

TIGHT MMSE BOUNDS FOR THE AGN CHANNEL UNDER KL DIVERGENCE CONSTRAINTS ON THE INPUT DISTRIBUTION

Michael Fauß, Abdelhak M. Zoubir

Alex Dytso, H. Vincent Poor*

Technische Universität Darmstadt
Signal Processing Group
64283 Darmstadt, Germany

Princeton University
Dept. of Electrical Engineering
Princeton, NJ 08544, USA

ABSTRACT

Tight bounds on the minimum mean square error for the additive Gaussian noise channel are derived, when the input distribution is constrained to be ε -close to a Gaussian reference distribution in terms of the Kullback–Leibler divergence. The distributions that attain the bounds are shown to be Gaussian whose means are identical to that of the reference distribution and whose covariance matrices are defined implicitly via systems of matrix equations. The estimator that attains the upper bound is identified as a minimax optimal estimator that is robust against deviations from the assumed prior. The lower bound is shown to provide a potentially tighter alternative to the Cramér–Rao bound. Both properties are illustrated with numerical examples.

Index Terms— MMSE bounds, robust estimation, minimax optimization, Cramér–Rao bound

1. INTRODUCTION AND PROBLEM FORMULATION

The mean square error (MSE) is a natural and commonly used measure for the accuracy of an estimator. The minimum MSE (MMSE) plays a central role in statistics [1, 2], information theory [3, 4], and signal processing [5, 6, 7] and has been shown to have close connections to entropy and mutual information [8, 9].

In this paper, lower and upper bounds on the MMSE are derived when the random variable of interest is contaminated by additive Gaussian noise and its distribution is constrained to be ε -close to a Gaussian reference distribution in terms of the Kullback–Leibler (KL) divergence. This problem is of interest in both information theory as well as robust statistics. More precisely, it is shown that the estimator that attains the upper bound is minimax robust in the sense that it minimizes the maximum MMSE over the set of feasible distributions. That is, within the specified KL divergence ball, it is robust against arbitrary deviations of the prior from the nominal Gaussian case. In addition, the lower bound provides a fundamental limit on the estimation accuracy and provides an alternative to the well-known Bayesian Cramér–Rao bound. However, since it uses additional information about the KL divergence, it can be significantly tighter in some cases.

More formally, let $(\mathbb{R}^K, \mathcal{B}^K)$ denote the K -dimensional Borel space. Consider an additive-noise channel

$$Y = X + N,$$

where X and N are independent $(\mathbb{R}^K, \mathcal{B}^K)$ -valued random variables. Without loss of generality, N is assumed to be zero-mean.

For the purpose of this paper, it is useful to define the MSE as a function of an estimator f and an input distribution P_X , i.e.,

$$\text{mse}_{X|Y}(f, P_X) := E_{P_{Y|X} P_X} [\|f(Y) - X\|^2].$$

The MMSE is accordingly defined as

$$\text{mmse}_{X|Y}(P_X) := \inf_{f \in \mathcal{F}} \text{mse}_{X|Y}(f, P_X),$$

where \mathcal{F} denotes the set of all all feasible estimators, i.e.,

$$\mathcal{F} = \left\{ f: (\mathbb{R}^K, \mathcal{B}^K) \rightarrow (\mathbb{R}^K, \mathcal{B}^K) \right\}.$$

The problems investigated in this paper are

$$\sup_{P_X} \text{mmse}_{X|Y}(P_X) \quad \text{s.t.} \quad P_X \in \mathcal{P}_\varepsilon, \quad (1)$$

$$\inf_{P_X} \text{mmse}_{X|Y}(P_X) \quad \text{s.t.} \quad P_X \in \mathcal{P}_\varepsilon. \quad (2)$$

where the set of all feasible distribution is defined as

$$\mathcal{P}_\varepsilon := \{ P_X \mid D_{\text{KL}}(P_X \parallel P_0) \leq \varepsilon \}.$$

Note that \mathcal{P}_ε is a KL ball centered at P_0 of radius ε . It is further assumed that

$$P_0 = \mathcal{N}(\mu_0, \Sigma_0),$$

where $\mathcal{N}(\mu, \Sigma)$ denotes the Gaussian distribution with mean μ and covariance Σ . All covariance matrices are assumed to be positive definite.

