

2

Local level model

2.1 Introduction

The purpose of this chapter is to introduce the basic techniques of state space analysis, such as filtering, smoothing, initialisation and forecasting, in terms of a simple example of a state space model, the local level model. This is intended to help beginners grasp the underlying ideas more quickly than they would if we were to begin the book with a systematic treatment of the general case. So far as inference is concerned, we shall limit the discussion to the classical standpoint; a Bayesian treatment of the local level model may be obtained as a special case of the Bayesian analysis of the linear Gaussian model in Chapter 8.

A *time series* is a set of observations y_1, \dots, y_n ordered in time. The basic model for representing a time series is the additive model

$$y_t = \mu_t + \gamma_t + \varepsilon_t, \quad t = 1, \dots, n. \quad (2.1)$$

Here, μ_t is a slowly varying component called the *trend*, γ_t is a periodic component of fixed period called the *seasonal* and ε_t is an irregular component called the *error* or *disturbance*. In general, the observation y_t and the other variables in (2.1) can be vectors but in this chapter we assume they are scalars. In many applications, particularly in economics, the components combine multiplicatively, giving

$$y_t = \mu_t \gamma_t \varepsilon_t. \quad (2.2)$$

By taking logs however and working with logged values model (2.2) reduces to model (2.1), so we can use model (2.1) for this case also.

To develop suitable models for μ_t and γ_t we need the concept of a *random walk*. This is a scalar series α_t determined by the relation $\alpha_{t+1} = \alpha_t + \eta_t$ where the η_t 's are independent and identically distributed random variables with zero means and variances σ_η^2 .

Consider a simple form of model (2.1) in which $\mu_t = \alpha_t$ where α_t is a random walk, no seasonal is present and all random variables are normally distributed. We assume that ε_t has constant variance σ_ε^2 . This gives the model

$$\begin{aligned} y_t &= \alpha_t + \varepsilon_t, & \varepsilon_t &\sim N(0, \sigma_\varepsilon^2), \\ \alpha_{t+1} &= \alpha_t + \eta_t, & \eta_t &\sim N(0, \sigma_\eta^2), \end{aligned} \quad (2.3)$$

for $t = 1, \dots, n$ where the ε_t 's and η_t 's are all mutually independent and are independent of α_1 . This model is called the *local level model*. Although it has a simple form, this model is not an artificial special case and indeed it provides the basis for the analysis of important real problems in practical time series analysis; for example, it provides the basis for exponentially weighted moving averages as discussed in §3.4. It exhibits the characteristic structure of state space models in which there is a series of unobserved values $\alpha_1, \dots, \alpha_n$ which represents the development over time of the system under study, together with a set of observations y_1, \dots, y_n which are related to the α_t 's by the state space model (2.3). The object of the methodology that we shall develop is to infer relevant properties of the α_t 's from a knowledge of the observations y_1, \dots, y_n .

We assume initially that $\alpha_1 \sim N(a_1, P_1)$ where a_1 and P_1 are known and that σ_ε^2 and σ_η^2 are known. Since random walks are non-stationary the model is non-stationary. By non-stationary here we mean that distributions of random variables y_t and α_t depend on time t .

For applications of the model to real series using classical inference, we need to compute quantities such as the mean of α_t given y_1, \dots, y_{t-1} or the mean of α_t given y_1, \dots, y_n , together with their variances; we also need to fit the model to data by calculating maximum likelihood estimates of the parameters σ_ε^2 and σ_η^2 . In principle, this could be done by using standard results from multivariate normal theory as described in books such as Anderson (1984). In this approach the observations y_t generated by the local level model are represented as an $n \times 1$ vector y such that

$$y \sim N(1a_1, \Omega), \quad \text{with } y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad 1 = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \quad \text{and } \Omega = 11'P_1 + \Sigma, \quad (2.4)$$

where the (i, j) th element of the $n \times n$ matrix Σ is given by

$$\Sigma_{ij} = \begin{cases} (i-1)\sigma_\eta^2, & i < j \\ \sigma_\varepsilon^2 + (i-1)\sigma_\eta^2, & i = j, \\ (j-1)\sigma_\eta^2, & i > j \end{cases} \quad i, j = 1, \dots, n, \quad (2.5)$$

which follows since the local level model implies that

$$y_t = \alpha_1 + \sum_{j=1}^{t-1} \eta_j + \varepsilon_t, \quad t = 1, \dots, n. \quad (2.6)$$

Starting from this knowledge of the distribution of y , estimation of conditional means, variances and covariances is in principle a routine matter using standard results in multivariate analysis based on the properties of the multivariate normal distribution. However, because of the serial correlation between the observations y_t , the routine computations rapidly become cumbersome as n increases. This naive

approach to estimation can be improved on considerably by using the filtering and smoothing techniques described in the next three sections. In effect, these techniques provide efficient computing algorithms for obtaining the same results as those derived by multivariate analysis theory. The remaining sections of this chapter deal with other important issues such as fitting the local level model and forecasting future observations.

The local level model (2.3) is a simple example of a *linear Gaussian state space model*. The variable α_t is called the *state* and is unobserved. The overall object of the analysis is to study the development of the state over time using the observed values y_1, \dots, y_n . Further examples of these models will be described in Chapter 3. The general form of the linear Gaussian state space model is given in equation (3.1); its properties will be considered in detail from the standpoint of classical inference in Chapters 4-7. A Bayesian treatment of the model will be given in Chapter 8.

2.2 Filtering

2.2.1 THE KALMAN FILTER

The object of filtering is to update our knowledge of the system each time a new observation y_t is brought in. We shall develop the theory of filtering for the local level model (2.3). Since all distributions are normal, conditional joint distributions of one set of observations given another set are also normal. Let Y_{t-1} be the set of past observations $\{y_1, \dots, y_{t-1}\}$ and assume that the conditional distribution of α_t given Y_{t-1} is $N(a_t, P_t)$ where a_t and P_t are to be determined. Given that a_t and P_t are known, our object is to calculate a_{t+1} and P_{t+1} when y_t is brought in. We do this by using some results from elementary regression theory.

Since $a_{t+1} = E(\alpha_{t+1}|Y_t) = E(\alpha_t + \eta_t|Y_t)$ and $P_{t+1} = \text{Var}(\alpha_{t+1}|Y_t) = \text{Var}(\alpha_t + \eta_t|Y_t)$ from (2.3), we have

$$a_{t+1} = E(\alpha_t|Y_t), \quad P_{t+1} = \text{Var}(\alpha_t|Y_t) + \sigma_\eta^2. \quad (2.7)$$

Define $v_t = y_t - a_t$ and $F_t = \text{Var}(v_t)$. Then

$$E(v_t|Y_{t-1}) = E(\alpha_t + \varepsilon_t - a_t|Y_{t-1}) = a_t - a_t = 0.$$

Thus $E(v_t) = E[E(v_t|Y_{t-1})] = 0$ and $\text{Cov}(v_t, y_j) = E(v_t y_j) = E[E(v_t|Y_{t-1})y_j] = 0$ so v_t and y_j are independent for $j = 1, \dots, t-1$. When Y_t is fixed, Y_{t-1} and y_t are fixed so Y_{t-1} and v_t are fixed and vice versa. Consequently, $E(\alpha_t|Y_t) = E(\alpha_t|Y_{t-1}, v_t)$ and $\text{Var}(\alpha_t|Y_t) = \text{Var}(\alpha_t|Y_{t-1}, v_t)$. Since all variables are normally distributed, the conditional expectation and variance are given by standard formulae from multivariate normal regression theory. For a general treatment of the results required see the regression lemma in §2.13.

It follows from equation (2.49) of this lemma that

$$E(\alpha_t|Y_t) = E(\alpha_t|Y_{t-1}, v_t) = E(\alpha_t|Y_{t-1}) + \text{Cov}(\alpha_t, v_t)\text{Var}(v_t)^{-1}v_t, \quad (2.8)$$

where

$$\begin{aligned}\text{Cov}(\alpha_t, v_t) &= E[\alpha_t(y_t - a_t)] = E[\alpha_t(\alpha_t + \varepsilon_t - a_t)] \\ &= E[\alpha_t(\alpha_t - a_t)] = E[\text{Var}(\alpha_t|Y_{t-1})] = P_t,\end{aligned}$$

and

$$\text{Var}(v_t) = F_t = \text{Var}(\alpha_t + \varepsilon_t - a_t) = \text{Var}(\alpha_t|Y_{t-1}) + \text{Var}(\varepsilon_t) = P_t + \sigma_\varepsilon^2.$$

Since $a_t = E(\alpha_t|Y_{t-1})$ we have from (2.8)

$$E(\alpha_t|Y_t) = a_t + K_t v_t. \quad (2.9)$$

where $K_t = P_t/F_t$ is the regression coefficient of α_t on v_t . From equation (2.50) of the regression lemma in §2.13 we have

$$\begin{aligned}\text{Var}(\alpha_t|Y_t) &= \text{Var}(\alpha_t|Y_{t-1}, v_t) \\ &= \text{Var}(\alpha_t|Y_{t-1}) - \text{Cov}(\alpha_t, v_t)^2 \text{Var}(v_t)^{-1} \\ &= P_t - P_t^2/F_t \\ &= P_t(1 - K_t).\end{aligned} \quad (2.10)$$

From (2.7), (2.9) and (2.10), we have the full set of relations for updating from time t to time $t + 1$,

$$\begin{aligned}v_t &= y_t - a_t, & F_t &= P_t + \sigma_\varepsilon^2, & K_t &= P_t/F_t, \\ a_{t+1} &= a_t + K_t v_t, & P_{t+1} &= P_t(1 - K_t) + \sigma_\eta^2,\end{aligned} \quad (2.11)$$

for $t = 1, \dots, n$. Note that a_1 and P_1 are assumed known here; however, more general initial specifications will be dealt with in §2.9. Relations (2.11) constitute the celebrated *Kalman filter* for this problem. It should be noted that P_t depends only on σ_ε^2 and σ_η^2 and does not depend on Y_{t-1} . We include the case $t = n$ in (2.11) for convenience even though a_{n+1} and P_{n+1} are not normally needed for anything except forecasting.

The notation employed in (2.11) is not the simplest that could have been used; we have chosen it in order to be compatible with the notation that we consider appropriate for the treatment of the general multivariate model in Chapter 4. A set of relations such as (2.11) which enables us to calculate quantities for $t + 1$ given those for t is called a *recursion*. Formulae (2.9) and (2.10) could be derived in many other ways. For example, a routine Bayesian argument for normal densities could be used in which the prior density is $p(\alpha_t|Y_{t-1}) \sim N(a_t, P_t)$, the likelihood is $p(y_t|\alpha_t)$ and we obtain the posterior density as $p(\alpha_t|Y_t) \sim N(a_t^*, P_t^*)$, where a_t^* and P_t^* are given by (2.9) and (2.10), respectively. We have used the regression approach because it is particularly simple and direct for deriving filtering and smoothing recursions for the general state space model in Chapter 4.

