

DOCUMENTOS DE TRABAJO

BILTOKI

D.T. 2005.03

Multiple imputation of time series:
an application to the construction of historical price indexes

Fernando TUSELL PALMER

eman ta zabal zazu



Universidad Euskal Herriko
del País Vasco Unibertsitatea

Facultad de Ciencias Económicas.
Avda. Lehendakari Aguirre, 83
48015 BILBAO.

Documento de Trabajo BILTOKI DT2005.03

Editado por el Departamento de Economía Aplicada III (Econometría y Estadística)
de la Universidad del País Vasco.

Depósito Legal No.: BI-2408-05

ISSN: 1134-8984

Multiple imputation of time series: an application to the construction of historical price indexes

Fernando TUSELL*

Abstract

Time series in many areas of application, and notably in the social sciences, are frequently incomplete. This is particularly annoying when we need to have complete data, for instance to compute indexes as a weighted average of values from a number of time series; whenever a single datum is absent, the index cannot be computed. This paper proposes to deal with such situations by creating multiple completed trajectories, drawing on state space modelling of time series, the simulation smoother and multiple imputation ideas.

Keywords: Multiple imputation; Time series; Missing data; Kalman filter; Index computation

1 Introduction

Missing data is a pervasive problem, afflicting not only the social sciences but also the physical and medical sciences.

Missing data in cross-sectional data is handled in a variety of ways, many admittedly *ad hoc* and making sense only in particular cases: cold deck and hot deck imputation, using only complete records, replacing the mean for missing data, using all available (possible incomplete) records, etc.

It has long been recognized that a sounder, more principled approach is desirable, and considerable effort has been expended in this direction. Much

*Departamento de Estadística y Econometría. Facultad de CC.EE. y Empresariales, Avda. Lehendakari Aguirre, 83, 48015 BILBAO. E-mail: fernando.tusel1@ehu.es.

of it stems from the seminal work in [11] (enlarged second edition [12]). The monograph [16] describes in considerable detail a methodology to deal with missing data in the cross sectional data, and [15] provides a useful overview of the ideas on multiple imputation and its impact on statistical practice.

The literature dealing with missing data in (multiple) time series is nonetheless sparse. Missing data in time series is considered in [12]; conceptually, the problem can be handled in the same way as in cross sectional data. However, the problem is both harder and more pressing. Harder, because an additional level of complexity exists when dealing with multivariate time series: both contemporaneous and lagged relationships between components need to be considered when imputing a missing data point. More pressing, because strategies like using only complete data records are no longer feasible. With cross sectional data, discarding records with data missing completely at random (MCAR) has no other effect than reducing the available sample. In a time series, each record is unique: dropping it would leave us with a series with holes, unusable for many purposes.

In the last fifteen years, state space modelling of time series has seen wider acceptance and is now a well established tool in the kit of the applied statistician. A number of theoretical breakthroughs like the simulation smoother (cf. [8], [3], [4]) and Markov Chain Monte Carlo (see for instance [6]), along with ever increasing computing power at the desktop, have improved our chances of dealing with missing data in multivariate time series. We describe in this paper a possible approach.

The rest of the paper is organized as follows. Section 2 reviews some of the basic theory on state space models and Kalman filtering and smoothing that we will be using. Section 3 discusses some of the models we have found useful for the purposes of this paper. Section 4 shows an example and discusses some further applications and the issues that they raise.

2 State space models

Let $\{\mathbf{y}_t\}$ be a p -variate time series, observed (perhaps partially) at times $t = 1, \dots, n$. We are concerned with the imputation of the missing values.

We will assume $\{\mathbf{y}_t\}$ generated by an state space model (see, for instance, [1] or [5]):

$$\mathbf{y}_t = \mathbf{Z}_t \boldsymbol{\alpha}_t + \boldsymbol{\epsilon}_t \quad (1)$$

$$\boldsymbol{\alpha}_{t+1} = \mathbf{T}_t \boldsymbol{\alpha}_t + \boldsymbol{\eta}_t \quad (2)$$

where $\boldsymbol{\epsilon}_t \sim N(\mathbf{0}, \mathbf{H}_t)$ and $\boldsymbol{\eta}_t \sim N(\mathbf{0}, \mathbf{Q}_t)$. Matrices $\mathbf{Z}_t, \mathbf{T}_t, \mathbf{Q}_t$ and \mathbf{H}_t are in general time-varying, and may depend on parameters that need to be

estimated.