The paper is organized as follows: the main result is stated in Section 2, followed by a brief discussion in Section 3. A proof is detailed in Section 4 and two illustrative numerical examples are presented in Section 5.

A remark on notation: In what follows, Σ_X and Σ_N denote the covariance matrices of X and N , respectively. Moreover, in order to keep the notation compact, the following matrices are introduced

$$W_X := \Sigma_X(\Sigma_X + \Sigma_N)^{-1},$$

$$D_X := \Sigma_N W_X^T W_X,$$

$$W_N := \Sigma_N(\Sigma_X + \Sigma_N)^{-1},$$

$$D_N := \Sigma_X W_N^T W_N,$$

where W^T denotes the transpose of W . Note that $W_X + W_N = I$.

*This work was supported in part by the U.S. National Science Foundation under Grants CCF-1420575 and ECCS-1549881.

2. RESULT

If (α^*, Σ_X^*) , with Σ_X^* positive definite and $\alpha^* \geq 0$, solve

$$\Sigma_X = (I + \alpha D_N) \Sigma_0, \quad (3)$$

$$\alpha \operatorname{tr}(D_N) - \log \det(I + \alpha D_N) = 2\varepsilon, \quad (4)$$

then

$$P_X^* = \mathcal{N}(\mu_0, \Sigma_X^*) \quad (5)$$

solves (1). Analogously, if $(\alpha^\dagger, \Sigma_X^\dagger)$, with Σ^\dagger positive definite and $\alpha^\dagger \leq 0$, solve (3) and (4), then

$$P_X^\dagger = \mathcal{N}(\mu_0, \Sigma_X^\dagger) \quad (6)$$

solves (2).

Since the optimal distributions are Gaussians of the form $P_X = \mathcal{N}(\mu_0, \Sigma_X)$, the MMSE estimators are in both cases given by

$$f(y) = W_X y + W_N \mu_0, \quad (7)$$

where $y \in \mathbb{R}^K$ denotes the observed realization of Y . The corresponding MMSE calculates to

$$\operatorname{mmse}_{X|Y}(P_X) = \operatorname{tr}(D_X) + \operatorname{tr}(D_N). \quad (8)$$

The lower and upper bound can be obtained by evaluating (8) at Σ_X^* and Σ_X^\dagger , respectively.

3. DISCUSSION

Before proceeding with the proof of the presented bounds, some clarifying remarks are in order.

1) The result does not make a statement about the existence or the uniqueness of Σ_X^* and Σ_X^\dagger . However, considering that the MMSE is a concave functional of the input distribution [10, Corollary 1] and that the KL divergence is strictly convex in both arguments [11, Theorem 2.7.2], we conjecture that the solution is guaranteed to exist. But, considering that the MMSE is not strictly concave [10, Corollary 2], it is likely not to be unique.

2) Solving (3) and (4) for Σ_X and α is non-trivial and, in general, requires the use of numerical solvers. For this purpose, it can be useful to rewrite (3) as

$$(I + \Sigma_N^{-1} \Sigma_X)^\top (I + \Sigma_N^{-1} \Sigma_X) (I - \Sigma_0^{-1} \Sigma_X) + \alpha \Sigma_X = 0 \quad (9)$$

in order to avoid numerical instabilities that might arise from inverting Σ_X . A detailed derivation of (9) is omitted due to space constraints. Our limited experimental results suggest that both (3) and (9) can usually be solved using off-the-shelf algorithms. However, a detailed discussion is beyond the scope of this work.

