

PROBABILISTIC FILTER AND SMOOTHER FOR VARIATIONAL INFERENCE OF BAYESIAN LINEAR DYNAMICAL SYSTEMS - APPENDIX

*Julian Neri** *Roland Badeau†* *Philippe Depalle**

*McGill University, CIRMMT, Montréal, Canada.

†LTCI, Télécom Paris, Institut Polytechnique de Paris, France.

Contents

A	Forward pass (filtering) derivation	2
A.1	Predictive distribution	2
A.2	Marginal likelihood	4
A.3	Marginal posterior	5
B	Backward pass (smoothing) derivation	8
C	Application: variational M step derivation	12
C.1	Joint distribution	12
C.2	Optimal distribution over parameters	12
D	Quadratic forms of random variables	14
D.1	General result	14
D.2	Case 1: Conjugate prior structure	14
D.3	Case 2: Non-conjugate prior structure	16
E	Schur complements	17
F	References	18

*Thanks to the Natural Sciences and Engineering Research Council of Canada (NSERC) Discovery Grant (RGPIN- 2018-05662) for funding.

A. FORWARD PASS (FILTERING) DERIVATION

The forward pass calculates the statistics of the marginal distribution

$$q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n}) = \frac{p(\mathbf{y}_n | \mathbf{x}_n) q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n-1})}{q_{\mathbf{x}}(\mathbf{y}_n | \mathbf{y}_{1:n-1})} \quad (1)$$

$$= \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, \mathbf{V}_n). \quad (2)$$

The predictive distribution $q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n-1})$, marginal likelihood $q_{\mathbf{x}}(\mathbf{y}_n | \mathbf{y}_{1:n-1})$, and marginal posterior $q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n})$ are derived in turn. These derivations involve quadratic forms of random variables, as detailed in Section D, and Schur complements, provided in Section E.

A.1. Predictive distribution

First, we evaluate the predictive distribution and show that it is Gaussian distributed with mean \mathbf{m}_{n-1} and covariance \mathbf{P}_{n-1} . We retrieve the predictive distribution by marginalizing out \mathbf{x}_{n-1} from the transition probability weighted by the previous marginal:

$$q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n-1}) = \int p(\mathbf{x}_n | \mathbf{x}_{n-1}) q_{\mathbf{x}}(\mathbf{x}_{n-1} | \mathbf{y}_{1:n-1}) d\mathbf{x}_{n-1} \quad (3)$$

$$= \int \mathcal{N}(\mathbf{x}_n | \mathbf{A}\mathbf{x}_{n-1} + \mathbf{B}\mathbf{u}_n, \mathbf{Q}) \mathcal{N}(\mathbf{x}_{n-1} | \boldsymbol{\mu}_{n-1}, \mathbf{V}_{n-1}) d\mathbf{x}_{n-1} \quad (4)$$

$$= \mathcal{N}(\mathbf{x}_n | \mathbf{m}_{n-1}, \mathbf{P}_{n-1}) \quad (5)$$

Collecting the quadratic terms over \mathbf{x}_n and \mathbf{x}_{n-1} , we get

$$\Omega_{11} = \mathbf{V}_{n-1}^{-1} + \langle \mathbf{A}^T \mathbf{Q}^{-1} \mathbf{A} \rangle \quad (6)$$

$$= \mathbf{V}_{n-1}^{-1} + \Sigma_{AQ} + \langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle \quad (7)$$

$$\Omega_{12} = \langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \quad (8)$$

$$\Omega_{21} = \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle \quad (9)$$

$$\Omega_{22} = \langle \mathbf{Q}^{-1} \rangle \quad (10)$$

Using the Schur complement (refer to Section E for complements) gives the following analytic inverse

$$\mathbf{P}_{n-1} = (\Omega_{22} - \Omega_{21}\Omega_{11}^{-1}\Omega_{12})^{-1} \quad (11)$$

$$= \Omega_{22}^{-1} + \Omega_{22}^{-1}\Omega_{21}\mathbf{F}_{11}^{-1}\Omega_{12}\Omega_{22}^{-1} \quad (12)$$

where

$$\mathbf{F}_{11}^{-1} = (\Omega_{11} - \Omega_{12}\Omega_{22}^{-1}\Omega_{21})^{-1}. \quad (13)$$

Substituting in and simplifying gives

$$\mathbf{F}_{11}^{-1} = (\mathbf{V}_{n-1}^{-1} + \Sigma_{AQ} + \langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle - \langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{Q}^{-1} \rangle^{-1} \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle)^{-1} \quad (14)$$

$$= (\mathbf{V}_{n-1}^{-1} + \Sigma_{AQ} + \cancel{\langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle} - \cancel{\langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle})^{-1} \quad (15)$$

$$= (\mathbf{V}_{n-1}^{-1} + \Sigma_{AQ})^{-1}. \quad (16)$$

It is computationally costly and numerically problematic to invert \mathbf{V}_{n-1}^{-1} at every time n . We apply the Woodbury matrix identity to retrieve the following preferred analytic expression for \mathbf{F}_{11}^{-1} ,

$$\mathbf{F}_{11}^{-1} = (\mathbf{V}_{n-1}^{-1} + \Sigma_{AQ})^{-1} \quad (17)$$

$$= \mathbf{V}_{n-1} - \mathbf{V}_{n-1}(\mathbf{I} + \Sigma_{AQ}\mathbf{V}_{n-1})^{-1}\Sigma_{AQ}\mathbf{V}_{n-1} \quad (18)$$

$$= (\mathbf{I} - \mathbf{V}_{n-1}(\mathbf{I} + \Sigma_{AQ}\mathbf{V}_{n-1})^{-1}\Sigma_{AQ})\mathbf{V}_{n-1} \quad (19)$$

$$= \mathbf{G}_{n-1}\mathbf{V}_{n-1} \quad (20)$$

where we have defined

$$\mathbf{G}_{n-1} = \mathbf{I} - \mathbf{V}_{n-1}(\mathbf{I} + \boldsymbol{\Sigma}_{AQA}\mathbf{V}_{n-1})^{-1}\boldsymbol{\Sigma}_{AQA}. \quad (21)$$

Substituting this result back into the expression for the covariance \mathbf{P}_{n-1} and simplifying gives

$$\mathbf{P}_{n-1} = \boldsymbol{\Omega}_{22}^{-1} + \boldsymbol{\Omega}_{22}^{-1}\boldsymbol{\Omega}_{21}\mathbf{F}_{11}^{-1}\boldsymbol{\Omega}_{12}\boldsymbol{\Omega}_{22}^{-1} \quad (22)$$

$$= \langle \mathbf{Q}^{-1} \rangle^{-1} + \cancel{\langle \mathbf{Q}^{-1} \rangle^{-1}}\cancel{\langle \mathbf{Q}^{-1} \rangle}\langle \mathbf{A} \rangle \mathbf{G} n - 1 \mathbf{V}_{n-1} \langle \mathbf{A}^\top \rangle \cancel{\langle \mathbf{Q}^{-1} \rangle}\cancel{\langle \mathbf{Q}^{-1} \rangle}^{-1} \quad (23)$$

$$= \langle \mathbf{Q} \rangle + \langle \mathbf{A} \rangle \mathbf{G}_{n-1} \mathbf{V}_{n-1} \langle \mathbf{A}^\top \rangle. \quad (24)$$

Next, we derive the expression for the mean \mathbf{m}_{n-1} . Completing the square for \mathbf{x}_n and \mathbf{x}_{n-1} , the linear terms over \mathbf{x}_{n-1} and \mathbf{x}_n are as follows.

$$\ell_1 = \mathbf{V}_{n-1}^{-1}\boldsymbol{\mu}_{n-1} - \langle \mathbf{A}^\top \mathbf{Q}^{-1} \mathbf{B} \rangle \mathbf{u}_n \quad (25)$$

$$= \mathbf{V}_{n-1}^{-1}\boldsymbol{\mu}_{n-1} - (\langle \mathbf{A}^\top \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle + \boldsymbol{\Sigma}_{AQB}) \mathbf{u}_n \quad (26)$$

$$\ell_2 = \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle \mathbf{u}_n \quad (27)$$

The mean is then given by

$$\mathbf{m}_{n-1} = \mathbf{P}_{n-1}(\ell_2 - \boldsymbol{\Omega}_{21}\boldsymbol{\Omega}_{11}^{-1}\ell_1). \quad (28)$$

Using the Schur complement, we have

$$\mathbf{P}_{n-1}\boldsymbol{\Omega}_{21}\boldsymbol{\Omega}_{11}^{-1} = \mathbf{F}_{22}^{-1}\boldsymbol{\Omega}_{21}\boldsymbol{\Omega}_{11}^{-1} \quad (29)$$

$$= \boldsymbol{\Omega}_{22}^{-1}\boldsymbol{\Omega}_{21}\mathbf{F}_{11}^{-1} \quad (30)$$

$$= \cancel{\langle \mathbf{Q}^{-1} \rangle^{-1}}\cancel{\langle \mathbf{Q}^{-1} \rangle}\langle \mathbf{A} \rangle \mathbf{G}_{n-1} \mathbf{V}_{n-1} \quad (31)$$

$$= \langle \mathbf{A} \rangle \mathbf{G}_{n-1} \mathbf{V}_{n-1}. \quad (32)$$

Substituting this term back into the expression for the mean and simplifying gives the following.

$$\mathbf{m}_{n-1} = \mathbf{P}_{n-1}(\ell_2 + \boldsymbol{\Omega}_{21}\boldsymbol{\Omega}_{11}^{-1}\ell_1) \quad (33)$$

$$= \mathbf{P}_{n-1}\ell_2 + \mathbf{P}_{n-1}\boldsymbol{\Omega}_{21}\boldsymbol{\Omega}_{11}^{-1}\ell_1 \quad (34)$$

$$= (\langle \mathbf{Q} \rangle + \langle \mathbf{A} \rangle \mathbf{G}_{n-1} \mathbf{V}_{n-1} \langle \mathbf{A}^\top \rangle) \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle \mathbf{u}_n \quad (35)$$

$$+ \langle \mathbf{A} \rangle \mathbf{G}_{n-1} \mathbf{V}_{n-1} (\mathbf{V}_{n-1}^{-1}\boldsymbol{\mu}_{n-1} - (\langle \mathbf{A}^\top \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle + \boldsymbol{\Sigma}_{AQB}) \mathbf{u}_n) \quad (36)$$

$$= \cancel{\langle \mathbf{Q} \rangle}\cancel{\langle \mathbf{Q}^{-1} \rangle}\langle \mathbf{B} \rangle \mathbf{u}_n + \langle \mathbf{A} \rangle \mathbf{G}_{n-1} \mathbf{V}_{n-1} \langle \mathbf{A}^\top \rangle \cancel{\langle \mathbf{Q}^{-1} \rangle} \langle \mathbf{B} \rangle \mathbf{u}_n \quad (37)$$

$$+ \langle \mathbf{A} \rangle \mathbf{G}_{n-1} \mathbf{V}_{n-1} \cancel{\mathbf{V}_{n-1}^{-1}\boldsymbol{\mu}_{n-1}} - \langle \mathbf{A} \rangle \mathbf{G}_{n-1} \mathbf{V}_{n-1} \langle \mathbf{A}^\top \rangle \cancel{\langle \mathbf{Q}^{-1} \rangle} \langle \mathbf{B} \rangle \mathbf{u}_n - \langle \mathbf{A} \rangle \mathbf{G}_{n-1} \mathbf{V}_{n-1} \boldsymbol{\Sigma}_{AQB} \mathbf{u}_n \quad (38)$$

$$= \langle \mathbf{B} \rangle \mathbf{u}_n + \langle \mathbf{A} \rangle \mathbf{G}_{n-1}\boldsymbol{\mu}_{n-1} - \langle \mathbf{A} \rangle \mathbf{G}_{n-1} \mathbf{V}_{n-1} \boldsymbol{\Sigma}_{AQB} \mathbf{u}_n \quad (39)$$

