

# Estimation for the Linear Model with Uncertain Covariance Matrices

Dave Zachariah, Nafiseh Shariati, Mats Bengtsson, Magnus Jansson and Saikat Chatterjee

**Abstract**—We derive a maximum a posteriori estimator for the linear observation model, where the signal and noise covariance matrices are both uncertain. The uncertainties are treated probabilistically by modeling the covariance matrices with prior inverse-Wishart distributions. The nonconvex problem of jointly estimating the signal of interest and the covariance matrices is tackled by a computationally efficient fixed-point iteration as well as an approximate variational Bayes solution. The statistical performance of estimators is compared numerically to state-of-the-art estimators from the literature and shown to perform favorably.

## I. INTRODUCTION

The linear observation model

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{w} \in \mathbb{R}^m, \quad (1)$$

is ubiquitous in signal processing, statistics and machine learning, cf. [1]–[7]. Applications include regression problems, model fitting, functional magnetic resonance imaging, finite impulse response identification, block data estimation, stochastic channel estimation, tracking, sensor fusion and multi-antenna receivers [2], [8]–[10]. Here  $\mathbf{x} \in \mathbb{R}^n$  denotes the unknown signal of interest,  $\mathbf{H} \in \mathbb{R}^{m \times n}$  denotes a given matrix with full column rank and  $\mathbf{w}$  is zero-mean noise. Many estimation procedures rely on prior knowledge of the statistical properties of  $\mathbf{x}$  and/or  $\mathbf{w}$ . In particular, the covariance matrices  $\mathbf{P} = \text{Cov}(\mathbf{x})$  and  $\mathbf{R} = \text{Cov}(\mathbf{w})$  are assumed to be known. In practice, however, these statistical properties may be subject to uncertainties. If assigned nominal covariance matrices,  $\mathbf{P}_0$  and  $\mathbf{R}_0$ , are based on prior knowledge where the statistics are only approximately stationary and/or prior estimates subject to errors, the resulting inaccuracies lead to degradation of estimation performance.

One approach is to treat the covariance uncertainties deterministically. This entails specifying a class of possible parameter values [11]. For instance, one could model  $\mathbf{P} = \mathbf{P}_0 + \delta\mathbf{P}$  and  $\mathbf{R} = \mathbf{R}_0 + \delta\mathbf{R}$  and assume that the errors  $\delta\mathbf{P}$  and  $\delta\mathbf{R}$  have known bounds on their spectral norms. In this case, [12] derived the linear estimator of  $\mathbf{x}$  that minimizes the worst-case mean square error (MSE) over the specified class of covariance matrices, drawing upon work in [13], [14]. The problem was shown to be convex and solved in closed form.

The ‘minimax’ MSE approach [15], however, was found to be overly conservative when evaluating its MSE performance. To compensate for this [12] also applied a different criterion based on the minimum attainable MSE over the covariance uncertainty class. The ‘minimax regret’ approach aims to minimize the maximum possible deviation from this MSE value. For the problem to be tractable, however, the uncertainty class was restricted such that the eigenvectors of  $\mathbf{P}$  and  $\mathbf{R}$  equal the right and left singular vectors of  $\mathbf{H}$ , respectively, and further, that their eigenvalues have known bounds. To circumvent this restriction, [16] generalized the minimax regret approach and applied it to a wider covariance uncertainty class with element-wise bounds, but only for the signal covariance  $\mathbf{P}$ . Further, unlike [12] the resulting estimator is not obtained in closed-form but requires solving a semidefinite program with quartic complexity in signal dimension  $n$ . In sum, a drawback of the deterministic approaches is the requirement of a restricted parametric class of covariance uncertainties. Further, they are formulated for a single snapshot and do not provide estimates of the signal and noise covariances, both of which are valuable statistical information in certain applications.

A different approach is to treat the covariance uncertainties probabilistically. This entails specifying distributions for the uncertain parameters [5]. For instance, [17] and [18] model  $\mathbf{w}$  as a Gaussian random variable and use various prior distributions on  $\mathbf{R}$ . The signal of interest  $\mathbf{x}$  is modeled with a noninformative prior distribution and therefore no signal covariance matrix  $\mathbf{P}$  is considered. In [17], the prior distribution of  $\mathbf{R}$  is noninformative resulting in closed-form solutions of the parameter estimates. By contrast, [18] consider informative priors for  $\mathbf{R}$  but require a sampling-based Markov chain Monte Carlo (MCMC) method for solving the problem, which becomes computationally intractable for larger signal dimensions.

In this paper we seek to generalize the probabilistic approach to jointly estimate the signal of interest, as well as the signal *and* noise covariance matrices. Both unknown matrices are modeled as random and independent quantities around the nominal ones, using tractable priors. To the best of the authors’ knowledge this has not been addressed and solved in a tractable way in the literature. In this work we use the inverse-Wishart distribution, which is a conjugate prior to the covariance matrix of a Gaussian distribution. A discussion on the use of this distribution is given in [5], [17], [18], where it is shown to be a modified version of the noninformative Jeffreys prior. The inverse-Wishart distribution has also been used in detection problems where the inaccuracies of the nominal covariance matrices arise due to environmental heterogeneity

The authors are with the ACCESS Linnaeus Centre, KTH Royal Institute of Technology, Stockholm. E-mail: {dave.zachariah, nafiseh, mats.bengtsson, magnus.jansson}@ee.kth.se and saikatchatt@gmail.com. This research has partly been funded by the Swedish Research Council under contracts 621-2011-5847 and 621-2012-4134. The research leading to these results has received funding from the European Research Council under the European Community’s Seventh Framework Programme (FP7/2007-2013) / ERC grant agreement n° 228044.

[19]–[21].

We show that the maximum a posteriori probability estimator results in a nonconvex optimization problem, but reveals certain connections with the standard estimators. To solve the problem in a computationally efficient manner we formulate a fixed-point iteration. Whilst proving convergence appears intractable, we prove that the iteration does not diverge and illustrate its converge properties empirically. Further, we derive a variational Bayes solution to the problem as a tractable but approximate alternative. Finally, the resulting estimators are evaluated in terms of average performance and robustness.

*Notation:*  $|\mathbf{A}|$  and  $\text{tr}\{\mathbf{A}\}$  denote the determinant and trace of  $\mathbf{A}$ , respectively.  $\mathbf{A} \otimes \mathbf{B}$  denotes the Kronecker product of matrices and  $\|\cdot\|_F$  denotes the Frobenius norm.  $\mathbf{E}_{ij}$  is the  $ij$ th standard basis matrix.  $\mathcal{N}(\boldsymbol{\mu}, \mathbf{P})$  denotes a Gaussian distribution with mean  $\boldsymbol{\mu}$  and covariance  $\mathbf{P}$ . The inverse-Wishart distribution with parameters  $\nu$  and  $\mathbf{C}$  is denoted  $\mathcal{W}^{-1}(\mathbf{C}, \nu)$ .

## II. PROBLEM FORMULATION

For generality we consider a set of  $N$  measurements  $\{\mathbf{y}_t\}_{t=1}^N$  and corresponding signals of interests  $\{\mathbf{x}_t\}_{t=1}^N$ . For notational simplicity we write  $\mathbf{Y} \triangleq [\mathbf{y}_1 \cdots \mathbf{y}_N] \in \mathbb{R}^{m \times N}$  and  $\mathbf{X} \triangleq [\mathbf{x}_1 \cdots \mathbf{x}_N] \in \mathbb{R}^{n \times N}$ . Then the linear observation model (1) is written as

$$\mathbf{Y} = \mathbf{H}\mathbf{X} + \mathbf{W}. \quad (2)$$

It is assumed that the signal and noise follow independent Gaussian distributions  $\mathbf{x}_t|\mathbf{P} \sim \mathcal{N}(\boldsymbol{\mu}_t, \mathbf{P})$  and  $\mathbf{w}_t|\mathbf{R} \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$ .

When the covariance matrices are known, the maximum a posteriori (MAP) estimator of  $\mathbf{X}$  coincides with the familiar linear minimum MSE estimator,

$$\begin{aligned} \hat{\mathbf{X}}_{\text{map}} &= \arg \max_{\mathbf{X} \in \mathbb{R}^{n \times N}} p(\mathbf{X}|\mathbf{Y}) \\ &= (\mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H} + \mathbf{P}^{-1})^{-1} (\mathbf{H}^\top \mathbf{R}^{-1} \mathbf{Y} + \mathbf{P}^{-1} \mathbf{U}), \end{aligned} \quad (3)$$

where  $\mathbf{U} \triangleq [\boldsymbol{\mu}_1 \cdots \boldsymbol{\mu}_N] \in \mathbb{R}^{n \times N}$  [3]. As the uncertainty or variance of the prior of  $\mathbf{x}_t$  increases, by setting  $\mathbf{P} = \sigma_x^2 \mathbf{I}_n$  and  $\sigma_x^2 \rightarrow \infty$ , the estimator coincides with the minimum variance unbiased (MVU) estimator,  $\hat{\mathbf{X}}_{\text{map}} \rightarrow \hat{\mathbf{X}}_{\text{mvu}} = (\mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{R}^{-1} \mathbf{Y}$  [2]. When  $\mathbf{P}$  and  $\mathbf{R}$  are not known precisely they are replaced by nominal matrices,  $\mathbf{P}_0$  and  $\mathbf{R}_0$ .

Henceforth the unknown covariance matrices are modeled as random and independent quantities around the nominal ones, using inverse-Wishart distributions:  $\mathbf{P} \sim \mathcal{W}^{-1}(\mathbf{C}_x, \nu_x)$  and  $\mathbf{R} \sim \mathcal{W}^{-1}(\mathbf{C}_w, \nu_w)$ . Assuming that  $\text{E}[\mathbf{P}] = \mathbf{P}_0$  and  $\text{E}[\mathbf{R}] = \mathbf{R}_0$ , we have  $\mathbf{C}_x = (\nu_x - n - 1)\mathbf{P}_0$  and  $\mathbf{C}_w = (\nu_w - m - 1)\mathbf{R}_0$ . The degrees of freedom,  $\nu_x > n + 1$  and  $\nu_w > m + 1$ , control the certainties of  $\mathbf{P}$  and  $\mathbf{R}$ . Extensions to the complex Gaussian and inverse-Wishart distributions [22] are straight-forward.

The goal is to estimate  $\mathbf{X}$ ,  $\mathbf{P}$  and  $\mathbf{R}$  from the set of observations  $\mathbf{Y}$ .

