

Dynamic State-Space Models*

Paul Karapanagiotidis[†]

Draft 6
June 3, 2014

Abstract

A review of the general state-space modeling framework. The discussion focuses heavily on the three prediction problems of forecasting, filtering, and smoothing within the state-space context. Numerous examples are provided detailing special cases of the state-space model and its use in solving a number of modeling issues. Independent sections are also devoted to both the topics of Factor models and Harvey's Unobserved Components framework.

Keywords: state-space models, signal extraction, unobserved components.
JEL: C10, C32, C51, C53, C58

1 Introduction

The *dynamic state-space model* was developed in the control systems literature, where physical systems are described mathematically as sets of inputs, outputs, and state variables, related by difference equations. The following, Section 2, describes the various versions of the linear state-space framework, discusses the relevant assumptions imposed, and provides examples encountered in economics and finance. Subsequently, Section 3 provides the analogous description of the more general nonlinear state-space framework. Section 4 then discusses some of the common terminologies related to the state-space framework within the different contexts they are encountered. Section 5 follows by discussing the general problems of state-space prediction, including forecasting, filtering, and smoothing. Moreover, it provides a number of simple applications to chosen models from Section 2. Sections 6 and 7 then go into more detail: Section

*I'd like to thank Christian Gourieroux for his helpful comments and suggestions.

[†]University of Toronto, Department of Economics, p.karapanagiotidis@utoronto.ca

6 briefly discusses the problem of prediction in the frequency domain and Section 7 outlines in detail the solutions to the forecasting, smoothing, and filtering problems, within the time domain. In particular, we interpret the solutions in terms of an orthogonal basis, and provide the MA and AR representations of the state-space model. Section 8 then details estimation of the state-space model parameters in the time domain. Finally, Section 9 discusses the equivalent Factor model representation, including the relationship between this representation, the VARMA, and the VECM models. It also discusses in more detail the Unobserved Components framework popularized by Harvey (1984,89).

2 Linear dynamic state-space model

2.1 The models

2.1.1 Weak linear state-space model

The weak form of the linear *dynamic state-space model* is as follows:

$$\mathbf{x}_t = \mathbf{F}_t \mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t, \quad (1a)$$

$$\mathbf{y}_t = \mathbf{H}_t \mathbf{x}_t + \boldsymbol{\eta}_t, \quad (1b)$$

with the moment restrictions:

$$E[\boldsymbol{\epsilon}_t] = \mathbf{0}, \quad Cov(\boldsymbol{\epsilon}_t, \boldsymbol{\epsilon}_{t-s}) = \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}_t} \mathbb{1}_{s=0}, \quad (2a)$$

$$E[\boldsymbol{\eta}_t] = \mathbf{0}, \quad Cov(\boldsymbol{\eta}_t, \boldsymbol{\eta}_{t-s}) = \boldsymbol{\Sigma}_{\boldsymbol{\eta}_t} \mathbb{1}_{s=0},$$

$$E[\boldsymbol{\epsilon}_{t-j} \boldsymbol{\eta}'_{t-s}] = \mathbf{0}, \quad \forall j, s \in \mathbb{Z},$$

$$\text{and } Cov(\mathbf{x}_0, \boldsymbol{\epsilon}_t) = Cov(\mathbf{x}_0, \boldsymbol{\eta}_t) = \mathbf{0}, \quad \forall t > 0,$$

where:

- \mathbf{y}_t is a $n \times 1$ vector of the observed values at time t .

- \mathbf{x}_t is a $p \times 1$ vector of state process values at time t .¹
- $\boldsymbol{\epsilon}_t$ and $\boldsymbol{\eta}_t$ are assumed uncorrelated with each other across all time lags, and their covariance matrices, $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}_t}$ and $\boldsymbol{\Sigma}_{\boldsymbol{\eta}_t}$, depend on t .
- \mathbf{F}_t is called the “system matrix” and \mathbf{H}_t the “observation matrix.” These matrices are assumed to be *non-stochastic* where \mathbf{F}_t is $p \times p$, and if we allow for more state processes than observed ones, \mathbf{H}_t is $n \times p$ where $n \geq p$.
- Equation (1a) is called the “state transition” equation and (1b) is called the “observation” or “measurement” equation.
- The state initial condition, \mathbf{x}_0 , is assumed stochastic with second order moments denoted $E[\mathbf{x}_0] = \boldsymbol{\mu}$ and $Var(\mathbf{x}_0) = \boldsymbol{\Sigma}_{\mathbf{x}_0}$. Finally, there exists zero correlation between the initial state condition and the observation and state error terms, for all dates $t > 0$.

2.1.2 The Gaussian linear state-space model

In the weak version of the linear dynamic state-space model, the assumptions concern only the first and second-order moments of the noise processes and initial state, or equivalently the first and second-order moments of the joint process $[(\mathbf{x}'_t, \mathbf{y}'_t)']$. We can also introduce a more restrictive version of the model by assuming, independent, and identically distributed, Gaussian white noises (IIN) for the errors of the state and measurement equations. The Gaussian linear state

¹Other terminology include the “signal” in the engineering context, “control variable” in the systems control literature, “latent factor” in the factor models approach, or “stochastic parameter” in the Bayesian context. See below for more details.

space model is therefore defined as:

$$\mathbf{x}_t = \mathbf{F}_t \mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t, \quad (1a)$$

$$\mathbf{y}_t = \mathbf{H}_t \mathbf{x}_t + \boldsymbol{\eta}_t, \quad (1b)$$

$$\text{where } \begin{pmatrix} \boldsymbol{\epsilon}_t \\ \boldsymbol{\eta}_t \end{pmatrix}_{t \geq 1} \sim IIN \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}_t} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{\boldsymbol{\eta}_t} \end{bmatrix} \right), \quad \mathbf{x}_0 \sim N(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{x}_0}), \quad (3a)$$

and \mathbf{x}_0 and the joint process $(\boldsymbol{\epsilon}_t, \boldsymbol{\eta}_t)$ are independent.

The Gaussian version of the state-space model is often used as a convenient intermediary tool. Indeed, under the assumption of Gaussian noise and initial state, we know that the joint process $[(\mathbf{x}'_t, \mathbf{y}'_t)']$ is Gaussian. This implies that all marginal and conditional distributions concerning the components of these processes are also Gaussian. If the distributions are easily derived, we get as a by-product the expression of the associated linear regressions and residual variances. Since these linear regressions and residual variances are functions of the first and second-order moments only, their expressions are valid even if the noises are not Gaussian—that is, for the weak linear state space model.

2.2 Examples

We will now discuss various examples of the state-space model. The first examples from 2.2.1-2.2.3 are descriptive models used for predicting the future; the second set of examples, 2.2.4-2.2.9 introduces some structure on the dynamics to capture measurement error, missing data, or aggregation. Finally, the last examples, 2.2.10-2.2.12, come from economic and financial applications.

2.2.1 The Vector Autoregressive model of order 1, VAR(1)

The (weak) Vector AutoRegressive model of order 1, $VAR(1)$, is defined as:

$$\mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t, \quad \text{with } \boldsymbol{\epsilon}_t \sim WWN(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon), \quad (4a)$$

$$\mathbf{y}_t = \mathbf{x}_t, \quad (4b)$$

with the condition of no correlation between initial state, \mathbf{x}_0 , and the error terms, $\boldsymbol{\epsilon}_t$, satisfied. Furthermore, WWN denotes “weak white noise” or a process with finite, constant, first and second-order moments which exhibits no serial correlation.

In this case the observation process coincides with the state process. This implies that:

$$\mathbf{y}_t = \mathbf{F}\mathbf{y}_{t-1} + \boldsymbol{\epsilon}_t, \quad \text{with } \boldsymbol{\epsilon}_t \sim WWN(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon), \quad (5)$$

which is the standard definition of the (weak) $VAR(1)$ process.

2.2.2 The univariate Autoregressive model of order p , AR(p)

The (weak) univariate AutoRegressive model of order p , $AR(p)$, is defined as:

$$x_t + b_1x_{t-1} + \dots + b_px_{t-p} = \epsilon_t, \quad \text{with } \epsilon_t \sim WWN(0, \sigma_\epsilon^2), \quad (6)$$

The model can be written in state-space form as:

$$\mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t, \quad (7a)$$

$$y_t = \mathbf{H}\mathbf{x}_t. \quad (7b)$$

The state vector includes the current and first $p - 1$ lagged values of \mathbf{x}_t :

$$\mathbf{x}_t = \begin{bmatrix} x_t & x_{t-1} & \dots & x_{t-p+1} \end{bmatrix}'_{p \times 1}, \quad (8)$$

with system matrices are given as:

$$\mathbf{F} = \begin{bmatrix} -b_1 & -b_2 & \dots & -b_{p-1} & -b_p \\ 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix}_{p \times p}, \quad (9a)$$

$$\mathbf{H} = \begin{bmatrix} 1 & 0 & \dots & 0 \end{bmatrix}_{1 \times p}, \quad (9b)$$

$$\text{and } \boldsymbol{\epsilon}_t = \begin{bmatrix} \epsilon_t & 0 & \dots & 0 \end{bmatrix}'_{p \times 1}. \quad (9c)$$

Since the $AR(p)$ process is completely observed, $\boldsymbol{\eta}_t = 0$ and $\boldsymbol{\Sigma}_{\boldsymbol{\eta}_t} = 0$ for all t . Moreover, $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}_t}$ is a singular matrix with zeroes in each element except the top diagonal element, which is equal to σ_ϵ^2 for all t .

2.2.3 The univariate Autoregressive-Moving Average model of order (p, q) , ARMA(p,q)

The (weak) univariate AutoRegressive-Moving Average model of order (p, q) , $ARMA(p, q)$, is defined as:

$$x_t + b_1 x_{t-1} + \dots + b_p x_{t-p} = \epsilon_t + a_1 \epsilon_{t-1} + \dots + a_q \epsilon_{t-q}, \quad \text{where } \epsilon_t \sim WWN(0, \sigma_\epsilon^2), \quad (10)$$

There are a number of possible state-space representations of an ARMA process. In the language of Akaike (1975), a “minimal” representation is a representation whose state vector elements represent the minimum collection of variables which contain all the information needed to produce forecasts given some forecast origin t . For ease of exposition in what follows we provide a non-minimal state-space representation, although the interested reader can consult Gouriéroux (1997, p. 607) for a minimal one.

Let the dimension of the state \mathbf{x}_t be $m = p + q$. We have:

$$\mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \mathbf{G}\epsilon_t, \quad (11a)$$

$$y_t = \mathbf{H}\mathbf{x}_t, \quad (11b)$$

The state vector is given as:

$$\mathbf{x}_t = \begin{bmatrix} x_t & x_{t-1} & \dots & x_{t-(p-2)} & x_{t-(p-1)} & \epsilon_t & \epsilon_{t-1} & \dots & \epsilon_{t-(q-2)} & \epsilon_{t-(q-1)} \end{bmatrix}'. \quad (12a)$$

The system matrices are given as:

$$\mathbf{F} = \begin{bmatrix} -b_1 & \dots & \dots & -b_{p-1} & -b_p & a_1 & \dots & \dots & a_{q-1} & a_q \\ 1 & 0 & \dots & 0 & 0 & 0 & \dots & \dots & \dots & 0 \\ 0 & 1 & \dots & 0 & 0 & 0 & \dots & \dots & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \dots & \dots & \dots & \vdots \\ 0 & 0 & \dots & 1 & 0 & 0 & \dots & \dots & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 & 0 & \dots & \dots & \dots & 0 \\ 0 & \dots & \dots & 0 & 0 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \dots & \dots & \vdots & \vdots & 0 & 1 & \dots & 0 & 0 \\ \vdots & \dots & \dots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & \dots & 0 & 0 & 0 & 0 & \dots & 1 & 0 \end{bmatrix}_{m \times m}, \quad (13a)$$

$$\mathbf{G} = \begin{bmatrix} 1 & 0 & \dots & 0 & 1 & 0 & \dots & 0 \end{bmatrix}'_{m \times 1}, \quad (13b)$$

$$\mathbf{H} = \begin{bmatrix} 1 & 0 & \dots & 0 \end{bmatrix}_{1 \times m}, \quad (13c)$$

$$\text{and } \epsilon_t \equiv \epsilon_t \sim WWN(0, \sigma_\epsilon^2). \quad (13d)$$

with the condition of no correlation between initial state, \mathbf{x}_0 , and the error terms, ϵ_t , satisfied.

2.2.4 Partially observed VARMA

In the discrete-time multivariate case, suppose that some linear, unobserved, state process \mathbf{x}_t is composed of filtered weak white noise $\boldsymbol{\epsilon}_t$. However, what we actually observe is \mathbf{y}_t , which has been corrupted by additive weak white noise $\boldsymbol{\eta}_t$. This model can be written as:

$$\mathbf{x}_t = \sum_{u=0}^{\infty} \mathbf{a}_u \boldsymbol{\epsilon}_{t-u} = \mathbf{A}(L)\boldsymbol{\epsilon}_t, \quad (14a)$$

$$\text{and } \mathbf{y}_t = \mathbf{x}_t + \boldsymbol{\eta}_t, \quad \text{for } t = 1, \dots, T, \quad (14b)$$

where:

- \mathbf{x}_t is an unobserved $P \times 1$ vector of state variables.
- the unobserved state \mathbf{x}_t (the “signal” in the engineering context) is corrupted with weak, white, additive noise, $\boldsymbol{\eta}_t$, with covariance matrix $\boldsymbol{\Sigma}_\eta$.
- $\boldsymbol{\epsilon}_t$ is an $M \times 1$ vector of weak white noise input processes with covariance $\boldsymbol{\Sigma}_\epsilon$.²
- \mathbf{y}_t is a $P \times 1$ vector of the observed “noisy” output process.
- $\boldsymbol{\Sigma}_\eta$ and $\boldsymbol{\Sigma}_\epsilon$ represent the covariance of the measurement noise and the state process noise, respectively, and are assumed time-invariant.

Note that $\mathbf{A}(L)$ is a $P \times M$ matrix infinite series where L denotes the lag operator; that is, $\mathbf{A}(L) = \mathbf{a}_0 L^0 + \mathbf{a}_1 L^1 + \mathbf{a}_2 L^2 + \dots$, where the individual $P \times M$ matrices, \mathbf{a}_u , collectively represent the impulse response function f where $f : \mathbb{Z} \rightarrow \mathbb{R}^{P \times M}$.

Given (14a), the infinite lag distribution makes working with this model in the time domain troublesome. However, the apparent multi-stage dependence can be reduced to first-order autoregressive dependence by means of a matrix representation if we assume that $\mathbf{A}(L)$ can be well-approximated by a ratio of finite lag matrix polynomials – the so called “transfer function”

²Note there is no loss of generality here since for any (possibly non-white) second-order stationary stochastic input process we might choose, \mathbf{z}_t , we can always represent it as $\mathbf{z}_t = \sum_{j=0}^{\infty} \boldsymbol{\Theta}_j \boldsymbol{\epsilon}_{t-j}$ by Wold’s theorem.

models of Box & Jenkins (1970). That is, suppose that we model instead:

$$\begin{aligned} C(L) &= B(L)^{-1}A(L) \\ &= (\mathbf{I} + \mathbf{b}_1L + \mathbf{b}_2L^2 + \cdots + \mathbf{b}_nL^n)^{-1} (\mathbf{a}_0 + \mathbf{a}_1L + \mathbf{a}_2L^2 + \cdots + \mathbf{a}_{n-1}L^{n-1}) \end{aligned} \quad (15)$$

where the inverse of the matrix lag polynomial $B(L)$ is assumed to exist.

The model in (14), with the ratio of finite order matrix lag polynomials $C(L)$ now replacing the infinite series $A(L)$, becomes:

$$\mathbf{x}_t + \mathbf{b}_1\mathbf{x}_{t-1} + \cdots + \mathbf{b}_n\mathbf{x}_{t-n} = \mathbf{a}_0\boldsymbol{\epsilon}_t + \cdots + \mathbf{a}_{n-1}\boldsymbol{\epsilon}_{t-(n-1)}, \quad \forall t = 1, \dots, T, \quad (16a)$$

$$\text{and } \mathbf{y}_t = \mathbf{x}_t + \boldsymbol{\eta}_t, \quad (16b)$$

(16) can now be reduced to first-order dependence by means of a redefinition of the state vector as:

$$\mathbf{x}_t^* = \mathbf{T}\hat{\mathbf{x}}_{t-1}, \quad (17a)$$

$$\text{where } \mathbf{T} = \begin{bmatrix} -\mathbf{b}_n & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ -\mathbf{b}_{n-1} & -\mathbf{b}_n & \mathbf{0} & \cdots & \mathbf{0} \\ -\mathbf{b}_{n-2} & -\mathbf{b}_{n-1} & -\mathbf{b}_n & \cdots & \mathbf{0} \\ \vdots & & & \ddots & \mathbf{0} \\ -\mathbf{b}_1 & -\mathbf{b}_2 & -\mathbf{b}_3 & \cdots & -\mathbf{b}_n \end{bmatrix}_{P_n \times P_n}, \quad (17b)$$

$$\text{and } \hat{\mathbf{x}}'_{t-1} = \begin{bmatrix} \mathbf{x}'_{t-1} & \mathbf{x}'_{t-2} & \cdots & \mathbf{x}'_{t-n} \end{bmatrix}_{1 \times P_n}, \quad (17c)$$

so that the new first-order autoregressive state-space model takes the form:

$$\mathbf{x}_t^* = \mathbf{F}\mathbf{x}_{t-1}^* + \mathbf{G}\hat{\boldsymbol{\epsilon}}_t, \quad (18a)$$

$$\text{and } \mathbf{y}_t = \mathbf{H}\mathbf{x}_t^* + \boldsymbol{\eta}_t \quad (18b)$$

where we have that:

$$\mathbf{F} = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & -\mathbf{b}_n \\ \mathbf{I} & \mathbf{0} & \dots & \mathbf{0} & -\mathbf{b}_{n-1} \\ \vdots & \vdots & & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{I} & -\mathbf{b}_1 \end{bmatrix}_{Pn \times Pn}, \quad (19a)$$

$$\mathbf{G} = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{a}_0 & \mathbf{a}_1 & \dots & \mathbf{a}_{n-1} \end{bmatrix}_{Pn \times Mn}, \quad (19b)$$

$$\mathbf{H} = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \dots & \mathbf{I} \end{bmatrix}_{P \times Pn}, \quad (19c)$$

$$\text{and } \tilde{\boldsymbol{\epsilon}}'_t = \begin{bmatrix} \boldsymbol{\epsilon}'_t & \boldsymbol{\epsilon}'_{t-1} & \dots & \boldsymbol{\epsilon}'_{t-(n-1)} \end{bmatrix}_{1 \times Mn} \quad (19d)$$

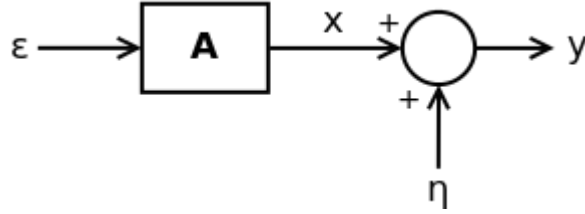
where again \mathbf{F} is the “system matrix,” \mathbf{G} the “input matrix,” and \mathbf{H} the “observation matrix.” Bear in mind that (18) and (19) represent only one possible state-space representation – in fact while the transfer function $\mathbf{C}(e^{-i\omega})$ in $\mathbf{h}_{xx}(\omega)$ (see Section 2.2.4, (20a)) implies an infinite number of possible state-space representations, any particular state-space representation has only one equivalent transfer function. Additionally, we can immediately see from (18b) that the observed process \mathbf{y}_t is a linear function of the unobserved “factors” since $\mathbf{y}_t = -\mathbf{b}_1 \mathbf{x}_{t-1} - \mathbf{b}_2 \mathbf{x}_{t-2} - \dots - \mathbf{b}_n \mathbf{x}_{t-n} + \mathbf{u}_t + \boldsymbol{\eta}_t$, where \mathbf{u}_t is equal to the right-hand side of (16a). See Akaike (1974) for a general treatment of finite order linear systems.

i) Spectral properties of the partially observed VARMA process

Note that from (14a), $\mathbf{A}(L) = \sum_{u=0}^{\infty} \mathbf{a}_u L^u$, so $\mathbf{A}(L)\boldsymbol{\epsilon}_t = \sum_{u=0}^{\infty} \mathbf{a}_u \boldsymbol{\epsilon}_{t-u}$. More generally, $\mathbf{A}(z) = \sum_{u=0}^{\infty} \mathbf{a}_u z^u$ where $z \in \mathbb{C}$ is known as the z-transform. Therefore, while $\mathbf{A}(L)$ is a polynomial function of an *operator* L , the z-transform, $\mathbf{A}(z)$, is a polynomial function of a complex variable. However, since both polynomials admit the same coefficients, we can solve for the *transfer function* of (14a) as $\mathbf{A}(z)$ where $z = e^{-i\omega}$, since this represents the Fourier

transform of the impulse response function.³ (Note that in continuous time this z-transform analogy is unnecessary since there is no need for defining the model in terms of lag operators, L). Therefore, the convolution observed in the time domain in (14a) is equivalent to a multiplication within the frequency domain, so that the Fourier transform of the impulse response, $\mathbf{A}(e^{-i\omega})$, disentangles the complicated interdependencies into a simple multiplicative relation between inputs and outputs given any frequency ω . Therefore, working with (14) in the frequency domain is often a useful approach. For clarity the frequency domain relationships are given diagrammatically in Figure 1.

Figure 1: Frequency domain relationships of the model in (14)



Since ϵ_t and η_t are jointly stationary and uncorrelated we have that:

$$\mathbf{h}_{yx}(\omega) = \mathbf{A}(e^{-i\omega})\mathbf{h}_{\epsilon\epsilon}(\omega)\mathbf{A}(e^{+i\omega})' = \frac{1}{2\pi}\mathbf{A}(e^{-i\omega})\Sigma_{\epsilon}\mathbf{A}(e^{+i\omega})' = \mathbf{h}_{xx}(\omega), \quad (20a)$$

$$\text{and } \mathbf{h}_{y\epsilon}(\omega) = \mathbf{A}(e^{-i\omega})\mathbf{h}_{\epsilon\epsilon}(\omega) = \frac{1}{2\pi}\mathbf{A}(e^{-i\omega})\Sigma_{\epsilon} \quad (20b)$$

represent the cross-spectral density matrices between \mathbf{y}_t and \mathbf{x}_t , \mathbf{y}_t and ϵ_t respectively. Therefore, from (20a) it is clear that \mathbf{x}_t represents “filtered” weak white noise, where the flat spectrum of ϵ_t (i.e. its variance) is given shape by $\mathbf{A}(e^{\pm i\omega})$.

Furthermore, the spectral-density matrix of \mathbf{y}_t is (from (16)):

$$\mathbf{h}_{yy}(\omega) = \mathbf{h}_{xx}(\omega) + \mathbf{h}_{\eta\eta}(\omega) = \frac{1}{2\pi} (\mathbf{A}(e^{-i\omega})\Sigma_{\epsilon}\mathbf{A}(e^{+i\omega})' + \Sigma_{\eta}). \quad (21)$$

³The system in (14) is constrained to be “physically realizable” by assuming the impulse response matrices are $\mathbf{a}_j = 0, \forall j < 0$. This form of impulse response exists, is unique, and is quadratically summable, with no zeros inside the unit circle as long as the integral from $-\pi$ to π of the log of ϵ_t ’s spectral density is finite – see Doob (1953), as cited in Priesley (1981, pg. 733). Note this condition is a very weak one and is satisfied here – in fact, the mentioned integral can only diverge to $-\infty$ if the spectral density vanishes in some interval in its domain.

Note that for the discrete time process all spectral densities are continuous in the frequency ω and periodic with period 2π .