In summary, the marginalizing over \mathbf{x}_{n-1} gives the predictive distribution

$$q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n-1}) = \mathcal{N}(\mathbf{x}_n | \mathbf{m}_{n-1}, \mathbf{P}_{n-1}) \quad (40)$$

where

$$\mathbf{G}_{n-1} = \mathbf{I} - \mathbf{V}_{n-1}(\mathbf{I} + \boldsymbol{\Sigma}_{AQA}\mathbf{V}_{n-1})^{-1}\boldsymbol{\Sigma}_{AQA} \quad (41)$$

$$\mathbf{m}_{n-1} = \langle \mathbf{A} \rangle \mathbf{G}_{n-1} (\boldsymbol{\mu}_{n-1} - \mathbf{V}_{n-1} \boldsymbol{\Sigma}_{AQB} \mathbf{u}_n) + \langle \mathbf{B} \rangle \mathbf{u}_n \quad (42)$$

$$\mathbf{P}_{n-1} = \langle \mathbf{A} \rangle \mathbf{G}_{n-1} \mathbf{V}_{n-1} \langle \mathbf{A}^\top \rangle + \langle \mathbf{Q} \rangle. \quad (43)$$

Initially at time $n = 1$, the predictive distribution is simply equal to the prior over the initial state $p(\mathbf{x}_1) = \mathcal{N}(\mathbf{x}_1 | \mathbf{m}_0, \mathbf{P}_0)$.

A.2. Marginal likelihood

Next, we derive the normalization term for the marginal posterior probability. The distribution is found by taking the integral over (marginalizing out) \mathbf{x}_n , resulting in a Gaussian distribution with mean denoted by $\hat{\mathbf{y}}_n$ and covariance denoted by \mathbf{S}_n ,

$$q_{\mathbf{x}}(\mathbf{y}_n | \mathbf{y}_{1:n-1}) = \int p(\mathbf{y}_n | \mathbf{x}_n) q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n-1}) d\mathbf{x}_n \quad (44)$$

$$= \int \mathcal{N}(\mathbf{y}_n | \mathbf{C}\mathbf{x}_n + \mathbf{D}\mathbf{u}_n, \mathbf{R}) \mathcal{N}(\mathbf{x}_n | \mathbf{m}_{n-1}, \mathbf{P}_{n-1}) d\mathbf{x}_n \quad (45)$$

$$= \mathcal{N}(\mathbf{y} | \hat{\mathbf{y}}_n, \mathbf{S}_n). \quad (46)$$

Collecting the quadratic terms over \mathbf{x}_n and \mathbf{y}_n :

$$\Omega_{11} = \mathbf{P}_{n-1}^{-1} + \langle \mathbf{C}^T \mathbf{R}^{-1} \mathbf{C} \rangle \quad (47)$$

$$= \mathbf{P}_{n-1}^{-1} + \Sigma_{CRC} + \langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{C} \rangle \quad (48)$$

$$\Omega_{12} = \langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \quad (49)$$

$$\Omega_{21} = \langle \mathbf{R}^{-1} \rangle \langle \mathbf{C} \rangle \quad (50)$$

$$\Omega_{22} = \langle \mathbf{R}^{-1} \rangle. \quad (51)$$

Using the Schur complement, we have

$$\mathbf{S}_n = (\Omega_{22} - \Omega_{21} \Omega_{11}^{-1} \Omega_{12})^{-1} \quad (52)$$

$$= \Omega_{22}^{-1} + \Omega_{22}^{-1} \Omega_{21} \mathbf{F}_{11}^{-1} \Omega_{12} \Omega_{22}^{-1} \quad (53)$$

where

$$\mathbf{F}_{11}^{-1} = (\Omega_{11} - \Omega_{12} \Omega_{22}^{-1} \Omega_{21})^{-1} \quad (54)$$

$$= (\mathbf{P}_{n-1}^{-1} + \Sigma_{CRC} + \langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{C} \rangle - \langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{C} \rangle)^{-1} \quad (55)$$

$$= (\mathbf{P}_{n-1}^{-1} + \Sigma_{CRC} + \cancel{\langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{C} \rangle} - \cancel{\langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{C} \rangle})^{-1} \quad (56)$$

$$= (\mathbf{P}_{n-1}^{-1} + \Sigma_{CRC})^{-1}. \quad (57)$$

The Woodbury matrix identity is used to get the following preferred form for \mathbf{F}_{11}^{-1} ,

$$\mathbf{F}_{11}^{-1} = (\mathbf{P}_{n-1}^{-1} + \Sigma_{CRC})^{-1} \quad (58)$$

$$= \mathbf{P}_{n-1} - \mathbf{P}_{n-1} (\mathbf{I} + \Sigma_{CRC} \mathbf{P}_{n-1})^{-1} \Sigma_{CRC} \mathbf{P}_{n-1} \quad (59)$$

$$= (\mathbf{I} - \mathbf{P}_{n-1} (\mathbf{I} + \Sigma_{CRC} \mathbf{P}_{n-1})^{-1} \Sigma_{CRC}) \mathbf{P}_{n-1} \quad (60)$$

$$= \mathbf{L}_{n-1} \mathbf{P}_{n-1} \quad (61)$$

where we have defined

$$\mathbf{L}_{n-1} = \mathbf{I} - \mathbf{P}_{n-1} (\mathbf{I} + \Sigma_{CRC} \mathbf{P}_{n-1})^{-1} \Sigma_{CRC}. \quad (62)$$

Substituting in and simplifying the expression for the covariance \mathbf{S}_n ,

$$\mathbf{S}_n = \Omega_{22}^{-1} + \Omega_{22}^{-1} \Omega_{21} \mathbf{F}_{11}^{-1} \Omega_{12} \Omega_{22}^{-1} \quad (63)$$

$$= \langle \mathbf{R}^{-1} \rangle^{-1} + \cancel{\langle \mathbf{R}^{-1} \rangle^{-1} \langle \mathbf{R}^{-1} \rangle \langle \mathbf{C} \rangle} \mathbf{L}_{n-1} \mathbf{P}_{n-1} \langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{C} \rangle \quad (64)$$

$$= \langle \mathbf{R} \rangle + \langle \mathbf{C} \rangle \mathbf{L}_{n-1} \mathbf{P}_{n-1} \langle \mathbf{C}^T \rangle. \quad (65)$$

Next, we derive the expression for the mean $\hat{\mathbf{y}}_n$. Completing the square for \mathbf{y}_n and \mathbf{x}_n , the linear terms are factored as

$$\ell_1 = \mathbf{P}_{n-1}^{-1} \mathbf{m}_{n-1} - \langle \mathbf{C}^T \mathbf{R}^{-1} \mathbf{D} \rangle \mathbf{u}_n \quad (66)$$

$$= \mathbf{P}_{n-1}^{-1} \mathbf{m}_{n-1} - (\langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{D} \rangle + \Sigma_{CRD}) \mathbf{u}_n \quad (67)$$

$$\ell_2 = \langle \mathbf{R}^{-1} \rangle \langle \mathbf{D} \rangle \mathbf{u}_n. \quad (68)$$

The mean is then given by

$$\hat{\mathbf{y}}_n = \mathbf{S}_n(\ell_2 - \Omega_{21}\Omega_{11}^{-1}\ell_1). \quad (69)$$

Using the Schur complement, we get

$$\mathbf{S}_n\Omega_{21}\Omega_{11}^{-1} = \mathbf{F}_{22}^{-1}\Omega_{21}\Omega_{11}^{-1} \quad (70)$$

$$= \Omega_{22}^{-1}\Omega_{21}\mathbf{F}_{11}^{-1} \quad (71)$$

$$= \cancel{\langle \mathbf{R}^{-1} \rangle^{-1}} \cancel{\langle \mathbf{R}^{-1} \rangle} \langle \mathbf{C} \rangle \mathbf{L}_{n-1} \mathbf{P}_{n-1} \quad (72)$$

$$= \langle \mathbf{C} \rangle \mathbf{L}_{n-1} \mathbf{P}_{n-1}. \quad (73)$$

Substituting this term back into the expression for $\hat{\mathbf{y}}_n$ and simplifying, we get

$$\hat{\mathbf{y}}_n = \mathbf{S}_n(\ell_2 + \Omega_{21}\Omega_{11}^{-1}\ell_1) \quad (74)$$

$$= \mathbf{S}_n\ell_2 + \mathbf{S}_n\Omega_{21}\Omega_{11}^{-1}\ell_1 \quad (75)$$

$$= (\langle \mathbf{R} \rangle + \langle \mathbf{C} \rangle \mathbf{L}_{n-1} \mathbf{P}_{n-1} \langle \mathbf{C}^T \rangle) \langle \mathbf{R}^{-1} \rangle \langle \mathbf{D} \rangle \mathbf{u}_n \quad (76)$$

$$+ \langle \mathbf{C} \rangle \mathbf{L}_{n-1} \mathbf{P}_{n-1} (\mathbf{P}_{n-1}^{-1} \mathbf{m}_{n-1} - (\langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{D} \rangle + \Sigma_{CRD}) \mathbf{u}_n) \quad (77)$$

$$= \langle \mathbf{D} \rangle \mathbf{u}_n + \langle \mathbf{C} \rangle \mathbf{L}_{n-1} \mathbf{m}_{n-1} - \langle \mathbf{C} \rangle \mathbf{L}_{n-1} \mathbf{P}_{n-1} \Sigma_{CRD} \mathbf{u}_n. \quad (78)$$

