

Kálmán's filtering technique in structural equation modeling

Marianna Bolla and Fatma Abdelkhalek

Dedicated to Professor Gheorghe Moroşanu on the occasion of his 70th anniversary.

Abstract. Structural equation modeling finds linear relations between exogenous and endogenous latent and observable random vectors. In this paper, the model equations are considered as a linear dynamical system to which the celebrated R. E. Kálmán's filtering technique is applicable. An artificial intelligence is developed, where the partial least squares algorithm of H. Wold and the block Cholesky decomposition of H. Kiiveri et al. are combined to estimate the parameter matrices from a training sample. Then the filtering technique introduced is capable to predict the latent variable case values along with the prediction error covariance matrices in the test sample. The recursion goes from case to case along the test sample, without having to re-estimate the parameter matrices. The algorithm is illustrated on real life sociological data.

Mathematics Subject Classification (2010): 62H05, 62P25, 68T99.

Keywords: Structural equation modeling, linear dynamical systems, Kálmán's filtering, artificial intelligence, application to social sciences.

1. Introduction

We consider structural equation model (SEM) for two latent random vectors that depend through a linear model on two observable random vectors, respectively (they usually include exogenous and endogenous variables). This kind of models was first investigated by T. Haavelmo [1], who obtained the Nobel Prize for it later. Unlike the traditional factor analysis, where latent variables were introduced and given a meaning based on the factor loadings, here the latent variables are organic parts of the model. The latent variables, e.g., alienation, ambition in [2] or mobility in our example, are given by the experts, and the observed (measurement) variables are indicators of them. In this way, so-called inner and outer relations are stated between the latent variables and between the observable and latent ones, respectively.

The estimation of the parameter matrices of this model was elaborated both in the Gaussian and distribution-free cases, former by K. G. Jöreskog (LISREL) [2], while the other by H. Wold (PLS) [8] in the 1970-1980s. These two approaches are sometimes called covariance-based and component-based SEM. However, we can consider the model equations as a linear dynamical system to which the celebrated R. E. Kálmán's filtering technique [3] is applicable. This technique was developed in the 1960s for time series to make predictions for the hidden state variables of a state space model, and was used in the lunar landing, for instance. We will show how to apply this technique in the more complicated dynamical system, containing two state and two observable equations, describing inner relations between the observable and latent variables, both for the exogenous and endogenous ones. Our contribution is that we connect these two approaches.

The parameter matrices are estimated from a training sample. We combine the first stage of the PLS algorithm of H. Wold to estimate the case values of the latent variables and the method of H. Kiiveri et al. [5] to decompose the inverse of the product moment matrix obtained with the latent case value estimates. At this point, we apply the block Cholesky decomposition. Then the filtering technique to be introduced is capable to make predictions for the endogenous variables based on the exogenous ones, through the latent variables. The driving force is that we propagate the error covariance matrices of the exogenous and endogenous latent variables in a recursion.

The test sample is a succession of observations coming one by one (like a time series or just subsequent observations), and the algorithm predicts their endogenous variables based on their own exogenous ones. Our contribution is that we combine existing methods for parameter estimation, and then apply filtering technique for prediction, so we develop an artificial intelligence. The computational gain is that the parameter matrices need not be estimated for every new-coming case in the test sample, but are estimated only once, in the training sample. The method is distribution-free (just second moments are used in the linear state equations) and applicable to small training, and not necessarily independent test samples.

The organization of the paper is as follows. In Section 2, the most important notions and facts about best linear predictions in Hilbert spaces are introduced. In Section 3, the prediction and propagation of the error covariance matrices are derived in two stages. The main results are summarized in Theorem 3.1 of Section 3.3. Then the proposed algorithm is illustrated on real life sociological data in Section 4. Eventually, the last Section 5 discusses the benefits of the proposed method together with some possible further perspectives.

2. Preliminaries

The following linear dynamical system that resembles the one to which R. E. Kálmán gave a recursive algorithm is considered:

$$\mathbf{B}\boldsymbol{\eta} = \mathbf{A}\boldsymbol{\xi} + \boldsymbol{\zeta},$$

where $\boldsymbol{\eta}$ is m - and $\boldsymbol{\xi}$ is n -dimensional latent vector, \mathbf{B} and \mathbf{A} are $m \times m$ and $m \times n$ coefficient matrices, and $\boldsymbol{\zeta}$ is a random vector of residuals of uncorrelated components.

It is also uncorrelated with $\boldsymbol{\xi}$, and \mathbf{B} is nonsingular. In the recursive models, \mathbf{B} is upper triangular, with 1s along its main diagonal.

Here $\boldsymbol{\eta}$ and $\boldsymbol{\xi}$ are not observed, but instead, the p -dimensional \mathbf{Y} and the q -dimensional \mathbf{X} are observed such that

$$\mathbf{X} = \mathbf{C}\boldsymbol{\xi} + \boldsymbol{\varepsilon}, \quad \mathbf{Y} = \mathbf{G}\boldsymbol{\eta} + \boldsymbol{\delta},$$

where $\boldsymbol{\varepsilon}$ and $\boldsymbol{\delta}$ are vectors of measurement errors in \mathbf{X} and \mathbf{Y} , respectively. They are uncorrelated with each other and $\boldsymbol{\zeta}$. Typically, $n \leq q$ and $m \leq p$.

For the estimation of the matrices \mathbf{A} , \mathbf{B} , \mathbf{C} , \mathbf{G} , and the covariance matrices of the errors, there is the LISREL algorithm of K. G. Jöreskog [2] (assuming multivariate Gaussian distribution of the measurement variables and large sample sizes) and component-wise SEM algorithms (not postulating normality and being able to treat small sample sizes), e.g., [7, 8], at our disposal.

In the first stage of his PLS algorithm, H. Wold [8] gives an iteration to find the case values of the latent variables. He states that this fixed point iteration converges. We use only this first stage to calculate the product moment estimate of the covariance matrix of the latent variables. Then we decompose the inverse of this matrix as \mathbf{LDL}^T with \mathbf{L} and \mathbf{D} having the form

$$\mathbf{L} = \begin{pmatrix} \mathbf{B}^T & \mathbf{O} \\ -\mathbf{A}^T & \mathbf{I} \end{pmatrix}, \quad \mathbf{D} = \begin{pmatrix} \mathbf{Q}^{-1} & \mathbf{O} \\ \mathbf{O} & \mathbf{F}^{-1} \end{pmatrix}, \quad (2.1)$$

where \mathbf{B} is $m \times m$ upper triangular matrix with 1s along its main diagonal, and \mathbf{A} is $m \times n$ matrix. The block-diagonal matrix \mathbf{D} comprises the inverse of the error covariance matrix \mathbf{Q} of $\boldsymbol{\zeta}$ and \mathbf{F} of $\boldsymbol{\xi}$, where \mathbf{Q} itself is a diagonal matrix. For this purpose we use the block Cholesky decomposition with block sizes $1, \dots, 1, n$ with number¹ m of 1s.

Wold's algorithm also provides the outer relation matrices \mathbf{C} and \mathbf{G} . In this way, we can estimate the parameter matrices \mathbf{A} , \mathbf{B} , \mathbf{C} , \mathbf{G} based on a training sample and on adjacency matrices that specify which latent variable is related to which observable one, both among the exogenous and endogenous variables. This is the point where the expert can intervene the system. For a detailed description, see Section 4.

In the heart of the estimation and the forthcoming filtering there is the simultaneous usage of OLS (ordinary least squares) regression, tracing back to the Gauss normal equations. We give a short summary of that.

Now we concentrate on linear estimates in Hilbert spaces that are the best whenever the underlying distribution is multivariate Gaussian, but is also applicable to second order processes.

Lemma 2.1. *Let $\mathbf{Y} \in \mathbb{R}^p$ and $\mathbf{X} \in \mathbb{R}^q$ be random vectors on a joint probability space with existing second moments and zero expectation. Then*

$\mathbb{E}\|\mathbf{Y} - \mathbf{A}^T\mathbf{X}\|^2$ is minimized with

$$\mathbf{A} = [\mathbb{E}\mathbf{X}\mathbf{X}^T]^{-1}[\mathbb{E}\mathbf{X}\mathbf{Y}^T], \quad (2.2)$$

¹Number of endogenous latent variables in the model.

where \mathbf{A} is a $q \times p$ matrix and we use generalized inverse $-$ if the covariance matrix $\mathbb{E}\mathbf{X}\mathbf{X}^T$ of \mathbf{X} is singular. If it is positive definite, then we get a unique estimate for \mathbf{A} with the unique inverse matrix $[\mathbb{E}\mathbf{X}\mathbf{X}^T]^{-1}$.

Note that the notation $-$ applies to any (not necessarily unique) generalized inverse, whereas $+$ will be used for the uniquely defined Moore–Penrose generalized inverse, see [6].

