W. WYSOCKI (Warszawa)

# MATHEMATICAL FOUNDATIONS OF MULTIVARIATE PATH ANALYSIS

*Abstract.* A generalization of path analysis to the multivariate case is proposed. The basic definitions which are not related to causality measurement are given independently of the dimension of the random vectors considered. Causality measurement is based on the concept of linear regression of a random vector w.r.t. a system of random vectors.

**1. Introduction.** Traditional path analysis was originated by works of the geneticist S. Wright [6]. Applications of the theory to genetics were also given by C. C. Li [3]. Nowadays path analysis is more frequently used when dealing with causal modelling in social sciences. This approach had often been abused, before the formal foundations of the theory were stated in Moran [4] and Carlin [1]. It seems that the main problem with using this statistical technique is a clear understanding what tacit assumptions are being made.

In this paper the formal foundations of path analysis are generalized to the multidimensional case. Multivariate path analysis is understood as the study of linear properties of the causal model $(\mathcal{X}, (\mathbf{P}_*, \mathbf{P}^*))$.

We define a causal system as $(\mathcal{X}, (I_*, I^*))$, where $\mathcal{X} = (X_1, \ldots, X_{m+n})$ is a system of $m + n$ $k$-dimensional random variables, and $I_*$, $I^*$ are incidence matrices describing the mutual dependences between $X_1, \ldots, X_{m+n}$. The random variables $X_1, \ldots, X_m$ are called *exogenous*, and $X_{m+1} \ldots X_{m+n}$ are called *endogenous*. Causality measurement for the causal system $(\mathcal{X}, (I_*, I^*))$ is based on the concept of linear regression of the vector $X_i$ $(i > m)$ w.r.t. the direct causes of $X_i$.

The corresponding regression coefficients form matrices $\mathbf{P}_*$ and $\mathbf{P}^*$ which are called path matrices. They are considered as measures of direct influence on endogenous variables exerted by their exogenous causes. The residuals pertaining to the linear regression are treated as endogenous variables in the model. The graphical representation of the linear causal model $(\mathcal{X}, (\mathbf{P}_*, \mathbf{P}^*))$ is given in the form of a complete path diagram. The diagram consists of $m + 2n$ points, each representing one variable (endogenous, exogenous or residual). All direct causes of the endogenous variable $X_i$ $(i > m)$ are linked on the diagram with $X_i$ by a single arrow with head at $X_i$. To each arrow there corresponds an element of the matrix $\mathbf{P}_*$ or $\mathbf{P}^*$. The residuals of $X_i$ $(i > m)$ are also linked by such an arrow with $X_i$, and in this case the identity matrix corresponds to this arrow. Any two points representing exogenous variables are linked by an arrow with two heads in opposite directions, provided the corresponding covariance matrix is not the zero matrix. The rules of moving in the complete path diagram are dealt with in Section 4.2.

In Section 2 the problem of causality measurement for multivariate random variables is stated. The relevant mathematical framework is contained in Section 3. In Section 4.1 simple conditions are given under which the construction of a linear causal model is possible. Moreover, the meaning of residuals in the context of causality measurement is explained.

In traditional path analysis all random variables of the system $\mathcal{X}$ are assumed to be normalized. Since for some applications this assumption appears to be too restrictive, we do not impose it in the paper.

## 2. Causal linear model

**2.1.** *Definition.* Let $\mathcal{H}(\Omega, \mathcal{A}, P; H)$ be the space of all random variables defined on a probability space $(\Omega, \mathcal{A}, P)$ with values in a normed space $H$ of finite dimension.

We distinguish a system $\mathcal{X}$ of $m + n$ elements of this space: $\mathcal{X} = (X_1, \ldots, X_m, X_{m+1}, \ldots, X_{m+n})$. The random variables $X_1, \ldots, X_m$ are called *exogenous*, and $X_{m+1}, \ldots, X_{m+n}$ are *endogenous*. Moreover, we distinguish two functions:

$$f_* : \{X_{m+1}, \ldots, X_{m+n}\} \times \{X_{m+1}, \ldots, X_{m+n}\} \to \{0, 1\}$$

where $f_*$ is such that $f_*(X_i, X_i) = 0$ for $i = m + 1, \ldots, m + n$, and

$$f^* : \{X_1, \ldots, X_m\} \times \{X_{m+1}, \ldots, X_{m+n}\} \to \{0, 1\}.$$

These functions can be represented in compact form by the matrices

(1)                     $I_* = [f_*(X_{m+j}, X_{m+i})], \quad i, j = 1, \ldots, n,$

and

(2) $$I^* = [f^*(X_j, X_{m+i})], \qquad i = 1, \ldots, n, \ j = 1, \ldots, m$$

(the indices $i$ and $j$ correspond to rows and columns, respectively).