Let $\mathcal{Y}_s \stackrel{\text{def}}{=} \{\mathbf{y}_1, \dots, \mathbf{y}_s\}$, i.e. the section of the time series up to and including time s . Given $\mathbf{Z}_t, \mathbf{T}_t, \mathbf{Q}_t, \mathbf{H}_t$ and \mathcal{Y}_t , the Kalman filter (see for instance [1], [5]) gives the conditional mean value and covariance matrix of state vector $\boldsymbol{\alpha}_t$ at each point in time, $\mathbf{a}_{t|t-1} = E[\boldsymbol{\alpha}_t | \mathcal{Y}_{t-1}]$ and $\mathbf{P}_{t|t-1} = \text{Cov}(\boldsymbol{\alpha}_t | \mathcal{Y}_{t-1})$. Defining,

$$\mathbf{F}_t = (\mathbf{Z}_t \mathbf{P}_{t|t-1} \mathbf{Z}_t^\top + \mathbf{H}_t) \quad (3)$$

$$\mathbf{K}_t = \mathbf{T}_t \mathbf{P}_{t|t-1} \mathbf{Z}_t^\top \mathbf{F}_t^{-1} \quad (4)$$

$$\mathbf{L}_t = \mathbf{T}_t - \mathbf{K}_t \mathbf{Z}_t \quad (5)$$

$$\mathbf{v}_t = \mathbf{y}_t - \mathbf{Z}_t \mathbf{a}_{t|t-1} \quad (6)$$

$$\mathbf{M}_t = \mathbf{P}_{t|t-1} \mathbf{Z}_t^\top \quad (7)$$

the formulae for one-step-ahead updating are:

$$\mathbf{a}_{t|t-1} = \mathbf{T}_{t-1} \mathbf{a}_{t-1|t-2} + \mathbf{K}_{t-1} \mathbf{v}_{t-1} \quad (8)$$

$$\mathbf{P}_{t|t-1} = \mathbf{T}_{t-1} \mathbf{P}_{t-1|t-2} \mathbf{L}_{t-1}^\top + \mathbf{Q}_{t-1}. \quad (9)$$

In order to start the iteration, either it is assumed that $\boldsymbol{\alpha}_0 \sim N(\mathbf{a}_{0|-1}, \mathbf{P}_{0|-1})$ or a diffuse prior is used (see [5], § 5.2).

Similar algorithms, collectively known as *smoothers*, give $\hat{\boldsymbol{\alpha}}_t \stackrel{\text{def}}{=} E[\boldsymbol{\alpha}_t | \mathcal{Y}_n]$ and $\mathbf{V}_t \stackrel{\text{def}}{=} \text{Cov}(\boldsymbol{\alpha}_t | \mathcal{Y}_n)$, i.e. conditional on the full length of time series. For instance, defining

$$\mathbf{r}_{t-1} \stackrel{\text{def}}{=} \mathbf{Z}_{t-1}^\top \mathbf{F}_{t-1}^{-1} \mathbf{v}_t + \mathbf{L}_{t-1}^\top \mathbf{r}_t \quad (10)$$

$$\mathbf{N}_{t-1} \stackrel{\text{def}}{=} \mathbf{Z}_t^\top \mathbf{F}_t^{-1} \mathbf{Z}_t + \mathbf{L}_t^\top \mathbf{N}_t \mathbf{L}_t, \quad (11)$$

we have

$$\hat{\boldsymbol{\alpha}}_t = \mathbf{a}_{t|t-1} + \mathbf{P}_{t|t-1} \mathbf{r}_{t-1} \quad (12)$$

$$\mathbf{V}_t = \mathbf{P}_{t|t-1} - \mathbf{P}_{t|t-1} \mathbf{N}_{t-1} \mathbf{P}_{t|t-1}. \quad (13)$$

The iteration is initialized with $\mathbf{N}_n = 0$ and $\mathbf{r}_n = 0$ (see [5], § 4.3.3).