Proof. Observe that minimizing

$$\mathbb{E}\|\mathbf{Y} - \mathbf{A}^T\mathbf{X}\|^2 = \sum_{i=1}^p (Y^i - \mathbf{a}_i^T\mathbf{X})^2$$

with respect to $\mathbf{A} = (\mathbf{a}_1 \dots \mathbf{a}_p)$ falls apart into the following p minimization tasks, with respect to the q -dimensional column vectors of \mathbf{A} :

$$\min_{\mathbf{a}_i} (Y^i - \mathbf{a}_i^T\mathbf{X})^2, \quad i = 1, \dots, p.$$

The solution (e.g., with the help of differentiation) gives the well known *Gauss normal equations* from the classical theory of multivariate regression:

$$[\mathbb{E}\mathbf{X}\mathbf{X}^T]\mathbf{a}_i = [\mathbb{E}\mathbf{X}\mathbf{Y}_i], \quad i = 1, \dots, p.$$

Since this system of linear equations is consistent (the vector $\mathbb{E}\mathbf{X}\mathbf{Y}_i$ is in the column space of $\mathbb{E}\mathbf{X}\mathbf{X}^T$), it always has a solution in the general form:

$$\mathbf{a}_i = [\mathbb{E}\mathbf{X}\mathbf{X}^T]^{-}[\mathbb{E}\mathbf{X}\mathbf{Y}_i], \quad i = 1, \dots, p.$$

Therefore the matrix \mathbf{A} giving the optimum is

$$\mathbf{A} = [\mathbb{E}\mathbf{X}\mathbf{X}^T]^{-}[\mathbb{E}\mathbf{X}\mathbf{Y}^T],$$

that is unique only if $\mathbb{E}\mathbf{X}\mathbf{X}^T$ is invertible (positive definite), otherwise (if $\mathbb{E}\mathbf{X}\mathbf{X}^T$ is singular, positive semidefinite) infinitely many versions of the generalized inverse give infinitely many convenient \mathbf{A} s. However, these always provide the same optimal linear prediction (projection) for \mathbf{Y} as follows:

$$\text{Proj}_{H(\mathbf{X})}\mathbf{Y} = \hat{\mathbf{Y}} = \begin{bmatrix} \hat{Y}_1 \\ \hat{Y}_2 \\ \vdots \\ \hat{Y}_p \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1^T\mathbf{X} \\ \mathbf{a}_2^T\mathbf{X} \\ \vdots \\ \mathbf{a}_p^T\mathbf{X} \end{bmatrix} = \mathbf{A}^T\mathbf{X},$$

where $\text{Proj}_{H(\mathbf{X})}$ denotes the projection onto the Hilbert space spanned by the linear combinations of the components of \mathbf{X} (the expectations are zeros and the inner product is the covariance). \square

Lemma 2.2. *Let $\mathbf{Y} \in \mathbb{R}^p$ and $\mathbf{X} \in \mathbb{R}^q$ be random vectors on a joint probability space with existing second moments and zero expectation, and let $\text{Proj}_{H(\mathbf{X})}\mathbf{Y}$ denotes the best linear prediction of \mathbf{Y} based on \mathbf{X} , as before. Then with any $p \times p$ matrix Φ ,*

$$\text{Proj}_{H(\mathbf{X})}(\Phi\mathbf{Y}) = \Phi\text{Proj}_{H(\mathbf{X})}\mathbf{Y}.$$

Proof. We saw that $\text{Proj}_{H(\mathbf{X})}\mathbf{Y} = \mathbf{A}^T\mathbf{X}$, where by (2.2), $\mathbf{A} = [\mathbb{E}\mathbf{X}\mathbf{X}^T]^{-1}[\mathbb{E}\mathbf{X}\mathbf{Y}^T]$, and we use the generalized inverse $^-$ if the covariance matrix $\mathbb{E}\mathbf{X}\mathbf{X}^T$ of \mathbf{X} is singular. Then

$$\begin{aligned} \text{Proj}_{H(\mathbf{X})}(\Phi\mathbf{Y}) &= \{[\mathbb{E}\mathbf{X}\mathbf{X}^T]^{-1}[\mathbb{E}\mathbf{X}(\Phi\mathbf{Y})^T]\}^T\mathbf{X} = [\mathbb{E}(\Phi\mathbf{Y}\mathbf{X}^T)][\mathbb{E}\mathbf{X}\mathbf{X}^T]^{-1}\mathbf{X} \\ &= \Phi[\mathbb{E}(\mathbf{Y}\mathbf{X}^T)][\mathbb{E}\mathbf{X}\mathbf{X}^T]^{-1}\mathbf{X} = \Phi\text{Proj}_{H(\mathbf{X})}\mathbf{Y}. \end{aligned}$$

□

The above lemma shows that this projection is linear in \mathbf{Y} and it commutes with Φ . In the Gaussian case, obviously, we have that

$$\text{Proj}_{H(\mathbf{X})}(\Phi\mathbf{Y}) = \mathbb{E}(\Phi\mathbf{Y} | \mathbf{X}) = \Phi\mathbb{E}(\mathbf{Y} | \mathbf{X}) = \Phi\text{Proj}_{H(\mathbf{X})}(\mathbf{Y})$$

by the properties of the conditional expectation.

The above setup is used for simultaneous (in other words, multiple response) regressions when we regress the components of a random vector (target) with all the components of the predictors.

3. The linear dynamical system for the prediction

Discrete time observations $\mathbf{X}_t, \mathbf{Y}_t$ arrive, whereas ξ_t and η_t are latent state variables corresponding to them. Starting at time 0, for $t = 1, 2, \dots$, the estimate of $\hat{\eta}_t$ is found, while observing $\mathbf{X}_1, \dots, \mathbf{X}_t$. Actually, to find $\hat{\eta}_t$, we only need the estimate $\hat{\xi}_t$ and the last observation \mathbf{X}_t . Then to find $\hat{\xi}_{t+1}$, the preceding estimate $\hat{\eta}_t$ and the last observation \mathbf{Y}_t are needed. In this way, a recursion is given via the propagation of the error covariance matrices. During the calculations, we use the linearity of the state equations and the predictions, for which we confine ourselves to the second moments of the underlying distributions (second order processes).

The linear dynamical system is

$$\begin{aligned} \mathbf{B}\eta_t &= \mathbf{A}\xi_t + \zeta_t \\ \mathbf{U}\xi_{t+1} &= \mathbf{V}\eta_t + \gamma_t \\ \mathbf{X}_t &= \mathbf{C}\xi_t + \varepsilon_t \\ \mathbf{Y}_t &= \mathbf{G}\eta_t + \delta_t, \end{aligned} \tag{3.1}$$

where \mathbf{A} is $m \times n$, \mathbf{B} is $m \times m$, \mathbf{V} is $n \times m$, \mathbf{U} is $n \times n$, \mathbf{C} is $q \times n$, and \mathbf{G} is $p \times m$ specified matrix; \mathbf{B} and \mathbf{U} are non-singular (in recursive models they are upper triangular with 1s along their main diagonals). Further, ζ_t is an orthogonal process with $\mathbb{E}\zeta_t\zeta_s^T = \delta_{st}\mathbf{Q}$ with diagonal covariance matrix \mathbf{Q} ; γ_t is an orthogonal process with $\mathbb{E}\gamma_t\gamma_s^T = \delta_{st}\mathbf{R}$ with diagonal covariance matrix \mathbf{R} ; $\mathbb{E}\xi_s^T\zeta_t = \mathbf{0}$ and $\mathbb{E}\eta_s^T\gamma_t = \mathbf{0}$ for $s \leq t$; ε_t is independent of ξ_t , δ_t is independent of η_t , they are also independent of each other and of ζ_t and γ_t . For simplicity, we assume that all the expectations are zeros.

The $\mathbf{A}, \mathbf{B}, \mathbf{U}, \mathbf{V}$ matrices are estimated from a training sample. Actually, the matrices \mathbf{A} and \mathbf{B} together with \mathbf{Q} and \mathbf{F} come from the block Cholesky decomposition (2.1), based on the product-moments of the estimated latent scores of the pairs ξ_s, η_s , where $s < 0$ is integer from the past (training sample). Likewise, the matrices

\mathbf{U} and \mathbf{V} together with \mathbf{R} and \mathbf{F}^* come from the block Cholesky decomposition (3.2) below, based on the product-moments of the estimated latent scores of the shifted pairs $\boldsymbol{\eta}_s, \boldsymbol{\xi}_{s+1}$ ($s < 0$). The inverse of this matrix is decomposed as $\mathbf{L}^* \mathbf{D}^* \mathbf{L}^{*T}$ with \mathbf{L}^* and \mathbf{D}^* having the form

$$\mathbf{L}^* = \begin{pmatrix} \mathbf{U}^T & \mathbf{O} \\ -\mathbf{V}^T & \mathbf{I} \end{pmatrix}, \quad \mathbf{D}^* = \begin{pmatrix} \mathbf{R}^{-1} & \mathbf{O} \\ \mathbf{O} & \mathbf{F}^{*-1} \end{pmatrix}, \quad (3.2)$$

where recall that, in the recursive model, \mathbf{U} is $n \times n$ upper triangular matrix with 1s along its main diagonal, and \mathbf{V} is $n \times m$. The block-diagonal matrix \mathbf{D}^* comprises the inverse of the error covariance matrix \mathbf{R} of $\boldsymbol{\gamma}$ and \mathbf{F}^* of $\boldsymbol{\eta}$, where \mathbf{R} itself is a diagonal matrix. For this purpose we use the block Cholesky decomposition with block sizes $1, \dots, 1, m$ with number² n of 1s.