In summary, the marginalization over \mathbf{x}_n gives the marginal likelihood (filtered output probability)

$$q_{\mathbf{x}}(\mathbf{y}_n | \mathbf{y}_{1:n-1}) = \mathcal{N}(\mathbf{y} | \hat{\mathbf{y}}_n, \mathbf{S}_n) \quad (79)$$

where

$$\mathbf{L}_{n-1} = \mathbf{I} - \mathbf{P}_{n-1}(\mathbf{I} + \Sigma_{CRC} \mathbf{P}_{n-1})^{-1} \Sigma_{CRC} \quad (80)$$

$$\hat{\mathbf{y}}_n = \langle \mathbf{C} \rangle \mathbf{L}_{n-1} (\mathbf{m}_{n-1} - \mathbf{P}_{n-1} \Sigma_{CRD} \mathbf{u}_n) + \langle \mathbf{D} \rangle \mathbf{u}_n \quad (81)$$

$$\mathbf{S}_n = \langle \mathbf{R} \rangle + \langle \mathbf{C} \rangle \mathbf{L}_{n-1} \mathbf{P}_{n-1} \langle \mathbf{C}^T \rangle. \quad (82)$$

A.3. Marginal posterior

Returning to the marginal posterior distribution, we substitute in the previously derived results and get

$$q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n}) = \frac{p(\mathbf{y}_n | \mathbf{x}_n) q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n-1})}{q_{\mathbf{x}}(\mathbf{y}_n | \mathbf{y}_{1:n-1})} \quad (83)$$

$$= \frac{\mathcal{N}(\mathbf{y}_n | \mathbf{C}\mathbf{x}_n + \mathbf{D}\mathbf{u}_n, \mathbf{R}) \mathcal{N}(\mathbf{x}_n | \mathbf{m}_{n-1}, \mathbf{P}_{n-1})}{\mathcal{N}(\mathbf{y} | \hat{\mathbf{y}}_n, \mathbf{S}_n)} \quad (84)$$

$$= \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, \mathbf{V}_n). \quad (85)$$

The covariance \mathbf{V}_n is found by completing the square for \mathbf{x}_n ,

$$\mathbf{V}_n = (\mathbf{P}_{n-1}^{-1} + \langle \mathbf{C}^T \mathbf{R}^{-1} \mathbf{C} \rangle)^{-1} \quad (86)$$

$$= (\mathbf{P}_{n-1}^{-1} + \Sigma_{CRC} + \langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{C} \rangle)^{-1}. \quad (87)$$

The Woodbury matrix identity can again be applied to get this result into the preferred form. Let

$$\Omega_{11} = -\langle \mathbf{R} \rangle \quad (88)$$

$$\Omega_{12} = \langle \mathbf{C} \rangle \quad (89)$$

$$\Omega_{21} = \langle \mathbf{C}^T \rangle \quad (90)$$

$$\Omega_{22} = \mathbf{P}_{n-1}^{-1} + \Sigma_{CRC} \quad (91)$$

then we re-write the expression for \mathbf{V}_n and use the identity

$$\mathbf{V}_n = (\Omega_{22} - \Omega_{21}\Omega_{11}^{-1}\Omega_{12})^{-1} \quad (92)$$

$$= \Omega_{22}^{-1} + \Omega_{22}^{-1}\Omega_{21}\mathbf{F}_{11}^{-1}\Omega_{12}\Omega_{22}^{-1} \quad (93)$$

where

$$\mathbf{F}_{11}^{-1} = (\Omega_{11} - \Omega_{12}\Omega_{22}^{-1}\Omega_{21})^{-1} \quad (94)$$

$$= (-\langle \mathbf{R} \rangle - \langle \mathbf{C} \rangle (\mathbf{P}_{n-1}^{-1} + \Sigma_{CRC})^{-1}ev\mathbf{C}^T)^{-1} \quad (95)$$

$$= -(\langle \mathbf{R} \rangle + \langle \mathbf{C} \rangle \mathbf{L}_{n-1}\mathbf{P}_{n-1}\langle \mathbf{C}^T \rangle)^{-1} \quad (96)$$

$$= -\mathbf{S}_n^{-1}. \quad (97)$$

Substituting back into the expression for \mathbf{V}_n gives

$$\mathbf{V}_n = \Omega_{22}^{-1} + \Omega_{22}^{-1}\Omega_{21}\mathbf{F}_{11}^{-1}\Omega_{12}\Omega_{22}^{-1} \quad (98)$$

$$= (\mathbf{P}_{n-1}^{-1} + \Sigma_{CRC})^{-1} - (\mathbf{P}_{n-1}^{-1} + \Sigma_{CRC})^{-1}\langle \mathbf{C}^T \rangle \mathbf{S}_n^{-1}\langle \mathbf{C} \rangle (\mathbf{P}_{n-1}^{-1} + \Sigma_{CRC})^{-1} \quad (99)$$

$$= \mathbf{L}_{n-1}\mathbf{P}_{n-1} - \mathbf{L}_{n-1}\mathbf{P}_{n-1}\langle \mathbf{C}^T \rangle \mathbf{S}_n^{-1}\langle \mathbf{C} \rangle \mathbf{L}_{n-1}\mathbf{P}_{n-1} \quad (100)$$

$$= (\mathbf{I} - \mathbf{L}_{n-1}\mathbf{P}_{n-1}\langle \mathbf{C}^T \rangle \mathbf{S}_n^{-1}\langle \mathbf{C} \rangle)\mathbf{L}_{n-1}\mathbf{P}_{n-1} \quad (101)$$

$$= (\mathbf{I} - \mathbf{K}_n\langle \mathbf{C} \rangle)\mathbf{L}_{n-1}\mathbf{P}_{n-1} \quad (102)$$

where we have defined

$$\mathbf{K}_n = \mathbf{L}_{n-1}\mathbf{P}_{n-1}\langle \mathbf{C}^T \rangle \mathbf{S}_n^{-1}. \quad (103)$$

In the non-Bayesian scenario (when the parameters are deterministic), \mathbf{K}_n is called the *Kalman gain*. Lastly, the mean is found by completing the square over \mathbf{x}_n . Collecting the linear terms over \mathbf{x}_n gives

$$\boldsymbol{\mu}_n = \mathbf{V}_n(\mathbf{P}_{n-1}^{-1}\mathbf{m}_{n-1} - (\langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{D} \rangle + \Sigma_{CRD})\mathbf{u}_n + \langle \mathbf{C}^T \mathbf{R}^{-1} \rangle \mathbf{y}_n). \quad (104)$$

Each term in the previous equation is addressed in order. First, we have

$$\mathbf{V}_n\mathbf{P}_{n-1}^{-1}\mathbf{m}_{n-1} = (\mathbf{I} - \mathbf{K}_n\langle \mathbf{C} \rangle)\mathbf{L}_{n-1}\mathbf{P}_{n-1}\mathbf{P}_{n-1}^{-1}\mathbf{m}_{n-1} \quad (105)$$

$$= (\mathbf{I} - \mathbf{K}_n\langle \mathbf{C} \rangle)\mathbf{L}_{n-1}\mathbf{m}_{n-1} \quad (106)$$

$$= \mathbf{L}_{n-1}\mathbf{m}_{n-1} - \mathbf{K}_n\langle \mathbf{C} \rangle \mathbf{L}_{n-1}\mathbf{m}_{n-1}. \quad (107)$$

Next, $\mathbf{V}_n\langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle$ can be simplified. Recall from the conditional distribution that $\Omega_{12} = \langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle$, and $\mathbf{V}_n = \Omega_{11}^{-1} = \mathbf{P}_{n-1}^{-1} + \langle \mathbf{C}^T \mathbf{R}^{-1} \mathbf{C} \rangle$. Then,

$$\mathbf{V}_n\langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle = \Omega_{11}^{-1}\Omega_{12} \quad (108)$$

$$= \mathbf{F}_{11}^{-1}\Omega_{12}\Omega_{22}^{-1}\mathbf{F}_{22} \quad (109)$$

$$= \mathbf{L}\mathbf{P}_{n-1}\langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{R}^{-1} \rangle^{-1}\mathbf{S}_n^{-1} \quad (110)$$

$$= \mathbf{L}\mathbf{P}_{n-1}\langle \mathbf{C}^T \rangle \mathbf{S}_n^{-1} \quad (111)$$

$$= \mathbf{K}_n. \quad (112)$$

Substituting this term back into the expression for $\boldsymbol{\mu}_n$ and simplifying, we get

$$\boldsymbol{\mu}_n = \mathbf{V}_n(\mathbf{P}_{n-1}^{-1}\mathbf{m}_{n-1} - (\langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{D} \rangle + \Sigma_{CRD})\mathbf{u}_n + \langle \mathbf{C}^T \mathbf{R}^{-1} \rangle \mathbf{y}_n) \quad (113)$$

$$= \mathbf{V}_n\mathbf{P}_{n-1}^{-1}\mathbf{m}_{n-1} - \mathbf{V}_n(\langle \mathbf{C}^T \rangle \langle \mathbf{R}^{-1} \rangle \langle \mathbf{D} \rangle + \Sigma_{CRD})\mathbf{u}_n + \mathbf{V}_n\langle \mathbf{C}^T \mathbf{R}^{-1} \rangle \mathbf{y}_n \quad (114)$$

$$= \mathbf{L}_{n-1}\mathbf{m}_{n-1} - \mathbf{K}_n\langle \mathbf{C} \rangle \mathbf{L}_{n-1}\mathbf{m}_{n-1} - \mathbf{K}_n\langle \mathbf{D} \rangle \mathbf{u}_n - \mathbf{V}_n\Sigma_{CRD}\mathbf{u}_n + \mathbf{K}_n\mathbf{y}_n \quad (115)$$

$$= \mathbf{L}_{n-1}\mathbf{m}_{n-1} + \mathbf{K}_n(\mathbf{y}_n - \langle \mathbf{C} \rangle \mathbf{L}_{n-1}\mathbf{m}_{n-1} - \langle \mathbf{D} \rangle \mathbf{u}_n) - \mathbf{V}_n\Sigma_{CRD}\mathbf{u}_n \quad (116)$$

$$= \mathbf{L}_{n-1}\mathbf{m}_{n-1} + \mathbf{K}_n(\mathbf{y}_n - \hat{\mathbf{y}}_n) - \mathbf{L}_{n-1}\mathbf{P}_{n-1}\Sigma_{CRD}\mathbf{u}_n. \quad (117)$$

In summary, the forward pass computes the statistics of the marginal posterior probability, $\forall n \in [1..N]$:

$$q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n}) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, \mathbf{V}_n) \quad (118)$$

where

$$\mathbf{K}_n = \mathbf{L}_{n-1} \mathbf{P}_{n-1} \langle \mathbf{C}^T \rangle \mathbf{S}_n^{-1} \quad (119)$$

$$\boldsymbol{\mu}_n = \mathbf{L}_{n-1} (\mathbf{m}_{n-1} - \mathbf{P}_{n-1} \boldsymbol{\Sigma}_{CRD} \mathbf{u}_n) + \mathbf{K}_n (\mathbf{y}_n - \hat{\mathbf{y}}_n) \quad (120)$$

$$\mathbf{V}_n = (\mathbf{I} - \mathbf{K}_n \langle \mathbf{C} \rangle) \mathbf{L}_{n-1} \mathbf{P}_{n-1}. \quad (121)$$

This concludes the derivation of the forward pass.

