

# New Routes from Minimal Approximation Error to Principal Components

Abhilash Alexander Miranda\*, Yann-Aël Le Borgne, Gianluca Bontempi

*Machine Learning Group, Département d'Informatique  
Université Libre de Bruxelles, Boulevard du Triomphe - CP212  
Brussels, Belgium 1050*

September 12, 2007

**Abstract.** We introduce two new methods of deriving the classical PCA in the framework of minimizing the mean square error upon performing a lower-dimensional approximation of the data. These methods are based on two forms of the mean square error function. One of the novelties of the presented methods is that the commonly employed process of subtraction of the mean of the data becomes part of the solution of the optimization problem and not a pre-analysis heuristic. We also derive the optimal basis and the minimum error of approximation in this framework and demonstrate the elegance of our solution in comparison with an existing solution in the framework.

**Keywords:** principal components analysis, eigenvalue, matrix trace

## 1. Introduction

The problem of approximating a given set of data using a weighted linear combination of a fewer number of vectors than the original dimensionality is classic. Many applications that require such a dimensionality reduction desire that the new representation retain the maximum variability in the data for further analysis. A popular method that attains simultaneous dimensionality reduction, minimum mean square error of approximation and retainment of maximum variance of the original data representation in the new representation is called the Principal Components Analysis (PCA) (Hotelling, 1933; Jolliffe, 2002).

The most popular framework for deriving PCA starts with the analysis of variance. A very common derivation of PCA in this framework generates the basis by iteratively finding the orthogonal directions of maximum retained variances (Hotelling, 1933; Jolliffe, 2002; Mardia et al., 1979; Johnson and Wichern, 1992). Since variance is implied in the statement of the problem here, the mean is subtracted from the data as a preliminary step. The second most predominant framework derives PCA by minimizing the mean square error of approximation (Duda et

---

\* abalexan@ulb.ac.be, Tel: +32 2 650 55 04, Fax: +32 2 650 56 09

al., 2001; Diamantaras and Kung, 1996; Bishop, 2006). Aided by the derivation in the variance-based framework above, it has become acceptable to resort to mean subtraction of the data prior to any analysis in this framework too in order to keep the analysis simple. In this letter our focus is on the latter framework within which we demonstrate two distinct and elegant analytical methods of deriving the PCA. In each of these methods of derivation, subtraction of data mean becomes part of the solution instead of being an initial assumption.

The letter is organized as follows: in Section 2 we describe the motivation behind the need for yet another derivation of the classical PCA. In particular, we highlight the issue of mean centering in Section 2.1. The notations are introduced in Section 2.2 and the PCA problem and its interpretations are discussed in Section 3. After reviewing an existing solution in Section 4, we make it evident in Section 5 that our two methods are due to two forms of the optimization function. Then we introduce our two methods of solving the PCA problem in Sections 6 and 7 and arrive at a simple common form of the optimization function in both the methods. This is analyzed further in Section 8 where we show the relation of the variance to the optimal basis in PCA as well as the minimum approximation error attained in PCA. In Section 8.3, we revisit the existing solution in our framework of PCA introduced in Section 4 and equate it with our approach.

## 2. Motivation

There are many standard textbooks of multivariate and statistical analysis (Jolliffe, 2002; Mardia et al., 1979; Johnson and Wichern, 1992) detailing PCA as a technique that seeks the best approximation of a given set of data points using a linear combination of a set of vectors which retain maximum variance along their directions. Since this framework of PCA starts by finding the covariances, the mean has to be subtracted from the data and becomes the *de facto* origin of the new coordinate system. The subsequent analysis is simple: find the eigenvector corresponding to the largest eigenvalue of the covariance matrix as the first basis vector. Then find the second basis vector on which the data components bear zero correlation with the data components on the first basis vector. This turns out to be the eigenvector corresponding to the second largest eigenvalue. In successively finding the basis vectors that have uncorrelated components as the eigenvectors of decreasing retained variances, the second order cross moments

between the components are successively eliminated<sup>1</sup>. Computationally, a widely employed trick in this framework finds the eigenvectors using singular value decomposition of the mean centered data matrix which effectively diagonalizes the covariance matrix without actually computing it (Jolliffe, 2002; Mardia et al., 1979). The set of orthogonal vectors corresponding to the largest few singular values proportional to the variances yield those directions which retain the maximum variance in the new representation of the data.

