

Comparison of Various Approaches for Joint Wiener/Kalman Filtering and Parameter Estimation with Application to BASS

Siouar Bensaïd and Dirk Slock
 Mobile Communications Department
 EURECOM, Sophia Antipolis, France
 Email: {bensaid, slock}@eurecom.fr

Abstract—In recent years, the Kalman filter (KF) has encountered renewed interest, due to an increasing range of applications. Even though in many cases the state-space model may be linear, it is often only known up to the values of some parameters, usually related to the vector autoregressive process of the state evolution equation. In this paper, after finding motivation in some applications, we review a number of approaches for adaptive Kalman filtering (AKF), in which state and parameters get estimated jointly. We propose an improved version of the Extended KF (EKF) in which the estimation error covariance matrix is computed exactly assuming overall joint Gaussianity. We also compare the performance and Cramer Rao bounds (CRBs) of joint Maximum A Posteriori Maximum Likelihood (MAP-ML) estimation of Bayesian state and deterministic parameters, and marginalized ML estimation of the parameters, and relate this to the Expectation-Maximization KF (EM-KF). The perspectives involve also the Variational Bayesian KF (VB-KF).

I. INTRODUCTION

Since Rudolf E. Kalman published his famous paper in 1960, the Kalman filter has become the work horse of many estimation processes in different application areas. The Kalman filter (KF) considers the estimation of a first-order vector autoregressive (AR(1)) (Markov) process from linear measurements in white noise. The KF performs this estimation recursively by alternating between filtering (measurement update) and (one step ahead) prediction (time update). An alternative viewpoint is that the Kalman filter recursively generates the innovations of the measurement signal (by a structured Gram-Schmidt approach that decorrelates the consecutive measurements). The KF corresponds to optimal (Minimum Mean Squared Error (MMSE) or Maximum A Posteriori (MAP)) Bayesian estimation of the state sequence if all random sources involved (measurement noise, state noise and state initial conditions) are Gaussian. In the non-Gaussian case, the KF performs Linear MMSE (LMMSE) estimation. The linear state-space model can be written as

$$\left\{ \begin{array}{l} \text{state update equation:} \\ \mathbf{x}_{k+1} = \mathbf{F}_k(\theta) \mathbf{x}_k + \mathbf{G}_k(\theta) \mathbf{w}_k \\ \text{measurement equation:} \\ \mathbf{y}_k = \mathbf{H}_k(\theta) \mathbf{x}_k + \mathbf{v}_k \end{array} \right. \quad (1)$$

for discrete time $k = 1, 2, \dots$, where the initial state $\mathbf{x}_0 \sim \mathcal{N}(\hat{\mathbf{x}}_0, \mathbf{P}_0(\theta))$, the measurement noise $\mathbf{v}_k \sim \mathcal{N}(0, \mathbf{R}_k(\theta))$,

EURECOMs research is partially supported by its industrial members: ORANGE, BMW Group, Swisscom, Cisco, SFR, ST Ericsson, Thales, Symantec, SAP, Monaco Telecom, and also by the EU FP7 project WHERE2.

and the state noise $\mathbf{w}_k \sim \mathcal{N}(0, \mathbf{Q}_k(\theta))$ and all these random quantities are mutually uncorrelated. As indicated, the state-space model is often specified up to the value of some parameters θ . Often the $\mathbf{F}_k(\theta)$, $\mathbf{G}_k(\theta)$, $\mathbf{H}_k(\theta)$ are linear in θ , which would correspond to the bilinear case. Although this signal model seems very simple, the applications are numerous. We will cite here three examples:

- Bayesian adaptive filtering [18] (or wireless channel estimation [9], [12]):
 in this case, \mathbf{x}_k = FIR filter response, and θ contains e.g. the Power Delay Profile (diagonal of a diagonal filter coefficient covariance matrix $\mathbf{P}_0 = \mathbf{P}_k$, and the AR(1) dynamics in e.g. diagonal \mathbf{F} and \mathbf{Q}).
- Position tracking (GPS) (see [5] and references therein):

$$\mathbf{x}_{t+1} = \begin{bmatrix} 1 & \Delta t & \frac{1}{2}\Delta t^2 \\ 0 & 1 & \Delta t \\ 0 & 0 & 1 \end{bmatrix} \cdot \mathbf{x}_t = \begin{bmatrix} x_t + \Delta t \cdot v_t + \frac{1}{2}\Delta t^2 a \\ v_t + \Delta t \cdot a \\ a \end{bmatrix}$$

the state contains position, velocity and possible acceleration and θ contains acceleration model parameters (e.g. white noise, AR(1))

- Blind Audio Source Separation (BASS) [4]: x_k = source signals, θ : (short+long term) AR parameters, reverb filters

In the literature, variations on the KF theme have been derived to handle the joint filtering and parameter estimation problem, such as e.g. the widely used EM-KF algorithm ([6], [8], [9]) which uses the famous Expectation Maximization technique (EM), and alternating optimization technique for ML estimation. Another well-known variation is the EKF algorithm, which can handle general nonlinear state space models. In this case, the state is extended with the unknown parameters, rendering the new state update equation nonlinear. A third derivation is the truncated Second-Order EKF (SOEKF) introduced by [3], [11] in which nonlinearities are expanded up to second order, third and higher order statistics being neglected. A corrected derivation of this filter is presented in [10]. In ([2], [11]), the Gaussian SOEKF is derived in which fourth-order terms in the Taylor series expansions are retained and approximated by assuming that the underlying joint probability distribution is Gaussian. In [21], Villares et al. introduced the Quadratic Extended Kalman Filter (QEKF) where they extend the EKF to a new algorithm using quadratic processing and incorporating fourth order statistics of the input signal. The problem of uncertainty about the process noise

and measurement noise covariance matrices was also tackled in [16] where a test of Kalman filter optimality is used in order to estimate the unknown noise covariance matrices. The performance of some of these Adaptive KF (AKF) approaches was studied in the literature. In [7], the EM approach is proved to converge to the ML performance. The asymptotic behavior of the EKF for AKF has been treated in [13] where it is proved that no global convergence is guaranteed. The performance analysis of linear and nonlinear KF has also been treated in terms of Cramer Rao Bound (CRB) computations. In [19], the Posterior CRB (PCRB) is developed for the discrete nonlinear KF. Recursive Bayesian CRBs were also developed for continuous and discrete nonlinear filtering for many problems. We can refer to [20] for an overview. This paper is organized as follows: a review of some AKF approaches is proposed in section II. An improved version of the EKF algorithm is developed in section III and some performance orderings are provided in section IV.