The matrices \mathbf{C} and \mathbf{G} are estimated by the PLS algorithm of H. Wold [8]. The details are given in Section 4.

Now a recursion is introduced for the following problem: starting the observations at time 0 in the test sample, we want to estimate $\boldsymbol{\eta}_t$ based on \mathbf{X}_t , and $\boldsymbol{\xi}_{t+1}$ based on \mathbf{Y}_t component-wise, with minimum mean square error. Former observations also play role, but only through the last one and through the propagation of the error covariance matrices. Here \mathbf{X}_0 and \mathbf{Y}_0 can be taken from the training sample.

3.1. First stage: $\mathbf{X}_t \rightarrow \hat{\boldsymbol{\eta}}_t$

For $t \geq 1$, let $H_{t-1}(\mathbf{X}) = \text{Span}(\mathbf{X}_0, \dots, \mathbf{X}_{t-1})$ consists of the linear combinations of all the components of $\mathbf{X}_0, \dots, \mathbf{X}_{t-1}$ over a common probability space. They are also in a Hilbert space (L_2 space) with the covariance as inner product. We denote the optimal prediction of $\boldsymbol{\xi}_t$ based on $\mathbf{X}_0, \dots, \mathbf{X}_{t-1}$ by $\hat{\boldsymbol{\xi}}_t$.

If $\mathbf{X}_0, \dots, \mathbf{X}_{t-1}$ are observed, i.e., $H_{t-1}(\mathbf{X})$ is known, then the newly observed (measured) \mathbf{X}_t can be orthogonally decomposed as

$$\mathbf{X}_t = \text{Proj}_{H_{t-1}(\mathbf{X})} \mathbf{X}_t + \tilde{\mathbf{X}}_t = \bar{\mathbf{X}}_t + \tilde{\mathbf{X}}_t, \quad (3.3)$$

where the orthogonal component $\tilde{\mathbf{X}}_t \in I_t(\mathbf{X})$, and $I_t(\mathbf{X})$ is the so-called innovation subspace (actually, the components of $\tilde{\mathbf{X}}_t$ generate $I_t(\mathbf{X})$). Assume that $I_t(\mathbf{X})$ is not the sole $\mathbf{0}$ vector, otherwise observing $\tilde{\mathbf{X}}_t$ does not give any additional information to $H_{t-1}(\mathbf{X})$. If $\{\mathbf{X}_t\}$ is weakly stationary, it means that the process is *regular*.

Equation (3.3) implies the decomposition of the corresponding subspaces like

$$H_t(\mathbf{X}) = H_{t-1}(\mathbf{X}) \oplus I_t(\mathbf{X}), \quad (3.4)$$

that is the analogue of multidimensional Wold decomposition when we make one-step ahead prediction based on finitely many past values. (The Wold decomposition applies to the stationary and infinite past case. Indeed, when $t \rightarrow \infty$, i.e., going to the future, we approach this situation in the stationary case).

²Number of exogenous latent variables in the model.

Assume that we have already found $\hat{\xi}_t$. We shall give a recursion to find $\hat{\eta}_t$ by using the new value of \mathbf{X}_t . In view of Equation (3.4), we proceed as follows:

$$\begin{aligned} B\hat{\eta}_t &= \text{Proj}_{H_t(\mathbf{X})}(B\eta_t) = \text{Proj}_{H_{t-1}(\mathbf{X})}(B\eta_t) + \text{Proj}_{I_t(\mathbf{X})}(B\eta_t) \\ &= A\text{Proj}_{H_{t-1}(\mathbf{X})}\xi_t + \text{Proj}_{H_{t-1}(\mathbf{X})}\zeta_t + \mathbf{K}_t\tilde{\mathbf{X}}_t \\ &= A\hat{\xi}_t + \mathbf{K}_t\tilde{\mathbf{X}}_t, \end{aligned} \quad (3.5)$$

where we utilized that $\zeta_t \perp H_{t-1}(\mathbf{X})$, Lemma 2.1 and the first state equation of (3.1). We refer to the linearity of the projection, see Lemma 2.2. Since $\text{Proj}_{I_t(\mathbf{X})}B\eta_t$ is the linear combination of the coordinates of the vector $\tilde{\mathbf{X}}_t \in I_t(\mathbf{X})$, its effect can be written as a matrix \mathbf{K}_t multiplied with $\tilde{\mathbf{X}}_t$. This $m \times q$ matrix \mathbf{K}_t is called *Kálmán gain matrix* after R. E. Kálmán (in fact, this notation was first used in the paper [4] of Kálmán and Bucy).

To specify the matrix \mathbf{K}_t , we have to write $\tilde{\mathbf{X}}_t$ in terms of $\hat{\xi}_t$ and \mathbf{X}_t . For this purpose, let us project both sides of the first observation equation of (3.1), i.e., of $\mathbf{X}_t = C\xi_t + \varepsilon_t$, onto $H_{t-1}(\mathbf{X})$. We get that

$$\bar{\mathbf{X}}_t = C\hat{\xi}_t.$$

Taking the orthogonal decomposition (3.3) of \mathbf{X}_t into consideration yields that

$$\tilde{\mathbf{X}}_t = \mathbf{X}_t - \bar{\mathbf{X}}_t = \mathbf{X}_t - C\hat{\xi}_t. \quad (3.6)$$

We substitute this into the last line of Equation (3.5) and obtain that

$$B\hat{\eta}_t = A\hat{\xi}_t + \mathbf{K}_t\tilde{\mathbf{X}}_t = (A - \mathbf{K}_tC)\hat{\xi}_t + \mathbf{K}_t\mathbf{X}_t.$$

With the notation

$$A_t^* = A - \mathbf{K}_tC \quad (3.7)$$

for the updated transition matrix, we get the new linear dynamics:

$$B\hat{\eta}_t = A_t^*\hat{\xi}_t + \mathbf{K}_t\mathbf{X}_t. \quad (3.8)$$

We also have the alternative expression

$$B\hat{\eta}_t = A\hat{\xi}_t + \mathbf{K}_t\tilde{\mathbf{X}}_t = A\hat{\xi}_t + \mathbf{K}_t(\mathbf{X}_t - C\hat{\xi}_t). \quad (3.9)$$

The estimation error is also governed by the linear dynamical system. This error term has two alternative forms. Using (3.8), the one is

$$\begin{aligned} B\tilde{\eta}_t &= B\eta_t - B\hat{\eta}_t = A\xi_t + \zeta_t - A^*\hat{\xi}_t - \mathbf{K}_tC\xi_t - \mathbf{K}_t\varepsilon_t \\ &= A_t^*(\xi_t - \hat{\xi}_t) + \zeta_t - \mathbf{K}_t\varepsilon_t = A_t^*\tilde{\xi}_t + \zeta_t - \mathbf{K}_t\varepsilon_t. \end{aligned}$$

Then, using (3.9), the other is

$$B\tilde{\eta}_t = A\xi_t + \zeta_t - A\hat{\xi}_t - \mathbf{K}_t(\mathbf{X}_t - C\hat{\xi}_t) = A\tilde{\xi}_t + \zeta_t - \mathbf{K}_t(\mathbf{X}_t - C\hat{\xi}_t).$$

From here, we get the following recursion for the covariance matrix

$$P_t = \mathbb{E}\tilde{\xi}_t\tilde{\xi}_t^T \quad (3.10)$$

of the optimal error (of predicting ξ_t) and so, of K_t :

$$\begin{aligned} B[\mathbb{E}\tilde{\eta}_t\tilde{\eta}_t^T]B^T &= \mathbb{E}[B\tilde{\eta}_t][B\tilde{\eta}_t]^T \\ &= \mathbb{E}[A^*\tilde{\xi}_t + \zeta_t][A\tilde{\xi}_t + \zeta_t - K_t(\mathbf{X}_t - C\hat{\xi}_t)]^T \\ &= A_t^*P_tA^T + Q, \end{aligned} \quad (3.11)$$

where recall that $Q = \mathbb{E}\zeta_t\zeta_t^T$, obtainable by (2.1). We used that ζ_t is uncorrelated with ξ_t and, therefore, with $\tilde{\xi}_t$ too. We also used that $\mathbf{X}_t - C\hat{\xi}_t$ is in $I_t(\mathbf{X})$, and ζ_t is uncorrelated with ε_t .

