \{ \sum_i r_{ii} a_{ii}^* + \sum_{i < j} r_{ij} a_{ij}^* - \psi(R) \}} \end{aligned}$$

where

$$b^* = \nu - m - 1$$

$$a_{ii}^* = -\frac{1}{2}$$

$$a_{ij}^* = 0$$

The posterior distribution of R is then given by

$$f(R | \mathbf{1}_U, \dots, \mathbf{n}_U, \nu, \tau) \propto e^{b\{\sum_i r_{ii} a_{ii} + \sum_{i < j} r_{ij} a_{ij}\} - \psi(R)}$$

where

$$b = n + b^*$$

$$a_{ii} = \frac{\bar{w}_{ii} + b^* a_{ii}^*}{n + b^*}$$

$$a_{ij} = \frac{\bar{w}_{ij} + b^* a_{ij}^*}{n + b^*}$$

$$\bar{w}_{ij} = \frac{1}{n} \sum_{t=1}^n w_{ij}^{(t_u)}$$

We note that the conditional, prior and posterior distributions have a natural parametrization. The condition necessary to write

$$\psi_i(R) = E[W_i(\underline{y}) | R]$$

is that we must pass a derivative under an integral sign which can be done [Barndorff-Nielsen (1973)] if the natural parameter space is open.

If we let $\mathcal{S} := \{\text{all symmetric } m \times m \text{ matrices}\}$ then \mathcal{S} may be viewed as $\mathbb{R}^{m(m+1)/2}$. We let

$\mathcal{V} := \{\text{all positive definite matrices}\}$. \mathcal{V} is a subset of \mathcal{S} so that we view \mathcal{V} as being a subset of $\mathbb{R}^{m(m+1)/2}$ and we need to show that \mathcal{V} is open in $\mathbb{R}^{m(m+1)/2}$.

Consider the following characterization of positive definiteness for matrices found in Bellman (1960).

Lemma A.3.1. Let A be an $m \times m$ symmetric matrix (a_{ij}) . Let $D_k := |(a_{ij})|_{i,j=1,\dots,k}$. A is positive definite if and only if $D_k > 0$, $k = 1, \dots, m$.

We note that $D_k : \mathcal{S} \rightarrow \mathbb{R}^1$ is a continuous function for $k = 1, \dots, m$ and that we may characterize \mathcal{V} as that subset of \mathcal{S} for which $D_k > 0$, $k = 1, \dots, m$. That is,

$$\begin{aligned} \mathcal{V} &= \{R \in \mathcal{S} : D_1(R) > 0, D_2(R) > 0, \dots, D_m(R) > 0\} \\ &= \{R \in \mathcal{S} : D_1(R) > 0\} \cap \{R \in \mathcal{S} : D_2(R) > 0\} \cap \dots \\ &\quad \cap \{R \in \mathcal{S} : D_m(R) > 0\} \\ &= \{D_1^{-1}(0, \infty)\} \cap \{D_2^{-1}(0, \infty)\} \cap \dots \cap \{D_m^{-1}(0, \infty)\} \end{aligned}$$

Now D_k being continuous for $k = 1, \dots, m$ implies that each of the sets $\{D_k^{-1}(0, \infty)\}$ is open and thus \mathcal{V} is the finite intersection of open sets and is therefore open.

The regularity condition given by Morgan (1969) sufficient for the validity of the statement $a_{ij} = E[\psi_{ij}(R)]$ that requires comment involves upper and lower limits of integration. The condition arises in the proof of his theorem where the integration by parts formula

$$\int u dv = uv - \int v du$$

is used. Roughly speaking the regularity condition requires that the uv term be zero when evaluated at the upper and lower limits of integration.

To be more precise in our situation, we consider a point R in \mathcal{V} as a vector r in $\mathbb{R}^{m(m+1)/2}$ with components r_{ij} , $1 \leq i \leq j \leq m$. We define the vector \mathcal{N}_{ij} by $\mathcal{N}_{ij} = r$ with the ij component deleted and view \mathcal{N}_{ij} as a vector in $\mathbb{R}^{m(m+1)/2 - 1}$. We define the functions corresponding to upper and lower limits by

$$u_{ij}(s) = \sup\{r_{ij} : R \in \mathcal{V}, \mathcal{N}_{ij} = s\}$$

$$l_{ij}(s) = \inf\{r_{ij} : R \in \mathcal{V}, \mathcal{N}_{ij} = s\}$$

We note that the domain of $u_{ij}(\cdot)$ and $l_{ij}(\cdot)$ is $\mathbb{R}^{\frac{m(m+1)}{2} - 1}$.

The required regularity condition is:

$$\lim_{r_{ij} \rightarrow l_{ij}(s)} e^{b^* \left\{ -\frac{1}{2} \sum_i r_{ii} + \frac{1}{2} \log |R| \right\}} = \lim_{r_{ij} \rightarrow u_{ij}(s)} e^{b^* \left\{ -\frac{1}{2} \sum_i r_{ii} + \frac{1}{2} \log |R| \right\}}$$

for $1 \leq i \leq j \leq m$.