II. ADAPTIVE KALMAN FILTERING APPROACHES

A. Basic Kalman Filter (KF)

In the following, we introduce the notation $\mathbf{y}_{1:k} = \{\mathbf{y}_1, \dots, \mathbf{y}_k\}$. The KF performs Gram-Schmidt orthogonalization (decorrelation) of the measurement variables \mathbf{y}_k . This is done by computing the LMMSE predictor $\hat{\mathbf{y}}_{k|k-1}$ of \mathbf{y}_k on the basis of $\mathbf{y}_{1:k-1}$, leading to the orthogonalized prediction error (or innovation) $\tilde{\mathbf{y}}_k = \tilde{\mathbf{y}}_{k|k-1} = \mathbf{y}_k - \hat{\mathbf{y}}_{k|k-1}$. We introduce the correlation matrix notation $R_{\mathbf{x}\mathbf{y}} = \mathbb{E} \mathbf{x} \mathbf{y}^T$ (correlation matrices will usually also be covariance matrices here since the processes \mathbf{y}_k and \mathbf{x}_k have zero mean and also various estimation errors will have (conditional) zero mean). We denote the covariance matrix $R_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k} = \mathbf{S}_k$. The idea of the innovations approach is that (linear) estimation in terms of $\mathbf{y}_{1:k}$ is equivalent to estimation in terms of $\tilde{\mathbf{y}}_{1:k}$ since one set is obtained from the other by an invertible linear transformation. Now, since the $\tilde{\mathbf{y}}_k$ are decorrelated, estimation in terms of $\tilde{\mathbf{y}}_{1:k}$ simplifies:

$$\hat{\mathbf{x}}_{|k} = \sum_{i=1}^k R_{\mathbf{x}\tilde{\mathbf{y}}_i} R_{\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i}^{-1} \tilde{\mathbf{y}}_i = \mathbf{x}_{|k-1} + R_{\mathbf{x}\tilde{\mathbf{y}}_k} \mathbf{S}_k^{-1} \tilde{\mathbf{y}}_k.$$

This will be used to obtain *predicted* estimates $\hat{\mathbf{x}}_{k|k-1}$ with estimation error $\tilde{\mathbf{x}}_{k|k-1} = \mathbf{x}_k - \hat{\mathbf{x}}_{k|k-1}$ with covariance matrix $\mathbf{P}_{k|k-1} = R_{\tilde{\mathbf{x}}_{k|k-1} \tilde{\mathbf{x}}_{k|k-1}}$ and also *filtered* estimates $\hat{\mathbf{x}}_{k|k}$ with estimation error $\tilde{\mathbf{x}}_{k|k} = \mathbf{x}_k - \hat{\mathbf{x}}_{k|k}$ with covariance matrix $\mathbf{P}_{k|k} = R_{\tilde{\mathbf{x}}_{k|k} \tilde{\mathbf{x}}_{k|k}}$.

Now exploiting the correlation structure in the signal model, this leads to the following two-step recursive procedure to go from $|k-1$ to $|k$:

Measurement Update

$$\begin{aligned} \hat{\mathbf{y}}_{k|k-1} &= \mathbf{H}_k \hat{\mathbf{x}}_{k|k-1} \\ \tilde{\mathbf{y}}_k &= \mathbf{y}_k - \hat{\mathbf{y}}_{k|k-1} \\ \mathbf{S}_k &= \mathbf{H}_k \mathbf{P}_{k|k-1} \mathbf{H}_k^T + \mathbf{R}_k \\ \mathbf{K}_k &= \mathbf{P}_{k|k-1} \mathbf{H}_k^T \mathbf{S}_k^{-1} \\ \hat{\mathbf{x}}_{k|k} &= \hat{\mathbf{x}}_{k|k-1} + \mathbf{K}_k \tilde{\mathbf{y}}_k \\ \mathbf{P}_{k|k} &= \mathbf{P}_{k|k-1} - \mathbf{K}_k \mathbf{H}_k \mathbf{P}_{k|k-1} \end{aligned} \quad (2)$$

Time Update (prediction)

$$\begin{aligned} \hat{\mathbf{x}}_{k+1|k} &= \mathbf{F}_k \hat{\mathbf{x}}_{k|k} \\ \mathbf{P}_{k+1|k} &= \mathbf{F}_k \mathbf{P}_{k|k} \mathbf{F}_k^T + \mathbf{G}_k \mathbf{Q}_k \mathbf{G}_k^T \end{aligned} \quad (3)$$

In the usual case of total absence of prior information on the initial state, one can choose $\hat{\mathbf{x}}_0 = 0$, $\mathbf{P}_0 = p_0 \mathbf{I}$ with p_0 a (very) large number.

B. Extended Kalman Filter (EKF)

For the case of a nonlinear state-space model, the idea of the EKF is to apply the KF to a linearized version of the state-space model, via a first-order Taylor series expansion. So we get

state update equation:

$$\mathbf{x}_{k+1} = f(\mathbf{x}_k, \mathbf{w}_k) \approx \mathbf{F}_k \mathbf{x}_k + \mathbf{G}_k \mathbf{w}_k \quad (4)$$

measurement equation:

$$\mathbf{y}_k = h(\mathbf{x}_k) + \mathbf{v}_k \approx \mathbf{H}_k \mathbf{x}_k + \mathbf{v}_k$$

where

$$\begin{aligned} \mathbf{F}_k &= \left. \frac{\partial f(\mathbf{x}, \mathbf{w})}{\partial \mathbf{x}^T} \right|_{(\mathbf{x}, \mathbf{w})=(\mathbf{x}_k, \mathbf{w}_k)} & \mathbf{G}_k &= \left. \frac{\partial f(\mathbf{x}, \mathbf{w})}{\partial \mathbf{w}^T} \right|_{(\mathbf{x}, \mathbf{w})=(\mathbf{x}_k, \mathbf{w}_k)} \\ \mathbf{H}_k &= \left. \frac{\partial h(\mathbf{x})}{\partial \mathbf{x}^T} \right|_{\mathbf{x}=\mathbf{x}_k} \end{aligned} \quad (5)$$

