

# A Corrected Criterion for Selecting the Optimum Number of Principal Components

Hannes Kazianka and Jürgen Pilz  
Institute of Statistics, University of Klagenfurt

**Abstract:** Determining the optimum number of components to be retained is a key problem in principal component analysis (PCA). Besides the rule of thumb estimates there exist several sophisticated methods for automatically selecting the dimensionality of the data. Based on the probabilistic PCA model Minka (2001) proposed an approximate Bayesian model selection criterion. In this paper we correct this criterion and present a modified version. We compare the novel criterion with various other approaches in a simulation study. Furthermore, we use it for finding the optimum number of principal components in hyper-spectral skin cancer images.

**Zusammenfassung:** Ein zentrales Problem der Hauptkomponentenanalyse (PCA) ist es, die Anzahl an wichtigen Komponenten zu wählen. Neben den bekannten Faustregeln gibt es verschiedene fortgeschrittene Methoden zur automatischen Wahl der optimalen Dimensionalität der Daten. Beispielsweise wurde von Minka (2001), basierend auf dem probabilistischen PCA Modell, ein approximativ Bayes'sches Kriterium zur Modellwahl vorgeschlagen. In der vorliegenden Arbeit korrigieren wir dieses Kriterium und stellen eine modifizierte Version vor. Das neue Kriterium wird mit anderen Verfahren in einer Simulationsstudie verglichen. Weiters wird es zum Finden der optimalen Anzahl an Hauptkomponenten in hyper-spektralen Hautkrebsdaten verwendet.

**Keywords:** PCA, Probabilistic PCA, Model Selection.

## 1 Introduction

Principal component analysis (PCA) was introduced by Pearson (1901) and is a linear orthogonal data transformation. The data is transformed in such a way that in the new coordinate system the components are uncorrelated and sorted in a descending order according to their variance. PCA is a standard tool for reducing multidimensional data sets to lower dimensions for further statistical analysis by omitting higher-order components. Determining the optimum number of components to be retained is a major problem in PCA. Besides certain rule of thumb estimates there exist numerous sophisticated methods for automatically selecting the dimensionality of the data. For the Gaussian probabilistic PCA model of Tipping and Bishop (1999) an approximate Bayesian model selection criterion was proposed by Minka (2001). Although mentioned in books and used in various applications, we found some mistakes in its proof. In the following we correct this criterion and present a modified version.

The paper is organized as follows. Section 2 presents the probabilistic PCA model and mentions its basic properties. In Section 3 the model selection criterion proposed by

Minka (2001) is corrected. Section 4 compares the corrected criterion with several other approaches in a simulation study while Section 5 presents the results when it is applied to hyper-spectral skin cancer image data. Section 6 is devoted to conclusions.

## 2 Probabilistic PCA

A probabilistic model for PCA was introduced by Tipping and Bishop (1999) and is a special case of the factor analysis model (see Basilevsky, 1994). It assumes that the random vector  $\mathbf{x}$  is a linear combination of basis vectors and an additive noise term,

$$\mathbf{x} = \sum_{j=1}^k \mathbf{h}_j w_j + \boldsymbol{\mu} + \boldsymbol{\varepsilon} = \mathbf{H}\mathbf{w} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}, \quad (1)$$

where  $\mathbf{x}$  has length  $d$ ,  $\mathbf{w}$  has smaller length  $k$ ,  $\boldsymbol{\mu}$  is the mean of  $\mathbf{x}$  and the noise model is isotropic Gaussian,  $p(\boldsymbol{\varepsilon}) \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$ . Moreover, the density of  $\mathbf{w}$  is assumed to be standard Gaussian,  $p(\mathbf{w}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$ . Therefore, the probability of observing  $\mathbf{x}$  follows a normal distribution,  $p(\mathbf{x} | \mathbf{w}, \mathbf{H}, \boldsymbol{\mu}, \sigma^2) \sim \mathcal{N}(\mathbf{H}\mathbf{w} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}_d)$ , and the marginal distribution of  $\mathbf{x}$  is  $p(\mathbf{x} | \mathbf{H}, \boldsymbol{\mu}, \sigma^2) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{H}\mathbf{H}^T + \sigma^2 \mathbf{I}_d)$ . With this, the likelihood of a data set,  $\mathcal{D} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ , is

$$p(\mathcal{D} | \mathbf{H}, \boldsymbol{\mu}, \sigma^2) = (2\pi)^{-\frac{nd}{2}} |\mathbf{H}\mathbf{H}^T + \sigma^2 \mathbf{I}_d|^{-\frac{n}{2}} \exp \left\{ -\frac{n}{2} \text{tr} [(\mathbf{H}\mathbf{H}^T + \sigma^2 \mathbf{I}_d)^{-1} \mathbf{S}] \right\}, \quad (2)$$

where  $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$ . Consequently, the maximum-likelihood estimate,  $\hat{\boldsymbol{\mu}}$ , for  $\boldsymbol{\mu}$  is the arithmetic mean. Defining  $\hat{\mathbf{S}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$  and  $\mathbf{C} = \mathbf{H}\mathbf{H}^T + \sigma^2 \mathbf{I}_d$  the log-likelihood evaluated at  $\hat{\boldsymbol{\mu}}$  is

$$\log p(\mathcal{D} | \mathbf{H}, \hat{\boldsymbol{\mu}}, \sigma^2) = -\frac{n}{2} \left( d \log(2\pi) + \log |\mathbf{C}| + \text{tr} [\mathbf{C}^{-1} \hat{\mathbf{S}}] \right). \quad (3)$$

Next we show how the maximum-likelihood estimates for  $\mathbf{H}$  and  $\sigma^2$  can be obtained. Using the symmetry of  $\mathbf{C}$  and  $\hat{\mathbf{S}}$  the gradient of (3) with respect to  $\mathbf{H}$  is

$$\begin{aligned} \frac{\partial \log p(\mathcal{D} | \mathbf{H}, \hat{\boldsymbol{\mu}}, \sigma^2)}{\partial \mathbf{H}} &= \frac{n}{2} \left( (\mathbf{C}^{-1} \hat{\mathbf{S}} \mathbf{C}^{-1})^T - \mathbf{C}^{-T} \right) 2\mathbf{H} \\ &= n \left( \mathbf{C}^{-1} \hat{\mathbf{S}} \mathbf{C}^{-1} \mathbf{H} - \mathbf{C}^{-1} \mathbf{H} \right). \end{aligned}$$

Setting the gradient to zero and multiplying with  $\mathbf{C}$  from the left yields

$$\hat{\mathbf{S}} \mathbf{C}^{-1} \mathbf{H} = \mathbf{H}. \quad (4)$$

We are interested in the solutions of (4) for  $\mathbf{H}$ . A possible solution would be  $\hat{\mathbf{H}} = \mathbf{0}$  which corresponds to a minimum of the likelihood. A second solution is  $\mathbf{C} = \hat{\mathbf{S}}$ , where the covariance model is exact. In this case,  $\mathbf{H}\mathbf{H}^T = \hat{\mathbf{S}} - \sigma^2 \mathbf{I}_d$ , which has the solution  $\hat{\mathbf{H}} = \mathbf{U}(\mathbf{L} - \sigma^2 \mathbf{I}_d)^{1/2} \mathbf{R}$ , where  $\mathbf{U} \in \mathbb{R}^{d \times d}$  is the orthogonal matrix with the eigenvectors of  $\hat{\mathbf{S}}$  in its columns,  $\mathbf{L} \in \mathbb{R}^{d \times d}$  is the diagonal matrix containing the corresponding eigenvalues and  $\mathbf{R} \in \mathbb{R}^{d \times d}$  is an arbitrary orthogonal matrix. As Tipping and Bishop (1999)

note, having an exact covariance model is generally undesirable because this information is effectively discarded in the dimensionality reduction process anyway.

