

Regression on non-orthogonal analytical base functions, orthogonal analytical base functions, and principal components.

R DIXON

1. Introduction

Much of the utility of principal components rests upon the theorem which states that a set of ordered eigenvectors will be more efficient on average at extracting variance from a dependent data set than any other set of linear functions of the data. Met O 11 Technical Note No 93 presents experimental evidence which, even though it relates to independent data sets, suggests that it may be desirable to study the conditions of validity and practical range of application for this theorem. Principal components have attracted wide interest since they were introduced into Meteorology in the fifties, but at a fundamental theoretical level there is a grey area containing questions which have remained largely undiscussed. This Note raises one such question, the assumptions as to the nature of the residuals when principal components are used in regression.

2. Estimation using non-orthogonal analytical base functions

Let there be a data vector $\underline{h} = \underline{h} (h_1, h_2, \dots, h_m)$ and a set of analytical base functions f_1, f_2, \dots, f_n . These analytical base functions can then each be evaluated over m discrete points of a physical domain to form the $m \times 1$ base vectors $\underline{f}_1, \underline{f}_2, \underline{f}_3, \dots, \underline{f}_n$. It may then be supposed that the "true" relationship between the data and the base vectors is given by

$$\underline{h}_{m \times 1} = \underline{F}_{m \times n} \cdot \underline{a}_n + \underline{r}_{m \times 1} \quad (1)$$

where \underline{F} is the matrix of column vectors

$$\underline{F} = (\underline{f}_1 | \underline{f}_2 | \dots | \underline{f}_n) \quad (2)$$

\underline{a} is the coefficient vector

\underline{r} is the vector of residuals which arises because we are considering the case where $m > n$ and so the functions cannot fit the data exactly.

The belief that (1) is the true structural model for the data may be based on physical reasoning, on some previous experimental evidence, or it may simply be postulated *faute de mieux*. Whatever the reason the problem arises of estimating \underline{a} , for \underline{a} itself is generally unknown. A well-known and much used estimate is the ordinary least squares estimate (OLS) $\hat{\underline{a}}$ given by

$$\hat{\underline{a}}_n = (\tilde{\underline{F}} \cdot \underline{F})^{-1} \cdot \tilde{\underline{F}} \cdot \underline{h}_{m \times 1} \quad (3)$$

If the residuals in (1) are such that, in the population

$$E(\underline{r}) = \underline{0} \quad (4)$$

and

$$E(\underline{r} \underline{r}^T) = \sigma^2 \underline{I}_{m \times m} \quad (5)$$

where \underline{I} is the unit diagonal matrix, E is the expectation operator, and σ is a scalar, then

$$E(\hat{\underline{a}}) = \underline{a} \quad (6)$$

and

$$E(\hat{\underline{a}} - \underline{a})(\hat{\underline{a}} - \underline{a})^T \text{ is a minimum} \quad (7)$$

(6) states that $\hat{\underline{a}}$ is unbiased, whilst (7) states that the sampling variance - co-variance of the difference between the estimate $\hat{\underline{a}}$ and the true unknown \underline{a} is a minimum amongst all such linear unbiased estimators. Clearly (6) and (7) are very desirable properties for an estimate of \underline{a} for they minimize the risk that one has a bad estimate.

If the residuals in (1) are such that (4) holds but in place of (5) we have

$$E(\underline{\hat{r}}\underline{\hat{r}}^T) = \sigma^2 \underline{\Omega} \quad (8)$$

where $\underline{\Omega}$ is a more general matrix than I then the OLS estimate $\hat{\underline{a}}$ given by (3) does not have the property (7). The solution vector which does have the desirable minimum sampling variance property (7) I will denote by \underline{a}^* and it is given by

$$\underline{a}^* = (\tilde{F} \cdot \underline{\Omega}^{-1} \cdot F)^{-1} \cdot \tilde{F} \cdot \underline{\Omega}^{-1} \cdot \underline{h} \quad (9)$$