We now show that this condition is satisfied. Our verification rests on Lemma A.3.1. For notation let $D_{i:j}$ be the determinant of the matrix obtained from R by the deletion of the i th row

and j th column and $D_{ij:j\ell}$ the determinant obtained from R by the deletion of the i th and k th rows and the j th and ℓ th columns.

Consider

$$u_{11}(s) = \sup\{r_{11} : R \in \mathcal{U}, r_{11} = s\}$$

Since $R \in \mathcal{U}$, $|R| > 0$. If we expand $|R|$ along the first column of R we have

$$r_{11}D_{1:1} - r_{12}D_{1:2} + r_{13}D_{1:3} - \dots + (-1)^{m+1}r_{1m}D_{1:m} = |R| > 0$$

We now note that r_{11} is involved in only the first term of the above sum. That is $r_{11}D_{1:1}$ is the only term involving r_{11} . Since $D_{1:1} > 0$, it now follows that with r_{ij} fixed for $(i,j) \neq (1,1)$, we have that $u_{11}(s) = \infty$. We now note that this same argument is valid for any diagonal element so that

$$u_{ii}(s) = \infty, \quad i = 1, \dots, m$$

Now consider

$$l_{mm}(s) = \inf\{r_{mm} : R \in \mathcal{U}, r_{mm} = s\}.$$

We note that the constraints D_k given by Lemma A.3.1 do not involve r_{mm} for $k = 1, \dots, m-1$. The only constraint on r_{mm} is the constraint $D_m = |R| > 0$. Therefore,

$$\inf\{r_{mm} : R \in \mathcal{V}, \mu_{mm} = s\} = \bar{r}_{mm}$$

where \bar{r}_{mm} is that value of r_{mm} for which $|R| = 0$ and we have

$$l_{mm}(s) = \bar{r}_{mm}.$$

By an even number of row and column interchanges any diagonal element can be thought of as the (m, m) element. Therefore,

$$l_{ii}(s) = \bar{r}_{ii}, \quad i = 1, \dots, m.$$

We now consider the "off-diagonal" terms. That is, consider

$$u_{ij}(s) = \sup\{r_{ij} : R \in \mathcal{V}, \mu_{ij} = s\}$$

and

$$l_{ij}(s) = \inf\{r_{ij} : R \in \mathcal{V}, \mu_{ij} = s\}.$$

We expand R by column j :

$$(A.1) \quad |R| = (-1)^{1+j} r_{ij} D_{1;j} + (-1)^{2+j} r_{2j} D_{2;j} + \dots + (-1)^{i+j} r_{ij} D_{i;j} \\ + \dots + (-1)^{m+j} r_{mj} D_{m;j} > 0.$$

We note that each matrix corresponding to $D_{t;j}$ for $t \neq i$ contains the component r_{ij} at most once. Therefore, the terms $r_{ij} D_{t;j}$ are linear or constant in r_{ij} for $t \neq i$ and quadratic in r_{ij} for $t = i$.

We now consider the quadratic term $(-1)^{i+j} r_{ij} D_{i;j}$. $D_{i;j}$ is itself a determinant. The matrix corresponding to this determinant contains the element r_{ji} . We now expand $D_{i;j}$ along the column containing r_{ji} and consider only that term involving r_{ji} . Namely, $(-1)^{i+j-1} r_{ji} D_{ij:ji}$. The quadratic term in (A. 1) now may be written as

$$(-1)^{i+j} r_{ij} (-1)^{i+j-1} r_{ji} D_{ij:ji} = -r_{ij}^2 D_{ij:ij}$$

Now, $D_{ij:ij}$ is the determinant of the matrix R with the i and j rows and columns deleted and is therefore positive.

We now have that the expansion given in (A. 1) is quadratic in r_{ij} with the coefficient of the quadratic term being negative. That is the quadratic is concave. Therefore,

$$\inf\{r_{ij} : R \in \mathcal{V}, \mathcal{N}_{ij} = s\} = \underline{r}_{ij}$$

and

$$\sup\{r_{ij} : R \in \mathcal{V}, \mathcal{N}_{ij} = s\} = \bar{r}_{ij}$$

where \underline{r}_{ij} and \bar{r}_{ij} are respectively the smallest and largest values of r_{ij} for which $|R| = 0$ when the remaining elements in R are fixed.

From the relationship,

$$e^{b^* \left\{ -\frac{1}{2} \sum_i r_{ii} + \frac{1}{2} \log |R| \right\}} = |R|^{b^*/2} e^{b^* \left\{ -\frac{1}{2} \sum_i r_{ii} \right\}}$$

and the values of $u_{ij}(s)$ and $l_{ij}(s)$ as defined above, the regularity condition follows. \square