# Smoothing Dynamic Systems with State-Dependent Covariance Matrices

Aleksandr Y. Aravkin and James V. Burke

*Abstract*— **Kalman filtering and smoothing algorithms are used in many areas, including tracking and navigation, medical applications, and financial trend filtering. One of the basic assumptions required to apply the Kalman smoothing framework is that error covariance matrices are known and given. In this paper, we study a general class of inference problems where covariance matrices can depend functionally on unknown parameters. In the Kalman framework, this allows modeling situations where covariance matrices may depend functionally on the state sequence being estimated. We present an extended formulation and generalized Gauss-Newton (GGN) algorithm for inference in this context. When applied to dynamic systems inference, we show the algorithm can be implemented to preserve the computational efficiency of the classic Kalman smoother. The new approach is illustrated with a synthetic numerical example.**

## I. INTRODUCTION

The Kalman filter [16] and smoother [19] are efficient algorithms to estimate the state of a dynamic system given noisy measurements. Over the last 10 years, the optimization perspective on the smoothing problem has produced many extensions to dynamic system estimation, including methods for smoothing systems with nonlinear process and measurement models [8], systems with nonlinear inequality constraints [10], robust Kalman smoothing [6], [5], [4], and smoothing of sparse systems [1].

In all of the above extensions, the variances of process and measurement errors are assumed to be fixed and known. In practice, these quantities are often not known, and may in fact depend on the state. For example, radar position errors are known to depend on the aspect angle as well as the position of the target [21]. In some applications [17], it may be of interest to do Kalman filtering in polar coordinates or other coordinates that induce a state dependence in measurement errors. Modeling of process error covariance may also be state dependent — for example, Bar-Shalom [7] suggests that the right choice of process noise level for flight tracking models depends on the turn rate range expected. Therefore if we are estimating turn rate as part of the state, the process noise level can be modeled as a function of (a portion of) the state.

These ideas motivate extensions of the standard Kalman smoothing formulation to situations where process and measurement variances have known functional dependence on the state. Several such extensions have already been considered. In [21], the Unscented Transform is used to fit models with state-dependent matrices acting on observation noise. Linear systems with additive observation noise where measurement error variance is a known function of the state are studied in [22]. Linear systems with control inputs transformed by state-dependent matrices are considered in [14]. Finally, adaptive system identification, as presented in [13], also falls into this class.

In this paper, we formulate the state-dependent covariance problem as a statistical estimation problem, and develop algorithms for obtaining the *maximum a posteriori* (MAP) estimate. The ideas presented here extend those developed in [11] for diagonal covariance matrices in kinetic tracer studies. In the theoretical development, we allow the process and measurement functions to be nonlinear, and we allow the functional dependence of covariance on the state to be nonlinear as well.

The paper proceeds as follows. In Section II, we review the statistical origins of the Kalman smoother, casting it as a structured nonlinear regression problem. We show that consideration of state-dependent variance in such a regression brings to the forefront terms that are usually ignored, and develop an extended MAP objective to optimize. The proposed formulation can be used for general nonlinear regression where variance depends in a known functional way on the parameters. In Section III we build a new algorithm for solving the resulting optimization problems, exploiting their *convex composite structure*. The key step is a special convex subproblem, which we solve in Section IV. In Section V, we show the necessary details required to implement this method for time series analysis, so as to preserve the computational complexity of the classic smoothing algorithms. In Section VI, we provide a numerical experiment using simulated data that demonstrates the performance of the new smoother and the potential modeling capabilities of the approach. We end with conclusions.

## II. KALMAN SMOOTHING WITH STATE-DEPENDENT UNCERTAINTY

The dynamic structure of the Kalman smoothing problem is specified as follows:

$$\begin{aligned}
\mathbf{x_1} &= g_1(x_0) + \mathbf{w_1}, \\
\mathbf{x_k} &= g_k(\mathbf{x_{k-1}}) + \mathbf{w_k} \quad k = 2, \ldots, N, \\
\mathbf{z_k} &= h_k(\mathbf{x_k}) + \mathbf{v_k} \qquad k = 1, \ldots, N,
\end{aligned} \quad (1)$$

where $g_k, h_k$ are known (nonlinear) process and measurement functions, and $\mathbf{w_k} \in \mathbb{R}^n$, $\mathbf{v_k} \in \mathbb{R}^{m(k)}$ are mutually independent Gaussian random variables with positive definite covariance matrices $Q_k$ and $R_k$, $\mathbf{x_k} \in \mathbb{R}^n$ are the unknown states, and $\mathbf{z_k} \in \mathbb{R}^{m(k)}$ are the observed measurements.

A. Y. Aravkin is with IBM T.J. Watson Research Center, Yorktown Heights, 10598, NY, USA saravkin@us.ibm.com

J. V. Burke is with the Department of Mathematics, University of Washington, Seattle, WA, USA jvburke@uw.edu

Considering model (1) and using Bayes' theorem, the conditional likelihood of the entire state sequence $\{x_k\}$ given the measurement sequence $\{z_k\}$ is given by

$$\mathbf{p}\left(\{x_k\}|\{z_k\}\right) \propto \mathbf{p}\left(\{z_k\}|\{x_k\}\right)\mathbf{p}\left(\{x_k\}\right), \qquad (2)$$

which in turn can be written in terms of the likelihood of state increments $\mathbf{p}(w_k)$ and measurement residuals $\mathbf{p}(v_k)$:

$$\mathbf{p}\left(\{z_k\}|\{x_k\}\right)\mathbf{p}\left(\{x_k\}\right) = \kappa \prod_{k=1}^{N} \mathbf{p}(v_k)\mathbf{p}(w_k)$$

$$= \kappa \prod_{k=1}^{N} \exp\Big(-\frac{1}{2}(z_k - h_k(x_k))^\top R_k^{-1}(z_k - h_k(x_k)) \qquad (3)$$
$$-\frac{1}{2}(x_k - g_k(x_{k-1}))^\top Q_k^{-1}(x_k - g_k(x_{k-1}))\Big),$$

where we define $g_1(x_0) = x_0$. The constant of proportionality $\kappa$ is usually ignored, since in classic models, the variance terms $Q_k$ and $R_k$ are fixed. It is given by

$$\kappa = \prod_{k=1}^{n} \frac{1}{(2\pi)^{n/2}\det(Q_k)}\frac{1}{(2\pi)^{m(k)/2}\det(R_k)}. \qquad (4)$$

Our main contribution here is to remove the assumption that $Q_k$ and $R_k$ are fixed and known, and instead model these covariance matrices as known $\mathcal{C}^2$ functions of the state. In this setting, $\kappa$ in (4) is no longer a constant, and must be accounted for. To design our approach, we assume that we are given the inverse Cholesky factors $Q_k^{-1/2}(x_k)$ and $R_k^{-1/2}(x_k)$ as functions of the state. For simple (e.g. diagonal variance) models, there is no loss of generality here; one can easily transform between different representations. When $Q_k$ and $R_k$ are full, however, the assumption that inverse Cholesky factors are available is essential for our approach. In addition to considerations of computational efficiency, the main motivation behind the assumption is the selection of an appropriate convex-composite model; this is explained in detail in the next section.