The second framework derives the PCA approximation by using its property of minimizing the mean square error. We think that this framework is more effective in introducing PCA to a novice because the two outcomes of optimal dimensionality reduction, viz. error minimization and retained variance maximization, are attained here simultaneously. Following the path of the retained variance maximization framework and to keep the analysis simple, many textbooks (Johnson and Wichern, 1992; Diamantaras and Kung, 1996; Hyvarinen et al., 2001; Ripley, 1996) advocate a mean subtraction for this framework too without sensible justification. Pearson stated in his now classical paper (Pearson, 1901):

*“The second moment of a system about a series of parallel lines is always least for the line going through the centroid. Hence: The best-fitting straight line for a system of points in a space of any order goes through the centroid of the system.”*

A procedure equivalent to rephrasing of this statement is followed in a much referenced textbook (Duda et al., 2001) which reasons that since the mean is the zero-dimensional hyperplane which satisfies the minimum average square error criterion, any higher dimensional hyperplane should be excused to pass through it too. In order to keep our analysis coherent with the concept of simultaneous dimensionality reduction, retained variance maximization and approximation error minimization, we do not invite the reader to such geometric intuitions. Note that the error minimization framework can also be viewed as a total least squares regression problem with all variables thought to be free so that the task is to fit a lower dimensional hyperplane that minimizes the perpendicular distances from the data points to the hyperplane (Van Huffel, 1997).

We will also be reviewing (Bishop, 2006) who derives PCA in the same framework as that of ours. Unlike in their approach, we neither undertake a complete decomposition nor force any basis vectors to bear a common statistic enticed by the prospect of an eventual mean

---

<sup>1</sup> Elimination of higher order cross moments is dealt in Independent Components Analysis (ICA) (Hyvarinen et al., 2001).

subtraction. Also for the benefit of practitioners who would like to deal data as realizations of a random variable, our treatment in the data samples domain can be readily extended to a population domain.

## 2.1. TO MEAN CENTER OR NOT

In the framework of finding the basis of a lower dimensional space which minimizes the mean square error of approximation, the process of mean subtraction has so far been part of the heuristics that the data need to be centered before installing the new low-dimensional coordinate system motivated by the philosophy according to (Pearson, 1901) that, had the mean of the data not been subtracted, the best fitting hyperplane would pass through the origin and not through the centroid. But there exist situations where a hyperplane is merely expected to partition the data space into orthogonal subspaces and as a result subtraction of mean is not desired. Note that in such situations, the term ‘principal component’ does not strictly hold as the basis vectors for the new space are not obtained from the data covariance matrix and the main concern there is the decomposition of the data rather than approximation.

One such set of situations are addressed by the Fukunaga-Koontz Transform (Fukunaga and Koontz, 1970; Miranda and Whelan, 2005) and it works by not requiring a subtraction of mean but instead finds the principal components of the autocorrelation matrices of two classes of data. It is widely used in automatic target rejection where eigenvalue decomposition generates basis for a target space orthogonal to the clutter space. But such is the issue of mean subtraction in using this transform that researchers of (Mahalanobis et al., 2004) and (Xuo et al., 2003) use autocorrelation and covariance matrices, respectively, for the same task without a justification of the impact of their choice to mean center or not. A similar approach called Eigenspace Separation Transformation (Plett et al., 1997) aimed at classification also does not involve mean subtraction. A family of techniques called Orthogonal Subspace Projection that is widely applied in noise rejection of signals use data that are not mean centered for the generalized PCA that follows (Harsanyi and Chang, 1994).

