# Derived variables calculated from similar joint responses: some characteristics and examples

Nanny Wermuth

*Universität Mainz, Germany*

D.R. Cox

*Nuffield College, Oxford, United Kingdom*

*Abstract:* A technique (Cox and Wermuth, 1992) is reviewed for finding linear combinations of a set of response variables having special relations of linear conditional independence with a set of explanatory variables. A theorem in linear algebra is used both to examine conditions in which the derived variables take a specially simple form and lead to reduced computations. Examples are discussed of medical and psychological investigations in which the method has aided interpretation.

## 1. Introduction

In situations in which the dependencies of several response variables $Y_1, Y_2, \ldots, Y_p$ on explanatory variables $X_1, X_2, \ldots, X_q$ are under study, it may be difficult to describe concisely the type of the dependence structure even if there are only linear relations. However if the joint responses are similar, in the sense that they are measured on comparable scales and they are thought of as capturing different aspects of an underlying phenomenon, then useful interpretations may be possible with new responses derived as linear combinations of the original response variables in such a way that each new response $Y_i^*$ has linear conditional independence of all explanatory variables except one, i.e. $Y_i^* \perp\!\!\!\perp (X_1, \ldots, X_{i-1}, X_{i+1}, \ldots, X_q) \mid X_i$. Typically, all variables are quantitative, but the explanatory variables may include binary indicator variables.

*Correspondence to:* N. Wermuth, Psychologisches Institut, Universität Mainz. D-55099 Mainz, Deutschland.

We review first the relation of such derived variables to canonical variables, we then give conditions under which derived response variables are likely to be ineffective in the sense of having weak relations to the explanatory variables and conditions under which they take on particularly simple forms. Finally we describe analyses of several sets of variables.

## 2. Some characteristics of derived variables

Let $Y$ and $X$ be $p \times 1$ and $q \times 1$ vector variables measured as deviations from their means, $p \geq q$ then their covariance matrix $\text{Cov}(Y, X) = \Sigma$ and their concentration matrix $\Sigma^{-1}$, assumed to be positive definite, can be written as

$$\Sigma = E\begin{pmatrix} YY^T & YX^T \\ . & XX^T \end{pmatrix} = \begin{pmatrix} \Sigma_{yy} & \Sigma_{yx} \\ & \Sigma_{xx} \end{pmatrix}, \quad \Sigma^{-1} = \begin{pmatrix} \Sigma^{yy} & \Sigma^{yx} \\ & \Sigma^{xx} \end{pmatrix},$$

and the matrices $\pi_{y|x}$ and $\pi_{x|y}$ which contain regression coefficients obtained by linear regression of $Y$ on $X$ and $X$ on $Y$, respectively, are

$$\pi_{y|x} = \Sigma_{yx}\Sigma_{xx}^{-1} = -(\Sigma^{yy})^{-1}\Sigma^{yx}, \quad \pi_{x|y} = \Sigma_{xy}\Sigma_{yy}^{-1} = -(\Sigma^{xx})^{-1}\Sigma^{xy}. \tag{1}$$

New variables defined with a $p \times p$ matrix $F$ and a $q \times q$ matrix $G$ by linear transformations as $F^T Y$ and $G^T X$ then have covariance matrix

$$\text{Cov}(F^T Y, G^T X) = \begin{pmatrix} F^T\Sigma_{yy}F & F^T\Sigma_{yx}G \\ . & G^T\Sigma_{xx}G \end{pmatrix}. \tag{2}$$

If $p = q$ a fairly direct argument determines the new vector $Y^*$ via the requirement that the matrix of regression coefficients of $Y^*$ on $X$ is the identity matrix, so that in particular the regression of $Y_i^*$ on $X$ involves only $X_i$. This, via (1) and (2), can be shown to imply that

$$Y^* = \pi_{y|x}^{-1}Y = \Sigma_{xx}\Sigma_{yx}^{-1}Y. \tag{3}$$

If however, $p > q$, $Y^*$ cannot be uniquely determined in this way. A unique solution, in a reasonable sense optimal, is achieved by first reducing $Y$ to the $q \times 1$ vector of variables in the canonical regression of $Y$ on $X$ before applying (3).

In defining canonical variables (Hotelling, 1936) nonsingular matrices $\tilde{F}$ and $\tilde{G}$ are chosen for (2) so that two new sets of uncorrelated variables $\tilde{Y}$ and $\tilde{X}$ are obtained in which all variables are standardized to have unit variance and each variable $\tilde{X}_i$ has maximal correlation with a single variable $\tilde{Y}_i$. This implies in particular that the covariance matrix coincides with the correlation matrix and can be written as

$$\text{Cov}(\tilde{Y}, \tilde{X}) = \begin{pmatrix} I_{p-q} & 0 & 0 \\ . & I_q & \Lambda \\ . & . & I_q \end{pmatrix}, \tag{4}$$

where $I_s$ denotes a $s \times s$ unit matrix and $\Lambda$ is a diagonal matrix with the so-called canonical correlations along the diagonal.