For the general and more interesting case consider all solutions of (4), but  $\hat{\mathbf{H}} \neq \mathbf{0}$  and  $\mathbf{C} \neq \hat{\mathbf{S}}$ . Using the singular value decomposition we know that the estimator can be written as  $\hat{\mathbf{H}} = \mathbf{U} \mathbf{D} \mathbf{R}^T$ , where now  $\mathbf{U} \in \mathbb{R}^{d \times k}$  is a matrix with orthonormal columns  $(\mathbf{u}_1, \dots, \mathbf{u}_k)$ ,  $\mathbf{D} \in \mathbb{R}^{k \times k}$  is the diagonal matrix containing the singular values  $d_1, \dots, d_k$ , and  $\mathbf{R} \in \mathbb{R}^{k \times k}$  is again an arbitrary orthogonal matrix. We deduce from (4) that

$$\hat{\mathbf{S}} (\mathbf{U} \mathbf{D}^2 \mathbf{U}^T + \sigma^2 \mathbf{I}_d)^{-1} \mathbf{U} \mathbf{D} \mathbf{R}^T = \mathbf{U} \mathbf{D} \mathbf{R}^T \implies \hat{\mathbf{S}} \mathbf{U} \mathbf{D} = \mathbf{U} (\mathbf{D}^2 + \sigma^2 \mathbf{I}_k) \mathbf{D}.$$

If  $d_i = 0$ ,  $i = 1, \dots, k$ , the vector  $\mathbf{u}_i$  can be chosen arbitrarily such that  $\mathbf{U}$  has orthonormal columns. If  $d_i \neq 0$ , we have  $\hat{\mathbf{S}} \mathbf{u}_i = (\sigma^2 + d_i^2) \mathbf{u}_i$  which means that each column of  $\mathbf{U}$  has to be an eigenvector of  $\hat{\mathbf{S}}$  corresponding to the eigenvalue  $l_i = \sigma^2 + d_i^2$ . Therefore,  $d_i = (l_i - \sigma^2)^{1/2}$  and all maximum-likelihood solutions for  $\mathbf{H}$  must have the form

$$\hat{\mathbf{H}} = \mathbf{U} (\mathbf{L} - \sigma^2 \mathbf{I}_k)^{1/2} \mathbf{R}, \quad (5)$$

where now  $\mathbf{L} \in \mathbb{R}^{k \times k}$  is the diagonal matrix with entries  $l_i$  and  $\mathbf{U}$  contains the corresponding eigenvectors. With this, we can simplify

$$\begin{aligned} \left| \hat{\mathbf{H}} \hat{\mathbf{H}}^T + \sigma^2 \mathbf{I}_d \right| &= \left| \mathbf{U} (\mathbf{L} - \sigma^2 \mathbf{I}_k) \mathbf{U}^T + \sigma^2 \mathbf{I}_d \right| = \sigma^{2(d-k)} \left| (\mathbf{L} - \sigma^2 \mathbf{I}_k) \mathbf{U}^T \mathbf{U} + \sigma^2 \mathbf{I}_k \right| \\ &= \sigma^{2(d-k)} \left| \mathbf{L} - \sigma^2 \mathbf{I}_k + \sigma^2 \mathbf{I}_k \right| = \sigma^{2(d-k)} |\mathbf{L}| = \sigma^{2(d-q)} \prod_{i=1}^q l_i, \end{aligned} \quad (6)$$

where  $q$  is the number of non-zero  $d_i$  and equals the rank of  $\hat{\mathbf{H}}$ . Therefore, the eigenvalues  $l_1, \dots, l_q$  correspond to those components retained in PCA and  $l_{q+1}, \dots, l_d$  correspond to those discarded. Using the Woodbury matrix identity we find that

$$\begin{aligned} \left( \hat{\mathbf{H}} \hat{\mathbf{H}}^T + \sigma^2 \mathbf{I}_d \right)^{-1} &= \sigma^{-2} \mathbf{I}_d - \sigma^{-2} \mathbf{U} \left( (\mathbf{L} - \sigma^2 \mathbf{I}_k)^{-1} + \sigma^{-2} \mathbf{I}_k \right)^{-1} \mathbf{U}^T \sigma^{-2} \\ &= \sigma^{-2} \mathbf{I}_d + \mathbf{U} (\mathbf{L}^{-1} - \sigma^{-2} \mathbf{I}_k) \mathbf{U}^T. \end{aligned} \quad (7)$$

Since  $\mathbf{U}^T \hat{\mathbf{S}} \mathbf{U} = \mathbf{L}$ , we get that

$$\begin{aligned} \text{tr} \left[ \left( \hat{\mathbf{H}} \hat{\mathbf{H}}^T + \sigma^2 \mathbf{I}_d \right)^{-1} \hat{\mathbf{S}} \right] &= \text{tr} \left[ \sigma^{-2} \hat{\mathbf{S}} \right] + \text{tr} \left[ \mathbf{U} (\mathbf{L}^{-1} - \sigma^{-2} \mathbf{I}_k) \mathbf{U}^T \hat{\mathbf{S}} \right] \\ &= \frac{1}{\sigma^2} \sum_{i=1}^d l_i + \sum_{i=1}^q \left( \frac{1}{l_i} - \frac{1}{\sigma^2} \right) l_i = q + \frac{1}{\sigma^2} \sum_{i=q+1}^d l_i. \end{aligned}$$

Plugging the latter result and (6) into the log-likelihood (3) we get

$$\log p \left( \mathcal{D} \mid \hat{\mathbf{H}}, \hat{\boldsymbol{\mu}}, \sigma^2 \right) = -\frac{n}{2} \left( \sum_{i=1}^q \log l_i + (d-q) \log \sigma^2 + d \log 2\pi + q + \frac{1}{\sigma^2} \sum_{i=q+1}^d l_i \right). \quad (8)$$

The maximum-likelihood estimate for  $\sigma^2$  turns out to be the average of the discarded eigenvalues:

$$\frac{n(d-q)}{\sigma} + \frac{n}{\sigma^3} \sum_{i=q+1}^d l_i = 0 \quad \implies \quad \hat{\sigma}^2 = \frac{1}{d-q} \sum_{i=q+1}^d l_i. \quad (9)$$

Substituting  $\hat{\sigma}^2$  into the log-likelihood (8) we arrive at

$$\log p(\mathcal{D} | \hat{\mathbf{H}}, \hat{\boldsymbol{\mu}}, \hat{\sigma}^2) = -\frac{n}{2} \left( \sum_{i=1}^q \log l_i + (d-q) \log \left( \frac{1}{d-q} \sum_{i=q+1}^d l_i \right) + d \log 2\pi + d \right). \quad (10)$$

In order to find those eigenvalues of  $\hat{\mathbf{S}}$  that correspond to discarded components we have to maximize the above log-likelihood with respect to their choice. The sum of all eigenvalues is constant implying that maximization of (10) with respect to the eigenvalues is equivalent to minimizing

$$\min_{l_i, i=q+1, \dots, d} \log \left( \frac{1}{d-q} \sum_{i=q+1}^d l_i \right) - \frac{1}{d-q} \sum_{i=q+1}^d \log l_i.$$

The above function is minimized when all  $l_i$ ,  $i = q+1, \dots, d$ , have the same value. Hence, the discarded eigenvalues are adjacent in the ordered spectrum of  $\hat{\mathbf{S}}$ . From the fact that  $l_i > \sigma^2$ ,  $\forall i = 1, \dots, q$  and (9) we know that the smallest eigenvalue of  $\hat{\mathbf{S}}$  is definitely discarded. This implies that the smallest  $d-q$  eigenvalues of  $\hat{\mathbf{S}}$  are discarded and the top  $q$  eigenvalues are retained. Seen from this angle, the estimate for  $\sigma^2$  can be interpreted as the average information loss when reducing the dimension of the data.

### 3 Choosing the Number of Principal Components Using Bayesian Model Selection

For the probabilistic PCA model with positive definite covariance matrix Minka (2001) developed a fast and efficient criterion to select the optimum number of principal components. It is based on Bayesian model selection and is therefore sometimes called the MIBS criterion (*Minka Bayesian model Selection*) in the literature. A variety of books refer to this criterion, for example Cichocki and Amari (2002), Smidl and Quinn (2005) and Bishop (2008). It showed good results for PCA and also for independent component analysis. Although the MIBS criterion is widely used, we found some errors in its proof. In this section we give a corrected proof and propose a modified criterion for selecting the optimum number of principal components.