B. BACKWARD PASS (SMOOTHING) DERIVATION

The backward pass computes the sufficient statistics of the marginal [1]

$$q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{Y}) = \frac{q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n})}{q_{\mathbf{x}}(\mathbf{y}_{n+1} | \mathbf{y}_{1:n})} \int \frac{p(\mathbf{y}_{n+1} | \mathbf{x}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{x}_n) q_{\mathbf{x}}(\mathbf{x}_{n+1} | \mathbf{Y})}{q_{\mathbf{x}}(\mathbf{x}_{n+1} | \mathbf{y}_{1:n+1})} d\mathbf{x}_{n+1} \quad (122)$$

$$= \frac{q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n})}{q_{\mathbf{x}}(\mathbf{y}_{n+1} | \mathbf{y}_{1:n})} \int \frac{\underline{p(\mathbf{y}_{n+1} | \mathbf{x}_{n+1})} p(\mathbf{x}_{n+1} | \mathbf{x}_n) p(\mathbf{x}_{n+1} | \mathbf{Y}) q_{\mathbf{x}}(\mathbf{y}_{n+1} | \mathbf{y}_{1:n})}{\underline{p(\mathbf{y}_{n+1} | \mathbf{x}_{n+1})} \int p(\mathbf{x}_{n+1} | \mathbf{x}_n) q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n}) d\mathbf{x}_n} d\mathbf{x}_{n+1} \quad (123)$$

$$= q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{y}_{1:n}) \int \frac{p(\mathbf{x}_{n+1} | \mathbf{x}_n) p(\mathbf{x}_{n+1} | \mathbf{Y})}{q_{\mathbf{x}}(\mathbf{x}_{n+1} | \mathbf{y}_{1:n})} d\mathbf{x}_{n+1} \quad (124)$$

$$= \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, \mathbf{V}_n) \int \frac{\mathcal{N}(\mathbf{x}_{n+1} | \mathbf{A}\mathbf{x}_n + \mathbf{B}\mathbf{u}_{n+1}, \mathbf{Q}) \mathcal{N}(\mathbf{x}_{n+1} | \hat{\boldsymbol{\mu}}_{n+1}, \hat{\mathbf{V}}_{n+1})}{\mathcal{N}(\mathbf{x}_{n+1} | \mathbf{m}_n, \mathbf{P}_n)} d\mathbf{x}_{n+1} \quad (125)$$

$$= \mathcal{N}(\mathbf{x}_n | \hat{\boldsymbol{\mu}}_n, \hat{\mathbf{V}}_n). \quad (126)$$

The square is completed over \mathbf{x}_n and \mathbf{x}_{n-1} to factor the quadratic terms

$$\Omega_{22} = \mathbf{V}_n^{-1} + \langle \mathbf{A}^T \mathbf{Q}^{-1} \mathbf{A} \rangle \quad (127)$$

$$= \mathbf{V}_n^{-1} + \Sigma_{AQ} \mathbf{A} + \langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle \quad (128)$$

$$\Omega_{21} = \langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \quad (129)$$

$$\Omega_{12} = \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle \quad (130)$$

$$\Omega_{11} = \hat{\mathbf{V}}_{n+1}^{-1} + \langle \mathbf{Q}^{-1} \rangle - \mathbf{P}_n^{-1} \quad (131)$$

$$= \hat{\mathbf{V}}_{n+1}^{-1} + \cancel{\langle \mathbf{Q}^{-1} \rangle} - (\cancel{\langle \mathbf{Q}^{-1} \rangle} - \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle) (\mathbf{V}_n^{-1} + \langle \mathbf{A}^T \mathbf{Q}^{-1} \mathbf{A} \rangle)^{-1} \langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \quad (132)$$

$$= \hat{\mathbf{V}}_{n+1}^{-1} + \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle \left(\mathbf{V}_n^{-1} + \Sigma_{AQ} \mathbf{A} + \langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle \right)^{-1} \langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \quad (133)$$

and the linear terms

$$\ell_2 = \mathbf{V}_n^{-1} \boldsymbol{\mu}_n - \langle \mathbf{A}^T \mathbf{Q}^{-1} \mathbf{B} \rangle \mathbf{u}_{n+1} \quad (134)$$

$$= \mathbf{V}_n^{-1} \boldsymbol{\mu}_n - (\langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle + \Sigma_{AQ} \mathbf{B}) \mathbf{u}_{n+1} \quad (135)$$

$$\ell_1 = \hat{\mathbf{V}}_{n+1}^{-1} \hat{\boldsymbol{\mu}}_{n+1} - \mathbf{P}_n^{-1} \mathbf{m}_n + \langle \mathbf{Q}^{-1} \mathbf{B} \rangle \mathbf{u}_{n+1}. \quad (136)$$

The covariance and mean are

$$\hat{\mathbf{V}}_n = (\Omega_{22} - \Omega_{21} \Omega_{11}^{-1} \Omega_{12})^{-1} \quad (137)$$

$$\hat{\boldsymbol{\mu}}_n = \hat{\mathbf{V}}_n (\ell_2 - \Omega_{21} \Omega_{11}^{-1} \ell_1). \quad (138)$$

We will use results from matrix theory to get these results into a preferred form, reminiscent of the Rauch-Tung-Striebel (RTS) smoother [2]. First the covariance. Using the Schur complement, we get

$$\hat{\mathbf{V}}_n = (\Omega_{22} - \Omega_{21} \Omega_{11}^{-1} \Omega_{12})^{-1} \quad (139)$$

$$= \Omega_{22}^{-1} + \Omega_{22}^{-1} \Omega_{21} \mathbf{F}_{11}^{-1} \Omega_{12} \Omega_{22}^{-1}. \quad (140)$$

The term \mathbf{F}_{11}^{-1} is reduced to

$$\mathbf{F}_{11}^{-1} = (\Omega_{11} - \Omega_{12} \Omega_{22}^{-1} \Omega_{21})^{-1} \quad (141)$$

$$= \left(\hat{\mathbf{V}}_{n+1}^{-1} + \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle \left(\mathbf{V}_n^{-1} + \Sigma_{AQ} \mathbf{A} + \langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle \right)^{-1} \langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \right. \quad (142)$$

$$\left. - \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle \left(\mathbf{V}_n^{-1} + \Sigma_{AQ} \mathbf{A} + \langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle \right)^{-1} \langle \mathbf{A}^T \rangle \langle \mathbf{Q}^{-1} \rangle \right)^{-1} \quad (143)$$

$$= (\hat{\mathbf{V}}_{n+1}^{-1})^{-1} \quad (144)$$

$$= \hat{\mathbf{V}}_{n+1}. \quad (145)$$

Next, the inverse of Ω_{22} is analytically expressed using the Woodbury matrix identity. Defining

$$\Lambda_{11} = -\langle \mathbf{Q} \rangle \quad (146)$$

$$\Lambda_{12} = \langle \mathbf{A} \rangle \quad (147)$$

$$\Lambda_{21} = \langle \mathbf{A}^\top \rangle \quad (148)$$

$$\Lambda_{22} = \mathbf{V}_n^{-1} + \Sigma_{AQ} \quad (149)$$

the inverse is then

$$\Omega_{22}^{-1} = (\Lambda_{22} - \Lambda_{21}\Lambda_{11}^{-1}\Lambda_{12})^{-1} \quad (150)$$

$$= \Lambda_{22}^{-1} + \Lambda_{22}^{-1}\Lambda_{21}\mathbf{O}_{11}^{-1}\Lambda_{12}\Lambda_{22}^{-1}. \quad (151)$$

The term Λ_{22}^{-1} is simplified using equation (20) from the forward pass:

$$\Lambda_{22}^{-1} = (\mathbf{V}_n^{-1} + \Sigma_{AQ})^{-1} \quad (152)$$

$$= \mathbf{G}_n \mathbf{V}_n. \quad (153)$$

Likewise, the term \mathbf{O}_{11}^{-1} is given by

$$\mathbf{O}_{11}^{-1} = (\Lambda_{11} - \Lambda_{12}\Lambda_{22}^{-1}\Lambda_{21})^{-1} \quad (154)$$

$$= (-\langle \mathbf{Q} \rangle - \langle \mathbf{A} \rangle \mathbf{G}_n \mathbf{V}_n \langle \mathbf{A}^\top \rangle)^{-1} \quad (155)$$

$$= -\mathbf{P}_n^{-1}. \quad (156)$$

These results are substituted into the equation for Ω_{22}^{-1} :

$$\Omega_{22}^{-1} = \Lambda_{22}^{-1} + \Lambda_{22}^{-1}\Lambda_{21}\mathbf{O}_{11}^{-1}\Lambda_{12}\Lambda_{22}^{-1} \quad (157)$$

$$= \mathbf{G}_n \mathbf{V}_n - \mathbf{G}_n \mathbf{V}_n \langle \mathbf{A}^\top \rangle \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \mathbf{V}_n. \quad (158)$$

Lastly, we can use another matrix identity to simplify the expression $\Omega_{22}^{-1}\Omega_{21}$.

$$\Omega_{22}^{-1}\Omega_{21} = \left(\mathbf{V}_n^{-1} + \Sigma_{AQ} + \langle \mathbf{A}^\top \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle \right)^{-1} \langle \mathbf{A}^\top \rangle \langle \mathbf{Q}^{-1} \rangle \quad (159)$$

$$= (\mathbf{V}_n^{-1} + \Sigma_{AQ})^{-1} \langle \mathbf{A}^\top \rangle \left(\langle \mathbf{Q} \rangle + \langle \mathbf{A} \rangle (\mathbf{V}_n^{-1} + \Sigma_{AQ})^{-1} \langle \mathbf{A}^\top \rangle \right)^{-1} \quad (160)$$

$$= \mathbf{G}_n \mathbf{V}_n \langle \mathbf{A}^\top \rangle \left(\langle \mathbf{Q} \rangle + \langle \mathbf{A} \rangle \mathbf{G}_n \mathbf{V}_n \langle \mathbf{A}^\top \rangle \right)^{-1} \quad (161)$$

$$= \mathbf{G}_n \mathbf{V}_n \langle \mathbf{A}^\top \rangle \mathbf{P}_n^{-1} \quad (162)$$