3) A special case for which an analytical solution of (3) and (4) can be given is a white signal in white noise, i.e., $\Sigma_0 = \sigma_0^2 I$ and $\Sigma_N = \sigma_N^2 I$. In this case it holds that

$$P_X^* = \mathcal{N}(\mu_0, s^* I) \quad \text{and} \quad P_X^\dagger = \mathcal{N}(\mu_0, s^\dagger I),$$

where the variances s^* and s^\dagger are unique and given by

$$s^* = -W_0 \left(-e^{-(1 + \frac{2}{K} \varepsilon)} \right) \sigma_0^2, \quad (10)$$

$$s^\dagger = -W_{-1} \left(-e^{-(1 + \frac{2}{K} \varepsilon)} \right) \sigma_0^2. \quad (11)$$

Here W_k denotes the k th branch of the Lambert W function [12]. The corresponding MMSE bounds calculate to

$$\frac{s^\dagger \sigma_N^2}{s^\dagger + \sigma_N^2} \leq \frac{\operatorname{mmse}_{X|Y}(P_X)}{K} \leq \frac{s^* \sigma_N^2}{s^* + \sigma_N^2} \quad \forall P_X \in \mathcal{P}_\varepsilon.$$

The proof of these results is straightforward and omitted for brevity.

4) The signal model can be extended to conditionally independent and identically distributed random variables $Y_1|X, \dots, Y_n|X$ by letting

$$\Sigma_N \leftarrow \frac{1}{n} \Sigma_N \quad \text{and} \quad y \leftarrow \frac{1}{n} \sum_{i=1}^n y_i,$$

where y_i denotes an observation of $Y_i|X$, $i = 1, \dots, n$.

4. PROOF

The idea underlying the proof is to reformulate (1) and (2) as nested optimizations over the estimator f and the distribution P_X . Only the solution of (1) is presented in detail since the solution of (2) can be given analogously.

Before detailing the proof of the main result, a result on the optimization of expected values under f -divergence constraints is derived. It simplifies the subsequent steps, but also constitutes a useful result in its own right.

4.1. Bounding expectations under f -divergence constraints

Consider the auxiliary problems

$$\sup_{P_X} E_{P_X} [h(X)] \quad \text{s.t.} \quad D_f(P_X \| P_0) \leq \varepsilon, \quad (12)$$

$$\inf_{P_X} E_{P_X} [h(X)] \quad \text{s.t.} \quad D_f(P_X \| P_0) \leq \varepsilon, \quad (13)$$

where $h: (\mathbb{R}^K, \mathcal{B}^K) \rightarrow (\mathbb{R}, \mathcal{B})$ is a measurable function. Assuming P_X and P_0 to be absolutely continuous with respect to a reference measure η , (12) and (13) can be reformulated in terms of the densities of P_X and P_0 w.r.t. η , namely,

$$\{\sup_{p_X}, \inf_{p_X}\} \int h(x) p_X(x) \eta(dx) \quad (14)$$

$$\text{s.t.} \quad \int_{\mathbb{R}^K} f\left(\frac{p_X(x)}{p_0(x)}\right) p_0(x) \eta(dx) \leq \varepsilon \quad (15)$$

$$\int_{\mathbb{R}^K} p_X(x) \eta(dx) = 1 \quad (16)$$

$$p_X(x) \geq 0. \quad (17)$$

The equivalent unconstrained optimization problems are given by

$$\{\sup_{p_X}, \inf_{p_X}\} L(p_X; \lambda, \nu), \quad (18)$$

with

$$L(p_X; \lambda, \nu) :=$$

$$\int_{\mathbb{R}} h(x) + \lambda f\left(\frac{p_X(x)}{p_0(x)}\right) p_0(x) + \nu p_X(x) \eta(dx), \quad (19)$$

and $\nu, \lambda \in \mathbb{R}$. Note that $\lambda \leq 0$ for (12) and $\lambda \geq 0$ for (13) and that the constraint (17) is neglected since it turns out to be redundant. The Fréchet derivative [13] of (19) w.r.t. p_X is given by

$$L'_{p_X}(p_X; \lambda, \nu) = h(x) + \lambda f'\left(\frac{p_X(x)}{q(x)}\right) + \nu,$$

where f' denotes the subderivative [14, p. 36] of f . A sufficient condition for p_X^* to solve (18) is that $L'_{p_X}(p_X; \lambda, \nu) = 0$, which yields

