

The Relationship between State Space Subspace Identification Methods and the EM Method

Stuart Gibson

Brett Ninness

Department of Electrical and Computer Engineering,
University of Newcastle, Australia.

Web: <http://www.ee.newcastle.edu.au>

Email: {shgibson, brett}@ee.newcastle.edu.au

Abstract

This paper exposes a close connection between subspace-type estimates and single iterates of the Expectation Maximisation (EM) method for Maximum Likelihood (ML) estimation. A key implication of this is that it suggests a means by which ML estimates for MIMO systems may be computed in a numerically robust and straightforward manner.

1 Introduction

Suppose that one is presented with observations $\{y_t\}$ of a stationary stochastic process and is faced with the task of estimating a state-space representation of this process in *innovations* form

$$x_{t+1} = Ax_t + Be_t, \quad (1)$$

$$y_t = Cx_t + De_t \quad (2)$$

where $\{e_t\}$ is an i.i.d. zero mean and unit variance white noise process.

One popular approach to this problem is to apply a subspace-based method [5] to estimate A , B , C and D .

A competing solution is the Maximum Likelihood (ML) method; the statistical superiority of which has been thoroughly examined [3].

As shown here, if the ML estimate is found via the EM scheme [2], then both this and the subspace-based approach employ identical projections on an approximate state-space in order to form their parameter estimates.

Thus, the only difference between a subspace-based estimate and one EM iteration seeking an ML estimate is the approximation used for the state space.

2 An Overview of Subspace Identification Methods

This overview is drawn from that in [1]. Define $Y_{t,p}^- \in \mathbf{R}^{k \times N}$ and $Y_{t,f}^+ \in \mathbf{R}^{k \times N}$ as $[Y_{t,p}^-]_{m,n} \triangleq y_{t-m+n-1}$ and $[Y_{t,f}^+]_{m,n} \triangleq y_{t+m+n}$, where $y_s = 0$ when $s \leq 0$ or $s > N$.

Then a so-called subspace method for estimating A , B , C and D is as follows:

1. Derive a matrix $\hat{\beta} = R_{fp}(R_{pp})^{-1}$ where

$$R_{fp} \triangleq \frac{1}{N} Y_{t,f}^+ (Y_{t,p}^-)^T \text{ and } R_{pp} \triangleq \frac{1}{N} Y_{t,p}^- (Y_{t,p}^-)^T.$$

2. Choose weighting matrices W_f , W_p and perform the SVD factorisation

$$W_f \hat{\beta} W_p = USV^T = [U_n \mid u] \begin{bmatrix} S_n & 0 \\ 0 & s \end{bmatrix} \begin{bmatrix} V_n^T \\ v^T \end{bmatrix}$$

3. Define further matrices

$$\hat{O}_f = W_f^{-1} U_n S_n^{\frac{1}{2}} \quad \hat{K}_p = S_n^{\frac{1}{2}} V_n^T W_p^{-1}.$$

4. Define the state estimate as $\hat{X}_{t,p}^- = \hat{K}_p Y_{t,p}^-$.

5. Estimate C as

$$\hat{C} = \left(\frac{1}{N} \sum y_t \hat{x}_t^T \right) \left(\frac{1}{N} \sum \hat{x}_t \hat{x}_t^T \right)^{-1} \quad (3)$$

6. Define a residual estimate as

$$\widehat{De}_t = [y_t^T, y_{t+1}^T, \dots, y_N^T] - \hat{C} \hat{X}_{t,p}^-. \quad (4)$$

7. Estimate A and BD^{-1} as the least-squares solution of

$$\hat{X}_{t+1,p}^- = [A, BD^{-1}] \begin{bmatrix} \hat{X}_{t,p}^- \\ \widehat{De}_t \end{bmatrix}. \quad (5)$$