Putting this altogether, and noting that $\Omega_{12}\Omega_{22}^{-1} = (\Omega_{22}^{-1}\Omega_{21})^\top$, the covariance from the backward pass at time n is given by

$$\hat{\mathbf{V}}_n = \Omega_{22}^{-1} + \Omega_{22}^{-1}\Omega_{21}\mathbf{F}_{11}^{-1}\Omega_{12}\Omega_{22}^{-1} \quad (163)$$

$$= \mathbf{G}_n \mathbf{V}_n - \mathbf{G}_n \mathbf{V}_n \langle \mathbf{A}^\top \rangle \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \mathbf{V}_n + \mathbf{G}_n \mathbf{V}_n \langle \mathbf{A}^\top \rangle \mathbf{P}_n^{-1} \hat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \mathbf{V}_n \quad (164)$$

$$= \mathbf{G}_n \mathbf{V}_n + \mathbf{G}_n \mathbf{V}_n \langle \mathbf{A}^\top \rangle \mathbf{P}_n^{-1} (\hat{\mathbf{V}}_{n+1} - \mathbf{P}_n) \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \mathbf{V}_n \quad (165)$$

$$= \mathbf{G}_n \mathbf{V}_n + \mathbf{J}_n (\hat{\mathbf{V}}_{n+1} - \mathbf{P}_n) \mathbf{J}_n^\top \quad (166)$$

where we have defined

$$\mathbf{J}_n = \mathbf{G}_n \mathbf{V}_n \langle \mathbf{A}^\top \rangle \mathbf{P}_n^{-1}. \quad (167)$$

Lastly, the mean $\hat{\boldsymbol{\mu}}_n$ is given by

$$\hat{\boldsymbol{\mu}}_n = \hat{\mathbf{V}}_n (\boldsymbol{\ell}_2 + \Omega_{21}\Omega_{11}^{-1}\boldsymbol{\ell}_1) \quad (168)$$

$$= \hat{\mathbf{V}}_n \mathbf{V}_n^{-1} \boldsymbol{\mu}_n - \hat{\mathbf{V}}_n (\langle \mathbf{A}^\top \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle + \Sigma_{AQ} \mathbf{B}) \mathbf{u}_{n+1} + \hat{\mathbf{V}}_n \Omega_{21}\Omega_{11}^{-1}\boldsymbol{\ell}_1 \quad (169)$$

$$= \hat{\mathbf{V}}_n \mathbf{V}_n^{-1} \boldsymbol{\mu}_n - \hat{\mathbf{V}}_n \langle \mathbf{A}^\top \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle \mathbf{u}_{n+1} - \hat{\mathbf{V}}_n \Sigma_{AQ} \mathbf{B} \mathbf{u}_{n+1} + \hat{\mathbf{V}}_n \Omega_{21}\Omega_{11}^{-1}\boldsymbol{\ell}_1. \quad (170)$$

The left-most term, $\widehat{\mathbf{V}}_n \mathbf{V}_n^{-1} \boldsymbol{\mu}_n$, is

$$\widehat{\mathbf{V}}_n \mathbf{V}_n^{-1} \boldsymbol{\mu}_n = \left(\mathbf{G}_n \mathbf{V}_n + \mathbf{J}_n (\widehat{\mathbf{V}}_{n+1} - \mathbf{P}_n) \mathbf{J}_n^\top \right) \mathbf{V}_n^{-1} \boldsymbol{\mu}_n \quad (171)$$

$$= \left(\mathbf{G}_n \mathbf{V}_n + \mathbf{J}_n (\widehat{\mathbf{V}}_{n+1} - \mathbf{P}_n) \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \mathbf{V}_n \right) \widehat{\mathbf{V}}_n \boldsymbol{\mu}_n \quad (172)$$

$$= \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n (\widehat{\mathbf{V}}_{n+1} - \mathbf{P}_n) \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n \quad (173)$$

$$= \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n - \mathbf{J}_n \mathbf{P}_n \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n \quad (174)$$

$$= \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n - \mathbf{J}_n \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n. \quad (175)$$

Next, $\widehat{\mathbf{V}}_n \boldsymbol{\Omega}_{21} \boldsymbol{\Omega}_{11}^{-1}$ can be evaluated using the Schur complement:

$$\widehat{\mathbf{V}}_n \boldsymbol{\Omega}_{21} \boldsymbol{\Omega}_{11}^{-1} = \mathbf{F}_{22}^{-1} \boldsymbol{\Omega}_{21} \boldsymbol{\Omega}_{11}^{-1} = \boldsymbol{\Omega}_{22}^{-1} \boldsymbol{\Omega}_{21} \mathbf{F}_{11}^{-1} \quad (176)$$

$$= \mathbf{G}_n \mathbf{V}_n \langle \mathbf{A}^\top \rangle \mathbf{P}_n^{-1} \widehat{\mathbf{V}}_{n+1} \quad (177)$$

$$= \mathbf{J}_n \widehat{\mathbf{V}}_{n+1}. \quad (178)$$

Likewise, we can use the Schur complement to evaluate $\widehat{\mathbf{V}}_n \langle \mathbf{A}^\top \rangle \langle \mathbf{Q}^{-1} \rangle$.

$$\widehat{\mathbf{V}}_n \langle \mathbf{A}^\top \rangle \langle \mathbf{Q}^{-1} \rangle = \mathbf{F}_{22}^{-1} \boldsymbol{\Omega}_{21} = \boldsymbol{\Omega}_{22}^{-1} \boldsymbol{\Omega}_{21} \mathbf{F}_{11}^{-1} \boldsymbol{\Omega}_{11} \quad (179)$$

$$= \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \boldsymbol{\Omega}_{11} \quad (180)$$

We expand the last term $\widehat{\mathbf{V}}_n \boldsymbol{\Omega}_{21} \boldsymbol{\Omega}_{11}^{-1} \boldsymbol{\ell}_1$ as follows.

$$\widehat{\mathbf{V}}_n \boldsymbol{\Omega}_{21} \boldsymbol{\Omega}_{11}^{-1} \boldsymbol{\ell}_1 = \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \left(\widehat{\mathbf{V}}_{n+1}^{-1} \widehat{\boldsymbol{\mu}}_{n+1} - \mathbf{P}_n^{-1} \mathbf{m}_n + \langle \mathbf{Q}^{-1} \mathbf{B} \rangle \mathbf{u}_{n+1} \right) \quad (181)$$

$$= \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \widehat{\mathbf{V}}_{n+1}^{-1} \widehat{\boldsymbol{\mu}}_{n+1} - \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \mathbf{m}_n + \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \langle \mathbf{Q}^{-1} \mathbf{B} \rangle \mathbf{u}_{n+1} \quad (182)$$

$$= \mathbf{J}_n \widehat{\boldsymbol{\mu}}_{n+1} - \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \mathbf{m}_n + \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \langle \mathbf{Q}^{-1} \mathbf{B} \rangle \mathbf{u}_{n+1} \quad (183)$$

The second term in the previous equation $\mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \mathbf{m}_n$ is further expanded to give

$$\mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \mathbf{m}_n = \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} (\langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n + (\langle \mathbf{B} \rangle - \langle \mathbf{A} \rangle \mathbf{G}_n \mathbf{V}_n \boldsymbol{\Sigma}_{AQ} \mathbf{B}) \mathbf{u}_{n+1}) \quad (184)$$

$$= \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \langle \mathbf{B} \rangle \mathbf{u}_{n+1} - \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \mathbf{V}_n \boldsymbol{\Sigma}_{AQ} \mathbf{B} \mathbf{u}_{n+1} \quad (185)$$

$$= \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \langle \mathbf{B} \rangle \mathbf{u}_{n+1} - \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{J}_n^\top \boldsymbol{\Sigma}_{AQ} \mathbf{B} \mathbf{u}_{n+1}. \quad (186)$$

Now we substitute all these terms back into equation (170) and get

$$\widehat{\boldsymbol{\mu}}_n = \widehat{\mathbf{V}}_n \mathbf{V}_n^{-1} \boldsymbol{\mu}_n - \widehat{\mathbf{V}}_n \langle \mathbf{A}^\top \rangle \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle \mathbf{u}_{n+1} - \widehat{\mathbf{V}}_n \boldsymbol{\Sigma}_{AQ} \mathbf{B} \mathbf{u}_{n+1} + \widehat{\mathbf{V}}_n \boldsymbol{\Omega}_{21} \boldsymbol{\Omega}_{11}^{-1} \boldsymbol{\ell}_1 \quad (187)$$

$$= \mathbf{G}_n \boldsymbol{\mu}_n + \underline{\mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n} - \mathbf{J}_n \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n - \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \boldsymbol{\Omega}_{11} \langle \mathbf{B} \rangle \mathbf{u}_{n+1} - \widehat{\mathbf{V}}_n \boldsymbol{\Sigma}_{AQ} \mathbf{B} \mathbf{u}_{n+1} + \mathbf{J}_n \widehat{\boldsymbol{\mu}}_{n+1} \quad (188)$$

$$- \underline{\mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n} - \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{P}_n^{-1} \langle \mathbf{B} \rangle \mathbf{u}_{n+1} + \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{J}_n^\top \boldsymbol{\Sigma}_{AQ} \mathbf{B} \mathbf{u}_{n+1} + \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \langle \mathbf{Q}^{-1} \mathbf{B} \rangle \mathbf{u}_{n+1} \quad (189)$$

$$= \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n (\widehat{\boldsymbol{\mu}}_{n+1} - \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n) - \widehat{\mathbf{V}}_n \boldsymbol{\Sigma}_{AQ} \mathbf{B} \mathbf{u}_{n+1} + \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{J}_n^\top \boldsymbol{\Sigma}_{AQ} \mathbf{B} \mathbf{u}_{n+1} \quad (190)$$

$$+ \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} (\langle \mathbf{Q}^{-1} \rangle - \mathbf{P}_n^{-1} - \boldsymbol{\Omega}_{11}) \langle \mathbf{B} \rangle \mathbf{u}_{n+1} \quad (191)$$

$$= \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n (\widehat{\boldsymbol{\mu}}_{n+1} - \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n) - \widehat{\mathbf{V}}_n \boldsymbol{\Sigma}_{AQ} \mathbf{B} \mathbf{u}_{n+1} + \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{J}_n^\top \boldsymbol{\Sigma}_{AQ} \mathbf{B} \mathbf{u}_{n+1} \quad (192)$$

$$+ \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} (-\widehat{\mathbf{V}}_{n+1}^\top) \langle \mathbf{B} \rangle \mathbf{u}_{n+1} \quad (193)$$

$$= \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n (\widehat{\boldsymbol{\mu}}_{n+1} - \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n) - \widehat{\mathbf{V}}_n \boldsymbol{\Sigma}_{AQ} \mathbf{B} \mathbf{u}_{n+1} + \mathbf{J}_n \widehat{\mathbf{V}}_{n+1} \mathbf{J}_n^\top \boldsymbol{\Sigma}_{AQ} \mathbf{B} \mathbf{u}_{n+1} - \mathbf{J}_n \langle \mathbf{B} \rangle \mathbf{u}_{n+1} \quad (194)$$

$$= \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n (\widehat{\boldsymbol{\mu}}_{n+1} - \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n - \langle \mathbf{B} \rangle \mathbf{u}_{n+1} + \widehat{\mathbf{V}}_{n+1} \mathbf{J}_n^\top \boldsymbol{\Sigma}_{AQ} \mathbf{B} \mathbf{u}_{n+1}) - \widehat{\mathbf{V}}_n \boldsymbol{\Sigma}_{AQ} \mathbf{B} \mathbf{u}_{n+1}. \quad (195)$$