2.2.2 ILLUSTRATION

In this chapter we shall illustrate the algorithms using observations from the river Nile. The data set consists of a series of readings of the annual flow volume at

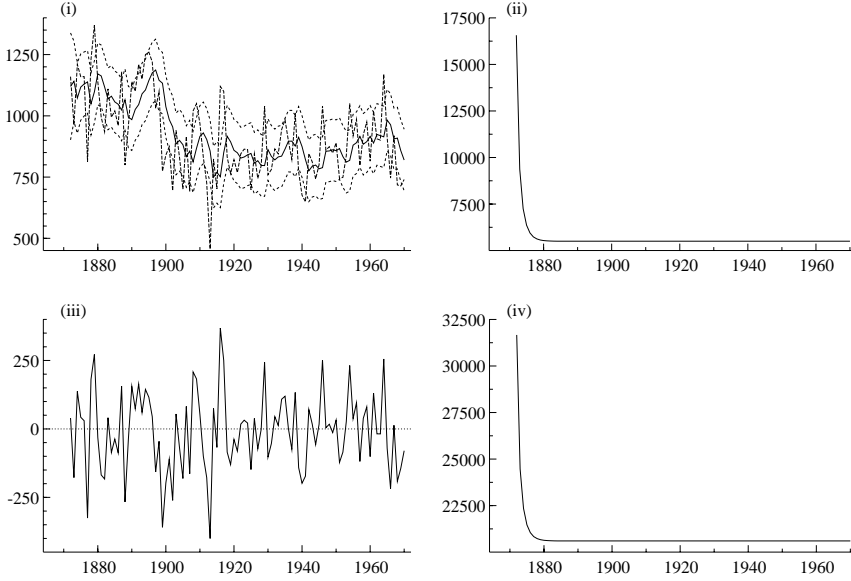


Fig. 2.1. Nile data and output of Kalman filter: (i) filtered state a_t and its 90% confidence intervals; (ii) filtered state variance P_t ; (iii) prediction errors v_t ; (iv) prediction variance F_t .

Aswan from 1871 to 1970. The series has been analysed by Cobb (1978) and Balke (1993). We analyse the data using the local level model (2.3) with $a_1 = 0$, $P_1 = 10^7$, $\sigma_\varepsilon^2 = 15099$ and $\sigma_\eta^2 = 1469.1$. The values for a_1 and P_1 were chosen arbitrarily for illustrative purposes. The values for σ_ε^2 and σ_η^2 are the maximum likelihood estimates which we obtain in §2.10.3. The output of the Kalman filter (that is v_t , F_t , a_t and P_t for $t = 2, \dots, n$) is presented graphically together with the raw data in Figure 2.1.

The most obvious feature of the four graphs is that F_t and P_t converge rapidly to constant values which confirms that the local level model has a steady state solution, that is, P_t and F_t converged to fixed values; for discussion of the concept of a steady state see §2.11. However, it was found that this local level model converged numerically to a steady state in around 25 updates of P_t although the graph of P_t seems to suggest that the steady state was obtained after around 10 updates.

2.3 Forecast errors

The Kalman filter residual $v_t = y_t - a_t$ and its variance F_t are the one-step-ahead forecast error and the one-step-ahead forecast error variance of y_t given Y_{t-1} as defined in §2.2. The forecast errors v_1, \dots, v_n are sometimes called *innovations* because they represent the new part of y_t that cannot be predicted from the past for $t = 1, \dots, n$. We shall make use of v_t and F_t for a variety of results in the next sections. It is therefore important to study them in detail.

2.3.1 CHOLESKY DECOMPOSITION

First we show that v_1, \dots, v_n are mutually independent. The joint density of y_1, \dots, y_n is

$$p(y_1, \dots, y_n) = p(y_1) \prod_{t=2}^n p(y_t | Y_{t-1}). \quad (2.12)$$

Transform from y_1, \dots, y_n to v_1, \dots, v_n . Since each v_t equals y_t minus a linear function of y_1, \dots, y_{t-1} for $t = 2, \dots, n$, the Jacobian is one. From (2.12) on making the substitution we have

$$p(v_1, \dots, v_n) = \prod_{t=1}^n p(v_t), \quad (2.13)$$

since $p(v_1) = p(y_1)$ and $p(v_t) = p(y_t | Y_{t-1})$ for $t = 2, \dots, n$. Consequently, the v_t 's are independently distributed.

We next show that the forecast errors v_t are effectively obtained from a Cholesky decomposition of the observation vector y . The Kalman filter recursions compute the forecast error v_t as a linear function of the initial mean a_1 and the observations y_1, \dots, y_t since

$$\begin{aligned} v_1 &= y_1 - a_1, \\ v_2 &= y_2 - a_1 - K_1(y_1 - a_1), \\ v_3 &= y_3 - a_1 - K_2(y_2 - a_1) - K_1(1 - K_2)(y_1 - a_1), \quad \text{and so on.} \end{aligned}$$

It should be noted that K_t does not depend on the initial mean a_1 and the observations y_1, \dots, y_n ; it depends only on the initial state variance P_1 and the disturbance variances σ_ε^2 and σ_η^2 . Using the definitions in (2.4), we have

$$v = C(y - 1a_1), \quad \text{with} \quad v = \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix},$$

where matrix C is the lower triangular matrix

$$C = \begin{bmatrix} 1 & 0 & 0 & 0 \\ c_{21} & 1 & 0 & 0 \\ c_{31} & c_{32} & 1 & 0 \\ & & \ddots & \vdots \\ c_{n1} & c_{n2} & c_{n3} & \cdots & 1 \end{bmatrix},$$

$$\begin{aligned} c_{i,i-1} &= -K_{i-1}, \\ c_{ij} &= -(1 - K_{i-1})(1 - K_{i-2}) \cdots (1 - K_{j+1})K_j, \end{aligned} \quad (2.14)$$

for $i = 2, \dots, n$ and $j = 1, \dots, i - 2$. The distribution of v is therefore

$$v \sim N(0, C\Omega C'), \quad (2.15)$$

where $\Omega = \text{Var}(y)$ as given by (2.4). On the other hand we know from (2.11) and (2.13) that $E(v_t) = 0$, $\text{Var}(v_t) = F_t$ and $\text{Cov}(v_t, v_j) = 0$, for $t, j = 1, \dots, n$ and $t \neq j$; therefore,

$$v \sim N(0, F), \quad \text{with} \quad F = \begin{bmatrix} F_1 & 0 & 0 & & 0 \\ 0 & F_2 & 0 & & 0 \\ 0 & 0 & F_3 & & 0 \\ & & & \ddots & \vdots \\ 0 & 0 & 0 & \dots & F_n \end{bmatrix},$$

where $C\Omega C' = F$. The transformation of a symmetric positive definite matrix (say Ω) into a diagonal matrix (say F) using a lower triangular matrix (say C) by means of the relation $C\Omega C' = F$ is known as the Cholesky decomposition of the symmetric matrix. The Kalman filter can therefore be regarded as essentially a Cholesky decomposition of the variance matrix implied by the local level model (2.3). This result is important for understanding the role of the Kalman filter and it will be used further in §§2.5.4 and 2.10.1. Note also that $F^{-1} = (C')^{-1}\Omega^{-1}C^{-1}$ so we have $\Omega^{-1} = C'F^{-1}C$.

2.3.2 ERROR RECURSIONS

Define the *state estimation error* as

$$x_t = \alpha_t - a_t, \quad \text{with} \quad \text{Var}(x_t) = P_t. \quad (2.16)$$

We now show that the state estimation errors x_t and forecast errors v_t are linear functions of the initial state error x_1 and the disturbances ε_t and η_t analogously to the way that α_t and y_t are linear functions of the initial state and the disturbances, for $t = 1, \dots, n$. It follows directly from the Kalman filter relations (2.11) that

$$\begin{aligned} v_t &= y_t - a_t \\ &= \alpha_t + \varepsilon_t - a_t \\ &= x_t + \varepsilon_t, \end{aligned}$$

and

$$\begin{aligned} x_{t+1} &= \alpha_{t+1} - a_{t+1} \\ &= \alpha_t + \eta_t - a_t - K_t v_t \\ &= x_t + \eta_t - K_t(x_t + \varepsilon_t) \\ &= L_t x_t + \eta_t - K_t \varepsilon_t, \end{aligned}$$

where

$$L_t = 1 - K_t = \sigma_\varepsilon^2 / F_t. \quad (2.17)$$

Thus analogously to the local level model relations

$$y_t = \alpha_t + \varepsilon_t, \quad \alpha_{t+1} = \alpha_t + \eta_t,$$

we have the error relations

$$v_t = x_t + \varepsilon_t, \quad x_{t+1} = L_t x_t + \eta_t - K_t \varepsilon_t, \quad t = 1, \dots, n, \quad (2.18)$$

with $x_1 = \alpha_1 - a_1$. These relations will be used in the next section. We note that P_t , F_t , K_t and L_t do not depend on the initial state mean a_1 or the observations y_1, \dots, y_n but only on the initial state variance P_1 and the disturbance variances σ_ε^2 and σ_η^2 . We note also that the recursion for P_{t+1} in (2.11) can alternatively be derived by

$$\begin{aligned} P_{t+1} &= \text{Var}(x_{t+1}) = \text{Cov}(x_{t+1}, \alpha_{t+1}) = \text{Cov}(x_{t+1}, \alpha_t + \eta_t) \\ &= L_t \text{Cov}(x_t, \alpha_t + \eta_t) + \text{Cov}(\eta_t, \alpha_t + \eta_t) - K_t \text{Cov}(\varepsilon_t, \alpha_t + \eta_t) \\ &= L_t P_t + \sigma_\eta^2. \end{aligned}$$

2.4 State smoothing

We now consider the estimation of $\alpha_1, \dots, \alpha_n$ given the entire sample Y_n . We shall find it convenient to use in place of the collective symbol Y_n its representation as the vector $y = (y_1, \dots, y_n)'$ defined in (2.4). Since all distributions are normal, the conditional density of α_t given y is $N(\hat{\alpha}_t, V_t)$ where $\hat{\alpha}_t = E(\alpha_t|y)$ and $V_t = \text{Var}(\alpha_t|y)$. We call $\hat{\alpha}_t$ the *smoothed state*, V_t the *smoothed state variance* and the operation of calculating $\hat{\alpha}_1, \dots, \hat{\alpha}_n$ *state smoothing*.