Finally, a natural non-parametric estimator of the transfer function matrix is given by:

$$\hat{\mathbf{A}}(e^{-i\omega}) = \hat{\mathbf{h}}_{y\epsilon}(\omega)\hat{\mathbf{h}}_{\epsilon\epsilon}^{-1}(\omega) = 2\pi\hat{\mathbf{h}}_{y\epsilon}(\omega)\widehat{\boldsymbol{\Sigma}}_{\epsilon}^{-1} \quad (22)$$

where the spectral densities in (22) can be estimated within the frequency domain. See Priestley (1981, Section 9.5) for more details.

Now, suppose we wish to establish the optimal manner of extracting the signal \mathbf{x}_t given only the noisy observations \mathbf{y}_t . That is, we wish to establish the optimal frequency response, or transfer function $\mathbf{C}(\omega)$, in Figure 2. It was Wiener that originally solved this frequency domain problem where he established the optimal frequency response as the ratio:⁴

$$\mathbf{C}(\omega) = \frac{\mathbf{h}_{xy}(\omega)}{\mathbf{h}_{yy}(\omega)} = \frac{\mathbf{h}_{xx}(\omega)}{\mathbf{h}_{xx}(\omega) + \mathbf{h}_{\eta\eta}(\omega)}. \quad (23)$$

Therefore, the Wiener filter attenuates those frequencies in which the signal to noise ratio is low and passes through those where it is high.

Figure 2: Wiener filter - the optimal transfer function $\mathbf{C}(\omega)$



⁴Noting of course that since $E[\mathbf{x}_t\mathbf{x}'_{t-s}]$ is symmetric in the time domain for all s , we have that $\mathbf{h}_{xx}(\omega)$ is real, and so $\mathbf{h}_{xy}(\omega) = \mathbf{h}_{xx}(\omega) = \mathbf{h}_{yx}(\omega)$ without the need of taking complex conjugates.

2.2.5 The VAR(1) with measurement error

A special case of the partially observed VARMA model in Section 2.2.4 arises as the (weak) VAR(1) with measurement error is defined as:

$$\mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t, \quad (24a)$$

$$\mathbf{y}_t = \mathbf{x}_t + \boldsymbol{\eta}_t, \quad (24b)$$

$$\begin{pmatrix} \boldsymbol{\epsilon}_t \\ \boldsymbol{\eta}_t \end{pmatrix}_{t \geq 1} \sim WWN \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_\epsilon & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_\eta \end{bmatrix} \right), \quad (24c)$$

with the condition of no correlation between initial state, \mathbf{x}_0 , and the error terms process, $(\boldsymbol{\epsilon}_t, \boldsymbol{\eta}_t)$, satisfied.

Therefore, the state-space process is a VAR(1) process, but measured with a multivariate error given by $\boldsymbol{\eta}_t$. The process (\mathbf{y}_t) is such that:

$$\begin{aligned} \mathbf{y}_t - \mathbf{F}\mathbf{y}_{t-1} &= \mathbf{x}_t + \boldsymbol{\eta}_t - \mathbf{F}(\mathbf{x}_{t-1} + \boldsymbol{\eta}_{t-1}) \\ &= \boldsymbol{\epsilon}_t + \boldsymbol{\eta}_t - \mathbf{F}\boldsymbol{\eta}_{t-1} \equiv \mathbf{v}_t \end{aligned} \quad (25a)$$

The process (\mathbf{v}_t) has serial covariances equal to zero for lags larger or equal to 2. Therefore, \mathbf{v}_t admits a Vector Moving Average, VMA(1), representation of order 1. Let $\mathbf{v}_t = \mathbf{u}_t - \boldsymbol{\Theta}\mathbf{u}_{t-1}$.⁵ We can therefore deduce that the process \mathbf{y}_t has a Vector Autoregressive-Moving Average of order (1,1), or VARMA(1,1), representation:

$$\mathbf{y}_t - \mathbf{F}\mathbf{y}_{t-1} = \mathbf{u}_t - \boldsymbol{\Theta}\mathbf{u}_{t-1}, \quad \text{with } \mathbf{u}_t \sim WWN(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad (26)$$

$\boldsymbol{\Theta}$ and $\boldsymbol{\Sigma}_u$ are functions of the initial parameters of the state-space representation: \mathbf{F} , $\boldsymbol{\Sigma}_\epsilon$, and

⁵And so we have that there exists no correlation between initial state, \mathbf{x}_0 , and the new error terms process, (\mathbf{u}_t) .

Σ_η . They are related by the matrix equations system:

$$-F\Sigma_\eta = \Theta\Sigma_u \quad (27a)$$

$$\Sigma_\eta + F\Sigma_\eta F' + \Sigma_\epsilon = \Sigma_u + \Theta\Sigma_u\Theta' \quad (27b)$$

which can be solved numerically for values of Θ and Σ_u .

2.2.6 Static state-space model

The (weak) static state-space model is defined as:

$$\mathbf{x}_t = \boldsymbol{\epsilon}_t, \quad (28a)$$

$$\mathbf{y}_t = \mathbf{H}\mathbf{x}_t + \boldsymbol{\eta}_t, \quad (28b)$$

$$\begin{pmatrix} \boldsymbol{\epsilon}_t \\ \boldsymbol{\eta}_t \end{pmatrix}_{t \geq 1} \sim WWN \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_\epsilon & \mathbf{0} \\ \mathbf{0} & \Sigma_\eta \end{bmatrix} \right), \quad (28c)$$

with the condition of no correlation between initial state, \mathbf{x}_0 , and the error terms, $(\boldsymbol{\epsilon}_t, \boldsymbol{\eta}_t)$, satisfied.

Therefore, from (28) the distribution of the state-space process is such that:

$$\begin{bmatrix} \mathbf{x}_t \\ \mathbf{y}_t \end{bmatrix}_{t \geq 1} = WWN \left[\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \Sigma_\epsilon & \Sigma_\epsilon \mathbf{H}' \\ \mathbf{H}\Sigma_\epsilon & \mathbf{H}\Sigma_\epsilon \mathbf{H}' + \Sigma_\eta \end{pmatrix} \right]. \quad (29)$$

In general, the state-space form is equivalent to the *factor model* representation, where we assume that some p factors, \mathbf{x}_t , influence the n observed processes \mathbf{y}_t , where $n > p$. Indeed, the goal of the factor model representation is to model the observed processes in terms of a smaller number of factor processes. Therefore, the particular form of the state-space model in (28) is equivalent to the *static* factor model representation, although it is clear that the factor may

instead be formulated in a dynamic manner as in (1a).

We can distinguish two cases in practice:

- Case 1: the factor \mathbf{x}_t is unobserved.

In this case, the unrestricted theoretical model above is unidentifiable. This is because the number of parameters exceeds the number of population moment conditions when \mathbf{x}_t is unobserved. Indeed, from (29), we have:

$$Var(\mathbf{y}_t)_{n \times n} = \mathbf{H}_{n \times p} \Sigma_{\epsilon_{p \times p}} \mathbf{H}'_{p \times n} + \Sigma_{\eta_{n \times n}}, \quad (30)$$

and so the $n(n+1)/2$ population second moment conditions are outnumbered by the $np + p(p+1)/2 + n(n+1)/2$ parameters.

Therefore, moment constraints on the parameters of the factor model are usually introduced to ensure identification. For example, we can assume without loss of generality that the factor covariance matrix, Σ_{ϵ} , is an identity matrix. Indeed, \mathbf{x}_t is unobserved and defined up to an invertible linear transform. That is, for any invertible matrix \mathbf{L} equation (28b) with $\mathbf{x}_t^* = \mathbf{L}\mathbf{x}_t$ and $\mathbf{H}^* = \mathbf{H}\mathbf{L}^{-1}$ is observationally equivalent. Therefore, for any Σ_{ϵ} we can always introduce the transformation $\mathbf{x}_t^* = \mathbf{D}^{-1/2} \mathbf{P}' \mathbf{x}_t$, where the matrix \mathbf{P} has an orthonormal basis of eigenvectors of Σ_{ϵ} as its columns and \mathbf{D} is a diagonal matrix of the respective eigenvalues, so as to make $Var(\mathbf{x}_t^*) = \mathbf{I}$. Moreover, it is often assumed in the literature that the observation error covariance matrix, Σ_{η} , is diagonal (or even scalar) so that $\Sigma_{\eta} = \sigma_{\eta}^2 \mathbf{I}$. This additional constraint is considered a “real” constraint since it reduces the model’s flexibility in favour of identification.

- Case 2: the factor \mathbf{x}_t is observed.

In this, case the unrestricted theoretical model becomes identifiable. Given the moment condition $Var(\mathbf{x}_t)$, Σ_{ϵ} is identified. We can then formulate the theoretical linear regression in (28b) to identify both \mathbf{H} and Σ_{η} .

The static factor model is popular in Finance, where the observed variable \mathbf{y}_t represents the returns of a set of assets. Under a (weak) efficient market hypothesis the returns are WWN. The observation equation (28b) thus decomposes the returns into a market component, $\mathbf{H}\mathbf{x}_t$, and a firm specific component, $\boldsymbol{\eta}_t$. Assuming an uncorrelated market component, the unobserved factors, \mathbf{x}_t , represent the returns on the market and $\boldsymbol{\eta}_t$ represent the firm specific returns, whose variability (or “idiosyncratic” risk) can be reduced through its addition to a well diversified portfolio. Of course, the assumption of uncorrelated market component can be generalized within the dynamic model. For more on the factor model representation, see Section 9.

2.2.7 The state-space model for “data aggregation”

Suppose that we assume the state vector \mathbf{x}_t represents some individual level components which we desire to aggregate in some way. In a model for aggregation you have to distinguish between both the behavioural equation which generally includes an error term, and the accounting relationship with no error term.

Therefore, let \mathbf{y}_t represent the observed aggregate variable, and let \mathbf{x}_t represent some possibly unobserved individual level variables. The state-space formulation defines both the behavioural equation for \mathbf{x}_t and the accounting equation for \mathbf{y}_t as:

$$\mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t, \quad \text{where } \boldsymbol{\epsilon}_t \sim WWN(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon), \quad (31a)$$

$$\mathbf{y}_t = \boldsymbol{\alpha}'\mathbf{x}_t, \quad (31b)$$

where $\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 & \alpha_2 & \dots & \alpha_p \end{bmatrix}'$ is a p vector of size adjustment parameters (which may possibly sum to 1) and so we can model the observed values y_t as the weighted aggregate of individual factors, the elements of \mathbf{x}_t .

Note that in an accounting relationship you can only add variables with the same unit. Therefore, we have first to transform the elements of \mathbf{x}_t into a common unit, which is usually done by considering some measure of value in common units, e.g. dollars.

The aggregation model can also be employed in the Finance context. In this case, as opposed to Section 2.2.6, the observation equation represents an accounting relationship between asset returns and the aggregate portfolio return, not a behavioural relationship. The returns may be weighted according to their contribution to the overall portfolio, where again the returns are written in the same domination, e.g. dollars.

2.2.8 The VAR(1) with “series selection” or “missing series”

The Vector AutoRegressive process of order 1, or VAR(1), with “series selection” or “missing series” is defined as:

$$\mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t, \quad \text{where } \boldsymbol{\epsilon}_t \sim WWN(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon), \quad (32a)$$

$$\text{and } y_t = x_{i,t}, \quad (32b)$$

where $x_{i,t}$ denotes the i 'th element of \mathbf{x}_t . Therefore the model can be interpreted in two ways, depending on whether or not \mathbf{x}_t is observed:

- Case 1: \mathbf{x}_t is observed.

The model is then interpreted as a method of selecting only that series from the state vector \mathbf{x}_t that is of interest. Notice that (31b) above is a special case of series selection, when \mathbf{x}_t is observed.

- Case 2: \mathbf{x}_t is not observed.

The model is interpreted as the case of “missing series.” That is, some of the elements of the series (\mathbf{x}_t) are missing.

2.2.9 The VAR(1) with “missing data”

The model in Section 2.2.8, with unobserved \mathbf{x}_t , can of course be generalized to the cases where not only are series missing, but perhaps individual data elements of some series are missing as

well. We call this the vector autoregressive process of order 1, or VAR(1), for “missing data”.

The state equation is the same for all cases below:

$$\mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t, \quad \text{where } \boldsymbol{\epsilon}_t \sim WWN(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon). \quad (32a)$$

- Case 1: the i 'th series is missing and some elements of the j 'th series are missing:

$$\mathbf{y}_t = \begin{cases} \left[x_{1,t} \ \dots \ x_{i-1,t} \ x_{i+1,t} \ \dots \ x_{j,t} \ \dots \ x_{p,t} \right]' & \text{if } t \neq m \\ \left[x_{1,t} \ \dots \ x_{i-1,t} \ x_{i+1,t} \ \dots \ x_{j-1,t} \ x_{j+1,t} \ \dots \ x_{p,t} \right]' & \text{if } t = m \end{cases}. \quad (33)$$

- Case 2: the i 'th and j 'th series both have missing data but the missing points occur at the same time:

$$\mathbf{y}_t = \begin{cases} \mathbf{x}_t & \text{if } t \neq m \\ \left[x_{1,t} \ \dots \ x_{i-1,t} \ x_{i+1,t} \ \dots \ x_{j-1,t} \ x_{j+1,t} \ \dots \ x_{p,t} \right]' & \text{if } t = m \end{cases}. \quad (34)$$

- Case 3: the i 'th and j 'th series both have missing data with no inherent pattern.

Where in each case, $m \in 0, \dots, T$, denotes a time period upon which some elements of the vector \mathbf{x}_t are missing.

2.2.10 The Unobserved Components model

Consider the special case of the state-space model for aggregation in Section 2.2.7, where the elements of $\boldsymbol{\alpha}$ are all equal to one, and we assume that the p elements of \mathbf{x}_t are independent of each other with specified marginal distributions or we at least specify their first two moments:

$$\mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t, \quad (35a)$$

$$y_t = \begin{bmatrix} 1 & \dots & 1 \end{bmatrix} \mathbf{x}_t, \quad (35b)$$

the observed series \mathbf{y}_t is therefore the sum of various *components*, generally unobserved.

2.2.10.1 The General Stochastic Trend

P.C. Young (2011, p.67) defines the *generalized random walk* or *general stochastic trend* as:

$$\mathbf{x}_t = \begin{bmatrix} x_{1,t} \\ \Delta x_{1,t} \end{bmatrix}, \quad (36a)$$

$$\mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \mathbf{G}\boldsymbol{\epsilon}_t, \quad \text{where } \mathbf{F} = \begin{bmatrix} \alpha & \beta \\ 0 & \gamma \end{bmatrix}, \quad \text{and } \begin{bmatrix} \delta & 0 \\ 0 & \epsilon \end{bmatrix}. \quad (36b)$$

$$\text{Also } \boldsymbol{\epsilon}_t \sim WWN(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon) \quad \text{where } \boldsymbol{\Sigma}_\epsilon \text{ is diagonal.} \quad (36c)$$

That is, we have defined the state process in such a manner as to allow us to modify the behaviour of the typical random walk in different ways. For example, if $\beta = \gamma = \epsilon = 0$ and $\alpha = \delta = 1$ the model represents the standard random walk. However, if $\alpha = \beta = \gamma = \epsilon = 1$ and $\delta = 0$ we have the *integrated random walk* which is smoother than the standard random walk. Moreover, if $0 < \alpha < 1$ and $\beta = \gamma = \epsilon = 1$ and $\delta = 0$ we have the case of the *smoothed random walk*. Also the case of $\beta = \gamma = \epsilon = 0$ and $0 < \alpha < 1$ and $\delta = 1$ is equivalent to the AR(1) model. Finally, both the *Local Linear Trend* (see Section 2.2.10.2) and *Damped Trend* from Harvey (1984,89) are both given by $\alpha = \beta = \gamma = \epsilon = \delta = 1$ (except in the latter case $0 < \gamma < 1$).

2.2.10.2 Harvey's Unobserved Components models: the "basic structural model"

Harvey (1984,89) attempts to decompose the series (\mathbf{y}_t) into a number of unobserved, orthogonal, components representing *trends, seasonals, other cycles, and irregular patterns*, all of which are

informed by the spectral properties of the observed series. For example, consider the model:

$$y_t = T_t + S_t + C_t + I_t \quad (37)$$

where y_t is the observed series, T_t is some trend component, S_t is a seasonal component, C_t is some other cyclical component, and I_t represents the irregular pattern.

Typically the trend component T_t is associated with the slowly changing, low frequency component of y_t (i.e. a spectral frequency close to zero, or equivalently a period close to ∞). It can be modeled by the stochastic counterpart of the linear time trend $\mu_t = \mu_0 + \beta t$, called the *Local Linear Trend* model:

$$T_t \equiv \mu_t = \mu_{t-1} + \beta_{t-1} + v_t, \quad \text{where } v_t \sim WWN(0, \sigma_v^2), \quad (38a)$$

$$\text{and } \beta_t = \beta_{t-1} + z_t, \quad \text{where } z_t \sim WWN(0, \sigma_z^2). \quad (38b)$$

Of course, the Local Linear Trend formulation is a special case of the general stochastic trend in Section 2.2.10.1.

Furthermore, the seasonal component S_t can be modeled as dummy intercepts which are constrained to sum to zero (with some small stochastic residual difference, ω). For example, suppose s is the number of “seasons” (say 12 for monthly data) and $z_{j,t}$ for $j = 1, 2, \dots, s$ is some set of dummy variables that take on the values:

$$z_{j,t} = \begin{cases} 1, & \text{if } t = j, j + s, j + 2s, \dots \\ 0, & \text{if } t \neq j, j + s, j + 2s, \dots \\ -1, & \text{if } t = s, 2s, 3s, \dots \end{cases} \quad (39)$$

then, if γ_j is the dummy intercept for time j , we have that at $t = xs$ for all $x \in \mathbb{N}^+$:

$$\sum_{j=1}^{s-1} z_{j,t} \gamma_j = - \sum_{j=1}^{s-1} \gamma_j \equiv \gamma_s \quad (40a)$$

$$\Leftrightarrow \sum_{j=1}^s \gamma_j = 0 \quad (40b)$$

and given a change in the notation, (40b) can be rewritten as $\sum_{j=0}^{s-1} \gamma_{t-j} = 0$. Adding a disturbance term with zero expectation to the right hand side allows the seasonal effect to vary stochastically:

$$\sum_{j=0}^{s-1} \gamma_{t-j} = \omega_t \quad \text{where} \quad \omega_t \sim WWN(0, \sigma_\omega^2), \quad (41a)$$

$$\Leftrightarrow (1 + L + L^2 + \dots + L^{s-1}) \gamma_t = \omega_t. \quad (41b)$$

Finally, the cyclical component C_t can be written as a sum of *stochastic harmonics*, where each component in the sum reflects some particular chosen frequency, $\lambda_j = 2\pi \frac{j}{s}$, where $j \leq \frac{s}{2}$. For example, given monthly data, let s be such that $s \pmod{12} = 0$, and let $j \in \mathbb{N}^+$ be chosen so that s/j represents the desired periodicity of the harmonic function. Therefore, we could choose that $s = 12$ and $j = 1$ so that the period is 12; that is, the cycle repeats every 12 months. Alternatively, if $j = 6$ then the period is 2 and the cycle repeats every 2 months, etc.

The cyclical component C_t can therefore be written as:

$$C_t \equiv \sum_{k \in J} c_{k,t} \quad (42a)$$

$$\text{where} \quad c_{k,t} = \rho_k \{c_{k,t-1} \cos \lambda_k + c_{k,t-1}^* \sin \lambda_k\} + \xi_{k,t}, \quad (42b)$$

$$\text{and} \quad c_{k,t}^* = \rho_k \{-c_{k,t-1} \sin \lambda_k + c_{k,t-1}^* \cos \lambda_k\} + \xi_{k,t}^*, \quad (42c)$$

where J is the set of chosen frequencies, ρ_k is a discount parameter, and ξ_k and ξ_k^* are zero mean *WWN* processes which are uncorrelated with each other, with common variance $\sigma_{\xi,k}^2$. For more

details on the stochastic harmonic cycles approach, see Hannan, Terrell and Tuckwell (1970).

Finally, the irregular component takes the form of a *WWN* innovation, $I_t \equiv \eta_t$. Putting all the components together into the state-space form with *WWN* innovations, we have the observation equation

$$\begin{aligned}
 y_t &= T_t + S_t + C_t + I_t \\
 \Leftrightarrow y_t &= \mu_t + \gamma_t + \sum_{k \in J} c_{k,t} + \eta_t \\
 &= \begin{bmatrix} 1 & 0 & 1 & 0 & 0 & \dots & 1 & 0 & 1 & \dots \end{bmatrix} \mathbf{x}_t + \eta_t \\
 &\equiv \mathbf{H}_t \mathbf{x}_t + \boldsymbol{\eta}_t,
 \end{aligned} \tag{43a}$$

and the state transition equation

$$\begin{aligned}
 \mathbf{x}_t = \begin{bmatrix} \mu_t \\ \beta_t \\ \gamma_t \\ \gamma_{t-1} \\ \vdots \\ c_{1,t} \\ c_{1,t}^* \\ c_{2,t} \\ c_{2,t}^* \\ \vdots \end{bmatrix} &= \begin{bmatrix} \mathbf{T} & \mathbf{0} & \mathbf{0} & \mathbf{0} & & \\ \mathbf{0} & \mathbf{S} & \mathbf{0} & \mathbf{0} & & \\ \mathbf{0} & \mathbf{0} & \mathbf{C}_1 & \mathbf{0} & \dots & \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{C}_2 & & \\ & & \vdots & & \ddots & \end{bmatrix} \begin{bmatrix} \mu_{t-1} \\ \beta_{t-1} \\ \gamma_{t-1} \\ \gamma_{t-2} \\ \vdots \\ c_{1,t-1} \\ c_{1,t-1}^* \\ c_{2,t-1} \\ c_{2,t-1}^* \\ \vdots \end{bmatrix} + \begin{bmatrix} v_t \\ z_t \\ \omega_t \\ 0 \\ \vdots \\ \xi_{1,t} \\ \xi_{1,t}^* \\ \xi_{2,t} \\ \xi_{2,t}^* \\ \vdots \end{bmatrix} \\
 &\equiv \mathbf{F} \mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t,
 \end{aligned} \tag{44a}$$

(44b)

such that

$$\mathbf{T} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix},$$

$$\mathbf{S} = \begin{bmatrix} -1 & -1 & -1 & -1 & \dots \\ 1 & 0 & 0 & 0 & \dots \\ 0 & 1 & 0 & 0 & \dots \\ 0 & 0 & 1 & 0 & \dots \\ \vdots & \vdots & \vdots & \ddots & \ddots \end{bmatrix},$$

and $\mathbf{C}_i = \begin{bmatrix} \cos \lambda_i & \sin \lambda_i \\ -\sin \lambda_i & \cos \lambda_i \end{bmatrix}.$

This state-space representation is known in Harvey (1989, pg.172) as the *Basic Structural Model*.

2.2.11 The CAPM

Another example of the state-space modeling framework is the capital asset pricing model (CAPM) with time-varying coefficients.