So, at this point, the basic KF can be applied to the thus obtained approximate linear state-space model. The EKF approach can be used to adapt some parameters in an otherwise linear state-space model $\mathbf{x}_{k+1} = \mathbf{F}' \mathbf{x}'_k + \mathbf{G}' \mathbf{w}_k$. For instance, consider the case in which one wants to adapt parameters appearing (e.g.) linearly in the matrix $\mathbf{F}' = \mathbf{F}'(\boldsymbol{\theta})$. One can jointly estimate the unknown constant parameter vector $\boldsymbol{\theta}$ by considering the following state update for them: $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k$. Then one can introduce the augmented state and system matrices

$$\mathbf{x}_k = \begin{bmatrix} \mathbf{x}'_k \\ \boldsymbol{\theta}_k \end{bmatrix}, \quad \mathbf{F}_k = \begin{bmatrix} \mathbf{F}'(\boldsymbol{\theta}_k) & \mathbf{C}(\mathbf{x}'_k) \\ 0 & \mathbf{I} \end{bmatrix}, \quad \mathbf{G}_k = \begin{bmatrix} \mathbf{G}' \\ 0 \end{bmatrix} \quad (6)$$

where $\mathbf{C}(\mathbf{x}'_k) = \frac{\partial \mathbf{F}'(\boldsymbol{\theta}) \mathbf{x}'_k}{\partial \boldsymbol{\theta}^T}$. When running the EKF, the state-dependent system matrices have to be filled with the latest state estimates, so in this case

$$\mathbf{F}_k = \begin{bmatrix} \mathbf{F}'(\hat{\boldsymbol{\theta}}_{k|k}) & \mathbf{C}(\hat{\mathbf{x}}'_{k|k}) \\ 0 & \mathbf{I} \end{bmatrix}. \quad (7)$$

The parameters $\boldsymbol{\theta}$ are often not really constant and hence need to be tracked adaptively. This can be done either by introducing some process noise in $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k$ (random walk time evolution) or by introducing exponential weighting (at least for the $\boldsymbol{\theta}$ portion) into the KF updates [1]. The EKF approach allows fairly straightforwardly to estimate parameters in \mathbf{F}_k , \mathbf{H}_k , or \mathbf{G}_k , but much less so in \mathbf{Q}_k , \mathbf{R}_k .

For adapting (parameters in) \mathbf{Q} and \mathbf{R} , one needs to consider the innovations representation $\hat{\mathbf{x}}_{k+1|k} = \mathbf{F}_k \hat{\mathbf{x}}_{k|k-1} + \mathbf{F}_k \mathbf{K}_k \tilde{\mathbf{y}}_k$ and consider gradients of the Kalman gain \mathbf{K}_k w.r.t. these matrices.

C. Recursive Prediction Error Method (RPEM-KF)

The RPEM [15], [14] is an adaptive implementation of Maximum Likelihood (ML) parameter estimation. The negative loglikelihood becomes a least-squares criterion in the prediction errors (innovations) and RPEM performs one iteration per new sample. Applied to KF, the RPEM can be seen as a more rigorous version of EKF and computes gradients more precisely [23]. Indeed, for the case of a state transition matrix $\mathbf{F}_k = \mathbf{F}_k(\boldsymbol{\theta})$, the EKF would consider the gradient

$$\frac{\partial \mathbf{x}_{k+1}}{\partial \boldsymbol{\theta}^T} = \frac{\partial \mathbf{F}_k(\boldsymbol{\theta}) \mathbf{x}_k}{\partial \boldsymbol{\theta}^T} \quad (8)$$

where only the explicit dependence of \mathbf{F} on $\boldsymbol{\theta}$ would be considered, whereas the RPEM would consider more correctly

$$\frac{\partial \mathbf{x}_{k+1}}{\partial \boldsymbol{\theta}^T} = \frac{\partial \mathbf{F}_k(\boldsymbol{\theta}) \mathbf{x}_k}{\partial \boldsymbol{\theta}^T} + \mathbf{F}_k(\boldsymbol{\theta}) \frac{\partial \mathbf{x}_k}{\partial \boldsymbol{\theta}^T}. \quad (9)$$

RPEM for KF can be found in the references above, but will not be pursued here further. One characteristic of the RPEM is a higher complexity.

D. Expectation-Maximization (EM-KF)

In EM [7], the parameters are estimated by minimizing expected values of negative loglikelihoods, see e.g. [22] for an application involving KF. For the state update, since \mathbf{G}_k is typically a tall matrix, $\mathbf{G}_k \mathbf{w}_k$ has a singular covariance matrix. The state update equation can be rewritten as

$$\mathbf{G}_k^+ \mathbf{x}_{k+1} = \mathbf{G}_k^+ \mathbf{F}_k \mathbf{x}_k + \mathbf{w}_k \quad (10)$$

where $\mathbf{G}_k^+ = (\mathbf{G}_k^T \mathbf{G}_k)^{-1} \mathbf{G}_k^T$ is the pseudo-inverse of \mathbf{G}_k . For the parameters involved in the state update equation, hence the following negative loglikelihood is applicable:

$$\sum_k \{ \ln \det(\mathbf{Q}_k) + (\mathbf{x}_{k+1} - \mathbf{F}_k \mathbf{x}_k)^T \mathbf{G}_k^{+T} \mathbf{Q}_k^{-1} \mathbf{G}_k^+ (\mathbf{x}_{k+1} - \mathbf{F}_k \mathbf{x}_k) \} \quad (11)$$

For the parameters involved in the measurement equation, the appropriate loglikelihood is

$$\sum_k \{ \ln \det(\mathbf{R}_k) + (\mathbf{y}_k - \mathbf{H}_k \mathbf{x}_k)^T \mathbf{R}_k^{-1} (\mathbf{y}_k - \mathbf{H}_k \mathbf{x}_k) \}. \quad (12)$$