2.4.1 SMOOTHED STATE

The forecast errors v_1, \dots, v_n are mutually independent and are a linear transformation of y_1, \dots, y_n , and v_t, \dots, v_n are independent of y_1, \dots, y_{t-1} with zero means. Moreover, when y_1, \dots, y_n are fixed, Y_{t-1} and v_t, \dots, v_n are fixed and vice versa. By (2.49) of the regression lemma in §2.13 we therefore have

$$\begin{aligned} \hat{\alpha}_t &= E(\alpha_t|y) = E(\alpha_t|Y_{t-1}, v_t, \dots, v_n) \\ &= E(\alpha_t|Y_{t-1}) + \text{Cov}[\alpha_t, (v_t, \dots, v_n)'] \text{Var}[(v_t, \dots, v_n)']^{-1} (v_t, \dots, v_n)' \\ &= a_t + \begin{pmatrix} \text{Cov}(\alpha_t, v_t) \\ \vdots \\ \text{Cov}(\alpha_t, v_n) \end{pmatrix}' \begin{bmatrix} F_t & 0 \\ & \ddots \\ 0 & F_n \end{bmatrix}^{-1} \begin{pmatrix} v_t \\ \vdots \\ v_n \end{pmatrix} \\ &= a_t + \sum_{j=t}^n \text{Cov}(\alpha_t, v_j) F_j^{-1} v_j. \end{aligned} \tag{2.19}$$

Now $\text{Cov}(\alpha_t, v_j) = \text{Cov}(x_t, v_j)$ for $j = t, \dots, n$, and

$$\begin{aligned} \text{Cov}(x_t, v_t) &= E[x_t(x_t + \varepsilon_t)] = \text{Var}(x_t) = P_t, \\ \text{Cov}(x_t, v_{t+1}) &= E[x_t(x_{t+1} + \varepsilon_{t+1})] = E[x_t(L_t x_t + \eta_t - K_t \varepsilon_t)] = P_t L_t. \end{aligned}$$

Similarly,

$$\begin{aligned} \text{Cov}(x_t, v_{t+2}) &= P_t L_t L_{t+1}, \\ &\vdots \\ \text{Cov}(x_t, v_n) &= P_t L_t L_{t+1} \dots L_{n-1}. \end{aligned} \tag{2.20}$$

Substituting in (2.19) gives

$$\begin{aligned}\hat{\alpha}_t &= a_t + P_t \frac{v_t}{F_t} + P_t L_t \frac{v_{t+1}}{F_{t+1}} + P_t L_t L_{t+1} \frac{v_{t+2}}{F_{t+2}} + \cdots \\ &= a_t + P_t r_{t-1},\end{aligned}$$

where

$$\begin{aligned}r_{t-1} &= \frac{v_t}{F_t} + L_t \frac{v_{t+1}}{F_{t+1}} + L_t L_{t+1} \frac{v_{t+2}}{F_{t+2}} + L_t L_{t+1} L_{t+2} \frac{v_{t+3}}{F_{t+3}} + \cdots \\ &\quad + L_t L_{t+1} \cdots L_{n-1} \frac{v_n}{F_n}\end{aligned}\tag{2.21}$$

is a weighted sum of innovations after $t - 1$. The value of this at time t is

$$\begin{aligned}r_t &= \frac{v_{t+1}}{F_{t+1}} + L_{t+1} \frac{v_{t+2}}{F_{t+2}} + L_{t+1} L_{t+2} \frac{v_{t+3}}{F_{t+3}} + \cdots \\ &\quad + L_{t+1} L_{t+2} \cdots L_{n-1} \frac{v_n}{F_n}.\end{aligned}\tag{2.22}$$

Obviously, $r_n = 0$ since no observations are available after time n . By substituting from (2.22) into (2.21), it follows that the values of r_{t-1} can be evaluated using the backwards recursion

$$r_{t-1} = \frac{v_t}{F_t} + L_t r_t,\tag{2.23}$$

with $r_n = 0$, for $t = n, n - 1, \dots, 1$. The smoothed state can therefore be calculated by the backwards recursion

$$r_{t-1} = F_t^{-1} v_t + L_t r_t, \quad \hat{\alpha}_t = a_t + P_t r_{t-1}, \quad t = n, \dots, 1,\tag{2.24}$$

with $r_n = 0$. The relations in (2.24) are collectively called the *state smoothing recursion*.

2.4.2 SMOOTHED STATE VARIANCE

The error variance of the smoothed state, $V_t = \text{Var}(\alpha_t | Y_n)$, is derived in a similar way. Using the properties of the innovations and the regression lemma in §2.13 with $x = \alpha_t$, $y = (y_1, \dots, y_{t-1})'$ and $z = (v_t, \dots, v_n)'$ in (2.50), we have

$$\begin{aligned}V_t &= \text{Var}(\alpha_t | y) = \text{Var}(\alpha_t | Y_{t-1}, v_t, \dots, v_n) \\ &= \text{Var}(\alpha_t | Y_{t-1}) - \text{Cov}[\alpha_t, (v_t, \dots, v_n)'] \text{Var}[(v_t, \dots, v_n)']^{-1} \\ &\quad \times \text{Cov}[\alpha_t, (v_t, \dots, v_n)'] \\ &= P_t - \begin{pmatrix} \text{Cov}(\alpha_t, v_t) \\ \vdots \\ \text{Cov}(\alpha_t, v_n) \end{pmatrix}' \begin{bmatrix} F_t & 0 \\ & \ddots \\ 0 & F_n \end{bmatrix}^{-1} \begin{pmatrix} \text{Cov}(\alpha_t, v_t) \\ \vdots \\ \text{Cov}(\alpha_t, v_n) \end{pmatrix} \\ &= P_t - \sum_{j=t}^n [\text{Cov}(\alpha_t, v_j)]^2 F_j^{-1},\end{aligned}\tag{2.25}$$

where the expressions for $\text{Cov}(\alpha_t, v_j)$ are given by (2.20). Substituting these into (2.25) leads to

$$\begin{aligned} V_t &= P_t - P_t^2 \frac{1}{F_t} - P_t^2 L_t^2 \frac{1}{F_{t+1}} - P_t^2 L_t^2 L_{t+1}^2 \frac{1}{F_{t+2}} - \cdots - P_t^2 L_t^2 L_{t+1}^2 \cdots L_{n-1}^2 \frac{1}{F_n} \\ &= P_t - P_t^2 N_{t-1}, \end{aligned} \quad (2.26)$$

where

$$\begin{aligned} N_{t-1} &= \frac{1}{F_t} + L_t^2 \frac{1}{F_{t+1}} + L_t^2 L_{t+1}^2 \frac{1}{F_{t+2}} + L_t^2 L_{t+1}^2 L_{t+2}^2 \frac{1}{F_{t+3}} + \cdots \\ &\quad + L_t^2 L_{t+1}^2 \cdots L_{n-1}^2 \frac{1}{F_n}, \end{aligned} \quad (2.27)$$

is a weighted sum of the inverse variances of innovations after time $t - 1$. Its value at time t is

$$N_t = \frac{1}{F_{t+1}} + L_{t+1}^2 \frac{1}{F_{t+2}} + L_{t+1}^2 L_{t+2}^2 \frac{1}{F_{t+3}} + \cdots + L_{t+1}^2 L_{t+2}^2 \cdots L_{n-1}^2 \frac{1}{F_n}, \quad (2.28)$$

and, obviously, $N_n = 0$ since no variances are available after time n . Substituting from (2.28) into (2.27) it follows that the value for N_{t-1} can be calculated using the backwards recursion

$$N_{t-1} = \frac{1}{F_t} + L_t^2 N_t, \quad (2.29)$$

with $N_n = 0$, for $t = n, n - 1, \dots, 1$. We see from (2.22) and (2.28) that $N_t = \text{Var}(r_t)$ since the forecast errors v_t are independent.

Combining these results, the error variance of the smoothed state can be calculated by the backwards recursion

$$N_{t-1} = F_t^{-1} + L_t^2 N_t, \quad V_t = P_t - P_t^2 N_{t-1}, \quad t = n, \dots, 1, \quad (2.30)$$

with $N_n = 0$. The relations in (2.30) are collectively called the *state variance smoothing recursion*. From the standard error $\sqrt{V_t}$ of $\hat{\alpha}_t$ we can construct confidence intervals for α_t for $t = 1, \dots, n$. It is also possible to derive the smoothed covariances between the states, that is, $\text{Cov}(\alpha_t, \alpha_s | y)$, $t \neq s$, using similar arguments. We shall not give them here but will derive them for the general case in §4.5.

2.4.3 ILLUSTRATION

We now show the results of state smoothing for the Nile data of §2.2.2 using the same local level model. The Kalman filter is applied first and the output v_t , F_t , a_t and P_t is stored for $t = 1, \dots, n$. Figure 2.2 presents the output of the backwards smoothing recursions (2.24) and (2.30); that is r_t , N_t , $\hat{\alpha}_t$ and V_t . The plot of $\hat{\alpha}_t$ includes the confidence bands for α_t . The graph of $\text{Var}(\alpha_t | y)$ shows that the

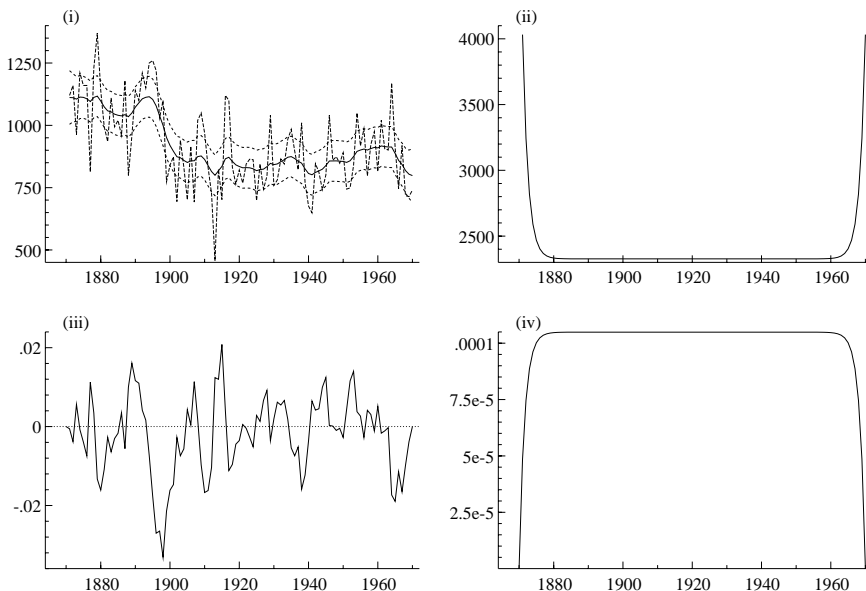


Fig. 2.2. Nile data and output of state smoothing recursion: (i) smoothed state $\hat{\alpha}_t$ and its 90% confidence intervals; (ii) smoothed state variance V_t ; (iii) smoothing cumulant r_t ; (iv) smoothing variance cumulant N_t .

conditional variance of α_t is larger at the beginning and end of the sample, as it obviously should be on intuitive grounds. Comparing the graphs of α_t and $\hat{\alpha}_t$ in Figures 2.1 and 2.2, we see that the graph of $\hat{\alpha}_t$ is much smoother than that of α_t , except at time points close to the end of the series, as it should be.

2.5 Disturbance smoothing

In this section we consider the calculation of the smoothed observation disturbance $\hat{\varepsilon}_t = E(\varepsilon_t|y) = y_t - \hat{\alpha}_t$ and the smoothed state disturbance $\hat{\eta}_t = E(\eta_t|y) = \hat{\alpha}_{t+1} - \hat{\alpha}_t$ together with their error variances. Of course, these could be calculated directly from a knowledge of $\hat{\alpha}_1, \dots, \hat{\alpha}_n$ and covariances $\text{Cov}(\alpha_t, \alpha_j|y)$ for $j \leq t$. However, it turns out to be computationally advantageous to compute them from r_t and N_t without first calculating $\hat{\alpha}_t$, particularly for the general model discussed in Chapter 4. The merits of smoothed disturbances are discussed in §4.4. For example, the estimates $\hat{\varepsilon}_t$ and $\hat{\eta}_t$ are useful for detecting outliers and structural breaks, respectively; see §2.12.2.