Assume that the system matrices $\mathbf{Z}_t, \mathbf{T}_t, \mathbf{Q}_t$ and \mathbf{H}_t , possibly depending on a parameter vector $\boldsymbol{\theta}$, are known. Algorithms referred to as *simulation smoothers* afford easy generation of trajectories of $\boldsymbol{\alpha}_t, \boldsymbol{\epsilon}_t$ and $\boldsymbol{\eta}_t$ conditional on both \mathcal{Y}_n and $\boldsymbol{\theta}$, by drawing from the distributions $p(\boldsymbol{\alpha}_t | \mathcal{Y}_n, \boldsymbol{\theta})$, $p(\boldsymbol{\epsilon}_t | \mathcal{Y}_n, \boldsymbol{\theta})$, $p(\boldsymbol{\eta}_t | \mathcal{Y}_n, \boldsymbol{\theta})$ (see [3] and a simpler algorithm in [4]).

In practice, $\mathbf{Z}_t, \mathbf{T}_t, \mathbf{Q}_t$ and \mathbf{H}_t are seldom known and need to be estimated, at least in part. If this is the case, the uncertainty introduced by

the use of estimated parameters in place of θ has to be accounted for. Very little work seems to have dealt with this issue: [17] gives some asymptotic results and [7] describes a technique to account for the influence of estimated parameters on the variance of $\hat{\alpha}_t$.

3 Models and strategy for imputation

There are no limitations in the choice of imputation models, other than the requirement to keep the number of estimated parameters down to a manageable size. Although, in principle, any model that fits reasonably well the data can be used, we have found simple structural models (see [9],[8] and [5] for instance) well suited for the task of imputation.

Local level full dimension multivariate model. Taking $\mathbf{T}_t = \mathbf{Z}_t = \mathbf{I}_p$ in equations (1)-(2) we have what may be the simplest multivariate model for \mathbf{y}_t : each component y_{it} ($1 \leq i \leq p$) fluctuates about a component α_{it} of α_t .

The choice of the covariance matrix \mathbf{Q}_t governs the degree of correlation among components α_{it} , α_{jt} , $i \neq j$. We can choose to have independent random walks (diagonal \mathbf{Q}_t), non-independent random walks (\mathbf{Q}_t with non null off-diagonal terms) or even a reduced rank model (\mathbf{Q}_t rank deficient; except for the possible influence of the prior distribution on α_1 this would be equivalent to a reduced dimension state vector).

Regarding \mathbf{H}_t , we can choose independent or correlated observation disturbances. We can also take $\mathbf{H}_t = \mathbf{0}$, effectively saying that the components α_{it} of the state can be observed without error, whenever the corresponding y_{it} is observed. In either case, interest normally centers in the generation of simulated trajectories $\tilde{\alpha}_t$ of the state vector.

Local level reduced rank multivariate model. An alternative to the full rank model consists in keeping the random walk dynamics for the vector state α_t while taking $p = \dim(\mathbf{y}_t) > \dim(\alpha_t) = s$. In that case, the observed time series $\{\mathbf{y}_t\}$ evolve as linear combinations of a small number s of unobserved component. In this case, matrix \mathbf{Z}_t will typically contain regression coefficients of the \mathbf{y}_t on the α_t . The model can be seen as a dynamic factor analysis.

Generation of imputed trajectories. A random sample of m trajectories from $(\alpha_t | \mathcal{Y}_n)$ can then be obtained as sketched in Algorithm 1.

The loop in steps 2 to 5 in the Algorithm 1 generates each time it is run a random $\hat{\boldsymbol{\theta}}$ from the posterior distribution of the parameters by using rejection sampling (see for instance [6], p. 85).

Algorithm 1 – Simulation of the state with prior information on $\boldsymbol{\theta}$

Require: $\hat{\boldsymbol{\theta}}_{\text{MLE}}, p(\boldsymbol{\theta}), m$

- 1: **for** $i = 1$ to m **do**
 - 2: **repeat**
 - 3: Draw $\hat{\boldsymbol{\theta}}$ from the prior distribution $p(\boldsymbol{\theta})$.
 - 4: Draw U from the uniform distribution on $[0, 1]$
 - 5: **until** $\ell(\hat{\boldsymbol{\theta}})/\ell(\hat{\boldsymbol{\theta}}_{\text{MLE}}) > U$
 - 6: $\hat{\boldsymbol{\theta}}^{(i)} \leftarrow \hat{\boldsymbol{\theta}}$
 - 7: Draw the i -th random trajectory $\tilde{\boldsymbol{\alpha}}_t^{(i)}$ from $p(\boldsymbol{\alpha}_t | \mathcal{Y}_n, \hat{\boldsymbol{\theta}}^{(i)})$.
 - 8: **end for**
-