## III. THE CMAP ESTIMATOR

The maximum a posterior estimator with random covariance matrices, henceforth denoted CMAP, is obtained by solving

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times N}, \mathbf{P} \succ \mathbf{0}, \mathbf{R} \succ \mathbf{0}} p(\mathbf{X}, \mathbf{P}, \mathbf{R}|\mathbf{Y}), \quad (4)$$

where  $p(\mathbf{X}, \mathbf{P}, \mathbf{R}|\mathbf{Y})$  denotes the joint posterior probability density function (pdf). By applying Bayes' rule and introducing

$$\begin{aligned} J(\mathbf{X}, \mathbf{P}, \mathbf{R}) &\triangleq \ln p(\mathbf{Y}|\mathbf{X}, \mathbf{P}, \mathbf{R}) + \ln p(\mathbf{X}, \mathbf{P}, \mathbf{R}) \\ &= \ln p(\mathbf{Y}|\mathbf{X}, \mathbf{R}) + \ln (p(\mathbf{X}|\mathbf{P})p(\mathbf{P})p(\mathbf{R})) \\ &= J_1(\mathbf{X}, \mathbf{R}) + J_2(\mathbf{X}, \mathbf{P}), \end{aligned}$$

where  $J_1(\mathbf{X}, \mathbf{R}) = [\ln p(\mathbf{Y}|\mathbf{X}, \mathbf{R}) + \ln p(\mathbf{R})]$  and  $J_2(\mathbf{X}, \mathbf{P}) = [\ln p(\mathbf{X}|\mathbf{P}) + \ln p(\mathbf{P})]$ , we can tackle the problem by first solving for  $\mathbf{R}$  and  $\mathbf{P}$ . Then

$$\hat{\mathbf{X}}_{\text{cmap}} = \arg \max_{\mathbf{X} \in \mathbb{R}^{n \times N}} \left[ \max_{\mathbf{R} \succ \mathbf{0}, \mathbf{P} \succ \mathbf{0}} J_1(\mathbf{X}, \mathbf{R}) + J_2(\mathbf{X}, \mathbf{P}) \right]. \quad (5)$$

We begin by finding the maximizing  $\mathbf{R}$  and  $\mathbf{P}$  below.

### A. Concentrated cost function

Let  $\tilde{\mathbf{y}}_t \triangleq \mathbf{y}_t - \mathbf{H}\mathbf{x}_t$ ,  $\tilde{\mathbf{Y}} \triangleq \mathbf{Y} - \mathbf{H}\mathbf{X}$  and  $\gamma_w \triangleq \nu_w + m + 1 + N$ , so that

$$\begin{aligned} J_1(\mathbf{X}, \mathbf{R}) &= \ln p(\tilde{\mathbf{Y}}|\mathbf{R}) + \ln p(\mathbf{R}) \\ &= \sum_{t=1}^N -\frac{1}{2} \ln |\mathbf{R}| - \frac{1}{2} \text{tr} \{ \mathbf{R}^{-1} \tilde{\mathbf{y}}_t \tilde{\mathbf{y}}_t^\top \} \\ &\quad - \frac{\nu_w + m + 1}{2} \ln |\mathbf{R}| - \frac{1}{2} \text{tr} \{ \mathbf{C}_w \mathbf{R}^{-1} \} + K \\ &= -\frac{\gamma_w}{2} \ln |\mathbf{R}| - \frac{1}{2} \text{tr} \{ (\mathbf{C}_w + \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^\top) \mathbf{R}^{-1} \} + K \\ &= \frac{\gamma_w}{2} \left( -\ln |\mathbf{R}| - \text{tr} \{ \tilde{\mathbf{R}} \mathbf{R}^{-1} \} \right) + K, \end{aligned} \quad (6)$$

where  $K$  denotes an unimportant constant and  $\tilde{\mathbf{R}} \triangleq \frac{1}{\gamma_w} (\mathbf{C}_w + \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^\top)$ . Then

$$\begin{aligned} \tilde{J}_1(\mathbf{X}, \mathbf{R}) &\triangleq -\ln |\mathbf{R}| - \text{tr} \{ \tilde{\mathbf{R}} \mathbf{R}^{-1} \} \\ &= -\ln |\tilde{\mathbf{R}} \tilde{\mathbf{R}}^{-1} \mathbf{R}| - \text{tr} \{ \tilde{\mathbf{R}} \mathbf{R}^{-1} \} \\ &= -\ln |\tilde{\mathbf{R}} (\mathbf{R}^{-1} \tilde{\mathbf{R}})^{-1}| - \text{tr} \{ \tilde{\mathbf{R}} \mathbf{R}^{-1} \} \\ &= -\ln |\tilde{\mathbf{R}}| + \ln |\mathbf{R}^{-1} \tilde{\mathbf{R}}| - \text{tr} \{ \mathbf{R}^{-1} \tilde{\mathbf{R}} \} \end{aligned}$$

attains its maximum when  $\mathbf{R}^{-1} \tilde{\mathbf{R}} = \mathbf{I}_m$ , or  $\mathbf{R}^* = \frac{1}{\gamma_w} (\mathbf{C}_w + \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^\top)$ . Similarly, let  $\tilde{\mathbf{X}} \triangleq \mathbf{X} - \mathbf{U}$  and  $\gamma_x \triangleq \nu_x + n + 1 + N$ , then  $\mathbf{P}^* = \frac{1}{\gamma_x} (\mathbf{C}_x + \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top)$ . Note that both  $\mathbf{P}^*$  and  $\mathbf{R}^*$  are functions of  $\tilde{\mathbf{X}}$ .

Plugging back the solution, and using the matrix determinant lemma, yields

$$\begin{aligned}
J_1(\mathbf{X}, \mathbf{R}^*) &= -\frac{\gamma_w}{2} \ln \left| \frac{1}{\gamma_w} (\mathbf{C}_w + \tilde{\mathbf{Y}}\tilde{\mathbf{Y}}^\top) \right| - \frac{\gamma_w}{2} \text{tr} \{ \mathbf{I}_m \} + K \\
&= -\frac{\gamma_w}{2} \ln \left( \frac{1}{\gamma_w^m} \left| \mathbf{C}_w + \tilde{\mathbf{Y}}\tilde{\mathbf{Y}}^\top \right| \right) + K' \\
&= -\frac{\gamma_w}{2} \ln \left| \mathbf{C}_w + \tilde{\mathbf{Y}}\tilde{\mathbf{Y}}^\top \right| + K'' \\
&= -\frac{\gamma_w}{2} \ln \left( \left| \mathbf{C}_w \right| \left| \mathbf{I}_N + \tilde{\mathbf{Y}}^\top \mathbf{C}_w^{-1} \tilde{\mathbf{Y}} \right| \right) + K'' \\
&= -\frac{\gamma_w}{2} \ln \left| \mathbf{I}_N + \tilde{\mathbf{Y}}^\top \mathbf{C}_w^{-1} \tilde{\mathbf{Y}} \right| + K'''.
\end{aligned}$$

Similarly,

$$J_2(\mathbf{X}, \mathbf{P}^*) = -\frac{\gamma_x}{2} \ln \left| \mathbf{I}_N + \tilde{\mathbf{X}}^\top \mathbf{C}_x^{-1} \tilde{\mathbf{X}} \right| + K.$$

In sum, the optimal estimator is given by

$$\hat{\mathbf{X}}_{\text{cmap}} = \arg \min_{\mathbf{X} \in \mathbb{R}^{n \times N}} V(\mathbf{X}), \quad (7)$$

where the concentrated cost function equals

$$\begin{aligned}
V(\mathbf{X}) &\triangleq \frac{\gamma_w}{2} \ln \left| \mathbf{I}_N + (\mathbf{Y} - \mathbf{H}\mathbf{X})^\top \mathbf{C}_w^{-1} (\mathbf{Y} - \mathbf{H}\mathbf{X}) \right| \\
&\quad + \frac{\gamma_x}{2} \ln \left| \mathbf{I}_N + (\mathbf{X} - \mathbf{U})^\top \mathbf{C}_x^{-1} (\mathbf{X} - \mathbf{U}) \right|.
\end{aligned} \quad (8)$$

Next, we study the properties of the cost function by writing it as  $V(\mathbf{X}) = V_1(\mathbf{X}) + V_2(\mathbf{X})$ , where

$$\begin{aligned}
V_1(\mathbf{X}) &= \frac{\gamma_w}{2} \ln |\mathbf{A}(\mathbf{X})| \\
V_2(\mathbf{X}) &= \frac{\gamma_x}{2} \ln |\mathbf{B}(\mathbf{X})|,
\end{aligned}$$

and  $\mathbf{A}(\mathbf{X}) \triangleq \mathbf{I}_N + (\mathbf{Y} - \mathbf{H}\mathbf{X})^\top \mathbf{C}_w^{-1} (\mathbf{Y} - \mathbf{H}\mathbf{X}) \succ \mathbf{0}$  and  $\mathbf{B}(\mathbf{X}) \triangleq \mathbf{I}_N + (\mathbf{X} - \mathbf{U})^\top \mathbf{C}_x^{-1} (\mathbf{X} - \mathbf{U}) \succ \mathbf{0}$ . While the inner matrices are quadratic functions of  $\mathbf{X}$ , the log-determinant makes  $V_1(\mathbf{X})$  and  $V_2(\mathbf{X})$  nonconvex functions. Their minima, however, provide the key for finding minima of  $V(\mathbf{X})$ .

The minimum of  $V_1(\mathbf{X})$  is  $\hat{\mathbf{X}}_{\text{mvu}}$  and can be verified by computing the gradient. Using the chain-rule,

$$\frac{\partial V_1}{\partial x_{it}} = \text{tr} \left\{ (\partial_A V_1)^\top \frac{\partial \mathbf{A}}{\partial x_{it}} \right\},$$

where the inner derivative equals

$$\begin{aligned}
\frac{\partial \mathbf{A}}{\partial x_{it}} &= \frac{\partial}{\partial x_{it}} (\mathbf{I}_N + (\mathbf{Y} - \mathbf{H}\mathbf{X})^\top \mathbf{C}_w^{-1} (\mathbf{Y} - \mathbf{H}\mathbf{X})) \\
&= -\mathbf{E}_{it}^\top \mathbf{H}^\top \mathbf{C}_w^{-1} (\mathbf{Y} - \mathbf{H}\mathbf{X}) - (\mathbf{Y} - \mathbf{H}\mathbf{X})^\top \mathbf{C}_w^{-1} \mathbf{H} \mathbf{E}_{it},
\end{aligned}$$

and the outer derivative is  $\partial_A V_1 = \frac{\gamma_w}{2} \mathbf{A}^{-1}$  due to symmetry. Hence

$$\begin{aligned}
\frac{\partial V_1}{\partial x_{it}} &= -\frac{\gamma_w}{2} \text{tr} \left\{ \mathbf{A}^{-1} \mathbf{E}_{it}^\top \mathbf{H}^\top \mathbf{C}_w^{-1} (\mathbf{Y} - \mathbf{H}\mathbf{X}) \right\} \\
&\quad - \frac{\gamma_w}{2} \text{tr} \left\{ \mathbf{A}^{-1} (\mathbf{Y} - \mathbf{H}\mathbf{X})^\top \mathbf{C}_w^{-1} \mathbf{H} \mathbf{E}_{it} \right\} \\
&= -\gamma_w \text{tr} \left\{ \mathbf{H}^\top \mathbf{C}_w^{-1} (\mathbf{Y} - \mathbf{H}\mathbf{X}) \mathbf{A}^{-1} \mathbf{E}_{it} \right\}
\end{aligned}$$

and

$$\partial_X V_1 = -\gamma_w \mathbf{H}^\top \mathbf{C}_w^{-1} (\mathbf{Y} - \mathbf{H}\mathbf{X}) \mathbf{A}^{-1}. \quad (9)$$

Setting  $\partial_X V_1(\mathbf{X}) = \mathbf{0}$  and solving for  $\mathbf{X}$  yields the stationary point  $\hat{\mathbf{X}}_{\text{mvu}}$  since  $\mathbf{C}_w \propto \mathbf{R}_0$ . Then  $\hat{\mathbf{X}}_{\text{mvu}}$  is the minimizer of

$V_1(\mathbf{X})$ , since  $\ln |\cdot|$  is a monotonically increasing function on the set of positive definite matrices and the quadratic function  $\mathbf{A}(\mathbf{X}) \succeq \mathbf{A}(\hat{\mathbf{X}}_{\text{mvu}})$ .

