

Efficient Linear Fusion of Partial Estimators

David Luengo^a, Luca Martino^b, Víctor Elvira^c, Mónica Bugallo^d

^a*Dep. of Signal Theory and Communications, Universidad Politécnica Madrid (Spain)*

^b*Image Processing Lab., Universitat de València (Spain)*

^c*IMT Lille Douai & CRISTAL (UMR CNRS 9189), Villeneuve d'Ascq (France)*

^d*Dep. of Electrical and Computer Eng., Stony Brook University, NY (USA)*

Abstract

Many signal processing applications require performing statistical inference on large datasets, where computational and/or memory restrictions become an issue. In this big data setting, computing an exact global centralized estimator is often either unfeasible or impractical. Hence, several authors have considered distributed inference approaches, where the data are divided among multiple workers (cores, machines or a combination of both). The computations are then performed in parallel and the resulting partial estimators are finally combined to approximate the intractable global estimator. In this paper, we focus on the scenario where no communication exists among the workers, deriving efficient linear fusion rules for the combination of the distributed estimators. Both a constrained optimization perspective and a Bayesian approach (based on the Bernstein-von Mises theorem and the asymptotic normality of the estimators) are provided for the derivation of the proposed linear fusion rules. We concentrate on finding the minimum mean squared error (MMSE) global estimator, but the developed framework is very general and can be used to combine any type of unbiased partial estimators (not necessarily MMSE partial estimators). Numerical results show the good performance of the algorithms developed, both in problems where analytical expressions can be obtained for the partial estima-

*Corresponding author: David Luengo

Email addresses: david.luengo@upm.es (David Luengo), luca.martino@uv.es (Luca Martino), victor.elvira@telecom-lille.fr (Víctor Elvira), monica.bugallo@stonybrook.edu (Mónica Bugallo)

tors, and in a wireless sensor network localization problem where Monte Carlo methods are used to approximate the partial estimators.

Keywords: distributed estimation, linear fusion, constrained minimization, big data, minimum mean squared error (MMSE) estimators, statistical inference, Bayesian estimation, Bernstein-von Mises theorem, Monte Carlo methods

1. Introduction

Estimation theory addresses the problem of inferring a set of unknown variables of interest given a collection of available data [1, 2]. This is a central problem in statistical signal processing, where a parametric model for the data is often
5 assumed and its parameters have to be inferred from the observations [3, 4, 5]. Indeed, even non-parametric approaches typically have a reduced set of hyper-parameters that need to be estimated from the data [6, 7, 8]. Unfortunately, determining the *global estimator* of these parameters using all the available in-
10 formation is often unfeasible or impractical for most real-world scenarios. Many current signal processing applications require performing statistical inference
on large datasets, where the amount of data at hand imposes computational and/or storage constraints that impede the global estimation process [9]. Fur-
thermore, even when approximate numerical solutions working directly on the whole dataset can be computed, they may not provide a satisfactory perfor-
15 mance either. For example, Monte Carlo (MC) methods are often used to attain asymptotically exact estimators when closed-form analytical expressions cannot be obtained [10, 11, 12]. However, large datasets also pose a challenge for MC-based estimators, since the posterior density tends to be concentrated on a
relatively small portion of the state space as the number of data increases [13].
20 Consequently, MC algorithms may have trouble locating this area (especially if the dimension of the state space is also large), and thus can lead to a poor performance in practice.

An alternative to *global estimation* consists of dividing the available data into groups of manageable information, and distributing them among multiple

25 workers (cores, machines or a combination of both). The computations are then performed in parallel (with or without communication among the different workers) and *distributed* or *partial* estimators of the unknown parameters are obtained. In this setting, two extreme situations may arise, namely the multi-core and the multi-machine scenarios. On the one hand, in the *multi-*
30 *core* case, the estimation is performed using several cores of a single machine (e.g., inside a graphics processing unit [GPU]) and communication among the cores can be considered costless [14, 15]. This approach allows for communication among workers, can provide significant speed-ups (if synchronization issues are properly addressed), and solves the computational cost problem, but not
35 the memory/disk storage bottleneck. On the other hand, in the *multi-machine* case, the estimation is distributed among several machines (typically lying inside a large cluster), and the cost of inter-machine communications cannot be ignored. This approach can alleviate all the issues associated to big data signal processing (i.e., both computational and memory/storage issues), but requires
40 each machine to work independently without any communication among workers (which typically communicate only to the central node at the beginning and the end of their tasks) [16]. Finally, note that a combination of both scenarios often occurs in practice (i.e., a large cluster where each machine may have several cores), thus resulting in situations where a moderate amount of
45 communications may be acceptable.

In this paper, we focus on the scenario where no communication exists among the workers, deriving efficient linear fusion rules for the combination of the partial estimators. Our main goal is finding an optimal combination of these partial estimators that allows us to achieve a performance which is as close as
50 possible to the performance of the global estimator that has access to all the information. We concentrate on minimum mean squared error (MMSE) global estimators, but the developed framework is very general and can be used to combine any type of unbiased partial estimators. In the following, we review related works (Section 1.1), detail our main contributions (Section 1.2) and
55 summarize the structure of the whole paper (Section 1.3).

1.1. Related Works

The fusion of different models or estimators has been studied in many different areas including control, signal processing, econometrics and digital communications. The literature on the subject is rather vast, so here we only mention the most important results related to the problem addressed. On the one hand, a related field in the statistical literature is the combination of forecasts [17]. Indeed, the optimal linear combination for the single parameter case was already derived in [18, 19], a Bayesian perspective was provided in [20], and a general procedure to combine estimators in the multiple parameter case has been proposed very recently in [21]. However, there are two important differences with respect to the scenario addressed here: (1) each forecaster is assumed to have access to the whole dataset; (2) the computational complexity issue is not addressed. Therefore, problems related to the scarcity of data per estimator (when the number of data is large but the ratio data/workers is not so large), such as the so-called *small sample bias* [22], the choice of the appropriate number of partial estimators or the feasibility of the optimal combination rules when the number of parameters to be estimated is also large, have never been investigated in this context as far as we know.

On the other hand, in wireless sensor networks the focus has been on distributed learning/estimation under communication constraints [23, 24]. The optimal linear fusion rule for the multi-dimensional case has also been derived in this context [24, 25], but the focus has been on developing optimal compression rules to restrict the amount of information being transmitted, rather than on obtaining efficient fusion schemes. Unfortunately, this compression is of limited use in the multi-machine learning scenario, since passing messages among multiple machines is expensive regardless of their size. Distributed fusion approaches, obtained by adapting methods developed for graphical models, have also been proposed [26], as well as many different consensus, gossip or diffusion algorithms [27, 28, 29]. However, all of these methods still require a non-negligible amount of communication that constitutes a burden for multi-machine signal processing.

Finally, there is currently a great deal of interest in parallel Bayesian computation using MC methods [30], and a few communication-free parallel Markov chain Monte Carlo (MCMC) algorithms working on disjoint partial datasets have been developed following the so-called *embarrassingly parallel* architecture [31]. In [32], four alternatives were proposed to combine the samples drawn from the partial posteriors using either a Gaussian approximation or importance resampling. Then, [33] derived the optimal linear combination of weights required to obtain samples approximately from the full posterior, noting that the approach is optimal when both the full and the partial posteriors are Gaussian. This was followed by [34], where three different approaches to approximate the full posterior from the partial posteriors were proposed: a simple parametric approach, a non-parametric estimator and a semi-parametric method. In [35], an improvement of the quality of the approximation to the full posterior was proposed by using the Weierstrass transform. A variational aggregation approach has been derived in [36], whereas an approach based on space partitioning and density aggregation has been introduced in [37]. However, none of these previous works addresses the potentially large dimension of the optimal combiners in practical problems, which demands the transmission of large matrices across the network. Furthermore, all of the aforementioned works focus on the generation of valid Monte Carlo samples at the fusion node (thus requiring the transmission of all the samples generated in the distributed nodes, which can be an excessive burden in some environments like wireless sensor networks), whereas we concentrate on the problem of obtaining a good global estimator using a single estimate from each distributed node.

1.2. Main Contributions

The main contribution of this work is the derivation of two novel efficient linear schemes for the fusion of the partial estimators. Although we focus on minimum mean squared error (MMSE) partial estimators throughout the paper, the proposed fusion schemes are independent from the specific approach followed to

115 obtain those partial estimators. Indeed, the only assumption is the unbiased-
ness of the partial estimators. Note that this requires both the total number of
data and the number of data per partial estimator (i.e., N and N_ℓ) to be large.
The case where N is large but N_ℓ is small requires either taking into account
120 the bias in the formulation of the fusion rules (a very hard task which is still
an open problem) or designing bias correction schemes, and falls outside of the
scope of this paper.

Our motivation comes from the optimal linear combination, which involves
the calculation of one weighting matrix per partial estimator and thus may be
too computationally demanding for large dimensional systems (both in the num-
125 ber of unknowns and observations), as it requires as many weighting matrices
(whose size depends quadratically on the number of unknowns) as partial esti-
mators (whose number is typically a fraction of the number of observations). For
instance, in a setting where the number of parameters to be estimated is D and
the N observations available are equally distributed among L partial estimators,
130 the optimal linear fusion approach requires computing one $D \times D$ matrix per
partial estimator (L matrices and LD^2 parameters in total), which must be es-
timated from the partial dataset composed of N/L samples. In order to reduce
the computational complexity, we propose two linear approaches: the single co-
efficient MMSE (SC-MMSE) fusion rule, that requires only a single weighting
135 coefficient per partial estimator (i.e., L weights in total), and the independent
linear MMSE (IL-MMSE) fusion, which requires one weighting coefficient per
parameter and partial estimator (i.e., LD weights in total), respectively.

Another important contribution of the paper is providing both a constrained
optimization perspective and a Bayesian point of view (based on the Bernstein-
140 von Mises theorem and the asymptotic normality of the estimators) for the
derivation of all the linear fusion rules considered. These two complementary vi-
sions help to explain the good performance of the derived fusion rules, even when
the normality assumption required by the Bayesian approach is not strictly ful-
filled. The optimal linear combination, derived first, provides the global MMSE
145 estimator only when the partial MMSE estimators have a Gaussian distribution.

Under certain regularity conditions, this is ensured by the Bernstein-von Mises theorem in the large-sample size limit for each partial estimator (i.e., when N/L is large). However, even when this theorem is not fulfilled and the partial estimators do not follow a Gaussian distribution, the optimal linear fusion rule provides the best linear unbiased estimator (in the sense of minimizing the MSE) given the unbiased partial estimates, as shown by the constrained optimization formulation. This explains the good performance of the optimal fusion rule even when the underlying distributions are not Gaussian. The efficient SC-MMSE and IL-MMSE linear fusion rules, derived next, can then be seen as the optimal restricted linear fusion rules corresponding to a single coefficient and a diagonal matrix, respectively. From a Bayesian point of view, this is equivalent to assuming independence among all the parameters and equal quality in the estimation of all the parameters, respectively. Again, even if these restrictive conditions are not fulfilled, the constrained optimization perspective ensures us that the derived estimators are the best ones that can be obtained given the restrictions. Moreover, it allows us to see that these fusion rules can still provide a good performance when the constraints are approximately met (e.g., if the dependence among the parameters is weak and the quality in the estimation of the different parameters is similar, respectively).

Finally, we analyze the performance of all the fusion rules on several numerical examples. Firstly, we perform a detailed study on two examples where the exact closed-form expressions for the partial and the global estimators can be obtained: a univariate Gamma distribution and a multi-variate Gaussian model. This allows us to rule out any approximation effects (e.g., due to slow convergence and poor mixing in MC methods) and to analyze the effect of the number of samples, the number of estimators, the prior, and the dimensionality of the state space. Then, we apply the proposed algorithms to the problem of target localization in a wireless sensor network using measurements acquired by several sensors with different noise characteristics. In this scenario, MC partial estimators are used to deal with the groups of measurements, showing that the performance of the novel fusion rules is close to that of the optimal fusion rule

(or even better in some cases) with only a fraction of its computational cost.

Note that a preliminary version of this work has been published in [38]. In this paper we elaborate on that work, introducing the following main novel contributions: a more exhaustive review of the literature, a full Bayesian analysis of the different fusion rules, all the technical derivations that could not be included in [38], and additional numerical examples. A related (but different) problem was also considered in [39]. In that paper we focused on developing a bias correction scheme for problems where the available number of data per partial estimator (N_ℓ for $\ell = 1, \dots, L$) was small and we tested it on the univariate Gamma example shown in Section 5.1. In short, this approach works by reserving a certain amount of data to perform a rough estimation of the desired parameter at the fusion center and then using this estimate to build a custom prior for the partial estimators. However, although this approach worked very well in this case, we have not tested it on more complex examples yet and it has some important disadvantages (e.g., when the information is distributed across several locations (as in a wireless sensor network) it may require the transmission of a large number of data to the fusion center).

1.3. Organization

The remainder of the paper is structured as follows. The notation and the problem statement are provided first in Section 2, which briefly recalls the statistical inference framework to derive parameter estimators (Section 2.1), and compares the global and partial estimators (Section 2.2). This is followed by Section 3, where the optimal linear combination method is obtained by solving a constrained minimization problem (Section 3.1), and two novel efficient linear fusion rules are also derived following this approach (Sections 3.2 and 3.3). An alternative approach is then pursued in Section 4, which provides a Bayesian perspective on linear fusion rules: the best linear unbiased fusion rule (i.e., the one that minimizes the MMSE) is derived for the Gaussian case (Section 4.1); the asymptotic optimality of this fusion rule in other cases is proved through

the asymptotic normality of the partial MMSE estimators as formulated by the Bernstein-von Mises theorem (Section 4.2); and some particular cases that lead to the proposed efficient fusion rules are finally discussed (Section 4.3). Several numerical experiments are analyzed and discussed in Section 5: a one-
210 dimensional problem using Gamma distributions, where analytical expressions for the partial MMSE estimators can be obtained and we analyze the effect of the bias (Section 5.1); a multi-variate Gaussian model, where closed-form estimators can still be found and we study the effect of increasing the number of parameters to be estimated (Section 5.2); and two localization problems in
215 wireless sensor networks, where MCMC methods have to be used to obtain the partial and global estimators (Section 5.3). Finally, some concluding remarks and future lines are provided in Section 6, whereas the appendices provide the detailed technical derivations.

2. Problem Statement

2.1. Statistical Inference

Many applications in statistical signal processing require inferring a set of variables of interest or unknowns given a collection of observations or measurements. Let us consider a D -dimensional vector of unknowns, $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^D$, and let $\mathbf{y} \in \mathcal{Y} \subseteq \mathbb{R}^N$ be a collection of N i.i.d. observed data. The two major
225 paradigms in statistical inference, the frequentist approach and the Bayesian perspective [3, 4], are briefly reviewed in the sequel.

2.1.1. Frequentist Parameter Estimation

From a *frequentist* point of view, the optimal estimator of \mathbf{x} is a function of the observations, $\hat{\mathbf{x}}^* = \mathbf{f}(\mathbf{y})$, that minimizes the expected value of some given *loss* or *cost function*, $C(\mathbf{x}, \hat{\mathbf{x}})$, i.e.,

$$\hat{\mathbf{x}}^* = \arg \min_{\hat{\mathbf{x}}} \mathbb{E}(C(\mathbf{x}, \hat{\mathbf{x}})), \quad (1)$$

where $\mathbb{E}(\cdot)$ denotes the mathematical expectation, in this case w.r.t. the PDF of the data, $p(\mathbf{y})$. Although many different cost functions have been considered, for statistical estimation the most common one is the squared loss, $C(\mathbf{x}, \hat{\mathbf{x}}) = (\hat{\mathbf{x}} - \mathbf{x})^\top (\hat{\mathbf{x}} - \mathbf{x})$, which leads to the well-known least squares (LS) or frequentist minimum mean squared error (MMSE) estimator:

$$\hat{\mathbf{x}}^{(\text{LS})} = \arg \min_{\hat{\mathbf{x}}} \text{MSE}(\hat{\mathbf{x}}), \quad (2)$$

where

$$\text{MSE}(\hat{\mathbf{x}}) = \mathbb{E}((\hat{\mathbf{x}} - \mathbf{x})^\top (\hat{\mathbf{x}} - \mathbf{x})) = \int_{\mathbf{y}} (\hat{\mathbf{x}} - \mathbf{x})^\top (\hat{\mathbf{x}} - \mathbf{x}) p(\mathbf{y}|\mathbf{x}) d\mathbf{y}. \quad (3)$$

Note that the MSE in (3) can be expressed as

$$\text{MSE}(\hat{\mathbf{x}}) = \text{Tr}(\mathbb{E}((\hat{\mathbf{x}} - \mathbf{x})(\hat{\mathbf{x}} - \mathbf{x})^\top)) = \text{Tr}(\mathbf{C}_{\hat{\mathbf{x}}}) + \mathbf{b}_{\hat{\mathbf{x}}}^\top \mathbf{b}_{\hat{\mathbf{x}}}, \quad (4)$$

with $\text{Tr}(\cdot)$ denoting the trace of a matrix (i.e., the sum of the elements along its main diagonal), $\mathbf{C}_{\hat{\mathbf{x}}} = (\hat{\mathbf{x}} - \mathbb{E}(\hat{\mathbf{x}}))(\hat{\mathbf{x}} - \mathbb{E}(\hat{\mathbf{x}}))^\top$ indicating the covariance matrix of $\hat{\mathbf{x}}$, and $\mathbf{b}_{\hat{\mathbf{x}}} = \mathbb{E}(\hat{\mathbf{x}}) - \mathbf{x}$ denoting the bias of $\hat{\mathbf{x}}$. Moreover, the estimator $\hat{\mathbf{x}}$ is usually required to be unbiased (i.e., $\mathbb{E}(\hat{\mathbf{x}}) = \mathbf{x}$), implying that $\text{MSE}(\hat{\mathbf{x}}) = \text{Tr}(\mathbf{C}_{\hat{\mathbf{x}}})$ and thus the LS estimator is simply the one that minimizes the trace of the covariance matrix of $\hat{\mathbf{x}}$.