DEFINITION 1. The function $f_*$ (resp. $f^*$) is called an *endogenous* (resp. *exogenous*) *causal function* for the system $\mathcal{X}$. The matrix $I_*$ (resp. $I^*$) is an *endogenous* (resp. *exogenous*) *incidence matrix* for $\mathcal{X}$.

If $f_*(X_j, X_i) = 1$ (resp. if $f^*(X_{j'}, X_{i'}) = 1$) then we say that $X_j$ (resp. $X_{j'}$) is a *direct cause* of $X_i$ (resp. of $X_{i'}$) and that $X_i$ (resp. $X_{i'}$) is a *direct effect* of $X_j$ (resp. of $X_{j'}$); symbolically, we write $X_j \to X_i$ (resp. $X_{j'} \to X_{i'}$). If $X_{j_i} \to X_{j_{i+1}}$ for $i = 1, \ldots, \nu - 1$, $\nu \geq 3$, then we say that $X_{j_1}$ is an *indirect cause* of $X_{j_\nu}$ (and $X_{j_\nu}$ is an *indirect effect* of $X_{j_1}$).

Suppose that the $i$th row of $I_*$ ($i = 1, \ldots, n$) has 1's in columns $j_1, \ldots$ $\ldots, j_{\nu(i)}$; in other words, $X_{m+j_1}, \ldots, X_{m+j_{\nu(i)}}$ are the endogenous direct causes of the endogenous variable $X_{m+i}$. On the other hand, let $X_{j'_1}, \ldots$ $\ldots, X_{j'_{\nu'(i)}}$ be the exogenous direct causes of $X_{m+i}$, where $j'_1, \ldots, j'_{\nu'(i)}$ indicate the columns of $I^*$ with 1's in the $i$th row. The set of all direct causes of the endogenous variable $X_i$ for $i = m + 1, \ldots, m + n$ is denoted by $\mathcal{X}_i$.

DEFINITION 2. The ordered pair $(\mathcal{X}, (I_*, I^*))$ is called a *causal system*.

Exogenous variables have no indirect causes. Any variable can be a direct or indirect effect of other variables; a direct or indirect cause (effect) will be shortly called a *cause* (*effect*). It can happen that an exogenous variable is an indirect cause of itself.

DEFINITION 3. The causal system $(\mathcal{X}, (I_*, I^*))$ is called *nonrecursive* if for some $i_0$ ($i_0 = m + 1, \ldots, m + n$) the variable $X_{i_0}$ of this system is an indirect cause of itself. Otherwise, the system is called *recursive*.

Recursivity of a causal system can be easily characterized by means of the endogenous incidence matrix $I_*$.

LEMMA 1. *The causal system* $(\mathcal{X}, (I_*, I^*))$ *is recursive if and only if* $I_*$ *is nilpotent of order* $n$ (*i.e.* $I_*^n = 0$).

The proof is analogous to that of a similar lemma of Kang and Seneta [2] for the one-dimensional case.

R e m a r k. It follows from the proof given by Kang and Seneta that recursivity of a causal system occurs if and only if the matrix $I_*$ is essentially lower triangular, possibly after some permutation of rows and columns. This permutation of rows and columns corresponds to a permutation $(X_{i_1}, \ldots, X_{i_n})$ of all endogenous variables such that $X_{i_j}$ is not a cause of $X_{i_{j'}}$ for $j > j'$, $j' = 1, \ldots, n$.

The condition $I_*^n = 0$ implies that $I_n - I_*$ is invertible, where $I_n$ is the identity matrix of order $n$. It is easy to check that $(I_n - I_*)^{-1} = \sum_{i=0}^{n-1} I_*^i$.

EXAMPLE. Let $\mathcal{X} = (X_1, \dots, X_5)$ be a system of random variables. Let

$$I_* = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \end{bmatrix}, \quad I^* = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \end{bmatrix}.$$

Thus, $m = 1$, $n = 4$. The causal system $(\mathcal{X}, (I_*, I^*))$ can be graphically represented by a path diagram (Fig. 1) which consists of $m + n$ points corresponding to the variables in $\mathcal{X}$.
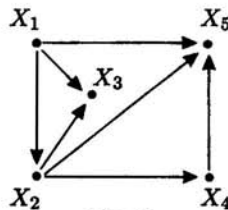


Fig. 1

Any endogenous variable $X_i$, $i = m+1, \dots, m+n$, is joined to any direct cause (endogenous or exogenous) of $X_i$ by an arrow with head in $X_i$. Such arrows are called *single step paths*.