#### 3.1 Choice of the Prior

In contrast to maximum-likelihood estimation performed in Section 2 we have to assign priors to the parameters  $\boldsymbol{\mu}$ ,  $\mathbf{H}$  and  $\sigma^2$  in the probabilistic PCA model (1) when we want

to go the Bayesian way. For the mean  $\boldsymbol{\mu}$  we use a non-informative, flat prior,  $p(\boldsymbol{\mu}) \propto 1$ . Integrating the likelihood (2) with respect to  $\boldsymbol{\mu}$  yields

$$p(\mathcal{D} | \mathbf{H}, \sigma^2) = n^{-\frac{d}{2}} (2\pi)^{-\frac{d(n-1)}{2}} |\mathbf{H}\mathbf{H}^T + \sigma^2 \mathbf{I}_d|^{-\frac{n-1}{2}} \exp \left\{ -\frac{n}{2} \text{tr} \left[ (\mathbf{H}\mathbf{H}^T + \sigma^2 \mathbf{I}_d)^{-1} \hat{\mathbf{S}} \right] \right\}. \quad (11)$$

The prior for  $\mathbf{H}$  is constructed similar to decomposition (5),  $\mathbf{H} = \mathbf{U}(\boldsymbol{\Lambda} - \sigma^2 \mathbf{I}_k)^{1/2} \mathbf{R}$ . In contrast to Section 2, here the orthogonal matrix  $\mathbf{R}$  and the matrix  $\mathbf{U}$ , which has orthonormal columns, are model parameters.  $\boldsymbol{\Lambda}$  is a diagonal matrix containing parameters  $\lambda_i$ ,  $i = 1, \dots, k$ . Minka (2001) found a conjugate prior for  $\mathbf{U}$ ,  $\boldsymbol{\Lambda}$ ,  $\mathbf{R}$  and  $\sigma^2$ . This prior is parameterized by a single parameter  $\alpha > 0$  and by applying simplifications similar to (6) and (7) it can be written as

$$\begin{aligned} p(\mathbf{H}, \sigma^2) &= p(\mathbf{U}, \boldsymbol{\Lambda}, \mathbf{R}, \sigma^2) \\ &\propto |\mathbf{H}\mathbf{H}^T + \sigma^2 \mathbf{I}_d|^{-\frac{\alpha+2}{2}} \exp \left\{ -\frac{\alpha}{2} \text{tr} \left[ (\mathbf{H}\mathbf{H}^T + \sigma^2 \mathbf{I}_d)^{-1} \right] \right\} \\ &\propto |\boldsymbol{\Lambda}|^{-\frac{\alpha+2}{2}} \sigma^{-(\alpha+2)(d-k)} \exp \left\{ -\frac{\alpha}{2} \left( \text{tr} [\boldsymbol{\Lambda}^{-1}] + \frac{d-k}{\sigma^2} \right) \right\}. \end{aligned} \quad (12)$$

With this definition of the prior the parameters are a-priori independent since their joint prior density (12) factors into the product of the marginal prior densities for each of the parameters:

$$p(\mathbf{U}, \boldsymbol{\Lambda}, \mathbf{R}, \sigma^2) = p(\sigma^2) p(\mathbf{U}) p(\mathbf{R}) \prod_{i=1}^k p(\lambda_i),$$

where  $p(\sigma^2) \sim \chi^{-2}(\alpha(d-k), (\alpha+2)(d-k)-2)$  and  $p(\lambda_i) \sim \chi^{-2}(\alpha, \alpha)$  are scaled inverse- $\chi^2$  priors. To be least informative the prior for  $\mathbf{U}$  is chosen as the reciprocal of the area of the orthonormally constrained subset of the Cartesian product of  $k$   $d$ -dimensional unit hyperballs, known as the Stiefel manifold (see Khatri and Mardia, 1977 and James, 1954):

$$p(\mathbf{U}) = 2^{-k} \pi^{\frac{1}{4}k(k-1-2d)} \prod_{i=1}^k \Gamma \left( \frac{d-i+1}{2} \right) \propto 1.$$

The prior for  $\mathbf{R}$  can also be assumed constant but as the variable does not appear in the likelihood (11) we can integrate it out.

### 3.2 The Corrected MIBS Criterion

Among the well-known statistical model selection criteria are Akaike's information criterion (AIC) and the minimum description length (MDL) criterion. These can be easily used for finding the optimum number of principal components. Assume we observe samples  $\mathcal{D} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$  from a zero-mean normal random vector. An appropriate number of principal components is the value of  $k = 1, \dots, d$  for which

$$\text{AIC}(k) = -2n(d-k) \log \rho(k) + 2k(2d-k), \quad (13)$$

$$\text{MDL}(k) = -n(d-k) \log \rho(k) + \frac{k}{2}(2d-k) \log n, \quad (14)$$

is minimized. In (13) and (14) the function  $\rho(k)$  is

$$\rho(k) = \frac{(l_{k+1}l_{k+2} \cdots l_d)^{\frac{1}{d-k}}}{\frac{1}{d-k}(l_{k+1} + l_{k+2} + \cdots + l_d)}.$$

The criterion that was recommended in Minka (2001) has shown reasonable performance in various applications and is fast to evaluate (see Cichocki and Amari, 2002 and Leonowicz, Karvanen, Tanaka, and Rezmer, 2004):

$$p(\mathcal{D} | k) \approx 2^{\frac{m-k}{2}} \rho(k)^{-n(d-k)} |\mathbf{A}|^{-\frac{1}{2}} n^{-\frac{k}{2}} \prod_{i=1}^k \Gamma\left(\frac{d-i+1}{2}\right) \prod_{i=1}^k l_i^{-\frac{n}{2}} \prod_{i=k+1}^d l_i^n, \quad (15)$$

where  $m = dk - k(k+1)/2$  and  $|\mathbf{A}|$  is defined as

$$|\mathbf{A}| = n^m \prod_{i=1}^k \prod_{j=i+1}^d (\tilde{l}_j^{-1} - \tilde{l}_i^{-1}) (l_i - l_j) \quad \text{where} \quad \tilde{l}_i = \begin{cases} l_i, & \text{if } i \leq k, \\ \frac{1}{d-k} \sum_{j=k+1}^d l_j, & \text{if } k < i \leq d. \end{cases}$$

The number of principal components to retain is taken to be the value of  $k = 1, \dots, d$  for which  $p(\mathcal{D} | k)$  is maximized. Dropping all the terms that do not grow with  $n$  a BIC approximation was also proposed in Minka (2001)

$$p(\mathcal{D} | k) \approx n^{-\frac{m+k}{2}} \rho(k)^{-n(d-k)} \prod_{i=1}^k l_i^{-\frac{n}{2}} \prod_{i=k+1}^d l_i^n. \quad (16)$$

The following theorem corrects the existing MIBS criterion (15) and serves as a guideline for easy implementation of the novel criterion proposed therein.

**Theorem:** Consider the probabilistic principal component model (1) with a prior for the parameters  $\mathbf{H}$ ,  $\boldsymbol{\mu}$ , and  $\sigma^2$  as discussed in Section 3.1

$$p(\mathbf{H}, \boldsymbol{\mu}, \sigma^2) \propto |\mathbf{H}\mathbf{H}^T + \sigma^2 \mathbf{I}_d|^{-\frac{\alpha+2}{2}} \exp\left\{-\frac{\alpha}{2} \text{tr}\left[(\mathbf{H}\mathbf{H}^T + \sigma^2 \mathbf{I}_d)^{-1}\right]\right\}, \quad (17)$$

where  $\alpha > 0$  is the prior parameter. By applying Laplace approximation the marginal likelihood of  $\mathcal{D} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$  for a model with  $k < d$  principal components is

$$p(\mathcal{D} | k) \approx 2^k c_k |\hat{\boldsymbol{\Lambda}}|^{-\frac{N}{2}+1} \hat{\sigma}^{-N(d-k)+2} \exp\left\{-\frac{Nd}{2} + k + 1\right\} (2\pi)^{\frac{m+k+1}{2}} (|\mathbf{A}_U| |\mathbf{A}_\Lambda| |\mathbf{A}_{\sigma^2}|)^{-\frac{1}{2}}, \quad (18)$$

where  $m = dk - k(k+1)/2$ ,  $N = n + 1 + \alpha$ ,  $\hat{\boldsymbol{\Lambda}}$  is a diagonal matrix with entries  $\hat{\lambda}_i$ ,  $i = 1, \dots, k$ , and

$$c_k = \frac{n^{-\frac{d}{2}} (2\pi)^{-\frac{(n-1)d}{2}} \pi^{\frac{1}{4}k(k-1-2d)}}{2^k \Gamma\left(\frac{(\alpha+2)(d-k)}{2} - 1\right) \Gamma\left(\frac{\alpha}{2}\right)^k} \left(\frac{\alpha(d-k)}{2}\right)^{\frac{(\alpha+2)(d-k)-2}{2}} \left(\frac{\alpha}{2}\right)^{\frac{k\alpha}{2}} \prod_{i=1}^k \Gamma\left(\frac{d-i+1}{2}\right). \quad (19)$$