2.1.2. Bayesian Inference

From a *Bayesian* perspective, the problem of finding an optimal estimator can be formulated as the minimization of the *Bayesian Expected Loss*,

$$r(\hat{\mathbf{x}}) = \int_{\mathbf{x}} C(\mathbf{x}, \hat{\mathbf{x}}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x}, \quad (5)$$

where $C(\mathbf{x}, \hat{\mathbf{x}})$ is again some suitable *cost function*, and $p(\mathbf{x}|\mathbf{y})$ is the posterior PDF. In the Bayesian framework, this posterior PDF contains all the information available to the user about the unknown variables \mathbf{x} , and is given by Bayes

theorem:

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})p_0(\mathbf{x})}{p(\mathbf{y})}, \quad (6)$$

where $p(\mathbf{y}|\mathbf{x})$ is the likelihood, $p_0(\mathbf{x})$ is the prior PDF, and $p(\mathbf{y})$ is the PDF of the data. Now, let us consider again the quadratic cost, $C(\mathbf{x}, \hat{\mathbf{x}}) = (\hat{\mathbf{x}} - \mathbf{x})^\top (\hat{\mathbf{x}} - \mathbf{x})$, which is also the most common cost function for regression problems within the Bayesian framework. Then, (5) becomes

$$r(\hat{\mathbf{x}}) = \text{MSE}(\hat{\mathbf{x}}|\mathbf{y}) = \int_{\mathcal{X}} (\hat{\mathbf{x}} - \mathbf{x})^\top (\hat{\mathbf{x}} - \mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x}, \quad (7)$$

where $\text{MSE}(\hat{\mathbf{x}}|\mathbf{y})$ refers to the *conditional Bayesian MSE* (i.e., the Bayesian MSE for a fixed set of data \mathbf{y}), that we will also denote in the sequel as MSE for the sake of simplicity.¹ The optimal estimator (i.e., the one that minimizes Eq. (7)) is the well-known Bayesian minimum mean squared error (MMSE) estimator, which corresponds to the conditional mean, i.e., the expected value of \mathbf{x} w.r.t. the posterior PDF [1, 3, 4, 5]:

$$\hat{\mathbf{x}}^{(\text{MMSE})} = \mathbb{E}(\mathbf{x}|\mathbf{y}) = \int_{\mathcal{X}} \mathbf{x} p(\mathbf{x}|\mathbf{y}) d\mathbf{x}. \quad (9)$$

2.2. Global vs. Distributed Partial Bayesian Estimators

All the estimators discussed in the previous section are *global estimators*, since they are assumed to have access to all the available data. Hence, their performance is optimal (from the point of view of minimizing their respective cost functions) whenever they can be computed exactly. However, in big data problems we cannot deal with the whole data set globally due to computational and/or memory restrictions. Furthermore, even when we can work with the

¹The *full Bayesian MSE* is obtained by performing a double integral on both the data and the parameters of interest using the joint PDF $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}|\mathbf{y})p(\mathbf{y})$:

$$\text{MSE}(\hat{\mathbf{x}}) = \int_{\mathcal{Y}} \int_{\mathcal{X}} (\hat{\mathbf{x}} - \mathbf{x})^\top (\hat{\mathbf{x}} - \mathbf{x}) p(\mathbf{x}|\mathbf{y}) p(\mathbf{y}) d\mathbf{x} d\mathbf{y}. \quad (8)$$

However, by assuming that the data are fixed (i.e., by conditioning on the data), the outer integral in (8) vanishes and the only integral remaining is the one on \mathbf{x} using $p(\mathbf{x}|\mathbf{y})$.

whole data set globally, splitting it into L data sets may be more efficient and lead to a better performance. For instance, in the Bayesian framework it is well-known that the posterior PDF tends to become more “peaky” as the amount of available data increases. Consequently, the inference process becomes harder (especially for high-dimensional scenarios) when closed form expressions for the Bayesian estimators cannot be found and we have to resort to numerical approximations (like Monte Carlo methods [11, 40]).

A natural solution in these cases is splitting the data into L disjoint groups or clusters, so that the ℓ -th cluster ($1 \leq \ell \leq L$) only has access to N_ℓ samples. Then, we can obtain a *partial estimator* for the ℓ -th cluster (i.e., a partial estimator of \mathbf{x} given only the partial data available within the ℓ -th cluster, \mathbf{y}_ℓ), $\hat{\mathbf{x}}_\ell$, by constructing a function of the ℓ -th data set that minimizes the corresponding cost function.² For instance, the partial MMSE estimator in the Bayesian setting would be given by

$$\hat{\mathbf{x}}_\ell^{(\text{MMSE})} = \mathbb{E}(\mathbf{x}|\mathbf{y}_\ell) = \int_{\mathcal{X}} \mathbf{x} p_\ell(\mathbf{x}|\mathbf{y}_\ell) d\mathbf{x}, \quad (10)$$

where $p_\ell(\mathbf{x}|\mathbf{y}_\ell)$ is the partial posterior associated to the ℓ -th dataset (see Table 1 for a summary of the notation used throughout the paper). The challenge now is trying to obtain the exact global MMSE estimator, $\hat{\mathbf{x}}^{(\text{MMSE})}$, from the set of partial MMSE estimators, $\{\hat{\mathbf{x}}_\ell^{(\text{MMSE})}\}_{\ell=1}^L$.

In this paper, we concentrate on a communications-free setup for the partial estimators, i.e., we assume that the partial estimators can only transmit their final estimates to the fusion center (FC), altogether with additional side information if needed, and are not allowed to communicate with each other during the estimation process. The FC will then be responsible for combining all the estimates in an efficient way in order to obtain the global MMSE estimator (if it is feasible) or at least the best possible approximation. We consider linear

²Note that we use the name partial estimator, instead of distributed or local estimator, to emphasize the fact that $\hat{\mathbf{x}}_\ell$ corresponds to the estimator of the complete set of variables of interest obtained using only the partial information available to the ℓ -th cluster.

Table 1: Summary of the Notation.

\mathbf{x}	Unknown parameters to be estimated.
$\hat{\mathbf{x}}$	Global estimator of \mathbf{x} .
D	Number of unknowns (i.e., dimension of \mathbf{x}).
\mathbf{y}	Vector of observations.
N	Number of observations.
L	Number of parallel (partial) estimators.
N_ℓ	Number of data for the ℓ -th estimator.
\mathbf{y}_ℓ	Data set for the ℓ -th estimator.
$\hat{\mathbf{x}}_\ell$	ℓ -th partial estimator of \mathbf{x} ($\ell = 1, \dots, L$).
$\mathbf{\Lambda}_\ell$	Weighting matrix used to combine the L partial estimators.
$p(\mathbf{x} \mathbf{y})$	Global posterior PDF.
$p(\mathbf{y} \mathbf{x})$	Global likelihood.
$p_0(\mathbf{x})$	Global prior PDF.
$p_\ell(\mathbf{x} \mathbf{y}_\ell)$	Partial posterior PDF for the ℓ -th estimator.
$p_\ell(\mathbf{y}_\ell \mathbf{x})$	Partial likelihood for the ℓ -th estimator.
$p_0^{(\ell)}(\mathbf{x})$	Local prior PDF for the ℓ -th estimator.

fusion rules, implying that the global estimator is given by the following general expression:

$$\hat{\mathbf{x}} = \sum_{\ell=1}^L \mathbf{\Lambda}_\ell \hat{\mathbf{x}}_\ell, \quad (11)$$

where the $\mathbf{\Lambda}_\ell$ ($\ell = 1, \dots, L$) are $D \times D$ weighting matrices. An important case where Eq. (11) leads to the exact global MMSE estimator occurs when both the global and the partial posteriors have Gaussian PDFs, as proved in Section 4.1. Indeed, when the conditions for the Bernstein-von Mises theorem are fulfilled, all the posterior PDFs are Gaussian and (11) becomes asymptotically optimal, as discussed in Section 4.2. Furthermore, even if (11) is not the global MMSE estimator, by choosing the $\mathbf{\Lambda}_\ell$ properly we obtain the best linear unbiased global estimator of \mathbf{x} given the $\hat{\mathbf{x}}_\ell$, as shown in Section 3.1. Imposing additional restrictions on the weighting matrices leads to the novel efficient linear fusion rules described in Sections 3.2 and 3.3, which can still be optimal under certain conditions (as discussed in Section 4.3) and provide good results in many other problems (as shown through numerical simulations in Section 5). In the following sections we discuss all these issues in detail.

3. Linear Fusion of Partial Estimators: A Constrained Minimization Approach

In this section, we first derive the optimal linear fusion rule (Section 3.1) from a constrained minimization perspective, and then provide two alternative efficient fusion rules by restricting the shape of the weighting matrix: the single coefficient fusion rule in Section 3.2 and the independent fusion rule in Section 3.3. These derivations are valid regardless of the approach followed to obtain the partial estimators, since the only assumption performed is their unbiasedness.

3.1. General Case: Optimal Linear Combination

Let us consider the most general linear fusion rule, which is given by

$$\hat{\mathbf{x}} = \sum_{\ell=1}^L \mathbf{\Lambda}_\ell \hat{\mathbf{x}}_\ell, \quad (12)$$

where $\hat{\mathbf{x}}_\ell$ can be any partial estimator (not necessarily the MMSE estimator) based on the ℓ -th partial dataset, \mathbf{y}_ℓ , and $\hat{\mathbf{x}}$ is the corresponding global estimator, obtained by linearly combining all those partial estimators. In this case, assuming that all the partial estimators are unbiased, the mean of the global estimator is given by

$$\mathbb{E}(\hat{\mathbf{x}}) = \sum_{\ell=1}^L \mathbf{\Lambda}_\ell \mathbb{E}(\hat{\mathbf{x}}_\ell) = \left(\sum_{\ell=1}^L \mathbf{\Lambda}_\ell \right) \mathbf{x}. \quad (13)$$

Thus, in order to obtain an unbiased global estimator (i.e., $\mathbb{E}(\hat{\mathbf{x}}) = \mathbf{x}$) we need to impose the following condition:

$$\sum_{\ell=1}^L \mathbf{\Lambda}_\ell = \mathbf{I}. \quad (14)$$

The covariance matrix of the global estimator is now given by

$$\mathbf{C}_{\mathbf{x}} = \mathbb{E}((\hat{\mathbf{x}} - \mathbf{x})(\hat{\mathbf{x}} - \mathbf{x})^\top) = \mathbb{E}\left(\left(\sum_{\ell=1}^L \mathbf{\Lambda}_\ell \hat{\mathbf{x}}_\ell - \mathbf{x}\right)\left(\sum_{\ell=1}^L \mathbf{\Lambda}_\ell \hat{\mathbf{x}}_\ell - \mathbf{x}\right)^\top\right), \quad (15)$$

where we have used (12) and assumed that the global estimator is unbiased.

Making use of (14), Eq. (15) can be expressed as

$$\begin{aligned} \mathbf{C}_{\mathbf{x}} &= \mathbb{E}\left(\sum_{\ell=1}^L \mathbf{\Lambda}_\ell (\hat{\mathbf{x}}_\ell - \mathbf{x}) \sum_{\ell=1}^L (\hat{\mathbf{x}}_\ell - \mathbf{x})^\top \mathbf{\Lambda}_\ell^\top\right) \\ &= \sum_{\ell,k=1}^L \mathbf{\Lambda}_\ell \mathbb{E}((\hat{\mathbf{x}}_\ell - \mathbf{x})(\hat{\mathbf{x}}_k - \mathbf{x})^\top) \mathbf{\Lambda}_k^\top. \end{aligned} \quad (16)$$

Finally, assuming that the partial estimators are independent and taking into account that $\mathbf{C}_{\mathbf{x}}^{(\ell)} = \mathbb{E}((\hat{\mathbf{x}}_\ell - \mathbf{x})(\hat{\mathbf{x}}_\ell - \mathbf{x})^\top)$, $\mathbf{C}_{\mathbf{x}}$ can be expressed as a function of the covariance of the partial estimators, $\mathbf{C}_{\mathbf{x}}^{(\ell)}$, and the weighting matrices, $\mathbf{\Lambda}_\ell$, as

$$\mathbf{C}_{\mathbf{x}} = \sum_{\ell=1}^L \mathbf{\Lambda}_\ell \mathbf{C}_{\mathbf{x}}^{(\ell)} \mathbf{\Lambda}_\ell^\top. \quad (17)$$

The MSE of the global estimator is then given by

$$\text{MSE}(\hat{\mathbf{x}}|\mathbf{y}) = \text{Tr}(\mathbf{C}_{\mathbf{x}}) = \sum_{\ell=1}^L \text{Tr}(\mathbf{\Lambda}_\ell \mathbf{C}_{\mathbf{x}}^{(\ell)} \mathbf{\Lambda}_\ell^\top), \quad (18)$$

275 where $\text{Tr}(\cdot)$ denotes the trace of a matrix.

Hence, in order to obtain the best linear unbiased global estimator (i.e., the linear combination of the partial estimators that minimizes the MSE), we need to solve the following constrained optimization problem:

$$\mathbf{\Lambda}^* = \arg \min_{\mathbf{\Lambda}} \sum_{\ell=1}^L \text{Tr}(\mathbf{\Lambda}_\ell \mathbf{C}_{\mathbf{x}}^{(\ell)} \mathbf{\Lambda}_\ell^\top), \quad (19a)$$

$$\text{s.t.} \quad \sum_{\ell=1}^L \mathbf{\Lambda}_\ell = \mathbf{I}, \quad (19b)$$

where $\mathbf{\Lambda} = [\mathbf{\Lambda}_1, \dots, \mathbf{\Lambda}_L]^\top$. Since (19a) and (19b) correspond to a convex optimization problem [41], by applying the method of the Lagrange multipliers it can be shown (see Appendix A) that the solution for each of the weighting matrices is given by

$$\mathbf{\Lambda}_\ell^* = \left[\sum_{k=1}^L (\mathbf{C}_x^{(k)})^{-1} \right]^{-1} (\mathbf{C}_x^{(\ell)})^{-1}. \quad (20)$$

Substituting this expression in (12), the optimal linear MMSE (L-MMSE) fusion rule is finally given by

$$\hat{\mathbf{x}}^{(\text{L-MMSE})} = \left[\sum_{k=1}^L (\mathbf{C}_x^{(k)})^{-1} \right]^{-1} \sum_{\ell=1}^L (\mathbf{C}_x^{(\ell)})^{-1} \hat{\mathbf{x}}_\ell, \quad (21)$$

regardless of the approach followed to derive the partial estimators.

3.2. Efficient Fusion Rule: Single Coefficient Fusion

Let us consider the particular case in which a single coefficient per estimator is used to construct the global estimator, i.e.,

$$\hat{\mathbf{x}} = \sum_{\ell=1}^L \alpha_\ell \hat{\mathbf{x}}_\ell, \quad (22)$$

which is obtained by setting $\mathbf{\Lambda}_\ell = \alpha_\ell \mathbf{I}$ in (12). Clearly, this will provide a suboptimal solution in general, but it is a fast and low-cost solution for the combination of estimators, and we can easily obtain a closed-form expression for the optimal weights.

On the one hand, since the partial estimators are unbiased, it is straightforward to see that the mean of the global estimator given by (22) is

$$\mathbb{E}(\hat{\mathbf{x}}) = \sum_{\ell=1}^L \alpha_\ell \mathbb{E}(\hat{\mathbf{x}}_\ell) = \left(\sum_{\ell=1}^L \alpha_\ell \right) \mathbf{x}. \quad (23)$$

Hence, in order to obtain an unbiased global estimator we need to have

$$\sum_{\ell=1}^L \alpha_{\ell} = 1. \quad (24)$$

On the other hand, the covariance matrix for the global estimator is given by

$$\mathbf{C}_{\mathbf{x}} = \sum_{\ell=1}^L \alpha_{\ell}^2 \mathbf{C}_{\mathbf{x}}^{(\ell)}, \quad (25)$$

and the MSE can thus be expressed as

$$\text{MSE}(\hat{\mathbf{x}}|\mathbf{y}) = \text{Tr}(\mathbf{C}_{\mathbf{x}}) = \sum_{\ell=1}^L \alpha_{\ell}^2 \text{Tr}(\mathbf{C}_{\mathbf{x}}^{(\ell)}), \quad (26)$$

where $\text{Tr}(\mathbf{C}_{\mathbf{x}})$ denotes the trace of the global covariance matrix:

$$\text{Tr}(\mathbf{C}_{\mathbf{x}}) = \sum_{d=1}^D \mathbf{C}_{\mathbf{x}}[d, d] = \sum_{d=1}^D \sigma_{x_d}^2, \quad (27)$$

with $\sigma_{x_d}^2 = \mathbb{E}((\hat{x}_d - x_d)^2)$, and $\text{Tr}(\mathbf{C}_{\mathbf{x}}^{(\ell)})$ denotes the trace of the ℓ -th partial covariance matrix:

$$\text{Tr}(\mathbf{C}_{\mathbf{x}}^{(\ell)}) = T_{\ell} = \sum_{d=1}^D \mathbf{C}_{\mathbf{x}}^{(\ell)}[d, d] = \sum_{d=1}^D \sigma_{\ell,d}^2, \quad (28)$$

with $\sigma_{\ell,d}^2 = \mathbb{E}((\hat{x}_d^{(\ell)} - x_d)^2)$.