Similarly, the trivial minimizer of  $V_2(\mathbf{X})$  is  $\mathbf{U}$ , and can be verified by

$$\partial_X V_2 = \gamma_x \mathbf{C}_x^{-1} (\mathbf{X} - \mathbf{U}) \mathbf{B}^{-1}. \quad (10)$$

When a realization  $\hat{\mathbf{X}}_{\text{mvu}}$  is far apart from  $\mathbf{U}$  then, in the vicinity of the minimizer of  $V_1(\mathbf{X})$ ,  $V_2(\mathbf{X})$  is approximately constant, and vice versa, due to the compressive property of the logarithm. In the extreme, therefore,  $V(\mathbf{X})$  may have at least two separated minima, located in the vicinity of  $\hat{\mathbf{X}}_{\text{mvu}}$  and  $\mathbf{U}$ , respectively, and the estimator is not amenable to closed-form solution. On the other hand, when  $\hat{\mathbf{X}}_{\text{mvu}}$  is close to  $\mathbf{U}$ , a single minimum of  $V(\mathbf{X})$  may result. These extreme scenarios are illustrated in Fig. 1.

Using  $\hat{\mathbf{X}}_{\text{mvu}}$  and  $\mathbf{U}$  as starting points, minima of  $V(\mathbf{X})$  can be found by gradient descent  $\hat{\mathbf{X}}^{\ell+1} = \hat{\mathbf{X}}^\ell - \mu \partial_X V(\hat{\mathbf{X}}^\ell)$ , where  $\mu > 0$  is the step size and  $\partial_X V = \partial_X V_1 + \partial_X V_2$  given by (9) and (10). The partial derivatives can be written in alternative forms that are computationally advantageous when  $N > n$  and  $N > m$ , using the matrix inversion lemma,

$$\begin{aligned}
\partial_X V_1 &= -\gamma_w \mathbf{H}^\top \mathbf{C}_w^{-1} \tilde{\mathbf{Y}} (\mathbf{I}_N + \tilde{\mathbf{Y}}^\top \mathbf{C}_w^{-1} \tilde{\mathbf{Y}})^{-1} \\
&= -\gamma_w \mathbf{H}^\top \mathbf{C}_w^{-1} \left( \mathbf{I}_N - \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^\top (\mathbf{C}_w + \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^\top)^{-1} \right) \tilde{\mathbf{Y}} \\
&= -\gamma_w \mathbf{H}^\top \mathbf{C}_w^{-1} (\mathbf{I}_N + \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^\top \mathbf{C}_w^{-1})^{-1} \tilde{\mathbf{Y}} \\
&= -\gamma_w \mathbf{H}^\top (\mathbf{C}_w + \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^\top)^{-1} \tilde{\mathbf{Y}}
\end{aligned}$$

and similarly

$$\begin{aligned}
\partial_X V_2 &= \gamma_x \mathbf{C}_x^{-1} \tilde{\mathbf{X}} (\mathbf{I}_N + \tilde{\mathbf{X}}^\top \mathbf{C}_x^{-1} \tilde{\mathbf{X}})^{-1} \\
&= \gamma_x (\mathbf{C}_x + \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top)^{-1} \tilde{\mathbf{X}}.
\end{aligned}$$

Thus

$$\begin{aligned}
\partial_X V &= -\gamma_w \mathbf{H}^\top \mathbf{C}_w^{-1} (\mathbf{Y} - \mathbf{H}\mathbf{X}) \mathbf{A}^{-1} + \gamma_x \mathbf{C}_x^{-1} (\mathbf{X} - \mathbf{U}) \mathbf{B}^{-1} \\
&= -\mathbf{H}^\top \hat{\mathbf{R}}^{-1} (\mathbf{Y} - \mathbf{H}\mathbf{X}) + \hat{\mathbf{P}}^{-1} (\mathbf{X} - \mathbf{U}),
\end{aligned}$$

where

$$\begin{aligned}
\hat{\mathbf{P}}(\mathbf{X}) &= \frac{1}{\gamma_x} (\mathbf{C}_x + (\mathbf{X} - \mathbf{U})(\mathbf{X} - \mathbf{U})^\top) \\
\hat{\mathbf{R}}(\mathbf{X}) &= \frac{1}{\gamma_w} (\mathbf{C}_w + (\mathbf{Y} - \mathbf{H}\mathbf{X})(\mathbf{Y} - \mathbf{H}\mathbf{X})^\top)
\end{aligned} \quad (11)$$

are the covariance matrix estimates. Note that their inverses can be computed recursively by a series of rank-1 updates, using the Sherman-Morrison formula [23]. The overall computational efficiency of the gradient decent method is, however, dependent on the user-defined step size  $\mu$ . To circumvent this limitation, we devise an alternative fixed-point iteration method.

### B. Fixed-point iteration

We attempt to find the local minima by iteratively fulfilling the condition for a stationary point. The solution to  $\partial_X V(\mathbf{X}) = \mathbf{0}$ , when holding the nonlinear functions  $\hat{\mathbf{P}}(\mathbf{X})$  and  $\hat{\mathbf{R}}(\mathbf{X})$  constant for a given estimate  $\hat{\mathbf{X}}^\ell$ , equals

$$\hat{\mathbf{X}}^{\ell+1} = (\mathbf{H}^\top \hat{\mathbf{R}}_\ell^{-1} \mathbf{H} + \hat{\mathbf{P}}_\ell^{-1})^{-1} (\mathbf{H}^\top \hat{\mathbf{R}}_\ell^{-1} \mathbf{Y} + \hat{\mathbf{P}}_\ell^{-1} \mathbf{U}) \quad (12)$$

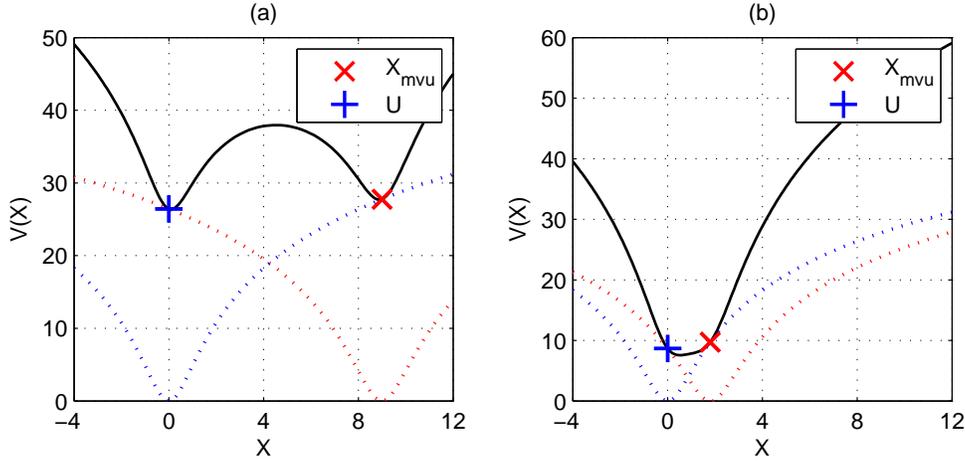


Fig. 1. Example of cost function  $V(\mathbf{X})$  where  $n = 1$ ,  $N = 1$  and  $m = 1$  for sake of illustration. Dotted lines show  $V_1(\mathbf{X})$  and  $V_2(\mathbf{X})$ . Here  $\mathbf{Y} = \mathbf{H}\mathbf{X} + \mathbf{W}$ , where  $\mathbf{X} = 1$  and  $\mathbf{H} = 1$ . Nominal variances  $\mathbf{P}_0 = 0.8$  and  $\mathbf{R}_0 = 1$  with minimum certainties. (a)  $\mathbf{W} = 8$  resulting in local minima of  $V(\mathbf{X})$ . (b)  $\mathbf{W} = 0.8$  resulting in a single minimum of  $V(\mathbf{X})$ . Note that the minima occur in the vicinity of  $\hat{\mathbf{X}}_{\text{mvu}}$  and  $\mathbf{U}$ .

and is iterated until convergence. Comparing (12) with (3) it is immediately recognized that the fixed-point method is an iterative application of the standard MAP estimator with covariance matrices  $\hat{\mathbf{P}}_\ell = \hat{\mathbf{P}}(\hat{\mathbf{X}}^\ell)$  and  $\hat{\mathbf{R}}_\ell = \hat{\mathbf{R}}(\hat{\mathbf{X}}^\ell)$ . Based on the analysis of the previous section, we propose using  $\hat{\mathbf{X}}_{\text{mvu}}$  and  $\mathbf{U}$  as two starting points, respectively. The resulting minimum with the lowest cost  $V(\mathbf{X})$  is then used as the estimate. When the costs happen to be equal, the estimator is indifferent and we can choose the solution that is closest to the MAP estimate, which assumes that the nominal covariances are true. Our numerical experiments show that the iterative solution is very likely to produce the optimal estimate, cf. section IV-D.

The CMAP estimator is summarized in Algorithm 1. The function  $\text{iter}(\cdot)$  iterates (12) until  $\|\hat{\mathbf{X}}^\ell - \hat{\mathbf{X}}^{\ell-1}\|_F < \varepsilon$ .

For a derivation of the conditions for convergence of (12) it would be sufficient to prove that the iteration is a contraction mapping [24]. Deriving these conditions appears intractable in general. However, it is possible to show that the iterative solution (12) does not diverge. Let  $\hat{\mathbf{Y}}^{\ell+1} = \mathbf{H}\hat{\mathbf{X}}^{\ell+1}$  denote the predicted observation, and  $\hat{\mathbf{y}}_t$  denote the  $t$ th column of  $\hat{\mathbf{Y}}^{\ell+1}$ . If  $\|\hat{\mathbf{y}}_t\|_2^2 = \hat{\mathbf{x}}_t^\top \mathbf{H}^\top \mathbf{H} \hat{\mathbf{x}}_t$  is bounded, then  $\|\hat{\mathbf{x}}_t\|_2^2$  is bounded since  $\mathbf{H}$  has full rank and  $\mathbf{H}^\top \mathbf{H} \succ \mathbf{0}$ . Hence  $\|\hat{\mathbf{Y}}^{\ell+1}\|_F^2 < \infty \Rightarrow \|\hat{\mathbf{X}}^{\ell+1}\|_F^2 < \infty$ . Next, consider  $\mathbf{U} = \mathbf{0}$ ,<sup>1</sup> so that (12) can be written as  $\hat{\mathbf{X}}^{\ell+1} = \hat{\mathbf{P}}_\ell \mathbf{H}^\top (\hat{\mathbf{R}}_\ell + \mathbf{H} \hat{\mathbf{P}}_\ell \mathbf{H}^\top)^{-1} \mathbf{Y}$ , and define  $\mathbf{\Gamma}_\ell \triangleq \hat{\mathbf{H}} \hat{\mathbf{P}}_\ell \mathbf{H}^\top \succ \mathbf{0}$  and  $\mathbf{\Phi}_\ell \triangleq \hat{\mathbf{R}}_\ell + \mathbf{\Gamma}_\ell \succ \mathbf{\Gamma}_\ell$ . Hence  $\|\mathbf{\Gamma}_\ell \mathbf{\Phi}_\ell^{-1}\|_2^2 < 1$  and it follows that  $\|\hat{\mathbf{y}}_t\|_2^2 = \|\mathbf{H} \hat{\mathbf{x}}_t\|_2^2 = \|\mathbf{\Gamma}_\ell \mathbf{\Phi}_\ell^{-1} \mathbf{y}_t\|_2^2 \leq \|\mathbf{\Gamma}_\ell \mathbf{\Phi}_\ell^{-1}\|_2^2 \|\mathbf{y}_t\|_2^2 < \|\mathbf{y}_t\|_2^2$ . Therefore  $\|\hat{\mathbf{Y}}^{\ell+1}\|_F^2$  is bounded and consequently  $\|\hat{\mathbf{X}}^{\ell+1}\|_F^2$  is bounded for all  $\ell$ . The iterative solution (12) must either converge or produce a bounded orbit. In fact, through extensive simulations the algorithm was always found to converge. In section IV-E we present an empirical convergence analysis of the fixed-point iteration.