**2.2.** *Path matrices.* The causal system $(\mathcal{X}, (I_*, I^*))$ indicates the links between causes and effects, but the strength of these links is not measured. Now, we present an intuitive background of such a measurement.

Consider an arbitrary endogenous variable $X_i$ from $(\mathcal{X}, (I_*, I^*))$, $i = m+1, \dots, m+n$. The direct causes (exogenous and endogenous) of $X_i$ form the set $\mathcal{X}_i$. The indices of these direct causes can be read off from the $i$th rows of $I_*$ and $I^*$. Let $X_{m+j_1}, \dots, X_{m+j_{\nu(i)}}$ and $X_{j'_1}, \dots, X_{j'_{\nu'(i)}}$ be the endogenous and exogenous direct causes of $X_i$, respectively.

If $H = \mathbb{R}^k$ then $X_i = (X_{i1}, \dots, X_{ik})^T$, $i = 1, \dots, m+n$.

We assume that the conditional expectation of any endogenous variable $X_i$ ($i = m+1, \dots, m+n$), conditioned on the set $\mathcal{X}_i$ of its direct causes, is of the form

$$(3) \qquad E(X_i \mid \mathcal{X}_i) = \sum_{l=1}^{\nu(i)} B_{*j_l} X_{m+j_l} + \sum_{l=1}^{\nu'(i)} B_{j'_l}^* X_{j'_l}$$

where $B_{*j_l}$, $l = 1, \dots, \nu(i)$, and $B_{j'_l}^*$, $l = 1, \dots, \nu'(i)$, are some $k \times k$ matrices.

We decompose the endogenous variable $X_i$, $i = m+1, \dots, m+n$, in $\mathcal{H}(\Omega, A, P; H)$ as

$$(4) \qquad X_i = E(X_i \mid \mathcal{X}_i) + U_i,$$

where $U_i \in \mathcal{H}(\Omega, A, P; H)$; $U_i$ is called the *residual* of $X_i$ (w.r.t. $\mathcal{X}_i$).

Geometrically, $E(X_i \mid \mathcal{X}_i)$ is the projection of $X_i$ on the subspace $\mathcal{H}_{\mathcal{X}_i}$, where

$$\mathcal{H}_{\mathcal{X}_i} = \Big\{ \sum_{l=1}^{\nu(i)} A_{*j_l} X_{m+j_l} + \sum_{l=1}^{\nu'(i)} A^*_{j'_l} X_{j'_l} : A_{*j_l}, A^*_{j'_l} \in \mathcal{M}_k \Big\}$$

and $\mathcal{M}_k$ is the set of all real $k \times k$ matrices.

Formula (4) is equivalent to

$$(5) \qquad X_i = \sum_{l=1}^{\nu(i)} B_{*j_l} X_{m+j_l} + \sum_{l=1}^{\nu'(i)} B^*_{j'_l} X_{j'_l} + I_k U_i \,,$$

where $I_k$ is the $k \times k$ identity matrix. To shorten notation, we introduce some changes and define some block matrices. Let

$$X_i = Z_i\,, \quad i = 1, \ldots, m\,, \quad X_{m+1} = Y_i\,, \quad i = 1, \ldots, n\,,$$

$$\mathbf{Z} = [Z_i]\,, \quad i = 1, \ldots, m\,, \qquad \mathbf{Y} = [Y_i]\,, \quad i = 1, \ldots, n\,,$$

$$(6) \qquad \mathbf{P}_* = [P_{*ij}]\,, \quad i, j = 1, \ldots, n\,,$$

$$(7) \qquad \mathbf{P}^* = [P^*_{ij}]\,, \quad i = 1, \ldots, n\,, \quad j = 1, \ldots, m\,.$$

The elements of $\mathbf{P}_*$ are $k \times k$ matrices. Similarly,

$$P^*_{ij} = \begin{cases} 0 & \text{if the } (i,j) \text{ element of } I^* \text{ is } 0, \\ B^*_{j_l}, & \text{where } j_l = j, \text{ otherwise.} \end{cases}$$