Although the theory of PCA demands mean subtraction for optimal low dimensional approximation, for many applications it is not without consequence. For example, the researchers of ecology and climate studies have extensively debated the purpose and result of mean centering for their PCA-based data analysis. In (Noy-Meir, 1973), the characteristics and apparent advantages of the principal components generated without mean subtraction are compared for data sampled homogeneously in the original space or otherwise. The claim made therein is that if data

form distinct clusters, the influence of variance within a cluster on another can be minimized by not subtracting the mean. Another ongoing debate named ‘Hockey Stick’ controversy (McIntyre and McKittrick, 2005) involves the appropriateness of mean subtraction for PCA in a much cited global warming study (Mann et al., 1998).

It should be borne in mind that this letter is neither solely about the aforementioned issue of mean centering that researchers using PCA often take it for granted nor does it change the results of PCA that is previously known to them. But we demonstrate in a new comprehensive framework that (i) the mean subtraction becomes a solution to the optimization problem in PCA and we reach this solution through two simple distinct methods that borrow little from traditional textbook derivations of PCA, and (ii) the derivation of the basis for the low dimensional space converges to minimum approximation error and maximum retained variance in the framework. Consequently, we believe that many problems which raise questions about their choice regarding mean subtraction can be revisited with ease using our proposed PCA framework.

## 2.2. NOTATIONS

While ensuring clarity to our analysis in this letter, we have tried to maintain its brevity through appropriate notations as summarized in the table below.

$J_q$	:	error function
$q$	:	new dimensionality
$p$	:	original dimensionality
$n$	:	number of samples
$x_k$	$\in$	$\mathbb{R}^p$ ; $k^{\text{th}}$ data sample
$\hat{x}_k$	$\in$	$\mathbb{R}^p$ ; approximation of $x_k$
$\theta$	$\in$	$\mathbb{R}^p$ ; new general origin
$\tilde{x}_k$	$=$	$x_k - \theta \in \mathbb{R}^p$
$e_i$	$\in$	$\mathbb{R}^p$ ; $i^{\text{th}}$ orthonormal basis vector of $\mathbb{R}^p$
$W$	$=$	$[e_1 \cdots e_q] \in \mathbb{R}^{p \times q}$
$B$	$=$	$I - WW^T \in \mathbb{R}^{p \times p}$
$\widetilde{W}$	$=$	$[e_{q+1} \cdots e_p] \in \mathbb{R}^{p \times p-q}$
$z_k$	$\in$	$\mathbb{R}^q$ ; dependent on $x_k$
$b$	$\in$	$\mathbb{R}^{p-q}$ ; a constant
$\text{Tr}(A)$	:	Trace of the matrix $A$
$\text{rank}(A)$	:	Rank of the matrix $A$
$\mu$	$\in$	$\mathbb{R}^p$ ; sample mean
$S$	$\in$	$\mathbb{R}^{p \times p}$ ; sample covariance matrix
$\lambda_i$	:	$i^{\text{th}}$ largest eigenvalue of $S$
$r$	$=$	$\text{rank}(S)$

### 3. Problem Definition in the Sample Domain

Let  $x_k \in \mathbb{R}^p, k = 1, \dots, n$  be a given set of data points. Suppose we are interested in orthonormal vectors  $e_i \in \mathbb{R}^p, i = 1, \dots, q \leq p$  whose resultant of weighted linear combination  $\hat{x}_k \in \mathbb{R}^p$  can approximate  $x_k$  with a minimum average (sample mean) square error or in other words minimize

$$J_q(\hat{x}_k) = \frac{1}{n} \sum_{k=1}^n \|x_k - \hat{x}_k\|^2. \quad (1)$$

The problem stated above means that we need an approximation  $x_k \simeq \hat{x}_k$  such that

$$\hat{x}_k = \sum_{i=1}^q (e_i^T x_k) e_i \quad (2)$$

so that we attain the minimum for  $J_q$ . This approximation assumes that the origin of all orthonormal  $e_i$  is the same as that of the coordinate

system in which the data is defined.