In order to economise on the amount of algebra in this chapter we shall present the required recursions for the local level model without proof, referring the reader to §4.4 for derivations of the analogous recursions for the general model.

2.5.1 SMOOTHED OBSERVATION DISTURBANCES

From (4.36) in §4.4.1, the smoothed observation disturbance $\hat{\varepsilon}_t = E(\varepsilon_t|y)$ is calculated by

$$\hat{\varepsilon}_t = \sigma_\varepsilon^2 u_t, \quad t = n, \dots, 1, \quad (2.31)$$

where

$$u_t = F_t^{-1} v_t - K_t r_t, \quad (2.32)$$

and where the recursion for r_t is given by (2.23). The scalar u_t is referred to as the *smoothing error*. Similarly, from (4.44) in §4.4.3, the smoothed variance $\text{Var}(\varepsilon_t|y)$ is obtained by

$$\text{Var}(\varepsilon_t|y) = \sigma_\varepsilon^2 - \sigma_\varepsilon^4 D_t, \quad t = n, \dots, 1, \quad (2.33)$$

where

$$D_t = F_t^{-1} + K_t^2 N_t, \quad (2.34)$$

and where the recursion for N_t is given by (2.29). Since from (2.22) v_t is independent of r_t , and $\text{Var}(r_t) = N_t$, we have

$$\text{Var}(u_t) = \text{Var}(F_t^{-1} v_t - K_t r_t) = F_t^{-2} \text{Var}(v_t) + K_t^2 \text{Var}(r_t) = D_t.$$

Consequently, from (2.31) we obtain $\text{Var}(\hat{\varepsilon}_t) = \sigma_\varepsilon^4 D_t$.

Note that the methods for calculating $\hat{\alpha}_t$ and $\hat{\varepsilon}_t$ are consistent since $K_t = P_t F_t^{-1}$, $L_t = 1 - K_t = \sigma_\varepsilon^2 F_t^{-1}$ and

$$\begin{aligned} \hat{\varepsilon}_t &= y_t - \hat{\alpha}_t \\ &= y_t - a_t - P_t r_{t-1} \\ &= v_t - P_t (F_t^{-1} v_t + L_t r_t) \\ &= F_t^{-1} v_t (F_t - P_t) - \sigma_\varepsilon^2 P_t F_t^{-1} r_t \\ &= \sigma_\varepsilon^2 (F_t^{-1} v_t - K_t r_t), \quad t = n, \dots, 1. \end{aligned}$$

Similar equivalences can be shown for V_t and $\text{Var}(\varepsilon_t|y)$.

2.5.2 SMOOTHED STATE DISTURBANCES

From (4.41) in §4.4.1, the smoothed mean of the disturbance $\hat{\eta}_t = E(\eta_t|y)$ is calculated by

$$\hat{\eta}_t = \sigma_\eta^2 r_t, \quad t = n, \dots, 1, \quad (2.35)$$

where the recursion for r_t is given by (2.23). Similarly, from (4.47) in §4.4.3, the smoothed variance $\text{Var}(\eta_t|y)$ is computed by

$$\text{Var}(\eta_t|y) = \sigma_\eta^2 - \sigma_\eta^4 N_t, \quad t = n, \dots, 1, \quad (2.36)$$

where the recursion for N_t is given by (2.29). Since $\text{Var}(r_t) = N_t$, we have that $\text{Var}(\hat{\eta}_t) = \sigma_\eta^4 N_t$. These results are interesting because they give an interpretation to the values r_t and N_t ; they are the scaled smoothed estimator of $\eta_t = \alpha_{t+1} - \alpha_t$ and its unconditional variance, respectively.

The method of calculating $\hat{\eta}_t$ is consistent with the definition $\eta_t = \alpha_{t+1} - \alpha_t$ since

$$\begin{aligned}\hat{\eta}_t &= \hat{\alpha}_{t+1} - \hat{\alpha}_t \\ &= a_{t+1} + P_{t+1}r_t - a_t - P_t r_{t-1} \\ &= a_t + K_t v_t - a_t + P_t L_t r_t + \sigma_\eta^2 r_t - P_t (F_t^{-1} v_t + L_t r_t) \\ &= \sigma_\eta^2 r_t.\end{aligned}$$

Similar consistencies can be shown for N_t and $\text{Var}(\eta_t|y)$.

2.5.3 ILLUSTRATION

The smoothed disturbances and their related variances for the Nile data and the local level model of §2.2.2 are calculated by the above recursions and presented in Figure 2.3. We note from the graphs of $\text{Var}(\varepsilon_t|y)$ and $\text{Var}(\eta_t|y)$ the extent that these conditional variances are larger at the beginning and end of the sample. Obviously, the plot of r_t in Figure 2.2 and the plot of $\hat{\eta}_t$ in Figure 2.3 are the same apart from a different scale.

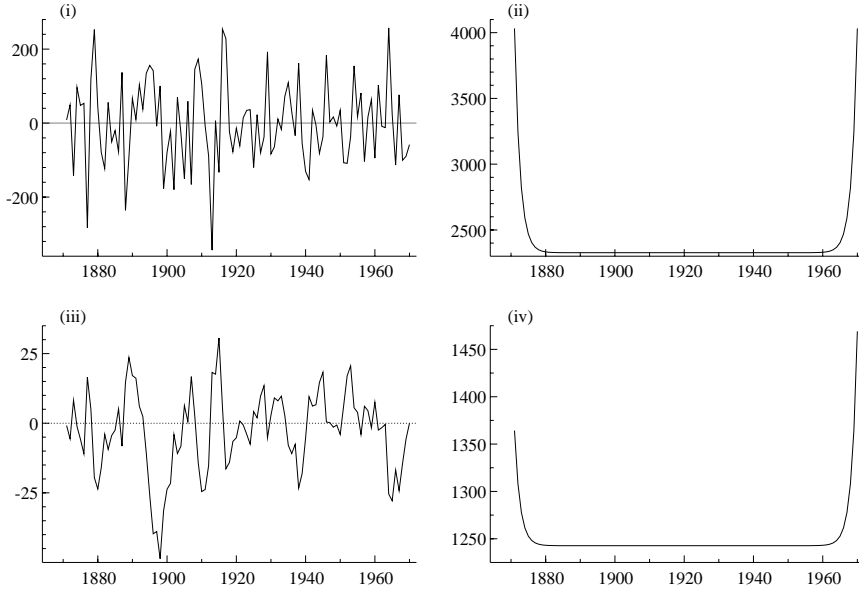


Fig. 2.3. Output of disturbance smoothing recursion: (i) observation error ε_t ; (ii) observation error variance $\text{Var}(\varepsilon_t|y)$; (iii) state error $\hat{\eta}_t$; (iv) state error variance $\text{Var}(\eta_t|y)$.

2.5.4 CHOLESKY DECOMPOSITION AND SMOOTHING

We now consider the calculation of $\hat{\varepsilon}_t = E(\varepsilon_t|y)$ by direct regression of $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)'$ on the observation vector y to obtain $\hat{\varepsilon} = (\hat{\varepsilon}_1, \dots, \hat{\varepsilon}_n)'$, that is,

$$\begin{aligned}\hat{\varepsilon} &= E(\varepsilon) + \text{Cov}(\varepsilon, y) \text{Var}(y)^{-1}[y - E(y)] \\ &= \text{Cov}(\varepsilon, y)\Omega^{-1}(y - 1a_1).\end{aligned}$$

It is obvious from (2.6) that $\text{Cov}(\varepsilon, y) = \sigma_\varepsilon^2 I_n$; also, from the Cholesky decomposition considered in §2.3.1 we have $\Omega^{-1} = C'F^{-1}C$ and $C(y - 1a_1) = v$. We therefore have

$$\hat{\varepsilon} = \sigma_\varepsilon^2 C'F^{-1}v,$$

which, by consulting the definitions of the lower triangular elements of C in (2.14), also leads to the disturbance equations (2.31) and (2.32). Thus

$$\hat{\varepsilon} = \sigma_\varepsilon^2 u, \quad u = \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix},$$

where

$$u = C'F^{-1}v \quad \text{with} \quad v = C(y - 1a_1).$$

It follows that

$$u = C'F^{-1}C(y - 1a_1) = \Omega^{-1}(y - 1a_1), \quad (2.37)$$

where $\Omega = \text{Var}(y)$ and $F = C\Omega C'$, as is consistent with standard regression theory.

2.6 Simulation

It is simple to draw samples generated by the local level model (2.3). We first draw the random normal deviates

$$\varepsilon_t^{(\cdot)} \sim N(0, \sigma_\varepsilon^2), \quad \eta_t^{(\cdot)} \sim N(0, \sigma_\eta^2), \quad t = 1, \dots, n.$$

Then we generate observations using the local level recursion as follows

$$y_t^{(\cdot)} = \alpha_t^{(\cdot)} + \varepsilon_t^{(\cdot)}, \quad \alpha_{t+1}^{(\cdot)} = \alpha_t^{(\cdot)} + \eta_t^{(\cdot)}, \quad t = 1, \dots, n,$$

for some starting value $\alpha_1^{(\cdot)}$.

For certain applications, which will be mainly discussed in Part II of this book, we may require samples generated by the local level model conditional on the observed time series y_1, \dots, y_n . Such samples can be obtained by use of the simulation smoother developed for the general linear Gaussian state space model by de Jong and Shephard (1995) which we derive in §4.7. For the local level model, a simulated sample for the disturbances ε_t , $t = 1, \dots, n$, given the observations

y_1, \dots, y_n can be obtained using (4.77) by the backwards recursion

$$\begin{aligned}\tilde{\varepsilon}_t &= d_t + \sigma_\varepsilon^2(v_t/F_t - K_t\tilde{r}_t), \\ \tilde{r}_{t-1} &= v_t/F_t - \tilde{W}_td_t/C_t + L_t\tilde{r}_t,\end{aligned}$$

where $d_t \sim N(0, C_t)$ with

$$\begin{aligned}C_t &= \sigma_\varepsilon^2 - \sigma_\varepsilon^4(1/F_t + K_t^2\tilde{N}_t), \\ \tilde{W}_t &= \sigma_\varepsilon^2(1/F_t - K_t\tilde{N}_tL_t), \\ \tilde{N}_{t-1} &= 1/F_t + \tilde{W}_t^2/C_t + L_t^2\tilde{N}_t,\end{aligned}$$

for $t = n, \dots, 1$ with $\tilde{r}_n = \tilde{N}_n = 0$. The quantities v_t, F_t, K_t and $L_t = 1 - K_t$ are obtained from the Kalman filter and they need to be stored for $t = 1, \dots, n$. We note that the recursions for \tilde{r}_t and \tilde{N}_t are similar to the recursions for r_t and N_t given by (2.23) and (2.29); in fact, they are equivalent if $\tilde{W}_t = 0$ for $t = 1, \dots, n$.