In order to develop a simpler notation for estimating the entire state sequence, we define functions $g : \mathbb{R}^{nN} \to \mathbb{R}^{nN}$ and $h : \mathbb{R}^{nN} \to \mathbb{R}^{M}$, with $M = \sum_k m_k$, from components $g_k$ and $h_k$ as follows:

$$g(x) = \begin{bmatrix} x_1 \\ x_2 - g_2(x_1) \\ \vdots \\ x_N - g_N(x_{N-1}) \end{bmatrix}, \quad h(x) = \begin{bmatrix} h_1(x_1) \\ h_2(x_2) \\ \vdots \\ h_N(x_N) \end{bmatrix}. \qquad (5)$$

Given a sequence of column vectors $\{u_k\}$ and matrices

$\{T_k\}$ we use the following notation:

$$\text{vec}(\{u_k\}) = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_N \end{bmatrix}, \; \text{diag}(\{T_k\}) = \begin{bmatrix} T_1 & 0 & \cdots & 0 \\ 0 & T_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & T_N \end{bmatrix},$$

$$\begin{aligned} R &= \text{diag}(\{R_k\}) & w &= \text{vec}(\{g_0, 0, \ldots, 0\}) \\ Q &= \text{diag}(\{Q_k\}) & z &= \text{vec}(\{z_1, \ldots, z_N\}) \\ x &= \text{vec}(\{x_k\}) & g_0 &= g_1(x_0). \end{aligned}$$
$$(6)$$

With this notation, and under Gaussian assumptions, the extended MAP object for the Kalman smoother, which incorporates state-dependent variance terms, is given by

$$\frac{1}{2}\|Q^{-1/2}(x)(g(x) - w)\|_2^2 + \frac{1}{2}\|R^{-1/2}(x)(h(x) - z)\|_2^2$$
$$-\log\det\left(Q^{-1/2}(x)\right) - \log\det\left(R^{-1/2}(x)\right). \qquad (7)$$

With (7) in front of us, we see why the log determinant terms play an important role. Without these terms, an optimization approach to minimize a weighted sum of squares will aim to drive $Q^{-1/2}(x)$ and $R^{-1/2}(x)$ to 0 if at all possible. The function $-\log(\cdot)$ acts as a barrier to prevent this from happening.

## III. CONVEX COMPOSITE FORMULATION AND ALGORITHM

We would like to apply the generalized Gauss-Newton methodology for minimizing convex composite functions [12] to the objective (7). The first step in this process is to write this objective in convex composite form, that is, in the form $f = \rho \circ F$, where $\rho$ is convex and $F$ is smooth. The choice of the functions $\rho$ and $F$ depend on how we wish to model the representation of the problem. The most straightforward way to rewrite (7) is the more general form

$$J(x) := \frac{1}{2}c(x)^T W(x)^{-1}c(x) + \frac{1}{2}\log\det(W(x)),$$

where $c : \mathbb{R}^{nN} \to \mathbb{R}^{M+nN}$ and $W : \mathbb{R}^{nN} \to \mathcal{S}_{++}^{M+nN}$ are smooth maps given by

$$c(x) = \begin{bmatrix} x_1 - g_0 \\ h_1(x_1) - z_1 \\ x_2 - g_2(x_1) \\ h_2(x_2) - z_2 \\ \vdots \\ x_N - g_N(x_{N-1}) \\ h_N(x_N) - z_N \end{bmatrix}, \qquad (8)$$

$$W(x) = \begin{bmatrix} Q_1(x_1) & 0 & & & \\ 0 & R_1(x_1) & & & \\ \vdots & & \ddots & \ddots & & \vdots \\ & & & Q_N(x_N) & 0 \\ & & & 0 & R_N(x_N) \end{bmatrix}$$
$$(9)$$

where $\mathcal{S}_{++}^{M+nN}$ is the cone of real symmetric $(M + nN) \times (M + nN)$ positive definite matrices.

Then $J = \hat{\rho} \circ \hat{F}$ with

$$\hat{\rho}(c, W) := \frac{1}{2}c^T W^{-1} c + \frac{1}{2}\log \det(W)$$
$$\hat{F}(x) = (c(x), W(x)).$$

Although the function $\hat{F}$ in this formulation can be assumed smooth, the function $\hat{\rho}$ is not convex. Indeed, $\hat{\rho}$ is the difference of two convex functions. When viewed as a function of $(c, W^{-1})$, it is still not jointly convex in these arguments.

Here, we propose an approach that applies in many practical settings and yields an efficient solution procedure. However, a price is paid in a more complex model for the covariance matrices. Specifically, we assume that the Cholesky factors for $Q_k^{-1}(x_k)$ and $R_k^{-1}(x_k)$ are given to us as explicit functions of the state. We denote these factors by $Q_k^{-1/2}(x_k)$ and $R_k^{-1/2}(x_k)$, respectively. In some settings, the matrices $Q_k(x_k)$ and $R_k(x_k)$ are modeled as diagonal matrices, in which case the inverse Cholesky factors are easily computed diagonal matrices. We provide an example of this type in the final section. Under this modeling assumption, the objective (7) can be abstracted to the more general form

$$K(x) = \frac{1}{2}c(x)^T V(x)^T V(x)c(x) - \log \circ \det[V(x)] , \quad (10)$$

where $c : \mathbb{R}^{nN} \to \mathbb{R}^{M+nN}$ is exactly as in (8) and $V : \mathbb{R}^{nN} \to \mathcal{L}^{M+nN}$ is given by

$$V(x) = \begin{bmatrix} Q_1^{-1/2} & 0 & & & & \\ 0 & R_1^{-1/2} & & & & \\ \vdots & & \ddots & \ddots & & \vdots \\ & & & & Q_N^{-1/2} & 0 \\ & & & & 0 & R_N^{-1/2} \end{bmatrix}, \quad (11)$$

where all blocks of V are functions of $x$, and $\mathcal{L}^{M+nN}$ is the subalgebra of $(M + nN) \times (M + nN)$ real lower triangular matrices. Throughout, we assume that both $c$ and $V$ are twice continuously differentiable and that $\text{dom}(K) := \{x \mid K(x) < +\infty\} = \{x \mid V(x) \in \mathcal{L}_{++}^{M+nN}\} \neq \emptyset$, where $\mathcal{L}_{++}^{M+nN}$ is the cone of $(M + nN) \times (M + nN)$ real lower triangular matrices with strictly positive entries on the diagonal. Now $K$ can be written in convex composite form $K(x) = \rho \circ F$ with

$$\rho(u, v) = \frac{1}{2}u^T u - \sum_i \log[v_i] \quad (12)$$

$$F(x) = \begin{bmatrix} F_1(x) \\ F_2(x) \end{bmatrix} = \begin{bmatrix} V(x)c(x) \\ \text{vec}[\{V_{ii}(x)\}] \end{bmatrix}. \quad (13)$$

Note that $\text{dom}(\rho) = \mathbb{R}^{M+nN} \times \mathbb{R}_{++}^{M+nN}$.