Some calculation then shows that

$$Y^* = \pi_{\tilde{y}|x}^{-1}\tilde{Y} = \Sigma_{xx}\left(C^T\Sigma_{yx}\right)^{-1}C^TY, \tag{5}$$

gives the new derived responses while explanatory variables remain untransformed because they are assumed to have strong individual identity. The matrix $C$ used to define $Y^*$ can be be viewed essentially as that part of the matrix $\tilde{F} = (\tilde{F}_1 \; \tilde{C})$ for (4) which corresponds to the $q \times q$ matrix $\Lambda$ of canonical correlations, i.e. each column of $C$ agrees with one of $\tilde{C}$ up to a factor of proportionality, the standard deviation of a derived response $Y_i^*$. The covariance matrix of $C^TY$ and $X$ can be expressed with (2) and (4) as

$$\text{Cov}(C^TY, X) = C^T\Sigma_{yx} = \text{diag}\left(\sqrt{\sigma_{11}^*}, \ldots, \sqrt{\sigma_{qq}^*}\right)\Lambda\tilde{G}^{-1},$$

and it is invertible whenever all canonical correlations are nonzero, i.e. it is possible to obtain $q$ derived responses if and only if all canonical correlations are nonzero.

The expression via the matrix $C$, i.e. in terms of canonical variables, is valuable, in particular because a check that all $q$ canonical variables are indeed nonzero is very desirable. For some purposes it is, however, simpler to use the equivalent expression

$$Y^* = \Sigma_{xx}\left(\Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx}\right)^{-1}\Sigma_{xy}\Sigma_{yy}^{-1}Y. \tag{6}$$

There will be one or more zero canonical correlations whenever there is a linear combination of the responses which is uncorrelated with all of the explanatory variables. From knowledge of the covariance matrix of the variables under study it may be possible to predict when such a case is likely to occur. One instance of a zero canonical correlation is if all responses have zero marginal correlations, zero regression coefficients, or zero concentrations with one explanatory variable, i.e. if one of the columns of $\Sigma_{yx}$, of $\pi_{y|x}$, or of $\Sigma^{yx}$ contains only zeros. Another instance of a zero canonical correlation is if the dependence of all responses on two explanatory variables agrees in the sense that two columns of $\Sigma_{yx}$, of $\pi_{y|x}$ or of $\Sigma^{yx}$ are identical or are multiples of each other. On the other hand, $q$ effective derived responses can be defined whenever all $q$ canonical correlations are sizeable, i.e. have values larger than 0.1 say.

When $p = q$ and there is the same number of sizeable canonical correlations then the form of the derived responses (3) depends just on the structure of $\pi_{y|x}$, the regression coefficients in the regression of $Y$ on $X$. In particular, the derived response $Y^*$ will coincide with $Y$ if $\pi_{y|x}$ can be permuted to be in diagonal form; a single derived response $Y_i^*$ will coincide with a single response $Y_i$ if row $i$ of $\pi_{y|x}$ can be permuted to contain just one nonzero element in position $(i,i)$ and subsets of derived responses will depend only on corresponding subsets of the original response variables if $\pi_{y|x}$ can be permuted to be in block-diagonal form.

When $p > q$ and there are $q$ sizeable canonical correlations the form of the derived responses (5) depends instead on the structure of $\Sigma_{xx}(C^T\Sigma_{yx})^{-1}C^T$, the inverse of $\pi_{\tilde{y}|x}$, i.e. the regression coefficients in the regression of the (unstandardized) canonical variable $\tilde{Y}$ on $X$. In that case $q$ effective derived variables may still be obtained without actually computing the eigenvalues and eigenvectors of the $p \times p$ matrix $\pi_{x|y}^T\pi_{y|x}^T$ which give $\Lambda^2$ and $\tilde{F}$, respectively, since at most an eigenvalue analysis of a possibly much smaller $q \times q$ matrix $\pi_{y|x}^T\pi_{x|y}^T$ is needed (compare Rao, 1973) and sometimes even closed form solutions can be given. Such results have been shown to hold for all symmetrizable matrix products (Wermuth and Rüssmann, 1993) and are restated here in a slightly specialized form without proof.