Given a sample for $\varepsilon_t, t = 1, \dots, n$, we obtain simulated samples for α_t and η_t via the relations

$$\begin{aligned}\tilde{\alpha}_t &= y_t - \tilde{\varepsilon}_t, & t &= 1, \dots, n, \\ \tilde{\eta}_t &= \tilde{\alpha}_{t+1} - \tilde{\alpha}_t, & t &= 1, \dots, n-1.\end{aligned}$$

2.6.1 ILLUSTRATION

To illustrate the difference between simulating a sample from the local level model unconditionally and simulating a sample conditional on the observations, we consider the Nile data and the local level model of §2.2.2. In Figure 2.4 (i) we present the smoothed state $\hat{\alpha}_t$ and a sample generated by the local level model unconditionally. The two series have seemingly nothing in common. In the next panel, again the smoothed state is presented but now together with a sample generated conditional on the observations. Here we see that the generated sample is much closer to $\hat{\alpha}_t$. The remaining two panels present the smoothed disturbances together with a sample from the corresponding disturbances conditional on the observations.

2.7 Missing observations

A considerable advantage of the state space approach is the ease with which missing observations can be dealt with. Suppose we have a local level model where observations y_j , with $j = \tau, \dots, \tau^* - 1$, are missing for $1 < \tau < \tau^* \leq n$. For the filtering stage, the most obvious way to deal with the situation is to define a new series y_t^* where $y_t^* = y_t$ for $t = 1, \dots, \tau - 1$ and $y_t^* = y_{t+\tau^*-\tau}$ for $t = \tau, \dots, n^*$ with $n^* = n - (\tau^* - \tau)$. The model for y_t^* with time scale $t = 1, \dots, n^*$ is then the same as (2.3) with $y_t = y_t^*$ except that $\alpha_\tau = \alpha_{\tau-1} + \eta_{\tau-1}$ where $\eta_{\tau-1} \sim N[0, (\tau^* - \tau)\sigma_\eta^2]$. Filtering for this model can be treated by the methods developed in Chapter 4 for the general state space model. The treatment is readily extended if more than one group of observations is missing.

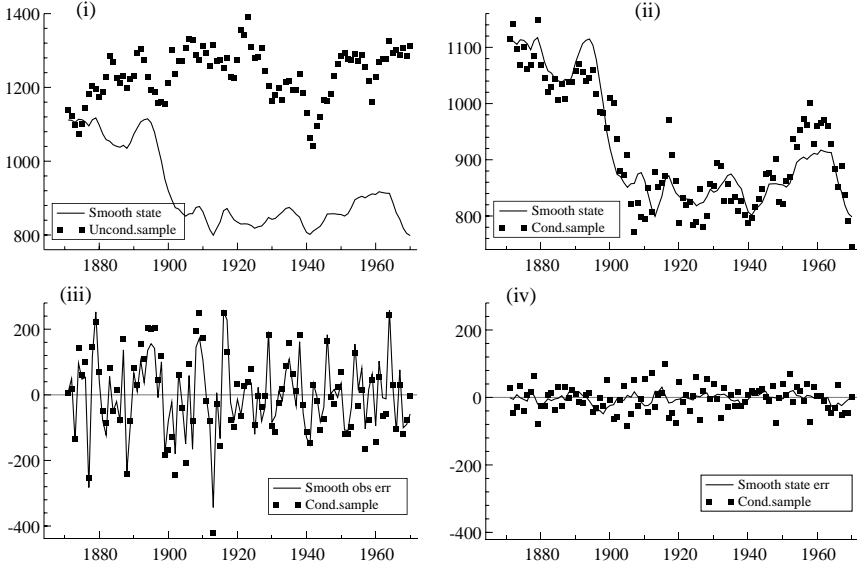


Fig. 2.4. Simulation: (i) smoothed state $\hat{\alpha}_t$ and sample $\alpha_t^{(\cdot)}$; (ii) smoothed state $\hat{\alpha}_t$ and sample $\tilde{\alpha}_t$; (iii) smoothed observation error $\hat{\varepsilon}_t$ and sample $\tilde{\varepsilon}_t$; (iv) smoothed state error $\hat{\eta}_t$ and sample $\tilde{\eta}_t$.

It is, however, easier and more transparent to proceed as follows, using the original time domain. For filtering at times $t = \tau, \dots, \tau^* - 1$, we have

$$E(\alpha_t | Y_{t-1}) = E(\alpha_t | Y_{\tau-1}) = E\left(\alpha_\tau + \sum_{j=\tau}^{t-1} \eta_j \middle| Y_{\tau-1}\right) = a_\tau$$

and

$$\text{Var}(\alpha_t | Y_{t-1}) = \text{Var}(\alpha_t | Y_{\tau-1}) = \text{Var}\left(\alpha_\tau + \sum_{j=\tau}^{t-1} \eta_j \middle| Y_{\tau-1}\right) = P_\tau + (t - \tau)\sigma_\eta^2.$$

giving

$$a_{t+1} = a_t, \quad P_{t+1} = P_t + \sigma_\eta^2, \quad t = \tau, \dots, \tau^* - 1, \quad (2.38)$$

the remaining values a_t and P_t being given as before by (2.11) for $t = 1, \dots, \tau$ and $t = \tau^*, \dots, n$. The consequence is that we can use the original filter (2.11) for all t by taking $v_t = 0$ and $K_t = 0$ at the missing time points. The same procedure is used when more than one group of observations is missing. It follows that allowing for missing observations when using the Kalman filter is extremely simple.

The forecast error recursions from which we derive the smoothing recursions are given by (2.18). These error-updating equations at the missing time points

become

$$v_t = x_t + \varepsilon_t, \quad x_{t+1} = x_t + \eta_t, \quad t = \tau, \dots, \tau^* - 1,$$

since $K_t = 0$ and therefore $L_t = 1$. The covariances between the state at the missing time points and the innovations after the missing period are given by

$$\text{Cov}(\alpha_t, v_{\tau^*}) = P_t,$$

$$\text{Cov}(\alpha_t, v_j) = P_t L_{\tau^*} L_{\tau^*+1} \dots L_{j-1}, \quad j = \tau^* + 1, \dots, n, \quad t = \tau, \dots, \tau^* - 1.$$

By deleting the terms associated with the missing time points, the state smoothing equation (2.19) for the missing time points becomes

$$\hat{\alpha}_t = a_t + \sum_{j=\tau^*}^n \text{Cov}(\alpha_t, v_j) F_j^{-1} v_j, \quad t = \tau, \dots, \tau^* - 1.$$

Substituting the covariance terms into this and taking into account the definition (2.21) leads directly to

$$r_{t-1} = r_t, \quad \hat{\alpha}_t = a_t + P_t r_{t-1}, \quad t = \tau, \dots, \tau^* - 1. \quad (2.39)$$

Again, the consequence is that we can use the original state smoother (2.24) for all t by taking $v_t = 0$ and $K_t = 0$, and hence $L_t = 1$, at the missing time points. This device applies to any missing observation within the sample period. In the same way the equations for the variance of the state error and the smoothed disturbances can be obtained by putting $v_t = 0$ and $K_t = 0$ at missing time points.

2.7.1 ILLUSTRATION

Here we consider the Nile data and the same local level model as before; however, we treat the observations at time points 21, ..., 40 and 61, ..., 80 as missing. The Kalman filter is applied first and the output v_t , F_t , a_t and P_t is stored for $t = 1, \dots, n$. Then, the state smoothing recursions are applied. The first two graphs in Figure 2.5 are the Kalman filter values of a_t and P_t , respectively. The last two graphs are the smoothing output $\hat{\alpha}_t$ and V_t , respectively.

Note that the application of the Kalman filter to missing observations can be regarded as extrapolation of the series to the missing time points, while smoothing at these points is effectively interpolation.

2.8 Forecasting

Let \bar{y}_{n+j} be the minimum mean square error forecast of y_{n+j} given the time series y_1, \dots, y_n for $j = 1, 2, \dots, J$ with J as some pre-defined positive integer. By minimum mean square error forecast here we mean the function \bar{y}_{n+j} of y_1, \dots, y_n which minimises $E[(y_{n+j} - \bar{y}_{n+j})^2 | Y_n]$. Then $\bar{y}_{n+j} = E(y_{n+j} | Y_n)$. This follows immediately from the well-known result that if x is a random variable with mean

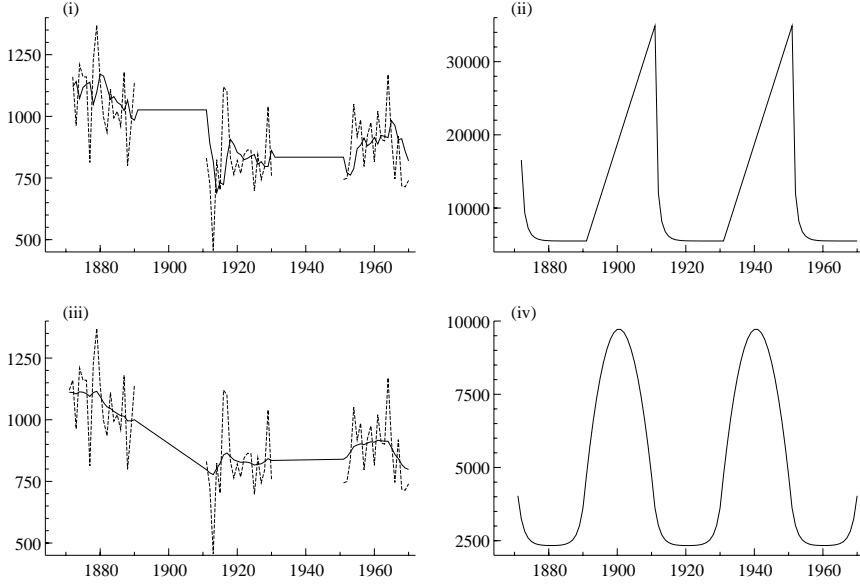


Fig. 2.5. Filtering and smoothing output when observations are missing: (i) filtered state a_t (extrapolation); (ii) filtered state variance P_t ; (iii) smoothed state $\hat{\alpha}_t$ (interpolation); (iv) smoothed state variance V_t .

μ the value of λ that minimises $E(x - \lambda)^2$ is $\lambda = \mu$. The variance of the forecast error is denoted by $\bar{F}_{n+j} = \text{Var}(y_{n+j}|Y_n)$. The theory of forecasting for the local level model turns out to be surprisingly simple; we merely regard forecasting as filtering the observations $y_1, \dots, y_n, y_{n+1}, \dots, y_{n+J}$ using the recursion (2.11) and treating the last J observations y_{n+1}, \dots, y_{n+J} as missing.