Moreover, if  $l_1, \dots, l_d$  denote the eigenvalues of  $\hat{\mathbf{S}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$ , where  $\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ , we have

$$\begin{aligned}
|\mathbf{A}_{\sigma^2}| &= \frac{N(d-k)-2}{2}, & \hat{\sigma}^2 &= \frac{n \sum_{i=k+1}^d l_i}{N(d-k)-2}, \\
|\mathbf{A}_{\Lambda}| &= \left(\frac{N}{2} - 1\right)^k, & \hat{\lambda}_i &= \frac{nl_i + \alpha}{N-2},
\end{aligned}$$

$$|\mathbf{A}_{\mathbf{U}}| = n^m \prod_{i=1}^k \prod_{j=i+1}^d \left(\tilde{\lambda}_j^{-1} - \tilde{\lambda}_i^{-1}\right) (l_i - l_j) \quad \text{where} \quad \tilde{\lambda}_i = \begin{cases} \hat{\lambda}_i, & \text{if } i \leq k, \\ \hat{\sigma}^2, & \text{if } k < i \leq d. \end{cases}$$

The value for  $\alpha$  is typically chosen very small to make the prior less informative. Large values for  $\alpha$  may lead to inferior results especially when the number of samples is small.

### 3.3 Proof of the Corrected MIBS Criterion

The marginal likelihood can be written as

$$\begin{aligned}
p(\mathcal{D} | k) &= c_k \int |\mathbf{H}\mathbf{H}^T + \sigma^2 \mathbf{I}_d|^{-\frac{N}{2}} \\
&\quad \times \exp \left\{ -\frac{n}{2} \text{tr} \left[ (\mathbf{H}\mathbf{H}^T + \sigma^2 \mathbf{I}_d)^{-1} \left( \hat{\mathbf{S}} + \frac{\alpha}{n} \mathbf{I}_d \right) \right] \right\} d\mathbf{U} d\Lambda d\sigma^2, \quad (20)
\end{aligned}$$

where  $N = n + 1 + \alpha$ . The multiplier  $c_k$  defined in (19) combines the constants of the prior (12) and the integrated likelihood (11). To evaluate the integral in (20) we use Laplace approximation as described in Lindley (1980). The matrix  $\mathbf{U} \in \mathbb{R}^{d \times k}$  has  $m = dk - k(k+1)/2$  free parameters since we have  $k(k+1)$  constraints because of the orthonormality condition for the columns of  $\mathbf{U}$ . Laplace approximation can therefore be written as

$$\int p(\mathcal{D}, \mathbf{U}, \Lambda, \sigma^2 | k) d\mathbf{U} d\Lambda d\sigma^2 \approx p(\mathcal{D}, \hat{\mathbf{U}}, \hat{\Lambda}, \hat{\sigma}^2 | k) (2\pi)^{\frac{k+m+1}{2}} |\mathbf{A}_{\mathbf{U}\Lambda\sigma^2}|^{-\frac{1}{2}}, \quad (21)$$

where  $(\hat{\mathbf{U}}, \hat{\Lambda}, \hat{\sigma}^2) = \arg \max_{\mathbf{U}, \Lambda, \sigma^2} p(\mathcal{D}, \mathbf{U}, \Lambda, \sigma^2 | k)$  and  $\mathbf{A}_{\mathbf{U}\Lambda\sigma^2}$  is the negative Hessian of the integrand at  $(\hat{\mathbf{U}}, \hat{\Lambda}, \hat{\sigma}^2)$ . It is important to know that Laplace approximation is a basis dependent method, so we have to choose an appropriate parameterization,  $(\mathbf{Z}, \hat{\Lambda}, \hat{\sigma}^2)$ , for  $(\mathbf{U}, \Lambda, \sigma^2)$ . In the case of  $\Lambda$  and  $\sigma^2$  it is easy to guess a good parameterization, since Laplace approximation works better when using integrals over the whole real axis than just over the positive real numbers. Therefore, one can use  $\hat{\sigma}^2 = \log \sigma^2$  and  $\hat{\lambda}_i = \log \lambda_i$  for all  $i = 1, \dots, k$ , where  $\hat{\lambda}_i$  are the diagonal elements of  $\hat{\Lambda}$ . The Jacobian matrices for the transformations are

$$\mathbf{J}_{\Lambda} = \frac{\partial(\lambda_1, \dots, \lambda_k)}{\partial(\hat{\lambda}_1, \dots, \hat{\lambda}_k)} = \begin{pmatrix} \exp\{\hat{\lambda}_1\} & & \\ & \ddots & \\ & & \exp\{\hat{\lambda}_k\} \end{pmatrix} = \Lambda, \quad \mathbf{J}_{\sigma^2} = \sigma^2.$$

The matrix  $U$  is expressed in Euler vector coordinates

$$U = U_d \exp\{Z\} \begin{bmatrix} I_k \\ \mathbf{0} \end{bmatrix}, \tag{22}$$

where  $Z \in \mathbb{R}^{d \times d}$  is a skew-symmetric matrix of parameters and  $U_d \in \mathbb{R}^{d \times d}$  is a fixed orthogonal matrix. The free parameters of  $Z$  are the first  $k$  rows of the upper triangle:  $d(d-1)/2 - (d-k)(d-k-1)/2 = m$ , as desired. If  $\Lambda$  and  $\sigma^2$  are fixed, we can use the result in (7) and write  $p(\mathcal{D}, U | k, \Lambda, \sigma^2)$  as

$$p(\mathcal{D}, U | k, \Lambda, \sigma^2) \propto \exp \left\{ -\frac{n}{2} \text{tr} \left[ (\Lambda^{-1} - \sigma^{-2} I_k) U^T \hat{S} U \right] \right\},$$

which is known to be the von Mises-Fisher matrix distribution (see Khatri and Mardia, 1977). This density is maximized for  $U$  at  $\hat{U}$ , which contains the first  $k$  eigenvectors of  $\hat{S}$ . In the parameterization (22) this corresponds to the case where  $\hat{Z} = \mathbf{0}$  and  $U_d$  is equal to the matrix of eigenvectors of  $\hat{S}$ . Note, that the density has the same value if the sign of a column of  $U$  is changed. This can happen  $2^k$  times, so the density has  $2^k$  extreme points with the same function values and we need to multiply (21) with  $2^k$ . The determinant of the Jacobian for the transformation (22) at  $\hat{Z}$  is obviously equal to 1.

Using the same argument as in (7) it follows that

$$\text{tr} \left[ (\mathbf{H}\mathbf{H}^T + \sigma^2 I_d)^{-1} \left( \hat{S} + \frac{\alpha}{n} I_d \right) \right] = \frac{1}{\sigma^2} \text{tr} \left[ \hat{S} - U^T \hat{S} U \right] + \text{tr} \left[ \Lambda^{-1} U^T \hat{S} U + \frac{\alpha}{n} \Lambda^{-1} \right], \tag{23}$$

and if  $L = \text{diag}(l_1, \dots, l_k)$  denotes the diagonal matrix containing the first  $k$  eigenvalues of  $\hat{S}$ , we can rewrite  $\log p(\mathcal{D}, \hat{Z}, \hat{\Lambda}, \hat{\sigma}^2 | k)$  as

$$\begin{aligned} \log p(\mathcal{D}, \hat{Z}, \hat{\Lambda}, \hat{\sigma}^2 | k) &\propto \left( -\frac{N}{2} + 1 \right) \log |\Lambda| + \left( -\frac{N(d-k)}{2} + 1 \right) \log \sigma^2 \\ &\quad - \frac{n}{2} \left( \frac{1}{\sigma^2} \left( \text{tr}[\hat{S}] - \text{tr}[L] \right) + \text{tr}[\Lambda^{-1} L] + \frac{\alpha}{n} \text{tr}[\Lambda^{-1}] \right). \end{aligned}$$