$$p_X^*(x) = q(x)g\left(-\frac{h(x) + \nu}{\lambda}\right) = q(x)g(\alpha h(x) + \beta), \quad (20)$$

where $\alpha = -1/\lambda$, $\beta = -\nu/\lambda$, and $g: \mathbb{R} \rightarrow \mathbb{R}$ is the generalized inverse [15] of f' , i.e.,

$$g(c) := \inf\{x \in \mathbb{R} \mid f'(x) \geq c\}.$$

Note that $\alpha \geq 0$ for (12) and $\alpha \leq 0$ for (13). Finally, in order for $p_X^*(x)$ to solve the constrained problem (12) or (13), α and β need to be chosen such that $p_X^*(x)$ is a valid density and the f -divergence constraint is fulfilled with equality. The latter follows directly from the fact that the complementary slackness constraint of the Karush–Kuhn–Tucker conditions [16] requires the f -divergence constraint to be satisfied with equality for all $\lambda \neq 0$ and hence for all $\alpha \in \mathbb{R}$.

4.2. Proof of the main result

Consider the maximization in (1) which can be written as the mini-max problem

$$\sup_{P_X \in \mathcal{P}_\varepsilon} \inf_{f \in \mathcal{F}} \text{mse}_{X|Y}(f, P_X). \quad (21)$$

A sufficient condition for P_X^* and f^* to solve (21), and hence (1), is that they satisfy the saddle point conditions [17, Exercise 3.14]

$$\text{mse}_{X|Y}(f^*, P_X^*) \leq \text{mse}_{X|Y}(f, P_X^*) \quad \forall f \in \mathcal{F}, \quad (22)$$

$$\text{mse}_{X|Y}(f^*, P_X^*) \geq \text{mse}_{X|Y}(f^*, P_X) \quad \forall P_X \in \mathcal{P}_\varepsilon. \quad (23)$$

The fact that f in (7) minimizes the right hand side of (22) follows directly from the definition of the MMSE [18, Chapter 10.4]. In the remainder of the proof, it is shown that P_X^* in (5) satisfies (23).

First, the right hand side of (23) is written as

$$\text{mse}_{X|Y}(f^*, P_X) = E_{P_X}[h(X)],$$

where $h: \mathbb{R}^K \rightarrow \mathbb{R}$ is independent of P_X and given by

$$\begin{aligned} h(x) &= E_{P_{Y|X=x}}[\|W_X Y + W_N \mu_0 - x\|_2^2] \\ &= E_{\mathcal{N}(x, \Sigma_N)}[\|W_X Y + W_N \mu_0 - x\|_2^2] \\ &= \text{tr}(\Sigma_N W_X^T W_X) + \|W_X x + W_N \mu_0 - x\|_2^2 \\ &= \text{tr}(D_X) + (x - \mu_0)^T W_N^T W_N (x - \mu_0). \end{aligned}$$

In order for $D_{\text{KL}}(P_X \| P_0)$ to be finite, P_X needs to be absolutely continuous w.r.t. P_0 . Therefore, the problem

$$\sup_{P_X \in \mathcal{P}_\varepsilon} \text{mse}_{X|Y}(f^*, P_X) = \sup_{P_X \in \mathcal{P}_\varepsilon} E_{P_X}[h(X)] \quad (24)$$

is of the form (12), with $f(x) = x \log x$. Since the latter is strictly convex, the inverse function of its derivative is unique and is given by $g(c) = e^{c-1}$. Inserting $g(c)$ into (20) yields

$$\begin{aligned} p_X^*(x) &= q(x)e^{\alpha h(x) + \beta - 1} \\ &\propto q(x)e^{\alpha(x - \mu_0)^T W_N^T W_N (x - \mu_0)} \\ &\stackrel{(1)}{\propto} e^{-\frac{1}{2}(x - \mu_0)^T (\Sigma_0^{-1} - \alpha W_N^T W_N)(x - \mu_0)} \\ &\propto e^{-\frac{1}{2}(x - \mu_0)^T \Sigma_X^{-1}(x - \mu_0)}, \end{aligned}$$

where $\alpha \geq 0$ and

$$\Sigma_X^{-1} = \Sigma_0^{-1} - \alpha W_N^T W_N. \quad (25)$$

Note that, without loss of generality, α has been scaled by 1/2 in (I). Multiplying (25) by Σ_X from the left, by Σ_0 from the right, and rearranging the terms yields the optimality condition in (3).