Letting $\bar{a}_{n+j} = E(\alpha_{n+j}|Y_n)$ and $\bar{P}_{n+j} = \text{Var}(\alpha_{n+j}|Y_n)$, it follows immediately from equation (2.38) with $\tau = n + 1$ and $\tau^* = n + J$ in §2.7 that

$$\bar{a}_{n+j+1} = \bar{a}_{n+j}, \quad \bar{P}_{n+j+1} = \bar{P}_{n+j} + \sigma_\eta^2, \quad j = 1, \dots, J - 1,$$

with $\bar{a}_{n+1} = a_{n+1}$ and $\bar{P}_{n+1} = P_{n+1}$ obtained from the Kalman filter (2.11). Furthermore, we have

$$\begin{aligned} \bar{y}_{n+j} &= E(y_{n+j}|Y_n) = E(\alpha_{n+j}|Y_n) + E(\varepsilon_{n+j}|Y_n) = \bar{a}_{n+j}, \\ \bar{F}_{n+j} &= \text{Var}(y_{n+j}|Y_n) = \text{Var}(\alpha_{n+j}|Y_n) + \text{Var}(\varepsilon_{n+j}|Y_n) = \bar{P}_{n+j} + \sigma_\varepsilon^2, \end{aligned}$$

for $j = 1, \dots, J$. The consequence is that the Kalman filter can be applied for $t = 1, \dots, n + J$ where we treat the observations at times $n + 1, \dots, n + J$ as missing. Thus we conclude that forecasts and their error variances are delivered by applying the Kalman filter in a routine way with $v_t = 0$ and $K_t = 0$ for

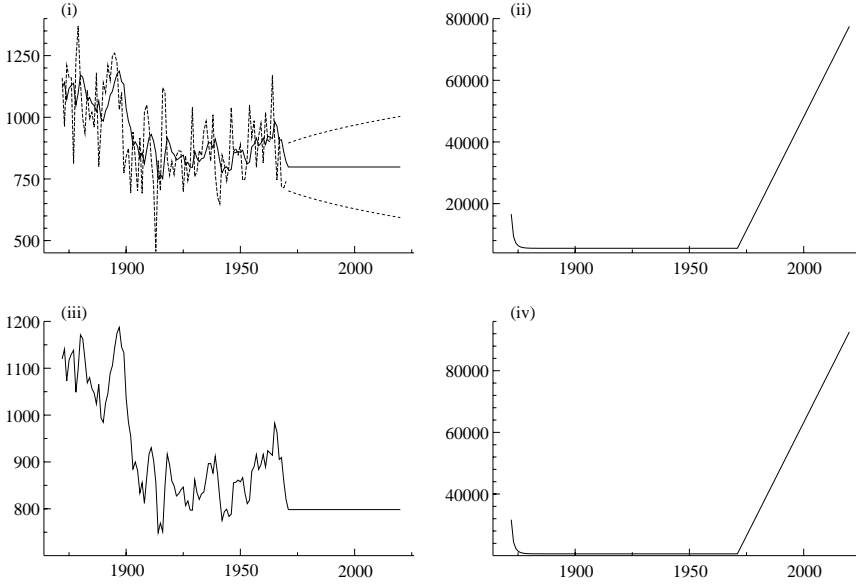


Fig. 2.6. Nile data and output of forecasting: (i) state forecast a_t and 50% confidence intervals; (ii) state variance P_t ; (iii) observation forecast $E(y_t|Y_{t-1})$; (iv) observation forecast variance F_t .

$t = n + 1, \dots, n + J$. The same property holds for the general linear Gaussian state space model as we shall show in §4.9.

2.8.1 ILLUSTRATION

The Nile data set is now extended by 30 missing observations allowing the computation of forecasts for the observations y_{101}, \dots, y_{130} . The Kalman filter only is required. The graphs in Figure 2.6 contain $\hat{y}_{n+j|n} = a_{n+j|n}$, $P_{n+j|n}$, $a_{n+j|n}$ and $F_{n+j|n}$, respectively, for $j = 1, \dots, J$ with $J = 30$. The confidence interval for $E(y_{n+j}|y)$ is $\hat{y}_{n+j|n} \pm k\sqrt{F_{n+j|n}}$ where k is determined by the required probability of inclusion; in Figure 2.6 this probability is 50%.

2.9 Initialisation

We assumed in previous sections that the distribution of the initial state α_1 is $N(a_1, P_1)$ where a_1 and P_1 are known. We now consider how to start up the filter (2.11) when nothing is known about the distribution of α_1 , which is the usual situation in practice. In this situation it is reasonable to represent α_1 as having a *diffuse prior density*, that is, fix a_1 at an arbitrary value and let $P_1 \rightarrow \infty$. From (2.11) we have

$$v_1 = y_1 - a_1, \quad F_1 = P_1 + \sigma_\varepsilon^2,$$

and, by substituting into the equations for a_2 and P_2 in (2.11), it follows that

$$a_2 = a_1 + \frac{P_1}{P_1 + \sigma_\varepsilon^2}(y_1 - a_1), \quad (2.40)$$

$$\begin{aligned} P_2 &= P_1 \left(1 - \frac{P_1}{P_1 + \sigma_\varepsilon^2} \right) + \sigma_\eta^2 \\ &= \frac{P_1}{P_1 + \sigma_\varepsilon^2} \sigma_\varepsilon^2 + \sigma_\eta^2. \end{aligned} \quad (2.41)$$

Letting $P_1 \rightarrow \infty$, we obtain $a_2 = y_1$, $P_2 = \sigma_\varepsilon^2 + \sigma_\eta^2$ and we can then proceed normally with the Kalman filter (2.11) for $t = 2, \dots, n$. This process is called *diffuse initialisation* of the Kalman filter and the resulting filter is called *the diffuse Kalman filter*. We note the interesting fact that the same values of a_t and P_t for $t = 2, \dots, n$ can be obtained by treating y_1 as fixed and taking $\alpha_1 \sim N(y_1, \sigma_\varepsilon^2)$. This is intuitively reasonable in the absence of information about the marginal distribution of α_1 since $(y_1 - \alpha_1) \sim N(0, \sigma_\varepsilon^2)$.

We also need to take account of the diffuse distribution of the initial state α_1 in the smoothing recursions. It is shown above that the filtering equations for $t = 2, \dots, n$ are not affected by letting $P_1 \rightarrow \infty$. Therefore, the state and disturbance smoothing equations are also not affected for $t = n, \dots, 2$ since these only depend on the Kalman filter output. From (2.24), the smoothed mean of the state α_1 is given by

$$\begin{aligned} \hat{\alpha}_1 &= a_1 + P_1 \left[\frac{1}{P_1 + \sigma_\varepsilon^2} v_1 + \left(1 - \frac{P_1}{P_1 + \sigma_\varepsilon^2} \right) r_1 \right], \\ &= a_1 + \frac{P_1}{P_1 + \sigma_\varepsilon^2} v_1 + \frac{P_1}{P_1 + \sigma_\varepsilon^2} \sigma_\varepsilon^2 r_1. \end{aligned}$$

Letting $P_1 \rightarrow \infty$, we obtain $\hat{\alpha}_1 = a_1 + v_1 + \sigma_\varepsilon^2 r_1$ and by substituting for v_1 we have

$$\hat{\alpha}_1 = y_1 + \sigma_\varepsilon^2 r_1.$$

The smoothed conditional variance of the state α_1 given y is, from (2.30)

$$\begin{aligned} V_1 &= P_1 - P_1^2 \left[\frac{1}{P_1 + \sigma_\varepsilon^2} + \left(1 - \frac{P_1}{P_1 + \sigma_\varepsilon^2} \right)^2 N_1 \right] \\ &= P_1 \left(1 - \frac{P_1}{P_1 + \sigma_\varepsilon^2} \right) - \left(\frac{P_1}{P_1 + \sigma_\varepsilon^2} \right)^2 \sigma_\varepsilon^4 N_1 \\ &= \left(\frac{P_1}{P_1 + \sigma_\varepsilon^2} \right) \sigma_\varepsilon^2 - \left(\frac{P_1}{P_1 + \sigma_\varepsilon^2} \right)^2 \sigma_\varepsilon^4 N_1. \end{aligned}$$

Letting $P_1 \rightarrow \infty$, we obtain $V_1 = \sigma_\varepsilon^2 - \sigma_\varepsilon^4 N_1$.

The smoothed means of the disturbances for $t = 1$ are given by

$$\hat{\varepsilon}_1 = \sigma_\varepsilon^2 u_1, \quad \text{with} \quad u_1 = \frac{1}{P_1 + \sigma_\varepsilon^2} v_1 - \frac{P_1}{P_1 + \sigma_\varepsilon^2} r_1,$$

and $\hat{\eta}_1 = \sigma_\eta^2 r_1$. Letting $P_1 \rightarrow \infty$, we obtain $\hat{\varepsilon}_1 = -\sigma_\varepsilon^2 r_1$. Note that r_1 depends on the Kalman filter output for $t = 2, \dots, n$. The smoothed variances of the disturbances for $t = 1$ depend on D_1 and N_1 of which only D_1 is affected by $P_1 \rightarrow \infty$; using (2.34),

$$D_1 = \frac{1}{P_1 + \sigma_\varepsilon^2} + \left(\frac{P_1}{P_1 + \sigma_\varepsilon^2} \right)^2 N_1.$$

Letting $P_1 \rightarrow \infty$, we obtain $D_1 = N_1$ and therefore $\text{Var}(\hat{\varepsilon}_1) = \sigma_\varepsilon^4 N_1$. The variance of the smoothed estimate of η_1 remains unaltered as $\text{Var}(\hat{\eta}_1) = \sigma_\eta^4 N_1$.

The initial smoothed state $\hat{\alpha}_1$ under diffuse conditions can also be obtained by assuming that y_1 is fixed and $\alpha_1 = y_1 - \varepsilon_1$ where $\varepsilon_1 \sim N(0, \sigma_\varepsilon^2)$. For example, for the smoothed mean of the state at $t = 1$, we have now only $n - 1$ varying y_t 's so that

$$\hat{\alpha}_1 = a_1 + \sum_{j=2}^n \frac{\text{Cov}(\alpha_1, v_j)}{F_j} v_j$$

with $a_1 = y_1$. It follows from (2.40) that $a_2 = a_1 = y_1$. Further, $v_2 = y_2 - a_2 = \alpha_2 + \varepsilon_2 - y_1 = \alpha_1 + \eta_1 + \varepsilon_2 - y_1 = -\varepsilon_1 + \eta_1 + \varepsilon_2$. Consequently, $\text{Cov}(\alpha_1, v_2) = \text{Cov}(-\varepsilon_1, -\varepsilon_1 + \eta_1 + \varepsilon_2) = \sigma_\varepsilon^2$. We therefore have from (2.19),

$$\begin{aligned} \hat{\alpha}_1 &= a_1 + \frac{\sigma_\varepsilon^2}{F_2} v_2 + \frac{(1 - K_2)\sigma_\varepsilon^2}{F_3} v_3 + \frac{(1 - K_2)(1 - K_3)\sigma_\varepsilon^2}{F_4} v_4 + \dots \\ &= y_1 + \sigma_\varepsilon^2 r_1, \end{aligned}$$

as before with r_1 as defined in (2.21) for $t = 1$. The equations for the remaining $\hat{\alpha}_t$'s are the same as previously.

Use of a diffuse prior for initialisation is the approach preferred by most time series analysts in the situation where nothing is known about the initial value α_1 . However, some workers find the diffuse approach uncongenial because they regard the assumption of an infinite variance as unnatural since all observed time series have finite values. From this point of view an alternative approach is to assume that α_1 is an unknown constant to be estimated from the data by maximum likelihood. The simplest form of this idea is to estimate α_1 by maximum likelihood from the first observation y_1 . Denote this maximum likelihood estimate by $\hat{\alpha}_1$ and its variance by $\text{Var}(\hat{\alpha}_1)$. We then initialise the Kalman filter by taking $a_1 = \hat{\alpha}_1$ and $P_1 = \text{Var}(\hat{\alpha}_1)$. Since when α_1 is fixed $y_1 \sim N(\alpha_1, \sigma_\varepsilon^2)$, we have $\hat{\alpha}_1 = y_1$ and $\text{Var}(\hat{\alpha}_1) = \sigma_\varepsilon^2$. We therefore initialise the filter by taking $a_1 = y_1$ and $P_1 = \sigma_\varepsilon^2$. But these are the same values as we obtain by assuming that α_1 is diffuse. It follows that we obtain the same initialisation of the Kalman filter by representing α_1 as a random variable with infinite variance as by assuming that it is fixed and unknown

and estimating it from y_1 . We shall show that a similar result holds for the general linear Gaussian state space model in §5.7.3.