Under the condition (8) it is the solution vector \underline{a}^* which has the property

$$E(\underline{a}^* - \underline{a})(\underline{a}^* - \underline{a})^T \text{ is a minimum} \quad (10)$$

\underline{a}^* is the fact the least-squares solution to a transformed version of the "true" relationship (1). The transformation is

$$\underline{\Omega}^{-\frac{1}{2}} \cdot \underline{h} = \underline{\Omega}^{-\frac{1}{2}} \cdot F \cdot \underline{a} + \underline{\Omega}^{-\frac{1}{2}} \cdot \underline{r} \quad (11)$$

Thus if it is known, assumed, inferred or believed that (8) holds then the least-squares fitting procedure must be applied to solve (11) not (1) if the minimum variance solution is to be obtained. This transformation of (1) by taking the scalar product of the square-root of $\underline{\Omega}^{-1}$ through (1) is known as prewhitening, and the whole process of finding \underline{a}^* is referred to as generalized least-squares (GLS) and \underline{a}^* is known as the GLS estimate of \underline{a} .

It should be noticed that getting the best solution vector is dependent upon making a correct assumption as to the nature of the residuals. Furthermore, it is the residuals in (1), the "true" model equation, about which assumptions have to be made, not the residuals which are obtained as a result of the fitting, i.e. as the result of using $\hat{\underline{a}}$ or \underline{a}^* . With laboratory data it may be that a scientist can be reasonably sure of his assumptions because the data have been produced by a well understood and controlled process, but with meteorological data often produced by processes which are only imperfectly understood and over which we have no control the difficulty is a very real one. The residuals from the fitted process have to be used somehow to infer the validity or otherwise of the assumptions on the true residuals. That is to say for example that if it is assumed that (1), (4), and (5) hold then the residuals $\hat{\underline{r}}$ in

$$\underline{h} = F \cdot \hat{\underline{a}} + \hat{\underline{r}} \quad (12)$$

have to be used to test the validity of the assumptions (4) and (5) about the unknown residuals \underline{r} in (1). It is something of a chicken and egg problem and it has a voluminous literature.

The OLS solution to (1) may be specified by the following two equations

$$\underline{h}_{mi} = \underset{mn}{F} \cdot \underset{ni}{\hat{a}} + \underset{mi}{\hat{r}} \quad (13)$$

$$\underline{0}_{ni} = \underset{nm}{\tilde{F}} \cdot \underset{mi}{\hat{r}} \quad (14)$$

This follows since, by virtue of (14), taking \tilde{F} through (13) leads in effect straight to the OLS solution (3). Equation (14) requires the residual vector \hat{r} to be orthogonal to the columns of F . The fact that (13) and (14) together specify the OLS solution to (1) is of central importance to the theme of this Note, as will emerge in Section 4.

3. Estimation using orthogonal analytical base functions

The model (1) is now replaced by

$$\underline{h}_{mi} = \underset{mn}{\Phi} \cdot \underset{ni}{b} + \underset{mi}{r} \quad (15)$$

where $\underset{mn}{\Phi}$ is the matrix of orthonormal column vectors

$$\underset{mn}{\Phi} = \left(\underset{mi}{\phi_1} \mid \underset{mi}{\phi_2} \mid \underset{mi}{\phi_3} \mid \dots \mid \underset{mi}{\phi_n} \right) \quad (16)$$

the $\underline{\phi}$ -vectors being orthonormal in the simple sense that

$$\underset{nm}{\tilde{\Phi}} \cdot \underset{mn}{\Phi} = \underset{nn}{I} \quad (17)$$

The OLS estimate for \underline{b} is then

$$\underset{ni}{\hat{b}} = \underset{nm}{\tilde{\Phi}} \cdot \underset{mi}{h} \quad (18)$$

which is computationally convenient in the sense that each coefficient in \hat{b} is obtained by a simple scalar product

$$b_i = \underset{im}{\tilde{\phi}_i} \cdot \underset{mi}{h} \quad i = 1, 2, \dots, n \quad (19)$$