Let $A = M_q L$, $B = LM_p$, where $L$ is a $q \times p$ matrix of rank $q$, $M_q$ a $q \times q$ positive definite matrix and $M_p$ a $p \times p$ positive definite matrix. Further, let $H(M)$ denote the matrix of normalized eigenvectors corresponding to the nonzero ordered eigenvalues of a matrix $M$, then

(i) the diagonal matrix $K$ of eigenvalues of $B^TA$ has only positive diagonal elements which coincide with the nonzero eigenvalues of $BA^T$ and of $AB^T$,

(ii) a matrix of eigenvectors of $B^TA$ is determined by the matrix of eigenvectors of $BA^T$ or of $AB^T$ as

$$H(B^TA)K^{\frac{1}{2}} = M_p L^T M_q H(BA^T) = M_p L^T H(AB^T),$$

(iii) matrices of eigenvectors of $BA^T$ and $AB^T$ can be found such that the product $H^T(BA^T)H(AB^T)$ is a diagonal matrix.

Applied to $\pi_{x|y}^T\pi_{y|x}^T$ with $A = \pi_{y|x}^T = \Sigma_{xx}^{-1}\Sigma_{xy}$ and $B = \pi_{x|y} = \Sigma_{xy}\Sigma_{yy}^{-1}$ (ii) gives for the matrix $C = H(\pi_{x|y}^T\pi_{y|x}^T)$ needed to compute the transformation matrix $C$ of the derived responses

$$\tilde{C}\Lambda = \Sigma_{yy}^{-1}\Sigma_{yx}\Sigma_{xx}^{-1}H(\pi_{x|y}\pi_{y|x}) = \Sigma_{yy}^{-1}\Sigma_{yx}H\left(\pi_{y|x}^T\pi_{x|y}^T\right), \tag{7}$$

i.e. column $j$ of $\tilde{C}$ is proportional to column $j$ of the matrices on the right-hand sides. If for instance $\pi_{x|y}\pi_{y|x}$ is a diagonal matrix of distinct elements then its diagonal elements are the squared canonical correlations and $H(\pi_{x|y}\pi_{y|x})$ is the identity matrix $I_q$.

The advantage of expression (7) may be computational if the number of responses is much larger than the number of explanatory variables, but mainly it is conceptual in that it permits study of the conditions under which the derived variables take on a form which possibly leads to simple interpretations. Some examples are as follows.

Two situations in which the rows of the transformation matrix $C^T$ in (5) for the derived responses are proportional to the rows of $\pi_{x|y}$ can be characterized as: (a) the explanatory variables are uncorrelated, i.e. their covariance matrix $\Sigma_{xx}$ is a diagonal matrix and all canonical correlations are distinct; (b) the product of the regression coefficient matrices is a diagonal matrix with identical diagonal elements, i.e. $\pi_{x|y}\pi_{y|x} = \Lambda^2 = cI_q$ where $c$ is some constant.

Examples are the following two covariance matrices with $p = 4$ response variables and $q = 2$ explanatory variables for which the squared canonical correlations are (0.4, 0.8) and (0.4, 0.4), respectively.

$$\Sigma^{(a)} = \begin{pmatrix} 1 & 0.2 & 0.2 & 0.2 & 0.4 & 0.4 \\ . & 1 & 0.2 & 0.2 & 0.4 & -0.4 \\ . & . & 1 & 0.2 & 0.4 & 0.4 \\ . & . & . & 1 & 0.4 & -0.4 \\ . & . & . & . & 1 & 0 \\ . & . & . & . & . & 1 \end{pmatrix},$$

$$\Sigma^{(b)} = \begin{pmatrix} 1 & 0.2 & 0.2 & 0.2 & 1.2 & 0.4 \\ . & 1 & 0.2 & 0.2 & 1.2 & 0.4 \\ . & . & 1 & 0.2 & 0.4 & 1.2 \\ . & . & . & 1 & 0.4 & 1.2 \\ . & . & . & . & 6 & 0.2 \\ . & . & . & . & . & 6 \end{pmatrix}$$

The relevant corresponding matrices of regression coefficients are

$$\pi_{x|y}^{(a)} = \begin{pmatrix} 0.25 & 0.25 & 0.25 & 0.25 \\ 0.50 & -0.50 & 0.50 & -0.50 \end{pmatrix}, \quad \pi_{x|y}^{(b)} = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{pmatrix},$$

so that the derived responses in these examples are just sums of sums or sums of differences of the variable pairs $(Y_1, Y_2)$ and $(Y_3, Y_4)$. It is conceivable that in some applications subject matter knowledge about the covariances of the variables under study is so strong that it can be judged directly whether conditions are likely to be satisfied which permit such a simple interpretation of the derived responses.

Different extensions for obtaining derived responses are possible if substantive knowledge suggests a chain of dependencies, for instance of the type where $U$-variables may depend on both $Y$- and $X$-variables and $Y$-variables only on the $X$-variables. We shall not explore this further here.