Knowing that P_X and P_0 are Gaussians with identical means, the KL divergence $D_{\text{KL}}(P_X \| P_0)$ is given by

$$\begin{aligned} D_{\text{KL}}(P_X \| P_0) &= \frac{1}{2} (\text{tr}(\Sigma_X \Sigma_0^{-1}) - K - \log \det \Sigma_X \Sigma_0^{-1}) \\ &\stackrel{(3)}{=} \frac{1}{2} (\alpha \text{tr}(D_N) + \log \det(I + \alpha D_N)). \quad (26) \end{aligned}$$

Equating (26) with ε yields the optimality condition (4). This concludes the proof.

The proof for the optimality of P_X^\dagger follows analogously, the only difference being that in (24) the supremum is replaced by the infimum and, consequently, the sign of α is reversed.

5. NUMERICAL EXAMPLES

In order to illustrate the usefulness of the bounds in Section 2, two examples are presented; the first in the context of robust MMSE estimation, the second in the context of estimation accuracy bounds.

5.1. Robust MMSE estimation

In this example, the nominal MMSE estimator f_0 with $\Sigma_X = \Sigma_0$ and the robust MMSE estimator f^* with $\Sigma_X = \Sigma_X^*$ are compared in terms of their performance under P_0 and under their respective least favorable distribution in the feasible set \mathcal{P}_ε . By definition, the least favorable distribution for the robust MMSE estimator f^* is P_X^* . From (25) it follows that the least favorable distribution for f_0 is given by a Gaussian with mean μ_0 and covariance

$$\Sigma_X^{-1} = \Sigma_0^{-1} - \alpha(\Sigma_0 + \Sigma_N)^{-1} \Sigma_N^2 (\Sigma_0 + \Sigma_N)^{-1},$$

where $\alpha \geq 0$ needs to be chosen such that the KL divergence constraint is fulfilled with equality. In a slight abuse of notation, the MMSE of f_0 under the corresponding least favorable distribution is denoted by $\text{mmse}_{X|Y}(f_0, P_X^*)$.

For the example, the matrix Σ_0 is assumed to be of size $K = 10$ and to admit a Toeplitz structure with entries

$$[\Sigma_0]_{ij} = e^{-0.9|i-j|}, \quad i, j = 1, \dots, 10.$$

This implies $\text{tr}(\Sigma_0) = K$. The noise is assumed to be white with variance $\sigma_N^2 = 1/\gamma$ so that $\gamma = \text{tr}(\Sigma_0)/\text{tr}(\Sigma_N)$ denotes the signal-to-noise ratio (SNR) under P_0 .

In Fig. 1, the MSE of the estimators f_0 and f^* under the nominal distribution P_0 and their respective least favorable distributions is plotted versus the SNR. The KL divergence tolerance is set to $\varepsilon = 2$. Especially in the low SNR regime, the robust estimator f^* yields a gain of 1-3 dB in terms of the worst case MSE. The price for this improvement is a loss of comparable magnitude in terms of nominal performance.

The increased robustness of f^* in comparison to f_0 becomes more apparent in Fig. 2, where the MSE is plotted versus the KL divergence tolerance ε at an SNR of 0 dB. It can clearly be seen how the worst case MSE scales differently for f^* and f_0 , which highlights the advantage of the robust estimator under large prior uncertainties. At the same time, the MSE of f^* under the nominal distribution P_0 deteriorates with increasing ε . However, in this example, the gain in worst case MSE is significantly larger than the loss in nominal performance.