It remains to find an explicit formula for K_t , and thus, also for A_t^* . Recall that K_t is the matrix of the linear operation $\text{Proj}_{I_t(\mathbf{X})}B\eta_t$, therefore by the projection principle (see Lemma 2.1):

$$K_t = [\mathbb{E}B\eta_t\tilde{X}_t^T][\mathbb{E}(\tilde{X}_t\tilde{X}_t^T)]^+,$$

where $^+$ denotes the Moore–Penrose generalized inverse (we use regular inverse if the underlying matrix is invertible).

Now we calculate the matrices in brackets. By the third equation of (3.1), that extends to $\tilde{X}_t = C\tilde{\xi}_t + \varepsilon_t$ and to their predictions, we get that

$$\mathbb{E}\tilde{X}_t\tilde{X}_t^T = \mathbb{E}(C\tilde{\xi}_t + \varepsilon_t)(C\tilde{\xi}_t + \varepsilon_t)^T = CP_tC^T + E,$$

where $E = \mathbb{E}\varepsilon_t\varepsilon_t^T$. E is obtainable by (2.1) in the following way:

$$\mathbb{E}X_tX_t^T = C(\mathbb{E}\xi_t\xi_t^T)C^T + E = CFC^T + E.$$

So E is the difference between $\mathbb{E}X_tX_t^T$ (estimated as $\hat{\Sigma}_{\mathbf{X}\mathbf{X}}$ from the training sample) and CFC^T , where F is the inverse of the second diagonal block of D in (2.1).

By the first and third equation of (3.1) and the orthogonality of $\hat{\xi}_t$ and $\tilde{\xi}_t$ we get that

$$\mathbb{E}(B\eta_t\tilde{X}_t^T) = A\mathbb{E}(\xi_t\tilde{X}_t^T) = A\mathbb{E}[(\hat{\xi}_t + \tilde{\xi}_t)(C\tilde{\xi}_t)^T] = AP_tC^T. \quad (3.12)$$

Therefore,

$$K_t = AP_tC^T[CP_tC^T + E]^+ \quad (3.13)$$

with the Moore–Penrose inverse.

With this matrix K_t of Equation (3.13) and using Equation (3.11), we are able to write the error covariance matrix in the form of a symmetric matrix:

$$\begin{aligned} B[\mathbb{E}\tilde{\eta}_t\tilde{\eta}_t^T]B^T &= A^*P_tA^T + Q = (A - K_tC)P_tA^T + Q \\ &= (A - AP_tC^T[CP_tC^T + E]^+C)P_tA^T + Q \\ &= A(I - P_tC^T[CP_tC^T + E]^+C)P_tA^T + Q \\ &= AP_tA^T - AP_tC^T[CP_tC^T + E]^+CP_tA^T + Q, \end{aligned}$$

so

$$BP_t^*B^T = AP_tA^T - AP_tC^T[CP_tC^T + E]^+CP_tA^T + Q, \quad (3.14)$$

where $P_t^* = \mathbb{E}(\tilde{\eta}_t\tilde{\eta}_t^T)$ is the covariance matrix of the error when predicting η_t . In the next stage, we use it to find P_{t+1} .

3.2. Second stage: $\mathbf{Y}_t \rightarrow \hat{\boldsymbol{\xi}}_{t+1}$

For $t \geq 1$, let $H_{t-1}(\mathbf{Y}) = \text{Span}(\mathbf{Y}_0, \dots, \mathbf{Y}_{t-1})$ consists of the linear combinations of all the components of $\mathbf{Y}_0, \dots, \mathbf{Y}_{t-1}$ over a common probability space. We denote the optimal prediction of $\boldsymbol{\eta}_t$ based on $\mathbf{Y}_0, \dots, \mathbf{Y}_{t-1}$ by $\check{\boldsymbol{\eta}}_t$.

If $\mathbf{Y}_0, \dots, \mathbf{Y}_{t-1}$ are observed, i.e., $H_{t-1}(\mathbf{Y})$ is known, then the newly observed (measured) \mathbf{Y}_t can be orthogonally decomposed as

$$\mathbf{Y}_t = \text{Proj}_{H_{t-1}(\mathbf{Y})} \mathbf{Y}_t + \tilde{\mathbf{Y}}_t = \bar{\mathbf{Y}}_t + \tilde{\mathbf{Y}}_t, \quad (3.15)$$

where the orthogonal component $\tilde{\mathbf{Y}}_t \in I_t(\mathbf{Y})$, and $I_t(\mathbf{Y})$ is the innovation subspace (actually, the components of $\tilde{\mathbf{Y}}_t$ generate $I_t(\mathbf{Y})$). Assume that $I_t(\mathbf{Y})$ is not the sole $\mathbf{0}$ vector, otherwise observing \mathbf{Y}_t does not give any additional information to $H_{t-1}(\mathbf{Y})$.

Equation (3.15) implies the decomposition of the corresponding subspaces like

$$H_t(\mathbf{Y}) = H_{t-1}(\mathbf{Y}) \oplus I_t(\mathbf{Y}). \quad (3.16)$$

Assume that we have already found $\check{\boldsymbol{\eta}}_t$. We shall give a recursion to find $\check{\boldsymbol{\xi}}_{t+1}$ by using the new value of \mathbf{Y}_t . In view of Equation (3.16):

$$\begin{aligned} U\check{\boldsymbol{\xi}}_{t+1} &= \text{Proj}_{H_t(\mathbf{Y})}(U\boldsymbol{\xi}_{t+1}) = \text{Proj}_{H_{t-1}(\mathbf{Y})}(U\boldsymbol{\xi}_{t+1}) + \text{Proj}_{I_t(\mathbf{Y})}(U\boldsymbol{\xi}_{t+1}) \\ &= V\text{Proj}_{H_{t-1}(\mathbf{Y})}\boldsymbol{\eta}_t + \text{Proj}_{H_{t-1}(\mathbf{Y})}\boldsymbol{\gamma}_t + \mathbf{M}_t\tilde{\mathbf{Y}}_t \\ &= V\check{\boldsymbol{\eta}}_t + \mathbf{M}_t\tilde{\mathbf{Y}}_t, \end{aligned} \quad (3.17)$$

where we utilized that $\boldsymbol{\gamma}_t \perp H_{t-1}(\mathbf{Y})$, Lemma 2.1 and the second state equation of (3.1). Furthermore, we refer to the linearity of the projection, see Lemma 2.2. Since $\text{Proj}_{I_t(\mathbf{Y})}U\boldsymbol{\xi}_{t+1}$ is the linear combination of the coordinates of the vector $\tilde{\mathbf{Y}}_t \in I_t(\mathbf{Y})$, its effect can be written as a matrix \mathbf{M}_t multiplied with $\tilde{\mathbf{Y}}_t$. This $n \times p$ matrix \mathbf{M}_t is another gain matrix.

To specify the matrix \mathbf{M}_t , we have to write $\tilde{\mathbf{Y}}_t$ in terms of $\check{\boldsymbol{\eta}}_t$ and \mathbf{Y}_t . For this purpose, let us project both sides of the second observation equation of (3.1), i.e., of $\mathbf{Y}_t = \mathbf{G}\boldsymbol{\eta}_t + \boldsymbol{\delta}_t$, onto $H_{t-1}(\mathbf{Y})$. We get that

$$\bar{\mathbf{Y}}_t = \mathbf{G}\check{\boldsymbol{\eta}}_t.$$

Taking the orthogonal decomposition (3.15) of \mathbf{Y}_t into consideration yields that

$$\tilde{\mathbf{Y}}_t = \mathbf{Y}_t - \bar{\mathbf{Y}}_t = \mathbf{Y}_t - \mathbf{G}\check{\boldsymbol{\eta}}_t. \quad (3.18)$$

We substitute this into the last line of Equation (3.17) and obtain that

$$U\check{\boldsymbol{\xi}}_{t+1} = V\check{\boldsymbol{\eta}}_t + \mathbf{M}_t\tilde{\mathbf{Y}}_t = (V - \mathbf{M}_t\mathbf{G})\check{\boldsymbol{\eta}}_t + \mathbf{M}_t\mathbf{Y}_t.$$

With the notation

$$\mathbf{V}_t^* = V - \mathbf{M}_t\mathbf{G} \quad (3.19)$$

for the updated transition matrix, we get the new linear dynamics:

$$U\check{\boldsymbol{\xi}}_{t+1} = \mathbf{V}_t^*\check{\boldsymbol{\eta}}_t + \mathbf{M}_t\mathbf{Y}_t. \quad (3.20)$$