The goal is finding the set of α_{ℓ} that minimizes (26), subject to Eq. (24) in order to obtain an unbiased estimator. Hence, the optimal selection of the weights can be formulated as the following constrained optimization problem:

$$\boldsymbol{\alpha}^* = \arg \min_{\boldsymbol{\alpha}} \sum_{\ell=1}^L \alpha_{\ell}^2 \text{Tr}(\mathbf{C}_{\mathbf{x}}^{(\ell)}), \quad (29a)$$

$$\text{s.t.} \quad \sum_{\ell=1}^L \alpha_{\ell} = 1, \quad (29b)$$

with $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_L]^\top$. Eqs. (29a) and (29b) correspond again to a convex optimization problem. Thus, by applying once more the method of the Lagrange multipliers, it can be shown (see Appendix A) that the optimal value of α_ℓ is

$$\alpha_\ell^* = \frac{[\text{MSE}(\hat{\mathbf{x}}_\ell|\mathbf{y}_\ell)]^{-1}}{\sum_{k=1}^L [\text{MSE}(\hat{\mathbf{x}}_k|\mathbf{y}_k)]^{-1}} = \frac{T_\ell^{-1}}{\sum_{k=1}^L T_k^{-1}}, \quad (30)$$

where we recall that $T_\ell = \text{Tr}(\mathbf{C}_\mathbf{x}^{(\ell)})$, and the single coefficient MMSE (SC-MMSE) fusion rule is then given by

$$\hat{\mathbf{x}}^{(\text{SC-MMSE})} = \sum_{\ell=1}^L \frac{T_\ell^{-1}}{\sum_{k=1}^L T_k^{-1}} \hat{\mathbf{x}}_\ell. \quad (31)$$

3.3. Efficient Fusion Rule: Independent Linear Fusion

The SC-MMSE fusion rule has a substantially reduced computational cost w.r.t. the L-MMSE approach, since it only requires the estimation of L parameters overall instead of the D^2L parameters of the L-MMSE rule. However, noting that the optimal weights in (31) involve the trace of the partial covariance matrices, we introduce an independent linear minimum mean squared error (IL-MMSE) fusion rule, where $\mathbf{A}_\ell = \mathbf{D}_\ell = \text{diag}(\alpha_{\ell,1}, \dots, \alpha_{\ell,D})$. This approach leads to an independent estimation of each of the D unknowns:

$$\hat{x}_d = \sum_{\ell=1}^L \alpha_{\ell,d} \hat{x}_{\ell,d}, \quad (32)$$

where $1 \leq d \leq D$, \hat{x}_d denotes the d -th component of the global estimator, and $\hat{x}_{\ell,d}$ indicates the d -th component of the ℓ -th partial estimator. The mean of the d -th component of the global estimator is then given by

$$\mathbb{E}(\hat{x}_d) = \sum_{\ell=1}^L \alpha_{\ell,d} \mathbb{E}(\hat{x}_{\ell,d}) = \left(\sum_{\ell=1}^L \alpha_{\ell,d} \right) x_d, \quad (33)$$

and the unbiasedness condition becomes

$$\sum_{\ell=1}^L \alpha_{\ell,d} = 1, \quad (34)$$

for $d = 1, \dots, D$. Hence, the covariance matrix for the global estimator is

$$\mathbf{C}_x = \sum_{\ell=1}^L \mathbf{D}_\ell \mathbf{C}_x^{(\ell)} \mathbf{D}_\ell^\top, \quad (35)$$

and the MSE can thus be expressed as

$$\text{MSE}(\hat{\mathbf{x}}|\mathbf{y}) = \text{Tr}(\mathbf{C}_x) = \sum_{\ell=1}^L \text{Tr}(\mathbf{D}_\ell \mathbf{C}_x^{(\ell)} \mathbf{D}_\ell^\top) = \sum_{\ell=1}^L \sum_{d=1}^D \alpha_{\ell,d}^2 \sigma_{\ell,d}^2. \quad (36)$$

From (34) and (36), it becomes apparent that the weights in (32) can be obtained by solving D single parameter constrained optimization problems:

$$\boldsymbol{\alpha}_d^* = \arg \min_{\boldsymbol{\alpha}_d} \sum_{\ell=1}^L \alpha_{\ell,d}^2 \sigma_{\ell,d}^2, \quad (37a)$$

$$\text{s.t.} \quad \sum_{\ell=1}^L \alpha_{\ell,d} = 1, \quad (37b)$$

where $\boldsymbol{\alpha}_d = [\alpha_{1,d}, \dots, \alpha_{L,d}]^\top$, and we recall that $\sigma_{\ell,d}^2 = C_{x_d}^{(\ell)}$ is the d -th element along the main diagonal of $\mathbf{C}_x^{(\ell)}$. The solution (see Appendix A) is now given by

$$\alpha_{\ell,d}^* = \frac{\left[\text{MSE}(\hat{x}_{\ell,d}^{(\text{MMSE})} | \mathbf{y}_\ell) \right]^{-1}}{\sum_{k=1}^L \left[\text{MSE}(\hat{x}_{k,d}^{(\text{MMSE})} | \mathbf{y}_k) \right]^{-1}} = \frac{\sigma_{\ell,d}^{-2}}{\sum_{k=1}^L \sigma_{k,d}^{-2}}, \quad (38)$$

and thus the IL-MMSE fusion rule is

$$\hat{x}_d^{(\text{IL-MMSE})} = \sum_{\ell=1}^L \frac{\sigma_{\ell,d}^{-2}}{\sum_{k=1}^L \sigma_{k,d}^{-2}} \hat{x}_{\ell,d}, \quad (39)$$

for $d = 1, \dots, D$. This approach requires the estimation of DL parameters overall, and thus it can be seen as an intermediate approach between the L-

MMSE and the SC-MMSE fusion rules, both in terms of computational cost and performance (as shown in Section 5).

4. Optimal Linear Fusion: Bayesian Perspective

In this section, we provide an alternative perspective of the optimal linear fusion
 290 problem from a Bayesian point of view. First of all, we derive the optimal fusion
 rule for the multi-variate Gaussian model in Section 4.1, and then we show
 that this rule is asymptotically optimal under mild conditions in Section 4.2.
 Finally, in Section 4.3 we address two relevant particular cases that correspond
 to the SC-MMSE and IL-MMSE fusion rules derived in Sections 3.2 and 3.3,
 295 respectively.

4.1. Gaussian Estimators: Optimal Fusion Rule

Let us consider the derivation of the optimal fusion rule at the Fusion Center
 (FC) from a Bayesian point of view. In this case, our observations are the
 outputs of each of the L partial estimators, $\hat{\mathbf{x}}_\ell$ for $\ell = 1, \dots, L$. Let us assume
 that these estimators are independent, unbiased (i.e., $\mathbb{E}(\hat{\mathbf{x}}_\ell) = \mathbf{x}$), and have
 Gaussian PDFs with means equal to the true parameter vector \mathbf{x} and covariance
 matrices $\mathbf{C}_\mathbf{x}^{(\ell)}$:

$$\begin{aligned} p(\hat{\mathbf{x}}_\ell|\mathbf{x}) &= \mathcal{N}(\hat{\mathbf{x}}_\ell|\mathbf{x}, \mathbf{C}_\mathbf{x}^{(\ell)}) \\ &= \left(2\pi|\mathbf{C}_\mathbf{x}^{(\ell)}|\right)^{-D/2} \exp\left(-\frac{1}{2}(\hat{\mathbf{x}}_\ell - \mathbf{x})^\top \left[\mathbf{C}_\mathbf{x}^{(\ell)}\right]^{-1} (\hat{\mathbf{x}}_\ell - \mathbf{x})\right). \end{aligned} \quad (40)$$

The full posterior is

$$p(\mathbf{x}|\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_L) \propto p(\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_L|\mathbf{x})p_0(\mathbf{x}), \quad (41)$$

with the likelihood function given by

$$p(\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_L | \mathbf{x}) = \prod_{\ell=1}^L p(\hat{\mathbf{x}}_\ell | \mathbf{x}) = \prod_{\ell=1}^L \mathcal{N}(\mathbf{x} | \hat{\mathbf{x}}_\ell, \mathbf{C}_\mathbf{x}^{(\ell)}), \quad (42)$$

and a Gaussian prior with mean $\hat{\mathbf{x}}_0$ and covariance matrix $\mathbf{C}_\mathbf{x}^{(0)}$,

$$p_0(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \hat{\mathbf{x}}_0, \mathbf{C}_\mathbf{x}^{(0)}). \quad (43)$$

Inserting (42) and (43) into (41), the full posterior can be finally expressed as

$$\begin{aligned} p(\mathbf{x} | \hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_L) &\propto \prod_{\ell=0}^L \mathcal{N}(\mathbf{x} | \hat{\mathbf{x}}_\ell, \mathbf{C}_\mathbf{x}^{(\ell)}) \\ &\propto \prod_{\ell=0}^L |\mathbf{C}_\mathbf{x}^{(\ell)}|^{-1/2} \exp\left(-\frac{1}{2} \sum_{\ell=1}^L (\mathbf{x} - \hat{\mathbf{x}}_\ell)^\top (\mathbf{C}_\mathbf{x}^{(\ell)})^{-1} (\mathbf{x} - \hat{\mathbf{x}}_\ell)\right). \end{aligned} \quad (44)$$

After some algebra, Eq. (44) can be expressed as a single multi-variate Gaussian:

$$p(\mathbf{x} | \hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_L) = (2\pi)^{-D/2} |\mathbf{C}_\mathbf{x}|^{-1/2} \exp\left(-\frac{1}{2} (\hat{\mathbf{x}} - \boldsymbol{\mu}_\mathbf{x})^\top \mathbf{C}_\mathbf{x}^{-1} (\hat{\mathbf{x}} - \boldsymbol{\mu}_\mathbf{x})\right), \quad (45)$$

where

$$\mathbf{C}_\mathbf{x} = \left[\sum_{\ell=0}^L (\mathbf{C}_\mathbf{x}^{(\ell)})^{-1} \right]^{-1}, \quad (46a)$$

$$\boldsymbol{\mu}_\mathbf{x} = \mathbf{C}_\mathbf{x} \sum_{\ell=0}^L (\mathbf{C}_\mathbf{x}^{(\ell)})^{-1} \hat{\mathbf{x}}_\ell. \quad (46b)$$

Hence, the global MMSE estimator for the multi-variate Gaussian model, that we denote as the G-MMSE estimator and corresponds to the mean of the full

posterior in Eq. (45), is finally given by

$$\hat{\mathbf{x}}^{(\text{G-MMSE})} = \boldsymbol{\mu}_{\mathbf{x}} = \sum_{\ell=0}^L \boldsymbol{\Lambda}_{\ell} \hat{\mathbf{x}}_{\ell}, \quad (47)$$

with

$$\boldsymbol{\Lambda}_{\ell} = \mathbf{C}_{\mathbf{x}} \cdot \left(\mathbf{C}_{\mathbf{x}}^{(\ell)} \right)^{-1} = \left[\sum_{k=0}^L \left(\mathbf{C}_{\mathbf{x}}^{(k)} \right)^{-1} \right]^{-1} \left(\mathbf{C}_{\mathbf{x}}^{(\ell)} \right)^{-1}. \quad (48)$$

Therefore, given any set of partial estimators $\{\hat{\mathbf{x}}_{\ell}\}_{\ell=1}^L$, the linear fusion is always optimal in the multi-variate Gaussian case, with an optimal weighting matrix given by Eq. (48). Furthermore, when the $\hat{\mathbf{x}}_{\ell}$ are the partial MMSE estimators (i.e., when $\hat{\mathbf{x}}_{\ell} = \hat{\mathbf{x}}_{\ell}^{(\text{MMSE})}$ for $\ell = 1, \dots, L$), then we have $\hat{\mathbf{x}}^{(\text{G-MMSE})} = \hat{\mathbf{x}}^{(\text{MMSE})}$, with $\hat{\mathbf{x}}^{(\text{MMSE})}$ denoting the global MMSE estimator that has access to all the data. Finally, note that the optimal fusion rule derived from the Bayesian perspective is equivalent to the optimal fusion rule obtained following the constrained optimization approach, since (48) is identical to (20).³

4.2. Asymptotically Optimal Fusion: Bernstein-von Mises Theorem

The Bernstein-von Mises (a.k.a. Bayesian central limit) theorem states that, under mild regularity conditions, the posterior PDF converges to a Gaussian PDF as the number of samples tends to infinity **regardless of the specific prior**

³Actually, there is a difference between the optimal fusion rules derived from a Bayesian and a constrained optimization point of view, since the sum in (48) starts at $\ell = 0$, while the sum in (20) starts at $\ell = 1$. This is due to the fact that the sum in (48) includes the prior, whose mean and covariance are $\hat{\mathbf{x}}_0$ and $\mathbf{C}_{\mathbf{x}}^{(0)}$, respectively. However, by using a non-informative prior (i.e., a prior such that both $\hat{\mathbf{x}}_0 = \mathbf{0}$ and $(\mathbf{C}_{\mathbf{x}}^{(0)})^{-1} = \mathbf{0}$), (48) and (20) become identical.

selected [42, 43].⁴ Applying this result to the partial posterior PDFs, we have

$$p_\ell(\mathbf{x}|\mathbf{y}_\ell) \rightarrow \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_\mathbf{x}^{(\ell)}, \mathbf{C}_\mathbf{x}^{(\ell)}) \quad \text{as } N_\ell \rightarrow \infty, \quad (49)$$

with $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_\mathbf{x}^{(\ell)}, \mathbf{C}_\mathbf{x}^{(\ell)})$ indicating that \mathbf{x} has a Gaussian PDF with a mean vector $\boldsymbol{\mu}_\mathbf{x}^{(\ell)} = \hat{\mathbf{x}}_\ell^{(\text{MMSE})}$ and a covariance matrix

$$\begin{aligned} \mathbf{C}_\mathbf{x}^{(\ell)} &= \mathbb{E} \left((\mathbf{x} - \hat{\mathbf{x}}_\ell^{(\text{MMSE})})(\mathbf{x} - \hat{\mathbf{x}}_\ell^{(\text{MMSE})})^\top \right) \\ &= \int_{\mathcal{X}} (\mathbf{x} - \hat{\mathbf{x}}_\ell^{(\text{MMSE})})(\mathbf{x} - \hat{\mathbf{x}}_\ell^{(\text{MMSE})})^\top p_\ell(\mathbf{x}|\mathbf{y}_\ell) d\mathbf{x}. \end{aligned} \quad (50)$$

Assuming that we have independent (though not necessarily identically distributed) observations and that each of them can only belong to one cluster (i.e., we have disjoint sets of samples such that $N = \sum_{\ell=1}^L N_\ell$), the global posterior PDF also converges to a Gaussian PDF as N tends to infinity, i.e.,

$$p(\mathbf{x}|\mathbf{y}) = \prod_{\ell=1}^L p_\ell(\mathbf{x}|\mathbf{y}_\ell) \rightarrow \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_\mathbf{x}, \mathbf{C}_\mathbf{x}) \quad \text{as } N \rightarrow \infty, \quad (51)$$

with $\mathbf{C}_\mathbf{x}$ and $\boldsymbol{\mu}_\mathbf{x}$ given by (46a) and (46b), respectively, using $\hat{\mathbf{x}}_\ell = \hat{\mathbf{x}}_\ell^{(\text{MMSE})}$. Eqs. (49) and (51) state that, even if $p(\mathbf{x}|\mathbf{y})$ and the $p_\ell(\mathbf{x}|\mathbf{y}_\ell)$ are not Gaussian, they will converge to Gaussian PDFs as $N_\ell \rightarrow \infty$ for $\ell = 1, \dots, L$ (and thus, ³¹⁰ $N = \sum_\ell N_\ell \rightarrow \infty$ also). This implies that, when the regularity conditions of

⁴The conditions for the Bernstein-von Mises theorem are remarkably weak (e.g., see Section 10.2 of [43]):

- The prior measure must be absolutely continuous in a neighborhood of the true parameter, θ_0 , with a continuous positive density at θ_0 .
- The assumed model must be differentiable in quadratic mean w.r.t. the parameter to be estimated.
- It must be possible to separate the true parameter θ_0 from the complements of balls centered at θ_0 .

Therefore, regarding the prior it is only required to be continuous and positive in a neighborhood of the true parameter, something satisfied by most prior densities. And regarding the other two conditions, the second one simply requires the differentiability of the model (fulfilled again by most models), whereas the last condition is automatically satisfied when the parameter space is compact and the model is continuous and identifiable (see again Section 10.2 of [43]).

the Bernstein-von Mises theorem are fulfilled, the linear fusion rule derived in Section 4.1 is asymptotically optimal, i.e., $\hat{\mathbf{x}}^{(\text{G-MMSE})} \rightarrow \hat{\mathbf{x}}^{(\text{MMSE})}$ when $\hat{\mathbf{x}}_\ell = \hat{\mathbf{x}}_\ell^{(\text{MMSE})}$ and $N_\ell \rightarrow \infty$ for $\ell = 1, \dots, L$.

4.3. Particular Cases

315 Note that the fusion rule of Eq. (47) requires the computation of a $D \times D$ weighting matrix, given by Eq. (48), for each of the L estimators. This implies calculating D^2L weights overall, which may be unfeasible (or at least very costly from a computational/storage point of view) when D and/or L is large. Moreover, even if the optimal weighting matrix can be computed, the numerical errors that arise in its computation (since it depends on the $\mathbf{C}_{\mathbf{x}}^{(\ell)}$, which
320 are usually unknown and must be estimated from the data) may hinder the performance of the optimal fusion rule, especially when D is large and D^2 is comparable to N_ℓ .