Now the expectation is taken, in principle with the conditional distribution given all data. Hence $\mathbb{E}_{|n}$ involving all data \mathbf{y}_k up to the last sample n . This leads to an iterative algorithm with in each iteration a whole fixed-interval smoothing operation. An adaptive version [22], [9] can be obtained by replacing fixed-interval smoothing by fixed-lag smoothing and performing one iteration per time sample. Since the state update equation corresponds to a vector AR(1) model, one may expect (as in [9]) that a lag of 1 should be enough (to guarantee convergence). In [22], complexity is reduced further by suggesting that filtering might be enough. In that case, the (presumably) slowly varying $\hat{\mathbf{Q}}_{k+1}$, $\hat{\mathbf{F}}_{k+1}$ (for use in the KF at time $k+1$) get determined by minimizing $\sum_{i=1}^k \lambda^{k-i} \mathbb{E}_{|i} \{ \text{Terms in (11)} \}$ w.r.t. \mathbf{Q} , \mathbf{F} (\mathbf{G} is known) where we introduced an exponential forgetting factor $\lambda \lesssim 1$. This is equivalent to

$$\gamma_k^{-1} \ln \det(\hat{\mathbf{Q}}) +$$

$$\sum_{i=1}^k \lambda^{k-i} \text{tr} \{ \mathbf{G}_i^{+T} \hat{\mathbf{Q}}^{-1} \mathbf{G}_i^+ \mathbb{E}_{|i} (\mathbf{x}_{i+1} - \hat{\mathbf{F}} \mathbf{x}_i) (\mathbf{x}_{i+1} - \hat{\mathbf{F}} \mathbf{x}_i)^T \} \quad (13)$$

where we introduced $\gamma_k^{-1} = \sum_{i=1}^k \lambda^{k-i} = \lambda \gamma_{k-1}^{-1} + 1$. γ_k^{-1} behaves initially as $1/k$ but saturates eventually at $\gamma_\infty^{-1} = 1 - \lambda$. We shall need

$$\begin{aligned} \mathbb{E}_{|i} \mathbf{x}_i \mathbf{x}_i^T &= \hat{\mathbf{x}}_{i|i} \hat{\mathbf{x}}_{i|i}^T + \mathbf{P}_{i|i} \\ \mathbb{E}_{|i} \mathbf{x}_{i+1} \mathbf{x}_i^T &= \mathbf{F}_i \hat{\mathbf{x}}_{i|i} \hat{\mathbf{x}}_{i|i}^T + \mathbf{F}_i \mathbf{P}_{i|i} \\ \mathbb{E}_{|i} \mathbf{x}_i \mathbf{x}_{i+1}^T &= \hat{\mathbf{x}}_{i|i} \hat{\mathbf{x}}_{i|i}^T \mathbf{F}_i^T + \mathbf{P}_{i|i} \mathbf{F}_i^T \\ \mathbb{E}_{|i} \mathbf{x}_{i+1} \mathbf{x}_{i+1}^T &= \mathbf{F}_i \hat{\mathbf{x}}_{i|i} \hat{\mathbf{x}}_{i|i}^T \mathbf{F}_i^T + \mathbf{P}_{i+1|i} \\ &= \mathbf{F}_i (\hat{\mathbf{x}}_{i|i} \hat{\mathbf{x}}_{i|i}^T + \mathbf{P}_{i|i}) \mathbf{F}_i^T + \mathbf{G}_i \mathbf{Q}_i \mathbf{G}_i^T \end{aligned} \quad (14)$$

In case of time-invariant $\mathbf{G}_k \equiv \mathbf{G}$, we can rewrite (13) as

$$\ln \det(\hat{\mathbf{Q}}) + \text{tr} \{ \mathbf{G}^{+T} \hat{\mathbf{Q}}^{-1} \mathbf{G}^+ (\mathbf{M}_k^{11} - \hat{\mathbf{F}} \mathbf{M}_k^{01} - \mathbf{M}_k^{10} \hat{\mathbf{F}}^T + \hat{\mathbf{F}} \mathbf{M}_k^{00} \hat{\mathbf{F}}^T) \} \quad \text{where} \quad (15)$$

$$\begin{aligned} \mathbf{M}_k^{00} &= (1 - \gamma_k) \mathbf{M}_{k-1}^{00} + \gamma_k (\hat{\mathbf{x}}_{k|k} \hat{\mathbf{x}}_{k|k}^T + \mathbf{P}_{k|k}) \\ \mathbf{M}_k^{10} &= (1 - \gamma_k) \mathbf{M}_{k-1}^{10} + \gamma_k \mathbf{F}_k (\hat{\mathbf{x}}_{k|k} \hat{\mathbf{x}}_{k|k}^T + \mathbf{P}_{k|k}) \\ \mathbf{M}_k^{01} &= (1 - \gamma_k) \mathbf{M}_{k-1}^{01} + \gamma_k (\hat{\mathbf{x}}_{k|k} \hat{\mathbf{x}}_{k|k}^T + \mathbf{P}_{k|k}) \mathbf{F}_k^T \\ \mathbf{M}_k^{11} &= (1 - \gamma_k) \mathbf{M}_{k-1}^{11} + \gamma_k (\hat{\mathbf{x}}_{k+1|k} \hat{\mathbf{x}}_{k+1|k}^T + \mathbf{P}_{k+1|k}) \end{aligned}$$

In case of furthermore time-invariant $\mathbf{F}_k \equiv \mathbf{F}$, $\mathbf{Q}_k \equiv \mathbf{Q}$, then

$$\begin{aligned} \mathbf{M}_k^{10} &= \mathbf{F} \mathbf{M}_k^{00} \\ \mathbf{M}_k^{01} &= \mathbf{M}_k^{00} \mathbf{F}^T \\ \mathbf{M}_k^{11} &= \mathbf{F} \mathbf{M}_k^{00} \mathbf{F}^T + \mathbf{G} \mathbf{G} \mathbf{G}^T \end{aligned} \quad (16)$$

As a result, (15) can be rewritten as $\ln \det(\hat{\mathbf{Q}}) + \text{tr} \{ \hat{\mathbf{Q}}^{-1} \mathbf{Q} \} + \text{tr} \{ \mathbf{G}^{+T} \hat{\mathbf{Q}}^{-1} \mathbf{G}^+ (\mathbf{F} - \hat{\mathbf{F}}) \mathbf{M}_k^{00} (\mathbf{F} - \hat{\mathbf{F}})^T \}$, the optimization of which now clearly leads to $\hat{\mathbf{F}} = \mathbf{F}$, $\hat{\mathbf{Q}} = \mathbf{Q}$. So we just get back the quantities that we use in the KF, without any additional information. Hence, just Kalman filtering in the EM-KF is not enough to adapt the state update parameters.