We also have the alternative expression

$$U\check{\boldsymbol{\xi}}_{t+1} = V\check{\boldsymbol{\eta}}_t + \mathbf{M}_t\tilde{\mathbf{Y}}_t = V\check{\boldsymbol{\eta}}_t + \mathbf{M}_t(\mathbf{Y}_t - \mathbf{G}\check{\boldsymbol{\eta}}_t). \quad (3.21)$$

The estimation error is also governed by the linear dynamical system. This error term has two alternative forms. Using (3.20), the one is

$$\begin{aligned} U\check{\xi}_{t+1} &= U\xi_{t+1} - U\check{\xi}_{t+1} = V\eta_t + \gamma_t - V_t^*\check{\eta}_t - M_t G\eta_t - M_t\delta_t \\ &= V_t^*(\eta_t - \check{\eta}_t) + \gamma_t - M_t\delta_t = V_t^*\check{\eta}_t + \gamma_t - M_t\delta_t. \end{aligned}$$

Then, using (3.21), the other is

$$U\check{\xi}_{t+1} = V\eta_t + \gamma_t - V\check{\eta}_t - M_t(\mathbf{Y}_t - G\check{\eta}_t) = V\check{\eta}_t + \gamma_t - M_t(\mathbf{Y}_t - M_t\check{\eta}_t).$$

From here, we get the following recursion for the covariance matrix

$$P_t^* = \mathbb{E}\check{\eta}_t\check{\eta}_t^T \quad (3.22)$$

of the optimal error (of predicting η_t) and so, of M_t :

$$\begin{aligned} U[\mathbb{E}\check{\xi}_{t+1}\check{\xi}_{t+1}^T]U^T &= \mathbb{E}[U\check{\xi}_{t+1}][U\check{\xi}_{t+1}]^T \\ &= \mathbb{E}[V_t^*\check{\eta}_t + \gamma_t - M_t\delta_t][V_t^*\check{\eta}_t + \gamma_t - M_t(\mathbf{Y}_t - G\check{\eta}_t)]^T \\ &= V_t^*P_t^*V^T + R, \end{aligned} \quad (3.23)$$

where recall that $R = \mathbb{E}\gamma_t\gamma_t^T$, obtainable by (3.2), and we used that γ_t is uncorrelated with η_t and, therefore, with $\check{\eta}_t$ too; further, $V_t^* = V - M_t G$. We also used that $\mathbf{Y}_t - G\check{\eta}_t$ is in $I_t(\mathbf{Y})$, and that γ_t is uncorrelated with δ_t .

Now an explicit formula is found for M_t , and thus, also for V_t^* . Recall that M_t is the matrix of the linear operation $\text{Proj}_{I_t(\mathbf{Y})}U\xi_{t+1}$, therefore by the projection principle (see Lemma 2.1):

$$M_t = [\mathbb{E}(U\xi_{t+1}\check{\mathbf{Y}}_t^T)][\mathbb{E}(\check{\mathbf{Y}}_t\check{\mathbf{Y}}_t^T)]^+,$$

where $^+$ denotes the Moore–Penrose generalized inverse (we use regular inverse if the underlying matrix is invertible). We calculate the matrices in brackets. By the last equation of (3.1), that extends to $\check{\mathbf{Y}}_t = G\check{\eta}_t + \delta_t$ and to their predictions, we get that

$$\mathbb{E}(\check{\mathbf{Y}}_t\check{\mathbf{Y}}_t^T) = \mathbb{E}[(G\check{\eta}_t + \delta_t)(G\check{\eta}_t + \delta_t)^T] = GP_t^*G^T + \Delta,$$

where $\Delta = \mathbb{E}\delta_t\delta_t^T$. Δ is obtainable by (3.2) in the following way:

$$\mathbb{E}\mathbf{Y}_t\mathbf{Y}_t^T = G(\mathbb{E}\eta_t\eta_t^T)G^T + \Delta = GF^*G^T + \Delta.$$

So Δ is the difference between $\mathbb{E}\mathbf{Y}_t\mathbf{Y}_t^T$ (estimated as $\hat{\Sigma}_{\mathbf{Y}\mathbf{Y}}$ from the training sample) and GF^*G^T , where F^* is the inverse of the second diagonal block of D^* in (3.2).

By the second and fourth equation of (3.1) and the orthogonality of $\check{\eta}_t$ and $\check{\eta}_t$ we get that

$$\mathbb{E}(U\xi_{t+1}\check{\mathbf{Y}}_t^T) = V\mathbb{E}(\eta_t\check{\mathbf{Y}}_t^T) = V\mathbb{E}[(\check{\eta}_t + \check{\eta}_t)(G\check{\eta}_t)^T] = VP_t^*G^T. \quad (3.24)$$

Therefore,

$$M_t = VP_t^*G^T[GP_t^*G^T + \Delta]^+ \quad (3.25)$$

with the Moore–Penrose inverse.

With this matrix \mathbf{M}_t of Equation (3.25) and using Equation (3.23), we are able to write the error covariance matrix in the form of a symmetric matrix:

$$\begin{aligned} \mathbf{U}[\mathbb{E}\check{\check{\xi}}_{t+1}\check{\check{\xi}}_{t+1}^T]\mathbf{U}^T &= \mathbf{V}_t^* \mathbf{P}_t^* \mathbf{V}_t^T + \mathbf{\Delta} = (\mathbf{V} - \mathbf{M}_t \mathbf{G}) \mathbf{P}_t^* \mathbf{V}_t^T + \mathbf{R} \\ &= (\mathbf{V} - \mathbf{V} \mathbf{P}_t^* \mathbf{G}^T [\mathbf{G} \mathbf{P}_t^* \mathbf{G}^T + \mathbf{\Delta}]^+ \mathbf{G} \mathbf{P}_t^* \mathbf{V}_t^T + \mathbf{R} \\ &= \mathbf{V} (\mathbf{I} - \mathbf{P}_t^* \mathbf{G}^T) [\mathbf{G} \mathbf{P}_t^* \mathbf{G}^T + \mathbf{\Delta}]^+ \mathbf{G} \mathbf{P}_t^* \mathbf{V}_t^T + \mathbf{R} \\ &= \mathbf{V} \mathbf{P}_t^* \mathbf{V}_t^T - \mathbf{V} \mathbf{P}_t^* \mathbf{G}^T [\mathbf{G} \mathbf{P}_t^* \mathbf{G}^T + \mathbf{\Delta}]^+ \mathbf{G} \mathbf{P}_t^* \mathbf{V}_t^T + \mathbf{R}, \end{aligned}$$

so

$$\mathbf{U} \mathbf{P}_{t+1} \mathbf{U}^T = \mathbf{V} \mathbf{P}_t^* \mathbf{V}_t^T - \mathbf{V} \mathbf{P}_t^* \mathbf{G}^T [\mathbf{G} \mathbf{P}_t^* \mathbf{G}^T + \mathbf{\Delta}]^+ \mathbf{G} \mathbf{P}_t^* \mathbf{V}_t^T + \mathbf{R}, \quad (3.26)$$

where we assumed that the error covariance matrix of $\check{\xi}_t$ and $\check{\xi}_t$, akin to that of $\check{\eta}_t$ and $\check{\eta}_t$ is the same. This fact gives rise to a recursion by connecting (3.14) and (3.26).

Finally, with (3.9) and (3.21) we are able to recursively estimate the latent state variables. During the $\mathbf{P}_1 \rightarrow \mathbf{P}_1^* \rightarrow \mathbf{P}_2 \rightarrow \mathbf{P}_2^* \dots$ recursion, from \mathbf{P}_t , we find \mathbf{K}_t by (3.13) and $\hat{\eta}_t$ by (3.9). Then, from \mathbf{P}_t^* , we find \mathbf{M}_t by (3.25) and $\check{\xi}_{t+1}$ by (3.21).