In these cases, structured weighting matrices (which require computing a reduced number of coefficients) can be used to obtain an approximation of the optimal case. Indeed, in certain cases the optimum weighting matrix may already contain a reduced number of different elements. As a first example, let us consider the situation where the parameters to be inferred are not interrelated. In this scenario, the covariance matrix for the ℓ -th estimator becomes diagonal,

$$\mathbf{C}_{\mathbf{x}}^{(\ell)} = \text{diag}(\sigma_{\ell,1}^2, \dots, \sigma_{\ell,D}^2), \quad (52)$$

with

$$\sigma_{\ell,d}^2 = \int_{\mathcal{X}_d} (\hat{x}_{\ell,d} - x_d)^2 p(x_d | \mathbf{y}_\ell) dx_d \quad (53)$$

for $d = 1, \dots, D$. Hence, the optimal weighting matrix becomes

$$\mathbf{\Lambda}_\ell = \text{diag}(\alpha_{\ell,1}, \dots, \alpha_{\ell,D}), \quad (54)$$

with

$$\alpha_{\ell,d} = \frac{\sigma_{\ell,d}^{-2}}{\sum_{k=1}^L \sigma_{k,d}^{-2}}. \quad (55)$$

Note that only D parameters are required now for each of the L estimators (i.e., 325 DL parameters in total).

As a second example, let us assume that, in addition to having independent parameters to be inferred, the uncertainty in their estimation is the same for a given partial estimator, i.e., $\sigma_{\ell,1}^2 = \sigma_{\ell,2}^2 = \dots = \sigma_{\ell,D}^2 = \sigma_{\ell}^2$. Then, the covariance matrix for the ℓ -th partial estimator becomes $\mathbf{C}_{\mathbf{x}}^{(\ell)} = \sigma_{\ell}^2 \mathbf{I}$, where \mathbf{I} denotes the $D \times D$ identity matrix, and the optimal weighting matrix is now given by

$$\mathbf{\Lambda}_{\ell} = \alpha_{\ell} \mathbf{I}, \quad (56)$$

with

$$\alpha_{\ell} = \frac{\sigma_{\ell}^{-2}}{\sum_{k=1}^L \sigma_k^{-2}}. \quad (57)$$

Note that this approach only requires a single parameter per estimator (i.e., L parameters in total), and thus it is very efficient from a computational point of view. Furthermore, if all the variances are equal (i.e., $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_L^2 = \sigma^2$) we obtain the so-called *equal weights fusion* (EWF) rule, which assigns the same weight to all the partial estimators,

$$\alpha_1 = \alpha_2 = \dots = \alpha_L = \frac{1}{L}, \quad (58)$$

and thus it does not require the transmission of any coefficient to perform the fusion at the FC.

Obviously, in real applications it is very unlikely that all the parameters to be inferred are completely independent, and even more unlikely that the uncertainty in their estimation is exactly equal. However, it can be shown (see 330 Appendix B) that the coefficients in Eqs. (54) and (56) correspond, respectively, to the best independent and isotropic Gaussian approximations (in terms of minimizing the Kullback-Leibler (KL) divergence) to the asymptotic Gaussian

PDFs of the partial estimators, $p_\ell(\mathbf{x}|\hat{\mathbf{x}}_\ell) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_\mathbf{x}, \mathbf{C}_\mathbf{x}^{(\ell)})$. Hence, even if the
 335 independence and/or equal quality assumptions are not fulfilled exactly, Eqs.
 (54) and (56) can still lead to good fusion rules (as shown in Section 5) if the
 approximations of the partial PDFs are accurate enough.

5. Numerical Experiments

In this section, we provide three different numerical examples where we compare
 340 the four linear fusion rules derived in the previous two sections:

1. **Linear MMSE (L-MMSE)** fusion rule, which uses the full weighting matrix and the global estimator is given by Eq. (21).
2. **Independent Linear MMSE (IL-MMSE)** fusion rule, which uses a diagonal weighting matrix and the global estimator is given by Eq. (31).
- 345 3. **Single Coefficient MMSE (SC-MMSE)** fusion rule, which uses a scaled identity matrix for the weights and the global estimator is given by Eq. (39).
4. **Equal Weights Fusion (EWF)**, which simply averages the results of all the partial estimators.

As a performance measure we always use the averaged MSE between the estimated parameters and the ground-truth (which is known in all the experiments),

$$\text{MSE}(\text{dB}) = 10 \cdot \log_{10} \left[\frac{1}{DN_s} \sum_{i=1}^{N_s} (\hat{\mathbf{x}}_i - \mathbf{x})^\top (\hat{\mathbf{x}}_i - \mathbf{x}) \right], \quad (59)$$

350 where D is the dimension of the parameter space, N_s is the number of simulations performed, $\hat{\mathbf{x}}_i$ denotes the estimation obtained in the i -th simulation, and \mathbf{x} is the true value of the parameter.

5.1. Univariate Gamma Distributions

Let us consider first a univariate example, where exact calculations may be
 355 performed. This allows us to rule out any potential issue with the underlying

MC methods typically used to approximate the MMSE estimators (e.g., slow convergence and poor mixing), and concentrate on the performance of the proposed fusion rules. Let us assume that we have N i.i.d. observations distributed according to a Gamma PDF with known shape parameter $\alpha > 0$ and unknown rate parameter $\beta > 0$. Although we do not consider any particular application, the Gamma distribution has been used, for instance, to represent the amount of rainfall in meteorological applications. In this context, the value of the shape parameter is related to the dryness of the area under study [44]: wet areas correspond to large values of α , while dry areas are associated to small values of α . Hence, this example might correspond to the estimation of the dryness of an area from distributed measurements obtained using a sensor network.

In this scenario, the likelihood is given by

$$p(\mathbf{y}|\beta) = \prod_{n=1}^N p(y_n|\beta), \quad (60)$$

with

$$p(y_n|\beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} y_n^{\alpha-1} \exp(-\beta y_n). \quad (61)$$

The conjugate prior is also a Gamma PDF over β with shape parameter $\alpha_0 > 0$ and rate parameter $\beta_0 > 0$, and thus the global posterior density is another Gamma PDF with parameters $\alpha^* = \alpha_0 + N\alpha$ and $\beta^* = \beta_0 + \sum_{n=1}^N y_n$ [45]. The global MMSE estimator is given by the mean of the posterior PDF,

$$\hat{\beta}^{(\text{MMSE})} = \frac{\alpha^*}{\beta^*} = \frac{\alpha_0 + N\alpha}{\beta_0 + \sum_{n=1}^N y_n}, \quad (62)$$

and its variance is given by

$$\sigma_\beta^2 = \frac{\alpha^*}{(\beta^*)^2} = \frac{\alpha_0 + N\alpha}{\left(\beta_0 + \sum_{n=1}^N y_n\right)^2}. \quad (63)$$

For the partial estimators, the MMSE estimates and their variances are still given by (62) and (63), respectively, but replacing N by N_ℓ and taking the sum

only over the N_ℓ samples available to each of the ℓ estimators:

$$\hat{\beta}_\ell^{(\text{MMSE})} = \frac{\alpha_0 + N_\ell \alpha}{\beta_0 + \sum_{n \in \mathcal{Y}_\ell} y_n}, \quad (64a)$$

$$\sigma_\ell^2 = \frac{\alpha_0 + N_\ell \alpha}{(\beta_0 + \sum_{n \in \mathcal{Y}_\ell} y_n)^2}, \quad (64b)$$

where \mathcal{Y}_ℓ denotes the set of data available to the ℓ -th partial estimator.

Finally, let us remark that in this simple case we can also obtain the optimal non-linear fusion rule analytically. First of all, from (64a) note that

$$\sum_{n \in \mathcal{Y}_\ell} y_n = \frac{\alpha_0 + N_\ell \alpha}{\hat{\beta}_\ell^{(\text{MMSE})}} - \beta_0. \quad (65)$$

Then, using (65) inside (62) it can be easily shown that

$$\hat{\beta}^{(\text{MMSE})} = \frac{\alpha_0 + N\alpha}{\beta_0 + \sum_{\ell=1}^L \sum_{n \in \mathcal{Y}_\ell} y_n} = \frac{\alpha_0 + N\alpha}{\sum_{\ell=1}^L \frac{\alpha_0 + N_\ell \alpha}{\hat{\beta}_\ell^{(\text{MMSE})}} - (L-1)\beta_0} \quad (66)$$

Hence, the global MMSE estimator can be expressed as a simple non-linear combination of the partial MMSE estimators in this particular example.

370 We are interested now in analyzing the effect of the sample size, the number of partial estimators and the number of samples per estimator. Therefore, we test $N \in \{10^3, 10^4, 10^5, 10^6, 5 \cdot 10^6\}$ with an equal number of samples per partial estimator ranging from $N_\ell = 1$ (i.e., as many partial estimators as observations) up to $N_\ell = N$ (i.e., a single global estimator that uses). Note that $N_\ell = N$
 375 (i.e., $L = 1$) corresponds to the optimal centralized MMSE estimator that has access to all the available data and is computed according to (62). For each case, $N_s = 1000$ simulations are performed to average the results.

Figure 1 shows the typical performance of the optimal linear fusion rule (since we only have one parameter, all the fusion rules discussed in the paper
 380 are equivalent), the EWF rule (that assigns the same weight to all the partial estimators) and an empirical estimator that combines the optimal and the EWF estimates at the fusion center, $\hat{\beta}^{(\text{Avg})} = \frac{\hat{\beta}^{(\text{L-MMSE})} + \hat{\beta}^{(\text{EWF})}}{2}$. The optimal non-

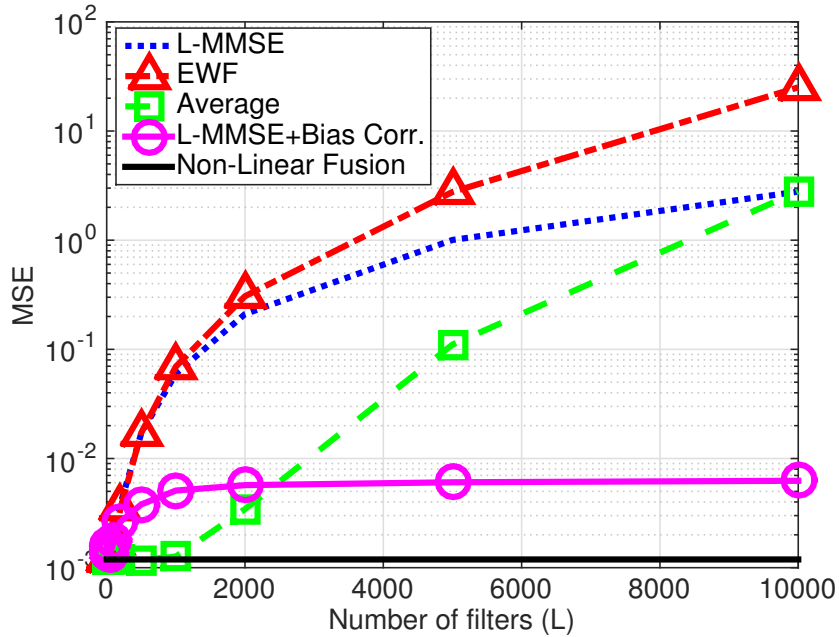


Figure 1: **Univariate Gamma Example:** MSE as a function of L using an improper prior for the three fusion rules considered: L-MMSE, EWF and their average. Note that $L = 1$ corresponds to the optimal centralized MMSE estimator given by (62). For the sake of comparison, we also include the MSE of the non-linear fusion rule of (66) and of the L-MMSE with the bias correction approach of [39].

linear fusion rule of Eq. (66) and the L-MMSE with the bias correction approach of [39] are also shown for comparison. In this example, the true parameters are $\alpha = 2$ and $\beta = 5$, and an improper prior is used by setting $\alpha_0 = \beta_0 = \epsilon$ with $\epsilon \rightarrow 0$.⁵ First of all, note that the performance of the optimal non-linear fusion rule is always the same regardless of the value of L . This is a direct consequence of the equivalence between (66) and (62). Then, note also that the optimal linear fusion rule always performs better than the EWF (as expected), especially when the number of partial estimators (a.k.a. filters) increases. The unexpected result is that combining the optimal fusion strategy and the EWF approach leads to a better performance than any of the two individual strategies. The reason for this good performance can be seen in Figure 2, which shows the estimated posterior

⁵In practice, we used $\epsilon = 0.01$ for the simulations.

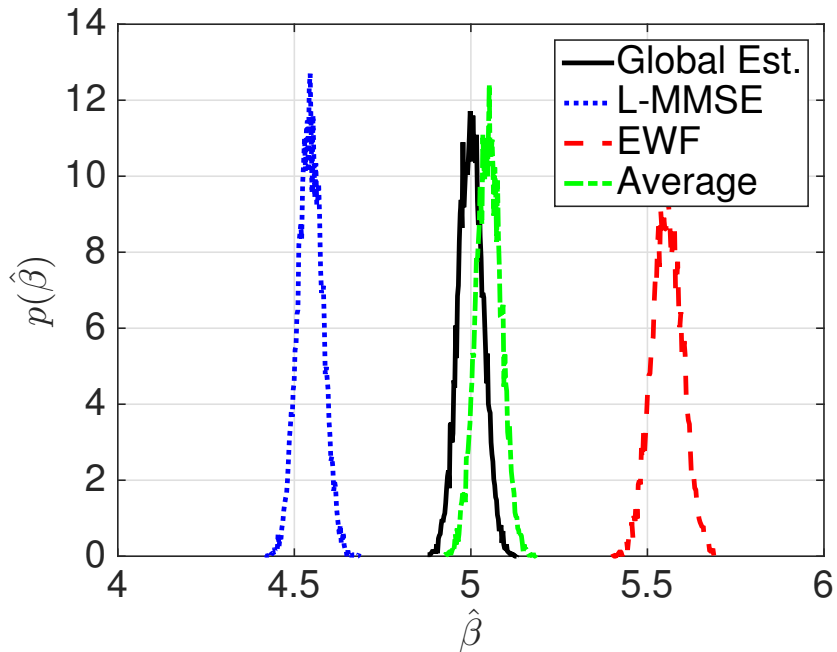


Figure 2: **Univariate Gamma Example:** Posterior density of the estimator for the three fusion rules considered (L-MMSE, EWF and their average), compared to the posterior of the global estimator.

PDFs of the estimators for the three fusion rules considered (L-MMSE, EWF
 395 and their average), compared to the posterior of the global estimator.⁶ It can be
 seen that the optimal linear fusion rule introduces a negative bias, whereas the
 EWF introduces a positive bias with a similar magnitude. Therefore, combining
 the two estimators leads to an average estimator with a reduced bias and thus
 a better performance. Unfortunately, although we can prove that the EWF
 400 tends to over-estimate β and the L-MMSE fusion rule tends to under-estimate
 it in this case (see Appendix C), this is not a general result and may not hold
 for other distributions. However, it illustrates how the bias may be reduced
 through the averaging of two simple fusion rules and maybe applicable in other
 problems (although a case by case analysis would have to be performed).

⁶Note that, given the large number of data available to the global estimator, the global posterior can be considered to be the “ground truth” for the evaluation of the PDFs obtained using the other fusion rules.

405 The second important issue in Figure 1 is related to the increase of the
MSE as the number of partial estimators increases. This is precisely due to
the fact that the bias increases as the number of samples per partial estimator
decreases (i.e., as the number of partial estimator increases for a fixed number
of data).⁷ This bias is caused by the mismatch between the “true” prior (in
410 the simulated scenario, an impulse centered around the true value $\beta = 5$) and
the prior assumed by the model (a Gamma PDF with parameters α_0 and β_0).
In order to see this, Figure 3 shows the evolution of the MSE with the number
of filters using a narrow prior (obtained setting $\beta_0 = \frac{\beta}{\varepsilon}$ and $\alpha_0 = \beta \times \beta_0$ for
 $\varepsilon = 0.01$) centered around the true value of β . In this case, all the estimators
415 are unbiased and the MSE decreases as we increase the number of partial es-
timators. These results, in an example where the exact MMSE estimator can
be obtained, highlight the importance of the prior in the Bayesian distributed
inference approach. Although this falls outside of the scope of this paper, let us
remark also that some approaches to reduce the bias in this type of distributed
420 inference problems have been recently proposed [46, 39].

Finally, Table 2 provides the complete picture regarding the evolution of the
MSE with the number of data and the number of data per partial estimator
for all the fusion rules considered. On the one hand, note that a minimum
amount of samples per estimator are required in order to attain a performance
425 that decreases linearly as a function of N for the optimal linear fusion and the
EWF. Otherwise, the bias dominates and nothing is gained by increasing N
(e.g., when $N_\ell = 10$ the MSE of the L-MMSE fusion rule is 0.0637 for $N = 10^3$
and 0.0568 for $N = 10^6$, implying that the MSE is only lowered by $\approx 10\%$ when
the total number of data increases by a factor of 1000). On the other hand,
430 note the excellent behaviour of the average fusion rule for all the cases: a linear
decrease in the MSE as a function of N is already clearly observed for $N_\ell = 20$.

⁷Note that we assumed unbiased partial estimators in the derivation of all the fusion rules.
However, when the number of samples per partial estimator is small this is no longer true.