The direction finding subproblem in a Gauss-Newton method takes the form

$$\min_d \rho(F(x) + F'(x)d) + \frac{\omega}{2}d^T d ,$$

for some $\omega \geq 0$. The quadratic term $\frac{\omega}{2}d^T d$ is a regularization term that both guarantees the uniqueness of the solution and regulates its magnitude. The convergence analysis of methods of this type rely heavily on the difference function

$$\Delta(x; d) = \rho( F(x) + F'(x)d ) - K(x) . \quad (14)$$

which is important for both convergence criteria and the line search in the overall method. In particular, [12, Lemma 2.3]

$$K'(x : d) = \inf_{t>0} t^{-1}\Delta(x; td) \text{ for all } d \in \mathbb{R}^{nN}, \ x \in \text{dom}(K) \quad (15)$$

since, whenever $F_2(x) > 0$, then, for all $d$, $F_2(x) + F_2'(x)(td) > 0$ for all $t$ sufficiently small.

Linearizing the functions $F_i(x)$ in (12) yields approximations $\tilde{F}_i(x; d) := F_i(x) + F_i'(x)d$, which in turn gives the approximation

$$\tilde{K}(x; d) = \rho[\tilde{F}_1(x; d), \tilde{F}_2(x; d)] . \quad (16)$$

This is the objective for the direction finding subproblem. Here,

$$\begin{aligned} \tilde{F}_1(x; d) &= \left(V(x)\partial_x c(x) + (c(x)^T \otimes I_N)\partial_x V(x)\right) d \\ &\quad + V(x)c(x) \\ \tilde{F}_2(x; d) &= \text{vec}\left(\{V_{ii}(x) + \partial_x V_{ii}(x)d\}\right) . \end{aligned} \quad (17)$$

Note that we must be sure that $\tilde{F}_2(x; d)$ is component-wise greater than zero. For details of these derivations, see [2]. The Gauss-Newton subproblem is now given by

$$\begin{aligned} \bar{\Delta}(x) &:= \min_{d \in \mathbb{R}^{nN}} \Delta(x; d) + \frac{\omega}{2}d^T d \\ &= \min_{d \in \mathbb{R}^{nN}} \frac{1}{2}\tilde{F}_1(x; d)^T \tilde{F}_1(x; d) + \frac{\omega}{2}d^T d - \sum_i \log[\tilde{F}_2(x; d)] . \end{aligned} \quad (18)$$

Due to our assumptions on $c$ and $V$, these subproblems are always well defined, are convex, and have a unique solution which must always exist. In addition, they provide an estimate for the first-order optimality for $K$.

*Theorem 3.1:* [12, Theorem 3.6] Let $x \in \text{dom}(K)$. Then the following three statements are equivalent:

(i) $\bar{\Delta}(x) = 0$,
(ii) $\bar{d} = 0$ solves (18), and
(ii) $0 \in \partial K(x)$, where $\partial K(x) = F'(x)^T \partial \rho(F(x))$ is the generalized subdifferential of $K$ at $x$ [20, Definition 8.3].

In particular, these conditions imply that $x$ is a first-order stationary point for $K$.

If we ignore dependence on $x$, the optimization problem in (18) can be rewritten as

$$\min_{d \in \mathbb{R}^{nN}} \frac{1}{2}d^T C d + a^T d - \sum_i \log[s_i] \quad (19)$$
$$\text{s.t.} \quad s = \text{vec}\{V_{ii}\} + \partial_x \text{vec}\{V_{ii}\}d,$$

where

$$\begin{aligned} C &= \omega I + \widetilde{V}^T \widetilde{V} \\ a &= \widetilde{V}^T V c \\ \widetilde{V} &= [V\partial_x c + (c^T \otimes I_N)\partial_x V] \end{aligned} \quad (20)$$

Note that $a$ is the gradient of the quadratic portion of the extended objective with respect to the state sequence $x$. The quantity $(c^{\mathrm{T}} \otimes I_N)\partial_x V$ that appears in (20) can be rewritten as

$$(c^{\mathrm{T}} \otimes I_N)\partial_x V \quad = \quad \sum_{i=1}^{M+nN} c_i \partial_x V_{i\cdot} . \tag{21}$$

The Lagrangian associated with the extended subproblem (19) is given by

$$L(d, s, \lambda) = \frac{1}{2} d^{\mathrm{T}} Cd + a^{\mathrm{T}} d - \sum_i \log[s_i]$$
$$+ \lambda^{\mathrm{T}} \left( s - \mathrm{vec}\{V_{ii}\} - \partial_x \mathrm{vec}\{V_{ii}\}d \right), \tag{22}$$

for $s > 0$ and $\lambda > 0$. The corresponding optimality conditions state that a direction $d$ solves (19) if and only if there exist $s, \lambda \in \mathbb{R}_{++}^{M+nN}$ such that

$$\begin{aligned}
\nabla_d L &= Cd + a - \partial_x \mathrm{vec}\{V_{ii}\}^{\mathrm{T}} \lambda = 0 \\
\nabla_s L &= -D(s)^{-1}\mathbf{1} + \lambda = 0 \\
\nabla_\lambda L &= s - \mathrm{vec}\{V_{ii}\} - \partial_x \mathrm{vec}\{V_{ii}\}d = 0 ,
\end{aligned} \tag{23}$$

where $D(s) := \mathrm{diag}(s)$.

We refer to (19) as the *extended subproblem*. In the next section, we show that this problem can be rapidly solved. This motivates the *Extended Gauss-Newton* method for (10).