Recall that the assumptions of the CAPM model imply that all investments should offer the same reward-to-risk ratio. If the ratio were better for one investment than another, investors would rearrange their portfolios towards the alternative featuring a better tradeoff. Such activity would put pressure on security prices until the ratios were equalized. Within the context of the CAPM, this ratio is known as the “Sharpe ratio” in honor of his pioneering work (Sharpe, 1966) and is defined in terms of excess returns over covariance:

$$\frac{E[R] - R_f}{Cov(R, R_m)} = \frac{E[R_m] - R_f}{\sigma_m^2}. \quad (46)$$

Of course, the Sharpe ratio directly implies a linear relationship between a) the covariance of

an asset's return with the market return; and b) the expected value of the asset's return itself:

$$E[R] - R_f = \frac{Cov(R, R_m)}{\sigma_m^2} (E[R_m] - R_f) = \beta (E[R_m] - R_f). \quad (47)$$

However, since it is clear that in the real world the assumptions of the CAPM may hold only approximately, some assets may deviate systematically from the Sharpe ratio relationship, by some amount α :

$$E[R] - R_f = \alpha + \beta (E[R_m] - R_f). \quad (48)$$

Moreover, each individual asset will be exposed to some form of idiosyncratic “micro” level risk, v , independent of what happens in the market as a whole. It is in fact this idiosyncratic risk that is minimized through the process of diversification. Therefore, we write:

$$E[r] = \alpha + \beta E[r_m] + v \quad (49)$$

where $r \equiv R - R_f$ is the observed excess return on some asset beyond the risk free rate, and $r_m \equiv R_m - R_f$ is the excess return on some market index (assumed to be completely diversified, so that it is orthogonal to the innovation or “idiosyncratic,” firm specific risk, v).

Therefore, we can treat the state transition equation as driving the dynamics of the *stochastic parameters* of the model, α_t and β_t . For example, consider the following model, given observations on r_t and $r_{m,t}$ for some $t = 1, \dots, T$ (which represents a linear regression with unobserved stochastic coefficients):

$$r_t = \alpha_t + \beta_t r_{m,t} + v_t, \quad \text{where } v_t \sim N(0, \sigma_v^2), \quad (50a)$$

$$\alpha_t = \gamma \alpha_{t-1} + u_t, \quad \text{where } u_t \sim N(0, \sigma_u^2), \quad (50b)$$

$$\text{and } \beta_t = \mu + \delta \beta_{t-1} + z_t, \quad \text{where } z_t \sim N(0, \sigma_z^2). \quad (50c)$$

Note that the nature of the equilibrium in the CAPM model suggests some reasonable restrictions on the dynamics of the stochastic parameters, α_t and β_t . First, it is safe to assume that $\frac{\mu}{1-\delta}$

will likely take on some value relatively close to, but not equal to, 1 and will depend directly on the long-run historical covariance between the asset and market returns. Moreover, δ and γ should take on values in the range $0 < x < 1$. That is, they should exhibit *mean reverting* behaviour since in the case of α_t , it is clear that arbitrage opportunities should eventually push α towards zero; and with β_t the relation between r and r_m in (47) should certainly be a bounded one.

Finally, the model in (50) can easily be put into state-space form:

$$\mathbf{x}_t \equiv \begin{bmatrix} \alpha_t \\ \beta_t \end{bmatrix} = \begin{bmatrix} 0 \\ \mu \end{bmatrix} + \begin{bmatrix} \gamma & 0 \\ 0 & \delta \end{bmatrix} \begin{bmatrix} \alpha_{t-1} \\ \beta_{t-1} \end{bmatrix} + \begin{bmatrix} u_t \\ z_t \end{bmatrix} \equiv \mathbf{c} + \mathbf{F}\mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t, \quad (51a)$$

$$\text{and } \mathbf{y}_t \equiv r_t = \begin{bmatrix} 1 & r_{m,t} \end{bmatrix} \begin{bmatrix} \alpha_t \\ \beta_t \end{bmatrix} + v_t \equiv \mathbf{H}\mathbf{x}_t + \boldsymbol{\eta}_t, \quad (51b)$$

where the covariance of the state transition equation is $\boldsymbol{\Sigma}_\epsilon \equiv \begin{bmatrix} \sigma_u^2 & 0 \\ 0 & \sigma_z^2 \end{bmatrix}$ and $\boldsymbol{\Sigma}_\eta \equiv \sigma_v^2$ is scalar. Given such a state-space representation, and observed values for both r_t and $r_{m,t}$ for all $t = 1, \dots, T$, we can now employ the techniques outlined in Section 7 to predict the values of the unobserved coefficients, α_t and β_t , across time.

Finally, the dynamic CAPM can also be interpreted as a *dynamic factor model*. Given this interpretation, the *stochastic slope coefficients* β_i now represent trending, unobserved, factors that follow their own stochastic factor dynamics. See Section 9 on “Factor models and common trends” for more details.

Of course, in the special case where \mathbf{x}_t is *observed* (through proxy) the above factor model representation is subject to Roll’s critique in that any empirical test of (51b) is really a test of the mean-variance efficiency of the proxy chosen for \mathbf{x}_t .

2.2.12 Stochastic Volatility

A popular method of modeling persistence in the second-moment of financial asset series is the ARCH model of Engle (1982):

$$y_t = \mu + \sigma_t u_t, \quad \text{where } u_t \sim N(0, 1), \quad (52a)$$

$$\text{and } \sigma_t^2 = E[\sigma_t^2 u_t^2 | Y_{t-1}] = \alpha + \beta \sigma_{t-1}^2 u_{t-1}^2. \quad (52b)$$

However, in the ARCH framework, the conditional volatility dynamics are driven in a completely deterministic fashion given past observations (that is, they are path dependent given information set $Y_t = \{y_t, y_{t-1}, \dots, y_1\}$, and the constraint imposed by (52b)). However, these second-moment dynamics may in fact be better modeled by imposing a specification that implies a strictly larger information set than Y_t . That is, we make the conditional volatility dynamics, σ_t^2 , stochastic by introducing the exogenous innovations v_t into (52b):

$$\ln(\sigma_t^2) = \ln(E[\sigma_t^2 u_t^2 | \Phi_{t-1}]) = \alpha + \beta \ln(\sigma_{t-1}^2) + v_t, \quad \text{where } v_t \sim N(0, \sigma_v^2), \quad (53)$$

and where we have taken logs to ensure positivity of the conditional volatility process.

Note that by enlarging the information set from Y_t in (52b) to $\Phi_t = \{Y_t, \sigma_t^2\}$ in (53), we are in an intuitive sense “increasing the types of questions” we can ask of our probability measure. That is, we are being more detailed about how outcomes in our probability space map to random variables in our model. However, note that the random variable σ_t^2 is latent or unobserved, and therefore, the information set we actually make inferences from will be an approximate one. In fact, it is this latent variable that makes this model amenable to state-space signal extraction methods.

Note that we can also impose a probability law on the conditional mean process, if we augment Φ_t again, to $\Phi_t = \{Y_t, \sigma_t^2, \mu_t\}$:

$$\mu_t = E[y_t | \Phi_{t-1}] = \gamma + \delta \mu_{t-1} + z_t, \quad \text{where } z_t \sim N(0, \sigma_z^2). \quad (54)$$

Therefore, the entire stochastic volatility and levels model can be written in state-space form, similar to (1), as:

$$\begin{aligned}\mathbf{y}_t &\equiv \ln(y_t^2) = \begin{bmatrix} 1 & 1 \end{bmatrix} \mathbf{x}_t + \ln(u_t^2) \\ &\equiv \mathbf{H} \mathbf{x}_t + \boldsymbol{\eta}_t,\end{aligned}\tag{55a}$$

$$\begin{aligned}\text{and } \mathbf{x}_t &\equiv \begin{bmatrix} \mu_t \\ \ln(\sigma_t^2) \end{bmatrix} = \begin{bmatrix} \gamma \\ \alpha \end{bmatrix} + \begin{bmatrix} \delta & 0 \\ 0 & \beta \end{bmatrix} \begin{bmatrix} \mu_{t-1} \\ \ln(\sigma_{t-1}^2) \end{bmatrix} + \begin{bmatrix} z_t \\ v_t \end{bmatrix} \\ &\equiv \mathbf{c} + \mathbf{F} \mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t,\end{aligned}\tag{55b}$$

where the observation equation in (52a) has been rewritten as $y_t = e^{\mu_t/2} \sigma_t u_t$ so that it is linear in logs. However, note that in this case we now have that $\boldsymbol{\eta}_t \equiv \ln(u_t^2)$ and so the observation equation innovations are not Gaussian. Therefore while the Kalman filter will represent the best unbiased *linear* predictor, it will not be as efficient as a nonlinear filtering method.

i) Factor GARCH

As a second example, we will consider augmenting the information set of the multivariate Factor GARCH model (Engle,1987). Note that this model has much in common with the material discussed in Section 9 which covers latent *dynamic factor models* and the use of principle components analysis to generate orthogonal factors.

First, consider the factor representation:

$$\mathbf{y}_t = \mathbf{B} \mathbf{f}_t + \boldsymbol{\delta}_t,\tag{56a}$$

$$\text{where } \mathbf{f}_t | \mathbf{F}_{t-1} \sim WWN(\mathbf{0}, \boldsymbol{\Omega}_t),$$

$$\boldsymbol{\delta}_t \sim WWN(\mathbf{0}, \boldsymbol{\Lambda}_\delta),$$

$$\text{and } \mathbf{F}_{t-1} = \{\mathbf{y}_{t-1}, \mathbf{y}_{t-2}, \dots, \mathbf{f}_{t-1}, \mathbf{f}_{t-2}, \dots\},$$

where \mathbf{B} is $n \times k$, and \mathbf{f}_t is $k \times 1$ where $k < n$.

From (56a) we have that the conditional covariance of \mathbf{y}_t is $\boldsymbol{\Sigma}_{\mathbf{y}_t} = \mathbf{B} \boldsymbol{\Omega}_t \mathbf{B}' + \boldsymbol{\Lambda}_\delta$. Of course,

assuming that Ω_t is diagonal we can write $\Sigma_{\mathbf{y}_t} = \sum_{i=1}^k \mathbf{b}_i \omega_{i,t} \mathbf{b}_i' + \Lambda_\delta$, where \mathbf{b}_i is the i 'th column of \mathbf{B} , and $\omega_{i,t}$ is the i 'th diagonal element of Ω_t .

In order to capture dynamic persistence in the second moment of the factors, Ω_t , we impose the following parsimonious GARCH(1,1) structure on each of the k diagonal elements:

$$\omega_{i,t} = \alpha_i \hat{f}_{i,t-1}^2 + \beta_i \omega_{i,t-1}, \quad \forall i = 1, \dots, k, \quad (57)$$

where $\hat{f}_{i,t}$ can be estimated by the i 'th element of $\mathbf{L}^{*T} \mathbf{y}_t$, and \mathbf{L}^* , $n \times k$, contains the first k columns of \mathbf{L} from the spectral decomposition of the unconditional variance of \mathbf{y}_t . That is, $\mathbf{L} \mathbf{D} \mathbf{L}' = \Sigma_{\mathbf{y}}$, so $\mathbf{L}' \Sigma_{\mathbf{y}} \mathbf{L}$ is diagonal and the k elements of $\mathbf{L}' \mathbf{y}_t$ represent an orthogonal set of random factors that account for all the variance of \mathbf{y}_t . These are the principle components (see Section 9).

Now, subbing (57) into $\Sigma_{\mathbf{y}_t} = \sum_{i=1}^k \mathbf{b}_i \omega_{i,t} \mathbf{b}_i' + \Lambda_\delta$ above yields:

$$\Sigma_{\mathbf{y}_t} = \sum_{i=1}^k \alpha_i \mathbf{b}_i \hat{f}_{i,t-1}^2 \mathbf{b}_i' + \sum_{i=1}^k \beta_i \mathbf{b}_i \omega_{i,t-1} \mathbf{b}_i' + \Lambda_\delta \quad (58a)$$

$$= \sum_{i=1}^k \alpha_i \mathbf{b}_i \left(\mathbf{l}_i^{*'} \mathbf{y}_{t-1} \right)^2 \mathbf{b}_i' + \sum_{i=1}^k \beta_i \mathbf{b}_i \left(\mathbf{l}_i^{*'} \Sigma_{\mathbf{y}_{t-1}} \mathbf{l}_i^* \right) \mathbf{b}_i' + \Lambda_\delta, \quad (58b)$$

where \mathbf{l}_i^* is the i 'th column of \mathbf{L}^* . Note that (58b) represents a first-order difference equation for $\Sigma_{\mathbf{y}_t}$ and is deterministic given the \mathbf{y}_t 's. Therefore, signal extraction methods are unnecessary as nothing is unobserved.

However, since the Factor GARCH model implies that the conditional heteroskedasticity is affecting the factors, \mathbf{f}_t , and not the innovations, δ_t , this is analogous to imposing a GARCH structure on (54) above, but where $\mu_t = z_t$ and $z_t \sim N(0, \sigma_{z,t}^2)$ is now path dependent. Of course, we could always allow for unobserved autoregressive dynamics on \mathbf{f}_t , implementing state-space framework prediction of this latent ‘‘state’’ variable and avoiding the need for principal components estimation. Another alternative would be to impose a ‘‘Factor Stochastic Volatility’’ specification, with unobserved stochastic processes driving the diagonal elements $\omega_{i,t}$.

3 Nonlinear dynamic state-space models

3.1 The nonlinear state-space model

Generally, the state-space representation requires two assumptions, namely that the process \mathbf{x}_t is Markov so that $f(\mathbf{x}_t|X_{t-1}, Y_{t-1}) = f(\mathbf{x}_t|\mathbf{x}_{t-1})$ and that the conditional distribution of \mathbf{y}_t only depends on the current value of the state, \mathbf{x}_t , or $g(\mathbf{y}_t|Y_{t-1}, X_t) = g(\mathbf{y}_t|\mathbf{x}_t)$, where $Y_t = \{\mathbf{y}_t, \dots, \mathbf{y}_0\}$ and $X_t = \{\mathbf{x}_t, \dots, \mathbf{x}_0\}$.

Therefore, the general state-space model considers the joint distribution of the process $[(\mathbf{x}'_t, \mathbf{y}'_t)']$, $l(\cdot)$:

$$l(\mathbf{y}_t, \mathbf{x}_t|X_{t-1}, Y_{t-1}) = f(\mathbf{y}_t|Y_{t-1}, X_t)g(\mathbf{x}_t|X_{t-1}, Y_{t-1}) \quad (59a)$$

$$= f(\mathbf{y}_t|\mathbf{x}_t)g(\mathbf{x}_t|\mathbf{x}_{t-1}). \quad (59b)$$

where the initial conditions of the process are defined by the marginal distribution of \mathbf{x}_0 .

3.1.1 Weak nonlinear state-space model

The weak form of the nonlinear dynamic state-space model is as follows:

$$\mathbf{x}_t = a(\mathbf{x}_{t-1}, \boldsymbol{\epsilon}_t), \quad (60a)$$

$$\mathbf{y}_t = c(\mathbf{x}_t, \boldsymbol{\eta}_t), \quad (60b)$$

with the moment restrictions:

$$E[\boldsymbol{\epsilon}_t] = \mathbf{0}, \quad Cov(\boldsymbol{\epsilon}_t, \boldsymbol{\epsilon}_{t-s}) = \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}_t} \mathbb{1}_{s=0}, \quad (2a)$$

$$E[\boldsymbol{\eta}_t] = \mathbf{0}, \quad Cov(\boldsymbol{\eta}_t, \boldsymbol{\eta}_{t-s}) = \boldsymbol{\Sigma}_{\boldsymbol{\eta}_t} \mathbb{1}_{s=0},$$

$$E[\boldsymbol{\epsilon}_{t-j} \boldsymbol{\eta}'_{t-s}] = \mathbf{0}, \quad \forall j, s \in \mathbb{Z},$$

$$\text{and } Cov(\mathbf{x}_0, \boldsymbol{\epsilon}_t) = Cov(\mathbf{x}_0, \boldsymbol{\eta}_t) = \mathbf{0}, \quad \forall t > 0,$$

where:

- \mathbf{y}_t is a $n \times 1$ vector of the observed values at time t .
- \mathbf{x}_t is a $p \times 1$ vector of state process values at time t .⁶
- $\boldsymbol{\epsilon}_t$ and $\boldsymbol{\eta}_t$ are assumed uncorrelated with each other across all time lags, and their covariance matrices, $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}_t}$ and $\boldsymbol{\Sigma}_{\boldsymbol{\eta}_t}$, depend on t .
- $a(\cdot)$ and $c(\cdot)$ are some nonlinear functions. These functions are assumed to be *non-stochastic*.
- Equation (1a) is called the “state transition” equation and (1b) is called the “observation” or “measurement” equation.
- The state initial condition, \mathbf{x}_0 , is assumed stochastic with second order moments denoted $E[\mathbf{x}_0] = \boldsymbol{\mu}$ and $Var(\mathbf{x}_0) = \boldsymbol{\Sigma}_{\mathbf{x}_0}$. Finally, there exists zero correlation between the initial state condition and the observation and state error terms, for all dates $t > 0$.

3.1.2 The Gaussian nonlinear state-space model

In the weak version of the nonlinear dynamic state-space model, the assumptions concern only the first and second-order moments of the noise processes, or equivalently the first and second-order moments of the joint process $[(\mathbf{x}'_t, \mathbf{y}'_t)']$. As in the case of the linear state-space model, we can introduce the restriction of, independent, and identically distributed, Gaussian white noises (IIN) for the errors of the state and measurement equations. The Gaussian nonlinear state space model is therefore defined as:

⁶Other terminology include the “signal” in the engineering context, “control variable” in the systems control literature, “latent factor” in the factor models approach, or “stochastic parameter” in the Bayesian context. See Section 4 for more details.

$$\mathbf{x}_t = a(\mathbf{x}_{t-1}, \boldsymbol{\epsilon}_t), \quad (60a)$$

$$\mathbf{y}_t = c(\mathbf{x}_t, \boldsymbol{\eta}_t), \quad (60b)$$

$$\begin{pmatrix} \boldsymbol{\epsilon}_t \\ \boldsymbol{\eta}_t \end{pmatrix} \sim IIN \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}_t} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{\boldsymbol{\eta}_t} \end{bmatrix} \right) \quad (62a)$$

$$E[\boldsymbol{\epsilon}_{t-j} \boldsymbol{\eta}'_{t-s}] = \mathbf{0}, \quad \forall j, s \in \mathbb{Z},$$

$$\text{with } \mathbf{x}_0 \sim N(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}_0}), \quad (62b)$$

$$\text{where } \mathbf{x}_0 \text{ and } (\boldsymbol{\epsilon}_t, \boldsymbol{\eta}_t) \text{ are independent.} \quad (62c)$$

However, when the functions $a(\cdot)$ and $c(\cdot)$ are nonlinear, under the assumption of Gaussian noise, it is no longer the case that the joint process $[(\mathbf{x}'_t, \mathbf{y}'_t)']$ is Gaussian. This implies that all marginal and conditional distributions concerning the components of these processes are also not necessarily Gaussian.

4 Terminologies

It is also interesting to note that given the widespread use of the state-space framework across different disciplines, a wide variety of interpretations have arisen regarding its implementation. For example, the examples illustrated in the previous section employ a number of different terminologies depending on the context:

- \mathbf{y}_t is equivalently referred to as the:
 - ◇ measure
 - ◇ endogenous variable

- ◇ output variable
- Likewise, \mathbf{x}_t is equivalently the:
 - ◇ state variable
 - ◇ signal
 - ◇ control variate
 - ◇ factor
 - ◇ latent variable
 - ◇ input variable
- And while, (1b) is typically called the “measurement equation,” we have that (1a) is equivalently called the:
 - ◇ state equation
 - ◇ factor dynamics

Consequently, the state-space formulation is often also referred to equivalently as the “dynamic factor model” representation, where (1a) defines the “factor dynamics” and the matrix \mathbf{H}_t is called the “factor loadings” matrix. See Section 9 for further discussion.

Moreover, the model also has the alternative Bayesian interpretation where \mathbf{x}_t represents instead a *stochastic parameter* and where the (1a) defines the dynamics of the prior distribution of the stochastic parameter. See, for example, Section 5.2.3 for more details.

Finally, with regards to the predictive problems from Section 5, we have equivalent terminologies describing the case where the state \mathbf{x}_t is unobserved. That is, $\mathcal{L}[\mathbf{x}_t | \mathbf{y}_t, \mathbf{y}_{t-1}, \dots]$, where \mathcal{L} is some optimal linear predictor, is equivalently referred to as:

- linear filtering
- linear signal extraction

- linear prediction of the latent factor
- the best linear predictor of \mathbf{x}_t given $Y_t = \{\mathbf{y}_t, \mathbf{y}_{t-1}, \dots\}$

5 The prediction problem

The problem of statistical inference can be broken into two equally important exercises. First, there is the problem of *estimation* of the theoretical model parameters, given observed random samples. Second, there is the problem of *prediction*, where we attempt to predict values of stochastic processes conditional on observed random samples, which we refer to as the *information set*. Therefore, in the first case what we are doing is drawing inferences on *fixed* parameters of the model, where in the latter case we are in fact inferring values of *random variables*. This section is exclusively concerned with the latter.

Of course, it is not always the case that we have access to random samples of the stochastic processes we have defined in our theoretical model. In this case we say that the stochastic process is *latent* or unobserved. Of course, this does not stop us from inferring the value of these stochastic processes any more than we are able to do so with those which we partially observe. In both cases, these inferences fall under the category of prediction, whether we are predicting the future value of a partially observed stochastic process, or we are predicting the value of one we never observed at all.

Therefore, given the state-space representations in Section 2, there are a number of particular prediction problems that can be formulated, which differ only in the choice of process we wish to predict (\mathbf{y}_t or \mathbf{x}_t) and the information set we have available.

The state-space models with \mathbf{x}_t unobservable imply the information set $Y_\tau \equiv \{\mathbf{y}_\tau, \mathbf{y}_{\tau-1}, \dots, \mathbf{y}_0\}$. Therefore, we say that the optimal *forecast* at time T , with horizon h , of either \mathbf{y}_{T+h} or \mathbf{x}_{T+h} employs the information set Y_T ; the optimal *filtering* of \mathbf{x}_t employs the set Y_t ; and finally the optimal *smoothing* of \mathbf{x}_t employs the set Y_T , such that $T > t$. Therefore, we have the problems:

1. Prediction of $\mathbf{y}_{T+h}|Y_T$ or $\mathbf{x}_{T+h}|Y_T$. These are the *forecasting* problems.

2. Prediction of $\mathbf{x}_t|Y_t$ for any $t = 1, \dots, T$. This is the *filtering* problem.
3. Prediction of $\mathbf{x}_t|Y_T$. This is the *smoothing* problem.

Notice that since the forecasting problem involves a choice of horizon, h , it implies an entire term structure of predictions as $h \rightarrow \infty$.

Moreover, so far we have only concerned ourselves with the *levels* of the processes. Therefore, we may also be interested in predicting how a process may behave, say in how much it varies or whether or not it tends to exhibit positive or negative outliers. Therefore, we may be interested in predicting:

1. Prediction of $(\mathbf{y}_{T+h} - \mathbf{y}_{T+h}|Y_T)^k|Y_T$ or $(\mathbf{x}_{T+h} - \mathbf{x}_{T+h}|Y_T)^k|Y_T$. These are the *forecasting* problems.
2. Prediction of $(\mathbf{x}_t - \mathbf{x}_t|Y_t)^k|Y_t$ for any $t = 1, \dots, T$. This is the *filtering* problem.
3. Prediction of $(\mathbf{x}_t - \mathbf{x}_t|Y_T)^k|Y_T$. This is the *smoothing* problem.

for any $k > 1 \in N$. Therefore, any prediction itself can be evaluated both on its “point” accuracy and its “interval” accuracy – thus we have the notions of both point and interval predictions.

5.1 Types of predictors

Of course, prediction itself can be divided into different types of predictors. Given a particular a problem we may approach the modeling exercise as an approximation, where this approximation may be some linear function, say $x\beta + \epsilon$, or some nonlinear function $f(x, \epsilon)$.

Let us now be more specific about how the prediction problem is formulated, given the general, mean zero, stochastic process $y_t = f(Y_{t-1}, \epsilon_t)$, where again Y_{t-1} represents the information set generated by the stochastic process at time $t - 1$.