Further simplification and rearrangement of terms gives

$$\hat{\boldsymbol{\mu}}_n = \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n (\hat{\boldsymbol{\mu}}_{n+1} - \langle \mathbf{A} \rangle \mathbf{G}_n \boldsymbol{\mu}_n - \langle \mathbf{B} \rangle \mathbf{u}_{n+1} + \hat{\mathbf{V}}_{n+1} \mathbf{J}_n^\top \boldsymbol{\Sigma}_{AQB} \mathbf{u}_{n+1}) - \hat{\mathbf{V}}_n \boldsymbol{\Sigma}_{AQB} \mathbf{u}_{n+1} \quad (196)$$

$$= \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n (\hat{\boldsymbol{\mu}}_{n+1} - \mathbf{m}_n - \langle \mathbf{A} \rangle \mathbf{G}_n \mathbf{V}_n \boldsymbol{\Sigma}_{AQB} \mathbf{u}_{n+1} + \hat{\mathbf{V}}_{n+1} \mathbf{J}_n^\top \boldsymbol{\Sigma}_{AQB} \mathbf{u}_{n+1}) - \hat{\mathbf{V}}_n \boldsymbol{\Sigma}_{AQB} \mathbf{u}_{n+1} \quad (197)$$

$$= \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n (\hat{\boldsymbol{\mu}}_{n+1} - \mathbf{m}_n) - \mathbf{J}_n \langle \mathbf{A} \rangle \mathbf{G}_n \mathbf{V}_n \boldsymbol{\Sigma}_{AQB} \mathbf{u}_{n+1} + \mathbf{J}_n \hat{\mathbf{V}}_{n+1} \mathbf{J}_n^\top \boldsymbol{\Sigma}_{AQB} \mathbf{u}_{n+1} - \hat{\mathbf{V}}_n \boldsymbol{\Sigma}_{AQB} \mathbf{u}_{n+1} \quad (198)$$

$$= \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n (\hat{\boldsymbol{\mu}}_{n+1} - \mathbf{m}_n) + \left(\mathbf{J}_n \hat{\mathbf{V}}_{n+1} \mathbf{J}_n^\top - \mathbf{J}_n \langle \mathbf{A} \rangle \mathbf{G}_n \mathbf{V}_n - \hat{\mathbf{V}}_n \right) \boldsymbol{\Sigma}_{AQB} \mathbf{u}_{n+1} \quad (199)$$

$$= \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n (\hat{\boldsymbol{\mu}}_{n+1} - \mathbf{m}_n) + \left(\mathbf{J}_n \hat{\mathbf{V}}_{n+1} \mathbf{J}_n^\top - \mathbf{J}_n \mathbf{P}_n \mathbf{J}_n^\top - \hat{\mathbf{V}}_n \right) \boldsymbol{\Sigma}_{AQB} \mathbf{u}_{n+1} \quad (200)$$

$$= \mathbf{G}_n \boldsymbol{\mu}_n + \mathbf{J}_n (\hat{\boldsymbol{\mu}}_{n+1} - \mathbf{m}_n) - \mathbf{G}_n \mathbf{V}_n \boldsymbol{\Sigma}_{AQB} \mathbf{u}_{n+1} \quad (201)$$

$$= \mathbf{G}_n (\boldsymbol{\mu}_n - \mathbf{V}_n \boldsymbol{\Sigma}_{AQB} \mathbf{u}_{n+1}) + \mathbf{J}_n (\hat{\boldsymbol{\mu}}_{n+1} - \mathbf{m}_n). \quad (202)$$

In summary, the backward pass computes the statistics of the smoothed marginal posterior probability:

$$q_{\mathbf{x}}(\mathbf{x}_n | \mathbf{Y}) = \mathcal{N}(\mathbf{x}_n | \hat{\boldsymbol{\mu}}_n, \hat{\mathbf{V}}_n) \quad (203)$$

where

$$\mathbf{J}_n = \mathbf{G}_n \mathbf{V}_n \langle \mathbf{A}^\top \rangle \mathbf{P}_n^{-1} \quad (204)$$

$$\hat{\boldsymbol{\mu}}_n = \mathbf{G}_n (\boldsymbol{\mu}_n - \mathbf{V}_n \boldsymbol{\Sigma}_{AQB} \mathbf{u}_{n+1}) + \mathbf{J}_n (\hat{\boldsymbol{\mu}}_{n+1} - \mathbf{m}_n) \quad (205)$$

$$\hat{\mathbf{V}}_n = \mathbf{G}_n \mathbf{V}_n + \mathbf{J}_n (\hat{\mathbf{V}}_{n+1} - \mathbf{P}_n) \mathbf{J}_n^\top. \quad (206)$$

This concludes the derivation of the backward pass (smoother) for the Bayesian linear dynamical system.

C. APPLICATION: VARIATIONAL M STEP DERIVATION

C.1. Joint distribution

The parameters of the Bayesian frequency estimation model are $\boldsymbol{\theta} = (\boldsymbol{\nu}, \boldsymbol{\tau}, \rho)$. The joint distribution over the variables and parameters is

$$p(\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}) = p(\mathbf{Y}|\mathbf{X}, \rho)p(\mathbf{X}|\boldsymbol{\nu}, \boldsymbol{\tau})p(\boldsymbol{\nu}|\boldsymbol{\tau})p(\boldsymbol{\tau})p(\rho) \quad (207)$$

where the conditional and marginal distributions are

$$p(\mathbf{Y}|\mathbf{X}, \rho) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{C}\mathbf{x}_n, R) \quad (208)$$

$$p(\mathbf{X}|\boldsymbol{\nu}, \boldsymbol{\tau}) = \mathcal{N}(\mathbf{x}_1 | \mathbf{m}_0, \mathbf{P}_0) \prod_{n=2}^N \mathcal{N}(\mathbf{x}_n | \mathbf{A}\mathbf{x}_{n-1}, \mathbf{Q}) \quad (209)$$

$$p(\boldsymbol{\nu}|\boldsymbol{\tau}) = \prod_{k=1}^K \mathcal{N}(\nu_k | 0, \alpha_k^{-1} \tau_k^{-1}) \quad (210)$$

$$p(\boldsymbol{\tau}) = \prod_{k=1}^K \text{Gam}(\tau_k | e_0, i_0) \quad (211)$$

$$p(\rho) = \text{Gam}(\rho | r_0, s_0) \quad (212)$$

C.2. Optimal distribution over parameters

The log optimal distribution over the parameters factors naturally as

$$\ln q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) = \langle \ln p(\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}) \rangle_{q_{\mathbf{x}}(\mathbf{x})} + \text{const.} \quad (213)$$

$$= \ln q_{\boldsymbol{\theta}}(\boldsymbol{\nu}, \boldsymbol{\tau}) + \ln q_{\boldsymbol{\theta}}(\rho) \quad (214)$$

The log optimal factors are derived in turn. First, we have

$$\ln q_{\boldsymbol{\theta}}(\boldsymbol{\nu}, \boldsymbol{\tau}) = \langle \ln p(\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}) \rangle_{q_{\mathbf{x}}(\mathbf{x})} + \text{const.} \quad (215)$$

$$= \langle \ln p(\mathbf{X}|\boldsymbol{\nu}, \boldsymbol{\tau}) \rangle_{q_{\mathbf{x}}(\mathbf{x})} + \ln p(\boldsymbol{\nu}|\boldsymbol{\tau}) + \ln p(\boldsymbol{\tau}) + \text{const.} \quad (216)$$

The expected log likelihood of the latent state sequence is given by

$$\begin{aligned} \langle \ln p(\mathbf{X}|\boldsymbol{\nu}, \boldsymbol{\tau}) \rangle_{q_{\mathbf{x}}(\mathbf{x})} &= \sum_{k=1}^K (N-1) \ln \tau_k - \frac{\tau_k}{2} \text{Tr} \left(\sum_{n=2}^N \langle \mathbf{x}_{n,k} \mathbf{x}_{n,k}^T \rangle \right) + \tau_k \text{Tr} \left(\mathbf{F} \sum_{n=2}^N \langle \mathbf{x}_{n-1,k} \mathbf{x}_{n,k}^T \rangle \right) \\ &\quad + \tau_k \nu_k \text{Tr} \left(\mathbf{E} \sum_{n=2}^N \langle \mathbf{x}_{n-1,k} \mathbf{x}_{n,k}^T \rangle \right) - \frac{\tau_k}{2} \text{Tr} \left(\mathbf{F} \sum_{n=2}^N \langle \mathbf{x}_{n-1,k} \mathbf{x}_{n-1,k} \rangle \mathbf{F}^T \right) \\ &\quad - \frac{\tau_k \nu_k^2}{2} \text{Tr} \left(\mathbf{E} \sum_{n=2}^N \langle \mathbf{x}_{n-1,k} \mathbf{x}_{n-1,k} \rangle \mathbf{E}^T \right) + \tau_k \nu_k \text{Tr} \left(\mathbf{E} \sum_{n=2}^N \langle \mathbf{x}_{n-1,k} \mathbf{x}_{n-1,k} \rangle \mathbf{F}^T \right) + \text{const.} \end{aligned} \quad (217)$$

The log prior distributions over $\boldsymbol{\nu}$ and $\boldsymbol{\tau}$ are, respectively,

$$\ln p(\boldsymbol{\nu}|\boldsymbol{\tau}) = \sum_{k=1}^K \frac{1}{2} \ln \tau_k - \frac{\tau_k}{2} \alpha_k \nu_k^2 + \text{const.} \quad (218)$$

$$\ln p(\boldsymbol{\tau}) = \sum_{k=1}^K (e_0 - 1) \ln \tau_k - i_0 \tau_k + \text{const.} \quad (219)$$