The estimate for  $\Lambda$  can be calculated by differentiation of the latter equation:

$$\begin{aligned} \frac{\partial \log p(\mathcal{D}, \hat{Z}, \hat{\Lambda}, \hat{\sigma}^2 | k)}{\partial \hat{\lambda}_i} &= -\frac{N}{2} + 1 + \frac{n l_i}{2 \exp\{\hat{\lambda}_i\}} + \frac{\alpha}{2 \exp\{\hat{\lambda}_i\}} \\ &= -\frac{N}{2} + 1 + \frac{n l_i + \alpha}{2 \hat{\lambda}_i} \quad \forall i = 1, \dots, k. \end{aligned}$$

If we set this derivative to zero, we get that  $\exp\{\hat{\lambda}_i\} = \frac{n l_i + \alpha}{n - 1 + \alpha}$  or equivalently

$$\hat{\lambda}_i = \frac{n l_i + \alpha}{N - 2} \quad \forall i = 1, \dots, k.$$

The cross derivatives are zero and the second derivatives with respect to  $\hat{\lambda}_i$  are

$$\frac{\partial^2 \log p(\mathcal{D}, \hat{Z}, \hat{\Lambda}, \hat{\sigma}^2 | k)}{\partial \hat{\lambda}_i^2} = -\frac{n l_i + \alpha}{2 \exp\{\hat{\lambda}_i\}} = -\frac{n l_i + \alpha}{2 \hat{\lambda}_i} \quad \forall i = 1, \dots, k.$$



The negative Hessian of  $\log p(\mathcal{D}, \hat{\mathbf{Z}}, \hat{\Lambda}, \hat{\sigma}^2 | k)$  at  $\Lambda = \hat{\Lambda} = \text{diag}(\hat{\lambda}_1, \dots, \hat{\lambda}_k)$  can now be calculated as

$$|\mathbf{A}_\Lambda| = \left(\frac{N}{2} - 1\right)^k. \quad (24)$$

The estimate of  $\sigma^2$  is obtained analogously:

$$\begin{aligned} \frac{\partial \log p(\mathcal{D}, \hat{\mathbf{Z}}, \hat{\Lambda}, \hat{\sigma}^2 | k)}{\partial \hat{\sigma}^2} &= -\frac{N(d-k)}{2} + 1 + \frac{n}{2\hat{\sigma}^2} \left( \text{tr}[\hat{\mathbf{S}}] - \text{tr}[\mathbf{L}] \right), \\ \hat{\sigma}^2 &= \frac{n \left( \text{tr}[\hat{\mathbf{S}}] - \text{tr}[\mathbf{L}] \right)}{N(d-k) - 2} = \frac{n \sum_{i=k+1}^d l_i}{N(d-k) - 2} \\ \frac{\partial^2 \log p(\mathcal{D}, \hat{\mathbf{Z}}, \hat{\Lambda}, \hat{\sigma}^2 | k)}{\partial \hat{\sigma}^4} &= -\frac{n}{2\hat{\sigma}^2} \left( \text{tr}[\hat{\mathbf{S}}] - \text{tr}[\mathbf{L}] \right), \\ \frac{\partial^2 \log p(\mathcal{D}, \hat{\mathbf{Z}}, \hat{\Lambda}, \hat{\sigma}^2 | k)}{\partial \hat{\sigma}^4} \Bigg|_{\sigma^2 = \hat{\sigma}^2} &= -\frac{N(d-k) - 2}{2} = -|\mathbf{A}_{\sigma^2}|. \end{aligned} \quad (25)$$

Substituting the estimates  $\hat{\mathbf{Z}}$ ,  $\hat{\Lambda}$  and  $\hat{\sigma}^2$  into (23) we get that

$$(23) |_{\hat{\mathbf{Z}}, \hat{\Lambda}, \hat{\sigma}^2} = \frac{N(d-k) - 2}{n} + \sum_{i=1}^k \frac{l_i(N-2)}{nl_i + \alpha} + \frac{\alpha}{n} \sum_{i=1}^k \frac{N-2}{nl_i + \alpha} = \frac{Nd - 2k - 2}{n},$$

which leads to the following equation for  $p(\mathcal{D}, \hat{\mathbf{Z}}, \hat{\Lambda}, \hat{\sigma}^2 | k)$ :

$$p(\mathcal{D}, \hat{\mathbf{Z}}, \hat{\Lambda}, \hat{\sigma}^2 | k) = c_k |\hat{\Lambda}|^{-\frac{N}{2}+1} \hat{\sigma}^{-N(d-k)+2} \exp \left\{ -\frac{Nd}{2} + k + 1 \right\}.$$

Next, we have to evaluate the Hessian of  $\log p(\mathcal{D}, \mathbf{Z}, \hat{\Lambda}, \hat{\sigma}^2 | k)$  at  $\hat{\mathbf{Z}}$ . Using

$$\exp\{\mathbf{Z}\} = \sum_{i=0}^{\infty} \frac{\mathbf{Z}^i}{i!} = \mathbf{I}_d + \mathbf{Z} + \frac{\mathbf{Z}^2}{2} + \frac{\mathbf{Z}^3}{6} + \dots$$

and certain matrix differential rules we obtain the first and second differential of  $\mathbf{U}$ :

$$d\mathbf{U} = \mathbf{U}_d \left( d\mathbf{Z} + \frac{1}{2} (d\mathbf{Z}\mathbf{Z} + \mathbf{Z}d\mathbf{Z}) + \dots \right) \begin{bmatrix} \mathbf{I}_k \\ \mathbf{0} \end{bmatrix},$$

$$d\mathbf{U} |_{\mathbf{Z}=\mathbf{0}} = \mathbf{U}_d d\mathbf{Z} \begin{bmatrix} \mathbf{I}_k \\ \mathbf{0} \end{bmatrix}, \quad (26)$$

$$d^2\mathbf{U} = \mathbf{U}_d \left( d(d\mathbf{Z}) + \frac{1}{2} (d\mathbf{Z}d\mathbf{Z} + d\mathbf{Z}d\mathbf{Z}) + \dots \right) \begin{bmatrix} \mathbf{I}_k \\ \mathbf{0} \end{bmatrix},$$

$$d^2\mathbf{U} |_{\mathbf{Z}=\mathbf{0}} = \mathbf{U}_d (d\mathbf{Z})^2 \begin{bmatrix} \mathbf{I}_k \\ \mathbf{0} \end{bmatrix}. \quad (27)$$

The differential of  $\log p(\mathcal{D}, \mathbf{Z}, \hat{\Lambda}, \hat{\sigma}^2 | k)$  can be evaluated as

$$\begin{aligned} d \log p \left( \mathcal{D}, \mathbf{Z}, \hat{\Lambda}, \hat{\sigma}^2 | k \right) &= -\frac{n}{2} \text{tr} \left[ d \left( \left( \hat{\Lambda}^{-1} - \hat{\sigma}^{-2} \mathbf{I}_k \right) \mathbf{U}^T \hat{\mathbf{S}} \mathbf{U} \right) \right] \\ &= -\frac{n}{2} \text{tr} \left[ \left( \hat{\Lambda}^{-1} - \hat{\sigma}^{-2} \mathbf{I}_k \right) \left( d \mathbf{U}^T \hat{\mathbf{S}} \mathbf{U} + \mathbf{U}^T \hat{\mathbf{S}} d \mathbf{U} \right) \right] \\ &= -n \text{tr} \left[ \left( \hat{\Lambda}^{-1} - \hat{\sigma}^{-2} \mathbf{I}_k \right) \mathbf{U}^T \hat{\mathbf{S}} d \mathbf{U} \right]. \end{aligned} \tag{28}$$

The second differential is obtained by taking the differential of (28):

$$d^2 \log p \left( \mathcal{D}, \mathbf{Z}, \hat{\Lambda}, \hat{\sigma}^2 | k \right) = -n \text{tr} \left[ \left( \hat{\Lambda}^{-1} - \hat{\sigma}^{-2} \mathbf{I}_k \right) \left( d \mathbf{U}^T \hat{\mathbf{S}} d \mathbf{U} + \mathbf{U}^T \hat{\mathbf{S}} d^2 \mathbf{U} \right) \right]. \tag{29}$$