It can be shown that in the general case, the optimal predictor $\hat{y}_{t|t-1} = f(Y_{t-1})$ under the minimum mean squared error (MMSE) criterion is the *conditional expectation*, $E[y_t|Y_{t-1}]$. That is, the argmin of the quadratic loss function, the expected mean squared error $E[(y_t - f(Y_{t-1}))^2]$,

Figure 3: Types of model approximations

		Problem	
		$X\beta + \varepsilon$	$f(x, \varepsilon)$
Approximation	$Xb + e$	Common	Common
	$g(x, e)$	Uncommon	Common

can be shown to be $f(Y_{t-1}) = E[y_t|Y_{t-1}]$. See see Priestley (1981, pg. 76) or Hamilton (1994, pg. 72) for proof.

However, suppose instead that we wish to restrict our choice to the class of *linear* predictors. That is, $\hat{y}_{t|t-1} = f(Y_{t-1}) = \mathbf{a}'\mathbf{y}_{t-1}^*$ where $\mathbf{y}_{t-1}^* = [y_{t-1} \ y_{t-2} \ y_{t-3} \ \dots \ y_0]'$. Interestingly, the value of \mathbf{a}' that satisfies $E[(y_t - \mathbf{a}'\mathbf{y}_{t-1}^*)\mathbf{y}_{t-1}^{*'}] = \mathbf{0}'$, that is exhibits zero covariance between forecast error and the independent variable, also minimizes the quadratic loss function above, given the constraint of having to choose amongst linear predictors.

We call this value of \mathbf{a} the *linear projection*:

$$\mathbf{a}' = E[y_t\mathbf{y}_{t-1}^{*'}]E[\mathbf{y}_{t-1}^*\mathbf{y}_{t-1}^{*'}]^{-1}, \quad (63)$$

and the definition becomes clear when we note that:

$$f(Y_{t-1}) = \mathbf{a}'\mathbf{y}_{t-1}^* = E[y_t\mathbf{y}_{t-1}^{*'}]E[\mathbf{y}_{t-1}^*\mathbf{y}_{t-1}^{*'}]^{-1}\mathbf{y}_{t-1}^* \equiv \mathcal{P}[y_t|Y_{t-1}] \quad (64)$$

and so $\mathcal{P}[y_t|Y_{t-1}]$ represents the scalar projection of y_t onto the Hilbert space spanned by the column \mathbf{y}_{t-1}^* , with the covariance norm imposed.

Interestingly, if we consider the special case of the linear, Gaussian, stochastic process with uncorrelated innovations, it turns out that the linear projection is equal to the conditional expectation since we know that for Gaussian random variables:

$$\begin{aligned} E[y_t|Y_{t-1}] &= Cov(y_t, Y_{t-1})Cov(Y_{t-1}, Y_{t-1})^{-1}Y_{t-1} \\ &\equiv \mathcal{P}[y_t|Y_{t-1}], \end{aligned} \tag{65a}$$

which is a standard result involving partitioning the set Y_t into y_t and Y_{t-1} , and without loss of generality we assume all the elements of Y_t are mean zero.

However, if the model is not linear in the stochastic variables, or if the stochastic processes are neither weak white nor Gaussian, then there is no reason to believe that the linear projection will prove the best predictor, although it is best amongst the class of linear predictors so that $MSE(\mathcal{P}[y_t|Y_{t-1}]) \geq MSE(E[y_t|Y_{t-1}])$. This is because the linear projection is a function only of the second moment properties of the stochastic processes.

Therefore, to summarize, we have three types of predictors (where one is a special case of the other):

- The linear projection, $\mathcal{P}[y_t|Y_\tau]$ $\tau \in \{1, \dots, T\}$, which is some linear function of the information set.
- The expectation, $E[y_t|Y_\tau]$, which may possibly be a nonlinear function of the information set.
- The expectation under the assumption of a linear model and Gaussianity co-incides with the linear projection.

5.2 Examples

5.2.1 Prediction of VAR(1) process (weak form)

As a simple example, consider the case where \mathbf{x}_t , the state process in (1), is observed and that the system matrices, including \mathbf{F} , are time-invariant. That is, we assume that $\mathbf{H} = \mathbf{I}$ and we know the covariance matrices $\Sigma_\epsilon \gg \mathbf{0}$ and $\Sigma_\eta = \mathbf{0}$. Therefore, the model is a VAR(1) and is a special case of the state-space model in (1) – see equation (4a).

If the eigenvalues of \mathbf{F} have modulus strictly less than 1, the VAR(1) can be inverted and written as an VMA(∞):

$$\mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t \quad (66a)$$

$$\begin{aligned} &= [\mathbf{I} - \mathbf{F}]^{-1} \boldsymbol{\epsilon}_t \\ &= \boldsymbol{\epsilon}_t + \mathbf{F}\boldsymbol{\epsilon}_{t-1} + \mathbf{F}^2\boldsymbol{\epsilon}_{t-2} + \dots \end{aligned} \quad (66b)$$

Therefore, we can consider the stochastic process as being a linear function of weak white noise, $\boldsymbol{\epsilon}_t$. That is, the information set is equivalently written as $Y_t \equiv \{\boldsymbol{\epsilon}_t, \boldsymbol{\epsilon}_{t-1}, \dots\}$.

Now, we would like to try and predict \mathbf{x}_t conditional on information set Y_{t-1} under the MMSE (minimum mean squared error) criterion. Let $f(Y_{t-1})$ represent our predictor function. We would like to minimize a quadratic loss function, the expected mean squared error $E [(\mathbf{x}_t - f(Y_{t-1})) (\mathbf{x}_t - f(Y_{t-1}))']$, where $f(Y_{t-1})$ is now a linear function of the information set at time $t-1$. The “best” (in the sense of using the information embodied in the information set efficiently, according to the MMSE criterion), linear, predictor of the VAR(1) process \mathbf{x}_t , given information set Y_{t-1} , is the *linear projection*, defined as $\mathcal{P}[\mathbf{x}_t|Y_{t-1}] \equiv \hat{\mathbf{x}}_{t|t-1}$.

However, it is clear that the linear projection must be:

$$\hat{\mathbf{x}}_{t|t-1} = \mathbf{F}\boldsymbol{\epsilon}_{t-1} + \mathbf{F}^2\boldsymbol{\epsilon}_{t-2} + \dots \quad (67)$$

since we have that:

$$E[(\mathbf{y}_t - \hat{\mathbf{x}}_{t|t-1})\mathbf{x}_{t-1}^{*'}] = \mathbf{0} \quad (68)$$

and from the discussion of the previous section we know that this value of $\hat{\mathbf{x}}_{t-1} = f(Y_{t-1})$ must be the argmin of the expected mean squared forecast error.

Note that we can also explicitly solve for the value of \mathbf{a} just as was done before, and this provides a moment condition that can be used to also *estimate* the system matrix \mathbf{F} given sample data by employing the sample moment counterpart, $\hat{\mathbf{a}}$.

Rewrite the model in (66b) as:

$$\mathbf{x}_t = \boldsymbol{\epsilon}_t + \begin{bmatrix} \mathbf{F} & \mathbf{F}^2 & \dots & \mathbf{F}^Q \end{bmatrix} \begin{bmatrix} \boldsymbol{\epsilon}_{t-1} \\ \boldsymbol{\epsilon}_{t-2} \\ \vdots \\ \boldsymbol{\epsilon}_{t-Q} \end{bmatrix} \quad (69a)$$

$$= \boldsymbol{\epsilon}_t + \mathbf{A}'\boldsymbol{\epsilon}_{t-1}^*, \quad (69b)$$

where we have truncated the lags at Q .

We can therefore solve for \mathbf{A}' as:

$$\begin{aligned} \mathbf{x}_t &= \boldsymbol{\epsilon}_t + \mathbf{A}'\boldsymbol{\epsilon}_{t-1}^* \\ \Leftrightarrow \mathbf{x}_t\boldsymbol{\epsilon}_{t-1}^{*'} &= \boldsymbol{\epsilon}_t\boldsymbol{\epsilon}_{t-1}^{*'} + \mathbf{A}'\boldsymbol{\epsilon}_{t-1}^*\boldsymbol{\epsilon}_{t-1}^{*'} \\ \Leftrightarrow E[\mathbf{x}_t\boldsymbol{\epsilon}_{t-1}^{*'}] &= E[\boldsymbol{\epsilon}_t\boldsymbol{\epsilon}_{t-1}^{*'}] + \mathbf{A}'E[\boldsymbol{\epsilon}_{t-1}^*\boldsymbol{\epsilon}_{t-1}^{*'}] \\ \Leftrightarrow \mathbf{A}' &= E[\mathbf{x}_t\boldsymbol{\epsilon}_{t-1}^{*'}]E[\boldsymbol{\epsilon}_{t-1}^*\boldsymbol{\epsilon}_{t-1}^{*'}]^{-1}. \end{aligned} \quad (70a)$$

Therefore, since the stochastic process is stationary and ergodic by assumption, we can esti-

mate the population moment condition above via the sample averages:

$$\frac{1}{T} \sum_{\tau=1}^t \mathbf{x}_\tau \boldsymbol{\epsilon}_{\tau-1}' \rightarrow^P E[\mathbf{x}_t \boldsymbol{\epsilon}_{t-1}'], \quad (71a)$$

$$\text{and } \frac{1}{T} \sum_{\tau=1}^t \boldsymbol{\epsilon}_{\tau-1}^* \boldsymbol{\epsilon}_{\tau-1}' \rightarrow^P E[\boldsymbol{\epsilon}_{t-1}^* \boldsymbol{\epsilon}_{t-1}']. \quad (71b)$$

Finally, consider the case where the state process \mathbf{x}_t in (1) is in fact *unobserved*.⁷ Now we must also predict \mathbf{x}_t conditional on the information set Y_τ (given different possible values for τ). This is where the concepts of smoothing and filtering become important – we leave the discussion of prediction given unobserved state processes to Section 7 and the estimation of the system matrices under these circumstances to Section 8 on estimation.

5.2.2 Prediction of static state-space model (strong form)

As another example, consider the case where \mathbf{x}_t , the state process in the strong form state-space model (3), is *unobserved* and that the system matrices, \mathbf{F} and \mathbf{H} , and covariance matrices, $\boldsymbol{\Sigma}_\epsilon$ and $\boldsymbol{\Sigma}_\eta$ are all time invariant. Moreover, in this case, the state process, \mathbf{x}_t , is i.i.d. Gaussian and does not exhibit serial correlation. We call this model the *static* state-space model and it is a special case of the strong form model in (3) – see formula (28) above.

From the model in (28) we know that:

$$\begin{bmatrix} \mathbf{x}_t \\ \mathbf{y}_t \end{bmatrix} = MVN \left[\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_\epsilon & \boldsymbol{\Sigma}_\epsilon \mathbf{H}' \\ \mathbf{H} \boldsymbol{\Sigma}_\epsilon & \mathbf{H} \boldsymbol{\Sigma}_\epsilon \mathbf{H}' + \boldsymbol{\Sigma}_\eta \end{pmatrix} \right]. \quad (72)$$

Moreover, from the discussion above, we know that given the linear Gaussian setting, the linear projection is equivalent to conditional expectation. Therefore, the best linear predictors of

⁷For example if $\boldsymbol{\Sigma}_\eta \gg \mathbf{0}$.

both observed and unobserved processes are given as:

$$E[\mathbf{y}_t|Y_{t-1}] = \mathbf{H}E[\mathbf{x}_t|Y_{t-1}] + \mathbf{0} = \mathbf{0}, \quad (73a)$$

$$\text{and } E[\mathbf{x}_t|Y_{t-1}] = \mathbf{0}. \quad (73b)$$

However consider that the expectation of the state variable, conditional on current time t information, is not trivial:

$$E[\mathbf{x}_t|Y_t] = E[\mathbf{x}_t\mathbf{y}_t']E[\mathbf{y}_t\mathbf{y}_t']^{-1}\mathbf{y}_t \quad (74a)$$

$$= E[\mathbf{x}_t(\mathbf{H}\mathbf{x}_t + \boldsymbol{\eta}_t)'] [\mathbf{H}\boldsymbol{\Sigma}_\epsilon\mathbf{H}' + \boldsymbol{\Sigma}_\eta]^{-1}\mathbf{y}_t \quad (74b)$$

$$= \boldsymbol{\Sigma}_\epsilon\mathbf{H}' [\mathbf{H}\boldsymbol{\Sigma}_\epsilon\mathbf{H}' + \boldsymbol{\Sigma}_\eta]^{-1}\mathbf{y}_t. \quad (74c)$$

Note the difference between this example and that given in Section 5.2.1. In the previous case the state process \mathbf{x}_t was *observable*, and so our prediction formula was really a prediction of the observation process. However, in this case it is not.

Therefore we are confronted with two problems: a) predict future values of the observed process \mathbf{y}_t given information set Y_{t-1} and also b) predict the state process \mathbf{x}_t given the information set Y_τ where τ can range across different values. If $\tau = t$, we call this the filtering problem and if $\tau > t$ we call this the smoothing problem. Under the static model framework, the solution is simple. It is simply the linear projection $\mathcal{P}[\mathbf{x}_t|Y_\tau]$ above. However, when the state process is not static, the problem becomes much more involved and involves updating the linear projection given new information. This is the essence of the Kalman filter and fixed point smoother algorithms discussed in Section 7. In fact, the Kalman filter solutions below collapse to the expressions above when the state process is static. For another example of a static state-space model, see the static Factor model of Section 9.

5.2.3 Bayesian interpretation of the state-space model

Within the Bayesian context we can always interpret the state transition equation (1a) as defining the dynamics of a prior distribution on the stochastic parameter \mathbf{x}_t . For example, from the linear state-space model in Section 2.1.2, if we view the “signal” \mathbf{x}_t to be a *stochastic parameter*, then the model is more akin to a linear regression of \mathbf{H}_t regressed on \mathbf{y}_t but where the slope coefficient \mathbf{x}_t is allowed to follow an unobserved stochastic process.

Given this interpretation we can work with the probability density functions implied by the respective model, to derive the *posterior* density $p(\mathbf{x}_t|Y_t)$, where $Y_t = \{\mathbf{y}_t, \mathbf{y}_{t-1}, \dots, \mathbf{y}_1\}$ (the posterior density is therefore the filtered density of the latent state, \mathbf{x}_t). This will be done by appealing to Bayes theorem, and in the process we will make use of the prior density $p(\mathbf{x}_t|Y_{t-1})$ and the likelihood $p(\mathbf{y}_t|\mathbf{x}_t)$, given observed values \mathbf{y}_t and stochastic parameter \mathbf{x}_t .

First, using Bayes theorem it can be shown that the posterior distribution $p(\mathbf{x}_t|Y_t)$ is derived as:

$$\begin{aligned}
 p(\mathbf{x}_t|Y_t) &= \frac{p(\mathbf{x}_t, Y_t)}{p(Y_t)} = \frac{p(\mathbf{x}_t, \mathbf{y}_t, Y_{t-1})}{p(\mathbf{y}_t, Y_{t-1})} \\
 \text{where } p(\mathbf{x}_t, \mathbf{y}_t, Y_{t-1}) &= p(\mathbf{y}_t|\mathbf{x}_t, Y_{t-1})p(\mathbf{x}_t|Y_{t-1})p(Y_{t-1}) \\
 &= p(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_t|Y_{t-1})p(Y_{t-1}) \\
 \Leftrightarrow p(\mathbf{x}_t|Y_t) &= \frac{p(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_t|Y_{t-1})}{p(\mathbf{y}_t|Y_{t-1})} \\
 &\propto p(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_t|Y_{t-1}) \tag{75a}
 \end{aligned}$$

Therefore, given the assumption of a linear model, from Section 2.1.2, with Gaussian inno-

vations, ϵ_t and η_t , we have that the prior density is defined as:

$$p(\mathbf{x}_t|Y_{t-1}) \sim MVN(\hat{\mathbf{x}}_{t|t-1}, \mathbf{P}_{t|t-1}) \quad (76a)$$

$$\text{where } \hat{\mathbf{x}}_{t|t-1} = \mathbf{F}\hat{\mathbf{x}}_{t-1|t-1} \quad (76b)$$

$$\text{and } \mathbf{P}_{t|t-1} = \mathbf{F}\mathbf{P}_{t-1|t-1}\mathbf{F}' + \Sigma_\epsilon, \quad (76c)$$

which follows directly from the linear projection of $\mathbf{F}\mathbf{x}_{t-1} + \epsilon_t$ into the space spanned by Y_{t-1} [see (89a) and (100a) below]. Therefore, $\hat{\mathbf{x}}_{t|t-1}$ is the forecast of the state variable \mathbf{x}_t , given information at time $t - 1$, which is a linear function of $\hat{\mathbf{x}}_{t-1|t-1}$, the filtered state variable.

Moreover, the likelihood given observed values \mathbf{y}_t is given as:

$$p(\mathbf{y}_t|\mathbf{x}_t) \sim MVN(\mathbf{H}\mathbf{x}_t, \Sigma_\eta) \quad (77)$$

which follows immediately from the definition of the state-space model in (1). Therefore, given the linear model, we have that the forecast of \mathbf{y}_t conditional on time $t - 1$ information is therefore $\mathbf{H}\hat{\mathbf{x}}_{t|t-1}$, which again follows from the linear projection.

Furthermore, by standard results on the moments of the product of two normal densities, we have that the posterior distribution is proportional to:

$$p(\mathbf{x}_t|Y_t) \sim MVN(\hat{\mathbf{x}}_{t|t}, \mathbf{P}_{t|t}) \quad (78a)$$

$$\text{where } \hat{\mathbf{x}}_{t|t} = \mathbf{P}_{t|t} \left(\mathbf{P}_{t|t-1}^{-1} \hat{\mathbf{x}}_{t|t-1} + \mathbf{H}'\Sigma_\eta^{-1}\mathbf{y}_t \right) \quad (78b)$$

$$\text{and } \mathbf{P}_{t|t} = \left[\mathbf{H}'\Sigma_\eta^{-1}\mathbf{H} + \mathbf{P}_{t|t-1}^{-1} \right]^{-1} \quad (78c)$$

where it can be shown that (78b) and (78c) are equal to (98b) and (100b), respectively. Therefore, the above expressions provide a recursive system of equations that together generate the filtered prediction of the state variable $\hat{\mathbf{x}}_{t|t}$ [to see this sub (76b) into (78b) and (76c) into (78c)]. This set of recursions is effectively the Kalman filter (1960) and so we can see that from the Bayesian framework the Kalman filter arises under the special case of the linear state-space model with

Gaussian innovations.

However, if the model is nonlinear we cannot expect the conditional expectation to be of a linear form. Generally, we must approximate the set of filter recursions by numerical methods such as Monte-Carlo. Essentially, we are able to generate the expectation of the filtered state density in (75a) by drawing from it and appealing to a law of large numbers.

Harrison and Stevens (1976) presents the general discussion of the Bayesian approach to state-space modeling.

6 Prediction in the frequency domain

Consider the engineering context, where the classical approach to linear time invariant system analysis involved spectral representations and “frequency response” to unit shocks passed through appropriate physical filters, the state, x_t , has kept its interpretation as the “signal” corrupted by noise. Since the signal is latent, we often wish to “extract” the signal from the noise and this defines the “signal extraction” problem:

$$\operatorname{argmin}_{\{a_j\}_{j=0}^{\infty}} E[(x_t - \hat{x}_t)^2] \quad (79a)$$

$$\text{where } \hat{x}_t = \sum_{j=0}^{\infty} a_j y_{t-j}, \quad (79b)$$

$$\text{and } y_t = x_t + \eta_t, \quad (79c)$$

and where η_t is a weak white noise, and we only observe the noisy signal y_t . Moreover, x_t is uncorrelated with the noise η_t , and both processes are stationary. That is, we wish to “filter” out the noise, η_t , to recover the signal, x_t , when we only observe, y_t . Under the MMSE criterion, the optimal solution will minimize (79a), by choosing the sequence of filter coefficients $\{a_j\}_{j=0}^{\infty}$.

That is, given observations of the output $Y_t = \{y_t, y_{t-1}, y_{t-2}, \dots\}$, we wish to find the most efficient predictor of the latent variable x_t in the sense of the discussion provided above in Section 5. The most famous solution, in the frequency domain, is the Wiener filter developed during the

mid-part of the last century. This solution is interpreted in terms of an optimal *transfer function*, $A(e^{-i\omega})$, which is the Fourier transform of the impulse response function, $\{a_j\}_{j=0}^{\infty}$:

$$A(e^{-i\omega}) = \frac{h_{xy}(\omega)}{h_{yy}(\omega)} = \frac{h_{xx}(\omega)}{h_{xx}(\omega) + h_{\eta\eta}(\omega)} \quad (80a)$$

$$\text{where } A(z) = \sum_{j=-\infty}^{\infty} a_j z^j, \quad (80b)$$

and where $h_{xy}(\omega)$ is the cross-spectral density between x_t and y_t , and $h_{xx}(\omega)$ is the spectral density of x_t . Therefore, we can see that the Wiener filter solution allows frequencies to pass through when the signal-to-noise ratio is high, and vice versa when low:

$$h_{\hat{x}\hat{x}} = |A(e^{-i\omega})|^2 h_{yy}(\omega). \quad (81)$$

See Priestley (1981), Chapter 10, for more details.

However, given the advent of digital computers, the more useful time domain counterpart is found in the work of Kalman (1960) and is discussed in more detail within Section 7.

Subsequently, the application of signal extraction methods to modern economics, and empirical time-series analysis more generally, was developed by Nerlove et al. (1979) and popularized by Harvey (1984,89) and the “Unobserved Components” framework which is discussed in Section 9.4 (see also example 2.2.10).

7 Prediction in the time domain

Let us turn now to the general solutions of the following prediction problems in the time domain, when the state-space representation is constrained to be linear:

1. Forecasting future values of \mathbf{y}_{T+h} for some horizon h given observed values from $t = 1, \dots, T$, or forecasting \mathbf{x}_{T+h} given predictions of the state $\hat{\mathbf{x}}_{t|t}, t = 1, \dots, T$.
2. Filtering \mathbf{x}_τ given observed values, $\mathbf{y}_t, t = 1, \dots, \tau \leq T$.

3. Smoothing x_τ given observed values, $\mathbf{y}_t, t = 1, \dots, T$, where $\tau < T$.

The model in (1) implies the information set $Y_t \equiv \{\mathbf{y}_t, \mathbf{y}_{t-1}, \dots, \mathbf{y}_1\}$. Therefore, the optimal forecast of either \mathbf{y}_t or x_t employs the information set Y_{t-1} ; the optimal filtering of x_t employs the set Y_t ; and finally the optimal smoothing of x_t employs the set Y_T .

It is of historical interest that the first of the three problems was solved almost simultaneously by both Wiener and Kolmogorov around the middle of the 20th century. However, while Kolmogorov worked in the time domain, Wiener solved the equivalent problem in the frequency domain [see Chapter 10 in Priestley (1981)]. Interestingly, the filtering problem is a special case of the smoother which was solved by Wiener in terms of an optimal transfer function (and frequency response) but where we have access to a doubly infinite set of observations on \mathbf{y}_t (i.e. our information set is $Y_\infty \equiv \{\dots, \mathbf{y}_{t+1}, \mathbf{y}_t, \mathbf{y}_{t-1}, \dots\}$) [see Section 6 above]. It was Kalman's celebrated work in (1960) that rephrased the three problems in terms of the state-space formulation in (1), allowing one to process the filtered state values through a recursive algorithm ideally suited to implementation on digital computers (as opposed to in an analog fashion using physically constructed electrical filters).

It can be shown that as $T \rightarrow \infty$ the steady-state solution to the Kalman state-space smoother is equivalent to the Wiener solution. However, the Kalman filter is superior in a number of ways: first, we need not work in the frequency domain; second, the Kalman optimal filter weights are easily updated at each time t , while the Wiener weights must be completely recomputed for each new time period; and finally, the Kalman filter allows us to easily extend the solution to the multivariate case, or to the linear time-variant system case.