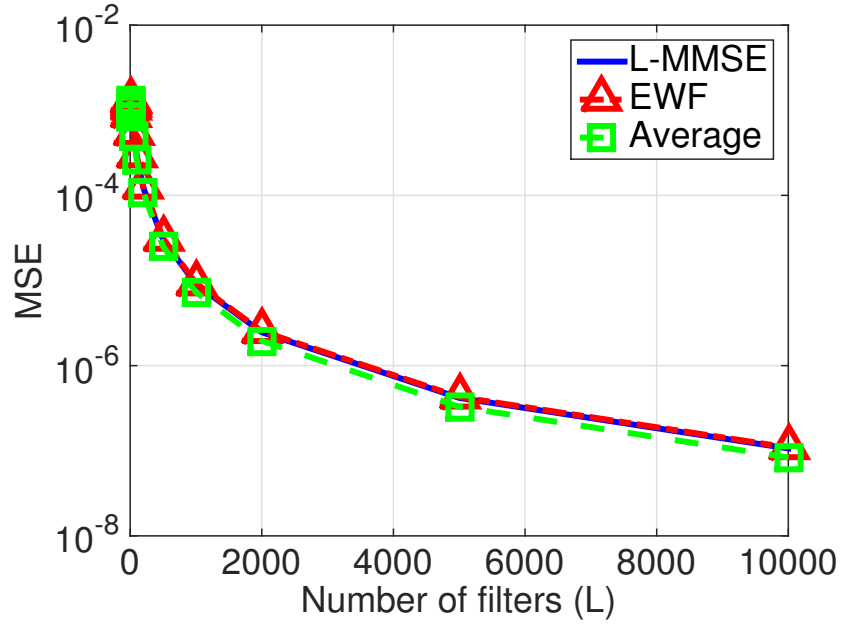


Figure 3: **Univariate Gamma Example:** MSE as a function of L using a narrow prior centered around the true value of β for the three fusion rules considered: L-MMSE, EWF and their average.

Table 2: **Univariate Gamma example:** Conditional MSE (averaged over 1000 independent runs) for the three fusion methods considered when $N \in \{10^3, 10^4, 10^5, 10^6\}$ and $N_\ell = N/L \in \{5, 10, 20, 50, 100, 200, N\}$.

Experiment		N_ℓ						
N	Estimator	5	10	20	50	100	200	N
10^3	EWF	0.3480	0.1011	0.0369	0.0193	0.0159	0.0148	0.0143
	L-MMSE	0.2042	0.0637	0.0243	0.0147	0.0141	0.0142	
	Average	0.0191	0.0152	0.0144	0.0143	0.0143	0.0143	
10^4	EWF	0.3067	0.0695	0.0172	0.0035	0.0017	0.0013	0.0012
	L-MMSE	0.2104	0.0598	0.0170	0.0038	0.0019	0.0014	
	Average	0.0034	0.0013	0.0012	0.0012	0.0012	0.0012	
10^5	EWF	0.3086	0.0695	0.0166	0.0027	0.0008	0.0003	0.0001
	L-MMSE	0.2058	0.0566	0.0149	0.0025	0.0007	0.0003	
	Average	0.0027	0.0003	0.0002	0.0001	0.0001	0.0001	
10^6	EWF	0.3081	0.0691	0.0164	0.0025	0.0006	0.0002	$1.2 \cdot 10^{-5}$
	L-MMSE	0.2069	0.0568	0.0149	0.0025	0.0006	0.0002	
	Average	0.0025	0.0002	$2 \cdot 10^{-5}$	$1.2 \cdot 10^{-5}$	$1.2 \cdot 10^{-5}$	$1.2 \cdot 10^{-5}$	

5.2. Multi-Variate Gaussian Distributions

Let us assume that we want to estimate a D -dimensional parameter vector, $\mathbf{x} \in \mathbb{R}^{D \times 1}$ for $D \geq 1$, from N independent D -dimensional observations that are divided into L groups, $\mathbf{Y} = [\mathbf{Y}_1, \dots, \mathbf{Y}_L]$, where $\mathbf{Y}_\ell = [\mathbf{y}_1^{(\ell)}, \dots, \mathbf{y}_{N_\ell}^{(\ell)}] \in \mathbb{R}^{D \times N_\ell}$ (with $\mathbf{y}_n^{(\ell)} \in \mathbb{R}^{D \times 1}$) is the set of observations available to the ℓ -th partial estimator. Let us consider a linear relationship among the $\mathbf{y}_n^{(\ell)}$ and \mathbf{x} :

$$\mathbf{y}_n^{(\ell)} = \mathbf{A}\mathbf{x} + \mathbf{w}_n^{(\ell)}, \quad (67)$$

where $\mathbf{w}_n^{(\ell)}$ follows a multi-variate Gaussian distribution, $\mathbf{w}_n^{(\ell)} \sim \mathcal{N}(\mathbf{w}_n^{(\ell)} | \mathbf{0}, \boldsymbol{\Sigma}_\ell)$. The likelihood is then another multi-variate Gaussian,

$$p(\mathbf{y}_1^{(\ell)}, \dots, \mathbf{y}_{N_\ell}^{(\ell)} | \mathbf{x}) = \prod_{n=1}^{N_\ell} \mathcal{N}(\mathbf{y}_n^{(\ell)} | \mathbf{A}\mathbf{x}, \boldsymbol{\Sigma}_\ell) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_\mathbf{y}^{(\ell)}, \mathbf{C}_\mathbf{y}^{(\ell)}), \quad (68)$$

where $\boldsymbol{\mu}_\mathbf{y}^{(\ell)} = \mathbf{A}^{-1} \cdot \left(\frac{1}{N_\ell} \sum_{n=1}^{N_\ell} \mathbf{y}_n^{(\ell)} \right)$ and $\mathbf{C}_\mathbf{y}^{(\ell)} = \frac{1}{N_\ell} (\mathbf{A}^\top \boldsymbol{\Sigma}_\ell \mathbf{A})^{-1}$. Let us consider also a multi-variate Gaussian prior, $p_0(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_\mathbf{x}^{(0)}, \mathbf{C}_\mathbf{x}^{(0)})$. The ℓ -th partial posterior is then another multi-variate Gaussian:

$$p(\mathbf{x} | \mathbf{y}_1^{(\ell)}, \dots, \mathbf{y}_{N_\ell}^{(\ell)}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_\mathbf{x}^{(\ell)}, \mathbf{C}_\mathbf{x}^{(\ell)}), \quad (69)$$

with

$$\mathbf{C}_\mathbf{x}^{(\ell)} = \left(\left(\mathbf{C}_\mathbf{x}^{(0)} \right)^{-1} + \left(\mathbf{C}_\mathbf{y}^{(\ell)} \right)^{-1} \right)^{-1}, \quad (70a)$$

$$\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}^{(\ell)} = \mathbf{C}_\mathbf{x}^{(\ell)} \left(\left(\mathbf{C}_\mathbf{x}^{(0)} \right)^{-1} \boldsymbol{\mu}_\mathbf{x}^{(0)} + \left(\mathbf{C}_\mathbf{y}^{(\ell)} \right)^{-1} \boldsymbol{\mu}_\mathbf{y}^{(\ell)} \right). \quad (70b)$$

Once more, we do not concentrate on any particular application, but this multi-variate Gaussian model arises in many statistical signal processing applications.

435 For the experiments, we set $\mathbf{A} = \mathbf{I}$ (implying that $\boldsymbol{\mu}_\mathbf{y}^{(\ell)} = \frac{1}{N_\ell} \sum_{n=1}^{N_\ell} \mathbf{y}_n^{(\ell)}$ and $\mathbf{C}_\mathbf{y}^{(\ell)} = \frac{1}{N_\ell} \boldsymbol{\Sigma}_\ell^{-1}$), and randomly generate the parameter vectors \mathbf{x} and the noise vectors $\mathbf{w}_n^{(\ell)}$. On the one hand, the elements of \mathbf{x} are randomly and

independently drawn from a zero-mean Gaussian with variance equal to 100, i.e., $\mathbf{x}_d \sim \mathcal{N}(0, 100)$. On the other hand, the covariance matrix for the noise is given by $\Sigma_\ell = \sigma_\ell^2 \mathbf{I} + \rho \mathbf{b}_\ell^\top \mathbf{b}_\ell$, with the elements of \mathbf{b}_ℓ randomly and independently drawn from a zero-mean and unit-variance Gaussian, $\sigma_\ell = |\varsigma_\ell|$ with $\varsigma_\ell \sim \mathcal{N}(0, 10)$, and $\rho \geq 0$ a user-dependent parameter that allows us to control the correlation among the outputs. In the simulations, we test the performance of the different fusion rules (EWF, SC-MMSE, IL-MMSE and L-MMSE) for different dimensions of the state space ($1 \leq D \leq 20$), correlation factors ρ ($\rho \in \{0, 0.2, 0.4\}$), and number of partial estimators ($L \in \{1, 2, 5, 10, 20, 50, 100, 200, 500, 1000, 2000, 5000\}$) in two different scenarios:

Sc1: The matrix $\mathbf{C}_y^{(\ell)}$ is assumed to be known.

Sc2: The matrix $\mathbf{C}_y^{(\ell)}$ is unknown and is substituted by the empirical estimate obtained from the data.

In all cases, we set $\boldsymbol{\mu}_0 = \mathbf{0}$ and $\mathbf{C}_0 = \mathbf{I}$ for the prior, which lead to the following expressions for the mean and covariance of the posterior PDF:

$$\mathbf{C}_x^{(\ell)} = \left(\mathbf{I} + \left(\mathbf{C}_y^{(\ell)} \right)^{-1} \right)^{-1}, \quad (71a)$$

$$\boldsymbol{\mu}_x^{(\ell)} = \mathbf{C}_x^{(\ell)} \left(\mathbf{C}_y^{(\ell)} \right)^{-1} \boldsymbol{\mu}_y^{(\ell)} = \boldsymbol{\mu}_y^{(\ell)} - \mathbf{C}_y^{(\ell)} \left(\mathbf{I} + \mathbf{C}_y^{(\ell)} \right)^{-1} \boldsymbol{\mu}_y^{(\ell)}, \quad (71b)$$

where we have used the matrix inversion lemma [6] in order to obtain the final expression of Eq. (71b).

The results (averaged over $N_s = 100$ realizations for each combination of parameters tested) are shown in Tables 3 and 4 for Sc1 and Sc2, respectively. The total number of data generated is $N = 10^5$, leading to a number of data per partial estimator $N_\ell \in \{10^5, 5 \cdot 10^4, 2 \cdot 10^4, 10^4, 5 \cdot 10^3, 2 \cdot 10^3, 10^3, 500, 200, 100, 50, 20\}$. Note that the first case (i.e., $L = 1$ and thus $N_\ell = N = 10^5$) corresponds to the global estimator that has access to all the data (i.e., the optimal centralized estimator). From these two tables we can extract the following conclusions:

- For $\rho = 0$ the IL-MMSE fusion rule shows the same performance as the

L-MMSE approach when the matrix $\mathbf{C}_y^{(\ell)}$ is known (Sc1), and even better when $\mathbf{C}_y^{(\ell)}$ is unknown (Sc2), due to the increased robustness in the estimation of D parameters per partial estimator instead of D^2 . This improvement can be appreciated especially when N_ℓ is small (e.g., when $N_\ell = 20$ and $\mathbf{C}_y^{(\ell)}$ has to be estimated from the data, the IL-MMSE scheme provides an improvement of 2.56 dB w.r.t. the L-MMSE rule).

- Both the IL-MMSE and the L-MMSE fusion rules outperform the SC-MMSE approach and the EWF substantially for $\rho = 0$. For example, when $\mathbf{C}_y^{(\ell)}$ is known (Sc1) the L-MMSE provides an improvement w.r.t. the EWF ranging from 1.25 dB when $N_\ell = 5 \cdot 10^4$ up to 6.79 dB for $N_\ell = 500$.
- For $\rho = 0.2$ the performance of all the fusion rules decreases in general, but the L-MMSE approach still outperforms the EWF by a similar amount as before (ranging from 1.41 dB when $N_\ell = 5 \cdot 10^4$ up to 6.31 dB for $N_\ell = 500$ in Sc1). The IL-MMSE rule decreases its performance w.r.t. the L-MMSE scheme, especially in Sc2, but its performance is still better than the EWF in general (up to 2.66 dB when $N_\ell = 20$ in Sc1, and up to 2.25 dB when $N_\ell = 10^4$ in Sc2).
- As the correlation increases (i.e., for $\rho = 0.4$), the L-MMSE fusion rule still outperforms the EWF by a similar amount as in the two previous cases. The performance of the IL-MMSE fusion rule deteriorates w.r.t. $\rho = 0.2$, but it can still be advantageous w.r.t. the EWF (e.g., 1.76 dB are gained in Sc1 when $N_\ell = 20$ and 1.19 dB are gained in Sc2 when $N_\ell = 10^4$).
- The SC-MMSE fusion rule attains a worse performance than the EWF in all cases. **Note that in this case we are replacing the optimal linear fusion rule (which requires $D \times D$ weighting matrices with $D = 10$ for each partial estimator) by weighting matrices of the form $\alpha_\ell \mathbf{I}$ with the α_ℓ estimated from the limited data available to each partial estimator. If the estimate of these coefficients is very sensitive to the data, then an increased MSE**

490

may be attained w.r.t. the EWF (which assumes that $\alpha_\ell = 1/L$ for all the partial estimators), as it happens here. However, in Section 5.3 we show that the SC-MMSE rule outperforms the EWF in the two localization examples considered.

495

- Regarding the behavior of the different approaches as the number of data per partial estimator (N_ℓ) decreases (i.e., the number of partial estimators (L) increases for a fixed total amount of data N), it seems to be rather constant for the SC-MMSE and EWF schemes, no pattern can be clearly seen for the IL-MMSE fusion rule, and there seems to be an optimal value around $N_\ell = 100$ (i.e., $L = 1000$) for the L-MMSE approach.

500

Finally, the performance of the different methods as the dimension of the state space (D) increases is shown in Figures 4 and 5. Note that the performance of all the methods remains quite stable (after an initial change for small values of D when $\rho = 0.2$) as the dimension of the state space increases.

Table 3: **Multi-variate Gaussian Example:** MSE in dB (averaged over $N_s = 100$ independent runs) for the four fusion methods considered when the matrix $\mathbf{\Sigma}_\ell$ is known, $N = 10^5$, $L \in \{2, 10, 50, 200, 1000, 5000\}$ (i.e., $N_\ell = N/L \in \{5 \cdot 10^4, 10^4, 2 \cdot 10^3, 500, 100, 20\}$), $D = 10$, and $\rho \in \{0, 0.2, 0.4\}$.

Experiment		N_ℓ					
Correlation	Estimator	$5 \cdot 10^4$	10^4	2000	500	100	20
$\rho = 0$	EWF	-46.43	-45.98	-46.07	-45.97	-46.07	-45.78
	SCMSE	-43.97	-41.59	-41.19	-41.07	-40.86	-41.01
	ILMSE	-47.68	-49.48	-51.59	-52.76	-52.45	-51.32
	LMSE	-47.68	-49.48	-51.59	-52.76	-52.45	-51.32
$\rho = 0.2$	EWF	-45.09	-45.35	-45.68	-45.61	-45.64	-45.81
	SCMSE	-42.75	-41.27	-40.89	-40.74	-40.80	-41.23
	ILMSE	-46.38	-47.60	-46.32	-44.07	-46.62	-48.47
	LMSE	-46.50	-48.72	-50.83	-51.92	-51.93	-51.15
$\rho = 0.4$	EWF	-45.16	-45.26	-45.45	-45.22	-45.34	-45.31
	SCMSE	-43.38	-41.21	-41.08	-40.49	-40.78	-40.65
	ILMSE	-45.77	-46.45	-45.66	-40.47	-44.16	-47.07
	LMSE	-45.99	-48.92	-50.58	-52.07	-52.12	-50.90

Table 4: **Multi-variate Gaussian Example:** MSE in dB (averaged over $N_s = 100$ independent runs) for the four fusion methods considered when the matrix Σ_ℓ is unknown, $N = 10^5$, $L \in \{2, 10, 50, 200, 1000, 5000\}$ (i.e., $N_\ell = N/L \in \{5 \cdot 10^4, 10^4, 2 \cdot 10^3, 500, 100, 20\}$), $D = 10$, and $\rho \in \{0, 0.2, 0.4\}$.