We reformulate the approximation

$$\hat{x}_k = \theta + \sum_{i=1}^q \left( e_i^T (x_k - \theta) \right) e_i \quad (3)$$

to assume that the new representation using basis vectors  $e_i$  has a general origin  $\theta \in \mathbb{R}^p$  and not the origin as in the approximation (2). We assume orthonormality here because (i) orthogonality guarantees linearly independent  $e_i$  so that they form a basis for  $\mathbb{R}^q$  (ii) normalizing  $e_i$  maintains notational simplicity in not having to divide the scalars  $e_i^T x_k$  in (2) by the norm  $\|e_i\|$  which is unity due to our assumption. Hence, the PCA problem may be defined as

$$\underset{e_i, \theta}{\operatorname{argmin}} \frac{1}{n} \sum_{k=1}^n \|x_k - \hat{x}_k\|^2 : \begin{aligned} & \hat{x}_k = \theta + \sum_{i=1}^q \left( e_i^T (x_k - \theta) \right) e_i; \\ & e_i^T e_j = 0, i \neq j; \quad e_i^T e_i = 1 \quad \forall i, j. \end{aligned} \quad (4)$$

which seeks a set of orthonormal basis vectors  $e_i$  with a new origin  $\theta$  which minimizes the error function in (1) in order to find a low-dimensional approximation  $W^T(x_k - \theta) \in \mathbb{R}^q$  for any  $x_k \in \mathbb{R}^p$ , where

$$W = [e_1 \cdots e_q]. \quad (5)$$

It is now easy to see that (3) becomes

$$\hat{x}_k = \theta + WW^T(x_k - \theta). \quad (6)$$

Hence the displacement vector directed from the approximation  $\hat{x}_k$  towards  $x_k$  is  $x_k - \hat{x}_k = (x_k - \theta) - WW^T(x_k - \theta)$ , which using  $\tilde{x}_k = x_k - \theta$  can be written concisely as  $x_k - \hat{x}_k = \tilde{x}_k - WW^T\tilde{x}_k$ . By setting  $B = I - WW^T$  for simplicity of notation, we write the displacement vector as

$$x_k - \hat{x}_k = B\tilde{x}_k. \quad (7)$$

#### 4. Review of an existing solution

The PCA solution in the framework of approximation error minimization is derived in (Bishop, 2006) is reviewed here. They derive PCA by undertaking a complete decomposition

$$\hat{x}_k = Wz_k + \widetilde{W}b \quad (8)$$

into basis vectors contained in the columns of matrix  $W$  of (5) and  $\widetilde{W} = [e_{q+1} \cdots e_p] \in \mathbb{R}^{p \times p-q}$  such that components of  $z_k \in \mathbb{R}^q$  depend

on  $x_k$ , whereas components of  $b \in \mathbb{R}^{p-q}$  are constants common for all data points.

By taking derivative of the error function with respect to  $b$ , they find that

$$b = \widetilde{W}^T \mu \quad (9)$$

so that the common components are those of the sample mean vector  $\mu$ . This implies that by subtracting the sample mean they are no longer obliged to retain the  $p - q$  dimensions corresponding to the columns of  $\widetilde{W}$  which preserve little information regarding the variation in the data. The first drawback of this approach is that it couples the process of dimensionality reduction with mean subtraction although the two are shown to be independent in our derivation. By taking derivative of the error function with respect to  $z_k$ , they also show that  $z_k = W^T x_k$ . Hence the approximation they are seeking is

$$\hat{x}_k = WW^T x_k + \widetilde{W}\widetilde{W}^T \mu. \quad (10)$$

The second drawback of their approach is the requirement of yet another constrained minimization of the error function before they reach the solution for the optimal columns of  $W$ .