The notions of forecasting, filtering, and smoothing are not as different as they first seem. All three involve determining the optimal prediction of some random variable, say Z_t , given a particular information set \mathcal{F}_τ . Therefore, despite the misleading nomenclature, forecasting, filtering and smoothing are really just methods to predict stochastic processes (whether observable or not). In fact, it is only the information set available to us that determines which of the three problems we face. Therefore, all three are examples of the problem of statistical inference of

random variables.

7.1 Linear projections

Recall from Section 5 that if we adopt a quadratic loss function (e.g. minimum mean squared error (MMSE)) criterion, then in general, the argmin function $f(Y_t)$ is the conditional expectation, given the relevant information set [see Section 5.1 and Priestley (1981, pg. 76) for a formal proof]. However, within the linear context, where we adopt a Hilbert space framework, with inner product defined as the covariance between random variables, the best linear unbiased estimator is the *linear projection* (also called the “orthogonal projector”).

Formally, for the first problem, given some $\mathcal{H}_T \subset \mathcal{H}_{T+h}$ where \mathcal{H} is a Hilbert space, let $\mathbf{y}_{T+h} \in \mathcal{H}_{T+h}$. It can be shown that there exists a unique $\hat{\mathbf{y}}_{T+h} \in \mathcal{H}_T$ where $\|\mathbf{y}_{T+h} - \hat{\mathbf{y}}_{T+h}\|^2$ is minimized, and where $\hat{\mathbf{y}}_{T+h}$ is uniquely determined by the property $(\mathbf{y}_{T+h} - \hat{\mathbf{y}}_{T+h}) \perp \forall \mathbf{y} \in \mathcal{H}_T$. Therefore, $\hat{\mathbf{y}}_{T+h}$ represents the *orthogonal projector*.

For the second problem, given some $\mathcal{L}_t \subseteq \mathcal{L}_T$ where \mathcal{L} is another Hilbert space (orthogonal to \mathcal{H}), let $\mathbf{x}_t \in \mathcal{L}_t$. It can be shown that there exists a unique $\hat{\mathbf{x}}_t \in \mathcal{H}_t$ where $\|\mathbf{x}_t - \hat{\mathbf{x}}_t\|^2$ is minimized, and where $\hat{\mathbf{x}}_t$ is uniquely determined by the property $(\mathbf{x}_t - \hat{\mathbf{x}}_t) \perp \forall \mathbf{y} \in \mathcal{H}_t$.

Finally, the third problem suggests that given some $\mathcal{L}_t \subseteq \mathcal{L}_T$, let $\mathbf{x}_t \in \mathcal{L}_t$. It can be shown that there exists a unique $\hat{\mathbf{x}}_t \in \mathcal{H}_T$ where $\|\mathbf{x}_t - \hat{\mathbf{x}}_t\|^2$ is minimized, and where $\hat{\mathbf{x}}_t$ is uniquely determined by the property $(\mathbf{x}_t - \hat{\mathbf{x}}_t) \perp \forall \mathbf{y} \in \mathcal{H}_T$.

7.2 The lag representation

The time-invariant version of the linear dynamic state-space model in (1) is obtained when the system matrices, \mathbf{F} and \mathbf{H} , and the noise covariance matrices, Σ_ϵ and Σ_η , do not depend on time. Therefore, we get:

$$\mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \epsilon_t, \quad (82a)$$

$$\text{and } \mathbf{y}_t = \mathbf{H}\mathbf{x}_t + \eta_t, \quad (82b)$$

$$\text{where } \begin{pmatrix} \boldsymbol{\epsilon}_t \\ \boldsymbol{\eta}_t \end{pmatrix} \sim WWN \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_\epsilon & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_\eta \end{bmatrix} \right). \quad (83a)$$

This time-invariant system can be written in terms of the lag operator (L). The operator is defined on the set of second-order processes.⁸ It transforms a given process (\boldsymbol{x}_t) into the process ($L\boldsymbol{x}_t$), where $(L\boldsymbol{x}_t) = (\boldsymbol{x}_{t-1})$; thus its components are deduced by drifting time by one lag. The operator is especially important in characterizing the second-order stationary process. Indeed, the process (\boldsymbol{x}_t) is second-order stationary if and only if the first and second-order moments of (\boldsymbol{x}_t) and ($L\boldsymbol{x}_t$) are the same for all t .

The system in (82) can be written as:

$$\boldsymbol{x}_t = \boldsymbol{F}L\boldsymbol{x}_t + \boldsymbol{\epsilon}_t \quad \text{and} \quad \boldsymbol{y}_t = \boldsymbol{H}\boldsymbol{x}_t + \boldsymbol{\eta}_t, \quad (84a)$$

$$\Leftrightarrow (\boldsymbol{I} - \boldsymbol{F}L)\boldsymbol{x}_t = \boldsymbol{\epsilon}_t \quad \text{and} \quad \boldsymbol{y}_t = \boldsymbol{H}\boldsymbol{x}_t + \boldsymbol{\eta}_t. \quad (84b)$$

If the operator $\boldsymbol{I} - \boldsymbol{F}L$ is invertible (i.e. the eigenvalues of \boldsymbol{F} are strictly less than 1 in absolute value), then we can write (84b) as:

$$\boldsymbol{x}_t = (\boldsymbol{I} - \boldsymbol{F}L)^{-1} \boldsymbol{\epsilon}_t. \quad (85)$$

Therefore, from both (84b) and (82b) the model can be written purely in terms of the observable variable \boldsymbol{y}_t and the error terms as:

$$\boldsymbol{y}_t = \boldsymbol{H}(\boldsymbol{I} - \boldsymbol{F}L)^{-1} \boldsymbol{\epsilon}_t + \boldsymbol{\eta}_t. \quad (86)$$

That is, if the model is stationary, we can always rewrite the state-space representation in terms of an *orthogonal basis*.

⁸That is, the lag operator, L , represents a mapping $L : \mathbb{R}^{Pt} \rightarrow \mathbb{R}^{P(t-1)}$ or $L^{-1} : \mathbb{R}^{Pt} \rightarrow \mathbb{R}^{P(t+1)}$ where $t \in \mathbb{Z}$. Moreover, the composition of lag operators implies that $LL^{-1} = 1$, since $L \circ L^{-1} \equiv LL^{-1} : \mathbb{R}^{Pt} \rightarrow \mathbb{R}^{P(t+1)} \rightarrow \mathbb{R}^{Pt}$. More generally of course, $L^k : \mathbb{R}^{Pt} \rightarrow \mathbb{R}^{P(t-k)}$.

7.3 Forecasting

Let us consider the model in (1), but with time invariant system and noise covariance matrices. That is, we write the system and observation matrices, \mathbf{F} and \mathbf{H} , without their time subscripts and the covariance matrices, Σ_η and Σ_ϵ , also do not depend on time – see (82). From (1b) we have that the optimal linear forecast of future values of \mathbf{y}_{T+h} for some horizon h given observed values from $t = 1, \dots, T$ is derived as:

$$\hat{\mathbf{y}}_{T+h} = \mathcal{P}[\mathbf{y}_{T+h}|Y_T] = \mathbf{H}\mathcal{P}[\mathbf{x}_{T+h}|Y_T] + \mathbf{0} \quad (87a)$$

$$= \mathbf{H} (\mathbf{F}\mathcal{P}[\mathbf{x}_{T+(h-1)}|Y_T] + \mathbf{0})$$

$$\vdots$$

$$= \mathbf{H}\mathbf{F}^h\mathcal{P}[\mathbf{x}_T|Y_T] \quad (87b)$$

where $\mathcal{P}[\cdot]$ denotes linear projection, and where $\mathcal{P}[\mathbf{y}_{T+h}|Y_T]$ represents the optimal *orthogonal projector* since we have that for any $\mathbf{y}_T = \mathbf{H}\mathbf{x}_T + \boldsymbol{\eta}_T$, $(\mathbf{y}_{T+h} - \hat{\mathbf{y}}_{T+h}) \perp \mathbf{y}_T$, according to the covariance norm.

Moreover, suppose we have already solved for the filtered values $\mathcal{P}[\mathbf{x}_T|\mathbf{y}_T, \mathbf{y}_{T-1}] \equiv \hat{\mathbf{x}}_{T|T}$, we have that the optimal forecast of \mathbf{x}_{T+h} given estimates of the state \mathbf{x}_t , $t = 1, \dots, T$, is derived as:

$$\hat{\mathbf{x}}_{T+h} = \mathcal{P}[\mathbf{x}_{T+h}|Y_T] = \mathbf{F}\mathcal{P}[\mathbf{x}_{T+(h-1)}|Y_T] + \mathbf{0} \quad (88a)$$

$$= \mathbf{F} (\mathbf{F}\mathcal{P}[\mathbf{x}_{T+(h-2)}|Y_T] + \mathbf{0})$$

$$\vdots$$

$$= \mathbf{F}^h\mathcal{P}[\mathbf{x}_T|Y_T]$$

$$= \mathbf{F}^h\mathcal{P}[\mathbf{x}_T|\mathbf{y}_T, \mathbf{y}_{T-1}] \equiv \mathbf{F}^h\hat{\mathbf{x}}_{T|T}. \quad (88b)$$

7.4 Filtering

The second prediction problem above is solved by establishing an expression for $\mathcal{P}[\mathbf{x}_t|Y_t]$, for any $t \in \{1, \dots, T\}$. That is, we wish to filter \mathbf{x}_τ given observed values, $\mathbf{y}_t, t = 1, \dots, \tau \leq T$. For simplicity of notation, let $\mathcal{P}[\mathbf{x}_t|Y_t] \equiv \hat{\mathbf{x}}_{t|t}$. The Kalman filter provides a direct way to compute $\hat{\mathbf{x}}_{t|t}$ recursively given any starting values assumed for $\hat{\mathbf{x}}_{1|0}$ and $\mathbf{P}_{1|0}$ and given values for the system and noise covariance matrices, Σ_ϵ and Σ_η .

First note that from (1a):

$$\hat{\mathbf{x}}_{t+1|t} = \mathcal{P}[\mathbf{F}\mathbf{x}_t + \epsilon_{t+1}|Y_t] = \mathbf{F}\hat{\mathbf{x}}_{t|t}, \quad (89a)$$

$$\text{and } \hat{\mathbf{y}}_{t+1|t} = \mathcal{P}[\mathbf{H}\mathbf{x}_{t+1}|Y_t] = \mathbf{H}\hat{\mathbf{x}}_{t+1|t}. \quad (89b)$$

Let:

$$\mathbf{f}_{t+1} = \mathbf{x}_{t+1} - \hat{\mathbf{x}}_{t+1|t}, \quad (90a)$$

$$\text{and } \mathbf{e}_{t+1} = \mathbf{y}_{t+1} - \hat{\mathbf{y}}_{t+1|t}, \quad (90b)$$

be the one-step ahead state and observation forecast errors, respectively. We can now define $\text{Var}[\mathbf{x}_t|\mathbf{y}_{t-1}] = E[\mathbf{f}_t\mathbf{f}_t'] \equiv \mathbf{P}_{t|t-1}$.

Our goal then is to update the linear projection $\mathcal{P}[\mathbf{x}_t|Y_{t-1}]$ with new information, as it arises at time t . It is shown in Section 7.4.1 below, that an expression for updating a linear projection $\mathcal{P}[\mathbf{x}_t|Y_{t-1}]$ with time t information is given as:

$$\mathcal{P}[\mathbf{x}_t|\mathbf{y}_t, \mathbf{y}_{t-1}] = \mathcal{P}[\mathbf{x}_t|\mathbf{y}_{t-1}] + E[\mathbf{f}_t\mathbf{e}_t']E[\mathbf{e}_t\mathbf{e}_t']^{-1}\mathbf{e}_t \quad (98a)$$

$$\equiv \hat{\mathbf{x}}_{t|t} = \hat{\mathbf{x}}_{t|t-1} + \mathbf{K}_t\mathbf{e}_t \quad (98b)$$

where \mathbf{K}_t is given from (99b), below, as:

$$\mathbf{K}_t = \mathbf{P}_{t|t-1}\mathbf{H}' [\mathbf{H}\mathbf{P}_{t|t-1}\mathbf{H}' + \Sigma_\eta]^{-1} \quad (99b)$$

and where from (99c), below, we have that:

$$\text{Var}[\mathbf{y}_t|\mathbf{y}_{t-1}] = E[\mathbf{e}_t\mathbf{e}_t'] = \Sigma_{e_t} = \mathbf{H}\mathbf{P}_{t|t-1}\mathbf{H}' + \Sigma_\eta. \quad (99c)$$

Finally from plugging (100b) into (100a) we get an expression for the conditional variance of \mathbf{x}_t in terms of a matrix difference equation. Equations (100b) and (100a), are introduced later but repeated here as

$$\begin{aligned} \text{Var}[\mathbf{x}_t|\mathbf{y}_{t-1}] &= E[\mathbf{f}_t\mathbf{f}_t'] = \mathbf{P}_{t|t-1} \\ &= \mathbf{F}\text{Var}[\mathbf{x}_{t-1}|\mathbf{y}_{t-1}]\mathbf{F}' + \Sigma_\epsilon = \mathbf{F}\mathbf{P}_{t-1|t-1}\mathbf{F}' + \Sigma_\epsilon, \end{aligned} \quad (100a)$$

$$\text{and } \mathbf{P}_{t|t} = \mathbf{P}_{t|t-1} - \mathbf{K}_t\Sigma_{e_t}\mathbf{K}_t'. \quad (100b)$$

Together they imply that:

$$\begin{aligned} \mathbf{P}_{t|t-1} &= \mathbf{F}\mathbf{P}_{t-1|t-1}\mathbf{F}' + \Sigma_\epsilon \\ &= \mathbf{F} [\mathbf{P}_{t-1|t-2} - \mathbf{K}_{t-1}\Sigma_{e_{t-1}}\mathbf{K}_{t-1}'] \mathbf{F}' + \Sigma_\epsilon \\ &= \mathbf{F} [\mathbf{P}_{t-1|t-2} - \mathbf{K}_{t-1} (\mathbf{H}\mathbf{P}_{t-1|t-2}\mathbf{H}' + \Sigma_\eta) \mathbf{K}_{t-1}'] \mathbf{F}' + \Sigma_\epsilon \\ &= \mathbf{F} [\mathbf{P}_{t-1|t-2} - \mathbf{P}_{t-1|t-2}\mathbf{H}'\mathbf{K}_{t-1}'] \mathbf{F}' + \Sigma_\epsilon \\ &= \mathbf{F}\mathbf{P}_{t-1|t-2} (\mathbf{F}(\mathbf{I} - \mathbf{K}_{t-1}\mathbf{H}))' + \Sigma_\epsilon \\ &= \mathbf{F}\mathbf{P}_{t-1|t-2}\mathbf{L}_{t-1}' + \Sigma_\epsilon \end{aligned} \quad (93a)$$

which is known as the discrete time algebraic Ricatti equation (ARE). Given sufficient conditions the ARE can be solved for a steady-state value which is often useful in embedded systems where computational memory is limited [see Simon (2006, pg.194-199)]. Note that a steady-state value for \mathbf{P}_∞ implies a steady-state value for \mathbf{K}_∞ .

Therefore, to summarize we have the following expressions that together suggest a recursive algorithm for computing the linear filtered solution $\mathcal{P}[\mathbf{x}_t|Y_t] \equiv \hat{\mathbf{x}}_{t|t}$ for each $t = 1, \dots, T$:

$$\hat{\mathbf{x}}_{t+1|t} = \mathbf{F}\hat{\mathbf{x}}_{t|t}, \quad (89a)$$

$$\hat{\mathbf{x}}_{t|t} = \hat{\mathbf{x}}_{t|t-1} + \mathbf{K}_t \mathbf{e}_t, \quad (98a)$$

$$\hat{\mathbf{y}}_{t+1|t} = \mathbf{H}\hat{\mathbf{x}}_{t+1|t}, \quad (89b)$$

$$\mathbf{f}_{t+1} = \mathbf{x}_{t+1} - \hat{\mathbf{x}}_{t+1|t}, \quad (90a)$$

$$\mathbf{e}_{t+1} = \mathbf{y}_{t+1} - \hat{\mathbf{y}}_{t+1|t}, \quad (90b)$$

$$\mathbf{K}_t = \mathbf{P}_{t|t-1} \mathbf{H}' \Sigma_{e_t}^{-1}, \quad (99b)$$

$$\Sigma_{e_t} = \mathbf{H} \mathbf{P}_{t|t-1} \mathbf{H}' + \Sigma_{\eta}, \quad (99c)$$

$$\mathbf{P}_{t|t-1} = \mathbf{F} \mathbf{P}_{t-1|t-2} \mathbf{L}'_{t-1} + \Sigma_{\epsilon}, \quad (93a)$$

$$\text{and } \mathbf{L}_{t-1} = \mathbf{F} (\mathbf{I} - \mathbf{K}_{t-1} \mathbf{H}). \quad (94a)$$

Kalman filter recursive algorithm:

1. Starting with $\mathbf{P}_{1|0}$, compute equation (99c) and (99b) to get \mathbf{K}_1 .
2. Plug \mathbf{K}_1 into (93a) to obtain $\mathbf{P}_{2|1}$.
3. Plug $\hat{\mathbf{x}}_{1|0}$ into (89b) to obtain $\hat{\mathbf{y}}_{1|0}$.
4. Plug $\hat{\mathbf{y}}_{1|0}$ into (90b) to obtain \mathbf{e}_1 .
5. Finally, plug \mathbf{K}_1 and \mathbf{e}_1 into (98a) to obtain the filtered value $\hat{\mathbf{x}}_{1|1}$.
6. Plugging $\hat{\mathbf{x}}_{1|1}$ into (89a) then yields $\hat{\mathbf{x}}_{2|1}$. Since we already have $\mathbf{P}_{2|1}$ from step (2) above, we can now continue from step (1), ultimately repeating all the steps until we solve for $\hat{\mathbf{x}}_{t|t}$ for some desired $t \leq T$.

7.4.1 Kalman filter in terms of orthogonal basis

Note that an intuitive way to approach the solution to $\hat{\mathbf{x}}_{t|t}$ is to write it in terms of an orthogonal basis, where the basis vectors are the forecast errors $\mathbf{e}_\tau = \mathbf{y}_\tau - \hat{\mathbf{y}}_{\tau|\tau-1}$, $\tau = 1, \dots, t$. This

allows us to appreciate how the Kalman filter is really the time domain analog to the frequency domain solution which attenuates frequencies where the signal-to-noise ratio is low.

In order to construct the appropriate orthogonal basis, we start with the requirement that the innovations $\mathbf{u}_t = \mathbf{x}_t - \mathcal{P}[\mathbf{x}_t|\mathbf{y}_t, \mathbf{y}_{t-1}]$, $\mathbf{e}_t = \mathbf{y}_t - \mathcal{P}[\mathbf{y}_t|\mathbf{y}_{t-1}]$, and finally \mathbf{y}_{t-1} itself should all be uncorrelated with each other. Since their joint covariance matrix is therefore diagonal this implies that:

$$E[\mathbf{b}\mathbf{b}'] = \mathbf{D}, \quad \text{where } \mathbf{b}' = \begin{bmatrix} \mathbf{y}'_{t-1} & \mathbf{e}'_t & \mathbf{u}'_t \end{bmatrix}. \quad (95a)$$

$$\text{Let } \mathbf{b}^* \equiv \mathbf{A}\mathbf{b}, \quad \text{where } \mathbf{A} = \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{A}_{21} & \mathbf{I} & \mathbf{0} \\ \mathbf{A}_{31} & \mathbf{A}_{32} & \mathbf{I} \end{bmatrix} \quad (95b)$$

$$\Leftrightarrow E[\mathbf{b}^*\mathbf{b}^{*'}] = \mathbf{A}\mathbf{D}\mathbf{A}' = \mathbf{\Omega} = \begin{bmatrix} \mathbf{\Omega}_{11} & \mathbf{\Omega}_{12} & \mathbf{\Omega}_{13} \\ \mathbf{\Omega}_{21} & \mathbf{\Omega}_{22} & \mathbf{\Omega}_{23} \\ \mathbf{\Omega}_{31} & \mathbf{\Omega}_{32} & \mathbf{\Omega}_{33} \end{bmatrix}. \quad (95c)$$

It can be shown that:

$$\mathbf{\Omega}' = \mathbf{\Omega}, \quad (96a)$$

$$\mathbf{A}_{i1} = \mathbf{\Omega}_{i1}\mathbf{\Omega}_{11}^{-1}, \quad \forall i = \{2, 3\}, \quad (96b)$$

$$\text{and } \mathbf{A}_{32} = (\mathbf{\Omega}_{32} - \mathbf{\Omega}_{31}\mathbf{\Omega}_{11}^{-1}\mathbf{\Omega}_{12}) (\mathbf{\Omega}_{22} - \mathbf{\Omega}_{21}\mathbf{\Omega}_{11}^{-1}\mathbf{\Omega}_{12})^{-1}. \quad (96c)$$

The first row of $\mathbf{A}\mathbf{b}$ is $\mathbf{y}_{t-1} = \mathbf{y}_{t-1}$; however, the second row is $\mathbf{\Omega}_{21}\mathbf{\Omega}_{11}^{-1}\mathbf{y}_{t-1} + \mathbf{e}_t = \mathbf{y}_t$, since $\mathbf{\Omega}_{21}\mathbf{\Omega}_{11}^{-1}\mathbf{y}_{t-1} = \mathcal{P}[\mathbf{y}_t|\mathbf{y}_{t-1}]$ and $E[\mathbf{e}_t\mathbf{y}'_{t-1}] = 0$ from (95a) implies it (alternatively, note that it must be true by the definition of \mathbf{e}_t given in (90b)). However, it is the third row of $\mathbf{A}\mathbf{b}$ that is of the most interest since it provides us with a way to “update” a linear projection given new

information. The third row is given as:

$$\Omega_{31}\Omega_{11}^{-1}\mathbf{y}_{t-1} + (\Omega_{32} - \Omega_{31}\Omega_{11}^{-1}\Omega_{12}) (\Omega_{22} - \Omega_{21}\Omega_{11}^{-1}\Omega_{12})^{-1} \mathbf{e}_t + \mathbf{u}_t = \mathbf{x}_t \quad (97a)$$

$$\Leftrightarrow \Omega_{31}\Omega_{11}^{-1}\mathbf{y}_{t-1} + (\Omega_{32} - \Omega_{31}\Omega_{11}^{-1}\Omega_{12}) (\Omega_{22} - \Omega_{21}\Omega_{11}^{-1}\Omega_{12})^{-1} \mathbf{e}_t = \mathcal{P}[\mathbf{x}_t|\mathbf{y}_t, \mathbf{y}_{t-1}], \quad (97b)$$

which follows from the definition of \mathbf{u}_t .