*Algorithm 3.1: Generalized Gauss-Newton Algorithm.*

The inputs to this algorithm are
- $x^0 \in \mathrm{dom}(K) := \{x : K(x) < \infty\} \subset \mathbb{R}^{Nn}$: initial estimate of state sequence
- $\varepsilon \geq 0$: overall termination criterion
- $\omega > 0$: regularization parameter
- $\beta \in (0, 1)$: step size selection parameter
- $\gamma \in (0, 1)$: line search step size factor

The steps are as follows:
1) Set the iteration counter $\nu = 0$.
2) (Generalized Gauss-Newton Step) Find descent direction $d^\nu$ solving (19) and set $\Delta_\nu := \bar{\Delta}(x^\nu) = \Delta(x^\nu; d^\nu)$. *Terminate if* $\Delta_\nu \geq -\varepsilon$.
3) (Line Search) Set

$$\begin{aligned}
t_\nu &= \max \gamma^i \\
&\text{s.t.} \quad i \in \{0, 1, 2, \cdots\} \text{ and} \\
&\text{s.t.} \quad \rho\left(F(x^\nu + \gamma^i d^\nu)\right) \leq \rho\left(F(x^\nu)\right) + \beta\gamma^i \Delta_\nu.
\end{aligned}$$

4) (Iterate) Set $x^{\nu+1} = x^\nu + t_\nu d^\nu$ and return to Step 2.

*Remark 3.2:* Note that the line search is well defined whenever $\Delta_\nu \neq 0$. Indeed, since whenever $\mathrm{diag}(V(x)) > 0$, then, for all $d$, $F_2(x; td) > 0$ for all $t$ sufficiently small. Consequently, since $x^0 \in \mathrm{dom}(K)$, we have $\{x^\nu\} \subset \mathrm{dom}(K)$. In addition, by (15), $K'(x^\nu; d^\nu) \leq \Delta_\nu < \bar{\Delta}(x^\nu)$, so that $\gamma^{-i}\Delta(x^\nu; \gamma^i d^\nu) < \beta\Delta_\nu$ for all $i$ sufficiently large.

*Theorem 3.2:* [Convergence] Let $\{x^\nu\}$ be generated by Algorithm 3.1 with $\epsilon = 0$. Then either the algorithm terminates finitely at a first-order stationary point for $K$ or the sequence $\{x^\nu\}$ is infinite and every cluster point of the sequence is a first-order stationary point for $K$. The proof is given in the appendix.

In the next section, we show how to solve the subproblem (19) for a general state-dependent covariance regression problem.

## IV. SOLVING THE EXTENDED SUBPROBLEM

To solve the direction finding subproblem, we apply a damped Newton method directly to the optimality conditions (23). We present the high-level method here, with details concerning Kalman smoothing given in the next section.

Let $E(s, \lambda, d)$ denote the KKT system given in (23), rearranged in a particular order:

$$E(s, \lambda, d) = \begin{bmatrix} s - \mathrm{vec}\{V_{ii}\} - \partial_x\mathrm{vec}\{V_{ii}\}d \\ D(s)D(\lambda)\mathbf{1} - \mathbf{1} \\ Cd + a - \partial_x\mathrm{vec}\{V_{ii}\}^{\mathrm{T}}\lambda \end{bmatrix} \tag{24}$$

Our goal is to find $(\bar{s}, \bar{\lambda}, \bar{d})$ for which $E(\bar{s}, \bar{\lambda}, \bar{d}) = 0$. We use damped Newton's method on E which requires solving the Newton equation

$$\nabla E(s, \lambda, d)\begin{bmatrix} \Delta s \\ \Delta\lambda \\ \Delta d \end{bmatrix} = -E(s, \lambda, d), \tag{25}$$

where $\nabla E(s, \lambda, d)$ is given by

$$\begin{bmatrix} I & 0 & -\partial_x\mathrm{vec}\{V_{ii}\} \\ D(\lambda) & D(s) & 0 \\ 0 & -\partial_x\mathrm{vec}\{V_{ii}\}^{\mathrm{T}} & C \end{bmatrix}. \tag{26}$$

Define

$$\mathcal{V} = \partial_x\mathrm{vec}\{V_{ii}\}.$$

Then, using row operations

$$\begin{aligned}
R_2 &= R_2 - D(\lambda)R_1 \\
R_3 &= R_3 + \mathcal{V}^{\mathrm{T}}D(s)^{-1}R_2,
\end{aligned}$$

we obtain the modified system

$$\begin{bmatrix} I & 0 & -\mathcal{V} \\ 0 & D(s) & D(\lambda)\mathcal{V} \\ 0 & 0 & \Phi \end{bmatrix}\begin{bmatrix} \Delta s \\ \Delta\lambda \\ \Delta d \end{bmatrix} = \begin{bmatrix} \alpha \\ \beta \\ \gamma \end{bmatrix}, \tag{27}$$

where

$$\Phi = C + \mathcal{V}^{\mathrm{T}}D(s)^{-1}D(\lambda)\mathcal{V} \tag{28}$$

and

$$\begin{aligned}
\alpha &= -s + \mathrm{vec}\{V_{ii}\} + \mathcal{V}d \\
\beta &= \mathbf{1} - D(\lambda)\left(\mathrm{vec}\{V_{ii}\} + \mathcal{V}d\right) \\
\gamma &= \mathcal{V}^{\mathrm{T}}\left(\lambda + D(s)^{-1}\left(\mathbf{1} - D(\lambda)(\mathrm{vec}\{V_{ii}\} + \mathcal{V}d)\right)\right) \\
&\quad - Cd - a .
\end{aligned}$$

By (20), the matrix $C$ is always positive definite and so $\Phi$ is always positive definite and hence invertible. This allows us to recover the Newton direction:

$$\begin{aligned}
\Delta d &= \Phi^{-1}\gamma \\
\Delta\lambda &= D(s)^{-1}\left(\mathbf{1} - D(\lambda)(\mathrm{vec}\{V_{ii}\} + \mathcal{V}(d + \Delta d))\right) \\
\Delta s &= -s + \mathrm{vec}\{V_{ii}\} + \mathcal{V}(d + \Delta d) .
\end{aligned} \tag{29}$$