E. Fixed-Lag Smoothing

Using the innovations approach, we have

$$\hat{\mathbf{x}}_{k-1|k} = \hat{\mathbf{x}}_{k-1|k-1} + \mathbf{R}_{\mathbf{x}_{k-1} \tilde{\mathbf{y}}_k} \mathbf{S}_k^{-1} \tilde{\mathbf{y}}_k. \quad (17)$$

After a few steps, we get the following lag-1 smoothing equations that need to be added to the basic Kalman Filter equations (to be inserted between the Measurement Update and the Time Update)

$$\begin{aligned} \mathbf{K}_{k;1} &= \mathbf{P}_{k-1|k-1} \mathbf{F}_{k-1}^T \mathbf{H}_k^T \\ \hat{\mathbf{x}}_{k-1|k} &= \hat{\mathbf{x}}_{k-1|k-1} + \mathbf{K}_{k;1} \mathbf{S}_k^{-1} \tilde{\mathbf{y}}_k \\ \mathbf{P}_{k-1|k} &= \mathbf{P}_{k-1|k-1} - \mathbf{K}_{k;1} \mathbf{S}_k^{-1} \mathbf{K}_{k;1}^T \end{aligned} \quad (18)$$

F. Adaptive EM-KF with Fixed-Lag Smoothing

Consider now the case in which the state-space model is essentially time-invariant (or slowly time-varying). In that case the time index of the system matrices \mathbf{F}_k etc. just reflects at

which time the (unknown) system matrices have been adapted. The resulting KF equations with lag-1 smoothing become

$$\begin{aligned}
\widehat{\mathbf{y}}_{k|k-1} &= \mathbf{H}_{k-1} \widehat{\mathbf{x}}_{k|k-1} \\
\widehat{\mathbf{y}}_k &= \mathbf{y}_k - \widehat{\mathbf{y}}_{k|k-1} \\
\mathbf{S}_k &= \mathbf{H}_{k-1} \mathbf{P}_{k|k-1} \mathbf{H}_{k-1}^T + \mathbf{R}_{k-1} \\
\mathbf{K}_{k;1} &= \mathbf{P}_{k-1|k-1} \mathbf{F}_{k-1}^T \mathbf{H}_{k-1}^T + \mathbf{R}_{k-1} \\
\widehat{\mathbf{x}}_{k-1|k} &= \widehat{\mathbf{x}}_{k-1|k-1} + \mathbf{K}_{k;1} \mathbf{S}_k^{-1} \widehat{\mathbf{y}}_k \\
\mathbf{P}_{k-1|k} &= \mathbf{P}_{k-1|k-1} - \mathbf{K}_{k;1} \mathbf{S}_k^{-1} \mathbf{K}_{k;1}^T \\
\mathbf{K}_k &= \mathbf{P}_{k|k-1} \mathbf{H}_{k-1}^T \mathbf{S}_k^{-1} \\
\widehat{\mathbf{x}}_{k|k} &= \widehat{\mathbf{x}}_{k|k-1} + \mathbf{K}_k \widehat{\mathbf{y}}_k \\
\mathbf{P}_{k|k} &= \mathbf{P}_{k|k-1} - \mathbf{K}_k \mathbf{H}_{k-1} \mathbf{P}_{k|k-1} \\
&\text{parameter update} \\
\widehat{\mathbf{x}}_{k+1|k} &= \mathbf{F}_k \widehat{\mathbf{x}}_{k|k} \\
\mathbf{P}_{k+1|k} &= \mathbf{F}_k \mathbf{P}_{k|k} \mathbf{F}_k^T + \mathbf{G}_k \mathbf{Q}_k \mathbf{G}_k^T
\end{aligned} \tag{19}$$

So, the system matrices (\mathbf{F} , \mathbf{G} , \mathbf{Q}) should be adapted after the smoothing step and before the filtering and prediction steps. We now adapt the system matrices \mathbf{F} , \mathbf{G} , \mathbf{Q} from the equivalent of (13) with $E_{|i}$ replaced by $E_{|i+1}$. This leads to the matrix updates

$$\begin{aligned}
\mathbf{M}_k^{00} &= (1 - \gamma_k) \mathbf{M}_{k-1}^{00} + \gamma_k (\widehat{\mathbf{x}}_{k-1|k} \widehat{\mathbf{x}}_{k-1|k}^T + \mathbf{P}_{k-1|k}) \\
\mathbf{M}_k^{10} &= (\mathbf{M}_k^{01})^T = (1 - \gamma_k) \mathbf{M}_{k-1}^{10} + \gamma_k (\widehat{\mathbf{x}}_{k|k} \widehat{\mathbf{x}}_{k-1|k}^T \\
&\quad + \mathbf{F}_{k-1} \mathbf{P}_{k-1|k} - \mathbf{G}_{k-1} \mathbf{Q}_{k-1} \mathbf{G}_{k-1}^T \mathbf{H}_{k-1}^T \mathbf{S}_k^{-1} \mathbf{K}_{k;1}^T) \\
\mathbf{M}_k^{11} &= (1 - \gamma_k) \mathbf{M}_{k-1}^{11} + \gamma_k (\widehat{\mathbf{x}}_{k|k} \widehat{\mathbf{x}}_{k|k}^T + \mathbf{P}_{k|k}) .
\end{aligned}$$

If for example \mathbf{G} would be fixed and invertible, minimization of the expected loglikelihood w.r.t. $\widehat{\mathbf{F}}$, $\widehat{\mathbf{Q}}$ would lead to the following minimizers

$$\begin{aligned}
\mathbf{F}_k &= \mathbf{M}_k^{10} (\mathbf{M}_k^{00})^{-1} \\
\mathbf{Q}_k &= \mathbf{G}^+ (\mathbf{M}_k^{11} - \mathbf{M}_k^{10} (\mathbf{M}_k^{00})^{-1} \mathbf{M}_k^{01}) \mathbf{G}^{+T} .
\end{aligned} \tag{20}$$