## 5. Methods of PCA

We have discussed the need for a new derivation of PCA by (i) explaining the lack of proper justification in the literature for subtracting the mean in a minimum mean square error framework, (ii) justifying its chronic necessity in many applications in Section 2, and (iii) reviewing a recent attempt to solve this problem in Section 4. Our derivations of the solution for the problem in (4) are due to two simple forms of the error function  $J_q$  of (1) which we state as follow:

$$\text{Form 1 : } J_q(\hat{x}_k) = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{x}_k)^T (x_k - \hat{x}_k) \quad (11)$$

$$\text{Form 2 : } J_q(\hat{x}_k) = \text{Tr} \left( \frac{1}{n} \sum_{k=1}^n (x_k - \hat{x}_k) (x_k - \hat{x}_k)^T \right) \quad (12)$$

We analyze Form 1 in (11) in Section 6 to arrive at a simplified  $J_q$  which is exactly the same as we get by following a different method of analyzing Form 2 in (12) in Section 7. These two methods take different paths towards the common error function, viz., the first using straightforward expansion of the terms in  $J_q$  and the second using the property of matrix trace. The common form of  $J_q$  is subsequently



treated in Section 8 to reveal the rest of the solution to our original problem.

## 6. Analysis of Form 1 of error function

Using (7), the error function  $J_q$  of Form 1 in (11) can be developed as

$$J_q(B, \theta) = \frac{1}{n} \sum_{k=1}^n \tilde{x}_k^T B^T B \tilde{x}_k. \quad (13)$$

The property that  $B = I - WW^T$  is idempotent and symmetric, i.e.,

$$B = B^2 = B^T, \quad (14)$$

or  $B$  is simply an orthogonal projector, may be used to reduce  $J_q$  further as

$$J_q(B, \theta) = \frac{1}{n} \sum_{k=1}^n \tilde{x}_k^T B \tilde{x}_k. \quad (15)$$

Expanding  $J_q$  above using  $\tilde{x}_k = x_k - \theta$  gives

$$J_q(B, \theta) = \frac{1}{n} \sum_{k=1}^n \left[ x_k^T B x_k - 2\theta^T B x_k + \theta^T B \theta \right] \quad (16)$$

In order to get the  $\theta$  which minimizes  $J_q$ , we find the partial derivative  $\partial J_q / \partial \theta = -2B \left[ \frac{1}{n} \sum_{k=1}^n x_k - \theta \right]$  and setting it to zero results in

$$\theta = \frac{1}{n} \sum_{k=1}^n x_k = \mu \quad (17)$$

which is as simple as regarding the sample mean of the data points as the new origin. Henceforth, we can assume that  $\tilde{x}_k$  is the data point  $x_k$  from which the sample mean has been subtracted.

### 6.1. SIMPLIFYING THE ERROR FUNCTION

We may analyze the error function in (15) as follow:

$$\begin{aligned} J_q(W) &= \frac{1}{n} \sum_{k=1}^n \tilde{x}_k^T (I - WW^T) \tilde{x}_k \\ &= \frac{1}{n} \sum_{k=1}^n \tilde{x}_k^T \tilde{x}_k - \frac{1}{n} \sum_{k=1}^n \tilde{x}_k^T WW^T \tilde{x}_k \\ &= \frac{1}{n} \sum_{k=1}^n \tilde{x}_k^T \tilde{x}_k - \text{Tr} \left( W^T \left[ \frac{1}{n} \sum_{k=1}^n \tilde{x}_k \tilde{x}_k^T \right] W \right). \end{aligned}$$

We have the sample covariance matrix

$$S = \frac{1}{n} \sum_{k=1}^n \tilde{x}_k \tilde{x}_k^T \Big|_{\theta=\mu} \quad (18)$$

so that the term  $\frac{1}{n} \sum_{k=1}^n \tilde{x}_k^T \tilde{x}_k \Big|_{\theta=\mu}$  equals  $\text{Tr}(S)$ , and we can write

$$J_q(W) = \text{Tr}(S) - \text{Tr}(W^T S W). \quad (19)$$

## 7. Analysis of Form 2 of error function

We now analyze the Form 2 of the error function  $J_q$  by substituting (7) in (12) as

$$J_q(B, \theta) = \text{Tr} \left( B \left[ \frac{1}{n} \sum_{k=1}^n \tilde{x}_k \tilde{x}_k^T \right] B^T \right). \quad (20)$$