Note that any damping scheme requires that $s > 0$ for the objective to be finite, and hence, in addition, we require that $\lambda > 0$ since we need $D(s)D(\lambda)\mathbf{1} = \mathbf{1}$ .
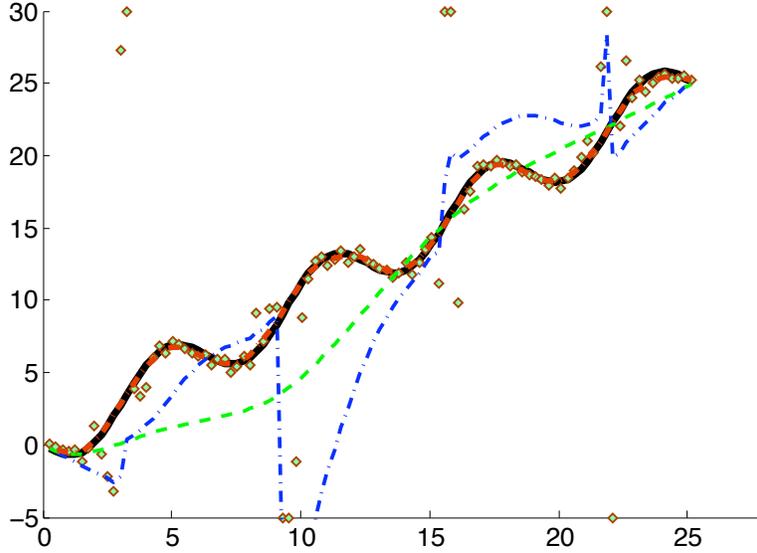
Fig. 1. True state $x_1$ (black curve), Extended Smoother estimate (thick red dash-dot), Kalman filter estimate (blue dash-dot) and Kalman Smoother estimate (green dashed curve). Measurements are displayed as diamonds, and those outside the axis range are displayed on the figure boundary.

## V. STRUCTURE OF THE EXTENDED KALMAN SMOOTHING OBJECTIVE

We now specify the method in the previous section to the Kalman smoothing problem, and demonstrate that the computational efficiency of the Kalman smoother can be preserved.

The functions $c(x)$ and $V(x)$ are given by (8) and (11).

With these definitions, objective $K(x)$ in (10) is exactly (7), and can be written explicitly as follows:

$$
\begin{aligned}
&\tfrac{1}{2}\left(c^{\mathrm{T}}(x)V(x)^{\mathrm{T}}V(x)c(x)\right) - \log \circ \det[V(x)] \\
=\ & \frac{1}{2}\sum_{k=1}^{N}\|[z_k - h_k(x_k)]\|^2_{R_k^{-T/2}(x_k)R_k^{-1/2}(x_k)} \\
&+ \frac{1}{2}\sum_{k=0}^{N}\|x_k - g_k(x_{k-1})\|^2_{Q_k^{-T/2}(x_k)Q_k^{-1/2}(x_k)} \\
&- \log\det(R_k^{-1/2}(x_k)) - \log\det(Q_k^{-1/2}(x_k)),
\end{aligned}
\tag{30}
$$

where, for any symmetric positive definite matrix $Q$, $\|u\|^2_Q := u^T Q u$.

We now derive the explicit forms for $C$ and $a$ in (20) for the Gauss-Newton subproblem (19). Recall that $C$ and $a$ are given by

$$
\begin{aligned}
C &= \omega I + [V\partial_x c + (c^{\mathrm{T}} \otimes I_N)\partial_x V]^{\mathrm{T}}[V\partial_x c + (c^{\mathrm{T}} \otimes I_N)\partial_x V], \\
a &= c^{\mathrm{T}}V^{\mathrm{T}}V\partial_x c + c^{\mathrm{T}}V^{\mathrm{T}}(c^{\mathrm{T}} \otimes I_N)\partial_x V.
\end{aligned}
$$

where

$$
(c^{\mathrm{T}} \otimes I_N)\partial_x V \;=\; \sum_{i=1}^{M+nN} c_i\partial_x V_{i\cdot},
$$

and

$$
\partial_x c(x) = \begin{bmatrix}
I & 0 & & & \\
H_1 & 0 & & \cdots & 0 \\
-G_2 & I & & \ddots & \\
0 & H_2 & 0 & & \cdots \\
& -G_3 & \ddots & & 0 \\
& & \ddots & H_{N-1} & 0 \\
& & & -G_N & I \\
& & & 0 & H_N
\end{bmatrix},
\tag{31}
$$

with $G_k = \partial_{x_k} g_{k+1}(x_k)$, $H_k = \partial_{x_k} h_k(x_k)$, and the dependence on $x$ has been suppressed to decrease the notational burden. Note that the matrix $G$ is invertible, and so $\partial_x c(x)$ is injective, that is, $\mathrm{Null}(\partial_x c(x)) = \{0\}$. In addition, since, we require $\mathrm{vec}\{V_{ii}\} > 0$ at every iteration, the matrix $V^T V$ is always positive definite. Consequently, the matrix $\partial_x c^T V^T V \partial_x c$ is always positive definite.

Define $\tilde{w}(x)$ and $\tilde{v}(x)$ in (6) by

$$
\begin{aligned}
\tilde{w}_k(x) &= x_k - g_k(x_{k-1}) & (32) \\
\tilde{v}_k(x) &= z_k - h_k(x_k) & (33)
\end{aligned}
$$

In the expressions below, we will use notation $\tilde{w}_{k,i}$ to mean the $i$th component of $\tilde{w}_k$.

The matrix $(c^{\mathrm{T}} \otimes I_N)\partial_x V$ has the following block struc-

ture:

$$[(c^{\mathrm{T}} \otimes I_N)\partial_x V] = \begin{bmatrix} \tilde{Q}_1 & 0 & 0 & 0 & 0 & \cdots & 0 \\ \tilde{R}_1 & 0 & 0 & 0 & 0 & \cdots & 0 \\ 0 & \tilde{Q}_2 & 0 & 0 & 0 & \cdots & 0 \\ 0 & \tilde{R}_2 & 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \ddots & \ddots & \cdots & 0 \\ 0 & 0 & 0 & \ddots & \ddots & Q_{N-1} & 0 \\ 0 & 0 & 0 & \ddots & \ddots & R_{N-1} & 0 \\ 0 & 0 & 0 & \cdots & \cdots & 0 & \tilde{Q}_N \\ 0 & 0 & 0 & \cdots & \cdots & 0 & \tilde{R}_N \end{bmatrix}$$

(34)

where

$$\tilde{Q}_i = \sum_{j=1}^{n} \tilde{w}_{i,j} \partial_{x_i} [Q_i^{-1/2}]_{j\cdot} \qquad (35)$$

$$\tilde{R}_i = \sum_{j=1}^{m(i)} \tilde{v}_{i,j} \partial_{x_i} [R_i^{-1/2}]_{j\cdot} \cdot \qquad (36)$$

Then we can write down the gradient $a$:

$$a = \begin{bmatrix} a_1^{\mathrm{T}} & a_2^{\mathrm{T}} & \cdots & a_N^{\mathrm{T}} \end{bmatrix}^{\mathrm{T}} ,$$

where

$$\begin{aligned} a_j = &-\tilde{v}_j^{\mathrm{T}} R_j^{-1}(x_j)\partial_{x(j)} h_j(x_j) + \tilde{w}_j^{\mathrm{T}} Q_j^{-1}(x_{j-1}) \\ &+ \tilde{w}_{j+1}^{\mathrm{T}} Q_{j+1}^{-1}(x_j)\partial_{x(j)} g_{j+1}(x_j) \\ &+ \tilde{w}_j^{\mathrm{T}} Q_j^{-T/2}\tilde{Q}_j + \tilde{v}_j^{\mathrm{T}} R_j^{-T/2}\tilde{R}_j . \end{aligned}$$