The likelihood for each $\hat{\boldsymbol{\theta}}$ needed at step 5 can be computed by running the Kalman filter, generating the set of innovations \mathbf{v}_t and their covariance matrices \mathbf{F}_t from (6) and (3) above and setting

$$\log \ell(\hat{\boldsymbol{\theta}}) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_t (\log |\mathbf{F}_t| + \mathbf{v}_t^T \mathbf{F}_t^{-1} \mathbf{v}_t);$$

the inversion of \mathbf{F}_t at each step may be avoided altogether by treating the time series \mathbf{y}_t as univariate ([5], § 6.4, [1], § 6.4), so the computation of the likelihood is reasonably fast.

Step 7 is handled with the simulation smoother (see [5], § 5.3).

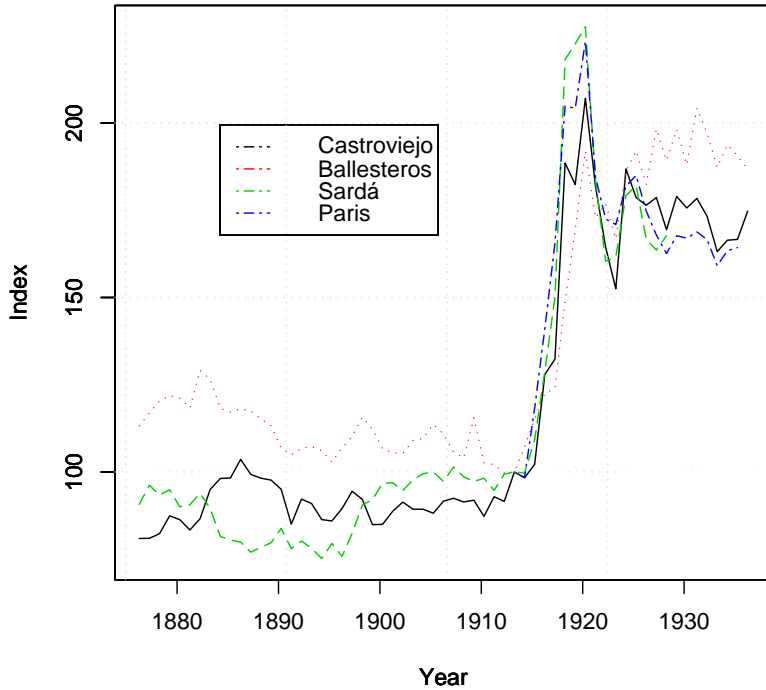
There are particular cases that can be handled faster. For certain *a priori* distributions $p(\boldsymbol{\theta})$ we may be able to exploit conjugacy and sample directly from the posterior.

On the other hand, if we have no prior information whatsoever, we may choose to sample $\hat{\boldsymbol{\theta}}^{(i)}$ from the asymptotic distribution $\mathcal{N}(\hat{\boldsymbol{\theta}}_{\text{MLE}}, \mathcal{I}(\hat{\boldsymbol{\theta}}_{\text{MLE}})^{-1})$, where $\mathcal{I}(\hat{\boldsymbol{\theta}}_{\text{MLE}})$ is the information matrix. (A similar approach in another context has been proposed by [7].) Thus, steps 2 to 7 in Algorithm 1 are replaced by a single draw from the asymptotic distribution of the estimates.

4 An illustration

Figure 1 shows four cost of living indexes, computed by four historians. They refer to the period 1876–1936 and different regions or the whole of Spain. (For a description of the indexes, see [13] and references therein.) As could

Figure 1: Cost of living indexes computed for the whole or part of Spain by four historians



be expected, the four indexes show a similar pattern: no long term trend before World War I, a phenomenal inflation during the war, then a drop of prices which nonetheless failed to return to the pre-war levels.

Not all four indexes are available for the whole period, the one computed by Sardá being the shortest.