The use of orthogonal base vectors has done nothing to change the considerations as to the nature of the residuals. The arguments of the previous Section still apply and \hat{b} will be the unbiased minimum variance estimate of b if and only if (4) and (5) hold. Indeed in most applications (1) and (15) will be the same model because the Φ -matrix will be obtained from the F matrix by some factorization process such as the Gram-Schmidt or Householder algorithms so that

$$F = \Phi \cdot C \quad (20)$$

and \underline{b} and \underline{a} , and \hat{b} and \hat{a} will be related by the simple linear transformation

$$C \cdot \underline{a} = \underline{b} \quad , \quad C \cdot \hat{a} = \hat{b} \quad (21)$$

However the use of orthogonal base vectors does give rise to a useful artifice. We can generate a full set of m $\underline{\phi}$ -vectors to form the $m \times m$ matrix

$$\Phi = (\underline{\phi}_1 | \underline{\phi}_2 | \underline{\phi}_3 | \dots | \underline{\phi}_m) \quad (22)$$

with the $\underline{\phi}$ -vectors orthonormal in the simple sense that

$$\tilde{\Phi} \cdot \Phi = I \quad (23)$$

The data vector \underline{h} can now be completely represented in terms of these $\underline{\phi}_s$ as

$$\underline{h} = b_1 \underline{\phi}_1 + b_2 \underline{\phi}_2 + \dots + b_m \underline{\phi}_m \quad (24)$$

There is no residual vector as the data is fitted exactly. We can now suppose that the terms in (24) have been ordered according to some principle and we can then decide to discard all but n of these terms and represent the data as

$$\underline{h} = b_1 \underline{\phi}_1 + b_2 \underline{\phi}_2 + \dots + b_n \underline{\phi}_n \quad (25)$$

This is quite legitimate mathematically as the orthogonality of the $\underline{\phi}_s$ makes their individual contributions independent of each other. Furthermore the ordering and selection principle which has decided the terms to be retained in (25) may have some apparently quite sound physical reasoning behind it and we may feel justified in arguing that the discarded $(m - n)$ terms can contribute little but noise. Nevertheless, this artifice needs to be examined more closely.

Although the representation of \underline{h} in terms of n vectors is often written as in (25) above, yet (25) is not strictly correct. Since $n < m$ the terms can no longer fit the data exactly and a vector of residuals must be introduced. Thus strictly speaking (25) has to be written as

$$\underline{h} = b_1 \underline{\phi}_1 + b_2 \underline{\phi}_2 + \dots + b_n \underline{\phi}_n + \underline{r} \quad (26)$$

where the use of the different diacritic Ψ over the \underline{r} merely indicates an absence of any commitment, at this stage, as to the nature of the \underline{r} in (26). But the nature of the \underline{r} in (26) can in fact be elicited, for if the representation (24) fits the data exactly whilst the representation (26) fits it leaving a residual \underline{r} it follows that this residual vector \underline{r} is itself represented by the discarded $\underline{\phi}$ terms and that

$$\underline{\overset{\Psi}{r}} = b_{n+1} \underline{\phi}_{n+1} + b_{n+2} \underline{\phi}_{n+2} + \dots + b_m \underline{\phi}_m \quad (27)$$

Now, since all the $\underline{\phi}$ -vectors in (24) are orthogonal to each other, all the $\underline{\phi}$ -vectors in (27) are orthogonal to all the $\underline{\phi}$ -vectors in (26). This means that \underline{r} itself in (26) is orthogonal to all the $\underline{\phi}$ -vectors in (26). Thus if we express this fact, and (26) itself, in matrix-vector form we have

$$\underline{h} = \underline{\Phi} \cdot \underline{b} + \underline{\overset{\Psi}{r}} \quad (28)$$

$$\underline{0} = \underline{\tilde{\Phi}} \cdot \underline{\overset{\Psi}{r}} \quad (29)$$