For adapting the parameters in the measurement equation on the other hand, Kalman filtering should be sufficient. A similar derivation from the expected measurement loglikelihood leads to

$$\begin{aligned}
\widehat{\mathbf{H}}_k &= \widehat{\mathbf{R}}_{\mathbf{y}\mathbf{x},k} \widehat{\mathbf{R}}_{\mathbf{x}\mathbf{x},k}^{-1} \\
\widehat{\mathbf{R}}_k &= \widehat{\mathbf{R}}_{\mathbf{y}\mathbf{y},k} - \widehat{\mathbf{R}}_{\mathbf{y}\mathbf{x},k} \widehat{\mathbf{R}}_{\mathbf{x}\mathbf{x},k}^{-1} \widehat{\mathbf{R}}_{\mathbf{x}\mathbf{y},k} \quad \text{where}
\end{aligned} \tag{21}$$

$$\begin{aligned}
\widehat{\mathbf{R}}_{\mathbf{y}\mathbf{y},k} &= (1 - \gamma_k) \widehat{\mathbf{R}}_{\mathbf{y}\mathbf{y},k-1} + \gamma_k \mathbf{y}_k \mathbf{y}_k^T \\
\widehat{\mathbf{R}}_{\mathbf{x}\mathbf{y},k} &= (1 - \gamma_k) \widehat{\mathbf{R}}_{\mathbf{x}\mathbf{y},k-1} + \gamma_k \widehat{\mathbf{x}}_{k|k} \mathbf{y}_k^T \\
\widehat{\mathbf{R}}_{\mathbf{y}\mathbf{x},k} &= \widehat{\mathbf{R}}_{\mathbf{x}\mathbf{y},k}^T \\
\widehat{\mathbf{R}}_{\mathbf{x}\mathbf{x},k} &= (1 - \gamma_k) \widehat{\mathbf{R}}_{\mathbf{x}\mathbf{x},k-1} + \gamma_k (\widehat{\mathbf{x}}_{k|k} \widehat{\mathbf{x}}_{k|k}^T + \mathbf{P}_{k|k})
\end{aligned} \tag{22}$$

For the initialization, in absence of any side information, one can take $\mathbf{M}_0^{00} = 1/p_0 \mathbf{I}$, $\mathbf{M}_0^{10} = \mathbf{0}$, $\mathbf{M}_0^{11} = \mathbf{0}$, $\widehat{\mathbf{R}}_{\mathbf{x}\mathbf{x},0} = 1/p_0 \mathbf{I}$, $\widehat{\mathbf{R}}_{\mathbf{x}\mathbf{y},0} = \mathbf{0}$, $\widehat{\mathbf{R}}_{\mathbf{y}\mathbf{y},0} = \mathbf{0}$ where again p_0 is a very large number.

Adaptive processing means one iteration per new sample. Now, esp. in BASS, the parameters to be adapted represent filters. For short+long term AR processes with narrow bandwidths (well separable signals), these filters have a long memory, hence convergence transients are long and estimation efficiency is low.

G. Alternating MAP-ML KF (AMAPMLKF)

Joint MAP estimate for the state sequence \mathbf{x}_k and ML estimate for the parameters θ , which is typically solved by alternating optimization between the two parts. The ML estimate of θ is then obtained by performing least-squares (LS) estimation, given the state sequence, which is replaced by its estimate. The resulting algorithm is similar to the EM-KF with only the $\widehat{\mathbf{x}}$ terms kept in the matrices \mathbf{M}_k^{ij} .

H. Variational Bayes KF (VB-KF)

This is again an application of alternating optimization, but this time applied to the Kullback-Leibler distance between the true joint posterior pdf of state and parameters, and an approximate product form using Gaussian pdfs for both state and parameters. In this case, not only does the estimation of the parameters account for the estimation errors in the state, but symmetrically the state estimation now also accounts for the estimation errors in the parameters, so the state estimation now uses a modified form of the KF. The fixed-interval vs. fixed-lag smoothing issue also applies to VB-KF (and AMAPMLKF).

I. Gaussian SOEKF for bilinear system

For a discussion of the Gaussian SOEKF, consider the bilinear state update equation $\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k) + \mathbf{G}_k \mathbf{w}_k = \mathbf{F}_k \mathbf{x}_k + \mathbf{G}_k \mathbf{w}_k$. where \mathbf{F}_k depends on the (augmented) $M \times 1$ state vector \mathbf{x}_k linearly. We get for the SOEKF time update equations:

$$\begin{aligned}
\widehat{\mathbf{x}}_{k+1|k} &= \mathbf{f}(\widehat{\mathbf{x}}_{k|k}) + \frac{1}{2} \sum_{i=1}^M \phi_i \text{tr} \{ \mathbf{H}_{i,k} \mathbf{P}_{k|k} \} \\
\mathbf{P}_{k+1|k} &= \mathbf{D}_k \mathbf{P}_{k|k} \mathbf{D}_k^T + \mathbf{G}_k \mathbf{Q}_k \mathbf{G}_k^T \\
&\quad + \frac{1}{2} \sum_{i,j=1}^M \phi_i \phi_j^T \text{tr} \{ \mathbf{H}_{i,k} \mathbf{P}_{k|k} \mathbf{H}_{j,k} \mathbf{P}_{k|k} \} \\
\mathbf{H}_{i,k} &= \frac{\partial^2 \mathbf{f}_i}{\partial \mathbf{x}_k \partial \mathbf{x}_k^T} \Big|_{\mathbf{x}_k = \widehat{\mathbf{x}}_{k|k}}, \quad \mathbf{D}_k = \frac{\partial \mathbf{f}_i}{\partial \mathbf{x}_k^T} \Big|_{\mathbf{x}_k = \widehat{\mathbf{x}}_{k|k}}
\end{aligned} \tag{23}$$

where \mathbf{f}_i is the i^{th} component of the vector \mathbf{f} and ϕ_i is the $M \times 1$ vector with all zeros except for 1 in the i^{th} element. The terms involving the ϕ 's represent correction terms w.r.t. the EKF.