The system of equations (5) is then given by

$$(8) \qquad \mathbf{Y} = \mathbf{P}_* \mathbf{Y} + \mathbf{P}^* \mathbf{Z} + \mathbf{U}\,, \quad \mathbf{U} = [U_i]\,.$$

Let $\mathbf{P}^i_*$ and $\mathbf{P}^{*i}$ be the $i$th rows of $\mathbf{P}_*$ and $\mathbf{P}^*$, respectively. Then

$$(9) \qquad Y_i = \mathbf{P}^i_* \mathbf{Y} + \mathbf{P}^{*i} \mathbf{Z} + I_k U_i\,.$$

This means that endogenous variable $Y_i$ is a "linear combination" of its direct causes (endogenous and exogenous) and of $U_i$. According to (9), the nonzero elements of $\mathbf{P}^i_*$ and $\mathbf{P}^{*i}$ can be treated as measures of direct influence on $Y_i$ (causality) of direct causes, endogenous and exogenous. Also, the direct influence of $U_i$ on $Y_i$ is represented by $I_k$.

DEFINITION 4. The elements $P_{*ij}$ and $P^*_{ij}$ of $\mathbf{P}_*$ and $\mathbf{P}^*$ are called *path matrices* which measure the direct influence ([1]) of the endogenous $X_{m+j}$ and exogenous $X_j$ on the endogenous $X_{m+i}$.

If the casual system $\mathcal{X}$ is recursive then equality (8) can be rewritten as

$$(10) \qquad \mathbf{Y} = (I - \mathbf{P}_*)^{-1} \mathbf{P}^* \mathbf{Z} + (I - \mathbf{P}_*)^{-1} \mathbf{U}\,.$$

---

([1]) This influence exists if and only if $P_{*ij}$ and $P^*_{ij}$ are not zero matrices.

This is due to the remark following Lemma 1. The block matrix $\mathbf{P}_*$ is nilpotent of order $n$. Hence, the matrix $I - \mathbf{P}_*$ is invertible.

Note that the matrices $\mathbf{P}_*$ and $\mathbf{P}^*$ appearing in (10) are unknown, since in formula (3) the coefficients corresponding to the direct causes $X_i$ are not specified. However, suppose that these matrices are somehow specified. Then

DEFINITION 5. The ordered pair $(\mathcal{X}, (\mathbf{P}_*, \mathbf{P}^*))$ is called a *linear causal model* for the causal system $(\mathcal{X}, (I_*, I^*))$.

## 3. Mathematical background

**3.1.** *The space $L^2(\Omega, \mathcal{A}, P; \mathbb{R}^k)$.* In Section 2 we have dealt with $\mathcal{H}(\Omega, \mathcal{A}, P; H)$. Setting $H = \mathbb{R}^k$ we obtain the space of $k$-dimensional random vectors which is considered in Section 4.

Assume that the covariance matrix exists for each vector and that the expectations are zero. Let $L^2(\Omega, \mathcal{A}, P; \mathbb{R}^k)$ be the subspace of $\mathcal{H}(\Omega, \mathcal{A}, P; \mathbb{R}^k)$ consisting of all centered random vectors.

For any given $k \times k$ matrix $\Lambda$ which is symmetric and positive definite we introduce a scalar product $\langle \, , \, \rangle_\Lambda$ in $L^2(\Omega, \mathcal{A}, P; \mathbb{R}^k)$ by

$$(11) \quad \langle Z_1, Z_2 \rangle_\Lambda = E(Z_1^T \Lambda^{-1} Z_2) = \mathrm{tr}(\Lambda^{-1} \mathrm{cov}(Z_1, Z_2)),$$
$$Z_1, Z_2 \in L^2(\Omega, \mathcal{A}, P; \mathbb{R}^k),$$

where $\mathrm{cov}(Z_1, Z_2)$ is the covariance matrix of the pair of vectors $Z_1, Z_2$. The corresponding norm, denoted by $\| \, \|_\Lambda$, is complete. For any $Z_1, Z_2$, let

$$(12) \quad \varrho_\Lambda(Z_1, Z_2) = \frac{\langle Z_1, Z_2 \rangle_\Lambda}{\|Z_1\|_\Lambda \|Z_2\|_\Lambda}$$
$$= \frac{\mathrm{tr}(\Lambda^{-1} \mathrm{cov}(Z_1, Z_2))}{\mathrm{tr}(\Lambda^{-1} \mathrm{cov}(Z_1, Z_1)))^{1/2} (\mathrm{tr}(\Lambda^{-1} \mathrm{cov}(Z_2, Z_2)))^{1/2}} \,.$$

This index was introduced by Sampson [5] as a measure of dependence between $Z_1$ and $Z_2$.