Since  $\mathbf{U}_d^T \hat{\mathbf{S}} \mathbf{U}_d = \mathbf{L}$  and  $\hat{\mathbf{U}}^T \hat{\mathbf{S}} \mathbf{U}_d = \begin{bmatrix} \mathbf{I}_k \\ \mathbf{0} \end{bmatrix}^T \mathbf{L}$ , just containing the top  $k$  rows of  $\mathbf{L}$ , we can rewrite (29) evaluated at  $\hat{\mathbf{Z}}$  by applying (26) and (27) as

$$-n \text{tr} \left[ \begin{bmatrix} \mathbf{I}_k \\ \mathbf{0} \end{bmatrix} \left( \hat{\Lambda}^{-1} - \hat{\sigma}^{-2} \mathbf{I}_k \right) \begin{bmatrix} \mathbf{I}_k \\ \mathbf{0} \end{bmatrix}^T \left( d \mathbf{Z}^T \mathbf{L} d \mathbf{Z} + \mathbf{L} d \mathbf{Z} d \mathbf{Z} \right) \right]. \tag{30}$$

If we define the diagonal matrix  $\mathbf{B} = \begin{bmatrix} \mathbf{I}_k \\ \mathbf{0} \end{bmatrix} \left( \hat{\Lambda}^{-1} - \hat{\sigma}^{-2} \mathbf{I}_k \right) \begin{bmatrix} \mathbf{I}_k \\ \mathbf{0} \end{bmatrix}^T$  and  $\mathbf{T} = \mathbf{B} d \mathbf{Z}^T + d \mathbf{Z} \mathbf{B}$ , we can simplify (30). Using the skew-symmetry of  $\mathbf{Z}$  we get  $\mathbf{T} = d \mathbf{Z} \mathbf{B} - \mathbf{B} d \mathbf{Z}$  and

$$d^2 \log p \left( \mathcal{D}, \mathbf{Z}, \hat{\Lambda}, \hat{\sigma}^2 | k \right) \Big|_{\mathbf{Z}=\mathbf{0}} = -n \text{tr} [\mathbf{T} \mathbf{L} d \mathbf{Z}]. \tag{31}$$

If we define a diagonal matrix  $\tilde{\Lambda}$  with entries  $\tilde{\lambda}_i, i = 1, \dots, d$ , as

$$\tilde{\Lambda} = \begin{pmatrix} \hat{\Lambda} & \mathbf{0} \\ \mathbf{0} & \hat{\sigma}^2 \mathbf{I}_k \end{pmatrix},$$

we get from the definition of  $\mathbf{B}$  and  $\mathbf{T}$  that the  $ij$ -th element of  $\mathbf{T}$  can be written as  $t_{ij} = (\tilde{\lambda}_j^{-1} - \tilde{\lambda}_i^{-1}) dz_{ij}$ . The next step is to show that (31) can be simplified to

$$d^2 \log p \left( \mathcal{D}, \mathbf{Z}, \hat{\Lambda}, \hat{\sigma}^2 | k \right) \Big|_{\mathbf{Z}=\mathbf{0}} = - \sum_{i=1}^k \sum_{j=i+1}^d n \left( \tilde{\lambda}_j^{-1} - \tilde{\lambda}_i^{-1} \right) (l_i - l_j) dz_{ij}^2. \tag{32}$$

To prove (32) we note that the  $ij$ -th element of  $\mathbf{T} \mathbf{L}$  is  $(\mathbf{T} \mathbf{L})_{ij} = (\tilde{\lambda}_j^{-1} - \tilde{\lambda}_i^{-1}) l_j dz_{ij}$ . This term is zero if  $i, j > k$  or  $i = j$ . Hence, the elements of  $\mathbf{T} \mathbf{L} d \mathbf{Z}$  can be written as  $(\mathbf{T} \mathbf{L} d \mathbf{Z})_{ij} = \sum_{u=1}^d (\mathbf{T} \mathbf{L})_{iu} dz_{uj}$  if  $i \leq k$  and  $(\mathbf{T} \mathbf{L} d \mathbf{Z})_{ij} = \sum_{u=1}^k (\mathbf{T} \mathbf{L})_{iu} dz_{uj}$  if  $i, j > k$ . Taking into account the skew-symmetry of  $d \mathbf{Z}$  the trace of  $\mathbf{T} \mathbf{L} d \mathbf{Z}$  is

$$\begin{aligned} \text{tr} [\mathbf{T} \mathbf{L} d \mathbf{Z}] &= - \sum_{i=1}^k \sum_{j=1}^d \left( \tilde{\lambda}_j^{-1} - \tilde{\lambda}_i^{-1} \right) l_j dz_{ij}^2 - \sum_{i=k+1}^d \sum_{j=1}^k \left( \tilde{\lambda}_j^{-1} - \tilde{\lambda}_i^{-1} \right) l_j dz_{ij}^2 \\ &= \sum_{i=1}^k \left( \sum_{j=k+1}^d \left( \tilde{\lambda}_j^{-1} - \tilde{\lambda}_i^{-1} \right) l_i dz_{ij}^2 - \sum_{j=1}^d \left( \tilde{\lambda}_j^{-1} - \tilde{\lambda}_i^{-1} \right) l_j dz_{ij}^2 \right). \end{aligned}$$

The latter equation can be rewritten as

$$\begin{aligned} \text{tr}[\mathbf{T}\mathbf{L}\mathbf{d}\mathbf{Z}] &= \sum_{i=1}^k \sum_{j=i+1}^d \left( \tilde{\lambda}_j^{-1} - \tilde{\lambda}_i^{-1} \right) (l_i - l_j) \mathbf{d}z_{ij}^2 \\ &\quad - \sum_{i=1}^k \left( \sum_{j=1}^{i-1} \left( \tilde{\lambda}_j^{-1} - \tilde{\lambda}_i^{-1} \right) l_j \mathbf{d}z_{ij}^2 + \sum_{j=i+1}^k \left( \tilde{\lambda}_j^{-1} - \tilde{\lambda}_i^{-1} \right) l_i \mathbf{d}z_{ij}^2 \right). \end{aligned} \quad (33)$$

It remains to show that the second term of (33) is zero and we prove it using induction. Let  $B(k)$  denote the value of this second term in (33). It is clear that  $B(1)$  and  $B(2)$  are both equal to zero. Let us assume that the claim holds for  $B(s)$  with a certain  $s \geq 2$ . For  $s + 1$  simple calculation yields

$$B(s+1) = B(s) + \sum_{j=1}^s \left( \tilde{\lambda}_j^{-1} - \tilde{\lambda}_{s+1}^{-1} \right) l_j \mathbf{d}z_{s+1,j}^2 + \sum_{i=1}^s \left( \tilde{\lambda}_{s+1}^{-1} - \tilde{\lambda}_i^{-1} \right) l_i \mathbf{d}z_{s+1,i}^2,$$

where the first term on the right side of the equation is zero due to the induction hypothesis and the sum of the last two terms is clearly also equal to zero.

With this we have shown that the cross derivatives are zero and therefore the Hessian of  $\log p(\mathcal{D}, \mathbf{Z}, \hat{\Lambda}, \hat{\sigma}^2 | k)$  is diagonal at  $\mathbf{Z} = \mathbf{0}$ . Using (32), the determinant of the negative Hessian with respect to  $\mathbf{Z}$ , required in (21), is

$$|\mathbf{A}_U| = n^m \prod_{i=1}^k \prod_{j=i+1}^d \left( \tilde{\lambda}_j^{-1} - \tilde{\lambda}_i^{-1} \right) (l_i - l_j). \quad (34)$$

Since all the cross derivatives between  $\hat{\Lambda}$ ,  $\hat{\sigma}^2$  and  $\mathbf{Z}$  are also zero, the negative Hessian of  $\log p(\mathcal{D}, \mathbf{Z}, \hat{\Lambda}, \hat{\sigma}^2 | k)$  is block diagonal. Its determinant equals  $|\mathbf{A}_{U\Lambda\sigma^2}| = |\mathbf{A}_U| |\mathbf{A}_\Lambda| |\mathbf{A}_{\sigma^2}|$ , where the factors are given in (24), (25), and (34), respectively.