(37)

Using (31) and (34), we can form

$$\Psi := V\partial_x c + (c^{\mathrm{T}} \otimes I_N)\partial_x V,$$

and obtain a closed form solution for $C = \Psi^T \Psi$:

$$C = \begin{bmatrix} C_1 & A_2^{\mathrm{T}} & 0 & \\ A_2 & C_2 & A_3^{\mathrm{T}} & 0 \\ 0 & \ddots & \ddots & \ddots \\ & 0 & A_N & C_N \end{bmatrix}$$

$$C_k = \omega I + [Q_k^{-1/2} + \tilde{Q}_k]^{\mathrm{T}}[Q_k^{-1/2} + \tilde{Q}_k] + \nabla g_{k+1}^{\mathrm{T}} Q_{k+1}^{-1} \nabla g_{k+1}$$

$$+ [-R_k^{-1/2}\nabla h_k + \tilde{R}_k]^{\mathrm{T}}[-R_k^{-1/2}\nabla h_k + \tilde{R}_k]$$

$$A_k = -(Q_k^{-1/2} + \tilde{Q}_k)^{\mathrm{T}} Q_k^{-1/2} \nabla g_k .$$

(38)

*Remark 5.1:* The matrix (38) is block tridiagonal, and so it can be inverted with effort $O(n^3 N)$ using any of the algorithms in [9], [3]. The proof sketch is given in the appendix.

Recall the direction finding equation (25) in the Extended GN algorithm:

$$\nabla E(s, \lambda, d) \begin{bmatrix} \Delta s \\ \Delta \lambda \\ \Delta d \end{bmatrix} = -E(s, \lambda, d).$$

The solution to this system is given by (29), with $\Phi$ as in (28), $C$ as in (38), and $[\mathcal{V}^{\mathrm{T}} D(s)^{-1} D(\lambda)\mathcal{V}]_k$ given by

$$\partial_{x(k)}\mathrm{diag}\{Q_k^{-1/2}\}^{\mathrm{T}} D(s_{Q_k})^{-1} D(\lambda_{Q_k})\partial_{x(k)}\mathrm{diag}\{Q_k^{-1/2}\}$$

$$+ \partial_{x(k)}\mathrm{diag}\{R_k^{-1/2}\}^{\mathrm{T}} D(s_{R_k})^{-1} D(\lambda_{R_k})\partial_{x(k)}\mathrm{diag}\{R_k^{-1/2}\} \cdot$$

## VI. NUMERICAL RESULTS

In this section, we present some numerical experiments to show the advantages and modeling possibilities of the new Kalman smoother. The simulation model we consider is similar to the one presented in [10]. The 'ground truth' time series for this simulated example is given by

$$x(t) = \begin{bmatrix} 1 - 2\cos(t) \\ t - 2\sin(t) \end{bmatrix} .$$

The time between measurements is a constant denoted by $\Delta t$. The models for the mean of $x_k$ given $x_{k-1}$ and for process covariance $Q_k$ [15], [18] are

$$g_k(x_{k-1}) = \begin{bmatrix} 1 & 0 \\ \Delta t & 1 \end{bmatrix} x_{k-1} , \quad Q_k = \begin{bmatrix} \Delta t & \Delta t^2/2 \\ \Delta t^2/2 & \Delta t^3/3 \end{bmatrix} .$$

The measurement model for the mean of $z_k$ given $x_k$ is $h_k(x_k) = x_{2,k}$ , where $x_{2,k}$ denotes the second component of $x_k$.

The main innovation of the example is in the measurement variance model. The smoother takes inverse Cholesky factors as input, and these are assumed to be $3 - x_{1,k}$. Then the variance model is given by $R_k(x_k) = (3 - x_{1,k})^{-2}$. The measurements were generated using the measurement model, from two full periods of the time series $x(t)$, with $N = 100$ discrete time points equally spaced over the interval $[0, 4\pi]$, and with noise sampled from $N(0, R_k(x_k))$. Since the true state $x_1$ varies in the interval $[-1, 3]$, the variance for the observations goes to infinity when $t$ is a multiple of $\pi$.

This simulation illustrates a situation where the measurements are very reliable for some state values, but completely unpredictable for others. This phenomenon may occur for example if sensors report garbage values when the attitude of a vehicle is in a particular configuration. The measurement model presented here can be easily adapted by the user to take their beliefs about the system into account. The main point is that as long as the inverse Cholesky factors for the variance can be coded as a smooth function of the state, smoothed estimates for state values can be obtained taking into account this bad behavior of the measurements.

The result of the simulation is shown in Figure 1. The extended Kalman smoother (thick red dash-dot) is able to recover the ground truth (shown in black) with no appreciable difference. The Kalman filter (thin blue dash-dot) is strongly affected by the outlying measurements, as expected. The Kalman smoother (green dashed) is able to smooth the measurements, but cannot pick up the oscillations of the ground truth, which are small in magnitude compared to the size of the errors.

This last point is the most important — it is not just the magnitude of the outliers that makes the Kalman smoother fail, although it can be seen to be rather far off the ground truth. The biggest challenge of the situation presented is knowing which measurements to trust, since this information depends on the state being estimated.

## VII. Conclusions

In this paper, we presented extended formulations for modeling dynamic systems in cases where the covariance matrices are known functions of the state. The formulation includes variance-control terms arising from statistical modeling assumptions. These terms give rise to an extended convex-composite structure, and we propose a new method, the extended Gauss-Newton, which repeatedly solves extended convex subproblems by exploiting their KKT optimality conditions. When applied to dynamic inference problems, the proposed approach preserves the complexity of the classic Kalman smoother.

## References

[1] D. Angelosante, S.I. Roumeliotis, and G.B. Giannakis. Lasso-kalman smoother for tracking sparse signals. In *Signals, Systems and Computers, 2009 Conference Record of the Forty-Third Asilomar Conference on*, pages 181–185, nov. 2009.

[2] A.Y. Aravkin. *Robust Methods with Applications to Kalman Smoothing and Bundle Adjustment*. PhD thesis, University of Washington, Seattle, WA, June 2010.

[3] A.Y. Aravkin, B M Bell, J V Burke, and G Pillonetto. New Stability Results and Algorithms for Block Tridiagonal Systems, with Applications to Kalman Smoothing. *http://arxiv.org/abs/1303.5237*, 2013.