As for the relation between $\check{\xi}_t$ and $\check{\xi}_t$, akin to that between $\hat{\eta}_t$ and $\check{\eta}_t$, we can estimate their cross-covariance matrices from the training sample, and then, linearly predict $\check{\eta}_t$ with $\hat{\eta}_t$ and linearly predict $\check{\xi}_t$ with $\hat{\xi}_t$ by Lemma 2.1 as follows:

$$\check{\eta}_t = \hat{\Sigma}_{\eta\mathbf{Y}} \hat{\Sigma}_{\mathbf{Y}\mathbf{Y}}^+ \hat{\Sigma}_{\mathbf{Y}\mathbf{X}} \hat{\Sigma}_{\mathbf{X}\mathbf{X}}^+ \Sigma_{\mathbf{X}\eta} [\mathbb{E}\hat{\eta}_t \hat{\eta}_t^T]^+ \hat{\eta}_t$$

and

$$\hat{\xi}_{t+1} = \hat{\Sigma}_{\xi\mathbf{X}} \hat{\Sigma}_{\mathbf{X}\mathbf{X}}^+ \hat{\Sigma}_{\mathbf{X}\mathbf{Y}} \hat{\Sigma}_{\mathbf{Y}\mathbf{Y}}^+ \hat{\Sigma}_{\mathbf{Y}\xi} [\mathbb{E}\check{\xi}_{t+1} \check{\xi}_{t+1}^T]^+ \check{\xi}_{t+1}.$$

Here, from Equation (3.11), we conclude that

$$\mathbb{E}\hat{\eta}_t \hat{\eta}_t^T = \hat{\Sigma}_{\eta\eta} - \mathbf{B}^{-1} [(\mathbf{A} - \mathbf{K}_t \mathbf{C}) \mathbf{P}_t \mathbf{A}^T + \mathbf{Q}] (\mathbf{B}^{-1})^T.$$

Likewise, from Equation (3.23), we conclude that

$$\mathbb{E}\check{\xi}_t \check{\xi}_t^T = \hat{\Sigma}_{\xi\xi} - \mathbf{U}^{-1} [(\mathbf{V} - \mathbf{M}_t \mathbf{G}) \mathbf{P}_t^* \mathbf{V}_t^T + \mathbf{R}] (\mathbf{U}^{-1})^T.$$

3.3. The main result

We assume that the system was at rest until time 0. The parameter matrices are estimated from the past, whereas newer and newer estimates for the latent variables are given, as observations arrive at time t ($t = 1, 2, \dots$ up to the end of the experimental time T). Thus, we can summarize the results in the subsequent theorem. It is important that in the derivation of the formulas we used the best linear prediction theory of Hilbert spaces.

Theorem 3.1. *In the linear dynamical system (3.1), the optimal estimate $\hat{\eta}_t$ of η_t and $\check{\xi}_{t+1}$ of ξ_{t+1} given $\mathbf{X}_1, \dots, \mathbf{X}_t$ and $\mathbf{Y}_1, \dots, \mathbf{Y}_t$ is generated by the new linear dynamical system*

$$\mathbf{B} \hat{\eta}_t = \mathbf{A}_t^* \hat{\xi}_t + \mathbf{K}_t \mathbf{X}_t$$

and

$$\mathbf{U} \check{\xi}_{t+1} = \mathbf{V}_t^* \check{\eta}_t + \mathbf{M}_t \mathbf{Y}_t.$$

The expected quadratic losses are $\text{tr} \mathbf{P}_t^*$ and $\text{tr} \mathbf{P}_{t+1}$, where \mathbf{P}_t^* and \mathbf{P}_{t+1} are the propagated covariance matrices of the estimation errors. The minimizing matrices and the

one-step ahead predictions $\hat{\eta}_t$ and $\hat{\xi}_{t+1}$ together with the error covariance matrices \mathbf{P}_t^* and \mathbf{P}_{t+1} are uniquely determined by the initial conditions

$$\hat{\xi}_1 = \text{Proj}_{\mathbf{X}_0} \xi_1, \quad \tilde{\xi}_1 = \xi_1 - \hat{\xi}_1, \quad \mathbf{P}_1 = \mathbb{E} \tilde{\xi}_1 \tilde{\xi}_1^T$$

and the recursions for $t = 1, 2, \dots$ as follows.

$$\mathbf{K}_t = \mathbf{A} \mathbf{P}_t \mathbf{C}^T [\mathbf{C} \mathbf{P}_t \mathbf{C}^T + \mathbf{E}]^+$$

$$\hat{\eta}_t = \mathbf{B}^{-1} [\mathbf{A} \hat{\xi}_t + \mathbf{K}_t (\mathbf{X}_t - \mathbf{C} \hat{\xi}_t)]$$

$$\hat{\mathbf{Y}}_t = \mathbf{G} \hat{\eta}_t$$

$$\tilde{\eta}_t = \hat{\Sigma}_{\eta\mathbf{X}} \hat{\Sigma}_{\mathbf{X}\mathbf{X}}^+ \hat{\Sigma}_{\mathbf{X}\mathbf{Y}} \hat{\Sigma}_{\mathbf{Y}\mathbf{Y}}^+ \hat{\Sigma}_{\mathbf{Y}\eta} \left[\hat{\Sigma}_{\eta\eta} - \mathbf{B}^{-1} ((\mathbf{A} - \mathbf{K}_t \mathbf{C}) \mathbf{P}_t \mathbf{A}^T + \mathbf{Q}) (\mathbf{B}^{-1})^T \right]^+ \hat{\eta}_t$$

$$\mathbf{P}_t^* = \mathbf{B}^{-1} [\mathbf{A} \mathbf{P}_t \mathbf{A}^T - \mathbf{A} \mathbf{P}_t \mathbf{C}^T [\mathbf{C} \mathbf{P}_t \mathbf{C}^T + \mathbf{E}]^+ \mathbf{C} \mathbf{P}_t \mathbf{A}^T + \mathbf{Q}] \mathbf{B}^{-1T}$$

$$\mathbf{M}_t = \mathbf{V} \mathbf{P}_t^* \mathbf{G}^T [\mathbf{G} \mathbf{P}_t^* \mathbf{G}^T + \Delta]^+$$

$$\tilde{\xi}_{t+1} = \mathbf{U}^{-1} [\mathbf{V} \tilde{\eta}_t + \mathbf{M}_t (\mathbf{Y}_t - \mathbf{G} \tilde{\eta}_t)]$$

$$\hat{\xi}_{t+1} = \hat{\Sigma}_{\xi\mathbf{Y}} \hat{\Sigma}_{\mathbf{Y}\mathbf{Y}}^+ \hat{\Sigma}_{\mathbf{Y}\mathbf{X}} \hat{\Sigma}_{\mathbf{X}\mathbf{X}}^+ \hat{\Sigma}_{\mathbf{X}\xi} \left[\hat{\Sigma}_{\xi\xi} - \mathbf{U}^{-1} ((\mathbf{V} - \mathbf{M}_t \mathbf{G}) \mathbf{P}_t^* \mathbf{V}^T + \mathbf{R}) (\mathbf{U}^{-1})^T \right]^+ \tilde{\xi}_{t+1}$$

$$\hat{\mathbf{X}}_{t+1} = \mathbf{C} \hat{\xi}_{t+1}$$

$$\mathbf{P}_{t+1} = \mathbf{U}^{-1} [\mathbf{V} \mathbf{P}_t^* \mathbf{V}^T - \mathbf{V} \mathbf{P}_t^* \mathbf{G}^T [\mathbf{G} \mathbf{P}_t^* \mathbf{G}^T + \Delta]^+ \mathbf{G} \mathbf{P}_t^* \mathbf{V}^T + \mathbf{R}] \mathbf{U}^{-1T},$$

where $^+$ denotes the Moore–Penrose generalized inverse (usual inverse if the matrix is invertible).

Note that $\hat{\xi}_1 = \text{Proj}_{\mathbf{X}_0} \xi_1 = \hat{\Sigma}_{\xi\mathbf{X}} \hat{\Sigma}_{\mathbf{X}\mathbf{X}}^+ \mathbf{X}_0$, by Lemma 2.1, where the last training sample entry can be chosen for \mathbf{X}_0 . To initialize \mathbf{P}_1 , the whole training sample can be used: if the L learning sample entries are indexed by ℓ , then $\hat{\xi}_\ell = \hat{\Sigma}_{\xi\mathbf{X}} \hat{\Sigma}_{\mathbf{X}\mathbf{X}}^+ \mathbf{X}_\ell$ and $\tilde{\xi}_\ell = \xi_\ell - \hat{\xi}_\ell$, where ξ_ℓ is the ℓ th case estimate of ξ , based on the forthcoming PLS algorithm in Section 4. Finally, the product moment estimate of \mathbf{P}_1 is $\frac{1}{L} \sum_{\ell=1}^L \tilde{\xi}_\ell \tilde{\xi}_\ell^T$ if the variables have zero expectation (otherwise, the sample means should be subtracted).

Note that the matrices \mathbf{U} , \mathbf{V} are also estimated from the learning sample with the shifted product moments, as discussed in the next section.

4. Application

Using data from three Egyptian villages, we applied our proposed algorithm to examine and predict to what extent parental views affect their daughters' thinking on two empowerment issues. Figure 1 visualizes the hypothesized outer and inner relations. The inner model examines the cause-effect structure between the latent variables (LVs): parental views on girls' participation in decision making and girls' mobility (exogenous: P-DM= ξ^1 and P-Mob= ξ^2) and the daughters' views on the same issues (endogenous: G-DM= η^1 and G-Mob= η^2), respectively. The outer model links LVs and observed variables (OVs) together. Mode A is used to construct all LVs in the model. To clarify, ξ^1 is composed of four independent x s (OVs) that measure parental views on girl's responsibility in making decisions related to marriage,

choosing a husband, entering and continuing schooling. ξ^2 is linked to three x s that ask parents whether girls can go alone to places such as market, field, and friends' home. In the same manner, η^1 and η^2 are connected to the same number of OVs, but the dependent y s that reflect the daughters' views on the same indicators. Note, all the OVs are on the same ordinal scale.