Now that we have all the terms and estimates which are needed in (21), Laplace approximation of the integral directly leads to

$$p(\mathcal{D} | k) \approx 2^k c_k |\hat{\Lambda}|^{-\frac{N}{2}+1} \hat{\sigma}^{-N(d-k)+2} \exp \left\{ -\frac{Nd}{2} + k + 1 \right\} (2\pi)^{\frac{m+k+1}{2}} (|\mathbf{A}_U| |\mathbf{A}_\Lambda| |\mathbf{A}_{\sigma^2}|)^{-\frac{1}{2}},$$

which is our recommended criterion. In this formula a few errors of Minka (2001) have been corrected. In the original paper the term  $\exp\{k + 1\}$  and the Jacobians of the transformations were missing.

## 4 Simulation Study

In this section we want to test the performance of the different criteria for selecting the optimum number of principal components using synthetic data. In a simulation study we compare four criteria that make use of the probabilistic PCA model: the original MIBS criterion (15), its corrected version (18), the BIC (16) and the orthogonal variational Bayes approximation, called OVPCA (see Smidl and Quinn, 2005). Additionally we consider the AIC (13) and MDL (14) criteria.

The prior parameter  $\alpha$  for the corrected MIBS criterion is chosen to be  $\alpha = 0.01$  to make the prior diffuse. The subsequently presented results for the simulation study would not change much as long as the value for  $\alpha$  is kept small, say smaller than 1. Large values for  $\alpha$ , however, would adversely affect the results, especially for small sample sizes when the prior information dominates the information obtained from the likelihood.

In the first experiment the data is sampled from a 10-dimensional zero-mean Gaussian distribution with variances [10 8 6 4 2] in the first five directions. The remaining five directions have a variance equal to 1. The results for sample sizes varying between 5 and 350 are shown in Figure 1(a). For every sample size we simulate 1000 data sets and count how often the true dimensionality is recovered by a certain criterion. One can see that AIC is the most accurate criterion for sample sizes smaller than 75, however, fails to increase the recovery rate when the number of samples gets large. The AIC is an inconsistent estimator for the true number of dimensionality which is visible in nearly all subsequent experiments. The corrected MIBS criterion, abbreviated by *corMIBS* in Figures 1(a) - (h), has excellent recovery rates for both small and large number of sampled data. Notably, it gives far better results than the original MIBS criterion for all different sample sizes. MDL turns out to be less accurate than the corrected and the original MIBS criterion but outperforms the BIC. The OVPCA shows the worst performance when there are less than 100 data samples but catches up to BIC and even MDL when we further increase the sample size.

The second experiment, displayed in Figure 1(b), is the same as the first except that the noise dimension is changed from 5 to 10. Again, the corrected MIBS criterion performs best and gives larger recovery rates than the original MIBS criterion. Compared to the first experiment, the MDL criterion and BIC need much more samples to reach a high recovery rate. This effect is far smaller for both MIBS criteria and also for the OVPCA criterion, which leads to superior results compared to MDL and BIC in this experiment.

The third experiment differs from the second experiment only by the fact that here we have a noise variance of 0.5 instead of 1. Compared to the second experiment, all criteria need less data to reach high recovery rates. The corrected MIBS criterion is the top performer for all sample sizes, however, the difference to the original MIBS criterion is smaller than in the first two experiments. Again, OVPCA is weak when the number of samples is small. Detailed results are shown in Figure 1(c).

In the fourth experiment the noise variance is reduced even more and set equal to 0.1. As can be seen in Figure 1(d), all criteria yield recovery rates higher than 95% when the sample size exceeds 25. The corrected MIBS criterion shows a far better performance than the original MIBS criterion for the very small sample sizes of 8 and 11. Note that this is the only experiment in which the BIC gives slightly better results than the MDL criterion. OVPCA is as good as the MIBS criteria when the number of samples is 17.

The fifth experiment is the same as the third except that here we have 20 noise dimensions instead of 10. The facts already mentioned in the second experiment also apply here. Figure 1(e) shows that the MDL criterion and BIC have a weak performance compared to the third experiment due to the added noise dimensions. The MIBS criteria and the OVPCA criterion are much less affected. Again, the corrected MIBS criterion gives the best results, especially for small sample sizes.

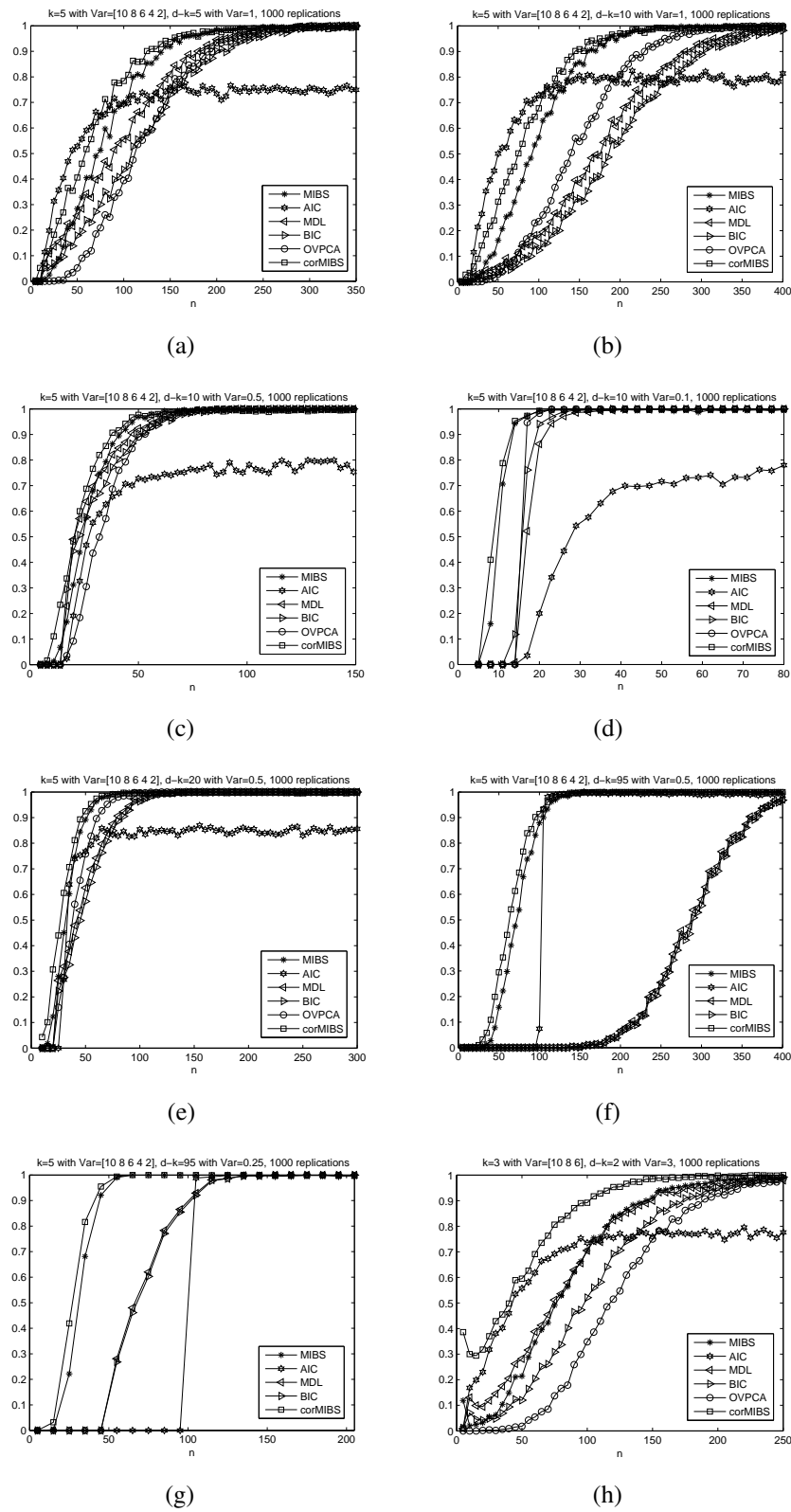


Figure 1: Recovery rates for the different criteria in the simulation study. The results of experiments 1-8 are displayed in (a)-(h).

The sixth and the seventh experiment, displayed in Figure 1(f)-(g), examine the case where the number of noise dimensions is 95 and, therefore, very large compared to the true dimensionality of the data which is equal to 5. When the data dimension is large the OVPCA is a very time-consuming method. This is the reason why we omit it for these two experiments. The corrected MIBS criterion is the most accurate method, followed by the original MIBS criterion. The AIC shows good results for sample sizes larger than 100. As we would expect, both MDL and BIC show a poor performance because of the large data dimensionality.