### 7.1. FINDING $\theta$

As in the previous section, we denote the sample mean and sample covariance matrix by  $\mu$  and  $S$ , respectively, and we may develop the term in (20):

$$\begin{aligned} \frac{1}{n} \sum_{k=1}^n \tilde{x}_k \tilde{x}_k^T &= \frac{1}{n} \sum_{k=1}^n (x_k - \theta)(x_k - \theta)^T \\ &= \frac{1}{n} \sum_{k=1}^n [x_k x_k^T - x_k \theta^T - \theta x_k^T + \theta \theta^T] \\ &= S + \mu \mu^T - \mu \theta^T - \theta \mu^T + \theta \theta^T, \end{aligned} \quad (21)$$

where we have used the sample autocorrelation matrix (Fukunaga, 1990)  $\frac{1}{n} \sum_{k=1}^n x_k x_k^T = S + \mu \mu^T$ . We get  $J_q(B) = \text{Tr} \left( B \left( S + \mu \mu^T - \mu \theta^T - \theta \mu^T + \theta \theta^T \right) B^T \right)$  upon substituting (21) in (20). Using (14) and the cyclic permutation property of trace of matrix products<sup>2</sup> we get

$$J_q(B) = \text{Tr} \left( B \left( S + \mu \mu^T - \mu \theta^T - \theta \mu^T + \theta \theta^T \right) \right) \quad (22)$$

and using the property of derivative of trace<sup>3</sup> and the chain rule of derivatives<sup>4</sup> we find that  $\partial J_q / \partial \theta = 2B(-\mu + \theta)$  which when equated

<sup>2</sup>  $\text{Tr}(\Upsilon \Phi \Psi) = \text{Tr}(\Psi \Upsilon \Phi) = \text{Tr}(\Phi \Psi \Upsilon)$

<sup>3</sup>  $\partial(\text{Tr}(\Psi \Phi^T)) / \partial \Phi = \Psi$

<sup>4</sup>  $\partial(\cdot) / \partial u = [\partial(\cdot) / \partial (uv^T)] v$

to zero results in

$$\theta = \mu \quad (23)$$

leading to the same solution of Form 1 in (17).

## 7.2. SIMPLIFYING THE ERROR FUNCTION

Having found  $\theta$ , we can substitute it in (22) to get  $J_q(B) = \text{Tr}(BS)$ . On substitution for  $B$  in terms of  $W$ , we may write  $J_q(W) = \text{Tr}(S) - \text{Tr}(WW^T S)$ . Utilizing the cyclic permutation property of matrix trace again, we get

$$J_q(W) = \text{Tr}(S) - \text{Tr}(W^T SW). \quad (24)$$

## 8. Optimal basis and minimum error

Note that we have arrived at the same set of equations in both (19) and (24) of Form 1 and Form 2, respectively, whereby substituting  $W$  as defined in (5) in either of them gives

$$J_q(e_i) = \text{Tr}(S) - \sum_{i=1}^q e_i^T S e_i. \quad (25)$$

### 8.1. RELATION OF VARIANCE TO OPTIMAL BASIS

Let us now find the variance  $\lambda_i$  of the data projected on the basis vector  $e_i$ . It is the average of the square of the difference between projections  $e_i^T x_k$  of the data points and the projection  $e_i^T \mu$  of the sample mean, i.e.,

$$\begin{aligned} \lambda_i &= \frac{1}{n} \sum_{k=1}^n (e_i^T x_k - e_i^T \mu)^2 \\ &= \frac{1}{n} \sum_{k=1}^n (e_i^T x_k - e_i^T \mu) (e_i^T x_k - e_i^T \mu)^T \\ &= e_i^T \left[ \frac{1}{n} \sum_{k=1}^n (x_k - \mu) (x_k - \mu)^T \right] e_i \\ &= e_i^T S e_i. \end{aligned} \quad (26)$$

Thus, the term  $\sum_{i=1}^q e_i^T S e_i$  in (25) gives the portion of the total variance  $\text{Tr}(S)$  retained along the directions of orthonormal  $e_i$ . Hence, we are looking for vectors  $e_i$  of the form  $\lambda_i = e_i^T (S e_i)$ , which is satisfied if

$Se_i = \lambda_i e_i$ . Such a relation implies  $(e_i, \lambda_i)$  form an eigen-pair of  $S$ . Note that since there is no unique basis for any nontrivial vector space, any basis that spans the  $q$ -dimensional space generated by the eigenvectors of  $S$  are solutions to  $e_i$  too. In (25), since

$$\operatorname{argmin}_{e_i} J_q = \operatorname{argmax}_{e_i} \sum_{i=1}^q e_i^T S e_i, \quad (27)$$

the vectors  $e_i$  have to be the eigenvectors corresponding to the  $q$  largest ('principal') eigenvalues of  $S$ . This is the classical result of the PCA.