**3.2.** *Linear regression.* Most of the facts stated here come from the author's paper [8].

Let $Y$ and $X_1, \dots, X_n$ be elements of $L^2(\Omega, \mathcal{A}, P; \mathbb{R}^k)$. Minimizing $\|Y - \sum A_i X_i\|_\Lambda$ over all "linear combinations" $A_1 X_1 + \dots + A_n X_n$, where the $A_i$ are $k \times k$ matrices, we get the vector $R(Y.1, \dots, n) = \sum_{i=1}^n B_i X_i$ with uniquely specified matrices $B_1, \dots, B_n$. The vector $R(Y.1, \dots, n)$ is called the *linear regression* of $Y$ on $(X_1, \dots, X_n)$. It satisfies the equalities

$$\sup \left\{ \varrho_\Lambda \left( Y, \sum A_i X_i \right) : A_i \in \mathcal{M}_k, \, i = 1, \dots, n \right\} = \left( \frac{\|R(Y.1, \dots, n)\|_\Lambda}{\|Y\|_\Lambda} \right)^2$$
$$= \varrho_\Lambda^2(Y; X_1, \dots, X_n).$$

The expression $\varrho_\Lambda(Y; X_1, \ldots, X_n)$ is called the *multiple correlation* between $Y$ and $(X_1, \ldots, X_n)$.

Let $Y = (Y_1, \ldots, Y_k)^T$, $X_i = (X_{i1}, \ldots, X_{ik})^T$, $i = 1, \ldots, n$, and let $\mathbf{X}$ be the column vector $\mathbf{X} = (X_1^T, \ldots, X_n^T)$. Finally, let $\Sigma = \mathrm{cov}(\mathbf{X}, \mathbf{X})$, $\widetilde{\Sigma} = \mathrm{cov}(\mathbf{X}, Y)$, and let $\mathbf{B}$ be the block matrix $\mathbf{B} = [B_i]$.

LEMMA 2. *If $\Sigma$ is nonsingular then the regression of $Y$ on $(X_1, \ldots, X_n)$ is given by*

(13) $$\mathbf{B}^T \mathbf{X} = \widetilde{\Sigma}^T \Sigma^{-1} \mathbf{X}.$$

R e m a r k. The geometrical interpretation of the linear regression of $Y$ on $(X_1, \ldots, X_n)$ is indicated in Wysocki [7]: under the scalar product (11), the linear regression is the orthogonal projection of $Y$ on the subspace $\mathcal{H}_{X_1, \ldots, X_n}$, defined analogously to $\mathcal{H}_{\mathcal{X}_i}$ from Section 2 ($\mathcal{X}_i = (X_1, \ldots, X_n)$).

The residual $U = Y - \mathbf{B}^T \mathbf{X}$ is orthogonal to $\mathcal{H}_{X_1, \ldots, X_n}$ and, in particular, to $AX_i$, where $A$ is any $k \times k$ matrix:

$$\langle U, AX_i \rangle_\Lambda = 0, \qquad i = 1, \ldots, n.$$

Finally, we introduce partial covariance matrices. Let $Y_1, Y_2, X_1, \ldots, X_n$ belong to $L^2(\Omega, \mathcal{A}, P; \mathbb{R}^k)$. The *partial covariance matrix* of $(Y_1, Y_2)$ on $(X_1, \ldots, X_n)$, denoted by $\mathrm{cov}(Y_1, Y_2.1, \ldots, n)$, is the covariance matrix for the pair of random vectors $Y_1 - R(Y_1.1, \ldots, n)$ and $Y_2 - R(Y_2.1, \ldots, n)$.

## 4. Path analysis

**4.1.** *Construction of a linear causal model.* Let us specify the assumptions under which a linear causal model exists. Let

$$\mathcal{X}_i = (X_{j_1'}, \ldots, X_{j_{\nu'(i)}'}, X_{m+j_1}, \ldots, X_{m+j_{\nu(i)}}).$$