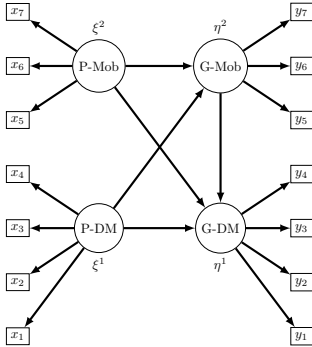


FIGURE 1. Schematic representation of the SEM that considers the effect of parental views on girls' views related to empowerment issues.

A total sample of 349 parents and their daughters are considered. Prior to the analysis, the data were standardized to have zero mean and unit variance. Then, it was divided randomly into a training sample of size 279 and a test sample of the remaining cases. We apply our proposed estimation algorithm on the training sample to obtain the specified parameter matrices ($\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{G}, \mathbf{Q}, \mathbf{F}, \mathbf{U}, \mathbf{V}, \mathbf{R}$, and \mathbf{F}^*). These matrices are used in the derivation of the proposed prediction recursion algorithm. Specifically, the filtering technique in Theorem 3.1 predicts the future values for the test sample observations that come sequentially.

Recall that our integrated estimation algorithm to obtain the model specified matrices combines the first stage of Wold's PLS technique and the block Cholesky decomposition of Kiiveri et al. The detailed explanation follows:

1. The first stage uses stage I of Wold's PLS algorithm, in which the outer relations and the LV case values are obtained from the training sample. It is an iterative process that consists of the following steps:
 - i. Initialize the LV scores for each case as the weighted sum of the observed indicators in the block that correspond to each LV:

$$\mathbf{H} = \mathbf{N}\mathbf{Z},$$

where \mathbf{H} is the exogenous and endogenous LV scores matrix, \mathbf{N} is the training sample data matrix of size 279×14 , and \mathbf{Z} is the 14×4 adjacency matrix of the measurement model. The entries z_{kj} are ones, if the indicator n_{kj} belongs to the block that defines the corresponding LV; and zeros, otherwise. After each step, the LV scores are standardized.

- ii. Update the obtained matrix \mathbf{H} with the inner weights

$$\tilde{\mathbf{H}} = \mathbf{H}\mathbf{W},$$

where \mathbf{W} is the LV inner weights matrix which is computed for each LV to indicate how strong it is connected to the other LVs in the model. There

are three schemes to obtain these weights, for more details, see [7] and [8]. We used the centroid scheme.

- iii. Use the obtained LV scores $\tilde{\mathbf{H}}$ to estimate the outer relations (loadings/weights). There are two modes of constructing the measurement model:
 - Mode A (reflective), where the arrows point outward from the LV to the OVs, as in our case, see Figure 1. The outer scores are called loadings (λ_{kj}) and estimated by OLS simple linear regression between each observed indicator of the measurement block and the corresponding LV score \tilde{h}_j .
 - Mode B (formative), where the arrows point inward from the observed variables to the corresponding LVs. The outer estimated scores called weights (w_{kj}) and are calculated by OLS multiple linear regression in which each LV score \tilde{h}_j is regressed on all the observed indicators of the corresponding block.

All the outer estimates λ_{kjs} and w_{kjs} are collected in the updated weight matrix $\hat{\mathbf{W}}$.

- iv. Using the obtained weight matrix $\hat{\mathbf{W}}$ to update the LV scores

$$\mathbf{H} = \mathbf{N}\hat{\mathbf{W}},$$

where \mathbf{H} contains the LV scores of the last iteration process.

This stage iterates sequentially from Step i. to Step iv. until convergence. At the convergence, the final LV case values \mathbf{H} are obtained as well as the estimated outer matrices \mathbf{C} and \mathbf{G}

$$\mathbf{C} = \begin{matrix} & \xi^1 & \xi^2 \\ \begin{matrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \end{matrix} & \begin{pmatrix} .6759 & 0 \\ .7477 & 0 \\ .7739 & 0 \\ .7093 & 0 \\ 0 & .5326 \\ 0 & .8217 \\ 0 & .7356 \end{pmatrix} \end{matrix}, \quad \mathbf{G} = \begin{matrix} & \eta^1 & \eta^2 \\ \begin{matrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \\ y_7 \end{matrix} & \begin{pmatrix} .6881 & 0 \\ .7408 & 0 \\ .7574 & 0 \\ .7287 & 0 \\ 0 & .4929 \\ 0 & .8132 \\ 0 & .7206 \end{pmatrix} \end{matrix}.$$

The matrices \mathbf{C} and \mathbf{G} contain the loadings that link the latent vectors ξ and η with \mathbf{X} and \mathbf{Y} , respectively.

2. The second stage runs the block Cholesky decomposition of Kiiveri et al. two times on the obtained LV case values. The first decomposition is applied on the inverse of the product moment of the covariance matrix of the LV final scores $\Sigma_{\mathbf{H}}^{-1}$ that is obtained at the convergence of the Wold algorithm, see Equation (2.1). The resulting block matrix \mathbf{L} gives the estimated path coefficient matrices of the inner relations,

$$\mathbf{B} = \begin{matrix} & \eta^1 & \eta^2 \\ \begin{matrix} \eta^1 \\ \eta^2 \end{matrix} & \begin{pmatrix} 1 & .107 \\ 0 & 1 \end{pmatrix} \end{matrix}, \quad \mathbf{A} = \begin{matrix} & \xi^1 & \xi^2 \\ \begin{matrix} \eta^1 \\ \eta^2 \end{matrix} & \begin{pmatrix} .337 & .156 \\ -.068 & .577 \end{pmatrix} \end{matrix};$$

whereas, the block matrix \mathbf{D} yields the error covariance matrices

$$\mathbf{Q} = cov(\zeta) = \begin{pmatrix} .869 & 0 \\ 0 & .663 \end{pmatrix}, \quad \mathbf{F} = cov(\xi) = \begin{pmatrix} 1 & -.024 \\ -.024 & 1 \end{pmatrix}.$$

The second decomposition is performed on the inverse of the product moments of the shifted LV score pairs $\Sigma_{H_{s,s+1}}^{-1}$, see Equation (3.2). The resulting block matrix \mathbf{L}^* contains

$$\mathbf{U} = \begin{matrix} & \xi^1 & \xi^2 \\ \xi^1 & \begin{pmatrix} 1 & .014 \\ 0 & 1 \end{pmatrix} \end{matrix}, \quad \mathbf{V} = \begin{matrix} & \eta^1 & \eta^2 \\ \xi^1 & \begin{pmatrix} .037 & -.074 \\ .051 & .145 \end{pmatrix} \\ \xi^2 & \end{matrix};$$

while, the matrix \mathbf{D}^* gives the error covariance matrices

$$\mathbf{R} = cov(\gamma) = \begin{pmatrix} .992 & 0 \\ 0 & .984 \end{pmatrix}, \quad \mathbf{F}^* = cov(\eta) = \begin{pmatrix} 1 & -.041 \\ -.041 & 1 \end{pmatrix}.$$

At this point, the specified parameter matrices are obtained from the training sample. The estimated matrices \mathbf{C} and \mathbf{G} show the loadings of each OV on the corresponding LV. The matrix \mathbf{B} displays the extent to which girls' views on mobility affect her views on participating in making decisions; while \mathbf{A} shows how parental views on girls' mobility and decision making influence their daughters' opinions on these issues. There is a direct effect of parental views on their daughters' opinions in the same domain, i.e., parents who reported a conservative view on girls' participation in making decisions tend to lead their daughters to think alike. The same scenario is true for mobility, where daughters tend to reproduce their parents' views. On the contrary, the effect is small when we consider parental opinions of one domain on their daughters' views of the other domain.

As new cases come one by one at a time sequence ($t = 1, 2, \dots, T$), instead of re-running the estimation algorithm, we give a recursion to predict the LV case values for the new observation. To do so, the estimated parameter matrices based on the training sample and the Kálmán filtering technique will be used. Theorem 3.1 discusses the recursion from which the optimal prediction of the latent case values is obtained along with the covariance matrix of the prediction error. Specifically, the prediction of $\hat{\eta}_t$ utilizes the estimated $\hat{\xi}_t$ and the new observation \mathbf{X}_t , while the new \mathbf{Y}_t and the obtained $\hat{\eta}_t$ are necessary to find $\hat{\xi}_{t+1}$. To start the recursion, the first propagated matrix \mathbf{P}_1 ought to be initialized from the training sample. Then, the Kálmán gain matrices \mathbf{K}_t and \mathbf{M}_t are obtained. The succession of calculations follows the order:

$$\mathbf{P}_t \rightarrow \mathbf{K}_t \rightarrow \hat{\eta}_t \rightarrow \check{\eta}_t \rightarrow \mathbf{P}_t^* \rightarrow \mathbf{M}_t \rightarrow \check{\xi}_{t+1} \rightarrow \hat{\xi}_{t+1} \rightarrow \hat{\mathbf{X}}_{t+1} \rightarrow \mathbf{P}_{t+1}.$$

For $t = 1$, we show the results of the highlighted matrices of the recursion as they are derived in Section 3.1 and Section 3.2.