III. IMPROVED EKF (IEKF)

The EKF performs AKF by using the parameter estimate $\widehat{\theta}$ without accounting for its estimation error $\widehat{\theta}$. We can correct that. With the extended state update equation being of the form $\mathbf{x}_{k+1} = \mathbf{F}_k \mathbf{x}_k + \mathbf{G}_k \mathbf{w}_k$, the predicted state $\widehat{\mathbf{x}}_{k+1|k}$ is computed as $\widehat{\mathbf{x}}_{k+1|k} = \mathbf{F}(\widehat{\mathbf{x}}_{k|k}) \widehat{\mathbf{x}}_{k|k}$. Hence, we get for the state prediction error $\widetilde{\mathbf{x}}_{k+1|k} = \mathbf{x}_{k+1} - \widehat{\mathbf{x}}_{k+1|k}$ the following evolution:

$$\begin{aligned}
\widetilde{\mathbf{x}}_{k+1|k} &= (\mathbf{F}(\widehat{\mathbf{x}}_{k|k}) + \mathbf{F}(\widetilde{\mathbf{x}}_{k|k})) \mathbf{x}_k - \mathbf{F}(\widehat{\mathbf{x}}_{k|k}) \widehat{\mathbf{x}}_{k|k} + \mathbf{G}_k \mathbf{w}_k \\
&= \mathbf{F}(\widehat{\mathbf{x}}_{k|k}) \widetilde{\mathbf{x}}_{k|k} + \mathbf{G}_k \mathbf{w}_k + \mathbf{F}(\widetilde{\mathbf{x}}_{k|k}) \mathbf{x}_k
\end{aligned} \tag{24}$$

The last term in the second line is not considered in the classic EKF. Considering the independence of the predicted state and the state error estimation (assuming Gaussian signals), the prediction error covariance gets updated as follows

$$\mathbf{P}_{k+1|k} = \mathbf{F}(\widehat{\mathbf{x}}_{k|k}) \mathbf{P}_{k|k} \mathbf{F}(\widehat{\mathbf{x}}_{k|k})^T + \mathbf{G}_k \mathbf{Q}_k \mathbf{G}_k^T + \mathbf{O}_k \tag{25}$$

where $F(\hat{\mathbf{x}}_{k|k})$ equals \mathbf{D}_k in (23) and \mathbf{O}_k is the covariance matrix of $F(\hat{\mathbf{x}}_{k|k}) \mathbf{x}_k$ and is a correction term w.r.t. the basic EKF. Since $F(\hat{\mathbf{x}}_{k|k}) \mathbf{x}_k$ is linear in $\tilde{\mathbf{x}}_{k|k}$, we can find a matrix function \mathbf{D} such that

$$F(\tilde{\mathbf{x}}_{k|k}) \mathbf{x}_k = \mathbf{D}(\mathbf{x}_k) \tilde{\mathbf{x}}_{k|k} = \mathbf{D}(\hat{\mathbf{x}}_{k|k}) \tilde{\mathbf{x}}_{k|k} + \mathbf{D}(\tilde{\mathbf{x}}_{k|k}) \tilde{\mathbf{x}}_{k|k}.$$

Consequently, \mathbf{O}_k can be computed as

$$\mathbf{O}_k = \mathbf{D}(\hat{\mathbf{x}}_{k|k}) \mathbf{P}_{k|k} \mathbf{D}(\hat{\mathbf{x}}_{k|k})^T + \mathbf{N}_k$$

where \mathbf{N}_k is the covariance of $\mathbf{D}(\tilde{\mathbf{x}}_{k|k}) \tilde{\mathbf{x}}_{k|k}$ which is a fourth order moment that can be computed assuming a Gaussian distribution for $\tilde{\mathbf{x}}_{k|k}$. \mathbf{O}_k represents increased (awareness of) estimation error covariance which tends to dampen the filters (rendering them less frequency-selective), hence potentially allowing faster convergence transients.

IV. PERFORMANCE ANALYSIS: BACK TO BASICS

The \mathbf{y} are the measurements, the state sequence \mathbf{x} are random (Gaussian) parameters (in some applications they are nuisance parameters, in others parameters of interest), and θ are the deterministic parameters. In a *joint* estimation approach, we maximize the likelihood

$$f(\mathbf{y}, \mathbf{x} | \theta) \quad (26)$$

which means MAP for \mathbf{x} and ML for θ , with (joint) ML error covariance matrix $C_{\theta\theta}^J$ and CRB $_{\theta}^J$. In a *marginalized* estimation approach, we maximize the likelihood

$$f(\mathbf{y} | \theta) \quad (27)$$

which means ML for θ , with (marginalized) ML error covariance matrix $C_{\theta\theta}^M$ and associated CRB $_{\theta}^M$. Indeed, since the state \mathbf{x} is random, it can be eliminated from the likelihood. Asymptotically (in the amount of data \mathbf{y}), we get

$$C_{\theta\theta}^J \stackrel{(i)}{\geq} C_{\theta\theta}^M \stackrel{(ii)}{=} \text{CRB}_{\theta}^M \stackrel{(iii)}{\geq} \text{CRB}_{\theta}^J \quad (28)$$

where (ii) is due to $\hat{\theta}_{ML}^M$ being consistent, (i) is due to the inconsistency of $\hat{\mathbf{x}}_{MAP}$ which prevents $\hat{\theta}_{ML}^J$ from reaching its CRB. (iii) on the other hand follows from

$$\text{CRB}_{\theta}^{-M} = \text{CRB}_{\theta}^{-J} - \mathbf{E}_{\mathbf{y}|\theta} \text{Cov}_{\mathbf{x}|\mathbf{y},\theta} \left(\frac{\partial}{\partial \theta} \ln f(\mathbf{y} | \mathbf{x}, \theta) \right).$$

(vectorized [17]) where $\text{CRB}_{\theta}^{-M} = -\mathbf{E}_{\mathbf{y}|\theta} \frac{\partial^2}{\partial \theta \partial \theta^T} \ln f(\mathbf{y} | \theta)$, $\text{CRB}_{\theta}^{-J} = -\mathbf{E}_{\mathbf{y}|\theta} \mathbf{E}_{\mathbf{x}|\mathbf{y},\theta} \frac{\partial^2}{\partial \theta \partial \theta^T} \ln f(\mathbf{y} | \mathbf{x}, \theta)$. In other words, even though the CRBs would indicate otherwise, in terms of actual performance, joint estimation of the state and the parameters leads to worse parameter estimates than when the parameters are estimated in a marginalized fashion. So we get

- AMAPMLKF: $\hat{\theta}$ from $\hat{\mathbf{x}}$ only, $\hat{\mathbf{x}}$ from $\hat{\theta}$.
Converges to joint MAP-ML (ML-KF).
- EM-KF: $\hat{\theta}$ from $\hat{\mathbf{x}}$ and $\tilde{\mathbf{x}}$, $\hat{\mathbf{x}}$ from $\hat{\theta}$.