Table 1
Observed marginal correlations (lower half), observed partial correlations given all remaining variables (upper half), means and standard deviations for 44 patients.

| Variable | $Y_1$ | $Y_2$ | $X_1$ | $X_2$ |
|---|---|---|---|---|
| $Y_1$: Log diast. bp | 1 | 0.723 | 0.321 | −0.239 |
| $Y_2$: Log syst. bp | 0.732 | 1 | −0.101 | −0.094 |
| $X_1$: Log weight | 0.283 | 0.165 | 1 | 0.483 |
| $X_2$: Log height | −0.094 | −0.034 | 0.430 | 1 |
| Mean | 4.29 | 4.74 | 4.142 | 5.123 |
| Standard deviation | 0.132 | 0.112 | 0.164 | 0.034 |

Table 2
Observed marginal correlations (lower half), and observed partial correlations given all remaining variables (upper half), means and standard deviations for 40 patients awaiting an operation

| Variable | $Y_1$ | $Y_2$ | $Y_3$ | $X_1$ | $X_2$ | $X_3$ |
|---|---|---|---|---|---|---|
| $Y_1$: Log palmitic acid | 1 | 0.07 | 0.73 | 0.20 | 0.14 | −0.03 |
| $Y_2$: Log oileic acid | 0.90 | 1 | 0.58 | 0.05 | 0.53 | 0.10 |
| $Y_3$: Log linoleic acid | 0.95 | 0.92 | 1 | −0.26 | −0.40 | 0.04 |
| $X_1$: Blood sugar | −0.25 | −0.27 | −0.32 | 1 | −0.01 | −0.10 |
| $X_2$: Gender | 0.28 | 0.43 | 0.23 | −0.03 | 1 | 0.05 |
| $X_3$: Anxiety | 0.35 | 0.39 | 0.37 | −0.21 | 0.20 | 1 |
| Mean | 4.91 | 4.26 | 4.88 | 80.93 | 0.05 | 41.75 |
| Stand. dev. | 0.37 | 0.47 | 0.40 | 9.05 | 1.01 | 11.22 |

## 3. Some examples

In all situations where we considered computing derived variables we checked first that the reported marginal correlations were sensible summaries of the relations under study, in the sense that neither theory nor empirical evidence suggested some strong form of nonlinearity, and that there were no extreme outliers in the data.

Tables 1 and 2 give summaries for two sets of data, in which inspection of the marginal correlations $s_{ij}/\sqrt{s_{ii}s_{jj}}$ and of the partial correlations given all of the remaining variables $-s^{ij}/\sqrt{s^{ii}s^{jj}}$ suggest that one canonical correlation will be close to zero. Here $s_{ij}$ and $s^{ij}$ denote elements of $S$ and $S^{-1}$, the observed covariance and concentration matrix. In Table 1 diastolic and systolic blood pressure are regarded as potential responses to weight and height, all measured in the logarithmic scale for 44 patients (Slangen et al., 1991, the raw data are given in the Appendix), and in Table 2 three kinds of free fatty acids, measured in the logarihmic scale for 40 patients on the morning before an operation on the jaw (Krohne et al. 1989), are seen as potential responses to anxiety, blood sugar and gender (coded by −1 for males and 1 for females). The data in Table 1 suggest that $(Y_1, Y_2)$ are marginally independent of $X_2$ since the marginal

Table 3
Observed marginal correlations (lower half), observed partial correlations given all remaining variables (upper half), means and standard deviations for 44 patients.

| Variable | $Y_1$ | $Y_2$ | $X_1$ | $X_2$ |
|---|---|---|---|---|
| $Y_1$: Log diast. bp | 1 | 0.657 | 0.186 | 0.098 |
| $Y_2$: Log syst. bp | 0.732 | 1 | −0.241 | 0.300 |
| $X_1$: Body mass | 0.336 | 0.188 | 1 | 0.572 |
| $X_2$: Age | 0.510 | 0.492 | 0.608 | 1 |
| Mean | 4.29 | 4.74 | 37.94 | 29.52 |
| Standard dev. | 0.13 | 0.11 | 5.98 | 10.59 |

Table 4

Observed marginal correlations (lower half), observed partial correlations given all remaining variables (upper half), with derived responses for the data of Table 1

| Variable | $Y_1^*$ | $Y_2^*$ | $X_1$ | $X_2$ |
|---|---|---|---|---|
| $Y_1^* = Y_2 - Y_1$ | 1 | $-0.566$ | $-0.241$ | 0.300 |
| $Y_2^* = Y_1$ | $-0.544$ | 1 | $-0.107$ | 0.491 |
| $X_1$: Body mass | $-0.253$ | 0.336 | 1 | 0.572 |
| $X_2$: Age | $-0.131$ | 0.510 | 0.608 | 1 |
| Mean | 0.453 | 4.29 | 37.94 | 29.52 |
| Standard dev. | 0.091 | 0.132 | 5.98 | 10.59 |

correlations ($-0.094$, $-0.034$) are near zero and those in Table 2 suggest that $(Y_1, Y_2, Y_3)$ are conditionally independent of $X_3$ given $(X_1, X_2)$ since the third column in $S^{yx}$ and hence in $\hat{\pi}_{y|x}$ contains values near zero. Consequently, there is one of the corresponding squared canonical correlations (0.0006, 0.15) and (0.004, 0.19, 0.39) near to zero. For the further analyses shown in Tables 3 and 4 the two explanatory variables in Table 1, weight and height, are replaced by the single explanatory variable weight relative to height, i.e. body mass, which is known to be more relevant to blood pressure levels. The explanatory variable $X_3$ is deleted before computing the derived variables in Table 5.