Notice, however, that $E[\mathbf{e}_t\mathbf{e}'_t] = \Omega_{22} - \Omega_{21}\Omega_{11}^{-1}\Omega_{12} = \mathbf{D}_{22}$ is simply the variance of the observed prediction error \mathbf{e}_t .⁹ Similarly, $E[\mathbf{f}_t\mathbf{e}'_t] = \Omega_{32} - \Omega_{31}\Omega_{11}^{-1}\Omega_{12}$ where $\mathbf{f}_t = \mathbf{x}_t - \mathcal{P}[\mathbf{x}_t|\mathbf{y}_{t-1}] = \mathbf{x}_t - \hat{\mathbf{x}}_{t|t-1}$ represents the state forecast error. Furthermore, $E[\mathbf{f}_{t|t}\mathbf{f}'_{t|t}] = \Upsilon_{33} - \Upsilon_{32}\Upsilon_{22}^{-1}\Upsilon_{23} = \mathbf{D}_{33}$, given $\Upsilon_{ij} = \Omega_{ij} - \Omega_{i1}\Omega_{11}^{-1}\Omega_{1j}$, represents the variance of the state filtered error $\mathbf{f}_{t|t} = \mathbf{x}_t - \mathcal{P}[\mathbf{x}_t|\mathbf{y}_t, \mathbf{y}_{t-1}] = \mathbf{x}_t - \hat{\mathbf{x}}_{t|t}$. Therefore, what we have done in essence, is to derive the Kalman filter by means of orthogonal projectors and this is what (97) represents:

$$\mathcal{P}[\mathbf{x}_t|\mathbf{y}_t, \mathbf{y}_{t-1}] = \mathcal{P}[\mathbf{x}_t|\mathbf{y}_{t-1}] + E[\mathbf{f}_t\mathbf{e}'_t]E[\mathbf{e}_t\mathbf{e}'_t]^{-1}\mathbf{e}_t \quad (98a)$$

$$\equiv \hat{\mathbf{x}}_{t|t} = \hat{\mathbf{x}}_{t|t-1} + \mathbf{K}_t\mathbf{e}_t \quad (98b)$$

To derive \mathbf{K}_t in terms of the system matrices, we have:

$$\mathbf{K}_t = E[\mathbf{f}_t\mathbf{e}'_t]E[\mathbf{e}_t\mathbf{e}'_t]^{-1} \quad (99a)$$

$$= E[\mathbf{f}_t(\mathbf{H}\mathbf{f}_t + \boldsymbol{\eta}_t)']E[(\mathbf{H}\mathbf{f}_t + \boldsymbol{\eta}_t)(\mathbf{H}\mathbf{f}_t + \boldsymbol{\eta}_t)']^{-1}$$

$$= E[\mathbf{f}_t\mathbf{f}'_t\mathbf{H}' + \mathbf{f}_t\boldsymbol{\eta}'_t]E[\mathbf{H}\mathbf{f}_t\mathbf{f}'_t\mathbf{H}' + \mathbf{H}\mathbf{f}_t\boldsymbol{\eta}'_t + \boldsymbol{\eta}_t\mathbf{f}'_t\mathbf{H}' + \boldsymbol{\eta}_t\boldsymbol{\eta}'_t]^{-1}$$

$$\equiv \mathbf{P}_{t|t-1}\mathbf{H}' [\mathbf{H}\mathbf{P}_{t|t-1}\mathbf{H}' + \boldsymbol{\Sigma}_\boldsymbol{\eta}]^{-1} \quad (99b)$$

$$= \mathbf{P}_{t|t-1}\mathbf{H}'\boldsymbol{\Sigma}_{\mathbf{e},t}^{-1} \quad (99c)$$

so $\mathbf{P}_{t|t-1}$ is the state forecast error covariance and $\boldsymbol{\Sigma}_{\mathbf{e},t} = E[\mathbf{e}_t\mathbf{e}'_t]$.

Note that the difference equation (98) represents the filtered value of the first-moment of the

⁹Or equivalently, the MSE of the projection.

state given information up to time t . Similarly, the filtered second-moment can be derived as:

$$\text{Var}[\mathbf{x}_t | \mathbf{y}_{t-1}] = \mathbf{P}_{t|t-1} = \mathbf{F} \text{Var}[\mathbf{x}_{t-1} | \mathbf{y}_{t-1}] \mathbf{F}' + \Sigma_\epsilon = \mathbf{F} \mathbf{P}_{t-1|t-1} \mathbf{F}' + \Sigma_\epsilon, \quad (100a)$$

$$\text{and } \mathbf{P}_{t|t} = \mathbf{P}_{t|t-1} - \mathbf{K}_t \Sigma_{e_t} \mathbf{K}_t' = \mathbf{D}_{33}, \quad (100b)$$

which when combined together represent a matrix difference equation (the ARE or algebraic Riccati equation) [see equation (93a)] which can be solved for a “steady-state” value of \mathbf{P}_∞ given sufficient stability conditions; see also Simon (2006, pg.194-199)].

Note that (100b) follows from (98b):

$$\begin{aligned} \mathbf{x}_t - \hat{\mathbf{x}}_{t|t} &= \mathbf{x}_t - (\hat{\mathbf{x}}_{t|t-1} + \mathbf{K}_t \mathbf{e}_t) & (101a) \\ \Leftrightarrow \text{MSE}[\hat{\mathbf{x}}_{t|t}] &= \mathbf{P}_{t|t} \\ &= E[\mathbf{f}_t \mathbf{f}_t'] - E[\mathbf{K}_t \mathbf{e}_t \mathbf{e}_t' \mathbf{K}_t'] \\ &= \mathbf{P}_{t|t-1} - \mathbf{K}_t \Sigma_{e_t} \mathbf{K}_t'. \end{aligned}$$

The usefulness of the orthogonal projections type of approach to deriving the Kalman filter is that we can immediately decompose the optimal filtered state $\hat{\mathbf{x}}_{t|t}$ according to its orthogonal basis vectors, $\mathbf{e}_\tau, \forall \tau = t, t-1, t-2, \dots$. First, consider starting from time t and recursively substituting $\hat{\mathbf{x}}_{t|t-1} = \mathbf{F} \hat{\mathbf{x}}_{t-1|t-1}$ into (98):

$$\hat{\mathbf{x}}_{t|t} = \hat{\mathbf{x}}_{t|t-1} + \mathbf{K}_t \mathbf{e}_t \quad (102a)$$

$$\begin{aligned} &= \mathbf{F} (\hat{\mathbf{x}}_{t-1|t-2} + \mathbf{K}_{t-1} \mathbf{e}_{t-1}) + \mathbf{K}_t \mathbf{e}_t \\ &= \mathbf{F} \mathbf{F} (\hat{\mathbf{x}}_{t-2|t-3} + \mathbf{K}_{t-2} \mathbf{e}_{t-2}) + \mathbf{F} \mathbf{K}_{t-1} \mathbf{e}_{t-1} + \mathbf{K}_t \mathbf{e}_t \\ &\vdots \\ &= \mathbf{F}^j (\hat{\mathbf{x}}_{t-j|t-j-1} + \mathbf{K}_{t-j} \mathbf{e}_{t-j}) + \sum_{u=0}^{j-1} \mathbf{F}^u \mathbf{K}_{t-u} \mathbf{e}_{t-u} \end{aligned} \quad (102b)$$

so as $j \rightarrow \infty$ in (102b), if the eigenvalues of F are less than 1 in modulus, we have that:

$$\hat{\mathbf{x}}_{t|t} = \sum_{u=0}^{\infty} F^u \mathbf{K}_{t-u} \mathbf{e}_{t-u}. \quad (103)$$

But this is nothing more than an exponentially weighted moving average (EWMA)! In other words, the Kalman filtered state is simply weak white noise, \mathbf{e}_t , passed through a low-pass filter (see Section 7.4.2 below for further discussion on this point).

Moreover, we can also represent $\hat{\mathbf{x}}_{t|t}$ in terms of the basis Y_t , which is useful for interpreting the observed values as the “input” of the filter. From (98b) we have:

$$\begin{aligned} \hat{\mathbf{x}}_{t|t} &= \hat{\mathbf{x}}_{t|t-1} + \mathbf{K}_t (\mathbf{y}_t - \mathbf{H} \hat{\mathbf{x}}_{t|t-1}) \\ &= (\mathbf{I} - \mathbf{K}_t \mathbf{H}) \hat{\mathbf{x}}_{t|t-1} + \mathbf{K}_t \mathbf{y}_t, \quad \text{so let } \hat{\mathbf{L}}_t = \mathbf{I} - \mathbf{K}_t \mathbf{H}. \end{aligned} \quad (104a)$$

$$\begin{aligned} \Leftrightarrow \hat{\mathbf{x}}_{t|t} &= \hat{\mathbf{L}}_t \mathbf{F} \hat{\mathbf{x}}_{t-1|t-1} + \mathbf{K}_t \mathbf{y}_t \\ &= \hat{\mathbf{L}}_t \mathbf{F} \left(\hat{\mathbf{L}}_{t-1} \hat{\mathbf{x}}_{t-1|t-2} + \mathbf{K}_{t-1} \mathbf{y}_{t-1} \right) + \mathbf{K}_t \mathbf{y}_t \\ &= \hat{\mathbf{L}}_t \mathbf{L}_{t-1} \mathbf{F} \left(\hat{\mathbf{L}}_{t-2} \hat{\mathbf{x}}_{t-2|t-3} + \mathbf{K}_{t-2} \mathbf{y}_{t-2} \right) + \hat{\mathbf{L}}_t \mathbf{F} \mathbf{K}_{t-1} \mathbf{y}_{t-1} + \mathbf{K}_t \mathbf{y}_t \\ &= \hat{\mathbf{L}}_t \mathbf{L}_{t-1} \mathbf{L}_{t-2} \hat{\mathbf{x}}_{t-2|t-3} + \hat{\mathbf{L}}_t \mathbf{L}_{t-1} \mathbf{F} \mathbf{K}_{t-2} \mathbf{y}_{t-2} + \hat{\mathbf{L}}_t \mathbf{F} \mathbf{K}_{t-1} \mathbf{y}_{t-1} + \mathbf{K}_t \mathbf{y}_t \\ &\vdots \\ &= \hat{\mathbf{L}}_t \left(\prod_{k=1}^j \mathbf{L}_{t-k} \right) \hat{\mathbf{x}}_{t-j|t-(j+1)} \\ &\quad + \hat{\mathbf{L}}_t \sum_{u=2}^j \left(\prod_{k=1}^{u-1} \mathbf{L}_{t-k} \right) \mathbf{F} \mathbf{K}_{t-u} \mathbf{y}_{t-u} + \hat{\mathbf{L}}_t \mathbf{F} \mathbf{K}_{t-1} \mathbf{y}_{t-1} + \mathbf{K}_t \mathbf{y}_t. \end{aligned} \quad (104b)$$

where the third equality after the implication holds since $\mathbf{L}_{t-1} = \mathbf{F} (\mathbf{I} - \mathbf{K}_{t-1} \mathbf{H}) = \mathbf{F} \hat{\mathbf{L}}_{t-1}$.

So, if as $j \rightarrow \infty$ we have that $\mathbf{P}_{t|t-1} \rightarrow \mathbf{P}_{\infty}$, and the eigenvalues of \mathbf{L}_{∞} are less than 1 in modulus, then we can write:

$$\hat{\mathbf{x}}_{t|t} = \hat{\mathbf{L}} \sum_{u=0}^{\infty} \mathbf{L}^u \mathbf{F} \mathbf{K} \mathbf{y}_{t-1-u} + \mathbf{K} \mathbf{y}_t \quad (105)$$

which implies directly from (89a) that:

$$\hat{\mathbf{x}}_{t+1|t} = \sum_{u=0}^{\infty} \mathbf{L}^u \mathbf{F} \mathbf{K} \mathbf{y}_{t-u}. \quad (106)$$

And so we can see that the steady-state Kalman solution represents a linear time-invariant filter, with the observations $\{\mathbf{y}_t, \mathbf{y}_{t-1}, \dots\}$ as the filter input.

Similarly, we can solve for the response to a shock to the input process vector, ϵ_t , by deriving the impulse response matrices:

$$\mathbf{y}_t = \mathbf{H} \mathbf{x}_t + \boldsymbol{\eta}_t \quad (107a)$$

$$= \mathbf{H} (\mathbf{F} \mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t) + \boldsymbol{\eta}_t$$

$$= \mathbf{H} (\mathbf{F} (\mathbf{F} \mathbf{x}_{t-2} + \boldsymbol{\epsilon}_{t-1}) + \boldsymbol{\epsilon}_t) + \boldsymbol{\eta}_t$$

\vdots

$$= \mathbf{H} \mathbf{F}^j (\mathbf{F} \mathbf{x}_{t-(j+1)} + \boldsymbol{\epsilon}_{t-j}) + \sum_{u=0}^{j-1} \mathbf{H} \mathbf{F}^u \boldsymbol{\epsilon}_{t-u} + \boldsymbol{\eta}_t \quad (107b)$$

so as $j \rightarrow \infty$ in (107b), if the eigenvalues of \mathbf{F} are less than 1 in modulus, we have that:

$$\mathbf{y}_t = \sum_{u=0}^{\infty} \mathbf{H} \mathbf{F}^u \boldsymbol{\epsilon}_{t-u} + \boldsymbol{\eta}_t \quad (108)$$

therefore, $\mathbf{H} \mathbf{F}^u$ represents the impulse response to a unit shock at time t , u periods later. Note that this result is equivalent to (86) above, and so we can see that \mathbf{y}_t can be interpreted as the result of filtering the weak white noise, ϵ_t , which drives the state, plus a current observation noise η_t .

7.4.2 Spectral properties of Kalman filter

Equation (108) provides an expression for the coefficient matrices making up the impulse response, $\mathbf{A}(L)$. The spectral density $\mathbf{h}_{yy}(\omega)$ and power spectrum $|\mathbf{A}(e^{-i\omega})|^2$ are given respectively

from the associated z-transform with the same coefficients, at $z = e^{-i\omega}$, as: ¹⁰

$$\begin{aligned} \mathbf{y}_t &= \mathbf{A}(L)\boldsymbol{\epsilon}_t + \boldsymbol{\eta}_t \\ &= (\mathbf{H} [\mathbf{I} - \mathbf{F}L]^{-1}) \boldsymbol{\epsilon}_t + \boldsymbol{\eta}_t \end{aligned} \quad (109a)$$

$$\Leftrightarrow \mathbf{h}_{yy}(\omega) = \frac{1}{2\pi} (\mathbf{A}(e^{-i\omega})\boldsymbol{\Sigma}_\epsilon\mathbf{A}(e^{+i\omega})' + \boldsymbol{\Sigma}_\eta) \quad (109b)$$

$$\Leftrightarrow |\mathbf{A}(e^{-i\omega})|^2 = \left[\mathbf{H} [\mathbf{I} - \mathbf{F}e^{-i\omega}]^{-1} \right] \left[\mathbf{H} [\mathbf{I} - \mathbf{F}e^{+i\omega}]^{-1} \right]' \quad (109c)$$

Furthermore, to be more specific about the nature of the Kalman filter's frequency response let us consider taking the Fourier transform of a z-transform with the same coefficients as the impulse response function from (103):

$$\begin{aligned} \hat{\mathbf{x}}_{t|t} &= \mathbf{A}_t(L)\mathbf{e}_t \\ &= [\mathbf{I} - \mathbf{F}L]^{-1} \mathbf{K}_t\mathbf{e}_t \end{aligned} \quad (110a)$$

$$\Leftrightarrow |\mathbf{A}_t(e^{-i\omega})|^2 = [\mathbf{I} - \mathbf{F}e^{-i\omega}]^{-1} \mathbf{K}_t\mathbf{K}_t' [\mathbf{I} - \mathbf{F}e^{+i\omega}]^{-1'} \quad (110b)$$

The Kalman filter attenuates the high frequencies and passes through the lower ones in predicting the state from the observation errors, given information Y_t . The shape of the power spectrum will depend on both the eigenvalues of \mathbf{K}_t and \mathbf{F} . Interestingly, the optimal filter weights (and thus frequency response) are computed easily for each t , since they depend only on \mathbf{K}_t in the linear time invariant case. This contrasts to the Wiener solution where the transfer function needs to be recomputed independently for every change in t .

In terms of the Y_t basis we can interpret the \mathbf{y}_t 's as "inputs" passing through the filter. From

¹⁰As in the example given in Section 2.2.4, we have that since the z-transform $\mathbf{A}(z) = \sum_{u=0}^{\infty} \mathbf{a}_u z^u$, where $z \in \mathbb{C}$, admits the same coefficients as $\mathbf{A}(L)$, we can solve for the transfer function as $\mathbf{A}(z)$ where $z = e^{-i\omega}$.

(105), and assuming the steady-state filter, we have:

$$\begin{aligned}\hat{\mathbf{x}}_{t|t} &= \mathbf{A}(L)\mathbf{y}_t \\ &= \left(\hat{\mathbf{L}} [\mathbf{I} - \mathbf{L}L]^{-1} \mathbf{F}\mathbf{K}L + \mathbf{K} \right) \mathbf{y}_t\end{aligned}\quad (111a)$$

$$\begin{aligned}\Leftrightarrow |\mathbf{A}(e^{-i\omega})|^2 &= \left(\hat{\mathbf{L}} [\mathbf{I} - \mathbf{L}e^{-i\omega}]^{-1} \mathbf{F}\mathbf{K}e^{-i\omega} + \mathbf{K} \right) \\ &\quad \left(\hat{\mathbf{L}} [\mathbf{I} - \mathbf{L}e^{+i\omega}]^{-1} \mathbf{F}\mathbf{K}e^{+i\omega} + \mathbf{K} \right)'\end{aligned}\quad (111b)$$

All the transfer functions and power spectrums above are periodic matrix functions, with period 2π in their argument ω .

7.5 The MA and AR representations of the state-space model

Another useful representation of (1) can be obtained by rewriting the observation and state transition equations in terms of the observed forecast errors, e_t . First, add \mathbf{y}_{t+1} to both sides of (89b) to get:

$$\mathbf{y}_{t+1} = \mathbf{H}\hat{\mathbf{x}}_{t+1|t} + e_{t+1}.\quad (112)$$

From (104a) we have that:

$$\begin{aligned}\hat{\mathbf{x}}_{t+1|t} &= \mathbf{L}_t\hat{\mathbf{x}}_{t|t-1} + \mathbf{F}\mathbf{K}_t\mathbf{y}_t \\ &= \mathbf{F}\hat{\mathbf{x}}_{t|t-1} - \mathbf{F}\mathbf{K}_t\mathbf{H}\hat{\mathbf{x}}_{t|t-1} + \mathbf{F}\mathbf{K}_t\mathbf{y}_t \\ &= \mathbf{F}\hat{\mathbf{x}}_{t|t-1} + \mathbf{F}\mathbf{K}_te_t.\end{aligned}\quad (113a)$$

Therefore, together, (112) and (113a) represent the state-space model in (1) where both the observation and state transition noise have been replaced by the one-step ahead observation forecast errors.

The form of the state-space model in (112) and (113a) suggests both an MA and AR representation, where the e_t 's represent weak white noise innovations. The MA representation is

derived by substituting (103) into (89a) and then plugging this into (112):

$$\begin{aligned} \mathbf{y}_t &= \mathbf{H}\hat{\mathbf{x}}_{t|t-1} + \mathbf{e}_t \\ &= \mathbf{H} \left(\sum_{u=0}^{\infty} \mathbf{F}^{u+1} \mathbf{K}_{t-1-u} \mathbf{e}_{t-1-u} \right) + \mathbf{e}_t \end{aligned} \quad (112)$$

and again assuming that $\mathbf{P}_{t|t-1} \rightarrow \mathbf{P}_{\infty}$ exists, we can write:

$$\begin{aligned} \mathbf{y}_t &= \mathbf{H} (\mathbf{I} - \mathbf{FL})^{-1} \mathbf{FKL} \mathbf{e}_t + \mathbf{e}_t \\ \Leftrightarrow \mathbf{y}_t &= \Phi(L) \mathbf{e}_t \end{aligned} \quad (115a)$$

$$\text{where } \Phi(L) = \mathbf{I} + \mathbf{H} (\mathbf{I} - \mathbf{FL})^{-1} \mathbf{FKL}.$$

Furthermore, the AR representation is derived from subbing (106) into (112):

$$\begin{aligned} \mathbf{y}_t &= \mathbf{H}\hat{\mathbf{x}}_{t|t-1} + \mathbf{e}_t \\ &= \mathbf{H} \left(\sum_{u=0}^{\infty} \mathbf{L}^u \mathbf{FK} \mathbf{y}_{t-1-u} \right) + \mathbf{e}_t \end{aligned} \quad (112)$$

so we can write

$$\begin{aligned} \mathbf{y}_t &= \mathbf{H} (\mathbf{I} - \mathbf{LL})^{-1} \mathbf{FKL} \mathbf{y}_t + \mathbf{e}_t \\ \Leftrightarrow \Theta(L) \mathbf{y}_t &= \mathbf{e}_t \end{aligned} \quad (117a)$$

$$\text{where } \Theta(L) = \mathbf{I} - \mathbf{H} (\mathbf{I} - \mathbf{LL})^{-1} \mathbf{FKL}.$$

Interestingly, if we assume that the MA representation is invertible, then we must have that:

$$\Phi^{-1}(L) = \Theta(L) \quad (118a)$$

$$\Leftrightarrow (\mathbf{I} + \mathbf{H} (\mathbf{I} - \mathbf{FL})^{-1} \mathbf{FKL})^{-1} = \mathbf{I} - \mathbf{H} (\mathbf{I} - \mathbf{LL})^{-1} \mathbf{FKL}. \quad (118b)$$

7.6 Smoothing

The Kalman smoother is a generalization of the Kalman filter, where instead of computing $\mathcal{P}[\mathbf{x}_t|Y_t]$, the best linear predictor given information up to time t , we wish rather to use all the more recent information available to us up to some time $T > t$. Therefore, we wish to compute the best linear predictor of \mathbf{x}_τ given observed values, \mathbf{y}_t , $t = 1, \dots, T$, where $\tau < T$; that is, we wish to compute $\mathcal{P}[\mathbf{x}_t|Y_T]$. Again, for simplicity of notation, let $\mathcal{P}[\mathbf{x}_t|Y_T] \equiv \hat{\mathbf{x}}_{t|T}$. The smoother is derived by a simple extension of (98a):

$$\mathcal{P}[\mathbf{x}_t|Y_t] = \mathcal{P}[\mathbf{x}_t|\mathbf{y}_{t-1}] + E[\mathbf{f}_t \mathbf{e}'_t] E[\mathbf{e}_t \mathbf{e}'_t]^{-1} \mathbf{e}_t \quad (98a)$$

$$\equiv \hat{\mathbf{x}}_{t|t} = \hat{\mathbf{x}}_{t|t-1} + \mathbf{K}_t \mathbf{e}_t \quad (98b)$$

which becomes:

$$\mathcal{P}[\mathbf{x}_t|Y_T] = \mathcal{P}[\mathbf{x}_t|\mathbf{y}_T, \dots, \mathbf{y}_t, Y_{t-1}] = \mathcal{P}[\mathbf{x}_t|\mathbf{y}_{t-1}] + \sum_{k=t}^T E[\mathbf{f}_k \mathbf{e}'_k] E[\mathbf{e}_k \mathbf{e}'_k]^{-1} \mathbf{e}_k \quad (120a)$$

$$\equiv \hat{\mathbf{x}}_{t|T} = \hat{\mathbf{x}}_{t|t-1} + \mathbf{K}_t \mathbf{e}_t + \sum_{k=t+1}^T \{ \mathbf{P}_{t|t-1} \left(\prod_{j=t}^{k-1} \mathbf{L}'_j \right) \mathbf{H}' \} \Sigma_{\mathbf{e}_k}^{-1} \mathbf{e}_k. \quad (120b)$$