The optimal distribution over these parameters is of the same form as the prior distribution (Normal-Gamma):

$$q_{\boldsymbol{\theta}}(\boldsymbol{\nu}, \boldsymbol{\tau}) = q_{\boldsymbol{\nu}}(\boldsymbol{\nu}|\boldsymbol{\tau})q_{\boldsymbol{\tau}}(\boldsymbol{\tau}) = \prod_{k=1}^K \mathcal{N}(\nu_k | \hat{\nu}_k, \sigma_k \tau_k^{-1}) \text{Gam}(\tau_k | \hat{e}_k, \hat{i}_k). \quad (220)$$

The statistics of the optimal distribution are found by completing the square for ν_k and collecting the $\ln \tau_k$ and τ_k terms from the log optimal distribution above. They are given by

$$\sigma_k = \left(\text{Tr} \left(\mathbf{E} \sum_{n=2}^N \langle \mathbf{x}_{n-1k} \mathbf{x}_{n-1k}^\top \rangle \mathbf{E}^\top \right) + \alpha_k \right)^{-1} \quad (221)$$

$$\hat{\nu}_k = \sigma_k \text{Tr} \left(\mathbf{E} \left(\sum_{n=2}^N \langle \mathbf{x}_{n-1k} \mathbf{x}_{nk}^\top \rangle - \sum_{n=2}^N \langle \mathbf{x}_{n-1k} \mathbf{x}_{n-1k}^\top \rangle \mathbf{F}^\top \right) \right) \quad (222)$$

$$\hat{e}_k = e_0 + N - 1 \quad (223)$$

$$\hat{i}_k = i_0 - \frac{1}{2} \nu_k^2 \sigma_k^{-1} + \frac{1}{2} \text{Tr} \left(\mathbf{F} \sum_{n=2}^N \langle \mathbf{x}_{n-1k} \mathbf{x}_{n-1k}^\top \rangle \mathbf{F}^\top - 2 \mathbf{F} \sum_{n=2}^N \langle \mathbf{x}_{n-1k} \mathbf{x}_{nk}^\top \rangle + \sum_{n=2}^N \langle \mathbf{x}_{nk} \mathbf{x}_{nk}^\top \rangle \right) \quad (224)$$

where \mathbf{x}_{nk} is the k th 2×1 sub-vector of \mathbf{x}_n .

Next, the optimal distribution over the observation noise precision ρ is given by

$$\ln q_\theta(\rho) = \langle \ln p(\mathbf{Y}, \mathbf{X}, \theta) \rangle_{q_\mathbf{x}(\mathbf{x})} + \text{const.} \quad (225)$$

$$= \langle \ln p(\mathbf{Y} | \mathbf{X}, \rho) \rangle_{q_\mathbf{x}(\mathbf{x})} + \ln p(\rho) + \text{const..} \quad (226)$$

The expected log likelihood of the observed sequence is given by

$$\langle \ln p(\mathbf{Y} | \mathbf{X}, \theta) \rangle_{q_\mathbf{x}(\mathbf{x})} = \frac{N}{2} \ln \rho - \frac{\rho}{2} \left(\sum_{n=1}^N \mathbf{y}_n \mathbf{y}_n^\top - 2\mathbf{C} \sum_{n=1}^N \langle \mathbf{x}_n \rangle \mathbf{y}_n^\top + \mathbf{C} \sum_{n=1}^N \langle \mathbf{x}_n \mathbf{x}_n^\top \rangle \mathbf{C}^\top \right) + \text{const.} \quad (227)$$

The log prior distribution over ρ is

$$\ln p(\rho) = (r_0 - 1) \ln \rho - s_0 \rho. \quad (228)$$

As expected, ρ is Gamma-distributed:

$$q_\rho(\rho) = \text{Gam}(\rho | \hat{r}, \hat{s}). \quad (229)$$

The statistics of the optimal distribution are found by collecting terms involving $\ln \rho$ and ρ that appear in $\ln q_\theta(\rho)$.

$$\hat{r} = r_0 + N/2 \quad (230)$$

$$\hat{s} = s_0 + \frac{1}{2} \left(\sum_{n=1}^N \mathbf{y}_n \mathbf{y}_n^\top - 2\mathbf{C} \sum_{n=1}^N \langle \mathbf{x}_n \rangle \mathbf{y}_n^\top + \mathbf{C} \sum_{n=1}^N \langle \mathbf{x}_n \mathbf{x}_n^\top \rangle \mathbf{C}^\top \right). \quad (231)$$

This concludes the variational M step.

D. QUADRATIC FORMS OF RANDOM VARIABLES

This section derives the decomposition of the quadratic form of parameters under the optimal distribution q_{θ} , which appears in the equations of the proposed filter and smoother. A comprehensive reference on quadratic forms of random variables is [3].

D.1. General result

The general decomposition is given by

$$\langle \Psi^T \Lambda^{-1} \Omega \rangle = \langle \Psi \rangle^T \langle \Lambda^{-1} \rangle \langle \Omega \rangle + \Sigma_{\Psi \Lambda \Omega} \quad (232)$$

where we define the covariance term as

$$\Sigma_{\Psi \Lambda \Omega} = \sum_i \sum_j \langle \Lambda^{-1} \rangle_{(i,j)} \text{cov}[\psi_{(i)}, \omega_{(j)}]. \quad (233)$$

where $\psi_{(i)}$ denotes the transposed i th row of matrix Ψ , and the expectations are taken with respect to the optimal distribution $q_{\theta}(\theta)$. In the following subsections, we show that this general expression holds for an arbitrary definition of $p(\theta)$. First, we consider the case when the prior over the noise covariance matrix Λ is conjugate to the priors over the matrices Ψ and Ω . Second, we show that it holds for non-conjugate priors as well, because the variational approximation must involve an induced factorization between non-conjugate parameters to retain analytic tractability.

D.2. Case 1: Conjugate prior structure

In the first case the prior $p(\theta)$ is constructed in such a way that the dynamics matrices are conjugate to the noise covariance. This is a common way of defining the parameter model as it simplifies the approximation, reducing the number of induced factorizations needed to keep the variational method tractable. Existing literature defines the priors in this way [4] [5] [6].

The prior over the parameters is $p(\theta) = p(\mathbf{A}, \mathbf{B}, \mathbf{Q})p(\mathbf{C}, \mathbf{D}, \mathbf{R})$ where

$$p(\mathbf{A}, \mathbf{B}, \mathbf{Q}) = p(\mathbf{A}|\mathbf{Q})p(\mathbf{B}|\mathbf{Q})p(\mathbf{Q}) \quad (234)$$

$$p(\mathbf{C}, \mathbf{D}, \mathbf{R}) = p(\mathbf{C}|\mathbf{R})p(\mathbf{D}|\mathbf{R})p(\mathbf{R}) \quad (235)$$

and the conditional and marginal distributions are

$$p(\mathbf{A}|\mathbf{Q}) = \prod_{h=1}^H \mathcal{N}(\mathbf{a}_{(h)} | \mathbf{0}, \text{diag}(\boldsymbol{\alpha}_h \tau_h)^{-1}) \quad (236)$$

$$p(\mathbf{B}|\mathbf{Q}) = \prod_{h=1}^H \mathcal{N}(\mathbf{b}_{(h)} | \mathbf{0}, \text{diag}(\boldsymbol{\beta}_h \tau_h)^{-1}) \quad (237)$$

$$p(\mathbf{Q}) = \prod_{h=1}^H \text{Gam}(\tau_h | e_0, g_0) \quad (238)$$

$$p(\mathbf{C}|\mathbf{R}) = \prod_{v=1}^V \mathcal{N}(\mathbf{c}_{(v)} | \mathbf{0}, \text{diag}(\boldsymbol{\gamma}_v \rho_v)^{-1}) \quad (239)$$

$$p(\mathbf{D}|\mathbf{R}) = \prod_{v=1}^V \mathcal{N}(\mathbf{d}_{(v)} | \mathbf{0}, \text{diag}(\boldsymbol{\delta}_v \rho_v)^{-1}) \quad (240)$$

$$p(\mathbf{R}) = \prod_{v=1}^V \text{Gam}(\rho_v | r_0, s_0). \quad (241)$$

The optimal distribution over the parameters that maximizes the lower bound, $q_{\theta}(\theta)$, factors naturally into

$$q_{\theta}(\theta) = q_{\theta}(\mathbf{A}, \mathbf{B}, \mathbf{Q})q_{\theta}(\mathbf{C}, \mathbf{D}, \mathbf{R}) \quad (242)$$

$$= q_{\theta}(\mathbf{A}, \mathbf{B}|\mathbf{Q})q_{\theta}(\mathbf{Q})q_{\theta}(\mathbf{C}, \mathbf{D}|\mathbf{R})q_{\theta}(\mathbf{R}) \quad (243)$$

where

$$q_{\theta}(\mathbf{A}, \mathbf{B} | \mathbf{Q}) = \prod_{h=1}^H \mathcal{N} \left(\begin{pmatrix} \mathbf{a}_{(h)} \\ \mathbf{b}_{(h)} \end{pmatrix} \mid \begin{pmatrix} \widehat{\mathbf{a}}_{(h)} \\ \widehat{\mathbf{b}}_{(h)} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{aa}^h & \boldsymbol{\Sigma}_{ab}^h \\ \boldsymbol{\Sigma}_{ba}^h & \boldsymbol{\Sigma}_{bb}^h \end{pmatrix} \tau_h^{-1} \right) \quad (244)$$

$$q_{\theta}(\mathbf{Q}) = \prod_{h=1}^H \text{Gam}(\tau_h | \widehat{e}_h, \widehat{g}_h) \quad (245)$$

$$q_{\theta}(\mathbf{C}, \mathbf{D} | \mathbf{R}) = \prod_{v=1}^V \mathcal{N} \left(\begin{pmatrix} \mathbf{c}_{(v)} \\ \mathbf{d}_{(v)} \end{pmatrix} \mid \begin{pmatrix} \widehat{\mathbf{c}}_{(v)} \\ \widehat{\mathbf{d}}_{(v)} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{cc}^v & \boldsymbol{\Sigma}_{cd}^v \\ \boldsymbol{\Sigma}_{dc}^v & \boldsymbol{\Sigma}_{dd}^v \end{pmatrix} \rho_v^{-1} \right) \quad (246)$$

$$q_{\theta}(\mathbf{R}) = \prod_{v=1}^V \text{Gam}(\rho_v | \widehat{r}_v, \widehat{s}_v). \quad (247)$$

Now we inspect an example of a quadratic form of these parameters under the optimal distribution that appears in the proposed filter and smoother equations. We decompose it using the definition of the quadratic. Note that in the optimal distribution, each row of \mathbf{A} and \mathbf{B} are independent from the other rows.