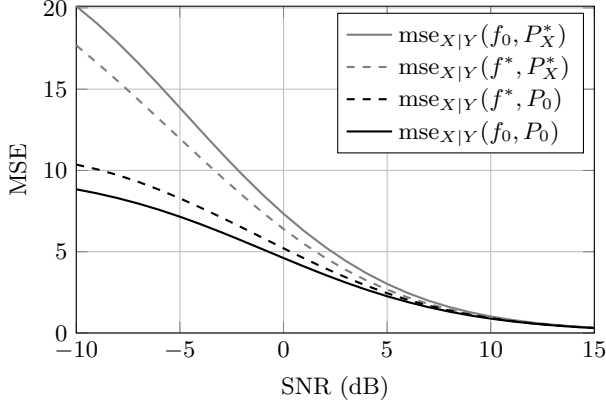


Fig. 1. MSE vs. SNR for the nominal MMSE estimator f_0 and the robust MMSE estimator f^* under P_0 and the respective least favorable distribution P_X^* for $\varepsilon = 2$.

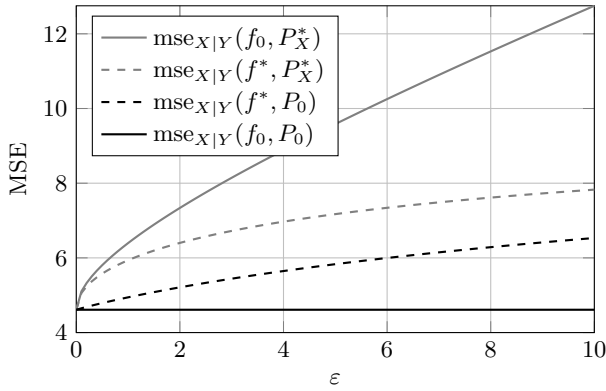


Fig. 2. MSE vs. ε for the nominal MMSE estimator f_0 and the robust MMSE estimator f^* under P_0 and the respective least favorable distribution P_X^* at SNR = 0 dB.

5.2. Comparison to the Cramér–Rao bound

In this example, the usefulness of the lower bound in Section 2 is illustrated by showing that it can provide a tighter alternative to the Cramér–Rao lower bound. Let $K = 1$ and consider the zero-mean generalized Gaussian distribution $\mathcal{G}(a, p)$ with density function

$$\mathcal{G}(x|a, p) = \frac{p}{2a\Gamma(1/p)} e^{-\left(\frac{|x|}{a}\right)^p},$$

where Γ denotes the gamma function [19], $a > 0$ is a scale parameter, and $p > 0$ determines the type of decay of the tails [20]. Note that the generalized Gaussian reduces to a regular Gaussian for $p = 2$ and to a Laplace distribution for $p = 1$. It is not difficult to verify that choosing $a = \sqrt{\frac{\Gamma(1/p)}{\Gamma(3/p)}} b$ implies that $E[X^2] = b^2$.

Calculating the exact MMSE for $X \sim \mathcal{G}(a, p)$ is non-trivial in general so that lower bounds are of interest that are either analytical or can be calculated with a low computational effort. A commonly used bound on the MMSE is the Bayesian Cramér–Rao bound [21], which for the transformation $Y = \sqrt{\gamma}X + N$ with $N \sim \mathcal{N}(0, 1)$ is

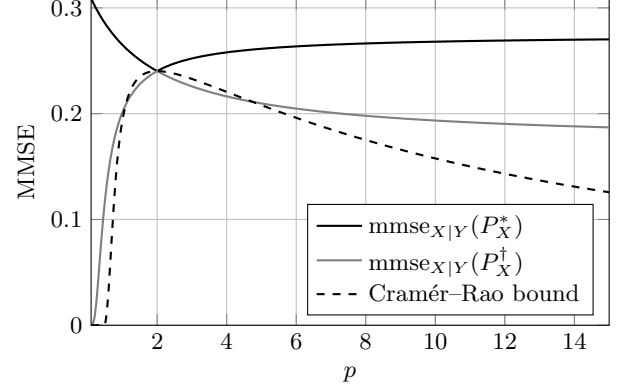


Fig. 3. MMSE bounds vs. the shape parameter p of the generalized Gaussian distribution at SNR = 5 dB.

given by

$$\text{mmse}_{X|Y}(\mathcal{G}(a, p)) \geq \frac{1}{\gamma + I(\mathcal{G}(a, p))},$$

where

$$I(\mathcal{G}(a, p)) = \begin{cases} \frac{p^2}{a^2} \frac{\Gamma(2-1/p)}{\Gamma(1/p)}, & p > \frac{1}{2} \\ \infty, & p \leq \frac{1}{2} \end{cases} \quad (27)$$

denotes the Fisher Information of the zero-mean generalized normal distribution [22, Chapter 3.2.1].