2.10 Parameter estimation

We now consider the fitting of the local level model to data from the standpoint of classical inference. In effect, this amounts to deriving formulae on the assumption that the parameters are known and then replacing these by their maximum likelihood estimates. Bayesian treatment will be considered for the general linear Gaussian model in Chapter 8. Parameters in state space models are often called *hyperparameters*, possibly to distinguish them from elements of state vectors which can plausibly be thought of as random parameters; however, in this book we shall just call them *parameters*, since with the usual meaning of the word parameter this is what they are. We will discuss methods for calculating the loglikelihood function and the maximisation of it with respect to the parameters, σ_ε^2 and σ_η^2 .

2.10.1 LOGLIKELIHOOD EVALUATION

Since

$$p(y_1, \dots, y_t) = p(Y_{t-1})p(y_t|Y_{t-1}),$$

for $t = 2, \dots, n$, the joint density of y_1, \dots, y_n can be expressed as

$$p(y) = \prod_{t=1}^n p(y_t|Y_{t-1}),$$

where $p(y_1|Y_0) = p(y_1)$. Now $p(y_t|Y_{t-1}) = N(a_t, F_t)$ and $v_t = y_t - a_t$ so on taking logs and assuming that a_1 and P_1 are known the loglikelihood is given by

$$\log L = \log p(y) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{t=1}^n \left(\log F_t + \frac{v_t^2}{F_t} \right). \quad (2.42)$$

The exact loglikelihood can therefore be constructed easily from the Kalman filter (2.11).

Alternatively, let us derive the loglikelihood for the local level model from the representation (2.4). This gives

$$\log L = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log|\Omega| - \frac{1}{2} (y - a_1 1)' \Omega^{-1} (y - a_1 1),$$

which follows from the multivariate normal distribution $y \sim N(a_1 1, \Omega)$. Since $\Omega = CFC'$, $|C| = 1$, $\Omega^{-1} = C'F^{-1}C$ and $v = C(y - a_1 1)$, it follows that

$$\log|\Omega| = \log|CFC'| = \log|C||F||C| = \log|F|,$$

and

$$(y - a_1 1)' \Omega^{-1} (y - a_1 1) = v' F^{-1} v.$$

Substitution and using the results $\log|F| = \sum_{t=1}^n \log F_t$ and $v'F^{-1}v = \sum_{t=1}^n F_t^{-1}v_t^2$ lead directly to (2.42).

The loglikelihood in the diffuse case is derived as follows. All terms in (2.42) remain finite as $P_1 \rightarrow \infty$ with y fixed except the term for $t = 1$. It thus seems reasonable to remove the influence of P_1 as $P_1 \rightarrow \infty$ by defining the *diffuse loglikelihood* as

$$\begin{aligned} \log L_d &= \lim_{P_1 \rightarrow \infty} \left(\log L + \frac{1}{2} \log P_1 \right) \\ &= -\frac{1}{2} \lim_{P_1 \rightarrow \infty} \left(\log \frac{F_1}{P_1} + \frac{v_1^2}{F_1} \right) - \frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{t=2}^n \left(\log F_t + \frac{v_t^2}{F_t} \right) \\ &= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{t=2}^n \left(\log F_t + \frac{v_t^2}{F_t} \right), \end{aligned} \quad (2.43)$$

since $F_1/P_1 \rightarrow 1$ and $v_1^2/F_1 \rightarrow 0$ as $P_1 \rightarrow \infty$. Note that v_t and F_t remain finite as $P_1 \rightarrow \infty$ for $t = 2, \dots, n$.

Since P_1 does not depend on σ_ε^2 and σ_η^2 , the values of σ_ε^2 and σ_η^2 that maximise $\log L$ are identical to the values that maximise $\log L + \frac{1}{2} \log P_1$. As $P_1 \rightarrow \infty$, these latter values converge to the values that maximise $\log L_d$ because first and second derivatives with respect to σ_ε^2 and σ_η^2 converge, and second derivatives are finite and strictly negative. It follows that the maximum likelihood estimators of σ_ε^2 and σ_η^2 obtained by maximising (2.42) converge to the values obtained by maximising (2.43) as $P_1 \rightarrow \infty$.

We estimate the unknown parameters σ_ε^2 and σ_η^2 by maximising expression (2.42) or (2.43) numerically according to whether a_1 and P_1 are known or unknown. In practice it is more convenient to maximise numerically with respect to the quantities $\psi_\varepsilon = \log \sigma_\varepsilon^2$ and $\psi_\eta = \log \sigma_\eta^2$. An efficient algorithm for numerical maximisation is implemented in the *STAMP* 6.0 package of Koopman, Harvey, Doornik and Shephard (2000). This optimisation procedure is based on the quasi-Newton scheme BFGS for which details are given in §7.3.2.

2.10.2 CONCENTRATION OF LOGLIKELIHOOD

It can be advantageous to re-parameterise the model prior to maximisation in order to reduce the dimensionality of the numerical search. For example, for the local level model we can put $q = \sigma_\eta^2/\sigma_\varepsilon^2$ to obtain the model

$$\begin{aligned} y_t &= \alpha_t + \varepsilon_t, & \varepsilon_t &\sim N(0, \sigma_\varepsilon^2), \\ \alpha_{t+1} &= \alpha_t + \eta_t, & \eta_t &\sim N(0, q\sigma_\varepsilon^2), \end{aligned}$$

and estimate the pair σ_ε^2, q in preference to $\sigma_\varepsilon^2, \sigma_\eta^2$. Put $P_t^* = P_t/\sigma_\varepsilon^2$ and $F_t^* = F_t/\sigma_\varepsilon^2$; from (2.11) and §2.9, the diffuse Kalman filter for the re-parameterised

Table 2.1. Estimation of parameters of local level model by maximum likelihood.

Iteration	q	ψ	Score	Log-likelihood
0	1	0	-3.32	-495.68
1	0.0360	-3.32	0.93	-492.53
2	0.0745	-2.60	0.25	-492.10
3	0.0974	-2.32	-0.001	-492.07
4	0.0973	-2.33	0.0	-492.07

local level model is then

$$\begin{aligned} v_t &= y_t - a_t, & F_t^* &= P_t^* + 1, \\ a_{t+1} &= a_t + K_t v_t, & P_{t+1}^* &= P_t^*(1 - K_t) + q, \end{aligned}$$

where $K_t = P_t^*/F_t^* = P_t/F_t$ for $t = 2, \dots, n$ and it is initialised with $a_2 = y_1$ and $P_2^* = 1 + q$. Note that F_t^* depends on q but not on σ_ε^2 . The loglikelihood (2.43) then becomes

$$\log L_d = -\frac{n}{2} \log(2\pi) - \frac{n-1}{2} \log \sigma_\varepsilon^2 - \frac{1}{2} \sum_{t=2}^n \left(\log F_t^* + \frac{v_t^2}{\sigma_\varepsilon^2 F_t^*} \right). \quad (2.44)$$

By maximising (2.44) with respect to σ_ε^2 , for given F_2^*, \dots, F_n^* , we obtain

$$\hat{\sigma}_\varepsilon^2 = \frac{1}{n-1} \sum_{t=2}^n \frac{v_t^2}{F_t^*}. \quad (2.45)$$

The value of $\log L_d$ obtained by substituting $\hat{\sigma}_\varepsilon^2$ for σ_ε^2 in (2.44) is called the *concentrated diffuse loglikelihood* and is denoted by $\log L_{dc}$, giving

$$\log L_{dc} = -\frac{n}{2} \log(2\pi) - \frac{n-1}{2} - \frac{n-1}{2} \log \hat{\sigma}_\varepsilon^2 - \frac{1}{2} \sum_{t=2}^n \log F_t^*. \quad (2.46)$$

This is maximised with respect to q by a one-dimensional numerical search.

2.10.3 ILLUSTRATION

The estimates of the variances σ_ε^2 and $\sigma_\eta^2 = q\sigma_\varepsilon^2$ for the Nile data are obtained by maximising the concentrated diffuse loglikelihood (2.46) with respect to ψ where $q = \exp(\psi)$. In Table 2.1 the iterations of the BFGS procedure are reported starting with $\psi = 0$. The relative percentage change of the loglikelihood goes down very rapidly and convergence is achieved after 4 iterations. The final estimate for ψ is -2.33 and hence the estimate of q is $\hat{q} = 0.097$. The estimate of σ_ε^2 given by (2.45) is 15099 which implies that the estimate of σ_η^2 is $\hat{\sigma}_\eta^2 = \hat{q}\hat{\sigma}_\varepsilon^2 = 0.097 \times 15099 = 1469.1$.

2.11 Steady state

We now consider whether the Kalman filter (2.11) converges to a *steady state* as $n \rightarrow \infty$. This will be the case if P_t converges to a positive value, \bar{P} say. Obviously,

we would then have $F_t \rightarrow \bar{P} + \sigma_\varepsilon^2$ and $K_t \rightarrow \bar{P}/(\bar{P} + \sigma_\varepsilon^2)$. To check whether there is a steady state, put $P_{t+1} = P_t = \bar{P}$ in (2.11) and verify whether the resulting equation in \bar{P} has a positive solution. The equation is

$$\bar{P} = \bar{P} \left(1 - \frac{\bar{P}}{\bar{P} + \sigma_\varepsilon^2} \right) + \sigma_\eta^2,$$

which reduces to the quadratic

$$x^2 - xh - h = 0, \quad (2.47)$$

where $x = \bar{P}/\sigma_\varepsilon^2$ and $h = \sigma_\eta^2/\sigma_\varepsilon^2$, with the solution

$$x = (h + \sqrt{h^2 + 4h})/2.$$

This is positive when $h > 0$ which holds for non-trivial models. The other solution to (2.47) is inapplicable since it is negative for $h > 0$. Thus all non-trivial local level models have a steady state solution.

The practical advantage of knowing that a model has a steady state solution is that, after convergence of P_t to \bar{P} has been verified to be close enough, we can stop computing F_t and K_t and the filter (2.11) reduces to the single relation

$$a_{t+1} = a_t + \bar{K} v_t,$$

with $\bar{K} = \bar{P}/(\bar{P} + \sigma_\varepsilon^2)$ and $v_t = y_t - a_t$. While this has little consequence for the simple local level model we are concerned with, it is a useful property for the more complicated models we shall consider in Chapter 4, where P_t can be a large matrix.