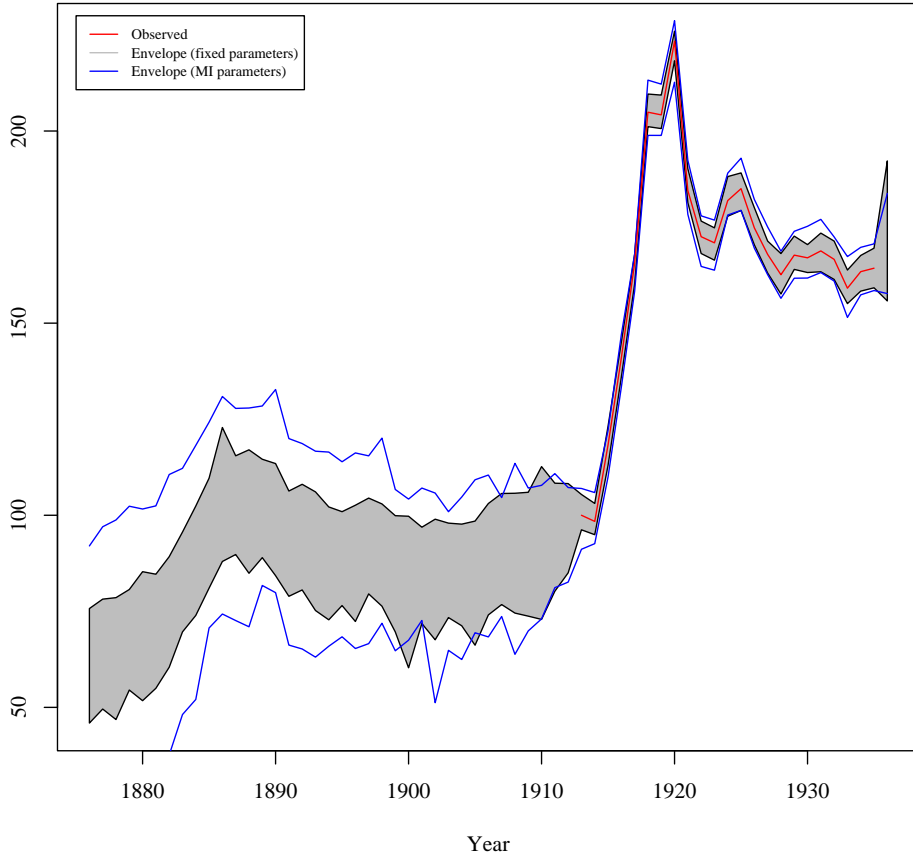
Given the similar patterns displayed by the indexes over the period where all four are observed, we can attempt to impute the missing years of an index using the past and future observations of itself and the others. To do so, we have set up the full dimension local level model described in Section 3 above. Thus, we consider:

$$\alpha_{t+1} = \alpha_t + \eta_t \quad (14)$$

$$\mathbf{y}_t = \mathbf{Z}_t \alpha_t + \epsilon_t \quad (15)$$

where $\{\mathbf{y}_t\}$ is the four dimensional time series and \mathbf{Z}_t is a matrix whose rows are a subset of the rows of the unit \mathbf{I}_4 matrix: those rows are taken that

Figure 2: Observed index and confidences bands for imputation. Grey band is the envelope of 100 trajectories from the conditional distribution $p(\boldsymbol{\alpha}_t | \mathcal{Y}_n, \hat{\boldsymbol{\theta}}_{\text{MLE}})$. The outer band is the envelope of 100 trajectories from $p(\boldsymbol{\alpha}_t | \mathcal{Y}_n, \hat{\boldsymbol{\theta}}^{(i)})$ with $\hat{\boldsymbol{\theta}}^{(i)}$ varying.



correspond to observed components of \mathbf{y}_t .

The covariance matrices are assumed time invariant; $\mathbf{H}_t = \mathbf{H}$ is chosen diagonal while $\mathbf{Q}_t = \mathbf{Q}$ is a full general covariance matrix, with no other restriction than being symmetric non-negative definite. Thus, we are assuming that what is observed are the “true” underlying indexes plus observation error, and the observation errors are unrelated for each of the four historians. On the other hand, the disturbances driving the state vector are correlated, as seems natural in this case.