Now, it is known that the EM approach converges to the marginalized ML approach, so the EM-KF would be one approach to get this optimal performance.

- VB-KF: $\hat{\theta}$ from $\hat{\mathbf{x}}$ and $\tilde{\mathbf{x}}$, $\hat{\mathbf{x}}$ from $\hat{\theta}$ and $\hat{\theta}$.

Open issues: Can VB-KF do better than EM-KF? In VB-KF, the state estimate is improved, but the question is whether this allows an improved estimation of the parameters. Relative performance of EKF, RPEM-KF etc? And finally, relation of the IEKF to SOEKF, which also handles the bias due to the bilinearity.

REFERENCES

- [1] B. Anderson and J. Moore, *Optimal Filtering*. Prentice Hall (or more recently: Dover, 2005), 1979.
- [2] M. Athans, R. Wishner, and A. Bertolini, "Suboptimal state estimation for continuous-time nonlinear systems from discrete noisy measurements," *Automatic Control, IEEE Transactions on*, vol. 13, no. 5, pp. 504 – 514, oct 1968.
- [3] R. D. Bass, V. D. Norum, and L. Swartz, "Optimal multichannel nonlinear filtering," *J. Mufh. Anal. Appl.*, vol. 16, pp. 152 – 164, 1966.
- [4] S. Bensaid, A. Schutz, and D. Slock, "Single Microphone Blind Audio Source Separation Using EM-Kalman Filter and Short+Long Term AR Modeling," in *Proc. Int'l Conf. on Latent Variable Analysis and Signal Separation (LVA-ICA)*, Saint-Malo, France, Sept. 2010.
- [5] Consortium Partners, "Deliverable D2.3 "Hybrid Localization Techniques";" EC FP7 project WHERE, Tech. Rep., May 2010, URL: http://www.kn-s.dlr.de/where/public_documents_deliverables.php.
- [6] C. Couvreur and Y. Bresler, "Decomposition of a mixture of gaussian ar processes," *Acoustics, Speech, and Signal Processing, IEEE International Conference on*, vol. 3, pp. 1605–1608, 1995.
- [7] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the em algorithm," *JOURNAL OF THE ROYAL STATISTICAL SOCIETY, SERIES B*, vol. 39, no. 1, pp. 1–38, 1977.
- [8] M. Feder and E. Weinstein, "Parameter estimation of superimposed signals using the em algorithm," *Acoustics, Speech and Signal Processing, IEEE Transactions on*, vol. 36, no. 4, pp. 477–489, Apr 1988.
- [9] W. Gao, T. S., and J. Lehnert, "Diversity combining for ds/ss systems with time-varying, correlated fading branches," *Communications, IEEE Transactions on*, vol. 51, no. 2, pp. 284–295, Feb 2003.
- [10] R. Henriksen, "The truncated second-order nonlinear filter revisited," *Automatic Control, IEEE Transactions on*, vol. 27, no. 1, pp. 247 – 251, feb 1982.
- [11] A. H. Jazwinski, *Stochastic processes and filtering theory*, 1970.
- [12] M. Lenardi and D. Slock, "Estimation of Time-Varying Wireless Channels and Application to the UMTS W-CDMA FDD Downlink," in *Proc. European Wireless (EW)*, Florence, Italy, Feb. 2002.
- [13] L. Ljung, "Asymptotic behavior of the extended kalman filter as a parameter estimator for linear systems," *IEEE Transactions on Automatic Control*, vol. 24, no. 1, pp. 36 – 50, feb 1979.
- [14] —, *System Identification: Theory for the User*. Upper Saddle River, NJ: Prentice-Hall, 2002, 2nd edition.
- [15] L. Ljung and T. Söderström, *Theory and Practice of Recursive Identification*. Cambridge, MA: MIT Press, 1983.
- [16] R. Mehra, "On the identification of variances and adaptive kalman filtering," *IEEE Transactions on Automatic Control*, vol. 15, no. 2, pp. 175 – 184, apr 1970.
- [17] M. Moeneclaey, "On the True and the Modified Cramer-Rao Bounds for the Estimation of a Scalar Parameter in the Presence of Nuisance Parameters," *IEEE Trans. Communications*, vol. 46, Nov. 1998.
- [18] T. Sadiki and D. Slock, "'Bayesian Adaptive Filtering: Principles and Practical Approaches";" in *Proc. 12th European Signal Processing Conf. (EUSIPCO)*, Vienna, Austria, Sept. 2004.
- [19] P. Tichavsky, C. Muravchik, and A. Nehorai, "Posterior cramer-rao bounds for discrete-time nonlinear filtering," *Signal Processing, IEEE Transactions on*, vol. 46, no. 5, pp. 1386 –1396, may 1998.
- [20] H. L. V. Trees and K. L. Bell, *Bayesian Bounds for Parameter Estimation and Nonlinear Filtering/Tracking*. Wiley-IEEE Press, 2007.
- [21] J. Villares and G. Vazquez, "The quadratic extended kalman filter," in *Sensor Array and Multichannel Signal Processing Workshop Proceedings, 2004*, July 2004, pp. 480 – 484.
- [22] E. Weinstein, A. Oppenheim, M. Feder, and J. Buck, "Iterative and sequential algorithms for multisensor signal enhancement," *Signal Processing, IEEE Transactions on*, vol. 42, no. 4, pp. 846–859, Apr 1994.
- [23] J. Wiklander, "Performance comparison of the Extended Kalman Filter and the Recursive Prediction Error Method," Master's thesis, Linköping Univ., 2003, reg nr: LiTH-ISY-EX-3351, <http://www.essays.se/essay/5eb48d45e1/>.