Equation (120a) is therefore a generalization of (98a), where we wish to update the best linear forecast of the state variable with not just information from the current time period t , but also with all future information, $\tau = t + 1, \dots, T$. Alternatively, we can see that we have replaced \mathbf{e}_t in (95a) with the stacked vector $\mathbf{e}_t^{*T} = \begin{bmatrix} \mathbf{e}'_t & \mathbf{e}'_{t+1} & \dots & \mathbf{e}'_{T-1} & \mathbf{e}'_T \end{bmatrix}'$, which suggests a generalized

form of (97b). Moreover, (120b) follows from:

$$\begin{aligned}
E[\mathbf{f}_t \mathbf{e}'_{t+1}] &= E[\mathbf{f}_t (\mathbf{H} \mathbf{f}_{t+1} + \boldsymbol{\eta}_{t+1})'] \\
&= E[\mathbf{f}_t \mathbf{f}'_{t+1} \mathbf{H}' + \mathbf{f}_t \boldsymbol{\eta}'_{t+1}] \\
&= E[\mathbf{f}_t \mathbf{f}'_{t+1}] \mathbf{H}' \\
&= \mathbf{P}_{t|t-1} \mathbf{L}'_t \mathbf{H}'
\end{aligned} \tag{121a}$$

$$\begin{aligned}
\Leftrightarrow E[\mathbf{f}_t \mathbf{e}'_{t+2}] &= \mathbf{P}_{t|t-1} \mathbf{L}'_t \mathbf{L}'_{t+1} \mathbf{H}' \\
&\vdots \\
E[\mathbf{f}_t \mathbf{e}'_T] &= \mathbf{P}_{t|t-1} \mathbf{L}'_t \mathbf{L}'_{t+1} \cdots \mathbf{L}'_{T-2} \mathbf{L}'_{T-1} \mathbf{H}',
\end{aligned} \tag{121b}$$

where equation (121a) follows from:

$$\begin{aligned}
\mathbf{f}_{t+1} &= \mathbf{x}_{t+1} - \hat{\mathbf{x}}_{t+1|t} \\
&= \mathbf{x}_{t+1} - \mathbf{F} \hat{\mathbf{x}}_{t|t} \\
&= \mathbf{x}_{t+1} - \mathbf{F} (\hat{\mathbf{x}}_{t|t-1} + \mathbf{K}_t \mathbf{e}_t) + \mathbf{F} \mathbf{x}_t - \mathbf{F} \mathbf{x}_t \\
&= \mathbf{F} (\mathbf{x}_t - \hat{\mathbf{x}}_{t|t-1}) + \boldsymbol{\epsilon}_{t+1} - \mathbf{F} \mathbf{K}_t \mathbf{e}_t \\
&= \mathbf{F} \mathbf{f}_t + \boldsymbol{\epsilon}_{t+1} - \mathbf{F} \mathbf{K}_t (\mathbf{H} \mathbf{f}_t + \boldsymbol{\eta}_t) \\
&= (\mathbf{F} (\mathbf{I} - \mathbf{K}_t \mathbf{H})) \mathbf{f}_t + \boldsymbol{\epsilon}_{t+1} - \mathbf{F} \mathbf{K}_t \boldsymbol{\eta}_t \\
&= \mathbf{L}_t \mathbf{f}_t + \boldsymbol{\epsilon}_{t+1} - \mathbf{F} \mathbf{K}_t \boldsymbol{\eta}_t.
\end{aligned} \tag{122a}$$

Note that from De Jong (1989) we can re-express (120b) in terms of the following recursion:

$$\mathbf{q}_{t-1} = \mathbf{L}'_t \mathbf{q}_t + \mathbf{H}' \boldsymbol{\Sigma}_{\mathbf{e}_t}^{-1} \mathbf{e}_t, \quad \text{where } t = T, \dots, 1. \tag{123}$$

That is, we can write (120b) as:

$$\hat{\mathbf{x}}_{t|T} = \hat{\mathbf{x}}_{t|t-1} + \mathbf{P}_{t|t-1} \mathbf{q}_{t-1}. \quad (124)$$

Moreover, from (124) (or alternatively from the same method employed in (101a)) we have that:

$$\mathbf{P}_{t|T} = \mathbf{P}_{t|t-1} - \mathbf{P}_{t|t-1} \text{Var}[\mathbf{q}_{t-1}] \mathbf{P}_{t|t-1} \quad (125a)$$

$$\text{where } \text{Var}[\mathbf{q}_{t-1}] \equiv \mathbf{M}_{t-1} = \mathbf{H}' \Sigma_{e_t}^{-1} \mathbf{H} + \mathbf{L}'_t \mathbf{M}_t \mathbf{L}_t \quad (125b)$$

noting that $\mathbf{P}_{t|T} = \mathbf{P}'_{t|T}$.

Therefore, (125) implies that:

$$\begin{aligned} \mathbf{P}_{t|T} &= \mathbf{P}_{t|t-1} - \mathbf{P}_{t|t-1} (\mathbf{H}' \Sigma_{e_t}^{-1} \mathbf{H} + \mathbf{L}'_t \mathbf{M}_t \mathbf{L}_t) \mathbf{P}_{t|t-1} \\ &= \mathbf{P}_{t|t-1} - \mathbf{P}_{t|t-1} \mathbf{H}' \Sigma_{e_t}^{-1} \mathbf{H} \mathbf{P}_{t|t-1} - \mathbf{P}_{t|t-1} \mathbf{L}'_t \mathbf{M}_t \mathbf{L}_t \mathbf{P}_{t|t-1} \\ &= \mathbf{P}_{t|t-1} - \mathbf{K}_t \Sigma_{e_t} \mathbf{K}'_t - \mathbf{P}_{t|t-1} \mathbf{L}'_t \mathbf{M}_t \mathbf{L}_t \mathbf{P}_{t|t-1} \\ &= \mathbf{P}_{t|t} - \mathbf{P}_{t|t-1} \mathbf{L}'_t \mathbf{M}_t \mathbf{L}_t \mathbf{P}_{t|t-1}. \end{aligned} \quad (126a)$$

So $\mathbf{P}_{t|t-1} \mathbf{L}'_t \mathbf{M}_t \mathbf{L}_t \mathbf{P}_{t|t-1}$ represents the reduction in MSE the smoother represents over the filter, since it uses the extra information embodied in $\{\mathbf{y}_T, \mathbf{y}_{T-1}, \dots, \mathbf{y}_{t+1}\}$.

What follows is known as the “fixed interval smoother” algorithm. This algorithm covers the case where T is fixed, and we desire the smoothed $\hat{\mathbf{x}}_{t|T}$ for any $t \leq T$. This is different from the “fixed-point smoother,” where t remains fixed and T increases, or the “fixed-lag smoother” where both t and T vary but their difference $T - t$ remains fixed.

Fixed Interval Smoother algorithm

1. Run the “Kalman filter algorithm” described in Section 7.4 to obtain $\hat{\mathbf{x}}_{t|t} = \hat{\mathbf{x}}_{t|t-1} + \mathbf{K}_t \mathbf{e}_t$.

2. Since we already have $\mathbf{P}_{t+1|t}$ from step (1), use equation (99c) to compute $\Sigma_{e_{t+1}}^{-1}$ and (99b) to get \mathbf{K}_{t+1} .
3. Plug \mathbf{K}_{t+1} into (93a) to obtain $\mathbf{P}_{t+2|t+1}$. \mathbf{K}_{t+1} can then be used to compute \mathbf{L}_{t+1} .
4. Plug $\hat{\mathbf{x}}_{t+1|t}$ [which we already have from step (1) and (89a)] into (89b) to obtain $\hat{\mathbf{y}}_{t+1|t}$.
5. Plug $\hat{\mathbf{y}}_{t+1|t}$ into (90b) to obtain \mathbf{e}_{t+1} .
6. Since we have $\mathbf{P}_{t+2|t+1}$ from step (3) above, we can now start again from step (2), ultimately repeating all the steps until we solve for $\mathbf{P}_{T|T-1}$. Once we've computed all the values of \mathbf{e}_τ , $\mathbf{L}_{\tau-1}$, and $\Sigma_{e_\tau}^{-1}$ for all $\tau = t+1, t+2, \dots, T$ we can then directly compute the sum in (120b).

7.6.1 Smoother in terms of orthogonal basis

Just as was done in Section 7.4.1, the fixed interval smoother can be represented in terms of an orthogonal basis. In fact, we have already done so—all that remains is to represent it here again in its entirety. This is useful since it allows us to appreciate the connection between the Kalman state-space approach to the prediction problem and that of the Wiener solution in the frequency domain.

The representation in terms of orthogonal basis is given from (89a) combined with (102b), and (120b):

$$\hat{\mathbf{x}}_{t|T} = \mathbf{F}^{j+1} \left(\hat{\mathbf{x}}_{t-1-j|t-2-j} + \mathbf{K}_{t-1-j} \mathbf{e}_{t-1-j} \right) + \sum_{u=0}^{j-1} \mathbf{F}^{u+1} \mathbf{K}_{t-1-u} \mathbf{e}_{t-1-u} \\ + \mathbf{K}_t \mathbf{e}_t + \sum_{k=t+1}^T \mathbf{P}_{t|t-1} \left(\prod_{j=t}^{k-1} \mathbf{L}'_j \right) \mathbf{H}' \Sigma_{e_k}^{-1} \mathbf{e}_k \quad (127)$$

so the first and second components in the sum represent past information, the third represents current information (at time t), and the last component represents future information.

Therefore, provided the eigenvalues of \mathbf{F} are less than 1 in modulus and $\mathbf{P}_{t|t-1} \rightarrow \mathbf{P}_\infty$,

allowing both j and T to approach ∞ results in the following expression of the impulse response function:

$$\begin{aligned}
\hat{\mathbf{x}}_{t|\infty} &= \sum_{u=0}^{\infty} \mathbf{F}^u \mathbf{F} \mathbf{K} e_{t-1-u} + \mathbf{K} e_t + \sum_{k=0}^{\infty} \mathbf{P} (\mathbf{L}')^{k+1} \mathbf{H}' \Sigma_e^{-1} e_{t+1+k} \\
&= \sum_{u=0}^{\infty} \mathbf{F}^u \mathbf{F} \mathbf{K} e_{t-1-u} + \sum_{k=0}^{\infty} \mathbf{P} (\mathbf{L}')^k \mathbf{H}' \Sigma_e^{-1} e_{t+k} \\
&= [\mathbf{I} - \mathbf{F} \mathbf{L}]^{-1} \mathbf{F} \mathbf{K} \mathbf{L} e_t + \mathbf{P} [\mathbf{I} - \mathbf{L}' \mathbf{L}^{-1}]^{-1} \mathbf{H}' \Sigma_e^{-1} e_t \\
&= \left([\mathbf{I} - \mathbf{F} \mathbf{L}]^{-1} \mathbf{F} \mathbf{K} \mathbf{L} + \mathbf{P} [\mathbf{I} - \mathbf{L}' \mathbf{L}^{-1}]^{-1} \mathbf{H}' \Sigma_e^{-1} \right) e_t \quad (128a) \\
&= \mathbf{C}(\mathbf{L}) e_t.
\end{aligned}$$

7.6.2 Spectral properties of the fixed interval smoother

The spectral properties of the filter $\mathbf{C}(\mathbf{L})$ in (128a) are available using the same methods as in 7.4.2 [see for example in equation (110a)] so they do not bear repeating.

However, what is of interest is the equivalence between the Wiener solution to the problem of signal extraction and that of the Kalman state-space smoother, as $T \rightarrow \infty$. I will not reproduce the entire equivalence proof here: interested readers can refer to Harvey & Proietti (2005). However, what can be said briefly is as follows.

First, take (128a) and replace e_t with its AR representation $\Theta(\mathbf{L}) \mathbf{y}_t$:

$$\hat{\mathbf{x}}_{t|\infty} = \mathbf{C}(\mathbf{L}) \Theta(\mathbf{L}) \mathbf{y}_t$$

Now, using the Woodbury matrix inversion lemma, the fact that the autocovariance generating matrix of \mathbf{y}_t is $\Gamma_{yy}(\mathbf{L}) = \Phi(\mathbf{L}) \Sigma_e \Phi(\mathbf{L}^{-1})'$ (from Section 7.5), the assumption of $\mathbf{P}_{t|t-1} \rightarrow \mathbf{P}_\infty$

and the solution to the discrete Riccati ARE, (93a), we have that:

$$\begin{aligned}
\hat{\mathbf{x}}_{t|\infty} &= \mathbf{C}(L)\boldsymbol{\Theta}(L)\mathbf{y}_t \\
&= \boldsymbol{\Gamma}_{xx}(L)\mathbf{H}'\boldsymbol{\Gamma}_{yy}^{-1}(L)\mathbf{y}_t \\
&= \boldsymbol{\Gamma}_{xy}(L)\boldsymbol{\Gamma}_{yy}^{-1}(L)\mathbf{y}_t
\end{aligned} \tag{130a}$$

$$\text{where } \boldsymbol{\Gamma}_{xx}(L) = \frac{1}{2\pi} [\mathbf{I} - \mathbf{F}L]^{-1} \boldsymbol{\Sigma}_\epsilon \left([\mathbf{I} - \mathbf{F}L^{-1}]^{-1} \right)' \tag{130b}$$

$$\text{and } \boldsymbol{\Gamma}_{yy}(L) = \frac{1}{2\pi} \mathbf{H}\boldsymbol{\Gamma}_{xx}(L)\mathbf{H}' + \boldsymbol{\Sigma}_\eta \tag{130c}$$

and where (130a) represents the Wiener solution in the time domain.

7.7 Conclusion

If we impose the stronger condition that the noise in (1), $\boldsymbol{\eta}_t$ and $\boldsymbol{\epsilon}_t$, is Gaussian (and not simply weak white noise), then we can replace all of the linear projections discussed above, $\mathcal{P}[\cdot]$, with expectations, since in the linear Gaussian case they coincide. Moreover, if the noise is Gaussian and uncorrelated, then it is also independent, and so these models are strong form.

In contrast, consider the case where the noise in (1) is weak white but the higher order moments rule out independence (e.g. $E[\epsilon^k \eta^j] \neq 0$ for some $k, j > 1$). Even in this case, the Kalman filter is still the best unbiased predictor amongst the class of *linear* predictions – see Simon (2006) pg.130. However, it is clear that in this case a *nonlinear* predictor will be more efficient.

Finally, if the noise has finite variance, but is coloured (and we are correct in assuming that all the relevant higher order moments are zero), then we can always modify the model in (1) to take account of this – see Simon (2006) sections 7.1 and 7.2 for more details.

8 Estimation in the time domain

The most straightforward way to estimate the parameters of the model in (1) is to employ Maximum Likelihood in the time domain. Given that the model is linear, and assuming the observation forecast errors, e_t , are Gaussian, their joint density can be factorized as:

$$p(\mathbf{e}_T, \mathbf{e}_{T-1}, \dots, \mathbf{e}_1 | \Theta) = p(\mathbf{y}_T, \mathbf{y}_{T-1}, \dots, \mathbf{y}_1 | \Theta) \quad (131a)$$

$$= p(\mathbf{y}_T | \mathbf{y}_{T-1}, \Theta) p(\mathbf{y}_{T-1} | \mathbf{y}_{T-2}, \Theta) \dots p(\mathbf{y}_2 | \mathbf{y}_1, \Theta) p(\mathbf{y}_1 | \Theta) \quad (131b)$$

where Θ is the parameter set. The first equality holds since $e_t = \mathbf{y}_t - \hat{\mathbf{y}}_{t|t-1}$ implies that the change of variables has unit Jacobian determinant (this can be proved by simply writing out $\hat{\mathbf{y}}_{t|t-1}$ recursively all the way back to $\hat{\mathbf{x}}_{1|0}$ – any e_t only depends on *present and past* values of \mathbf{y}_t and thus the Jacobian matrix will be block lower triangular with identity matrices on the diagonal).

Therefore, (131) implies that the log-likelihood of the model can be derived as:

$$p(\mathbf{e}_T, \mathbf{e}_{T-1}, \dots, \mathbf{e}_1 | \Theta) = \prod_{t=1}^T f(\mathbf{e}_t | \Theta) \quad (132a)$$

$$\propto \prod_{t=1}^T \frac{1}{\det(\mathbf{P}_{t|t-1})^{1/2}} \exp\left\{-\frac{1}{2} \mathbf{e}_t' \mathbf{P}_{t|t-1}^{-1} \mathbf{e}_t\right\} \quad (132b)$$

$$\Leftrightarrow \mathcal{L}(\Theta) = \sum_{t=1}^T -\frac{1}{2} \ln(\det(\mathbf{P}_{t|t-1})) - \frac{1}{2} \mathbf{e}_t' \mathbf{P}_{t|t-1}^{-1} \mathbf{e}_t \quad (132c)$$

where the first equality holds since the predicted observation errors are Gaussian and thus their uncorrelatedness implies independence.

Therefore, given (132) and some initial starting value for $\mathbf{x}_{1|0}$, we can optimize the log-likelihood with respect to Θ given the recursive state estimates $\hat{\mathbf{x}}_{t|t-1}$, $t = 1, \dots, T$, the Kalman gain \mathbf{K}_t and the estimated predicted state MSE, $\mathbf{P}_{t|t-1}$, for each choice of Θ .

If the predicted observation errors are not truly Gaussian, then this reduces to a Quasi-Maximum likelihood estimator which is still consistent asymptotically.

9 Factor models and common trends

Historically, the use of Factor models in Economics has largely involved their application to financial data. However, what may not be obvious to most practitioners, is that the Factor model is really just a special case of the general state-space model. This Section will briefly touch on the role of Factor models in finance and discuss both the traditional static Factor model and the newer dynamic Factor model representations. Moreover, we will show that the dynamic Factor model representation is intimately connected to the notion of co-integrated processes. In fact, the dynamic Factor model is equivalent to the co-integrated Vector Autoregressive Moving Average process (VARMA) or its equivalent Vector Error Correction (VECM) process representation.

9.1 Factor models in finance

The impetus for the popularity of Factor models in finance can probably be attributed to the seminal paper on Arbitrage Pricing Theory by Ross (1976). The CAPM model of Sharpe (1970) generated a linear security market line representation for excess returns, based on Markowitz mean-variance optimization and the implication of an efficient “market” portfolio. The “Sharpe ratio” implied that under equilibrium the beta of an asset, β , entirely determined its price proportional to the hypothetical market excess return. Concerned over the plausibility of the assumptions imposed by the CAPM and Markowitz portfolio theory, Ross (1976) aimed to reproduce the same linear representation for excess returns but where the model made no assumption of equilibrium behaviour. Since the model assumed only arbitrage as the force driving prices, the linear relationship could exist even under disequilibrium in the asset markets. In fact, it required no assumption of a “market” portfolio at all.

In order to formulate this linear model, Ross (1976) chose the Factor representation, where excess returns are explained as a linear function of some common factor and where an appropriately chosen arbitrage portfolio could still induce full diversification and elimination of the idiosyncratic risk component. This not only freed us from the need to model excess returns as a function of the hypothetical market return, but it allowed us to test for the predictive power of

other potential factors. Moreover, it presented us with a tractible way to explain excess returns using a limited number of factors, as opposed to calculating the direct covariance matrix between a multitude of assets as under Markowitz theory.

For example, Chen, Roll, and Ross (1986) introduced a multifactor model for stock returns, where the factors consist of unexpected, yet observable, changes in macroeconomic variables. Another famous example is the so-called 3-Factor model by Fama and French (1992), where they showed that beta only explained part of excess returns, and that book-to-market price ratios and market capitalization were also important explanatory factors.

In the previous examples, the factors chosen were observed (either as proxies or directly). However, factors can either be observed or unobserved (sometimes referred to as latent). If the factors chosen are unobserved then we must employ statistical methods to reconstitute the factors from the observable data. For the simple i.i.d. factor model, the traditional method of Principle Component Analysis can be used. For a dynamic factor model, the methods of state-space prediction discussed above in Section 7 are available.

9.2 The static Factor model

Traditionally, Factor models had treated the factors as static i.i.d. processes. The static Factor model for the observed variables \mathbf{y}_t is given as:

$$\mathbf{y}_t = \mathbf{B}\mathbf{f}_t + \boldsymbol{\eta}_t \quad (133a)$$

$$\text{and } \mathbf{f}_t = \boldsymbol{\epsilon}_t, \quad (133b)$$

$$\text{where } \boldsymbol{\eta}_t \sim IIN(0, \boldsymbol{\Sigma}_\eta),$$

$$\boldsymbol{\epsilon}_t \sim IIN(0, \boldsymbol{\Sigma}_\epsilon),$$

$$\text{and } 0 = E[\boldsymbol{\eta}_{t-j}\boldsymbol{\epsilon}'_{t-s}], \quad \forall j, s \in \mathbb{Z},$$

which is the same as the (strong form) model in (28) but where $\mathbf{H} \equiv \mathbf{B}$ and $\mathbf{x}_t \equiv \mathbf{f}_t$ and where we assume that:

1. The observed data has $E[\mathbf{y}_t] = 0$ and \mathbf{y}_t is $n \times 1$ and weakly stationary with no serial correlation.
2. $E[\mathbf{f}_t] = 0$ and $E[\mathbf{f}_t \mathbf{f}_t'] = \mathbf{\Omega}$, such that $\mathbf{\Omega}$ is diagonal and \mathbf{f}_t is a $k \times 1$ vector of common factors and $k < n$.
3. \mathbf{B} is $n \times k$ and represents the factor loadings.
4. $\boldsymbol{\eta}_t$ is such that $E[\boldsymbol{\eta}_t] = 0$ and $E[\boldsymbol{\eta}_t \boldsymbol{\eta}_t'] = \mathbf{\Lambda}_\eta$ which is a diagonal matrix, $\mathbf{\Lambda}_\eta = \text{diag}\{\sigma_1^2, \dots, \sigma_n^2\}$. Moreover, $\boldsymbol{\eta}_t$ is also $n \times 1$ and represents the specific error (often call the “idiosyncratic error”).
5. Finally, \mathbf{f}_t and $\boldsymbol{\eta}_t$ are independent.

Therefore, given the discussion in Section 5 we know that the optimal predictor of the unknown factor \mathbf{f}_t is given as:

$$E[\mathbf{f}_t | \mathbf{y}_t] = E[\mathbf{f}_t \mathbf{y}_t'] E[\mathbf{y}_t \mathbf{y}_t']^{-1} \mathbf{y}_t \quad (134a)$$

$$= E[\mathbf{f}_t (\mathbf{B} \mathbf{f}_t + \boldsymbol{\eta}_t)'] [\mathbf{B} \mathbf{\Omega} \mathbf{B}' + \mathbf{\Lambda}_\eta]^{-1} \mathbf{y}_t \quad (134b)$$

$$= \mathbf{\Omega} \mathbf{B}' [\mathbf{B} \mathbf{\Omega} \mathbf{B}' + \mathbf{\Lambda}_\eta]^{-1} \mathbf{y}_t \quad (134c)$$

which is equivalent to the static state-space model prediction, (74c), provided in the Section 5 examples, or analogously, the steady-state Kalman gain solution in (99b) where the state covariance matrix is time invariant, $\mathbf{\Omega} \equiv \mathbf{P}_\infty$.

9.2.1 Estimation in the time domain: the static Factor model and PCA

Unobserved static factors can be statistically reconstituted by employing the method of principle components analysis (PCA). The aim of PCA is to identify, from the observed data, statistically

orthogonal factors that can account for most of the variation in the covariance matrix of \mathbf{y}_t . To do so, first let us write the covariance matrix of \mathbf{y}_t :

$$E[\mathbf{y}_t \mathbf{y}_t'] \equiv \Sigma_{\mathbf{y}} = \mathbf{B} \Omega \mathbf{B}' + \Lambda_{\eta}. \quad (135)$$

Since $\Sigma_{\mathbf{y}}$ is symmetric, real, and positive-definite, it admits a spectral decomposition $\mathbf{L} \mathbf{D} \mathbf{L}'$ where \mathbf{L} has as its columns the orthonormal eigenvectors of $\Sigma_{\mathbf{y}}$ and \mathbf{D} is diagonal with $\mathbf{D} = \text{diag}\{\lambda_1, \dots, \lambda_n\}$, the eigenvalues of $\Sigma_{\mathbf{y}}$, arranged in decreasing order from largest to smallest. This implies that $E[\mathbf{L}' \mathbf{y}_t \mathbf{y}_t' \mathbf{L}] = \mathbf{D}$ is diagonal since $\mathbf{L}' \mathbf{L} = \mathbf{I}$.