With the model thus specified, two sets of one hundred trajectories of the

state α_t have been generated with the simulation smoother. In one case, the parameters were kept fixed at the values estimated by maximum likelihood, while in the other each trajectory was generated with a vector of parameters $\hat{\theta}^{(i)}$ sample from $N(\hat{\theta}_{\text{MLE}}, \mathcal{I}(\hat{\theta}_{\text{MLE}})^{-1})$. The envelopes for each set of trajectories are represented in Figure 2 and can be interpreted as approximate (simultaneous) confidence bands for the state. It can be noticed that taking into account the variability of the parameters increases substantially the width of the band, which nearly doubles in the regions where no observations were available and the Sardá index had to be imputed with information from the other three.

All computations were programmed in R (see [14]). Software is available from the author.

5 Conclusion

An approach has been proposed for the imputation of multivariate time series, and its use illustrated imputing a price index with unavailable information. The approach is general enough and able to cope, in particular, with the general situation in which several, partially overlapping sources of information are available and we need to construct an index.

We also may notice that the Kalman filter and smoother can deal with series with disparate observation periods, i.e., some series could be observed monthly and other quarterly. The use of the Kalman filter in such situations is demonstrated in [10], where the emphasis is in prediction or benchmarking while in our case is imputation.

Finally, we would like to point out that the purpose of multiple imputation in the example shown is to account for the uncertainty in the estimation of parameters, and hence be able to produce “honest” confidence intervals or bands for the estimands of interest. This issue is all too often neglected in time series analysis: [2] is one of the rare monographs to discuss this issue in his Chapter 8. Perhaps in many applications the uncertainty due to imprecise estimation of the parameters is likely to be negligible, given large enough sample sizes. That this is not always so is well epitomized by the example above.

Our approach still leaves unaccounted the uncertainty due to the choice of the model: we are assuming the model known, which will rarely be the case. A further step would be to consider a set of models and perform Bayesian model averaging.

Acknowledgments. Support from UPV/EHU (Grupo Consolidado 9/UPV 00038.321-13631/2001) and MCyT (grant BEC2003-02273) is gratefully acknowledged.

References

- [1] B.D.O. Anderson and J.B. Moore. *Optimal Filtering*. Prentice-Hall, 1979.
- [2] C. Chatfield. *Time-Series Forecasting*. Chapman and Hall/CRC, 2001.
- [3] P. de Jong. The simulation smoother for time series models. *Biometrika*, 82(2):339–350, 1995.
- [4] J. Durbin and S.J. Koopman. A simple and efficient simulation smoother for state space time series analysis. *Biometrika*, 89(3):603–615, 2002.
- [5] J. Durbin and S.J. Koopman. *Time Series Analysis by State Space Methods*. Oxford Univ. Press, 2004.
- [6] D. Gamerman. *Markov chain Monte Carlo*. Chapman and Hall, 1997.
- [7] J.D. Hamilton. A standard error for the estimated state vector of a state-space model. *Journal of Econometrics*, 33:387–397, 1986.
- [8] A. Harvey, S.J. Koopman, and N. Sheppard, editors. *State Space and Unobserved Components Models*. Cambridge Univ. Press, 2004.
- [9] A.C. Harvey. *Forecasting, structural time series models and the Kalman filter*. Cambridge Univ. Press, 1989.
- [10] A.C. Harvey and R.G. Pierse. Estimating missing observations in economic time series. *Journal of the American Statistical Association*, 79(385):125–131, 1984.
- [11] R.J.A. Little and D.B. Rubin. *Statistical Analysis with Missing Data*. John Wiley and Sons, New York, 1987. Second edition appeared 2002.
- [12] R.J.A. Little and D.B. Rubin. *Statistical Analysis with Missing Data*. Wiley, second edition, 2002.
- [13] P.M. Pérez Castroviejo. Poder adquisitivo y calidad de vida de los trabajadores vizcainos, 1876–1936. *Revista de Historia Industrial*, submitted for publication.

- [14] R Development Core Team. *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria, 2004. ISBN 3-900051-07-0.
- [15] D.B. Rubin. Multiple imputation after 18+ years (with discussion). *Journal of the American Statistical Association*, 91:473–489, 1996.
- [16] J.L. Schafer. *Analysis of Incomplete Multivariate Data*. Chapman and Hall, London, 1997.
- [17] N. Watanabe. Note on the Kalman filter with estimated parameters. *Journal of Time Series Analysis*, 6(4):269–278, 1985.