[4] A.Y. Aravkin, B.M. Bell, J.V. Burke, and G. Pillonetto. An $\ell_1$-laplace robust kalman smoother. *Automatic Control, IEEE Transactions on*, 56(12):2898–2911, dec. 2011.

[5] A.Y. Aravkin, J.V. Burke, and G. Pillonetto. Robust and trend-following kalman smoothers using student's t. In *IFAC, 16th Symposium of System Identification*, oct. 2011.

[6] A.Y. Aravkin, J.V. Burke, and G. Pillonetto. A statistical and computational theory for robust and sparse Kalman smoothing. In *IFAC, 16th Symposium of System Identification*, oct. 2011.

[7] Y. Bar-Shalom, X. Rong Li, and T. Kirubarajan. *Estimation with Applications to Tracking and Navigation*. John Wiley and Sons, 2001.

[8] B.M. Bell. The iterated Kalman smoother as a Gauss-Newton method. *SIAM J. Optimization*, 4(3):626–636, August 1994.

[9] B.M. Bell. The marginal likelihood for parameters in a discrete Gauss-Markov process. *IEEE Transactions on Signal Processing*, 48(3):626–636, August 2000.

[10] B.M. Bell, J.V. Burke, and G. Pillonetto. An inequality constrained nonlinear kalman-bucy smoother by interior point likelihood maximization. *Automatica*, 45(1):25–33, January 2009.

[11] B.M. Bell, J.V. Burke, and A. Schumitzky. A relative weighting method for estimating parameters and variances in multiple data sets. *Computational Statistics and Data Analysis*, 22:119–135, 1996.

[12] J.V. Burke. Descent methods for composite nondifferentiable optimization problems. *Mathematical Programming*, 33:260–279, 1985.

[13] C.K. Chui and G. Chen. *Kalman Filtering: with Real-Time Applications*. Springer series in information sciences. Springer, 2008.

[14] A. Dutka, H. Javaherian, and M.J. Grimble. State-dependent Kalman filters for robust engine control. In *Proceedings of the 2006 American Control Conference*, pages 1185–1190, 2006.

[15] Andrew Jazwinski. *Stochastic Processes and Filtering Theory*. Dover Publications, Inc, 1970.

[16] R. E. Kalman. A new approach to linear filtering and prediction problems. *Transactions of the AMSE - Journal of Basic Engineering*, 82(D):35–45, 1960.

[17] G.A. McIntyre and K.J. Hintz. A comparison of several maneuvering target tracking models. In *SPIE Conference on Signal Processing, Sensor Fusion, and Target Recognition VII*, volume 3374, pages 48–63, April 1998.

[18] Bernt Oksendal. *Stochastic Differential Equations*. Springer, sixth edition, 2005.

[19] H. E. Rauch, F. Tung, and C. T. Striebel. Maximum likelihood estimates of linear dynamic systems. *AIAA J.*, 3(8):1145–1150, 1965.

[20] R. Tyrrell Rockafellar and Roger J-B. Wets. *Variational Analysis*, volume 317 of *A Series of Comprehensive Studies in Mathematics*. Springer, 1998.

[21] M. Stakkeland, O. Overrein, and E.F. Brekke. Tracking of targets with state dependent measurement errors using recursive BLUE filters. In *12th International Conference on Information Fusion*, pages 2052–2061, 2009.

[22] B. Zehnwirth. A generalization of the Kalman filter for models with state-dependent observation variance. *Journal of the Amer. Stat. Association*, 83(401):164–167, March 1988.

## VIII. Appendix

### A. Proof sketch of Remark 5.1

Recall the matrix inversion lemma:

*Lemma 8.1 (Matrix Inversion Lemma):* Assume matrices $M_1 \in \mathbb{R}^{m_1 \times m_1}$ and $M_2 \in \mathbb{R}^{m_2 \times m_2}$ are symmetric positive definite. Then for any matrix $U \in \mathbb{R}^{m_1 \times m_2}$, the following matrix is also symmetric positive definite:

$$S = M_1^{-1} - M_1^{-1} U (M_2 + U^T M_1^{-1} U)^{-1} U^T M_1^{-1}. \quad (39)$$

From this, we get an immediate and useful corollary:

*Corollary 8.2:* For a positive definite matrix $M_2$, and any matrix $U$, we have

$$\|U^T(M_2 + UU^T)^{-1}U\|_2 < 1,$$

where $\|\cdot\|_2$ denotes the spectral norm.

Simply take $M_1 = I$ in lemma 8.1. The conclusion of the lemma gives the result.

Corollary 8.2 can be applied to show that the algorithms in [3] yield invertible blocks at every iteration. Application of these algorithms to $C$ in (38) requires inverting matrices of form

$$H^+ + V^T(I - U^T(H + UU^T)^{-1}U)V,$$

where $H$ and $H^+$ are always positive definite. The second term in the sum is clearly seen to be positive semidefinite by Corollary 8.2. To be specific, at the first iteration,

$$
\begin{aligned}
V &= Q_2^{-1/2} + \tilde{Q}_2, \\
H &= \omega I + (Q_1^{-1/2} + \tilde{Q}_1)^T(Q_1^{-1/2} + \tilde{Q}_1) \\
&\quad + (-R_1^{-1/2}\nabla h_1 + \tilde{R}_1)^{\mathrm{T}}(-R_1^{-1/2}\nabla h_1 + \tilde{R}_1), \\
U &= \nabla g_2^T Q_2^{-1/2}, \\
H^+ &= \omega I + (-R_2^{-1/2}\nabla h_2 + \tilde{R}_2)^{\mathrm{T}}(-R_2^{-1/2}\nabla h_2 + \tilde{R}_2)
\end{aligned}
$$

The full argument can be made by induction, but we do not include it here.

### B. Proof of Theorem 3.2

The algorithm can only terminate if

$$0 = \Delta_\nu = \Delta(x^\nu; d^\nu) \leq \Delta(x^\nu; d^\nu) + \tfrac{\omega}{2}\|d^\nu\|^2 = \bar{\Delta}(x^\nu) \leq 0,$$

i.e., $\bar{\Delta}(x^\nu) = 0$, or equivalently, $x^\nu$ is a first-order stationary point for $K$ by Theorem 3.1.