In Table 3 correlations of diastolic and of systolic blood pressure are shown for the same collective of 44 patients but for two different potential explanatory variables: body mass (weight relative to height multiplied by 100) and age. The derived responses are $Y_1^* = Y_2 - Y_1$ and $Y_2^* = Y_1$ as approximations to the precise weights of $(1, -0.95)$ for $Y_1^*$ and $(1, .26)$ for $Y_2^*$. The approximation is judged to be well compatible with the desired independencies since a test statistic comparing $S^*$ with the matrix $\hat{\Sigma}^*$ estimated under the hypothesis $Y_1^* \perp\!\!\!\perp X_2 \mid X_1$ and $Y_2^* \perp\!\!\!\perp X_1 \mid X_2$: $-n \log\{\det(S^*)/\det(\hat{\Sigma}^*)\}$ has a value of 0.17 on 2 degrees of freedom. Under the assumption of normally distributed variables this is the likelihood-ratio statistic and can be tested as chi-squared, but more generally, being much smaller than its expected value, it is an indication

Table 5

Observed marginal correlations (lower half), observed partial correlations given all remaining variables (upper half), with derived responses for the data of Table 2

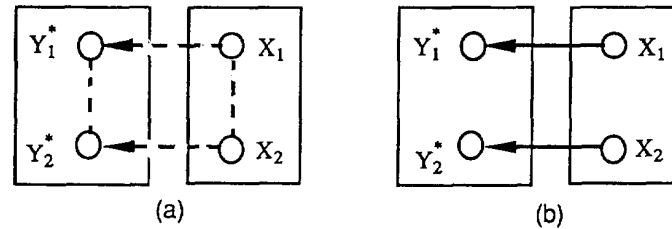| Variable | $Y_1^*$ | $Y_2^*$ | $X_1$ | $X_2$ |
|---|---|---|---|---|
| $Y_1^* = Y_1 - Y_3$ | 1 | 0.04 | 0.32 | 0.06 |
| $Y_2^* = Y_2 - Y_3$ | 0.09 | 1 | 0.02 | 0.57 |
| $X_1$: Blood sugar | 0.32 | 0.01 | 1 | $-0.06$ |
| $X_2$: Gender | 0.08 | 0.57 | $-0.03$ | 1 |
| Mean | 0.03 | $-0.62$ | 80.93 | 0.05 |
| Stand. dev. | 0.12 | 0.19 | 9.05 | 1.01 |

Fig. 1. Graphical representation of the independencies to which correlations with the derived responses are close: (a) for the data of Table 4 we have the typical structure for two derived responses with $Y_1^* \perp\!\!\!\perp X_2 \mid X_1$ and $Y_2^* \perp\!\!\!\perp X_1 \mid X_2$, (b) for the data of Table 5 satisfying nearly $X_1 \perp\!\!\!\perp X_2$ we get $(Y_1^*, X_1) \perp\!\!\!\perp (Y_2^*, X_2)$.

for a satisfactory fit. The matrices of observed standardized regression coefficients of $Y$ on $X$ compare with those of $Y^*$ on $X$ as:

$$\hat{\pi}_{y\mid x} = \begin{pmatrix} 0.040 & 0.486 \\ -0.177 & 0.600 \end{pmatrix}, \quad \hat{\pi}_{y^*\mid x} = \begin{pmatrix} -0.275 & 0.037 \\ 0.040 & 0.486 \end{pmatrix}.$$

These derived variables suggest a plausible interpretation: diastolic blood pressure increases with age after controlling for an increase in body mass, the ratio of systolic to diastolic blood pressure is higher the lower the body mass for persons of the same age. The graphical representation of the independencies (Cox and Wermuth 1993; Lauritzen and Wermuth, 1989) among derived responses and explanatory variables is shown in Figure 1(a) while Figure 1(b) gives the corresponding graph for the variables in Table 5.