<sup>1</sup>This is no restriction as it is possible to define an equivalent problem with zero-mean variables,  $\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{U}$  and  $\tilde{\mathbf{Y}} = \mathbf{Y} - \mathbf{H}\mathbf{U}$ , and then shift the estimate of  $\tilde{\mathbf{X}}$ .

---

#### Algorithm 1 CMAP estimator

---

- 1: Input:  $\mathbf{Y}$ ,  $\mathbf{H}$ ,  $\mathbf{C}_x$ ,  $\mathbf{C}_w$ ,  $\gamma_x$ ,  $\gamma_w$ ,  $\varepsilon$
  - 2:  $\hat{\mathbf{X}}_1^0 = \hat{\mathbf{X}}_{\text{mvu}}$  and  $\hat{\mathbf{X}}_2^0 = \mathbf{U}$
  - 3:  $\hat{\mathbf{X}}_1 = \text{iter}(\mathbf{Y}, \hat{\mathbf{X}}_1^0, \mathbf{H}, \mathbf{C}_x, \mathbf{C}_w, \gamma_x, \gamma_w, \varepsilon)$
  - 4:  $\hat{\mathbf{X}}_2 = \text{iter}(\mathbf{Y}, \hat{\mathbf{X}}_2^0, \mathbf{H}, \mathbf{C}_x, \mathbf{C}_w, \gamma_x, \gamma_w, \varepsilon)$
  - 5: **if**  $V(\hat{\mathbf{X}}_1) < V(\hat{\mathbf{X}}_2)$  **then**
  - 6:      $\hat{\mathbf{X}} := \hat{\mathbf{X}}_1$
  - 7: **else if**  $V(\hat{\mathbf{X}}_1) > V(\hat{\mathbf{X}}_2)$  **then**
  - 8:      $\hat{\mathbf{X}} := \hat{\mathbf{X}}_2$
  - 9: **else**
  - 10:     $\hat{\mathbf{X}} := \arg \min_{\mathbf{X} \in \{\hat{\mathbf{X}}_1, \hat{\mathbf{X}}_2\}} \|\hat{\mathbf{X}}_{\text{map}} - \mathbf{X}\|_F$
  - 11: **end if**
  - 12:  $\hat{\mathbf{P}} = (\mathbf{C}_x + (\hat{\mathbf{X}} - \mathbf{U})(\hat{\mathbf{X}} - \mathbf{U})^\top) / \gamma_x$
  - 13:  $\hat{\mathbf{R}} = (\mathbf{C}_w + (\mathbf{Y} - \mathbf{H}\hat{\mathbf{X}})(\mathbf{Y} - \mathbf{H}\hat{\mathbf{X}})^\top) / \gamma_w$
  - 14: Output:  $\hat{\mathbf{X}}$ ,  $\hat{\mathbf{P}}$  and  $\hat{\mathbf{R}}$
- 

#### C. Marginalized MAP

In certain applications the covariance matrices  $\mathbf{P}$  and  $\mathbf{R}$  may not be of interest and can be treated as nuisance parameters that are marginalized out from the prior and likelihood pdfs. Utilizing the conjugacy of the inverse-Wishart distribution to the Gaussian distribution,

$$p(\mathbf{X}) = \int p(\mathbf{X}|\mathbf{P})p(\mathbf{P})d\mathbf{P} \propto \left| \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top + \mathbf{C}_x \right|^{-(\nu_x + N)/2}$$

and

$$p(\mathbf{Y}|\mathbf{X}) = \int p(\mathbf{Y}|\mathbf{X}, \mathbf{R})p(\mathbf{R})d\mathbf{R} \propto \left| \tilde{\mathbf{Y}}\tilde{\mathbf{Y}}^\top + \mathbf{C}_w \right|^{-(\nu_w + N)/2}.$$

Then taking the negative logarithm of the marginalized pdf,  $p(\mathbf{X}|\mathbf{Y}) \propto p(\mathbf{Y}|\mathbf{X})p(\mathbf{X})$ , results in a cost function of the same

form as (8) and the marginalized MAP estimator is given by

$$\hat{\mathbf{X}}_{\text{mmap}} = \arg \min_{\mathbf{X} \in \mathbb{R}^{n \times N}} V'(\mathbf{X}) \quad (13)$$

where

$$V'(\mathbf{X}) \triangleq \frac{\gamma'_w}{2} \ln |\mathbf{I}_N + (\mathbf{Y} - \mathbf{H}\mathbf{X})^\top \mathbf{C}_w^{-1} (\mathbf{Y} - \mathbf{H}\mathbf{X})| \\ + \frac{\gamma'_x}{2} \ln |\mathbf{I}_N + (\mathbf{X} - \mathbf{U})^\top \mathbf{C}_x^{-1} (\mathbf{X} - \mathbf{U})|,$$

and the weights are  $\gamma'_w = \nu_w + N$  and  $\gamma'_x = \nu_x + N$ . Thus we can apply the same solution methods as used for CMAP but with different weights.

#### D. Variational MAP

We note that the sought variables follow the conditional distributions:  $[\mathbf{X}]_i | \mathbf{P}, \mathbf{R}, \mathbf{Y} \sim \mathcal{N}([\mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H} + \mathbf{P}^{-1}]^{-1} [\mathbf{H}^\top \mathbf{R}^{-1} \mathbf{Y} + \mathbf{P}^{-1} \mathbf{U}]_i, (\mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H} + \mathbf{P}^{-1})^{-1})$ ,  $\mathbf{P} | \mathbf{X}, \mathbf{R}, \mathbf{Y} \sim \mathcal{W}^{-1}(\mathbf{C}_x + (\mathbf{X} - \mathbf{U})(\mathbf{X} - \mathbf{U})^\top, \nu_x + N)$  and  $\mathbf{R} | \mathbf{X}, \mathbf{P}, \mathbf{Y} \sim \mathcal{W}^{-1}(\mathbf{C}_w + (\mathbf{Y} - \mathbf{H}\mathbf{X})(\mathbf{Y} - \mathbf{H}\mathbf{X})^\top, \nu_w + N)$ , where  $[\mathbf{X}]_i$  denotes the  $i$ th column of  $\mathbf{X}$ . This enables a numerical computation of the mean of the posterior pdf  $p(\mathbf{X}, \mathbf{P}, \mathbf{R} | \mathbf{Y})$  in (4) by means of Markov Chain Monte Carlo methods, e.g., Gibbs sampling [7]. Whilst the posterior mean is the MSE-optimal estimate, the dimensionality of the problem requires a very large number of samples for accurate computation, rendering the sampling methods intractable. For completeness we consider a variational approximation of the posterior pdf [25], and derive the corresponding MAP estimator. The solution to this approximated problem results in an iteration that converges to a local minimum.

The pdf  $p(\mathbf{X}, \mathbf{P}, \mathbf{R} | \mathbf{Y})$  is approximated by conditionally independent pdfs  $q(\mathbf{X} | \mathbf{Y})q(\mathbf{P} | \mathbf{Y})q(\mathbf{R} | \mathbf{Y})$ . The distributions that minimize the Kullback-Leibler divergence to  $p(\mathbf{X}, \mathbf{P}, \mathbf{R} | \mathbf{Y})$  are given by [7]

$$q(\mathbf{X} | \mathbf{Y}) \propto e^{\mathbb{E}_{P, R | Y} [\ln p(\mathbf{X}, \mathbf{P}, \mathbf{R}, \mathbf{Y})]} \\ q(\mathbf{P} | \mathbf{Y}) \propto e^{\mathbb{E}_{X, R | Y} [\ln p(\mathbf{X}, \mathbf{P}, \mathbf{R}, \mathbf{Y})]} \\ q(\mathbf{R} | \mathbf{Y}) \propto e^{\mathbb{E}_{X, P | Y} [\ln p(\mathbf{X}, \mathbf{P}, \mathbf{R}, \mathbf{Y})]}. \quad (14)$$

Using the chain rule and introducing  $\mathbf{V} = \mathbb{E}_{P | Y} [\mathbf{P}^{-1}]$  and  $\mathbf{W} = \mathbb{E}_{R | Y} [\mathbf{R}^{-1}]$  for notational simplicity, we have

$$\ln q(\mathbf{X} | \mathbf{Y}) = \mathbb{E}_{P, R | Y} [\ln p(\mathbf{X}, \mathbf{P}, \mathbf{R}, \mathbf{Y})] + K_1 \\ = \mathbb{E}_{P, R | Y} [\ln p(\mathbf{Y} | \mathbf{X}, \mathbf{R}) + \ln p(\mathbf{X} | \mathbf{P})] + K_2 \\ = -\frac{1}{2} \text{tr}\{(\mathbf{Y} - \mathbf{H}\mathbf{X})^\top \mathbb{E}_{R | Y} [\mathbf{R}^{-1}] (\mathbf{Y} - \mathbf{H}\mathbf{X})\} \\ - \frac{1}{2} \text{tr}\{(\mathbf{X} - \mathbf{U})^\top \mathbb{E}_{P | Y} [\mathbf{P}^{-1}] (\mathbf{X} - \mathbf{U})\} + K_3 \\ = -\frac{1}{2} \text{tr}\{\mathbf{Y}^\top \mathbf{W} \mathbf{Y} - \mathbf{Y}^\top \mathbf{W} \mathbf{H} \mathbf{X} - \mathbf{X}^\top \mathbf{H}^\top \mathbf{W} \mathbf{Y} \\ + \mathbf{X}^\top \mathbf{H}^\top \mathbf{W} \mathbf{H} \mathbf{X} + \mathbf{X}^\top \mathbf{V} \mathbf{X} \\ - \mathbf{X}^\top \mathbf{V} \mathbf{U} - \mathbf{U}^\top \mathbf{V} \mathbf{X} + \mathbf{U}^\top \mathbf{V} \mathbf{U}\} + K_4 \\ = -\frac{1}{2} \text{tr}\{\mathbf{X}^\top (\mathbf{H}^\top \mathbf{W} \mathbf{H} + \mathbf{V}) \mathbf{X} \\ - (\mathbf{Y}^\top \mathbf{W} \mathbf{H} + \mathbf{U}^\top \mathbf{V}) \mathbf{X} \\ - \mathbf{X}^\top (\mathbf{H}^\top \mathbf{W} \mathbf{Y} + \mathbf{V} \mathbf{U})\} + K_5 \\ = -\frac{1}{2} \text{tr}\{(\mathbf{X} - \tilde{\mathbf{U}})^\top \tilde{\mathbf{P}}^{-1} (\mathbf{X} - \tilde{\mathbf{U}})\} + K_6,$$

which equals the functional form of  $N$  independent Gaussians with mean and covariance

$$\tilde{\mathbf{U}} = (\mathbf{H}^\top \mathbb{E}_{R | Y} [\mathbf{R}^{-1}] \mathbf{H} + \mathbb{E}_{P | Y} [\mathbf{P}^{-1}])^{-1} \\ \times (\mathbf{H}^\top \mathbb{E}_{R | Y} [\mathbf{R}^{-1}] \mathbf{Y} + \mathbb{E}_{P | Y} [\mathbf{P}^{-1}]^{-1} \mathbf{U}) \quad (15) \\ \tilde{\mathbf{P}} = (\mathbf{H}^\top \mathbb{E}_{R | Y} [\mathbf{R}^{-1}] \mathbf{H} + \mathbb{E}_{P | Y} [\mathbf{P}^{-1}])^{-1}.$$