2.12 Diagnostic checking

2.12.1 DIAGNOSTIC TESTS FOR FORECAST ERRORS

The assumptions underlying the local level model are that the disturbances ε_t and η_t are normally distributed and serially independent with constant variances. On these assumptions the standardised one-step forecast errors

$$e_t = \frac{v_t}{\sqrt{F_t}}, \quad t = 1, \dots, n, \quad (2.48)$$

(or for $t = 2, \dots, n$ in the diffuse case) are also normally distributed and serially independent with unit variance. We can check that these properties hold by means of the following large-sample diagnostic tests:

- Normality

The first four moments of the standardised forecast errors are given by

$$m_1 = \frac{1}{n} \sum_{t=1}^n e_t,$$

$$m_q = \frac{1}{n} \sum_{t=1}^n (e_t - m_1)^q, \quad q = 2, 3, 4,$$

with obvious modifications in the diffuse case. Skewness and kurtosis are denoted by S and K , respectively, and are defined as

$$S = \frac{m_3}{\sqrt{m_2^3}}, \quad K = \frac{m_4}{m_2^2},$$

and it can be shown that when the model assumptions are valid they are asymptotically normally distributed as

$$S \sim N\left(0, \frac{6}{n}\right), \quad K \sim N\left(3, \frac{24}{n}\right);$$

see Bowman and Shenton (1975). Standard statistical tests can be used to check whether the observed values of S and K are consistent with their asymptotic densities. They can also be combined as

$$N = n \left\{ \frac{S^2}{6} + \frac{(K - 3)^2}{24} \right\},$$

which asymptotically has a χ^2 distribution with 2 degrees of freedom on the null hypothesis that the normality assumption is valid. The *QQ plot* is a graphical display of ordered residuals against their theoretical quantiles. The 45 degree line is taken as a reference line (the closer the residual plot to this line, the better the match).

- Heteroscedasticity

A simple test for heteroscedasticity is obtained by comparing the sum of squares of two exclusive subsets of the sample. For example, the statistic

$$H(h) = \frac{\sum_{t=n-h+1}^n e_t^2}{\sum_{t=1}^h e_t^2},$$

is $F_{h,h}$ -distributed for some preset positive integer h , under the null hypothesis of homoscedasticity. Here, e_t is defined in (2.48) and the sum of h squared forecast errors in the denominator starts at $t = 2$ in the diffuse case.

- Serial correlation

When the local level model holds, the standardised forecast errors are serially uncorrelated as we have shown in §2.3.1. Therefore, the correlogram of the forecast errors should reveal serial correlation insignificant. A standard portmanteau test statistic for serial correlation is based on the Box-Ljung statistic suggested by Ljung and Box (1978). This is given by

$$Q(k) = n(n+2) \sum_{j=1}^k \frac{c_j^2}{n-j},$$

for some preset positive integer k where c_j is the j th correlogram value

$$c_j = \frac{1}{nm_2} \sum_{t=j+1}^n (e_t - m_1)(e_{t-j} - m_1).$$

More details on diagnostic checking will be given in §7.5.

2.12.2 DETECTION OF OUTLIERS AND STRUCTURAL BREAKS

The standardised smoothed residuals are given by

$$u_t^* = \hat{\varepsilon}_t / \sqrt{\text{Var}(\hat{\varepsilon}_t)} = D_t^{-\frac{1}{2}} u_t,$$

$$r_t^* = \hat{\eta}_t / \sqrt{\text{Var}(\hat{\eta}_t)} = N_t^{-\frac{1}{2}} r_t, \quad t = 1, \dots, n;$$

see §2.5 for details on computing the quantities u_t , D_t , r_t and N_t . Harvey and Koopman (1992) refer to these standardised residuals as *auxiliary residuals* and they investigate their properties in detail. For example, they show that the auxiliary residuals are autocorrelated and they discuss their autocorrelation function. The auxiliary residuals can be useful in detecting outliers and structural breaks in time series because $\hat{\varepsilon}_t$ and $\hat{\eta}_t$ are estimators of ε_t and η_t . An outlier in a series that we postulate as generated by the local level model is indicated by a large (positive or negative) value for ε_t and a break in the level α_t is indicated by a large (positive or negative) value for η_t . A discussion of use of auxiliary residuals for the general model will be given in §7.5.

2.12.3 ILLUSTRATION

We consider the fitted local level model for the Nile data as obtained in §2.10.3. A plot of e_t is given in Figure 2.7 together with the histogram, the QQ plot and the correlogram. These plots are satisfactory and they suggest that the assumptions

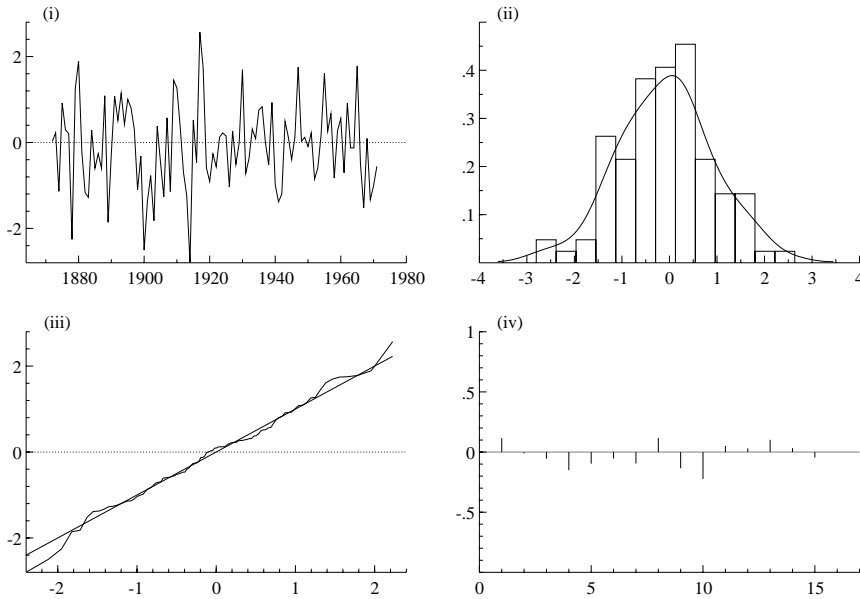


Fig. 2.7. Diagnostic plots for standardised prediction errors: (i) standardised residual; (ii) histogram plus estimated density; (iii) ordered residuals; (iv) correlogram.

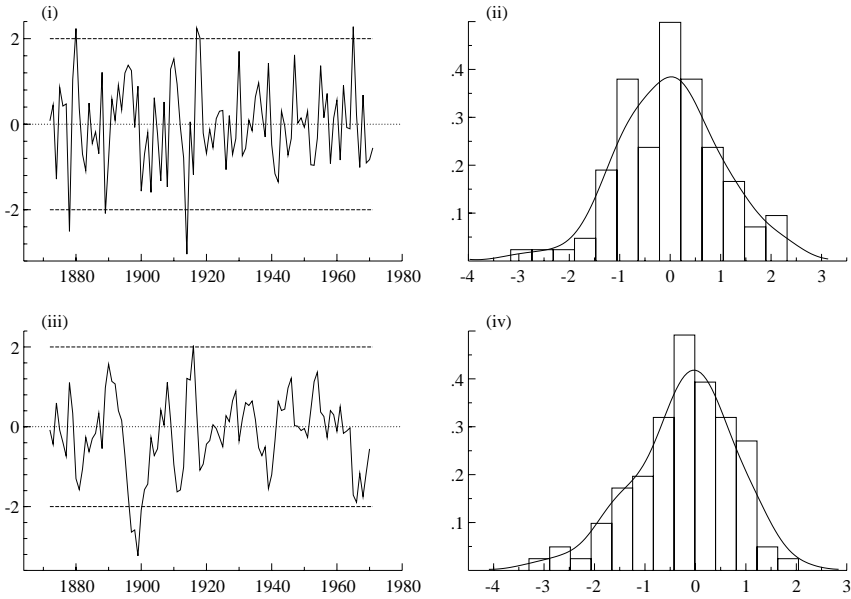


Fig. 2.8. Diagnostic plots for auxiliary residuals: (i) observation residual u_t^* ; (ii) histogram and estimated density for u_t^* ; (iii) state residual r_t^* ; (iv) histogram and estimated density for r_t^* .

underlying the local level model are valid for the Nile data. This is largely confirmed by the following diagnostic test statistics

$$S = -0.03, \quad K = 0.09, \quad N = 0.05, \quad H(33) = 0.61, \quad Q(9) = 8.84.$$

The low value for the heteroscedasticity statistic H indicates a degree of heteroscedasticity in the residuals. This is apparent in the plots of u_t^* and r_t^* together with their histograms in Figure 2.8. These diagnostic plots indicate outliers in 1913 and 1918 and a level break in 1899. The plot of the Nile data confirms these findings.

2.13 Appendix: Lemma in multivariate normal regression

We present here a simple lemma in multivariate normal regression theory which we use extensively in the book to derive results in filtering, smoothing and related problems.

Suppose that x , y and z are random vectors of arbitrary orders that are jointly normally distributed with means μ_p and covariance matrices $\Sigma_{pq} = E[(p - \mu_p)(q - \mu_q)']$ for $p, q = x, y$ and z with $\mu_z = 0$ and $\Sigma_{yz} = 0$. The symbols x , y , z , p and q are employed for convenience and their use here is unrelated to their use in other parts of the book.

Lemma

$$E(x|y, z) = E(x|y) + \Sigma_{xz} \Sigma_{zz}^{-1} z, \quad (2.49)$$

$$\text{Var}(x|y, z) = \text{Var}(x|y) - \Sigma_{xz} \Sigma_{zz}^{-1} \Sigma'_{xz}. \quad (2.50)$$

PROOF. By standard multivariate normal regression theory we have

$$E(x|y) = \mu_x + \Sigma_{xy} \Sigma_{yy}^{-1} (y - \mu_y), \quad (2.51)$$

$$\text{Var}(x|y) = \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma'_{xy}; \quad (2.52)$$

see, for example, Anderson (1984, Theorem 2.5.1).

Applying (2.51) to vector $\begin{pmatrix} y \\ z \end{pmatrix}$ in place of y gives

$$\begin{aligned} E(x|y, z) &= \mu_x + [\Sigma_{xy} \quad \Sigma_{xz}] \begin{bmatrix} \Sigma_{yy}^{-1} & 0 \\ 0 & \Sigma_{zz}^{-1} \end{bmatrix} \begin{pmatrix} y - \mu_y \\ z \end{pmatrix} \\ &= \mu_x + \Sigma_{xy} \Sigma_{yy}^{-1} (y - \mu_y) + \Sigma_{xz} \Sigma_{zz}^{-1} z, \end{aligned}$$

since $\mu_z = 0$ and $\Sigma_{yz} = 0$. This proves (2.49).

Applying (2.52) to vector $\begin{pmatrix} y \\ z \end{pmatrix}$ in place of y gives

$$\begin{aligned} \text{Var}(x|y, z) &= \Sigma_{xx} - [\Sigma_{xy} \quad \Sigma_{xz}] \begin{bmatrix} \Sigma_{yy}^{-1} & 0 \\ 0 & \Sigma_{zz}^{-1} \end{bmatrix} \begin{bmatrix} \Sigma'_{xy} \\ \Sigma'_{xz} \end{bmatrix} \\ &= \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma'_{xy} - \Sigma_{xz} \Sigma_{zz}^{-1} \Sigma'_{xz}, \end{aligned}$$

since $\Sigma_{yz} = 0$. This proves (2.50). This simple lemma provides the basis for the treatment of the Kalman filter and smoother in this book.