In a next step of the analysis of the data in Table 2 the explanatory variable $X_3$ was discarded and derived responses were computed for $p = 3$ and $q = 2$. Then the very simple structure displayed in Table 5 and in Figure 1(b) results: the ratio of palmitic to linoleic acid relates only to the level of blood sugar and the ratio of oileic acid to linoleic acid relates only to gender. This approximates the precise weights of $(0.67\ 0.11\ -1)$ for $Y_1^*$ and $(0.31\ 0.82\ -1)$ for $Y_2^*$ and is judged to be well compatible with the independencies $(Y_1^*, X_1) \perp\!\!\!\perp (Y_2^*, X_2)$ since the above mentioned test statistic has a value of 0.59 on 4 degrees of freedom. The matrices of standardized regression coefficients of $Y^*$ on $X$ compare with those estimated under the independence assumption as:

$$\hat{\pi}_{y^*\mid x} = \begin{pmatrix} 0.322 & 0.097 \\ 0.032 & 0.575 \end{pmatrix}, \quad \hat{\pi}_{y^*\mid x} = \begin{pmatrix} 0.317 & 0 \\ 0 & 0.571 \end{pmatrix}.$$

The last set of data in Table 6 contains information for 72 students on brain activity (CNV: Contingent negative variation) measured at three locations (frontal, central and parietal) and under two experimental conditions, one in which the participant has to prepare for a motor activity at a stimulus four seconds after a signal has forwarned him ($g$: the 'go'-situation) and another in which he is not to react ($n$: the 'nogo'-situation). This defines the six responses $Y_1, \ldots, Y_6$ in Table 6, where all measurements are taken during a fixed first part of an

Table 6
Observed marginal correlations (lower half), observed partial correlations given all remaining variables (upper half) for 72 students

| Variable | $Y_1$ | $Y_2$ | $Y_3$ | $Y_4$ | $Y_5$ | $Y_6$ | $X_1$ | $X_2$ |
|---|---|---|---|---|---|---|---|---|
| $Y_1$: CNV-frontal, $g$ | 1 | 0.67 | −0.14 | 0.69 | −0.34 | −0.03 | −0.13 | −0.21 |
| $Y_2$: CNV-central, $g$ | 0.80 | 1 | 0.66 | −0.55 | 0.73 | −0.50 | −0.09 | 0.29 |
| $Y_3$: CNV-parietal, $g$ | 0.47 | 0.78 | 1 | 0.13 | −0.47 | 0.76 | −0.03 | −0.20 |
| $Y_4$: CNV-frontal, $n$ | 0.70 | 0.56 | 0.37 | 1 | 0.67 | −0.22 | 0.19 | 0.12 |
| $Y_5$: CNV-central, $n$ | 0.61 | 0.78 | 0.68 | 0.78 | 1 | 0.72 | 0.07 | −0.15 |
| $Y_6$: CNV-parietal, $n$ | −0.23 | −0.52 | −0.77 | 0.46 | 0.79 | 0.1 | 0.06 | 0.23 |
| $X_1$: EOG, $g - n$ | −0.12 | −0.13 | −0.04 | 0.24 | 0.19 | 0.22 | 1 | −0.03 |
| $X_2$: Trait anxiety | 0.10 | 0.30 | 0.28 | 0.17 | 0.32 | 0.36 | 0.03 | 1 |
| Mean | −7.18 | −3.61 | 5.63 | −6.25 | −2.49 | 5.95 | −0.54 | 39.81 |
| Stand. dev. | 4.96 | 5.76 | 5.03 | 4.83 | 5.97 | 5.64 | 8.21 | 9.94 |

early interval after the stimulus (Glanzmann and Fröhlich, 1986; Hänsel, 1992). The potential explanatory variables are the personality characteristic anxiety ($X_1$) and attention ($X_2$) measured as the difference in eye movements under the two experimental conditions (EOG, $g$-$n$: Electrooculogram, difference between 'go'- and 'nogo'-situation). Extreme values in EOG are taken as an indication that no unconfounded measurement of brain activities is possible, hence persons with such values are excluded.

Derived responses were calculated for three corresponding sets of four variables $[Y_1, Y_4, X_1, X_2]$, $[Y_2, Y_5, X_1, X_2]$, and $[Y_3, Y_6, X_1, X_2]$, in which we have the same explanatory variables and the same kind of responses under the two experimental conditions; only the location of the measurement is different for the three variable sets (compare Table 6). The squared canonical correlations are (0.02, 0.24), (0.10, 0.23), and (0.07, 0.22), respectively, and the calculated transformation matrices ((1, 0.54) & (1, −0.65); (1, 0.78) & (1, −0.91); (1, 0.21) & (1, −0.71)) suggest for all three locations of measurement that we can take as derived responses the sums and the differences of the measurements under the two experimental conditions. This leads to the observed association structure displayed in the lower part of Table 7.