Thus, $\mathbf{L}' \mathbf{y}_t$ makes a perfect candidate for the unobserved factor \mathbf{f}_t since it satisfies points (2) and (5) above in Section 9.2, and the i 'th row of $\mathbf{L}' \mathbf{y}_t$ represents a linear combination of the observed data which accounts for some proportion $\frac{\lambda_i}{\sum_{j=1}^n \lambda_j}$ of the variance in \mathbf{y}_t , since $\text{tr}(\Sigma_{\mathbf{y}}) = \sum_{j=1}^n \lambda_j$.

Now, suppose that we choose to only employ some k factors, such that $k < n$. What is the optimal choice of \mathbf{B} (i.e. what value of \mathbf{B} satisfies the optimal predictor of \mathbf{y}_t in the sense of Section 5) given that $\mathbf{f}_t = \mathbf{L}^* \mathbf{y}_t$ (where \mathbf{L}^* denotes the matrix composed of only the first k eigenvectors of $\hat{\Sigma}_{\mathbf{y}}$, the estimated unconditional covariance of \mathbf{y}_t) included as its columns? Under the constraint that $\mathbf{L}^* \mathbf{L}^{*'} = \mathbf{I}$ and the MMSE criterion, we have that the optimal value of \mathbf{B} is the linear projection of \mathbf{y}_t onto the space spanned by the columns of $\mathbf{f}_t = \mathbf{L}^* \mathbf{y}_t$:

$$E[\mathbf{y}_t (\mathbf{y}_t' \mathbf{L}^*)] E[(\mathbf{L}^* \mathbf{y}_t) (\mathbf{y}_t' \mathbf{L}^*)]^{-1} = \mathbf{B}, \quad (136a)$$

which can be estimated by sample moment counterpart. Moreover, since we have that $\boldsymbol{\eta}_t$ must lie in the space orthogonal to \mathbf{B} , we have that $\boldsymbol{\eta}_t = \mathbf{y}_t - E[\mathbf{y}_t \mathbf{f}_t'] E[\mathbf{f}_t \mathbf{f}_t']^{-1} \mathbf{f}_t$. Therefore the estimate for the residual error is $\hat{\Lambda}_{\eta} = \hat{\Sigma}_{\mathbf{y}} - \mathbf{B} \hat{\Omega} \mathbf{B}'$.

However, since simple PCA relies on the assumption that the dependent variable process \mathbf{y}_t is weakly stationary with no serial correlation, the method's popularity amongst cross-sectional researchers has not been matched within the time series community.

Attempts at reconciling PCA with dynamic factor models can be made by employing a two-step procedure where first some static factors are estimated by PCA as above, and, in the second step, a VAR model is specified as the factor process – see Breitung & Eickmeier (2006). Another method involves estimating the dynamic principal components within the frequency domain by eigenvalue decomposition of the spectral density of \mathbf{y}_t , then an estimate of the common components is obtained through its inversion – see Brillinger (1981, ch.9).

The nice feature of the state-space approach is that it allows us to easily predict latent, dynamic, factors, by means of the Kalman filter. Moreover, the dynamic Factor model for “common trends,” implied by the Engle-Granger representation theorem (1987) discussed below, manifests in a state-space representation. This makes the state-space framework useful for modeling unobserved stochastic trends that may drive the dynamics of co-integrated multivariate series.

9.3 The dynamic Factor model for common trends

The modern *dynamic Factor model* fits naturally within the state-space framework since it explicitly defines a parametric model for the dynamics of the unobserved factors. Moreover, within the multivariate context, the existence of co-integrating relationships amongst dependent variables leads directly to the “common trends” representation.

The *dynamic Factor model* for common trends is given in (137). Pre-multiplying the stochastic trend component $\boldsymbol{\mu}_t$, is the $N \times K$ matrix Θ of *factor loadings*, where $\boldsymbol{\mu}_t$ is now $K \times 1$ and represents the unobserved *common factors* where $0 \leq K \leq N$ and:

$$\mathbf{y}_t = \Theta \boldsymbol{\mu}_t + \boldsymbol{\mu}_0 + \boldsymbol{\eta}_t, \quad \text{for } t = 1, \dots, T \quad (137a)$$

$$\text{and } \boldsymbol{\mu}_t = \boldsymbol{\mu}_{t-1} + \boldsymbol{\beta} + \boldsymbol{\epsilon}_t. \quad (137b)$$

Additionally, $\boldsymbol{\mu}_0$ is included as a constant on the last $N - K$ series. That is, the first N elements of $\boldsymbol{\mu}_0$ are zero, and the last $N - K$ are $\tilde{\boldsymbol{\mu}}$.

As it stands the model in (137) is unidentifiable since for any singular $N \times N$ matrix \mathbf{H} , we

could redefine $\Theta^* = \Theta \mathbf{H}^{-1}$ and $\boldsymbol{\mu}_t^* = \mathbf{H} \boldsymbol{\mu}_t$ so that $\mathbf{y}_t = \Theta^* \boldsymbol{\mu}_t^* + \boldsymbol{\mu}_0 + \boldsymbol{\eta}_t$ and $\boldsymbol{\mu}_t^* = \boldsymbol{\mu}_{t-1}^* + \boldsymbol{\beta}^* + \boldsymbol{\epsilon}_t^*$ are observationally equivalent to (137), given $\boldsymbol{\epsilon}_t^* = \mathbf{H} \boldsymbol{\epsilon}_t$ and $\boldsymbol{\beta}_t^* = \mathbf{H} \boldsymbol{\beta}_t$.

Therefore, restrictions must be placed on the covariance matrix $\boldsymbol{\Sigma}_\epsilon$ and the factor loadings Θ . One way to solve the problem is to set $\Theta_{ij} = 0$ for $j > i$ and $\Theta_{ii} = 1$ where $i = 1, \dots, K$, while $\boldsymbol{\Sigma}_\epsilon$ is set as a diagonal matrix. Interestingly, if $K = N$, then the model is equivalent to (142) below, the simplest form of the so-called “structural” models popularized by Harvey (1984). Rewriting $\boldsymbol{\mu}_t^* = \Theta \boldsymbol{\mu}_t$, we have that the variance of $\boldsymbol{\epsilon}_t$ is the positive semi-definite matrix $\Theta \boldsymbol{\Sigma}_\epsilon \Theta'$.

As it stands the identification restriction placed above on Θ is a rather arbitrary one and does not suggest any reasonable economic intuition. However, this problem is not insurmountable since once the model in (137) is estimated under the suggested restriction, we can then subsequently introduce a *factor rotation* by including an orthogonal matrix \mathbf{H} , as above, in between Θ and $\boldsymbol{\mu}_t$. The new common trends are then $\boldsymbol{\mu}_t^*$. This approach may yield useful interpretations since we would be able to associate certain subsets of the elements our vector of observed variables, \mathbf{y}_t , with certain subsets of the common factors. Moreover, it may be possible to assign the observed series into various *clusters*, allowing the factor loadings matrix Θ to be block-diagonal.

9.3.1 Co-integration of levels process \mathbf{y}_t

The model in (137) also has the interesting feature that $N - K$ linear combinations of \mathbf{y}_t are stationary, despite the fact that all the elements of \mathbf{y}_t individually are nonstationary, $I(1)$, processes. Consider a matrix \mathbf{C} which is $(N - K) \times N$ and which has the property that $\mathbf{C} \Theta = 0$. We have from (137) that $\mathbf{C} \mathbf{y}_t = \mathbf{C} \boldsymbol{\mu}_0 + \mathbf{C} \boldsymbol{\eta}_t$, and thus the elements of $\mathbf{C} \mathbf{y}_t$ now represent $(N - K)$ stationary series, where \mathbf{C} has as its rows the $(N - K)$ *co-integrating relations*.

It turns out that this result is one implication of the Engle-Granger representation theorem (1987) which states that for any multivariate series which are co-integrated, there exist three equivalent representations: the VARMA, the VECM (vector error correction model), and the “common trends” Factor model representation.

To illustrate the connection suppose that:

$$\mathbf{y}_t = \mathbf{A}\mathbf{y}_{t-1} + \boldsymbol{\eta}_t \quad (138)$$

so that \mathbf{y}_t , $N \times 1$, follows a $VAR(1)$ process (and so \mathbf{A} is not generally symmetric).

Now, suppose that \mathbf{A} has a spectral decomposition $\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$, where \mathbf{P} contains the eigenvectors of \mathbf{A} as its columns at that \mathbf{P} has full column rank (i.e. there exists N distinct linearly independent eigenvectors). Moreover, \mathbf{D} is diagonal and contains the eigenvalues of \mathbf{A} .

The conditional forecast for \mathbf{y}_{t+h} , given some horizon $h \in \mathbb{N}^+$ is from (138):

$$\begin{aligned} E[\mathbf{y}_{T+h}|\mathbf{y}_T] &= \mathbf{A}^h \mathbf{y}_T \\ &= [\mathbf{P}\mathbf{D}\mathbf{P}^{-1}]^h \mathbf{y}_T \\ &= [\mathbf{P}\mathbf{D}^h \mathbf{P}^{-1}] \mathbf{y}_T \end{aligned} \quad (139a)$$

Without any loss of generality, let the first K eigenvalues in \mathbf{D} be equal to 1 and the other $N - K$ eigenvalues have modulus strictly less than 1. Furthermore, let $\mathbf{P}^{-1}\mathbf{y}_t = \mathbf{z}_t$. Equation (139a) implies that in the limit, as $h \rightarrow \infty$, the first K elements of $E[\mathbf{z}_{t+h}|\mathbf{z}_T]$ will approach ∞ while the last $N - K$ elements will approach 0. That is, the last $N - K$ elements of \mathbf{z}_t are mean reverting¹¹ while for the first K elements no unconditional mean exists. Therefore, there exist linear combinations of \mathbf{y}_t which are $I(0)$ (i.e. the last $N - K$ elements of \mathbf{z}_t) despite the fact that all the elements \mathbf{y}_t are $I(1)$ (by virtue of the K unit eigenvalues of \mathbf{A}). We say that there exist $N - K$ co-integrating relations, or equivalently, that there exist K common stochastic trends driving \mathbf{y}_t .

Given that the last $N - K$ diagonal elements of \mathbf{D}^h approach 0 in the limit, we can always

¹¹That is, the conditional mean reverts back to the unconditional mean of 0 as $h \rightarrow \infty$. Or in other words, $E[\mathbf{z}_{t+h}|\mathbf{z}_T] \rightarrow E[\mathbf{z}_t] = 0$ as $h \rightarrow \infty$ for those last $N - K$ stable elements of \mathbf{z}_t .

rewrite (139a) as:

$$\begin{aligned}
\lim_{h \rightarrow \infty} E[\mathbf{y}_{T+h} | \mathbf{y}_T] &= \mathbf{A}^\infty \mathbf{y}_T \\
&= \mathbf{P} \begin{bmatrix} \bar{\mathbf{z}}_T \\ \mathbf{0} \end{bmatrix} \\
&= \bar{\mathbf{P}} \bar{\mathbf{z}}_T
\end{aligned} \tag{140a}$$

where $\bar{\mathbf{z}}_T$ is the first K rows of \mathbf{z}_T , and $\bar{\mathbf{P}}$, $N \times K$, is the first K columns of \mathbf{P} .

Therefore, we have that the “common trends” representation of the co-integrated system \mathbf{y}_t is given as:

$$\begin{aligned}
\mathbf{y}_t &= \lim_{h \rightarrow \infty} E[\mathbf{y}_t | \mathbf{y}_{t-h}] + \boldsymbol{\eta}_t \\
&= \bar{\mathbf{P}} \bar{\mathbf{z}}_t + \boldsymbol{\eta}_t
\end{aligned} \tag{141a}$$

so if we allow $\mathbf{C} \equiv \bar{\mathbf{P}}^{-1}$ to be the $(N - K) \times N$ matrix with the co-integrating relations (i.e. the last $N - K$ rows of \mathbf{P}^{-1}) as its rows, then we have the result that $\mathbf{C}\mathbf{y}_t = \mathbf{C}\boldsymbol{\eta}_t$ is $I(0)$, since $\bar{\mathbf{P}}^{-1} \bar{\mathbf{P}} = \mathbf{0}$. Note the similarity to (137) where $\bar{\mathbf{z}}_t \equiv \boldsymbol{\mu}_t$, $\boldsymbol{\Theta} \equiv \bar{\mathbf{P}}$, and of course, $\bar{\mathbf{z}}_t$ is $I(1)$.

9.4 Unobserved Components

The *dynamic Factor model* for common trends in (137) also represents the simplest form of the structural (or “Unobserved Components”) framework popularized by Harvey (1984), wherein it is referred to as the *Local level random walk plus noise* model:¹²

$$\mathbf{y}_t = \boldsymbol{\mu}_t + \boldsymbol{\eta}_t, \quad \text{for } t = 1, \dots, T, \tag{142a}$$

$$\text{and } \boldsymbol{\mu}_t = \boldsymbol{\mu}_{t-1} + \boldsymbol{\epsilon}_t, \tag{142b}$$

¹²See Harvey (1989) pg. 429 and Nerlove et al. (1979) for a review of the unobserved components models more generally.

where $\boldsymbol{\mu}_t$ is an $N \times 1$ vector of local level components, $\boldsymbol{\eta}_t$ and $\boldsymbol{\epsilon}_t$ are $N \times 1$ vectors of multivariate white noise with zero mean and covariance matrices $\boldsymbol{\Sigma}_\eta$ and $\boldsymbol{\Sigma}_\epsilon$ respectively, and \mathbf{y}_t represents an $N \times 1$ vector of series which may be contemporaneously related (e.g. cross-sectional data).

The common trends framework, with $I(1)$ trends, models the local level component, $\boldsymbol{\mu}_t$, as a random walk (i.e. stochastic trend), although it is free to take on any number of different structural representations [see (148a) below or example 2.2.10 for the general representation]. Not only do the $\boldsymbol{\eta}_t$'s exhibit contemporaneous correlation, but the innovations of the trend components themselves, $\boldsymbol{\epsilon}_t$, do as well. The N series, \mathbf{y}_t , are thus linked via the off-diagonals of the two innovation covariance matrices, each of which include $N(N + 1)/2$ parameters to be estimated. It is worth mentioning that (142) is similar to the time-varying SUR model of Zellner (1962), with the exception that we replace the regressors $\mathbf{X}_t\boldsymbol{\beta}$ with the time-varying intercept vector $\boldsymbol{\mu}_t$.

9.4.1 Alternative representations

The model in (142) can also be rewritten in stationary single equation form as:

$$(\mathbf{I} - \mathbf{I}L)\mathbf{y}_t = \boldsymbol{\epsilon}_t + (\mathbf{I} - \mathbf{I}L)\boldsymbol{\eta}_t, \quad \text{where } t = 1, \dots, T, \quad (143)$$

so that $E[\Delta\mathbf{y}_t] = 0$, $Var[\Delta\mathbf{y}_t] = \boldsymbol{\Sigma}_\epsilon + 2\boldsymbol{\Sigma}_\eta$, and $Covar[\Delta\mathbf{y}_t, \Delta\mathbf{y}_{t-1}] = -\boldsymbol{\Sigma}_\eta$.¹³ Therefore, the reduced form stationary model of (143) is of the VARIMA(0,1,1) type since it matches its autocorrelation representation. That is, if $\mathbf{x}_t = \mathbf{u}_t + \boldsymbol{\Phi}\mathbf{u}_{t-1}$ we have that $\boldsymbol{\Gamma}_x(0) = E[\mathbf{x}_t\mathbf{x}'_{t-1}] = \boldsymbol{\Sigma}_u + \boldsymbol{\Phi}\boldsymbol{\Sigma}_u\boldsymbol{\Phi}'$ and $\boldsymbol{\Gamma}_x(1) = E[\mathbf{x}_t\mathbf{x}'_{t-1}] = \boldsymbol{\Phi}\boldsymbol{\Sigma}_u$ which is of the same form if $\boldsymbol{\Sigma}_\epsilon = 0$ and $\boldsymbol{\Phi} = -\mathbf{I}$.

The structural form of the model in (142) implies restrictions on the parameter space of the reduced form VARIMA(0,1,1) representation of (143).¹⁴ Take for simplicity the example of the

¹³All other autocovariances are zero.

¹⁴Note, this is similar to translating between the state-space form and the transfer-function form of a system within the engineering context. While the state-space form implies a unique transfer function form, the transfer function form has any number of respective state-space representations. So in this context, the specified state-space form implies restrictions on the transfer function form of the model. By transfer function form, we mean the model written purely in terms of impulse response to exogenous shocks. In this sense, we can think of the multiplicity of state-space forms within the context of Akaike (1975) and the ‘‘minimal’’ representation.

univariate analog to (143), the ARIMA(0,1,1) model:

$$x_t = u_t + \phi u_{t-1}, \quad \text{where } u_t \sim N(0, \sigma_u^2). \quad (144)$$

Matching autocorrelations of (144) with the univariate counterpart to (143) results in:

$$\rho(1) = \frac{-\sigma_\eta^2}{\sigma_\epsilon^2 + 2\sigma_\eta^2} = \frac{\phi\sigma_u^2}{(1 + \phi^2)\sigma_u^2} \quad (145a)$$

$$\Leftrightarrow 0 = \phi^2\sigma_\eta^2 + \phi(\sigma_\epsilon^2 + 2\sigma_\eta^2) + \sigma_\eta^2 \quad (145b)$$

which is a second-order polynomial in ϕ . One of the roots is dropped since it implies non-invertibility. Solving (145b) we find that:

$$\phi = -(q + 2)/2 + \sqrt{q^2/4 + q}, \quad \text{where } q = \frac{\sigma_\epsilon^2}{\sigma_\eta^2} \quad (146a)$$

$$\Rightarrow \text{if } 0 \leq q \leq \infty \text{ we have that } -1 \leq \phi \leq 0. \quad (146b)$$

And thus the standard stationary region of the parameter space is cut in half. Note that the order of the polynomial is equal to one plus the order of the MA process. Therefore solving for the restrictions of an $MA(q)$ process would generally require solving a Yule-Walker type system of equations, which include, at most, a $q + 1$ order polynomial. Given this, Nerlove et al. (1979, pp. 70-78) provides a more general algorithm applicable to complicated univariate structural forms. As for the multivariate case, the problem is even more complex. In the univariate case, the structural form implies restrictions on the parameter space of the reduced form, but in the multivariate case it also implies a reduction in the dimension of the model. See Harvey (1989, pp. 432) for more details.

Interestingly, the model in (142) can also be viewed as the multivariate *UCARIMA* type (Engle, 1978). The UCARIMA representation decomposes the process $\mathbf{y}_t = T_t + S_t + C_t + I_t$ into individual ARIMA processes (where T_t is the trend component, S_t is some seasonal component, C_t a non-seasonal cyclical component, and I_t is some residual irregular component.) For example

in (142) we have only defined the trend component, $T_t = \boldsymbol{\mu}_t$ and the irregular $I_t = \boldsymbol{\eta}_t$, so that the second and third components are effectively zero. However, this need not be the case—we could choose to define other components for S_t and C_t . For example, we could have chosen stochastic harmonic processes.

Either way, the first trend component can manifest as an ARIMA process by representing (142) as an ARIMA(0,1,1) process, with its parameter, $\boldsymbol{\Phi}$, characterized in the discussion above in determining the reduced form of (143):

$$\mathbf{y}_t = \Delta^{-1} \boldsymbol{\epsilon}_t + \boldsymbol{\eta}_t \quad (147a)$$

$$\equiv \Delta^{-1} (\mathbf{I} + \boldsymbol{\Phi}L) \mathbf{u}_t, \quad (147b)$$

where $\Delta = (\mathbf{I} - \mathbf{I}L)$. Therefore, \mathbf{y}_t is composed of a singular ARIMA(0,1,1) process (where we have subsumed the irregular component into the trend for clarity).

More generally, we can write the multivariate *UCARIMA* form under any particular set of components $\{T_t, S_t, C_t, \dots, I_t\}$ as:

$$\mathbf{y}_t = \sum_{m=0}^M \Delta_m^{-1}(L) \boldsymbol{\Theta}_m^{-1}(L) \boldsymbol{\Phi}_m(L), \quad (148a)$$

$$\text{where } \Delta(L) = \prod_{m=0}^M \Delta_m(L), \quad (148b)$$

so that pre-multiplying (148a) by (148b) makes the process stationary. For reference, consider the example in 2.2.10 where we presented the *Basic Structural Model*; it also has a UCARIMA representation – see Harvey (1989, pg.74).

10 References

H. Akaike (1974) "Markovian representation of stochastic processes and its application to the analysis of autoregressive moving average processes," *Ann. Inst. Statist. Math.*, 26, 363-387

————— (1975) "Markovian representation of stochastic processes by canonical variables," *SIAM Journal on Control* 13, 162-173

G. Box and G. Jenkins (1970) *Time series analysis: Forecasting and control*, San Francisco: Holden-Day

J. Breitung, S. Eickmeier (2006) "Dynamic factor models," in *Modern Econometric Analysis: Surveys on Recent Developments*, eds. O. Hubler, J. Frohn, Springer

D.R. Brillinger (1981) *Time Series: Data Analysis and Theory: Expanded edition*, Holden-Day, San Francisco

N.F. Chen, R. Roll, S.A. Ross (1986) "Economic forces and the stock market," *The Journal of Business*, 59, 383-404

De Jong (1989) "Smoothing and interpolation with the state space model," *Journal of the American Statistical Association*, 84, 1085-1088

R. Engle (1978) "Estimating structural models of seasonality," in *Seasonal Analysis of Economic Time Series*, A. Zellner ed., NBER

————— (1982) "Autoregressive conditional heteroskedasticity with estimates of the variance of the United Kingdom inflation," *Econometrica*, 50(4), 987-1007

————— (1987) "Multivariate ARCH with factor structures – cointegration in variance," University of California Press, San Diego

R.F. Engle, C.W.J. Granger (1987) "Co-integration and error correction: representation, estimation, and testing," *Econometrica*, 55, 251-276

E.F. Fama, K.R. French (1992) "The cross-section of expected stock returns," *Journal of Finance*, 47(2), 427-465

- J.D. Hamilton (1994) *Time Series Analysis*, Princeton University Press, New Jersey
- E.J. Hannan, R.D. Terrell, N. Tuckwell (1970) "The seasonal adjustment of economic time series," *International Economic Review*, 11, 24-52
- P.J. Harrison, C.F. Stevens (1976) "Bayesian forecasting," *Journal of the Royal Statistical Society Series B*, 38(3), 205-247
- A.C. Harvey (1984) "A unified view of statistical forecasting procedures," *Journal of Forecasting*, 3(3), 245-275
- (1989) *Forecasting, structural time series models, and the Kalman filter*, Cambridge University Press
- A.E. Kalman (1960) "A new approach to linear filtering and prediction problems," *Journal of Basic Engineering*, 82(1), 35-45
- Nerlove, Grether, and Carvalho (1979) *Analysis of economic time series: A synthesis*, New York: Academic Press
- M.B. Priestley (1981) *Spectral Analysis and Time Series, 2 volumes*, London: Elsevier Academic Press
- S. Ross (1976) "The arbitrage theory of capital asset pricing," *Journal of Economic Theory*, 13(3), 341-360
- W.F. Sharpe (1966) "Mutual fund performance," *The Journal of Business*, 39(1), 119-138
- (1970) *Portfolio theory and capital markets*, McGraw-Hill, New York
- D. Simon (2006) *Optimal State Estimation*, New Jersey: Wiley
- P.C. Young (2011) *Recursive Estimation and Time-Series Analysis*, 2nd ed., Berlin: Springer-Verlag

Zellner (1962) "An efficient method of estimating seemingly unrelated regression equations and tests of aggregation bias," *Journal of the American Statistical Association*, 57, 500-509