Experiment		N_ℓ					
Correlation	Estimator	$5 \cdot 10^4$	10^4	2000	500	100	20
$\rho = 0$	EFWF	-46.43	-45.98	-46.07	-45.97	-46.07	-45.78
	SCMSE	-43.97	-41.59	-41.19	-41.07	-40.83	-40.70
	ILMSE	-47.68	-49.47	-51.59	-52.84	-53.63	-51.62
	LMSE	-47.68	-49.45	-51.57	-52.69	-53.15	-49.06
$\rho = 0.2$	EFWF	-45.09	-45.35	-45.68	-45.61	-45.64	-45.81
	SCMSE	-42.75	-41.27	-40.89	-40.74	-40.77	-40.94
	ILMSE	-46.38	-47.60	-46.22	-42.54	-37.37	-45.20
	LMSE	-46.50	-48.72	-50.78	-51.84	-52.56	-48.75
$\rho = 0.4$	EFWF	-45.16	-45.26	-45.45	-45.22	-45.34	-45.31
	SCMSE	-43.38	-41.21	-41.08	-40.49	-40.74	-40.21
	ILMSE	-45.77	-46.45	-45.61	-38.65	-33.75	-42.90
	LMSE	-45.99	-48.90	-50.56	-52.08	-52.78	-48.22

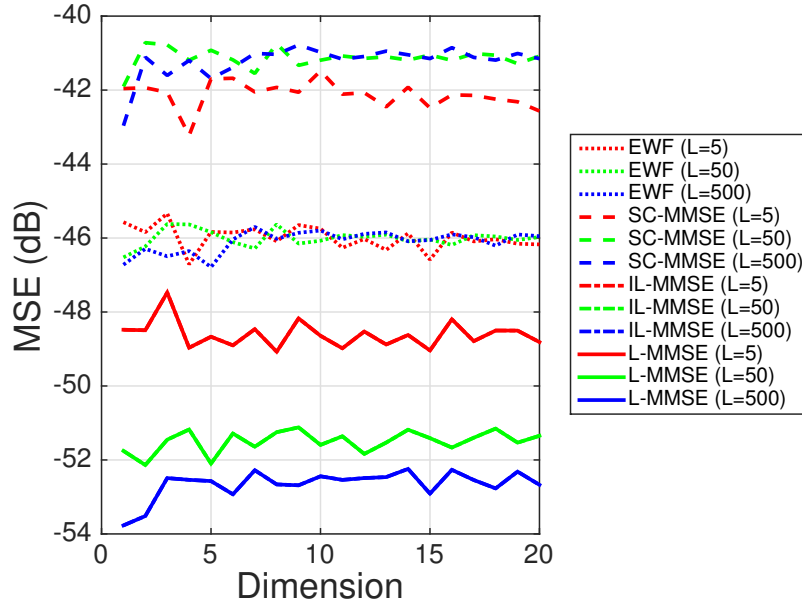
5.3. Localization in a Wireless Sensor Network

505 5.3.1. Case 1: Estimation of the Target's Position

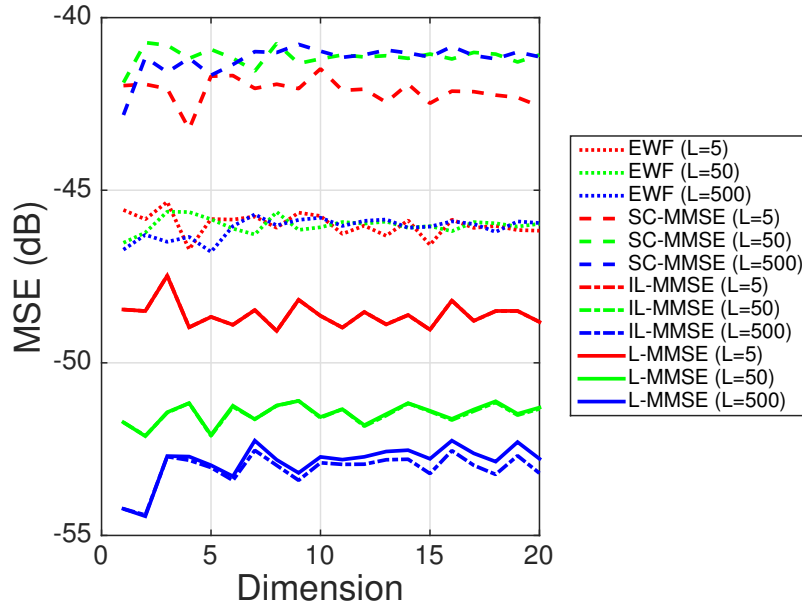
In this section, we address the problem of positioning a static target in the two-dimensional space of a wireless sensor network using only range measurements. More specifically, we consider a random vector $\mathbf{X} = [X_1, X_2]^\top$ to denote the target's position in the \mathbb{R}^2 plane. The position of the target is then a specific realization \mathbf{x} . The measurements are obtained from 6 range sensors located at $\mathbf{h}_1 = [1, -8]^\top$, $\mathbf{h}_2 = [8, 10]^\top$, $\mathbf{h}_3 = [-15, -7]^\top$, $\mathbf{h}_4 = [-8, 1]^\top$, $\mathbf{h}_5 = [10, 0]^\top$, and $\mathbf{h}_6 = [0, 10]^\top$. The measurement equations are

$$Y_j = -10 \log (\|\mathbf{x} - \mathbf{h}_j\|_2^2) + \Theta_j, \quad j = 1, \dots, 6, \quad (72)$$

where $\Theta_j \sim \mathcal{N}(\theta_j | \mathbf{0}, \omega_j^2 \mathbf{I})$, with $\omega_j = 5$ for $j \in \{1, 2, 3\}$ and $\omega_j = 20$ for $j \in \{4, 5, 6\}$. We simulate $N = 6000$ observations from the model ($\frac{N}{6} = 1000$ observations from each sensor) fixing $x_1 = x_2 = 3.5$. We consider a varying number of partial estimators L with $N_\ell = N/L$ for $1 \leq \ell \leq L$, and three

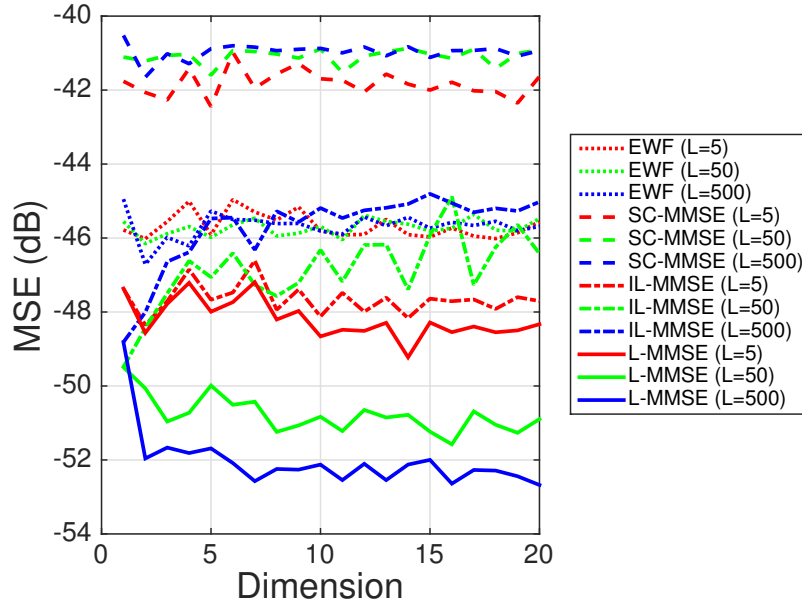


(a)

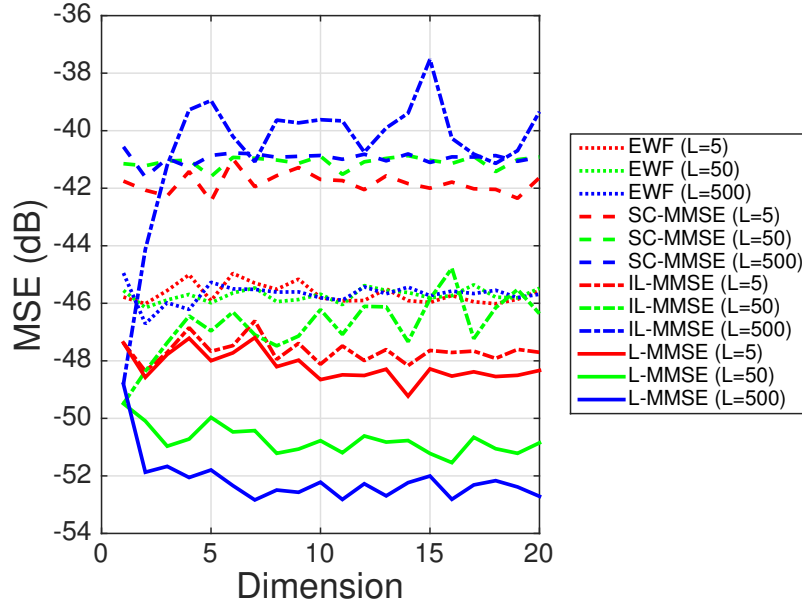


(b)

Figure 4: **Multi-variate Gaussian Example:** MSE as a function of D for $\rho = 0$, $N = 10^5$ and $N_\ell = N/L$. (a) Σ_ℓ known. (b) Σ_ℓ estimated from the data.



(a)



(b)

Figure 5: **Multi-variate Gaussian Example:** MSE as a function of D for $\rho = 0.2$, $N = 10^5$ and $N_\ell = N/L$. (a) Σ_ℓ known. (b) Σ_ℓ estimated from the data.

510 scenarios for splitting the data:

Sc1: Exactly $\frac{N}{6L}$ measurements from each sensor are provided to each partial estimator.

Sc2: The first $L/2$ estimators contain an equal number of observations from the first 3 sensors (the best ones), whereas the remaining $L/2$ estimators
515 work with measurements from the last 3 sensors (the noisiest ones).

Sc3: Measurements are randomly assigned to the estimators.

For each scenario, we run $M_C^{(\ell)} = 100$ MCMC independent parallel chains with length $T_C^{(\ell)} = 5000$, compute the MMSE estimates $\hat{x}_1^{(\ell)}$ and $\hat{x}_2^{(\ell)}$, and fuse these estimates into the final result. We compare the Equal Weights Fusion (EWF)
520 method, where each estimator is given the same weight, $1/L$, and the three fusion methods described in the paper. We repeat the experiments 50 times and average the results. The results, shown in Table 5 and Figures 6–8, confirm the good performance of the SC-MMSE and IL-MMSE estimators, which outperform the naive EWF and show an MSE similar to the optimal and more costly
525 L-MMSE. Regarding the three scenarios considered, we note that the best performance is obtained in the second case (with $\text{MSE}(\hat{\mathbf{x}}^{(\text{L-MMSE})}|\mathbf{y}) = 0.0021$), i.e., splitting the data into separate filters according to their quality. This opens up the possibility of performing a “smart” division of the data in order to optimize the performance.

530 Finally, in order to study the scaling behaviour of the fusion rules as N increases, we have also simulated the three scenarios for $N = 30000$, as well as Scenario 2 for $N = 600000$. The results, displayed in Table 6 and Figure 9, respectively, show that the performance of all the fusion rules scales roughly as a function of the number of samples, N .

535 **5.3.2. Case 2: Estimation of the Target’s Position and the Sensors’ Noise Variance**

In this section, we address a more complicated version of the previous case. We still consider the localization of a single target (placed at $x_1 = x_2 = 3.5$,

Table 5: **Localization Example (Case 1)**: MSE (averaged over $N_s = 50$ independent runs) for the three scenarios and the four fusion methods considered when $N = 6000$, $L \in \{5, 10, 25, 100, 200, 500, 1000\}$, and $N_\ell = N/L \in \{6, 12, 30, 60, 240, 600, 1200\}$.

Experiment		N_ℓ						
Scenario	Estimator	6	12	30	60	240	600	1200
Sc1	EWF	0.0041	0.0049	0.0065	0.0090	0.0167	0.0590	0.1192
	SCMSE	0.0039	0.0046	0.0063	0.0089	0.0166	0.0587	0.1191
	ILMSE	0.0038	0.0046	0.0063	0.0089	0.0166	0.0586	0.1188
	LMSE	0.0037	0.0045	0.0062	0.0088	0.0165	0.0584	0.1183
Sc2	EWF	0.0087	0.0053	0.0064	0.0104	0.0343	0.0648	0.1681
	SCMSE	0.0057	0.0034	0.0047	0.0092	0.0328	0.0628	0.1623
	ILMSE	0.0052	0.0031	0.0043	0.0085	0.0304	0.0588	0.1521
	LMSE	0.0037	0.0021	0.0028	0.0057	0.0210	0.0410	0.1107
Sc3	EWF	0.0078	0.0061	0.0068	0.0092	0.0169	0.0587	0.1181
	SCMSE	0.0060	0.0053	0.0066	0.0091	0.0168	0.0584	0.1180
	ILMSE	0.0055	0.0051	0.0065	0.0090	0.0168	0.0583	0.1177
	LMSE	0.0051	0.0048	0.0064	0.0090	0.0167	0.0582	0.1174

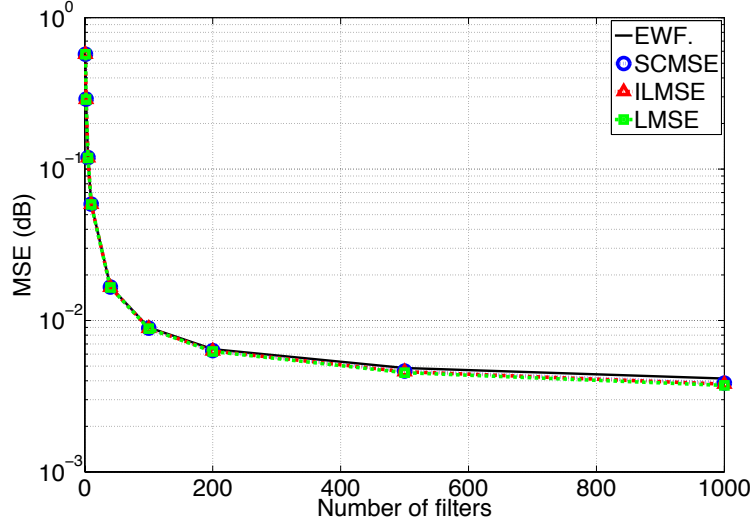


Figure 6: **Localization Example (Case 1)**: MSE as a function of L for Scenario 1 (Sc1) when $N = 6000$.

as before), using range-only measurements from 6 sensors placed at the same
540 locations as in case 1: $\mathbf{h}_1 = [1, -8]^\top$, $\mathbf{h}_2 = [8, 10]^\top$, $\mathbf{h}_3 = [-15, -7]^\top$, $\mathbf{h}_4 =$
 $[-8, 1]^\top$, $\mathbf{h}_5 = [10, 0]^\top$ and $\mathbf{h}_6 = [0, 10]^\top$. However, in addition to the estimation

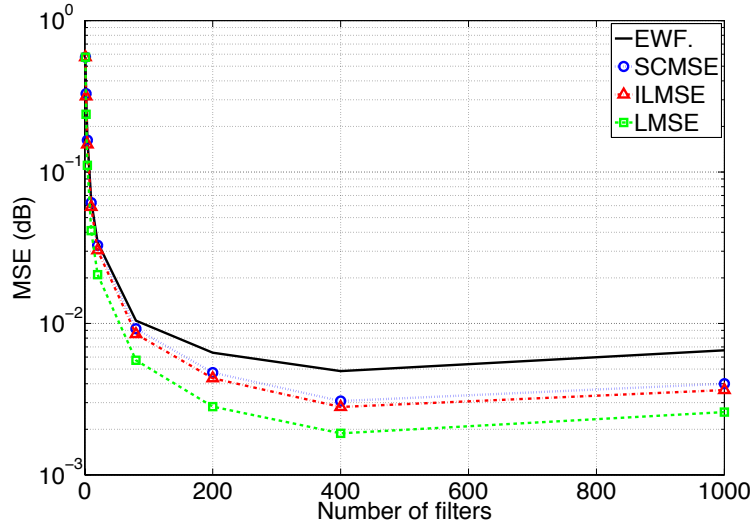


Figure 7: **Localization Example (Case 1):** MSE as a function of L for Scenario 2 (Sc2) when $N = 6000$.

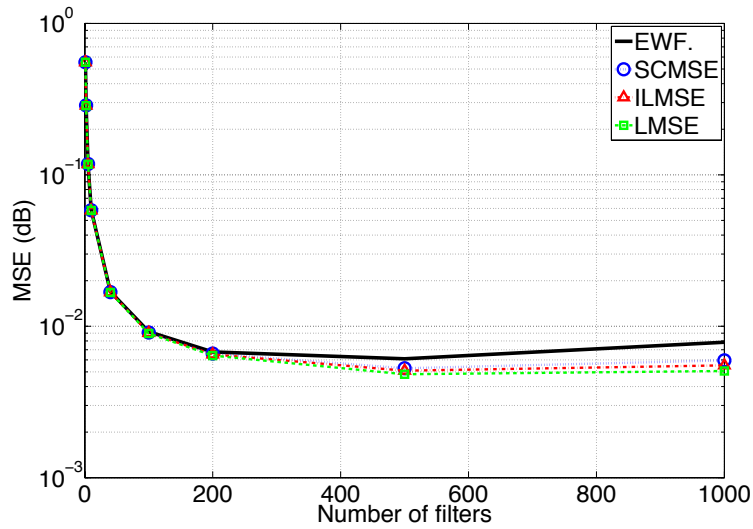


Figure 8: **Localization Example (Case 1):** MSE as a function of L for Scenario 3 (Sc3) when $N = 6000$.

of the target's location, now we also estimate the noise variance of the different sensors, implying that the dimension of the state space in this case is $D = 8$. The

Table 6: **Localization Example (Case 1)**: MSE (averaged over $N_s = 50$ independent runs) for the three scenarios and the four fusion methods considered when $N = 30000$, $L \in \{25, 50, 125, 500, 1000, 2500, 5000\}$, and $N_\ell = N/L \in \{6, 12, 30, 60, 240, 600, 1200\}$.