While the correlations of the original responses in Table 6 are not close to any easily interpretable structure a rather simple interpretation arises from the derived responses. The level of the brain activity under the two experimental condition, i.e. the sum of CNV, relates only to trait anxiety, while the difference in brain activity under the two conditions relates only to the attention or arousal of the participant. This is reflected in Table 7 in how close most of the correlations estimated under the hypothesis $(Y_1^*, Y_2^*, Y_3^*, X_2) \perp\!\!\!\perp (Y_4^*, Y_5^*, Y_6^*, X_1)$ are to those observed and in the overall chi square statistic of 20.31 on 16 degrees of freedom in a likelihood test of goodness of fit. The correlations in Table 7 indicate further that site of measurement plays a different role for the two explanatory variables: the weakest nonzero correlation of $X_1$, the difference

Table 7

Observed marginal correlations (lower half), marginal correlations estimated under hypothesis $(Y_1^*, Y_2^*, Y_3^*, X_2) \perp\!\!\!\perp (Y_4^*, Y_5^*, Y_6^*, X_1)$ (upper half) with derived responses for the data of Table 6

| Variable | $Y_1^*$ | $Y_2^*$ | $Y_3^*$ | $Y_4^*$ | $Y_5^*$ | $Y_6^*$ | $X_1$ | $X_2$ |
|---|---|---|---|---|---|---|---|---|
| $Y_1^* = Y_1 + Y_4$ | 1 | 0.79 | 0.44 | 0 | 0 | 0 | 0 | 0.15 |
| $Y_2^* = Y_2 + Y_5$ | 0.79 | 1 | 0.78 | 0 | 0 | 0 | 0 | 0.33 |
| $Y_3^* = Y_3 + Y_6$ | 0.44 | 0.78 | 1 | 0 | 0 | 0 | 0 | 0.34 |
| $Y_4^* = Y_1 - Y_4$ | 0.04 | 0.07 | -0.09 | 1 | 0.82 | 0.65 | -0.48 | 0 |
| $Y_5^* = Y_2 - Y_5$ | -0.06 | -0.06 | -0.19 | 0.82 | 1 | 0.82 | -0.48 | 0 |
| $Y_6^* = Y_3 - Y_6$ | 0.05 | -0.01 | -0.18 | 0.65 | -0.82 | 1 | -0.40 | 0 |
| $X_1$: EOG, $g - n$ | 0.06 | 0.03 | 0.10 | -0.48 | -0.48 | -0.40 | 1 | 0 |
| $X_2$: Trait anxiety | 0.15 | 0.33 | 0.34 | -0.09 | -0.06 | -0.17 | 0.03 | 1 |
| Mean | -13.53 | -6.10 | 11.57 | -0.94 | -1.12 | -0.32 | -0.54 | 39.81 |
| Stand. dev. | 9.03 | 11.05 | 10.05 | 3.77 | 3.91 | 3.63 | 8.21 | 9.94 |

in eye movements, is to the parietal measurement, $Y_6$; the weakest nonzero correlation of $X_1$, trait anxiety, is to the frontal measurement $Y_1$, but that there is nevertheless a replication of the results at the different sites: the correlations $(-0.48, -0.48, -0.40)$ and $(0.15, 0.33, 0.34)$ of each of the explanatory variables with one of the derived responses differing in site show the same direction of dependence just different strengths.

If a further condensation of the results were desired, i.e. a summary irrespective of site of measurement, it could be based on the correlation matrix of the derived responses estimated under this independence hypothesis. The regression coefficient matrices for $\hat{Y}^*$ and $X$ have a particularly simple form, so that the canonical correlations are given by the diagonal elements of $\pi_{x|\hat{y}^*}\pi_{\hat{y}^*|x}$ as 0.14 and 0.25 and as transformation $C^T$ we take the regression coefficient matrix obtained by regressing $\hat{Y}^*$ on $X$, i.e.

$$\pi_{x|\hat{y}^*} = \begin{pmatrix} 0 & 0 & 0 & -0.54 & -0.49 & -0.09 \\ -0.55 & 0.79 & 0.21 & 0 & 0 & 0 \end{pmatrix}.$$

If this transformation is applied to the observed sums and differences in CNV all except two correlations are nearly zero, the weighted sum of differences in CNV, i.e. $(-0.54Y_4^* -0.49Y_5^* -0.09Y_6^*)$, has correlation about 0.50 with $X_1$ the difference in EOG and the weighted sum of sums of CNV has correlation about 0.38 with trait anxiety. If one person with somewhat extreme values for EOG is removed, the first correlation reduces to 0.25, all other conclusions remaining qualitatively unchanged. Except for signs essentially the same correlation structure results if only approximate weights $(0\ 0\ 0\ 1\ 1\ 0)\ \&\ (1\ -1\ 0\ 0\ 0\ 0)$ are used, i.e. if we define derived variables $Y_1^{**} = Y_4^* + Y_5^*$ and $Y_2^{**} = Y_1^* - Y_2^*$ we get correlations close to $(Y_1^{**}, X_1) \perp\!\!\!\perp (Y_2^{**}, X_2)$.