The mean and mode coincide and the variational MAP estimator equals

$$\hat{\mathbf{X}}_{\text{vmap}} = \arg \max_{\mathbf{X} \in \mathbb{R}^{n \times N}} q(\mathbf{X} | \mathbf{Y}) = \tilde{\mathbf{U}}. \quad (16)$$

Further,

$$\ln q(\mathbf{P} | \mathbf{Y}) = \mathbb{E}_{X, R | Y} [\ln p(\mathbf{X}, \mathbf{P}, \mathbf{R}, \mathbf{Y})] + K_1 \\ = \mathbb{E}_{X, R | Y} [\ln p(\mathbf{X} | \mathbf{P}) + \ln p(\mathbf{P})] + K_2 \\ = \mathbb{E}_{X, R | Y} \left[ -\frac{\nu_x + N + n + 1}{2} \ln |\mathbf{P}| \right. \\ \left. - \frac{1}{2} \text{tr}\{(\mathbf{C}_x + \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top) \mathbf{P}^{-1}\} \right] + K_3 \\ = -\frac{\gamma'_x + n + 1}{2} \ln |\mathbf{P}| - \frac{1}{2} \text{tr}\{\tilde{\mathbf{C}}_x \mathbf{P}^{-1}\} + K_4.$$

This is the functional form of an inverse-Wishart with parameters  $\gamma'_x = \nu_x + N$  and

$$\tilde{\mathbf{C}}_x = \mathbf{C}_x + \mathbb{E}_{X, R | Y} [(\mathbf{X} - \mathbf{U})(\mathbf{X} - \mathbf{U})^\top] \\ = \mathbf{C}_x + N \tilde{\mathbf{P}} + N \tilde{\mathbf{U}} \tilde{\mathbf{U}}^\top - \tilde{\mathbf{U}} \mathbf{U}^\top - \mathbf{U} \tilde{\mathbf{U}}^\top + \mathbf{U} \mathbf{U}^\top. \quad (17)$$

Thus  $\mathbf{P}^{-1}$  follows a Wishart distribution and  $\mathbb{E}_{P | Y} [\mathbf{P}^{-1}] = \gamma'_x \tilde{\mathbf{C}}_x^{-1}$ . Similarly,

$$\ln q(\mathbf{R} | \mathbf{Y}) = \mathbb{E}_{X, R | Y} [\ln p(\mathbf{X}, \mathbf{P}, \mathbf{R}, \mathbf{Y})] + K_1 \\ = -\frac{\gamma'_w + m + 1}{2} \ln |\mathbf{R}| - \frac{1}{2} \text{tr}\{\tilde{\mathbf{C}}_w \mathbf{R}^{-1}\} + K_2$$

has the functional form of an inverse-Wishart distribution with parameters  $\gamma'_w = \nu_w + N$  and

$$\tilde{\mathbf{C}}_w = \mathbf{C}_w + \mathbb{E}_{X, P | Y} [(\mathbf{Y} - \mathbf{H}\mathbf{X})(\mathbf{Y} - \mathbf{H}\mathbf{X})^\top] \\ = \mathbf{C}_w + \mathbf{Y} \mathbf{Y}^\top - \mathbf{Y} \tilde{\mathbf{U}}^\top \mathbf{H}^\top - \mathbf{H} \tilde{\mathbf{U}} \mathbf{Y}^\top \\ + N \mathbf{H} \tilde{\mathbf{P}} \mathbf{H}^\top + N \mathbf{H} \tilde{\mathbf{U}} \tilde{\mathbf{U}}^\top \mathbf{H}^\top. \quad (18)$$

Thus  $\mathbf{R}^{-1}$  follows a Wishart distribution and  $\mathbb{E}_{R | Y} [\mathbf{R}^{-1}] = \gamma'_w \tilde{\mathbf{C}}_w^{-1}$ .

Inserting these results into (16) the variational MAP estimator is computed iteratively as

$$\hat{\mathbf{X}}_{\text{vmap}} = (\gamma'_w \mathbf{H}^\top \tilde{\mathbf{C}}_w^{-1} \mathbf{H} + \gamma'_x \tilde{\mathbf{C}}_x^{-1})^{-1} \\ \times (\gamma'_w \mathbf{H}^\top \tilde{\mathbf{C}}_w^{-1} \mathbf{Y} + \gamma'_x \tilde{\mathbf{C}}_x^{-1} \mathbf{U}).$$

The parameters  $\tilde{\mathbf{C}}_x$  and  $\tilde{\mathbf{C}}_w$  are subsequently updated using (17) and (18). The iteration is initialized by setting the parameters  $\tilde{\mathbf{C}}_x = \mathbf{C}_x$  and  $\tilde{\mathbf{C}}_w = \mathbf{C}_w$ . Experimentally we find that using more informative initialization points, i.e., initializing (17) and (18) with  $\tilde{\mathbf{U}} = \hat{\mathbf{X}}_{\text{mvu}}$  and  $\mathbf{P} = \mathbf{P}_0$  produces virtually identical results.

#### IV. EXPERIMENTAL RESULTS

In this section we compare the statistical performance of  $\widehat{\mathbf{X}}_{\text{cmap}}$  with other estimators using the distribution of normalized squared errors,  $\text{NSE} \triangleq \|\mathbf{X} - \widehat{\mathbf{X}}\|_F^2 / \mathbb{E}[\|\mathbf{X}\|_F^2]$ . The expectation is over all random variables. In particular, we will use the normalized mean square error  $\text{NMSE} \equiv \mathbb{E}[\text{NSE}]$  and the complementary cumulative distribution function (ccdf),  $\Pr\{\text{NSE} > \kappa\}$ . The former measures the average performance of the estimators and the latter quantifies their robustness to noise and covariance uncertainties.

We also evaluate the NMSE of the covariance matrix estimates  $\widehat{\mathbf{P}}$  and  $\widehat{\mathbf{R}}$  in comparison with the nominal matrices  $\mathbf{P}_0$  and  $\mathbf{R}_0$ .

The statistical measures are estimated by means of Monte Carlo simulations.

##### A. Estimators

We compare  $\widehat{\mathbf{X}}_{\text{cmap}}$  with  $\widehat{\mathbf{X}}_{\text{map}}$  and  $\widehat{\mathbf{X}}_{\text{mvu}}$ , that use nominal covariance matrices  $\mathbf{P}_0$  and  $\mathbf{R}_0$ . For CMAP we set the tolerance parameter  $\varepsilon = 10^{-6}$ .

For comparison of robustness properties with respect to both signal and noise covariance uncertainties, we also apply the state of the art difference regret estimator (DRE) given in [12],  $\widehat{\mathbf{X}}_{\text{dre}}$ . This estimator assumes that  $\mathbf{X}$  is zero mean and is derived on assumption that the spatial correlations of the signal and noise are structured by the singular vectors of  $\mathbf{H}$ . Nevertheless, in [12] it is suggested that the estimator can be implemented whether or not this correlation structure is satisfied. The covariance uncertainties are treated deterministically as bounds on the eigenvalues of  $\mathbf{P}_0$ , i.e.,  $l_i^x \leq \lambda_i^x \leq u_i^x$  for  $i = 1, \dots, n$ , and  $\mathbf{R}_0$ , i.e.,  $l_j^w \leq \lambda_j^w \leq u_j^w$  for  $j = 1, \dots, m$ .

The estimator has the form

$$\widehat{\mathbf{X}}_{\text{dre}} = \mathbf{D}_x \mathbf{H}^\top (\mathbf{H} \mathbf{D}_x \mathbf{H}^\top + \mathbf{D}_w)^{-1} \mathbf{Y}. \quad (19)$$

The input covariance matrices are set as  $\mathbf{D}_x = \mathbf{V} \mathbf{\Delta}_x \mathbf{V}^\top$  and  $\mathbf{D}_w = \mathbf{W} \mathbf{\Delta}_w \mathbf{W}^\top$ , where  $\mathbf{V}$  and  $\mathbf{W}$  are eigenvector matrices of  $\mathbf{P}_0$  and  $\mathbf{R}_0$ , respectively. Further,  $\mathbf{\Delta}_x = \text{diag}(\delta_1^x, \dots, \delta_n^x)$  and  $\mathbf{\Delta}_w = \text{diag}(\delta_1^w, \dots, \delta_m^w)$ , where

$$\begin{aligned} \delta_i^x &= \alpha_i l_i^x + (1 - \alpha_i) u_i^x, & i = 1, \dots, n \\ \delta_i^w &= \alpha_i l_i^w + (1 - \alpha_i) u_i^w, & i = 1, \dots, m \end{aligned}$$

and  $\delta_i^w = \lambda_i^w$  for all  $i = n + 1, \dots, m$ . Here

$$\alpha_i = \frac{\sqrt{l_i^w + u_i^x \sigma_i^2}}{\sqrt{l_i^w + u_i^x \sigma_i^2} + \sqrt{u_i^w + l_i^x \sigma_i^2}}, \quad (20)$$

where  $\sigma_i$  are the singular values of  $\mathbf{H}$ .

Since the covariance uncertainties are treated probabilistically in this work, selecting deterministic bounds on the eigenvalues can only be done heuristically. Here we have selected,  $l_i = (1 - \nu^0/\nu)\lambda_i$  and  $u_i = (1 + \nu^0/\nu)\lambda_i$ , where  $\nu^0$  denotes the minimum integer value of  $\nu$ , i.e.,  $\nu_x^0 = n + 2$  and  $\nu_w^0 = m + 2$ . Thus with minimum certainty of the covariances, the lower bound is 0 and upper bound is  $2\lambda_i$ . As  $\nu \rightarrow \infty$ , the bounds become tight.

##### B. Signal setup

For sake of illustration, we consider the problem of estimating a stochastic  $2 \times 2$  multiple-input multiple output (MIMO) channel  $\mathbf{A} \in \mathbb{R}^{2 \times 2}$  from observed signals  $\mathbf{z}_k = \mathbf{A} \mathbf{s}_k + \mathbf{n}_k$ . As is common in wireless communications, this is achieved by transmitting a known training sequence  $\mathbf{S} = [\mathbf{s}_1 \cdots \mathbf{s}_K]$  [26], [27]. Collecting  $K$  snapshots and vectorizing, the observation is rewritten as  $\mathbf{y} = \mathbf{H} \mathbf{x} + \mathbf{w}$ , where  $\mathbf{x} = \text{vec}(\mathbf{A})$  and  $\mathbf{H} = (\mathbf{S}^\top \otimes \mathbf{I}_2)$ . The vectorized channel coefficients  $\mathbf{x}$  and noise  $\mathbf{w}$  follow independent, zero-mean, conditionally Gaussian distributions.  $\mathbf{S}$  is chosen as a deterministic white sequence with constrained power,  $\|\mathbf{S}\|_2^2 \equiv 10$ . We set  $\mathbf{P}_0 = \frac{1}{n} \mathbf{I}_n$ , and  $\mathbf{R}_0 = \sigma_w^2 \mathbf{I}_m$ . The covariance matrices are drawn according to inverse-Wishart distributions. We consider estimating  $N$  channel realizations  $\mathbf{X} \in \mathbb{R}^{n \times N}$  from observations  $\mathbf{Y} = \mathbf{H} \mathbf{X} + \mathbf{W} \in \mathbb{R}^{m \times N}$ .