But these two equations have a familiar look and referring back it is seen that they are the analogues, for the orthogonal case, of equations (13) and (14). Equations (28) and (29) are in fact the two conditions which specify the OLS solution to (15). This follows because taking $\underline{\tilde{\Phi}} \cdot$ through (28) leads by virtue of (29) and (17) straight to (18). Thus we may drop the non-committal diacritic Ψ and rewrite (28) and (29) with the appropriate diacritic as

$$\underline{h} = \underline{\Phi} \cdot \underline{\hat{b}} + \underline{\hat{r}} \quad (30)$$

$$\underline{0} = \underline{\tilde{\Phi}} \cdot \underline{\hat{r}} \quad (31)$$

This Section brings out the fact that the artifice of selecting from a complete set of orthogonal base vectors, or ordering and truncating the set, and forming the coefficients of a data set representation by taking scalar products (19) is fully equivalent to the OLS process and implies the same assumptions on the residuals. This last point is crucial, for it is not obvious and may easily be overlooked.

4. Estimation using Principal Components

In this case instead of using base vectors derived by evaluating a set of analytical functions over the discrete points of the domain we derive the base vectors directly from a set of data vectors. Taking N data vectors $\underline{h}_1, \underline{h}_2, \underline{h}_3, \dots, \underline{h}_N$ and viewing these N vectors as $m \times 1$ column vectors we assemble them to form the dependent data matrix

$$H = (\underline{h}_1 | \underline{h}_2 | \underline{h}_3 | \dots | \underline{h}_N) \quad (32)$$

Now by taking the scalar product $\cdot \tilde{H}$ we can form the symmetric matrix C, where

$$NC = H \cdot \tilde{H} \quad (33)$$

$\begin{matrix} m & m \\ m & N \end{matrix}$

The effect of forming the scalar product has been to eliminate the N-space. Also if the components of the individual \underline{h} vectors have been measured from the individual means of each vector then C is the covariance matrix of the data. C can now be factored into

$$C = U \cdot \Lambda \cdot \tilde{U} \quad (34)$$

$\begin{matrix} m & m \\ m & m \end{matrix}$

where

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_m) \quad (35)$$

$\begin{matrix} m & m \end{matrix}$

is a diagonal matrix of eigenvalues of $\frac{1}{N} H \cdot \tilde{H}$ and

$$U = \left(\begin{array}{c|c|c|c|c} \underline{u}_1 & \underline{u}_2 & \underline{u}_3 & \dots & \underline{u}_m \end{array} \right) \quad (36)$$

$\begin{matrix} m & m & m & m & m \end{matrix}$

is a modal matrix, the columns of which are the eigenvectors corresponding to $\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_m$. Without loss of generality we can assume that $\lambda_1 > \lambda_2 > \lambda_3 > \dots > \lambda_m$. The modal matrix U has the property

$$\tilde{U} \cdot U = I \quad (37)$$

$\begin{matrix} m & m \\ m & m \end{matrix}$

Since the \underline{u}_i are m-vectors they can be used as base-vectors for the representation of any other m-vector. In particular they may be used as base vectors for the representation of the individual vectors of the dependent data set, so that we have

$$\underline{h}_i = b_{i1} \underline{u}_1 + b_{i2} \underline{u}_2 + b_{i3} \underline{u}_3 + \dots + b_{im} \underline{u}_m \quad i = 1, 2, \dots, N \quad (38)$$

$\begin{matrix} m & m & m & m & m \end{matrix}$

where the b_{ij} are sets of coefficients (not, of course, the same as the b_i in Section 3).

The m equations (38) may be rewritten in matrix-vector notation as

$$H = U \cdot B \quad (39)$$

$\begin{matrix} m & N \\ m & m \end{matrix}$

From (39) and (37) it follows that

$$\underset{mm}{\tilde{U}} \cdot \underset{mN}{H} = \underset{mN}{B} \quad (40)$$

and that therefore

$$\underset{Nm}{\tilde{B}} = \underset{Nm}{\tilde{H}} \cdot \underset{mm}{U} \quad (41)$$

and consequently

$$\underset{mN}{B} \cdot \underset{Nm}{\tilde{B}} = \underset{mm}{\tilde{U}} \cdot \underset{mN}{H} \cdot \underset{Nm}{\tilde{H}} \cdot \underset{mm}{U} \quad (42)$$