Recall that for any endogenous $X_i$, $i = m+1, \ldots, m+n$, the primed variables indicate its direct exogenous causes, and the nonprimed ones its endogenous causes. Define

$$\mathbf{X}_i = (X^T, \ldots, X_{j_{\nu'(i)}'}^T, X_{m+j_1}^T, \ldots, X_{m+j_{\nu(i)}}^T)^T, \qquad i = 1, \ldots, n,$$

$$\Sigma_i = \mathrm{cov}(\mathbf{X}_i, \mathbf{X}_i), \qquad i = 1, \ldots, n,$$

$$\widetilde{\Sigma}_i = \mathrm{cov}(\mathbf{X}_i, Y_i), \qquad i = 1, \ldots, n,$$

$$\mathbf{B}_i = [B_l]^T, \quad B_l \in \mathcal{M}_k, \qquad l = j_1', \ldots, j_{\nu'(i)}', m+j_1, \ldots, m+j_{\nu(i)}.$$

LEMMA 3. *Suppose that the variables from $(\mathcal{X}, (I_*, I^*))$ satisfy the following conditions:*

(i) *all variables belong to $L^2(\Omega, \mathcal{A}, P; \mathbb{R}^k)$;*

(ii) *the conditional expectation of any endogenous $X_i$ ($i = m + 1, \ldots, m + n$) given its direct causes $\mathcal{X}_i$, is a.e. equal to the linear regression of $X_i$ given $\mathcal{X}_i$;*

(iii) *the matrices $\Sigma_i$ $(i = 1, \ldots, n)$ are nonsingular.*

*Then there exists exactly one linear causal model $(\mathcal{X}, (\mathbf{P}_*, \mathbf{P}^*))$.*

P r o o f. The matrices (6) and (7) have to be uniquely determined. We use Lemma 2 substituting $X = X_i$, $\Sigma = \Sigma_i$, $B = B_i$. By (13), $B_i^T = \widetilde{\Sigma}_i^T \Sigma_i^{-1}$. From (i) and (iii) it follows that $B_i$ exists and is uniquely determined. Moreover, for any $i = 1, \ldots, n$,

$$P_{*ij} := \begin{cases} B_{m+j} & \text{if } j = j_l, \ l = 1, \ldots, \nu_i, \\ 0 & \text{if } j \in \{1, \ldots, n\} \backslash \{j_1, \ldots, j_{\nu'(i)}\}, \end{cases}$$

$$P_{ij}^* := \begin{cases} B_j & \text{if } j = j'_l, \ l = 1, \ldots, \nu'_i, \\ 0 & \text{if } j \in \{1, \ldots, m\} \backslash \{j'_1, \ldots, j'_{\nu'(i)}\}. \end{cases}$$

This ends the proof.

Now, let us change the notation in the model $(\mathcal{X}, (I_*, I^*))$ to make it consistent with Section 2.2.

By (10), the vector $\mathbf{U} = [U_1, \ldots, U_n]$ is uniquely determined. The vectors $U_1, \ldots, U_n$ can be treated in $(\mathcal{X}, (\mathbf{P}_*, \mathbf{P}^*))$ as direct exogenous causes of $X_{m+i}$, $i = 1, \ldots, n$. Let $\widetilde{\mathcal{X}}$ denote the system $\mathcal{X}$ augmented by the residuals $U_1, \ldots, U_n$. For convenience, the set of all exogenous variables in $\widetilde{\mathcal{X}}$ (including residuals) is denoted by $(U_1, \ldots, U_m)$. Then

$$\widetilde{\mathcal{X}} = (U_1, \ldots, U_m, \ X_{m+1}, \ldots, X_{m+n}).$$

So we have a new block matrix

$$\mathbf{P} = [\widetilde{\mathbf{P}}^* \mathbf{P}_*] = [P_{ij}], \quad i = m + 1, \ldots, m + n, \ j = 1, \ldots, m + n,$$

where $\widetilde{\mathbf{P}}^*$ is a suitable modification of $\mathbf{P}^*$. The linear causal model corresponding to $\widetilde{\mathcal{X}}$ is called an *augmented linear causal model* for $(\mathcal{X}, (I_*, I^*))$ and is denoted by $(\widetilde{\mathcal{X}}, \mathbf{P})$.

To construct the linear causal model $(\mathcal{X}, (\mathbf{P}_*, \mathbf{P}^*))$ in practice, we need the incidence matrices $(I_*, I^*)$ for $\mathcal{X}$ and the covariance matrix for this system. The starting point is to determine $\mathcal{X} = (X_1, \ldots, X_{m+n})$ in such a way that the $X_i$, $i = 1, \ldots, m + n$, have a real interpretation. Then we determine the incidence matrices $I_*$ and $I^*$ which represent direct causal relations between the real entities. Mathematically, these matrices can be arbitrary.

The path matrices (elements of $\mathbf{P}_*$ and $\mathbf{P}^*$) "measure" the direct influence of direct endogenous and exogenous causes on endogenous variables.