Experiment		N_ℓ						
Scenario	Estimator	6	12	30	60	240	600	1200
Sc1	EWF	0.0008	0.001	0.0013	0.0018	0.0033	0.0117	0.0231
	SC-MMSE	0.0008	0.0009	0.0013	0.0018	0.0033	0.0117	0.0230
	IL-MMSE	0.0008	0.0009	0.0013	0.0018	0.0033	0.0117	0.0230
	L-MMSE	0.0007	0.0009	0.0012	0.0017	0.0033	0.0116	0.0229
Sc2	EWF	0.0007	0.0009	0.0012	0.0018	0.0036	0.0131	0.0335
	SC-MMSE	0.0004	0.0006	0.0009	0.0015	0.0033	0.0125	0.0323
	IL-MMSE	0.0004	0.0005	0.0009	0.0014	0.0031	0.0118	0.0304
	L-MMSE	0.0003	0.0003	0.0006	0.0009	0.0021	0.0082	0.0214
Sc3	EWF	0.0018	0.0011	0.0013	0.0018	0.0033	0.0118	0.0229
	SC-MMSE	0.0014	0.001	0.0013	0.0018	0.0033	0.0118	0.0228
	IL-MMSE	0.0013	0.001	0.0013	0.0018	0.0033	0.0118	0.0228
	L-MMSE	0.0012	0.001	0.0013	0.0018	0.0033	0.0117	0.0227

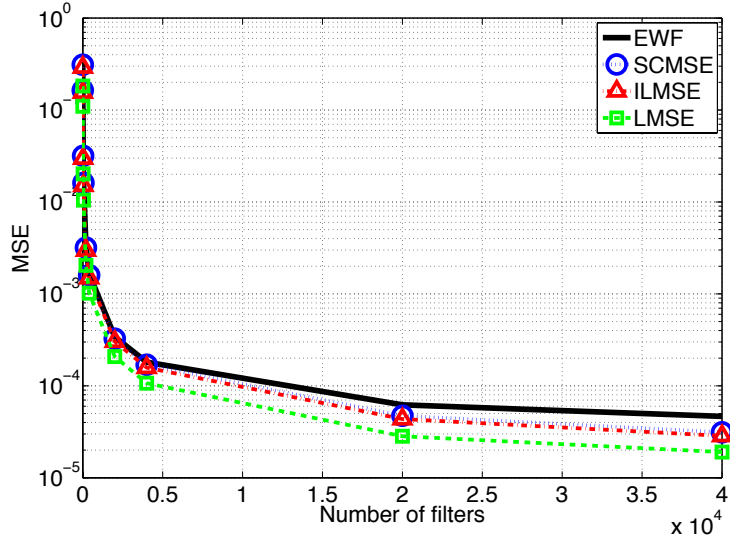


Figure 9: **Localization Example (Case 1)**: MSE as a function of L for Scenario 2 (Sc2) when $N = 600000$.

ground-truth values of the noise variances used in the simulations are $\omega_1 = 1$,
 $\omega_2 = 2$, $\omega_3 = 1$, $\omega_4 = 0.5$, $\omega_5 = 3$, and $\omega_6 = 0.2$. Furthermore, we investigate
545 the performance of the proposed approaches when a different number of data is

available to each the partial estimators. Namely, we generate 87 data from each sensor (i.e., $N = 87 \times 6 = 522$ data overall) according to the model in Eq. (72) and split them randomly among $L = 5$ partial estimators in such a way that
550 $N_1 = 12$, $N_2 = 30$, $N_3 = 60$, $N_4 = 120$, and $N_5 = 300$.

We test the performance of the optimal L-MMSE estimator, comparing it to the EWF approach and the two novel efficient fusion methods proposed (IL-MMSE and SC-MMSE). Moreover, we also test two heuristic fusion rules that take into account the unequal number of data available to each of the partial
555 estimators: a linearly proportional weights fusion (L-PWF) scheme that assigns weights to the partial estimates as $w_\ell = N_\ell/N$ (i.e., in our case $w_1 = 2/87$, $w_2 = 5/87$, $w_3 = 10/87$, $w_4 = 20/87$, and $w_5 = 50/87$); and a square root PWF (SR-PWF) scheme that assigns weights as $w_\ell = \sqrt{N_\ell}/S$, with $S = \sum_{\ell=1}^L \sqrt{N_\ell}$ (i.e., in our case, $S = \sqrt{2} + \sqrt{5} + \sqrt{10} + \sqrt{20} + \sqrt{50} \approx 18.3558$, $w_1 = \sqrt{2}/S \approx 0.0770$,
560 $w_2 = \sqrt{5}/S \approx 0.1218$, $w_3 = \sqrt{10}/S \approx 0.1723$, $w_4 = \sqrt{20}/S \approx 0.2436$, and $w_5 = \sqrt{50}/S \approx 0.3852$). The rationale behind these two heuristic fusion rules is that the estimates provided by estimators that have more data available should be trusted more, since the precision in the estimation is directly proportional to the number of data (L-PWF implicitly assumes that the precision is proportional
565 to N_ℓ , whereas SR-PWF assumes that it is proportional to $\sqrt{N_\ell}$).

Figure 10 shows the MSE for all the different fusion rules considered as a function of the number of samples of the Markov chain. First of all, it is remarkable the increase in performance as the number of iterations of the chain increases: around $5 \cdot 10^3$ iterations are required in order to attain a good per-
570 formance in all cases. Then, note the similar performance of the three heuristic approaches (EWF, L-PWF and SR-PWF), with EWF performing slightly better than the other two (showing that having weights directly proportional to N_ℓ or $\sqrt{N_\ell}$ is not a good idea in this case). Finally, note also that the IL-MMSE fusion rule always outperforms the optimal L-MMSE fusion, whereas the SC-MMSE approach performs better than the L-MMSE fusion for less than 10^4
575 iterations. This shows the increased robustness of the two proposed efficient fusion rules, which require the estimation of a reduced number of parameters

from the noisy data (in this case, $D = 8$ parameters per partial estimator for the IL-MMSE approach and just one parameter for the SC-MMSE fusion rule, instead of the $D^2 = 64$ parameters per partial estimator of the L-MMSE rule).
 Furthermore, both of the efficient fusion rules proposed are able to provide an improved performance w.r.t. the naive EWF. For instance, using $5 \cdot 10^4$ iterations of the chain the SC-MMSE fusion rule achieves ≈ 1 dB reduction in MSE w.r.t. the EWF with the transmission of a single weight per partial estimator, whereas the IL-MMSE approach attains a reduction of 1.69 dB by transmitting $D = 8$ coefficients per partial estimator.

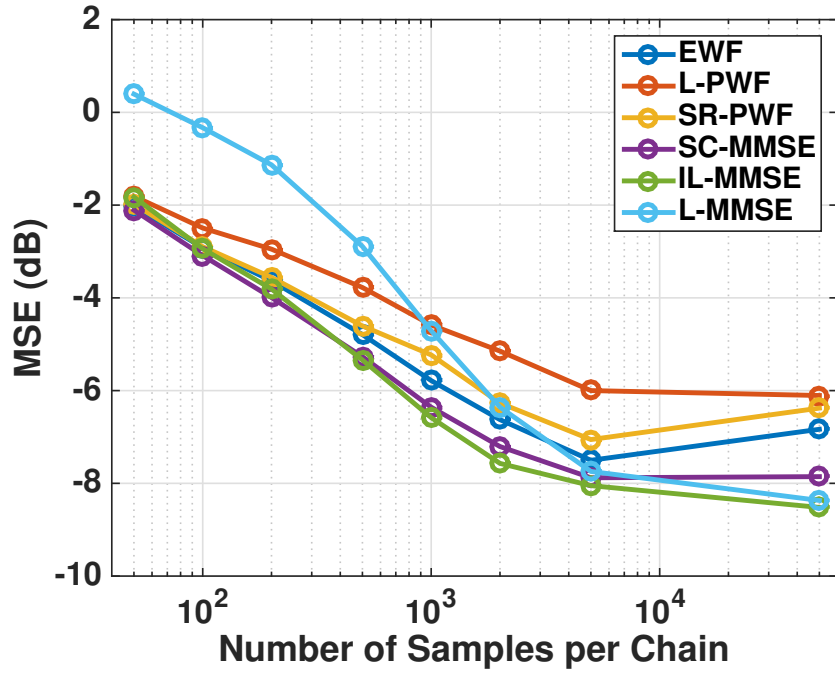


Figure 10: **Localization Example (Case 2)**: MSE as a function of the number of samples per chain for all the different fusion rules considered.

6. Conclusions

In this paper, we have addressed the linear fusion of unbiased and independent partial estimators, focusing on the scenario where no communication is allowed

590 among them and a large number of data is available to the partial estimators.

The optimal linear fusion rule has been derived and two novel efficient linear combination schemes, that achieve an excellent trade-off between the amount of information that has to be transmitted to the fusion center and the performance in terms of mean squared error (MSE) of the final estimator, have been

595 proposed. Both a constrained optimization point of view and a Bayesian perspective have been provided for these three fusion rules, allowing us to explore their connections and to explain their good performance even in situations where the assumptions that lead to them are not strictly fulfilled. All the methods were tested through computer simulations by applying them to three problems.

600 First of all, we considered a univariate problem, where all the posterior densities followed a Gamma PDF and all the estimators could be computed analytically. We use this example to analyze the performance of the fusion rules as the number of samples per estimator decreases and the partial estimators cannot be considered unbiased any more. Then, we address a parameter estimation prob-

605 lem in a multi-variate Gaussian model, where we increase the dimension of the state space up to $D = 20$, showing the good performance and robustness of the two efficient fusion rules proposed. Finally, we tackle a localization problem with one target and six sensors whose measurements were processed using several parallel filters. Here we consider two cases: localization of the target

610 when the noise characteristics of the sensors are known (dimension of the state space $D = 2$), and joint localization of the target and estimation of the noise characteristics of the sensors (dimension of the state space $D = 8$). The new fusion methods show a performance equivalent to the optimal linear combination (sometimes even better) with a reduced computational cost. Furthermore,

615 it has been shown that splitting the data can be advantageous in terms of attaining a good MSE with a reduced computational cost, but only when the bias in the partial estimators can be controlled. In future works we plan to address bias correction approaches, extending our previous work of [39], to deal with the small sample bias that appears when the number of samples per partial

620 estimator is small (even if the total amount of data may be very large). Alterna-

tively, the bias might be included in the formulation of the optimal linear fusion schemes, although this is a very challenging open issue. Some other interesting areas of research are non-linear fusion techniques and the development of fusion schemes where the partial Monte Carlo estimators are allowed to exchange a
625 reduced amount of information.

Acknowledgements

This work has been supported by Ministerio de Economía y Competitividad of Spain through the MIMOD-PLC (TEC2015-64835-C3-3-R) and KERMES (TEC2016-81900-REDT/AEI) projects; Ministerio de Educación, Cultura y De-
630 porte of Spain under CAS15/00350 grant; Universidad Politécnica de Madrid through a mobility grant for a short visit of D. Luengo to Stony Brook University; the BBVA Foundation through project MG-FIAR (“I Convocatoria de Ayudas Fundación BBVA a Investigadores, Innovadores y Creadores Culturales”); the National Science Foundation under Award CCF-1617986; and
635 the European Research Council (ERC) through the ERC Consolidator Grant SEDAL ERC-2014-CoG 647423.

References

- [1] H. L. Van Trees, Detection, Estimation and Modulation Theory, John Wiley and Sons, Hoboken, NJ (USA), 1968.
- 640 [2] G. Casella, R. L. Berger, Statistical inference, Duxbury, 2002.
- [3] L. L. Scharf, Statistical Signal Processing, Addison-Wesley, Reading, MA (USA), 1991.
- [4] S. M. Kay, Fundamentals of Statistical Signal Processing: Estimation Theory, Prentice Hall, Upper Saddle River, NJ (USA), 1993.

- 645 [5] J. M. Mendel, *Lessons in estimation theory for signal processing, communications, and control*, Pearson Education, New York, NY (USA), 1995.
- [6] C. Rasmussen, C. Williams, *Gaussian Processes for Machine Learning*, MIT Press, 2006.
- [7] N. L. Hjort, C. Holmes, P. Müller, S. G. Walker, *Bayesian nonparametrics*,
650 Vol. 28, Cambridge University Press, Cambridge (UK), 2010.
- [8] J. D. Gibbons, S. Chakraborti, *Nonparametric statistical inference*, Springer, 2011.
- [9] G. B. Giannakis, F. Bach, R. Cendrillon, M. Mahoney, J. Neville, *Signal processing for big data (special issue)*, *IEEE Signal Processing Magazine*
655 31 (5) (2014) 15–111.
- [10] P. M. Djuric, S. J. Godsill, *Special issue on Monte Carlo methods for statistical signal processing*, *IEEE Transactions on Signal Processing* 50 (2) (2002) 173–173.
- [11] C. P. Robert, G. Casella, *Monte Carlo Statistical Methods*, Springer, 2004.
- 660 [12] A. Doucet, X. Wang, *Monte Carlo methods for signal processing: a review in the statistical signal processing context*, *Signal Processing Magazine, IEEE* 22 (6) (2005) 152–170.
- [13] N. Chopin, *A sequential particle filter method for static models*, *Biometrika* 89 (3) (2002) 539–552.
- 665 [14] M. A. Suchard, Q. Wang, C. Chan, J. Frelinger, A. Cron, M. West, *Understanding GPU programming for statistical computation: Studies in massively parallel massive mixtures*, *Journal of Computational and Graphical Statistics* 19 (2) (2010) 419–438.
- [15] A. Lee, C. Yau, M. B. Giles, A. Doucet, C. C. Holmes, *On the utility of graphics cards to perform massively parallel simulation of advanced Monte*
670

- Carlo methods, *Journal of computational and graphical statistics* 19 (4) (2010) 769–789.
- [16] J. Dean, S. Ghemawat, Mapreduce: simplified data processing on large clusters, *Communications of the ACM* 51 (1) (2008) 107–113.
- 675 [17] K. F. Wallis, Combining forecasts – forty years later, *Applied Financial Economics* 21 (1–2) (2011) 33–41.
- [18] J. M. Bates, C. W. Granger, The combination of forecasts, *Operational Research Quarterly* 20 (4) (1969) 451–468.
- [19] J. P. Dickinson, Some statistical results in the combination of forecasts,
680 *Operational Research Quarterly* 24 (1975) 253–260.
- [20] R. F. Bordley, The combination of forecasts: A Bayesian approach, *Journal of the Operational Research Society* 33 (2) (1982) 171–174.
- [21] F. Lavancier, P. Rochet, A general procedure to combine estimators, arXiv preprint arXiv:1401.6371 (2014) 1–38.
- 685 [22] G. M. Allenby, Cross-validation, the Bayes theorem, and small-sample bias, *Journal of Business & Economic Statistics* 8 (2) (1990) 171–178.
- [23] J. B. Predd, S. R. Kulkarni, H. V. Poor, Distributed learning in wireless sensor networks, *IEEE Signal Processing Magazine* 23 (4) (2006) 56–69.
- [24] J.-J. Xiao, A. Ribeiro, Z.-Q. Luo, G. B. Giannakis, Distributed
690 compression-estimation using wireless sensor networks, *IEEE Signal Processing Magazine* 23 (4) (2006) 27–41.
- [25] A. Swami, Q. Zhao, Y.-W. Hong, L. Tong (Eds.), *Wireless Sensor Networks: Signal Processing and Communications Perspectives*, John Wiley and Sons, 2007.
- 695 [26] M. Cetin, L. Chen, J. W. F. III, A. T. Ihler, R. L. Moses, M. J. Wainwright, A. S. Willsky, Distributed fusion in sensor networks, *IEEE Signal Processing Magazine* 23 (4) (2006) 42–55.

- [27] R. Olfati-Saber, J. A. Fax, R. M. Murray, Consensus and cooperation in networked multi-agent systems, *Proceedings of the IEEE* 95 (1) (2007) 215–233.
- 700
- [28] A. G. Dimakis, S. Kar, J. F. Moura, M. G. Rabbat, A. Scaglione, Gossip algorithms for distributed signal processing, *Proceedings of the IEEE* 98 (11) (2010) 1847–1864.
- [29] F. S. Cattivelli, A. H. Sayed, Diffusion LMS strategies for distributed estimation, *IEEE Transactions on Signal Processing* 58 (3) (2010) 1035–1048.
- 705
- [30] D. J. Wilkinson, *Parallel Bayesian computation*, *Statistics Textbooks and Monographs* 184.
- [31] B. Wilkinson, M. Allen, *Parallel programming: techniques and applications using networked workstations and parallel computers*, Prentice-Hall, Inc., 1999.
- 710
- [32] Z. Huang, A. Gelman, *Sampling for Bayesian computation with large datasets*, Available at SSRN 1010107.
- [33] S. L. Scott, A. W. Blocker, F. V. Bonassi, H. A. Chipman, E. I. George, R. E. McCulloch, Bayes and big data: The consensus Monte Carlo algorithm, *International Journal of Management Science and Engineering Management* 11 (2) (2016) 78–88.
- 715
- [34] W. Neiswanger, C. Wang, E. Xing, Asymptotically exact, embarrassingly parallel MCMC, *arXiv:1311.4780v2* (21 Mar. 2014) 1–16.
- [35] X. Wang, D. B. Dunson, Parallelizing MCMC via Weierstrass sampler, *arXiv:1312.4605v2* (25 May 2014) 1–35.
- 720
- [36] M. Rabinovich, E. Angelino, M. I. Jordan, Variational consensus monte carlo, in: *Advances in Neural Information Processing Systems (NIPS)*, 2015, pp. 1207–1215.

- [37] X. Wang, F. Guo, K. A. Heller, D. B. Dunson, Parallelizing mcmc with
725 random partition trees, in: *Advances in Neural Information Processing Systems (NIPS)*, 2015, pp. 451–459.
- [38] D. Luengo, L. Martino, V. Elvira, M. Bugallo, Efficient linear combination of partial Monte Carlo estimators, in: *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2015, pp. 4100–4104.
- 730 [39] D. Luengo, L. Martino, V. Elvira, M. Bugallo, Bias correction for distributed Bayesian estimators, in: *IEEE 6th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, IEEE, 2015, pp. 253–256.
- [40] A. Doucet, X. Wang, Monte Carlo methods for signal processing, IEEE
735 *Signal Processing Magazine* 22 (6) (2005) 152–170.
- [41] S. Boyd, L. Vandenberghe, *Convex optimization*, Cambridge university press, 2004.
- [42] L. Le Cam, *Asymptotic Methods in Statistical Decision Theory*, Springer-Verlag, New York, NY (USA), 1986.
- 740 [43] A. W. van der Vaart, *Asymptotic Statistics*, Cambridge University Press, 1998.
- [44] G. J. Husak, J. Michaelson, C. Funk, Use of the gamma distribution to represent monthly rainfall in africa for drought monitoring applications, *International Journal of Climatology* 27 (7) (2007) 935–944.
- 745 [45] D. Fink, A compendium of conjugate priors, Technical Report (May 1997) 1–47.
- [46] H. Strathmann, D. Sejdinovic, M. Girolami, Unbiased Bayes for big data: Paths of partial posteriors, arXiv preprint arXiv:1501.03326.
- [47] T. M. Cover, J. A. Thomas, *Elements of information theory*, John Wiley & Sons, 2012.
750

- [48] J. L. Fleiss, The teacher's corner: A note on the expectation of the reciprocal of a random variable, *The American Statistician* 20 (1) (1966) 25–25.
- [49] A. N. Kolmogorov, S. V. Fomin, *Introductory real analysis*, Dover Publications Inc., 1975.
- 755 [50] N. L. García, J. L. Palacios, On inverse moments of nonnegative random variables, *Statistics & Probability Letters* 53 (3) (2001) 235–239.