The signal to noise ratio,

$$\text{SNR} \triangleq \frac{\mathbb{E}[\|\mathbf{H} \mathbf{X}\|_F^2]}{\mathbb{E}[\|\mathbf{W}\|_F^2]} = \frac{\text{tr}\{\mathbf{H} \mathbf{P}_0 \mathbf{H}^\top\}}{\text{tr}\{\mathbf{R}_0\}},$$

is varied in the experiments, i.e., setting  $\sigma_w^2 = \text{tr}\{\mathbf{H} \mathbf{H}^\top\} / (mn \times \text{SNR})$ . We consider  $K = 8$  snapshots so that  $m = 16$  and  $n = 4$ . For  $m > n$ , the resulting low-rank signal structure enables CMAP to estimate parts of both covariances. When  $m = n$ , the loss of parameter identifiability makes CMAP rely less on the prior signal statistics at higher SNR levels, thus performing closer to the MVU estimator.

Throughout the experiments we ran  $10^5$  Monte Carlo simulations for each signal setup.

##### C. Results for single observation

In the following experiments we consider  $N = 1$ . First, the average performance of the estimators are compared. Fig. 2 shows the NMSE as a function of SNR when the degrees of freedom for  $\mathbf{P}$  and  $\mathbf{R}$  are set to their minimum integer values,  $\nu_x^0 = n + 2$  and  $\nu_w^0 = m + 2$ , respectively. This yields the minimum certainties of the random quantities. CMAP is capable of reducing the NMSE by up to approximately 2 dB compared to the standard MAP. As SNR increases, MAP converges faster to MVU than does CMAP. The average performance of DRE is initially similar to MAP but the gap increases with SNR as it injects a larger bias.

Next, the statistical performance of the estimators is compared using the ccdf,  $\Pr\{\text{NSE} > \kappa\}$ , at SNR=0 dB. The curves in Fig. 3 illustrate the relative robustness of the estimators to covariance uncertainties. Estimators that produce a lower fraction of poor estimates will have lower ccdfs. Note that  $\text{NSE} > 1$  are estimates that have errors greater than the average NSE of the mean,  $\widehat{\mathbf{X}} = \mathbf{0}$ . As expected, MAP and DRE perform similarly at this SNR level, while MVU is slightly worse but declines at a similar rate. CMAP declines more rapidly, with a ccdf that is approximately one order of magnitude lower than MVU at  $\kappa = 10$ .

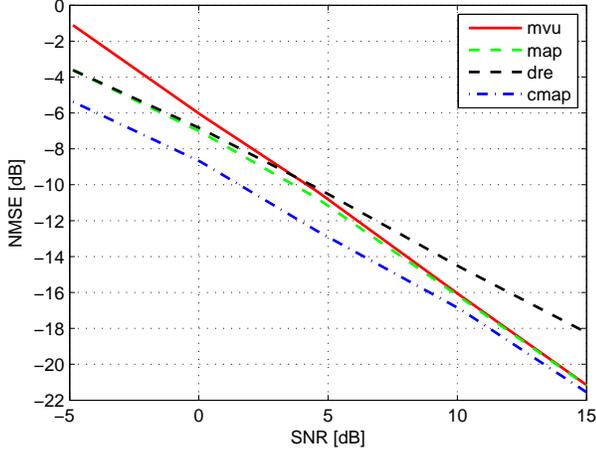


Fig. 2. NMSE versus SNR with minimum certainties of  $\mathbf{P}$  and  $\mathbf{R}$ .  $N = 1$ .

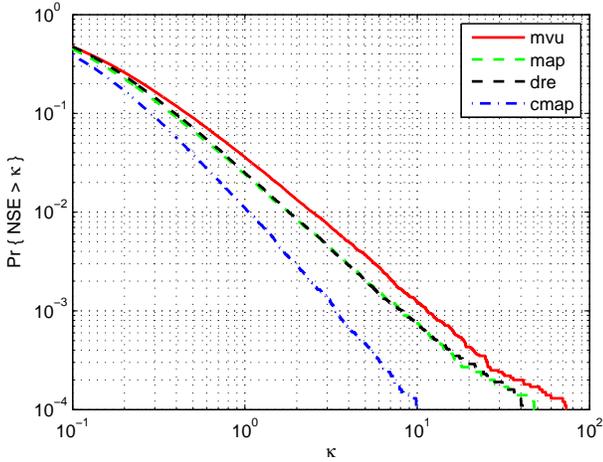


Fig. 3. Ccdf of NSE at SNR=0 dB, with  $\nu_x = \nu_x^0$  and  $\nu_w = \nu_w^0$ .  $N = 1$ .

#### D. Results for multiple observations

The previous experiments are repeated for  $N = 4$ . We use the Gibbs sampling approximation of the posterior mean which provides a bound on the NMSE, cf. Fig. 4. In this scenario we see that CMAP is very close to the optimum. When  $m > N$ , the computational complexity of the Gibbs sampler and CMAP is of the order  $\mathcal{O}(m^3 N_{\text{iter}})$ , where  $m^3$  is the complexity of matrix inversion and  $N_{\text{iter}}$  is the number of repetitions. For the Gibbs sampler,  $N_{\text{iter}} = 2 \times 10^4$  is about 100 times the number of parameters to estimate and provides a good approximation of the mean. For CMAP, the expected number of iterations is approximately three orders of magnitude less, cf. Sec. IV-E.

Further, we vary the certainties of the covariance matrices by setting  $\nu$  to the extremes,  $\nu^0$  and  $\infty$ . (For  $\infty$ , we set  $\nu$  numerically to  $10^5$ .) The relative difference in average performance between CMAP and MAP is denoted by  $\Delta \text{NMSE}$ , where a negative value means reduction in NMSE in decibel using CMAP. The results are shown in Fig. 5. When  $(\nu_x, \nu_w) = (\infty, \infty)$ , CMAP is identical to MAP but for  $(\nu_x, \nu_w) = (\nu_x^0, \infty)$  the estimator relies primarily on

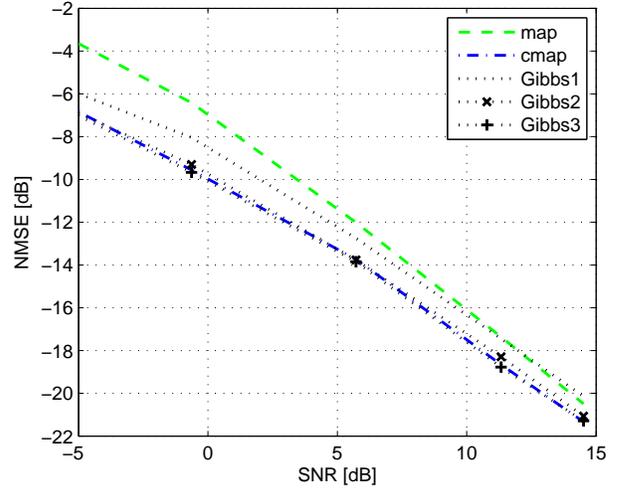


Fig. 4. NMSE versus SNR, with  $\nu_x = \nu_x^0$  and  $\nu_w = \nu_w^0$ .  $N = 4$ . Gibbs 1, 2 and 3 use 200, 2 000 and 20 000 samples, respectively.

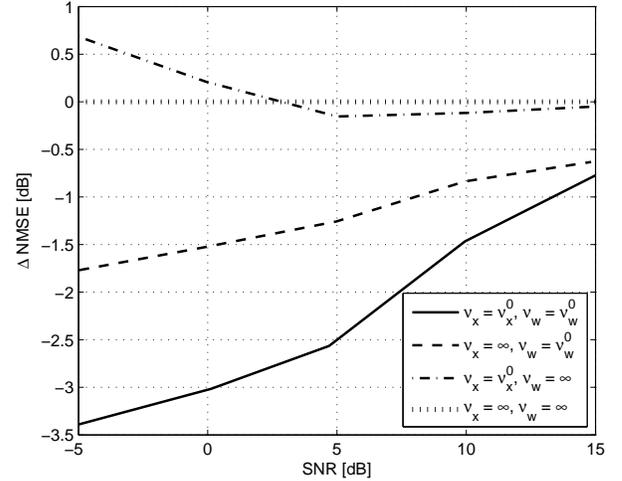


Fig. 5. Difference NMSE between  $\hat{\mathbf{X}}_{\text{cmap}}$  and  $\hat{\mathbf{X}}_{\text{map}}$  versus SNR, for different certainties of  $\mathbf{P}$  and  $\mathbf{R}$ .  $N = 4$ .

the noise statistics and CMAP approaches MVU. As both covariance matrices become less certain  $(\nu_x, \nu_w) \rightarrow (\nu_x^0, \nu_w^0)$ , the advantage of CMAP increases, illustrated by the dashed and solid lines. The improvement for  $N = 4$  snapshots is above 3 dB for low SNR.

Next, the statistical performance is assessed using the cdf at SNR=0 dB. Figs. 6, 7 and 8 illustrate robustness at various covariance uncertainties. A comparison between Fig. 3 and 6 shows how the cdf of CMAP is reduced when the number of samples increases from  $N = 1$  to 4. When CMAP relies primarily on the noise statistics, as in Fig. 7, it tends towards MVU. While the average NSE of CMAP rises slightly above MAP in this case at low SNR (Fig. 5), its cdf exhibits a sharp decline relative to MAP. When only the signal statistics are reliable, the differences in decline are more pronounced, see Fig. 8. In all three cases the fraction of poor estimates cuts

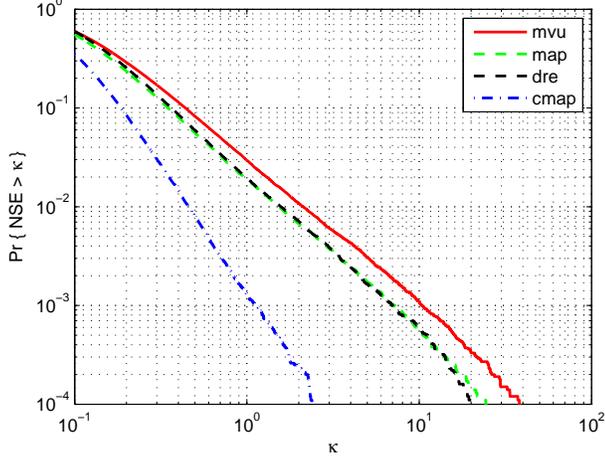


Fig. 6. Ccdf of NSE at SNR=0 dB, with  $\nu_x = \nu_x^0$  and  $\nu_w = \nu_w^0$ .  $N = 4$ .