### 8.2. RELATION OF VARIANCE TO MINIMUM APPROXIMATION ERROR

It follows from (26) that the term  $\sum_{i=1}^q e_i^T S e_i = \sum_{i=1}^q \lambda_i$  of (25) is the sum of the  $q$  principal eigenvalues of  $S$ ; this is the maximum variance that could be retained upon approximation using any  $q$  basis vectors. Also,  $\operatorname{Tr}(S) = \sum_{i=1}^r \lambda_i$ ,  $r = \operatorname{rank}(S)$  is the total variance in the data. Substituting these in  $J_q$  in (25) gives the difference of the total variance and the maximum retained variance; the result is the minimum of the eliminated variance. Hence, for  $\lambda_i \geq \lambda_j, j > i$ , the minimum mean square approximation error can be expressed as

$$J_q = \underbrace{\sum_{i=1}^r \lambda_i}_{\text{total variance}} - \underbrace{\sum_{i=1}^q \lambda_i}_{\text{retained variance}} = \underbrace{\sum_{i=q+1}^r \lambda_i}_{\text{eliminated variance}}. \quad (28)$$

### 8.3. COMPARISON OF THE REVIEWED SOLUTION WITH THE PRESENT WORK

In order to compare the solution of (Bishop, 2006) reviewed in Section 4, let us first write the approximation in (6) as  $\hat{x}_k = WW^T x_k + B\theta$ . We know from (17) and (23) that  $\theta = \mu$  and, hence,

$$\hat{x}_k = WW^T x_k + B\mu. \quad (29)$$

If  $\widetilde{W}\widetilde{W}^T = B$ , we have the approximation according to (Bishop, 2006) in (10) of Section 4 equivalent to the approximation in (29).

While the drawbacks of (6) highlighted in Section 4 exist, let us outline the difference in these two approaches: we have demonstrated in our solutions that the new origin  $\theta \in \mathbb{R}^p$  of the low dimensional coordinate system should be the mean  $\mu \in \mathbb{R}^p$  so that the error of the approximation is reduced. But (Bishop, 2006) necessitates an orthogonal projection of certain data-independent components  $b \in \mathbb{R}^{p-q}$

to  $\mu \in \mathbb{R}^p$  to achieve the same objective. Our approach has shown that such a dimensionality reduction coupled with mean subtraction is unnecessary for deriving PCA.

#### 8.4. POPULATION PCA

For population PCA (Mardia et al., 1979; Johnson and Wichern, 1992), where the samples that form the data are assumed to be realizations of a random variable, we have made it easy for the reader to follow our analysis by just replacing all occurrences of  $\frac{1}{n} \sum_{k=1}^n \rightarrow \mathcal{E}$ , the expectation operator; and bold faces for random variables as in  $x_k \rightarrow \mathbf{x}$ ,  $\hat{x}_k \rightarrow \hat{\mathbf{x}}$ , and  $\tilde{x}_k \rightarrow \tilde{\mathbf{x}}$ .

### 9. Conclusion

Motivated by the need to justify the heuristics of pre-analysis mean centering in PCA and related questions, we have demonstrated through two distinct methods that the mean subtraction becomes part of the solution of the standard PCA problem in an approximation error minimization framework. We based these two methods on two subtly different forms of the error function. We have also derived the optimal basis and the minimum error of approximation in this framework and have compared our results with an existing solution.