First stage: $\mathbf{X}_1 \rightarrow \hat{\eta}_1$

$$\mathbf{P}_1 = \begin{matrix} & \xi^1 & \xi^2 \\ \begin{pmatrix} .007 & .005 \\ .005 & .003 \end{pmatrix} \end{matrix},$$

$$\mathbf{K}_1 = \begin{pmatrix} x_1 & x_2 & x_3 & x_4 & x_5 & x_6 & x_7 \\ -.018 & -.036 & -.036 & -.022 & .022 & .064 & .051 \\ -.013 & -.025 & -.025 & -.015 & .015 & .045 & .036 \end{pmatrix},$$

$$\hat{\boldsymbol{\eta}}_1 = \begin{pmatrix} \eta^1 & \eta^2 \\ .601 & -.377 \end{pmatrix}, \quad \mathbf{P}_1^* = \begin{pmatrix} \eta^1 & \eta^2 \\ .878 & -.070 \\ -.070 & .664 \end{pmatrix}.$$

where $\hat{\boldsymbol{\eta}}_1$ show the predicted case values based on the information \mathbf{X}_1 of the new observation; and \mathbf{P}_1^* is the covariance matrix of the prediction error of $\hat{\boldsymbol{\eta}}_1$.

Second stage: $\mathbf{Y}_1 \rightarrow \hat{\boldsymbol{\xi}}_2$

$$\mathbf{M}_1 = \begin{pmatrix} y_1 & y_2 & y_3 & y_4 & y_5 & y_6 & y_7 \\ .005 & .016 & .019 & .009 & -.022 & -.042 & -.034 \\ .032 & .011 & .007 & .018 & .041 & .086 & .064 \end{pmatrix},$$

$$\hat{\boldsymbol{\xi}}_2 = \begin{pmatrix} \xi^1 & \xi^2 \\ -.595 & -1.460 \end{pmatrix}, \quad \mathbf{P}_2 = \begin{pmatrix} \xi^1 & \xi^2 \\ .993 & -.014 \\ -.014 & .985 \end{pmatrix}.$$

where $\hat{\boldsymbol{\xi}}_2$ presents the estimated case values for the exogenous LVs at $t = 2$ based on the new information \mathbf{Y}_1 . Then the covariance matrix of the prediction error $\hat{\boldsymbol{\xi}}_2$ is obtained. From this we can calculate the propagation matrix \mathbf{P}_2 to start the recursion once again at $t = 2$ for the next new observation.

The root mean square error (*RMSE*) statistic measures the prediction errors. For the test sample of size 70 observations, we compared the predicted LV case values that are obtained from the Wold algorithm and the filtering technique, simultaneously. Table 1 shows the values of *RSME*. It indicates that the prediction capability of the filtering technique and that of the Wold algorithm are quite homogeneous. Moreover, the difference (in Frobenius norm) between the error covariance and gain matrices in the t th and $(t + 1)$ th consecutive steps of the recursion are displayed in Table 2. This shows that these matrices are stabilized after the first few steps.

TABLE 1. RMSE for the predictions of the test sample.

Test Obs. at Sequence “ t ”	Wold Prediction (W)				Filtering Prediction (F)			
	$\hat{\xi}_t^1$	$\hat{\xi}_t^2$	$\hat{\eta}_t^1$	$\hat{\eta}_t^2$	$\hat{\xi}_t^1$	$\hat{\xi}_t^2$	$\hat{\eta}_t^1$	$\hat{\eta}_t^2$
1	-.44005	-.28108	-.25305	-.46325	-.33149	-.52816	-.32865	-.14788
2	.68377	.28108	.84937	-.46325	.89402	.30331	.44740	-.34648
3	.68377	.28108	.84937	-.46325	.35488	.18442	.46531	-.36156
4	.59157	1.50271	.18313	1.18603	.59467	2.43277	.73569	2.05939
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
70	-.53202	-.28108	-.69213	-.46325	-.47413	-.54895	-.59782	-.31266
<i>RMSE:</i>	$\sqrt{\frac{1}{T} \sum_{t=1}^T (W(LV_t) - F(LV_t))^2}$.3075	.4033	.3117	.3229

TABLE 2. Consecutive norm for gain and propagation of predictions.

Sequence "t"	Frobenius Norm $\ \cdot\ _F$			
	$P_{t+1} - P_t$	$K_{t+1} - K_t$	$P_{t+1}^* - P_t^*$	$M_{t+1} - M_t$
1	-	-	-	-
2	1.391595	.5148828	.008475387	8.145129e-5
3	7.044589e-7	7.491063e-9	1.094965e-10	1.138605e-12
4	9.341999e-15	1.487013e-16	4.388542e-17	1.357636e-17
5	1.734723e-18	0	0	0
6	0	0	0	0
\vdots	\vdots	\vdots	\vdots	\vdots
70	0	0	0	0

In sum, the numerical results show a good performance of our proposed algorithm. Once the specified matrices are obtained from the training sample, the Kálmán filtering technique yields an optimal prediction for the LV case values along with the error covariance matrices for the test sample.

5. Discussion and Conclusion

It should be emphasized that the PLS method of Wold is applicable to a given sample, where estimates for the endogeneous variables are given through the exogenous latent ones, and the case values of the LVs are also estimated. The algorithm uses a lot of OLS regressions and so, the estimation of the coefficient matrices is time demanding akin to the block Cholesky decomposition we use. This is the case when we have a long time series with small time intervals or data when the observations come frequently in subsequent order. Our point is that for the new observations, we need not to repeat the whole estimation procedure to obtain the model parameters, but instead we can update the latent variable scores with the help of the new observable data, the estimated matrices, and the Kálmán filtering technique.

In this way, an artificial intelligence is developed. The parameter matrices are estimated from a training sample at the beginning, and the latent variable scores are estimated as observable variables arrive one by one from the test sample. Moreover, there is no need for any distribution assumptions and the data are not necessarily independent. It should be noted that in the possession of a stationary time series, the matrix sequences K_t and M_t (as $t \rightarrow \infty$) tend to fixed points of an iteration finding the solution of a matrix Riccati equation (see [4]), but this is the topic of a further research.

Acknowledgement. The research reported in this paper which was carried out at the Budapest University of Technology and Economics has been supported by the National Research Development and Innovation Fund based on the charter of bolster issued by the National Research Development and Innovation Office under the auspices of the Ministry for Innovation and Technology; also supported by the National Research, Development and Innovation Fund (TUDFO/51757/2019-ITM, Thematic Excellence Program). The research was supported by the following project too:

EFOP-3.6.2-16-2017-00015, HU-MATHS-IN for deepening the activity of the Hungarian Industrial and Innovation Network.

References

- [1] Haavelmo, T., *The statistical implications of a system of simultaneous equations*, *Econometrica*, **11**(1943), 1-12.
- [2] Jöreskog, K.G., *Structural equation models in the social sciences specification, estimation and testing*, In: "Applications of Statistics", (P.R. Krishnaiah, Ed), North-Holland Publishing Co., 1977, 265-287.
- [3] Kálmán, R.E., *A new approach to linear filtering and prediction problems*, *Trans. ASME J. Basic. Eng.*, **82D**(1960), 35-45.
- [4] Kálmán, R.E., Bucy, R.S., *New results in linear filtering and prediction theory*, *Trans. Amer. Soc. Mech. Eng., J. Basic Eng.*, **83**(1961), 95-108.
- [5] Kiiveri, H., Speed, T.P., Carlin, J.B., *Recursive casual models*, *J. Aust. Math. Soc.*, **36**(1984), 30-52.
- [6] Rao, C.R., *Linear Statistical Inference and its Applications*, Wiley (1973).
- [7] Tenenhaus, M., Esposito Vinzi, V., Chatelinc, Y-L., Lauro, C., *PLS path modeling*, *Comput. Statist. Data Anal.*, **48.1**(2005), 159-205.
- [8] Wold, H., *Partial least squares*, In: "Encyclopedia of Statistical Sciences", (Kotz, S., Johnson, N. L., Ed), Wiley, New York, **6**(1985), 581-591.

Marianna Bolla

Budapest University of Technology and Economics,
Institute of Mathematics,
Műgyetem rkp. 3. Budapest 1111, Hungary
e-mail: marib@math.bme.hu

Fatma Abdelkhalek

Budapest University of Technology and Economics,
Institute of Mathematics,
Műgyetem rkp. 3. Budapest 1111, Hungary
Also at Faculty of Commerce,
Assiut University, Egypt
e-mail: fatma@math.bme.hu