8. Estimate D as the lower triangular Cholesky factor of the sample covariance of \widehat{De}_t .

3 The EM Algorithm for Likelihood Maximisation

An essential feature of the EM algorithm is the postulate of an unobserved ‘complete data set’ $Z_N \triangleq (X_N, Y_N)$ that contains what is actually observed $Y_N \triangleq \{y_1, \dots, y_N\}$, plus other observations X_N which one might wish were available, but in fact are not. Then approximating $\log p_\theta(Z_N \mid Y_N)$ as (subscripts denote conditional dependence)

$$\log p_\theta(Z_N \mid Y_N) \approx Q(\theta, \theta') \triangleq \mathbf{E}_{\theta'} \{ \log p_\theta(Z_N \mid Y_N) \mid Y_N \}$$

leads to the EM algorithm:

1. **E Step** : Calculate $Q(\theta, \hat{\theta}_n)$.
2. **M Step** : Maximise $\hat{\theta}_{n+1} = \arg \max_\theta Q(\theta, \hat{\theta}_n)$.

4 Application of the EM Algorithm for State-Space Estimation

The most obvious choice for the incomplete data set X_N is the state sequence, i.e. $X_N = \{x_0, x_1, \dots, x_N\}$ (see [4]). Via repeated applications of Bayes' Rule

$$p_\theta(Z_N) = p_\theta(x_0) \prod_{k=0}^N p_\theta(x_k | x_{k-1}) \prod_{k=0}^N p_\theta(y_k | x_k).$$

Furthermore, suppose that the data is generated by (2) and $x_{t+1} = Ax_t + Be_t + w_t$, where $e_t \sim \mathcal{N}(0, I_m)$ is independent of $w_t \sim \mathcal{N}(0, \epsilon I)$ and $\epsilon > 0$ is arbitrarily small. This system will be 'close' to (1), (2) yet avoid the mathematical complication of singular distributions.

For this new system the M-step with $\theta = (A, B, C, D)$ becomes the closed form expressions (all expectations in (6)-(9) are to be understood as conditional on the observed data Y_N)

$$\hat{A} = \left[\sum_{k=1}^N \mathbf{E}_{\theta'} \{x_k x_{k-1}^T\} \right] \left[\sum_{k=1}^N \mathbf{E}_{\theta'} \{x_{k-1} x_{k-1}^T\} \right]^{-1}, \quad (6)$$

$$\hat{B}\hat{B}^T = \frac{1}{N} \left[\sum_{k=1}^N \mathbf{E}_{\theta'} \{x_k x_k^T\} - \left[\sum_{k=1}^N \mathbf{E}_{\theta'} \{x_k x_{k-1}^T\} \right] \times \left[\sum_{k=1}^N \mathbf{E}_{\theta'} \{x_{k-1} x_{k-1}^T\} \right]^{-1} \left[\sum_{k=1}^N \mathbf{E}_{\theta'} \{x_k x_k^T\} \right]^T \right] - \epsilon I, \quad (7)$$

$$\hat{C} = \bar{Y}_N \hat{X}_N^T \left[\sum_{k=0}^N \mathbf{E}_{\theta'} \{x_k x_k^T\} \right]^{-1}, \quad (8)$$

$$\hat{D}\hat{D}^T = \frac{1}{N} \bar{Y}_N \left[I - \hat{X}_N^T \left[\sum_{k=0}^N \mathbf{E}_{\theta'} \{x_k x_k^T\} \right]^{-1} \hat{X}_N \right] \bar{Y}_N^T, \quad (9)$$

where \hat{x}_k and $\mathbf{E}_{\theta'} \{x_k x_k^T | Y_N\}$ may be extracted from a Kalman smoother and $\hat{X}_N \triangleq [\hat{x}_0, \dots, \hat{x}_N]$, $\bar{Y}_N \triangleq [y_1, \dots, y_N]$.