LEMMA 4. *Let $(\widetilde{\mathcal{X}}, \mathbf{P})$ be the augmented linear causal model for the recursive causal system $(\mathcal{X}, (I_*, I^*))$. The following conditions are equivalent:*

(i) *for any $(i, j)$ such that $X_i$ is an indirect cause of $X_j$, the partial*

covariance matrix of $(X_i, X_j)$ on all direct causes $X_j$ is the zero matrix:

$$\operatorname{cov}(X_i, X_j . \mathcal{X}_j) = 0 \, ;$$

(ii) *for any* $(i', j')$, $i' \neq j'$, *such that* $X'_i$ *is a cause of* $X'_j$, *the covariance matrix of* $(U_{i'}, U_{j'})$ *is the zero matrix.*

The proof is omitted as it is analogous to that of Kang and Seneta [2] in the one-dimensional case.

If $(\mathcal{X}, \mathbf{P})$ is a linear causal model then the covariance matrix of $U$ is easily obtained from (8):

$$\operatorname{cov}(\mathbf{U}, \mathbf{U}) = (I - \mathbf{P}_*) \operatorname{cov}(\mathbf{Y}, \mathbf{Y})(I - \mathbf{P}_*)^T - (I - \mathbf{P}_*) \operatorname{cov}(\mathbf{Y}, \mathbf{Z}) \mathbf{P}^{*T}$$
$$- \mathbf{P}^* \operatorname{cov}(\mathbf{Z}, \mathbf{Y})(I - \mathbf{P}_*) + \mathbf{P}^* \operatorname{cov}(\mathbf{Z}, \mathbf{Z}) \mathbf{P}^{*T} \, .$$

**4.2.** *Decomposition of the covariance matrix.* An augmented linear causal model $(\widetilde{\mathcal{X}}, \mathbf{P})$ for a causal system $(\mathcal{X}, (I_*, I^*))$ can be graphically represented as the so-called *complete path diagram*. It consists of the causal system $(\mathcal{X}, (I_*, I^*))$ and of $n$ points which represent the residuals $U_1, \ldots, U_n$ appearing in the linear causal model. Each point $U_i$ $(i = 1, \ldots, n)$ is linked with each point representing an endogenous variable $X_{m+i}$ by a single arrow with head at $X_{m+i}$. A pair of points from the union of points representing exogenous variables and residuals is linked by an arrow with two heads in opposite directions provided that the respective covariance matrix is not the zero matrix.

To any single arrow starting at a point representing an endogenous variable or an exogenous variable which is not a residual we attach the respective path matrix belonging to $\mathbf{P}_*$ or $\mathbf{P}^*$; to any single arrow linking $U_i$ and $X_{m+i}$ we attach the respective identity matrix; to any double arrow we attach the respective covariance matrix.

Now we generalize the famous Wright rules of moving in the complete path diagram of the augmented linear causal model.

Let $(\widetilde{\mathcal{X}}, \mathbf{P})$ be the augmented linear causal model for the recursive causal system $(\mathcal{X}, (I_*, I^*))$. Suppose that $X_i$ and $X_j$ are two endogenous variables, represented as

$$X_i = \sum_{X_{i'} \in \mathcal{X}_i} P_{ii'} X_{i'} \, , \qquad X_j = \sum_{X_{j'} \in \mathcal{X}_j} P_{jj'} X_{j'} \, .$$

Then

$$(14) \quad \operatorname{cov}(X_i, X_j) = \sum_{X_l \in \mathcal{X}_i \cap \mathcal{X}_j} P_{il} \operatorname{cov}(X_l, X_l) P_{jl}^T$$
$$+ \sum_{X_{i'} \in \mathcal{X}_i, \, X_{j'} \in \mathcal{X}_j, \, X_{i'} \neq X_{j'}} P_{ii} \operatorname{cov}(X_{i'}, X_{j'}) P_{jj}^T \, .$$

Putting $X_i = X_j$ we get the so-called *complete determination equation*.

Equation (14) has an interesting interpretation. Note that each term in the first sum of the right-hand side of (14) corresponds to a two-step path: $X_i \leftarrow X_l \rightarrow X_j$. This can be interpreted as moving in the diagram in the following manner: from $X_i$ to $X_l$ (this direction is opposite to the direction of the arrow) and from $X_l$ to $X_j$ as indicated by the arrow. The variable $X_l$ is a common direct cause of $X_i$ and $X_j$.