Assume that the algorithm does not terminate finitely, and let $\hat{x}$ be a cluster point of the sequence of iterates $\{x^\nu\}$. Since this is a descent algorithm, it is necessarily the case that $K(x^\nu) \downarrow K(\hat{x})$. Let $J \subset \mathbb{N}$ and $\hat{x} \in \mathbb{R}^{nN}$ be such that $x^\nu \xrightarrow{J} \hat{x}$, and suppose to the contrary that $\hat{x}$ is not a first-order staionary point for $K$, i.e. $\bar{\Delta}(\hat{x}) < 0$. We now use the optimality conditions (23) to show that the subsequence of search directions $\{d^\nu\}_J$ is bounded. Let

$(s^\nu, \lambda^\nu, d^\nu)$ denote the triple satisfying these conditions for each $x^\nu$. Then multiplying the second condition in (23) by $s^\nu$ and the third condition in (23) by $\lambda^\nu$ and combining, we find that

$$M + nN = (\lambda^\nu)^{\mathrm{T}}(\mathrm{vec}\{V_{ii}(x^\nu)\} + \partial_x \mathrm{vec}\{V_{ii}(x^\nu)\}d^\nu).$$

By combining this with the first condition in (23), we find that

$$\begin{aligned} M + nN &= (\lambda^\nu)^{\mathrm{T}}\mathrm{vec}\{V_{ii}(x^\nu)\} + (C(x^\nu)d^\nu + a(x^\nu))^{\mathrm{T}}d^\nu \\ &\geq a(x^\nu)^{\mathrm{T}}d^\nu + (d^\nu)^{\mathrm{T}}C(x^\nu)d^\nu \\ &\geq \omega\|d^\nu\|_2^2 - \|a(x^\nu)\|_2\|d^\nu\|_2, \end{aligned}$$

where the first inequality follows since $\lambda^\nu > 0$ and $\mathrm{diag}(V(x^\nu)) > 0$, and the second inequality follows from (20). Consequently, the subsequence $\{d^\nu\}_J$ is bounded due to the continuity of

$$a(x) = \nabla\left(\frac{1}{2}c(x)^{\mathrm{T}}V(x)^{\mathrm{T}}V(x)c(x)\right).$$

With no loss in generality, we can now assume that there is a $\hat{d}$ such that $d^\nu \xrightarrow{J} \hat{d}$. By continuity,

$$\Delta_\nu = \Delta(x^\nu; d^\nu) \to \Delta(\hat{x}; \hat{d}).$$

Moreover, for all $d \in \mathbb{R}^{M+nN}$,

$$\Delta(x^\nu; d^\nu) + \tfrac{\omega}{2}(d^\nu)^{\mathrm{T}}d^\nu \leq \Delta(x^\nu; d) + \tfrac{\omega}{2}d^{\mathrm{T}}d .$$

Taking the limit over $\nu \in J$ gives

$$\Delta(\hat{x}; \hat{d}) + \tfrac{\omega}{2}\hat{d}^{\mathrm{T}}\hat{d} \leq \Delta(\hat{x}; d) + \tfrac{\omega}{2}d^{\mathrm{T}}d .$$

Therefore, $\Delta(\hat{x}; \hat{d}) + \tfrac{\omega}{2}\hat{d}^{\mathrm{T}}\hat{d} = \bar{\Delta}(\hat{x})$.

Recall our working assumption that $\bar{\Delta}(\hat{x}) < 0$. Since we have just shown that $\Delta_\nu \to \bar{\Delta}(\hat{x})$, we must therefore have $\xi := \sup_{\nu \in J}\Delta_\nu < 0$. Since $K(x^{\nu+1}) - K(x^\nu) \leq \beta t_\nu \Delta_\nu$ with $K(x^\nu)$ convergent, we must have $t_\nu \to 0$. Again, with no loss in generality, $1 > t_\nu \downarrow_J 0$. By continuity, there are $\delta > 0$ and $\mu > 0$ such that

$$\mathrm{diag}(V(x)) \in \mathrm{diag}(V(\hat{x})) + \mu\mathbb{B} \subset \mathbb{R}_{++}^{M+nN} \quad \text{and}$$
$$\mathrm{diag}(V(x) + V'(x)d) \in \mathrm{diag}(V(\hat{x})) + \mu\mathbb{B}$$

for all $x \in \hat{x} + \gamma^{-1}\delta\mathbb{B}$ and $d \in \gamma^{-1}\delta\mathbb{B}$. Since $\{d^\nu\}_J$ is bounded and $t_\nu \downarrow_J 0$, we can assume with no loss in generality that $x^\nu \in \hat{x} + \delta\mathbb{B}$ and $t_\nu d^\nu \in \delta\mathbb{B}$ with $1 > t_\nu$ for all $\nu \in J$. Let

$$\kappa_1 := \sup\left\{\frac{1}{2}\|V(x)c(x)\|_2^2 \ : \ x \in \hat{x} + \delta\mathbb{B}\right\}.$$

Since $\kappa_1\mathbb{B} \times \mu\mathbb{B} \subset \mathrm{intr}(\mathrm{dom}(\rho))$, $\rho$ is Lipschitz continuous on $\kappa_1\mathbb{B} \times \mu\mathbb{B}$ with Lipschitz constant $\kappa_2 > 0$. Also, the function $F$ defined in (13) is such that $F'$ is Lipschitz continuous on $\hat{x} + \gamma^{-1}\delta\mathbb{B}$ with Lipschitz constant $\kappa_3 > 0$.

Due to the way the step sizes $t_\nu$ are chosen and the fact that $t_\nu d^\nu \in \delta\mathbb{B}$ and $1 > t_\nu$ for all $\nu \in J$, we have

$$\begin{aligned} \gamma^{-1}t_\nu\beta\bar{\Delta}(x^\nu) &< K(x^\nu + \gamma^{-1}t_\nu d^\nu) - K(x^\nu) \\ &\leq \Delta(x^\nu; \gamma^{-1}t_\nu d^\nu) \\ &\quad + \kappa_2\|F(x^\nu + \gamma^{-1}t_\nu d^\nu) - F(x^\nu) - F'(x^\nu)(\gamma^{-1}t_\nu d^\nu)\| \\ &\leq \gamma^{-1}t_\nu\bar{\Delta}(x^\nu) + \frac{\kappa_2\kappa_3}{2}(\gamma^{-1}t_\nu)^2\|d^\nu\|_2^2. \end{aligned}$$

Consequently,

$$\begin{aligned} 0 &< (1-\beta)\bar{\Delta}(x^\nu) + \frac{\kappa_2\kappa_3}{2}(\gamma^{-1}t_\nu)\|d^\nu\|_2^2 \\ &\leq (1-\beta)\xi + \frac{\kappa_2\kappa_3}{2}(\gamma^{-1}t_\nu)\|d^\nu\|_2^2. \end{aligned}$$

Taking the limit over $\nu \in J$ in this inequality gives the contradiction $0 \leq (1-\beta)\xi < 0$.