A confirmation of the usefulness of the derived variables would be the applicability to different but similar sets of data and ultimately the demonstration of direct substantive importance.

# Appendix

Table 8
Raw data for the data summaries in Tables 1, 3 and 4 on 44 healthy female patients expecting cosmetic surgery

| Systolic blood pr. | Diastolic blood pr. | Age in years | Height in cm | Weight in kg |
|---|---|---|---|---|
| 115 | 80 | 35 | 163 | 58 |
| 120 | 80 | 31 | 162 | 58 |
| 120 | 80 | 22 | 170 | 59 |
| 110 | 70 | 23 | 164 | 50 |
| 110 | 70 | 18 | 179 | 58 |
| 105 | 60 | 28 | 167 | 62 |
| 150 | 95 | 51 | 163 | 67 |
| 110 | 75 | 25 | 164 | 54 |
| 110 | 60 | 22 | 170 | 54 |
| 115 | 70 | 19 | 183 | 60 |
| 100 | 80 | 32 | 165 | 63 |
| 100 | 60 | 26 | 160 | 57 |
| 110 | 80 | 26 | 170 | 58 |
| 120 | 80 | 25 | 170 | 67 |
| 110 | 60 | 22 | 167 | 55 |
| 110 | 60 | 21 | 174 | 69 |
| 110 | 60 | 22 | 167 | 53 |
| 120 | 70 | 36 | 168 | 75 |
| 140 | 80 | 57 | 161 | 70 |
| 110 | 80 | 43 | 165 | 75 |
| 110 | 70 | 17 | 157 | 41 |
| 90 | 60 | 22 | 163 | 63 |
| 115 | 80 | 30 | 160 | 55 |
| 140 | 80 | 25 | 165 | 57 |
| 110 | 70 | 20 | 170 | 65 |
| 115 | 70 | 26 | 168 | 55 |
| 100 | 70 | 24 | 170 | 69 |
| 125 | 80 | 29 | 165 | 72 |
| 140 | 90 | 44 | 163 | 75 |
| 120 | 75 | 22 | 168 | 58 |
| 110 | 70 | 28 | 170 | 58 |
| 120 | 80 | 39 | 158 | 52 |
| 130 | 70 | 25 | 174 | 65 |
| 105 | 65 | 25 | 174 | 65 |
| 110 | 70 | 23 | 168 | 55 |
| 110 | 80 | 34 | 170 | 93 |
| 110 | 70 | 36 | 164 | 60 |
| 130 | 80 | 23 | 165 | 59 |
| 100 | 60 | 27 | 169 | 70 |
| 120 | 80 | 51 | 172 | 85 |
| 115 | 80 | 23 | 176 | 65 |
| 100 | 60 | 49 | 173 | 75 |
| 150 | 100 | 55 | 180 | 90 |
| 100 | 70 | 18 | 173 | 82 |

## Acknowledgement

## References

Cox, D.R. and Wermuth, N. (1992), On the calculation of derived variables in the analyses of multivariate responses, *J. Multiv. Analysis.* **42**, 167–172.

Cox, D.R. and Wermuth, N. (1993), Linear dependencies represented by chain graphs, *J. Statist. Science*, **8**, 204–218.

Glanzmann, P.G. and Fröhlich, W.D. (1986), Anxiety and covert changes of attention control, In: Hentschel, U., Smith, J.G., Draguns, J.G. (Eds.) *The Roots of Perception*, 381–400, Amsterdam: Elsevier Science.

Hänsel, F. (1992) Kortikale Aufmerksamkeitsregulation, Angstneigung und subjektive Erwartung, Psychological Dissertation. Universität Mainz.

Hotelling, H. (1936), Relations between two sets of variates, *Biometrika*, **28**, 321–377.

Krohne, H.W., Kleemann, P.P., Hardt, J. and Theisen, A. (1989), Beziehungen zwischen Bewältigungsstrategien und präoperativen Streßreaktionen, *Z. Klinische Psychologie*, **18**, 350–364.

Lauritzen, S.L. and Wermuth, N. (1989), Graphical models for associations between variables, some of which are qualitative and some quantitative, *Ann. Statist.*, **17**, 31–57.

Rao, C.R. (1973) *Linear Statistical Inference and its Applications*, New York: Wiley.

Slangen, K., Kleemann, P.P. & Krohne, H.W, Coping with surgical stress, In: Krohne, H.W. (ed.) *Attention and avoidance*; *strategies in coping with aversiveness*, New York, Heidelberg: Springer, pp. 321–348.

Wermuth, N. and Rüssmann, H. (1993); Eigenanalysis of symmetrizeable matrix products: a result with statistical applications. *Scand. J. Statist.*, **20**, 361–367.