5 Relationship between EM and Subspace Methods

Both the EM and the subspace methods estimate the state sequence. In the case of subspace methods the estimate $\hat{X}_{t,p}^-$ can be thought of as an approximate Kalman filtered estimate [5]. For the EM algorithm a Kalman smoother is employed to estimate the state sequence. Both the subspace and EM method then calculate parameter estimates using a series of projections which involve the state sequence estimate.

In particular, given an estimate of the state sequence the subspace method calculates A , $\bar{B} \triangleq BD^{-1}$ as in (5), that is,

$$\hat{X}_{t+1,p}^- = [A, \bar{B}] \begin{bmatrix} \hat{X}_{t,p}^- \\ \bar{Y}_N \hat{X}_{t,p}^\perp \end{bmatrix}, \quad (10)$$

where $\hat{X}_{t,p}^\perp = I - (\hat{X}_{t,p}^-)^T \left[\hat{X}_{t,p}^- (\hat{X}_{t,p}^-)^T \right]^{-1} \hat{X}_{t,p}^-$. Solving this leads to

$$\hat{A} = \hat{X}_{t+1,p}^- (\hat{X}_{t,p}^-)^T \left[\hat{X}_{t,p}^- (\hat{X}_{t,p}^-)^T \right]^{-1}, \quad (11)$$

$$\hat{B} = \hat{X}_{t+1,p}^- \hat{X}_{t,p}^\perp \bar{Y}_N^T \left[\bar{Y}_N \hat{X}_{t,p}^\perp \bar{Y}_N^T \right]^{-1}. \quad (12)$$

Comparing (11) and (6) indicates that \hat{A} is calculated by both the subspace and EM methods using an identical projection. Similarly, equations (3) and (8) immediately indicate that both methods employ the same projection to find \hat{C} .

The subspace method's equation for \widehat{De}_t , (4), yields

$$\hat{D}\hat{D}^T = \frac{1}{N} \bar{Y}_N \left[I - (\hat{X}_{t,p}^-)^T \left[\hat{X}_{t,p}^- (\hat{X}_{t,p}^-)^T \right]^{-1} \hat{X}_{t,p}^- \right] \bar{Y}_N^T,$$

which is identical to the EM algorithm projection (9).

Lastly, the subspace method estimates $\hat{B} = \hat{B}\hat{D}^{-1}$ as per (12), implying that

$$\hat{B}\hat{B}^T = \frac{1}{N} \hat{X}_{t+1,p}^- \hat{X}_{t,p}^\perp \left(\hat{X}_{t+1,p}^- \right)^T,$$

which, again, is identical to the EM scheme (7) except for the particular method of state estimation.

6 Conclusion

Subspace methods have attracted interest because of their numerical robustness and (because of their lack of imposed state space parameterisation) their suitability for MIMO systems. The fact shown here, that a subspace method is essentially an iteration towards an ML estimate (once Kalman smoother state estimates are used), means that ML estimates can be obtained with the same numerical robustness and with the same parameterisation free advantages. Preliminary results towards this goal are provided in an extended report available via <http://www.ee.newcastle.edu.au>.

References

- [1] M. DEISTLER, K. PETERNELL, AND W. SCHERRER, *Consistency and relative efficiency of subspace methods*, *Automatica*, 31 (1995).
- [2] A. DEMPSTER, N. LAIRD, AND D. RUBIN, *Maximum likelihood from incomplete data via the em algorithm*, *J. R. Stat. Soc., B*, 39 (1977), pp. 1-38.
- [3] L. LJUNG, *System Identification: Theory for the User*, Prentice-Hall, Inc., New Jersey, 1987.
- [4] R. SHUMWAY AND D. STOFFER, *An approach to time series smoothing and forecasting using the em algorithm*, *J. Time Ser. Anal.*, 3 (1982).
- [5] P. VAN OVERSCHEE AND B. DE MOOR, *Subspace Identification for Linear Systems*, Kluwer Academic Publishers, 1996.