The last experiment is designed to test the case in which the true dimensionality is larger than the number of noise dimensions. The data is generated from a five-dimensional Gaussian distribution with variances [10 8 6] in the first three directions. The remaining two directions have a variance equal to 3. Again, the corrected MIBS criterion leads to the best results. However, in this experiment the difference in terms of recovery rate between the corrected MIBS criterion and all other criteria is more striking. The original MIBS criterion and the MDL criterion perform equally well. The BIC and especially the OVPCA give poor results.

## 5 Analysis of Hyper-Spectral Skin Cancer Data

If spectral measurements using hundreds of narrow contiguous wavelength intervals are performed, the resulting image is called a hyper-spectral image and is often represented as a hyper-spectral image cube. In contrast to RGB-images, where every pixel can be represented as a three-dimensional vector with entries corresponding to the red, green and blue channels, a hyper-spectral image contains pixels represented as multidimensional vectors with elements indicating the reflectivity at a specific wavelength. Thus, these vectors correspond to spectra in the physical meaning and are equal to spectra measured with e.g. spectrometers.

A set of 310 hyper-spectral images ( $171 \times 170$  pixels and 270 spectral bands after preprocessing steps) of malign and benign lesions were taken in clinical studies at the Medical University Graz, Austria. They are classified as melanoma or mole by human experts on the basis of a histological examination. Kazianka, Leitner, and Pilz (2008) used this data set and reported on the classification results for unobserved skin cancer images. In their paper they mention that as the dimensionality of the data equals the number of spectral bands, using the full spectral information in classification and clustering approaches leads to computational complexity. Moreover, the spectral bands are highly correlated and contain noise. To overcome the curse of dimensionality PCA is used to reduce the dimensions of the data, and inherently also the unwanted noise. Preceding analysis and inspection of scores and loadings showed that the optimum number of principal components to be retained is 7 (see Kazianka, 2007).

To test the performance of the AIC, BIC, MDL, MIBS and corrected MIBS criterion we selected 500 pixels from the training set. The prior parameter  $\alpha$  for the corrected MIBS criterion is again set to  $\alpha = 0.01$ . The dimensionality picked by each method is shown in Table 1. The original MIBS and the corrected MIBS criterion choose the number of components that was suggested by the preceding analysis.

Table 1: Number of components picked by the different criteria for the skin cancer data.

Estimator	MIBS	corMIBS	MDL	AIC	BIC
No. Features	7	7	6	8	6

The values of the criteria for varying number of retained components are shown in Figure (2). As can be seen in Figure 2(a), the corrected MIBS criterion gives less probability mass for larger dimensionalities than the original formulation.

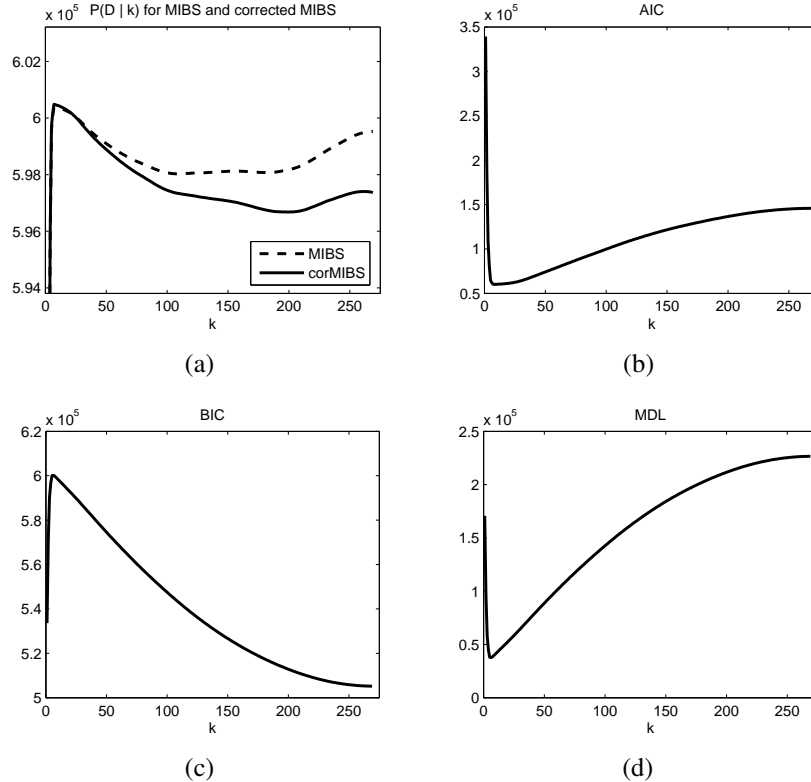


Figure 2: Values of the criteria for varying dimensionality in the skin cancer study.

## 6 Conclusion

The corrected version of the MIBS criterion for automatically selecting the optimum data dimensionality in the probabilistic PCA model shows excellent results in a simulation study. It is the constant top performer for all eight different experiments and works especially well with small sample sizes. The proposed criterion not only outperforms the original MIBS criterion but also other well-known methods such as BIC, AIC, MDL and OVPCA. The study also reveals that the performance of the BIC and the MDL criterion strongly depends on the data dimension. Both criteria give poor results when the number of dimensions is large. Besides the promising results for synthetic data the application to the hyper-spectral skin cancer images shows that the corrected MIBS criterion is also suitable for real world data which do not necessarily follow a Gaussian distribution.

## Acknowledgements

The authors are grateful to Vaclav Smidl from the UTIA Prague for providing the code for the OVPCA method. We would also like to express our gratitude to CTR Carinthian Tech Research for providing the hyper-spectral skin cancer data. Furthermore, our special thanks go to the referee whose comments and suggestions helped a lot in preparing this final version of the paper.

## References

- Basilevsky, A. (1994). *Statistical Factor Analysis and Related Methods*. New York: Wiley.
- Bishop, C. (2008). *Pattern Recognition and Machine Learning*. Berlin: Springer.
- Cichocki, A., and Amari, S. (2002). *Adaptive Blind Signal and Image Processing*. Chichester: Wiley.
- James, A. (1954). Normal Multivariate Analysis and the Orthogonal Group. *Annals of Mathematical Statistics*, 25, 40-75.
- Kazianka, H. (2007). *Classification Techniques for Hyper-Spectral Medical Image Data*. Unpublished master's thesis, University of Klagenfurt.
- Kazianka, H., Leitner, R., and Pilz, J. (2008). Segmentation and Classification of Hyper-Spectral Skin Data. In C. Preisach, H. Burkhardt, L. Schmidt-Thieme, and R. Decker (Eds.), *Data Analysis, Machine Learning and Applications* (p. 245-252). Berlin: Springer.
- Khatri, C., and Mardia, K. (1977). The von Mises-Fisher Distribution in Orientation Statistics. *Journal of the Royal Statistical Society, Series B*, 39, 95-106.
- Leonowicz, Z., Karvanen, J., Tanaka, T., and Rezmer, J. (2004). Model Order Selection Criteria: Comparative Study and Applications. In *Proceedings of the VIth International Workshop CPEE 2004* (p. 193-196). Warsaw: University of Technology.
- Lindley, D. (1980). Approximate Bayesian Statistics Methods. In J. Bernardo, M. de Groot, D. Lindley, and A. Smith (Eds.), *Bayesian Statistics* (p. 223-237). Valencia: University Press.
- Minka, T. (2001). Automatic Choice of Dimensionality for PCA. *Advances in Neural Information Processing Systems*, 13, 598-604.
- Pearson, K. (1901). On Lines and Planes of Closest Fit to Systems of Points in Space. *Philosophical Magazine*, 2, 559-572.
- Smidl, V., and Quinn, A. (2005). *The Variational Bayes Method in Signal Processing*. Berlin: Springer.
- Tipping, M., and Bishop, C. (1999). Probabilistic Principal Component Analysis. *Journal of the Royal Statistical Society, Series B*, 61, 611-622.

Authors' address:

Hannes Kazianka and Jürgen Pilz  
Institute of Statistics, University Klagenfurt  
Universitätsstraße 65-67  
9020 Klagenfurt  
E-Mails: Hannes.Kazianka@uni-klu.ac.at