But from (33) this is

$$\underset{mN}{B} \cdot \underset{Nm}{\tilde{B}} = \underset{mm}{\tilde{U}} \cdot \underset{mm}{N} \cdot \underset{mm}{C} \cdot \underset{mm}{U} \quad (43)$$

and now, using (34) and (37) again, this comes to

$$\underset{mN}{B} \cdot \underset{Nm}{\tilde{B}} = \underset{mm}{N} \Lambda \quad (44)$$

Λ is the diagonal matrix of eigenvalues and so it transpires that the row vectors of B are orthogonal.

These are the basic equations and relationships of principal component analysis and using them one can prove the main well-known theorem to the effect that such a set of ordered eigenvectors will be more efficient on average at extracting variance from the dependent data set than any other set of linear functions of the data. Note the words "on average" in this statement of the theorem. It does not follow that the theorem is true for any particular individual data vector taken from the set.

In the light of this theorem it is natural to attempt to obtain an economical representation of the data vectors of the set by truncating (38) on some criterion. Thus dropping the subscript and taking \underline{h} as a typical data vector we will have

$$\underline{h}_{m_i} = b_1 \underline{u}_1 + b_2 \underline{u}_2 + \dots + b_n \underline{u}_n + \frac{\psi}{m_i} \quad (45)$$

or in matrix-vector form

$$\underline{h}_{m_i} = \underset{mm}{U} \cdot \underset{ni}{b} + \frac{\psi}{m_i} \quad (46)$$

the coefficients being found by taking the scalar products

$$b_j = \tilde{u}_j \cdot \underline{h} \quad (47)$$

Notice that since the n terms cannot fit the data exactly it has been necessary in (45) to introduce the unspecified residual vector \underline{r} . In the general literature of principal components, including the meteorological literature, this residual vector goes virtually unacknowledged in the theoretical discussion. Here the main theoretical point to be made is that once placed plainly in evidence as in (45) or (46) it does not long remain unspecified. It is clear that \underline{r} , being composed of all the discarded \underline{u} terms, is orthogonal to each of the column vectors of U and we are in fact dealing with precisely the same situation as in the case of equation (26) in Section 3 on orthogonal analytical base functions. The argument following (26) applies equally in the case of (45) and leads in exactly the same way to

$$\underline{h} = U \cdot \underline{\hat{b}} + \underline{\hat{r}} \quad (48)$$

$$\underline{0} = \tilde{U} \cdot \underline{\hat{r}} \quad (49)$$

In other words the truncation of the complete set of eigenvectors and the formation of the coefficients as scalar products (47) leads to the OLS solution. We have in fact found the OLS solution to a model

$$\underline{h} = U \cdot \underline{b} + \underline{r} \quad (50)$$

which has not been explicitly postulated. The general argument and terminology of the principal component approach somewhat obscures the fact that a model has been postulated, but such is the case and the OLS solution $\underline{\hat{b}}$ found by the above procedure will be the desirable "best" (i.e. minimum variance) solution only if (5) applies, just as in the cases of the analytical base functions dealt with in Sections 2 and 3. In other words a principal component enthusiast has as much of an obligation to state and test his assumptions concerning the residuals as does an enthusiast for polynomials, Fourier terms, exponentials, etc.

The principal components concept appears to have come over to meteorology from psychology and biology and in these sciences the basic theory, equations (32) to (44), ~~has been presented~~ has been presented with a subtle difference of emphasis. In the presentation as given in Section 4 the condition (44), revealing the orthogonality of the rows of the coefficient matrix B , emerges as an almost incidental consequence of the use of eigenvectors. In the presentation favoured by psychologists, biologists, and some meteorologists the condition (44) is made a central requirement. Indeed, it is sometimes referred to as the "principal component property". It is then shown that the eigenvectors of the sample covariance matrix are the only base vectors which permit this. Mathematically there is nothing to choose between the two approaches. The argument and terminology accompanying either of them tends to obscure the point made in this Note.

R Dixon

R Dixon
Met 0 11