The terms in the second sum of the right-hand side of (14) correspond to a three-step path of the form

$$X_i \leftarrow X_{i'} \leftrightarrow X_{j'} \rightarrow X_j.$$

Here, $X_{i'}$ is a direct cause of $X_i$, while $X_{j'}$ is a direct cause of $X_j$.

Paths of the first and second kinds correspond to the matrices $P_{il} \operatorname{cov}(X_l, X_l) P_{jl}^T$ and $P_{il'} \operatorname{cov}(X_{i'}, X_{j'}) P_{jj'}^T$. Thus, the covariance matrix of $(X_i, X_j)$ is the sum of all such products of matrices over both kinds of paths linking $X_i$ and $X_j$.

Equation (14) and the complete determination equation can be applied to covariance matrices appearing on the right-hand side of (14), yielding a further decomposition of $\operatorname{cov}(X_i, X_j)$.

LEMMA 5. *If $X_i$ and $X_j$ are two arbitrary variables of an augmented linear causal model $(\widetilde{\mathcal{X}}, \mathbf{P})$ for a recursive causal system $(\mathcal{X}, (I_*, I^*))$ then*

$$(15) \quad \operatorname{cov}(X_i, X_j)$$

$$= {\sum}' P_{i_1 i_2} P_{i_2 i_3} \ldots P_{i_{r-1} i_r} \operatorname{cov}(X_{i_r}, X_{i_r}) P_{i_{r+1} i_r}^T P_{i_{r+2} i_{r+1}}^T \cdots P_{i_{r+s} i_{r+s-1}}^T$$

$$+ {\sum}'' P_{j_1 j_2} P_{j_2 j_3} \ldots P_{j_{q-1} j_q} \operatorname{cov}(X_{j_q}, X_{j_{q+1}}) P_{j_{q+2} j_{q+1}}^T \cdots P_{j_{q+t} j_{q+t-1}}^T$$

*where the summation $\sum'$ is taken over all paths without loops which have the form*

$$(16) \quad X_i = X_{i_1} \leftarrow X_{i_2} \leftarrow \ldots \leftarrow X_{i_r} \rightarrow X_{i_{r+1}} \rightarrow \ldots \rightarrow X_{i_{r+s}} = X_j,$$

*and $\sum''$ over all paths without loops which have the form*

$$(17) \quad X_i = X_{j_1} \leftarrow X_{j_2} \leftarrow \ldots \leftarrow X_{j_q} \leftrightarrow X_{j_{q+1}} \rightarrow \ldots \rightarrow X_{j_{q+t}} = X_j.$$

The proof is analogous to that of Kang and Seneta [2] for the univariate case and is omitted.

Formula (15) has a clear interpretation in the complete path diagram. In case of a path described by (16) we start from $X_i$ and move against the direction of the arrows till the point $X_{i_r}$ is reached and then proceed to $X_j$ according to the direction of the arrows. The variable $X_{i_r}$ is a common cause of $X_i$ and $X_j$. On a path of the form (17) we move analogously; the only difference is that instead of one changing point $X_{i_r}$ we have two variables $X_{j_q}$ and $X_{j_{q+1}}$ linked by a double arrow. Here, $X_{j_q}$ and $X_{j_{q+1}}$ are exogenous causes of $X_i$ and $X_j$, respectively.

### References

[1]  S. Carlin, *Causal models: an attempt at a unified approach*, Honours Essay, Dept. Math., University of Western Australia, Nedlands 1977.

[2]  K. M. Kang and E. Seneta, *Path analysis: an exposition*, in: Developments in Statistics, Vol. 3, Academic Press, New York 1980, 217–246.

[3]  C. C. Li, *Path Analysis—A Primer*, Boxwood Press, Pacific Grove, CA, 1975.

[4]  P. A. P. Moran, *Path coefficients reconsidered*, Austral. J. Statist. 3 (1961), 87–93.

[5]  A. R. Sampson, *Positive dependence properties of elliptically symmetric distributions*, J. Multivariate Anal. 13 (1983), 375–381.

[6]  S. Wright, *Correlation and causation*, J. Agric. Res. 20 (1921), 557–585.

[7]  W. Wysocki, *Geometrical aspects of measures of dependence for random vectors*, Zastos. Mat. 21 (1991), 211–224.

[8]  —, *Maximal correlation in path analysis*, ibid., 225–233.

WŁODZIMIERZ WYSOCKI
INSTITUTE OF COMPUTER SCIENCE
POLISH ACADEMY OF SCIENCES