### Acknowledgements

This work was funded by the Project TANIA (WALEO II) of the Walloon Region, Belgium. The authors thank their colleague Olivier Caelen for his appreciable comments. Thanks are also due to Dr. P.P. Mohanlal of ISRO Inertial Systems Unit, India for his valuable insights. The authors are very grateful to the editor-in-chief and three anonymous reviewers for their excellent suggestions on an earlier version of this letter.

### References

- C. M. Bishop. Pattern Recognition and Machine Learning. Information Science and Statistics. Springer, New York, August 2006.
- K. I. Diamantaras and S. Y. Kung. Principal Component Neural Networks: Theory and Applications. New York, February, 1996.
- R. O. Duda, P. E. Hart, and D. G. Stork. Pattern Classification. Wiley Interscience, New York, second edition, 2001.

- K. Fukunaga. Introduction to Statistical Pattern Recognition. Computer Science and Scientific Computing. Academic Press, San Diego, second edition, 1990.
- K. Fukunaga and W.L.G. Koontz. Application of the Karhunen-Loeve expansion to feature selection and ordering. *IEEE Transactions on Computers*, C-19(4):311-318, 1970.
- J. C. Harsanyi and C.-I. Chang. Hyperspectral Image Classification and Dimensionality Reduction: An Orthogonal Subspace Projection Approach. *IEEE Transactions on Geoscience and Remote Sensing*, 32(4):779-785, July, 1994.
- H. Hotelling. Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology*, 24:417-441, 1933.
- X. Huo, M. Elad, A. G. Flesia, B. Muiise, R. Stanfill, A. Mahalanobis, et al. Optimal Reduced-Rank Quadratic Classifiers using the Fukunaga-Koontz Transform with Applications to Automated Target Recognition, *Proceedings of SPIE*, 5094:59-72, September, 2003.
- A. Hyvarinen, J. Karhunen, and E. Oja. Independent Component Analysis, volume 27 of *Adaptive and Learning Systems for Signal Processing, Communications and Control*. Wiley-Interscience, New York, June, 2001.
- R. A. Johnson and D. W. Wichern. *Applied Multivariate Statistical Analysis*. Prentice-Hall, Inc., Upper Saddle River, New Jersey, third edition, 1992.
- I. T. Jolliffe. *Principal Component Analysis*. Springer, 2nd edition, New York, 2002.
- A. Mahalanobis, R. R. Muiise, S. R. Stanfill and A. Van Nevel. Design and Application of Quadratic Correlation Filters for Target Detection. *IEEE Transaction on Aerospace and Electronic Systems*, 40(3):837-850, July, 2004.
- M. E. Mann, R. S. Bradley and M. K. Hughes. Global-scale Temperature Patterns and Climate Forcing over the Past Six Centuries. *Nature*, 392:779-788, April, 1998.
- K. Mardia, J. Kent, and J. Bibby. *Multivariate Analysis*. Academic Press, London, 1979.
- S. McIntyre and R. McKittrick. Reply to comment by Huybers on "Hockey sticks, Principal Components, and Spurious Significance". *Geophysical Research Letters*, 32, L20713, 2005.
- A. A. Miranda and P. F. Whelan. Fukunaga-Koontz Transform for Small Sample Size Problems. *Proceedings of the IEE Irish Signals and Systems Conference*, pp. 156-161, Dublin, 2005.
- I. Noy-Meir. Data Transformations in Ecological Ordination: I. Some Advantages of Non-Centering. *The Journal of Ecology*, 61(2):329-341, July, 1973.
- K. Pearson. On lines and planes of closest fit to systems of points in space. *Philosophical Magazine*, 2:559-572, 1901.
- G. L. Plett, T. Doi, and D. Torrieri. Mine Detection using Scattering Parameters and an Artificial Neural Network. *IEEE Transactions on Neural Networks*, 8(6):1456-1467, November, 1997.
- B. D. Ripley. *Pattern Recognition and Neural Networks*. Cambridge University Press, Cambridge, 1996.
- S. Van Huffel (Ed.), *Recent Advances in Total Least Squares Techniques and Errors-in-Variables Modeling*, SIAM, Philadelphia, PA, 1997.