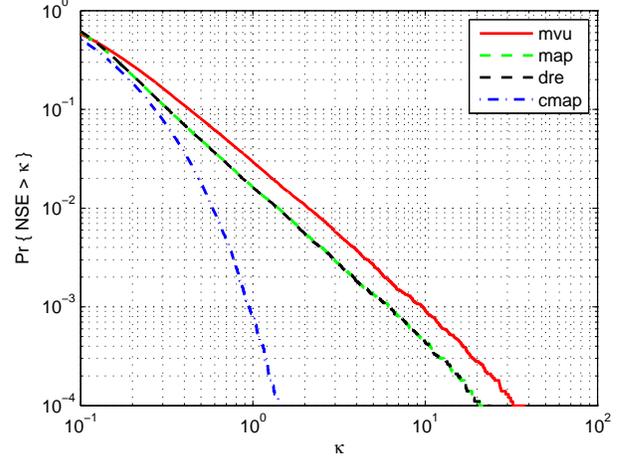


Fig. 8. Ccdf of NSE at SNR=0 dB, with  $\nu_x = \infty$  and  $\nu_w = \nu_w^0$ .  $N = 4$ .

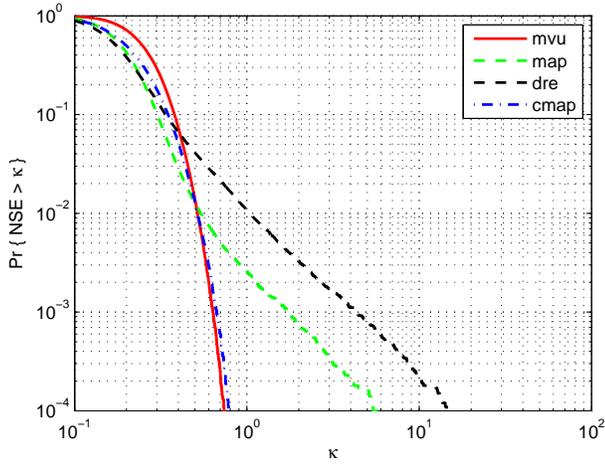


Fig. 7. Ccdf of NSE at SNR=0 dB, with  $\nu_x = \nu_x^0$  and  $\nu_w = \infty$ .  $N = 4$ .

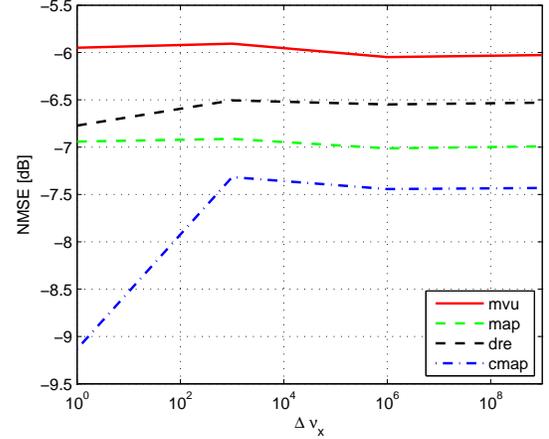


Fig. 10. NMSE versus  $\Delta\nu_x$  at SNR=0 dB.  $N = 4$ .

off faster for CMAP than MAP.

Further, we investigate the estimation errors of the covariance matrix estimates  $\hat{\mathbf{P}}$  and  $\hat{\mathbf{R}}$ . More specifically, we compute the difference NMSE between the estimates and the priors, which quantifies the information gain, as  $N$  increases. The results are shown in Fig. 9 for various SNR levels. Note that there is a measurable gain even at  $N < n$  and  $N < m$ . Thus CMAP is useful also as a covariance estimator for applications in which signal statistics are of importance.

In practical scenarios with uncertain covariances, CMAP would be implemented with the minimal integer values  $\nu_x^0$  and  $\nu_w^0$ . We now investigate the effect of mismatches from this conservative prior knowledge by setting the true values to  $\nu^0 + \Delta\nu$ . In Figures 10 and 11 we increase  $\Delta\nu_x$  and  $\Delta\nu_w$ , respectively. At SNR=0 dB, we see increases in NMSE for CMAP but the advantage of the estimator is still robust with respect to mismatches for either distributions of  $\mathbf{P}$  or  $\mathbf{R}$ .

Finally, we investigate how CMAP performs when increasing the signal dimensions. We now set  $m = 64$  and  $n = 16$ , for SNR=0 dB, with  $\nu_x = \nu_x^0$ ,  $\nu_w = \nu_w^0$  and  $N = 16$ . The NMSE

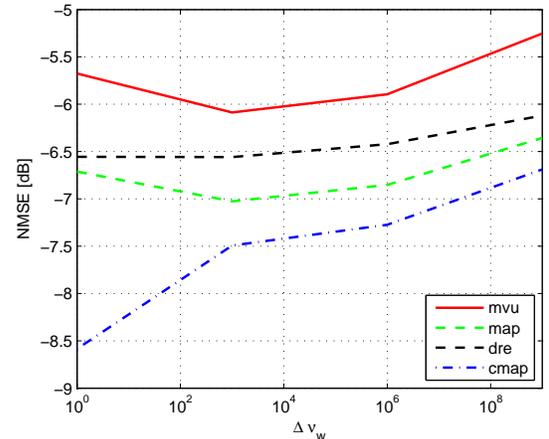


Fig. 11. NMSE versus  $\Delta\nu_w$  at SNR=0 dB.  $N = 4$ .

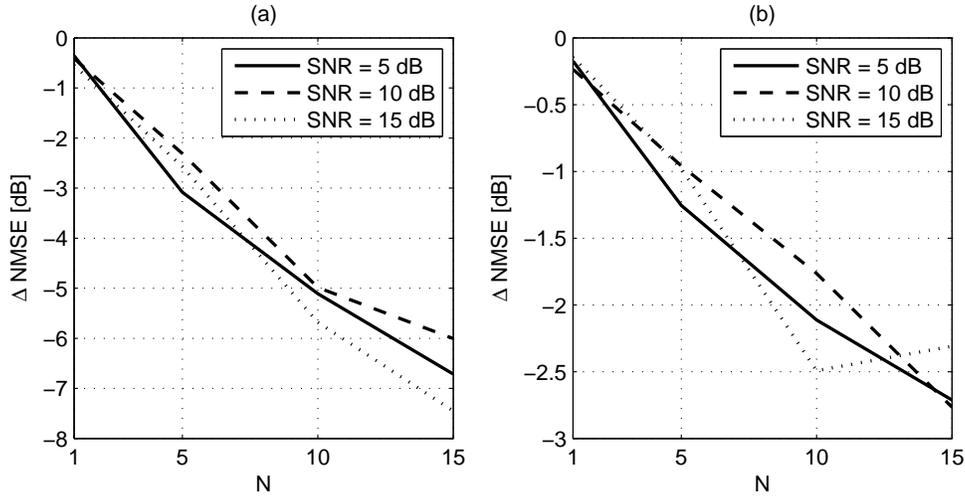


Fig. 9. Difference NMSE between (a)  $\hat{\mathbf{P}}$  and  $\mathbf{P}_0$  and (b)  $\hat{\mathbf{R}}$  and  $\mathbf{R}_0$  versus  $N$ . Here  $\nu_x = \nu_x^0$  and  $\nu_w = \nu_w^0$ .

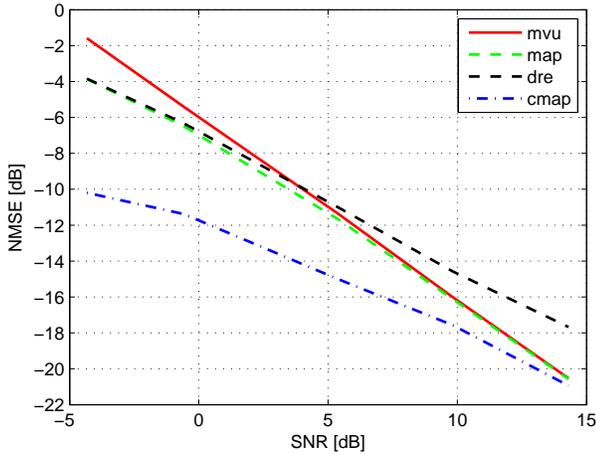


Fig. 12. NMSE versus SNR for  $m = 64$  and  $n = 16$ .  $N = 16$ .

performance as a function of SNR is illustrated in Fig. 12 which shows a gain of CMAP over MAP greater than in the setup considered in Fig. 4, where  $m = 16$ ,  $m = 4$  and  $N = 4$ .

### E. Empirical convergence properties

We now turn to the convergence properties of the iterative solution (12) of CMAP for the same scenario as considered in the previous section, i.e., SNR=0 dB, with  $\nu_x = \nu_x^0$ ,  $\nu_w = \nu_w^0$  and  $N = 4$ . Fig. 13 shows a comparison of the convergence rate of the fixed-point iteration and the gradient descent solution, for a typical realization. Both solutions exhibit similar rates once the estimates are sufficiently close to a minimum as both are based on the gradient. But the fixed-point iteration reaches this region within a few iterations without the need for a user-defined step length.

Next, we study the statistical convergence properties. Let  $N_{\text{iter}}$  denote the total number of iterations until (12) fulfills  $\|\hat{\mathbf{X}}^{\ell+1} - \hat{\mathbf{X}}^\ell\|_F < \varepsilon$ . Then we can estimate the ccdf  $\Pr\{N_{\text{iter}} >$

$k\}$ , as displayed in Fig. 14. When  $\varepsilon = 10^{-6}$  we see that the probability of  $N_{\text{iter}}$  exceeding 300 iterations is less than  $10^{-3}$ , and the mean of  $N_{\text{iter}}$  is 24.7. For  $\varepsilon = 10^{-3}$ , the mean of  $N_{\text{iter}}$  is reduced by more than a half, while the NMSE is virtually the same. For  $\varepsilon = 10^{-1}$ , the entire ccdf is substantially reduced while incurring an increase in NMSE of only 0.24 dB.

We also estimate the proportion of instances in which the fixed-point iteration converges to two different minima starting from  $\hat{\mathbf{X}}_{\text{mvu}}$  and  $\mathbf{U}$ , respectively. We quantify this as when the convergence points,  $\hat{\mathbf{X}}_1$  and  $\hat{\mathbf{X}}_2$ , differ substantially from the numerical tolerance, i.e.,  $\Pr\{\|\hat{\mathbf{X}}_1 - \hat{\mathbf{X}}_2\|_F > 10^{-2} \times nN\}$ . For the given scenario, the probability was estimated to 0.03. In 98% out of those instances  $\hat{\mathbf{X}}_1$  produced a lower cost  $V(\mathbf{X})$  than  $\hat{\mathbf{X}}_2$ .

In all our simulations we did never encounter a single case when the fixed-point iteration failed to converge to the tolerance.

### F. Comparison between alternative MAP estimators

Finally, we compare a scenario in which the covariance matrices are not of interest and can be marginalized out, resulting in the marginalized MAP estimator (13) using the same initial points as CMAP. The difference between the estimators is marginal in terms of NSE performance (see Fig. 15). The NMSE is marginally better for MMAP as it estimates fewer parameters than CMAP;  $-9.98$  and  $-9.94$  dB for MMAP and CMAP, respectively. The variational MAP estimator performs better than the standard MAP but is inferior to MMAP and CMAP. The NMSE is  $-6.83$  and  $-8.10$  dB for MAP and VMAP, respectively.

We also evaluate the significance and robustness of the choice of starting points for CMAP and MMAP. Tests were performed using initial points  $\hat{\mathbf{X}}^0$  randomized by a Gaussian distribution with covariance  $\mathbf{P}_0$  and a given mean. For each observation  $\mathbf{Y}$  we then form 10 random initial points  $\hat{\mathbf{X}}^0$  resulting in 10 search paths. The convergence point that yields the lowest cost,  $V(\mathbf{X})$ , is retained as the estimate. We denote this randomized MAP-based estimator as ‘RMAP’.