$$\langle \mathbf{A}^T \mathbf{Q}^{-1} \mathbf{B} \rangle = \left\langle \sum_{i=1}^H \sum_{j=1}^H \mathbf{Q}_{(i,j)}^{-1} \mathbf{a}_{(i)} \mathbf{b}_{(j)}^T \right\rangle \quad \text{Definition of quadratic.} \quad (248)$$

$$= \sum_{i=1}^H \sum_{j=1}^H \langle \mathbf{Q}_{(i,j)}^{-1} \mathbf{a}_{(i)} \mathbf{b}_{(j)}^T \rangle \quad \text{Linearity of expected value.} \quad (249)$$

Since the noise covariance is diagonal, the summand is non-zero only when $i = j = h$:

$$\langle \mathbf{A}^T \mathbf{Q}^{-1} \mathbf{B} \rangle = \sum_{h=1}^H \langle \tau_h \mathbf{a}_{(h)} \mathbf{b}_{(h)}^T \rangle. \quad (250)$$

Using standard results of the Normal-Gamma distribution [7], this expected value is

$$\langle \mathbf{A}^T \mathbf{Q}^{-1} \mathbf{B} \rangle = \sum_{h=1}^H \langle \tau_h \rangle (\langle \mathbf{a}_{(h)} \rangle \langle \mathbf{b}_{(h)} \rangle^T + \text{cov}[\mathbf{a}_{(h)}, \mathbf{b}_{(h)}^T]) \quad (251)$$

$$= \sum_{h=1}^H \langle \tau_h \rangle (\langle \mathbf{a}_{(h)} \rangle \langle \mathbf{b}_{(h)} \rangle^T + \sum_{h=1}^H \langle \tau_h \rangle \text{cov}[\mathbf{a}_{(h)}, \mathbf{b}_{(h)}^T]) \quad (252)$$

$$= \langle \mathbf{A} \rangle^T \text{diag}(\langle \boldsymbol{\tau} \rangle) \langle \mathbf{B} \rangle + \sum_{h=1}^H \langle \tau_h \rangle \text{cov}[\mathbf{a}_{(h)}, \mathbf{b}_{(h)}^T]. \quad (253)$$

Finally, since $\langle \mathbf{Q}^{-1} \rangle = \text{diag}(\langle \boldsymbol{\tau} \rangle)$, we substitute in the covariance from the optimal distribution and simplify:

$$\langle \mathbf{A}^T \mathbf{Q}^{-1} \mathbf{B} \rangle = \langle \mathbf{A} \rangle^T \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle + \sum_{h=1}^H \langle \tau_h \rangle \text{cov}[\mathbf{a}_{(h)}, \mathbf{b}_{(h)}^T] \quad (254)$$

$$= \langle \mathbf{A} \rangle^T \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle + \sum_{h=1}^H \langle \tau_h \rangle \boldsymbol{\Sigma}_{ab}^h \langle \tau_h \rangle^{-1} \quad (255)$$

$$= \langle \mathbf{A} \rangle^T \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle + \sum_{h=1}^H \boldsymbol{\Sigma}_{ab}^h \quad (256)$$

$$= \langle \mathbf{A} \rangle^T \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle + \boldsymbol{\Sigma}_{AQB}. \quad (257)$$

where $\boldsymbol{\Sigma}_{AQB} = \sum_{h=1}^H \boldsymbol{\Sigma}_{ab}^h$. An identical result is obtained by starting from the general expression in equation (233).

D.3. Case 2: Non-conjugate prior structure

More generally, we can assume any prior $p(\boldsymbol{\theta})$ that is not necessarily constructed from a conjugate pair of Normal and Gamma distributions. In this case, the optimal distribution $q_{\boldsymbol{\theta}}(\boldsymbol{\theta})$ will need to be factorized to keep the variational procedure tractable. Here we assume that only the noise needs to be factored from the dynamics matrices (depending on the prior, it might also be the case where the dynamics matrices need to be factored from one another):

$$q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) = q_{\boldsymbol{\theta}}(\mathbf{A}, \mathbf{B})q_{\boldsymbol{\theta}}(\mathbf{Q})q_{\boldsymbol{\theta}}(\mathbf{C}, \mathbf{D})q_{\boldsymbol{\theta}}(\mathbf{R}) \quad (258)$$

Without detailing any particular prior distributions or resulting optimal distribution forms, it is simple to show that the quadratic form of these parameters decomposes as follows. Again, we start with the definition of the quadratic.

$$\langle \mathbf{A}^T \mathbf{Q}^{-1} \mathbf{B} \rangle = \left\langle \sum_{i=1}^H \sum_{j=1}^H \mathbf{Q}_{(i,j)}^{-1} \mathbf{a}_{(i)} \mathbf{b}_{(j)}^T \right\rangle \quad \text{Definition of quadratic.} \quad (259)$$

$$= \sum_{i=1}^H \sum_{j=1}^H \langle \mathbf{Q}_{(i,j)}^{-1} \mathbf{a}_{(i)} \mathbf{b}_{(j)}^T \rangle \quad \text{Linearity of expected value.} \quad (260)$$

Now, since the optimal distribution of the noise covariance matrix is independent of the other parameters, we get

$$\langle \mathbf{A}^T \mathbf{Q}^{-1} \mathbf{B} \rangle = \sum_{i=1}^H \sum_{j=1}^H \langle \mathbf{Q}^{-1} \rangle_{(i,j)} \langle \mathbf{a}_{(i)} \mathbf{b}_{(j)}^T \rangle \quad \mathbf{Q} \text{ is independent of } \mathbf{A} \text{ and } \mathbf{B}. \quad (261)$$

$$= \sum_{i=1}^H \sum_{j=1}^H \langle \mathbf{Q}^{-1} \rangle_{(i,j)} (\langle \mathbf{a}_{(i)} \rangle \langle \mathbf{b}_{(j)} \rangle^T + \text{cov}[\mathbf{a}_{(i)}, \mathbf{b}_{(j)}]) \quad \text{Expand the second moment.} \quad (262)$$

$$= \langle \mathbf{A} \rangle^T \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle + \sum_{i=1}^H \sum_{j=1}^H \langle \mathbf{Q}^{-1} \rangle_{(i,j)} \text{cov}[\mathbf{a}_{(i)}, \mathbf{b}_{(j)}] \quad \text{Simplify first term with definition of quadratic.} \quad (263)$$

In summary, the quadratic forms of the random parameters under their optimal distributions can be expressed as

$$\langle \mathbf{A}^T \mathbf{Q}^{-1} \mathbf{A} \rangle = \langle \mathbf{A} \rangle^T \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{A} \rangle + \Sigma_{AQ} \quad (264)$$

$$\langle \mathbf{A}^T \mathbf{Q}^{-1} \mathbf{B} \rangle = \langle \mathbf{A} \rangle^T \langle \mathbf{Q}^{-1} \rangle \langle \mathbf{B} \rangle + \Sigma_{AQ} \quad (265)$$

$$\langle \mathbf{C}^T \mathbf{R}^{-1} \mathbf{C} \rangle = \langle \mathbf{C} \rangle^T \langle \mathbf{R}^{-1} \rangle \langle \mathbf{C} \rangle + \Sigma_{CR} \quad (266)$$

$$\langle \mathbf{C}^T \mathbf{R}^{-1} \mathbf{D} \rangle = \langle \mathbf{C} \rangle^T \langle \mathbf{R}^{-1} \rangle \langle \mathbf{D} \rangle + \Sigma_{CD} \quad (267)$$

where we have defined

$$\Sigma_{AQ} = \sum_{i=1}^H \sum_{j=1}^H \langle \mathbf{Q}^{-1} \rangle_{(i,j)} \text{cov}[\mathbf{a}_{(i)}, \mathbf{a}_{(j)}] \quad (268)$$

$$\Sigma_{AQ} = \sum_{i=1}^H \sum_{j=1}^H \langle \mathbf{Q}^{-1} \rangle_{(i,j)} \text{cov}[\mathbf{a}_{(i)}, \mathbf{b}_{(j)}] \quad (269)$$

$$\Sigma_{CR} = \sum_{i=1}^H \sum_{j=1}^V \langle \mathbf{R}^{-1} \rangle_{(i,j)} \text{cov}[\mathbf{c}_{(i)}, \mathbf{c}_{(j)}] \quad (270)$$

$$\Sigma_{CD} = \sum_{i=1}^H \sum_{j=1}^V \langle \mathbf{R}^{-1} \rangle_{(i,j)} \text{cov}[\mathbf{c}_{(i)}, \mathbf{d}_{(j)}]. \quad (271)$$

These have the same form as the general expression given in equation (233).

E. SCHUR COMPLEMENTS

The Schur complements were used extensively to derive the forward-backward equations. They are given here for reference.

$$\begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{pmatrix}^{-1} = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \quad (272)$$

$$\Sigma_{11} = (\Lambda_{11} - \Lambda_{12}\Lambda_{22}^{-1}\Lambda_{21})^{-1} \quad (273)$$

$$= \Lambda_{11}^{-1} + \Lambda_{11}^{-1}\Lambda_{12}\Sigma_{22}\Lambda_{21}\Lambda_{11}^{-1} \quad (274)$$

$$\Sigma_{22} = (\Lambda_{22} - \Lambda_{21}\Lambda_{11}^{-1}\Lambda_{12})^{-1} \quad (275)$$

$$= \Lambda_{22}^{-1} + \Lambda_{22}^{-1}\Lambda_{21}\Sigma_{11}\Lambda_{12}\Lambda_{22}^{-1} \quad (276)$$

$$\Sigma_{12} = -\Lambda_{11}^{-1}\Lambda_{12}\Sigma_{22} \quad (277)$$

$$= -\Sigma_{11}\Lambda_{12}\Lambda_{22}^{-1} \quad (278)$$

$$\Sigma_{21} = -\Sigma_{22}\Lambda_{21}\Lambda_{11}^{-1} \quad (279)$$

$$= -\Lambda_{22}^{-1}\Lambda_{21}\Sigma_{11} \quad (280)$$

F. REFERENCES

- [1] C. Bishop, *Pattern Recognition and Machine Learning*, Springer-Verlag New York, 2006.
- [2] H. Rauch, F. Tung, and C. Striebel, “Maximum likelihood estimates of linear dynamical systems,” *AIAA Journal*, vol. 3, no. 8, pp. 1445–1450, Aug. 1965.
- [3] A. Mathai and S. Provost, *Quadratic forms in random variables: Theory and applications*, Marcel Dekker, 1992.
- [4] M. Beal, *Variational Algorithms for Approximate Bayesian Inference*, Ph.D. thesis, University College London, May 2003.
- [5] D. Barber, A. Cemgil, and S. Chiappa, Eds., *Bayesian Time Series Models*, Cambridge University Press, 2011.
- [6] J. Luttinen, “Fast variational Bayesian linear state-space model,” in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, Springer, Berlin, Heidelberg, 2013, pp. 305–320.
- [7] J. Bernardo and A. Smith, *Bayesian theory*, John Wiley & Sons, Inc., 2000.