Appendices

A. Solution of the Constrained Optimization

A.1. General Case: L-MMSE Fusion Rule

Let us consider the optimal L-MMSE fusion rule. From Eqs. (19a) and (19b), and applying the method of the Lagrange multipliers, the cost function that has to be minimized is

$$J_{\text{L-MMSE}} = \sum_{k=1}^L \text{Tr} \left(\mathbf{\Lambda}_k \mathbf{C}_{\mathbf{x}}^{(k)} \mathbf{\Lambda}_k^{\top} \right) + \lambda \left(\sum_{k=1}^L \mathbf{\Lambda}_k - \mathbf{I} \right), \quad (73)$$

where \mathbf{I} denotes the $D \times D$ identity matrix. Differentiating w.r.t. $\mathbf{\Lambda}_\ell$ ($\ell = 1, \dots, L$) and equating the result to zero we obtain a set of LD^2 equations,

$$\frac{\partial J_{\text{L-MMSE}}}{\partial \mathbf{\Lambda}_\ell} = 2\mathbf{\Lambda}_\ell \left[\mathbf{C}_{\mathbf{x}}^{(\ell)} \right]^{\top} + \lambda \mathbf{I} = \mathbf{0}, \quad (74)$$

whereas differentiating w.r.t. λ and equating the result to zero we obtain one additional equation:

$$\frac{\partial J_{\text{L-MMSE}}}{\partial \lambda} = \sum_{k=1}^L \mathbf{\Lambda}_k - \mathbf{I} = \mathbf{0}. \quad (75)$$

And now, it can be easily checked that the weighting matrix given by (48),

$$\mathbf{\Lambda}_\ell^* = \left[\sum_{k=1}^L (\mathbf{C}_x^{(k)})^{-1} \right]^{-1} (\mathbf{C}_x^{(\ell)})^{-1},$$

and the regularization parameter,

$$\lambda^* = -2 \left[\sum_{k=1}^L (\mathbf{C}_x^{(k)})^{-1} \right]^{-1},$$

760 fulfill both (74) and (75), and are thus the unique solution of the convex optimization problem posed by Eqs. (19a) and (19b).⁸

A.2. Particular Cases: IL-MMSE and SC-MMSE Fusion

For the IL-MMSE fusion rule we have $\mathbf{\Lambda}_\ell = \mathbf{D}_\ell = \text{diag}(\alpha_{\ell,1}, \dots, \alpha_{\ell,D})$. Hence, (74) and (75) become, respectively,

$$\frac{\partial J_{\text{IL-MMSE}}}{\partial \alpha_{\ell,d}} = 2\sigma_{\ell,d}^2 \alpha_{\ell,d} + \lambda_d = 0, \quad \text{for } \ell = 1, \dots, L \quad \text{and } d = 1, \dots, D; \quad (76a)$$

$$\frac{\partial J_{\text{IL-MMSE}}}{\partial \lambda_d} = \sum_{k=1}^L \alpha_{k,d} - 1 = 0, \quad \text{for } d = 1, \dots, D, \quad (76b)$$

and the optimal solution of the convex estimation problem is now given by (38):

$$\alpha_{\ell,d}^* = \frac{\sigma_{\ell,d}^{-2}}{\sum_{k=1}^L \sigma_{k,d}^{-2}}. \quad (77)$$

⁸Note that $\mathbf{C}_x^{(k)} = (\mathbf{C}_x^{(k)})^\top$ (for $k = 1, \dots, L$), since $\mathbf{C}_x^{(k)}$ is a covariance matrix, and thus $[(\mathbf{C}_x^{(k)})^{-1}]^\top = (\mathbf{C}_x^{(k)})^{-1}$.

Finally, for the SC-MMSE fusion rule we have $\mathbf{A}_\ell = \alpha_\ell \mathbf{I}$, implying that

$$\frac{\partial J_{\text{SC-MMSE}}}{\partial \alpha_{\ell,d}} = 2\alpha_\ell T_\ell + \lambda = 0, \quad \text{for } \ell = 1, \dots, L; \quad (78a)$$

$$\frac{\partial J_{\text{SC-MMSE}}}{\partial \lambda} = \sum_{k=1}^L \alpha_k - 1 = 0, \quad (78b)$$

with T_ℓ denoting the trace of $\mathbf{C}_\mathbf{x}^{(\ell)}$, and the optimal solution now corresponds to $\lambda^* = \frac{2}{\sum_{k=1}^L T_k^{-1}}$ and

$$\alpha_\ell^* = \frac{T_\ell^{-1}}{\sum_{k=1}^L T_k^{-1}}, \quad (79)$$

which are the optimal coefficients given in Eq. (30).

B. Optimal Independent and Isotropic Multivariate Gaussian Approximations

765

Let us define $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\hat{\mathbf{x}}_\ell, \mathbf{\Sigma}_p)$ and $q(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\hat{\mathbf{x}}_\ell, \mathbf{\Sigma}_q)$. The Kullback-Leibler (KL) divergence, $D_{KL} \geq 0$, is a standard measure of the discrepancy between two distributions [47]: the larger the value of D_{KL} the more different the two distributions, with $D_{KL} = 0$ indicating that they are equal almost everywhere. The KL divergence between $p(\mathbf{x})$ and $q(\mathbf{x})$ is defined as

$$D_{KL}(p||q) = \int_{\mathcal{X}} p(\mathbf{x}) \ln \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}. \quad (80)$$

Since both $p(\mathbf{x})$ and $q(\mathbf{x})$ are Gaussian PDFs, we have

$$\ln \frac{p(\mathbf{x})}{q(\mathbf{x})} = \frac{D}{2} \ln |\mathbf{\Sigma}_q| - \frac{D}{2} \ln |\mathbf{\Sigma}_p| - \frac{1}{2} (\mathbf{x} - \hat{\mathbf{x}}_\ell)^\top (\mathbf{\Sigma}_p^{-1} - \mathbf{\Sigma}_q^{-1}) (\mathbf{x} - \hat{\mathbf{x}}_\ell), \quad (81)$$

and the KL divergence becomes

$$D_{KL}(p||q) = K + \frac{D}{2} \ln |\mathbf{\Sigma}_q| + \frac{1}{2} \int_{\mathcal{X}} (\mathbf{x} - \hat{\mathbf{x}}_\ell)^\top \mathbf{\Sigma}_q^{-1} (\mathbf{x} - \hat{\mathbf{x}}_\ell) \mathcal{N}(\mathbf{x}|\hat{\mathbf{x}}_\ell, \mathbf{\Sigma}_p) d\mathbf{x}, \quad (82)$$

where K is a constant that includes all the terms that do not depend on Σ_q . Now, let us consider the two particular cases that we are interested in: the approximation of an arbitrary multi-variate Gaussian PDF using an independent and an isotropic Gaussian PDF, respectively.

770 B.1. Isotropic Approximation

Let us consider $\Sigma_p = \mathbf{C}_x^{(\ell)}$ and $\Sigma_q = \sigma_\ell^2 \mathbf{I}$, so that $|\Sigma_q| = \sigma_\ell^{2D}$ and $\Sigma_q^{-1} = \frac{1}{\sigma_\ell^2} \mathbf{I}$. Then, the KL divergence in (82) becomes

$$D_{KL}(p||q) = K + \frac{D^2}{2} \ln(\sigma_\ell^2) + \frac{1}{2\sigma_\ell^2} \text{Tr} \left(\mathbf{C}_x^{(\ell)} \right). \quad (83)$$

In order to minimize the KL divergence in (83), we take the derivative w.r.t. σ_ℓ^2 and equate it to zero:

$$\frac{\partial D_{KL}(p||q)}{\partial \sigma_\ell^2} = \frac{D^2}{2\sigma_\ell^2} - \frac{1}{2\sigma_\ell^4} \text{Tr} \left(\mathbf{C}_x^{(\ell)} \right) = 0. \quad (84)$$

The solution of this equation is

$$\sigma_\ell^2 = \frac{1}{D^2} \text{Tr} \left(\mathbf{C}_x^{(\ell)} \right), \quad (85)$$

and it can be easily checked (by taking the second derivative of $D_{KL}(p||q)$) that it corresponds to a minimum. Hence, the best isotropic approximation of $p(\mathbf{x}|\hat{\mathbf{x}}_\ell) = \mathcal{N}(\mathbf{x}|\hat{\mathbf{x}}_\ell, \mathbf{C}_x^{(\ell)})$ is $q(\mathbf{x}|\hat{\mathbf{x}}_\ell) = \mathcal{N}(\mathbf{x}|\hat{\mathbf{x}}_\ell, \sigma_\ell^2 \mathbf{I})$, with σ_ℓ^2 given by Eq. (85). Using this approximation in the expression of the posterior of Eq. (44), it is straightforward to see that the weighting matrix for the linear fusion in this case is $\mathbf{\Lambda}_\ell = \alpha_\ell \mathbf{I}$, with

$$\alpha_\ell = \frac{\text{Tr} \left(\mathbf{C}_x^{(\ell)} \right)^{-1}}{\sum_{k=1}^L \text{Tr} \left(\mathbf{C}_x^{(k)} \right)^{-1}} = \frac{\sigma_\ell^{-2}}{\sum_{k=1}^L \sigma_k^{-2}}. \quad (86)$$

Note that Eq. (86) corresponds to the SC-MMSE fusion rule, obtained as a solution of a constrained optimization problem in Section 3.2 and given by (30),

and also derived in Section 4.3 from a Bayesian perspective and given by (57).

B.2. Independent Approximation

Let us consider again $\Sigma_p = \mathbf{C}_x^{(\ell)}$, but now let us use $\Sigma_q = \text{diag}(\hat{\sigma}_{\ell,1}^2, \dots, \hat{\sigma}_{\ell,D}^2)$. Then, $|\Sigma_q| = \prod_{d=1}^D \hat{\sigma}_{\ell,d}^2$, $\ln |\Sigma_q| = \sum_{d=1}^D \ln(\hat{\sigma}_{\ell,d}^2)$, and $\Sigma_q^{-1} = \text{diag}(\hat{\sigma}_{\ell,1}^{-2}, \dots, \hat{\sigma}_{\ell,D}^{-2})$. In this scenario, the KL divergence in (82) becomes

$$D_{KL}(p||q) = K + \frac{D}{2} \sum_{d=1}^D \ln(\hat{\sigma}_{\ell,d}^2) + \frac{1}{2} \sum_{d=1}^D \hat{\sigma}_{\ell,d}^{-2} \sigma_{\ell,d}^2, \quad (87)$$

where $\sigma_{\ell,d}^2 = \mathbf{C}_x^{(\ell)}[d, d]$. In order to minimize the KL divergence in (87), we take again the derivative w.r.t. $\hat{\sigma}_{\ell,d}^2$ and equate it to zero:

$$\frac{\partial D_{KL}(p||q)}{\partial \hat{\sigma}_{\ell,d}^2} = \frac{D}{2\hat{\sigma}_{\ell,d}^2} - \frac{\sigma_{\ell,d}^2}{2\hat{\sigma}_{\ell,d}^4} = 0. \quad (88)$$

The solution of this equation is

$$\hat{\sigma}_{\ell,d}^2 = \frac{1}{D} \sigma_{\ell,d}^2, \quad (89)$$

and it can be easily checked (by taking the second derivative of $D_{KL}(p||q)$) that it corresponds to a minimum. Hence, the best independent approximation of $p(\mathbf{x}|\hat{\mathbf{x}}_\ell) = \mathcal{N}(\mathbf{x}|\hat{\mathbf{x}}_\ell, \mathbf{C}_x^{(\ell)})$ is $q(\mathbf{x}|\hat{\mathbf{x}}_\ell) = \mathcal{N}(\mathbf{x}|\hat{\mathbf{x}}_\ell, \text{diag}(\hat{\sigma}_{\ell,1}^2, \dots, \hat{\sigma}_{\ell,D}^2))$, with $\hat{\sigma}_{\ell,d}^2$ given by Eq. (89). Using this approximation in the expression of the posterior of Eq. (44), it is straightforward to see that the weighting matrix for the linear fusion in this case is $\mathbf{A}_\ell = \mathbf{D}_\ell = \text{diag}(\alpha_{\ell,1}, \dots, \alpha_{\ell,D})$, with

$$\alpha_{\ell,d} = \frac{\sigma_{\ell,d}^{-2}}{\sum_{k=1}^L \sigma_{k,d}^{-2}}. \quad (90)$$

⁷⁷⁵ Note that Eq. (90) corresponds now to the IL-MMSE fusion rule, obtained as a solution of a constrained optimization problem in Section 3.3 and given by (38), and also derived in Section 4.3 from a Bayesian perspective and given by (55).

C. Bounds for the Univariate Gamma Example

Let us consider the univariate Gamma example described in Section 5.1. In this case, the EWF is given by

$$\hat{\beta}^{(\text{EWF})} = \frac{1}{L} \sum_{\ell=1}^L \hat{\beta}_{\ell}^{(\text{MMSE})} = \frac{1}{L} \sum_{\ell=1}^L \frac{\alpha_0 + N_{\ell}\alpha}{\beta_0 + \sum_{n \in \mathcal{Y}_{\ell}} y_n}.$$

Assuming $N_{\ell} = N/L$ for $\ell = 1, \dots, L$ (as done in Section 5.1), the expectation of $\hat{\beta}^{(\text{EWF})}$ can be lower bounded as

$$\mathbb{E} \left(\hat{\beta}^{(\text{EWF})} \right) = \left(\alpha_0 + \frac{N\alpha}{L} \right) \mathbb{E} \left(\frac{1}{\beta_0 + \sum_{n \in \mathcal{Y}_{\ell}} y_n} \right) > \frac{\alpha_0 + N\alpha/L}{\mathbb{E}(\beta_0 + \sum_{n \in \mathcal{Y}_{\ell}} y_n)}, \quad (91)$$

where we have used the fact that $\mathbb{E} \left(\frac{1}{X} \right) > \frac{1}{\mathbb{E}(X)}$ for any non-degenerate positive random variable [48]. After some algebra, (91) can be expressed as

$$\mathbb{E} \left(\hat{\beta}^{(\text{EWF})} \right) > \frac{\alpha_0 + N\alpha/L}{\beta_0\beta + N\alpha/L} \times \beta. \quad (92)$$

Finally, note that the right-hand side of (92) is larger or equal than β whenever $\alpha_0 \geq \beta_0\beta$ (as it happens for both of the priors considered in Section 5.1, the 780 improper prior and the delta-like prior centered on the true value of β), thus showing that the EWF tends to over-estimate β in this particular example.

Finding an exact bound for the L-MMSE fusion rule is much more complicated, but an approximate bound can be obtained by using Hölder's inequality and an approximation on the expectation of the reciprocal of a random variable. Assuming, for the sake of simplicity, that $N_{\ell} = N/L$ for $\ell = 1, \dots, L$ and $\alpha_0, \beta_0 \rightarrow 0$ (which is the configuration used to obtain Figures 1 and 2), we can write the L-MMSE fusion rule as

$$\hat{\beta}^{(\text{L-MMSE})} = \sum_{\ell=1}^L \frac{\sigma_{\ell}^{-2}}{\sum_{k=1}^L \sigma_k^{-2}} \hat{\beta}_{\ell}^{(\text{MMSE})} = \frac{N\alpha}{L} \times \frac{V}{1 + Z_N}, \quad (93)$$

where

$$V = \sum_{n=1}^N y_n \quad (94a)$$

$$Z_N = \sum_{\ell=1}^L \left(\sum_{n \in \mathcal{Y}_\ell} y_n \right)^2 - 1 \quad (94b)$$

Now, using Hölder's inequality [49], the expected value of $\hat{\beta}^{(\text{L-MMSE})}$ can be upper bounded as

$$\mathbb{E} \left(\hat{\beta}^{(\text{L-MMSE})} \right) = \frac{N\alpha}{L} \mathbb{E} \left(\frac{V}{1 + Z_N} \right) \leq \frac{N\alpha}{L} [\mathbb{E}(V^2)]^{1/2} \left[\mathbb{E} \left(\frac{1}{(1 + Z_N)^2} \right) \right]^{1/2}, \quad (95)$$

where

$$\mathbb{E}(V^2) = \frac{N\alpha(1 + N\alpha)}{\beta^2} \approx \left(\frac{N\alpha}{\beta} \right)^2, \quad (96a)$$

$$\mathbb{E} \left(\frac{1}{(1 + Z_N)^2} \right) \approx \frac{1}{(1 + \mathbb{E}(Z_N))^2} = \left(\frac{L\beta^2}{N\alpha(L + N\alpha)} \right)^2, \quad (96b)$$

and the approximation in (96b) comes from the fact that $\frac{1/(1+\mathbb{E}(Z_N))^\alpha}{\mathbb{E}[1/(1+Z_N)^\alpha]} \rightarrow 1$ for a non-negative sequence of random variables Z_N , any $\alpha > 0$ and $N \rightarrow \infty$ [50].

Inserting (96a) and (96b) in (95), we finally obtain

$$\mathbb{E} \left(\hat{\beta}^{(\text{L-MMSE})} \right) \lesssim \frac{N\alpha}{L + N\alpha} \beta, \quad (97)$$

showing that the L-MMSE fusion rule tends to under-estimate β in this particular case.