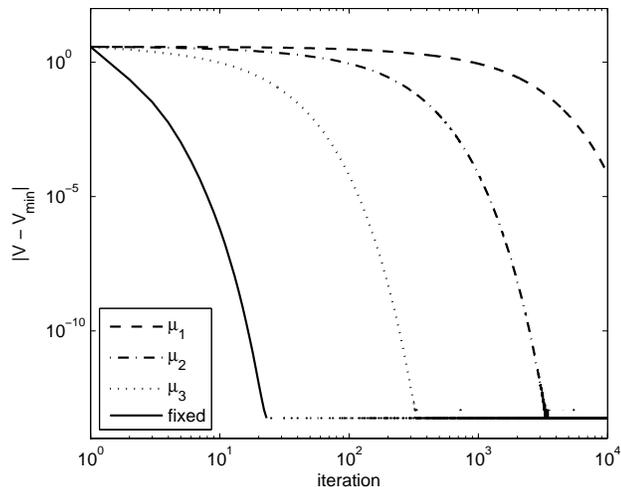


Fig. 13. Convergence to minimum  $V_{\min}$  of gradient descent and fixed-point iteration, starting from  $\hat{\mathbf{X}}_{\text{mvu}}$ . Step sizes  $\mu_1$ ,  $\mu_2$  and  $\mu_3$  were set to  $10^{-5}$ ,  $10^{-4}$  and  $10^{-3}$ , respectively. For step-size  $10^{-2}$ , the gradient descent became unstable. Based on a realization of signal setup given in section IV-B with  $N = 4$ .

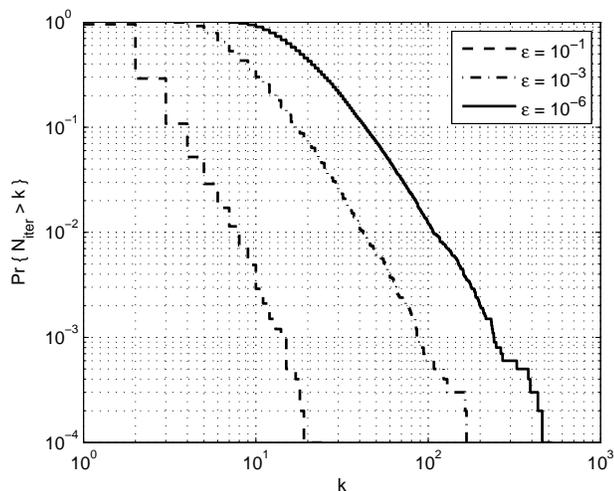


Fig. 14. Ccdf of  $N_{\text{iter}}$  at SNR=0 dB, with  $\nu_x = \nu_x^0$ ,  $\nu_w = \nu_w^0$  and  $N = 4$ . Tolerances  $\varepsilon_1 = 10^{-1}$ ,  $\varepsilon_2 = 10^{-3}$  and  $\varepsilon_3 = 10^{-6}$ . The corresponding  $E[N_{\text{iter}}]$  was estimated to 2.5, 10.3 and 24.7, respectively, and NMSE was  $-9.55$ ,  $-9.78$  and  $-9.78$  dB, respectively.

The different means tested were based on starting points for CMAP, i.e. prior mean  $\mathbf{U}$  and  $\hat{\mathbf{X}}_{\text{mvu}}$ , as well as  $\hat{\mathbf{X}}_{\text{map}}$  and  $\hat{\mathbf{X}}_{\text{dre}}$ . Randomizing  $\hat{\mathbf{X}}^0$  around the  $\hat{\mathbf{X}}_{\text{mvu}}$ , as well as the proximate values  $\hat{\mathbf{X}}_{\text{map}}$  and  $\hat{\mathbf{X}}_{\text{dre}}$ , is found to produce near identical performance to CMAP. Randomizing  $\hat{\mathbf{X}}^0$  around  $\mathbf{U}$ , on the other hand, leads to significantly reduced NSE performance for the worst estimates. These results corroborate the choice of initial points described in section III-A. Fig. 15 shows the performance for RMAP when using  $\hat{\mathbf{X}}_{\text{dre}}$  as a mean, and the NMSE equals  $-9.93$  dB.

*Reproducible research:* Code for reproducing Figs. 2, 5 and 12 is available at [www.ee.kth.se/~davez/rr-cmap](http://www.ee.kth.se/~davez/rr-cmap).

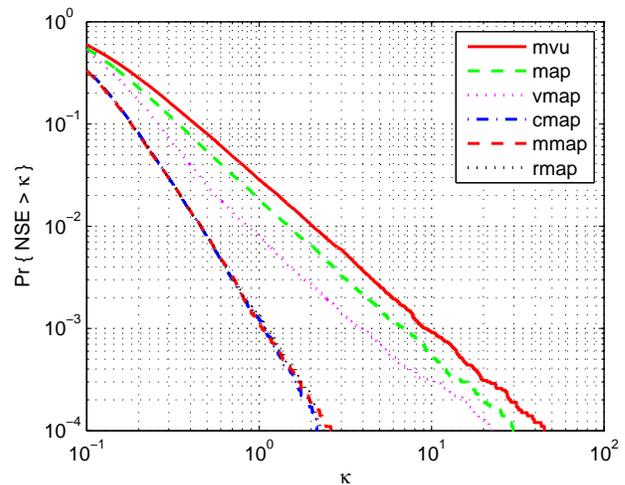


Fig. 15. Ccdf of NSE at SNR=0 dB, with  $\nu_x = \nu_x^0$  and  $\nu_w = \nu_w^0$ ,  $N = 4$ . Comparison with the marginalized estimator, ‘MMAP’, the variational estimator, ‘VMAP’, and estimator with randomized initial points ‘RMAP’.

## V. CONCLUSION

We have derived a joint signal and covariance maximum a posteriori estimator for the linear observation model, where the signal and noise covariance matrices are modeled as random quantities. We formulated a solution of the nonconvex problem as a fixed-point iterations. The resulting estimator, CMAP, exhibits robustness properties relative to the standard MAP and MVU estimators as well as the minimax difference regret estimator in low-rank signal estimation problems. In this scenario CMAP also shows near MSE-optimal performance. As the number of samples increases, the performance gains of CMAP can be quite substantial.

## REFERENCES

- [1] L. Scharf and C. Demeure, *Statistical Signal Processing: detection, estimation, and time series analysis*. Addison-Wesley Series in Electrical and Computer Engineering, Addison-Wesley Pub. Co., 1991.
- [2] S. M. Kay, *Fundamentals of Statistical Signal Processing, Vol.1—Estimation theory*. Prentice Hall, 1993.
- [3] T. Kailath, A. Sayed, and B. Hassibi, *Linear Estimation*. Prentice Hall, 2000.
- [4] C. Rao, *Linear Statistical Inference and its Applications*. Wiley Series in Probability and Statistics, Wiley, 1973.
- [5] S. J. Press, *Applied multivariate analysis—using Bayesian and frequentist methods of inference*. Dover, 2005 [1972].
- [6] C. Rao, H. Toutenburg, Salabh, and C. Heumann, *Linear Models and Generalizations: Least Squares and Alternatives*. Springer series in statistics, Springer, 2007.
- [7] C. Bishop, *Pattern Recognition and Machine Learning*. Information Science and Statistics, Springer, 2006.
- [8] K. J. Friston, A. P. Holmes, K. J. Worsley, J.-P. Poline, C. D. Frith, and R. S. Frackowiak, “Statistical parametric maps in functional imaging: a general linear approach,” *Human brain mapping*, vol. 2, no. 4, pp. 189–210, 1994.
- [9] A. Sayed, *Fundamentals of Adaptive Filtering*. Wiley, 2003.
- [10] H. Van Trees and K. Bell, *Detection Estimation and Modulation Theory, part 1*. Wiley, 2013.
- [11] S. A. Kassam and H. V. Poor, “Robust techniques for signal processing: A survey,” *Proc. IEEE*, vol. 73, pp. 433–481, Mar. 1985.
- [12] Y. C. Eldar, “Robust competitive estimation with signal and noise covariance uncertainties,” *IEEE Trans. Inf. Theory*, vol. 52, pp. 4532–4547, Oct. 2006.

- [13] Y. C. Eldar and N. Merhav, "A competitive minimax approach to robust estimation of random parameters," *IEEE Trans. Signal Processing*, vol. 52, pp. 1931–1946, July 2004.
- [14] Y. C. Eldar and N. Merhav, "Minimax MSE-ratio estimation with signal covariance uncertainties," *IEEE Trans. Signal Processing*, vol. 53, pp. 1335–1347, Apr. 2005.
- [15] S. Verdú and H. Poor, "On minimax robustness: A general approach and applications," *IEEE Trans. Inf. Theory*, vol. 30, pp. 328–340, Mar. 1984.
- [16] R. Mittleman and E. L. Miller, "Robust estimation of a random parameter in a gaussian linear model with joint eigenvalue and element-wise covariance uncertainties," *IEEE Trans. Signal Processing*, vol. 58, pp. 1001–1011, Mar. 2010.
- [17] G. C. Tiao and A. Zellner, "On the Bayesian estimation of multivariate regression," *J. Royal Statistical Soc. Series B*, vol. 26, pp. 277–285, Apr. 1964.
- [18] L. Svensson and M. Lundberg, "On posterior distributions for signals in gaussian noise with unknown covariance matrix," *IEEE Trans. Signal Processing*, vol. 53, pp. 3554–3571, Sept. 2005.
- [19] S. Bidon, O. Besson, and J.-Y. Tournéret, "A Bayesian approach to adaptive detection in nonhomogeneous environments," *IEEE Trans. Signal Processing*, vol. 56, pp. 205–217, Jan. 2008.
- [20] S. Bidon, O. Besson, and J.-Y. Tournéret, "The adaptive coherence estimator is the generalized likelihood ratio test for a class of heterogeneous environments," *IEEE Signal Process Lett.*, vol. 15, pp. 281–284, 2008.
- [21] P. Wang, H. Li, and B. Himed, "A Bayesian parametric test for multichannel adaptive signal detection in nonhomogeneous environments," *IEEE Signal Processing Lett.*, vol. 17, pp. 351–354, Apr. 2010.
- [22] D. Maiwald and D. Kraus, "Calculation of moments of complex Wishart and complex inverse Wishart distributed matrices," *IEE Proc.: Radar, Sonar and Nav.*, vol. 147, pp. 162–168, Aug. 2000.
- [23] H. Hager, "Updating the inverse of a matrix," *SIAM Rev.*, vol. 32, no. 2, pp. 221–239, 1989.
- [24] B. Hasselblatt and A. Katok, *A First Course in Dynamics: with a Panorama of Recent Developments*. Cambridge University Press, 2003.
- [25] V. Šmídl and A. Quinn, *The Variational Bayes Method in Signal Processing*. Signals and Communication Technology, Springer, 2010.
- [26] J. Kotecha and A. Sayeed, "Transmit signal design for optimal estimation of correlated MIMO channels," *IEEE Trans. Signal Processing*, vol. 52, pp. 546–557, Feb. 2004.
- [27] E. Björnson and B. Ottersten, "A framework for training-based estimation in arbitrarily correlated Rician MIMO channels with Rician disturbance," *IEEE Trans. Signal Processing*, vol. 58, pp. 1807–1820, Mar. 2010.