Using the result in Section 2, an alternative lower bound can be obtained via a detour over the KL divergence between a Gaussian and a generalized Gaussian. The latter is given by [23, eq. (17)]

$$D_{\text{KL}}(\mathcal{G}(a, p) \parallel \mathcal{G}(a_0, 2)) = \log \frac{pa_0}{2a} \frac{\Gamma(1/2)}{\Gamma(1/p)} + \frac{a^2}{a_0^2} \frac{\Gamma(3/p)}{\Gamma(1/p)} - \frac{1}{p}.$$

In order to find the best Gaussian reference distribution, this distance needs to be minimized w.r.t. a_0 . It is not hard to show that this is accomplished by choosing $a_0 = \sqrt{\frac{2\Gamma(3/p)}{\Gamma(1/p)}} a$ (i.e., the KL divergence is minimized if the second moments agree), so that the following divergence can be defined:

$$\begin{aligned} D_{\text{KL}}(\mathcal{G}(a, p) \parallel \mathcal{N}) &:= \min_{a_0 > 0} D_{\text{KL}}(\mathcal{G}(a, p) \parallel \mathcal{G}(a_0, 2)) \\ &= \log \frac{p}{\sqrt{2}} \sqrt{\frac{\Gamma(3/p)}{\Gamma(1/p)} \frac{\Gamma(1/2)}{\Gamma(1/p)}} + \frac{1}{2} - \frac{1}{p} \\ &=: d_{\text{KL}}(p), \end{aligned}$$

where $d_{\text{KL}}: \mathbb{R}_{>0} \rightarrow \mathbb{R}_{>0}$ is defined implicitly and only depends on p . Using this relation, bounds on the MMSE for any $X \sim \mathcal{G}(a, p)$ can be obtained by evaluating (10) and (11) at $\varepsilon = d_{\text{KL}}(p)$.

An example of bounds on $\text{mmse}_{X|Y}(\mathcal{G}(a, p))$ that were obtained via the KL divergence approach is shown in Fig. 3. The SNR was set to 5 dB and the Cramér–Rao bound is plotted for comparison. Interestingly, neither of the bounds is uniformly tighter than the other. While the Cramér–Rao bound is more accurate on the interval $\approx [1.1, 4.7]$, the bound in Section 2 is tighter for all p outside this interval. In particular for very large and very small values of p , the proposed bound is a significant improvement. Moreover, it can also be calculated for values $p \leq 0.5$, for which the Fisher Information in (27) is infinite so that the Cramér–Rao bound becomes trivial.

6. REFERENCES

- [1] E. L. Lehmann and G. Casella, *Theory of Point Estimation*, Springer, New York City, New York, USA, 2 edition, 1998.
- [2] Y. Dodge, *The Concise Encyclopedia of Statistics*, chapter Criterion of Total Mean Squared Error, pp. 141–144, Springer, New York City, New York, USA, 2008.
- [3] D. Guo, Y. Wu, S. Shamai (Shitz), and S. Verdú, “Estimation in Gaussian Noise: Properties of the Minimum Mean-Square Error,” *IEEE Transactions on Information Theory*, vol. 57, no. 4, pp. 2371–2385, April 2011.
- [4] A. Dytso, R. Bustin, D. Tuninetti, N. Devroye, H. V. Poor, and S. Shamai (Shitz), “On Communication Through a Gaussian Channel With an MMSE Disturbance Constraint,” *IEEE Transactions on Information Theory*, vol. 64, no. 1, pp. 513–530, 2018.
- [5] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*, Prentice-Hall, Upper Saddle River, NJ, USA, 1993.
- [6] L. A. Dalton and E. R. Dougherty, “Exact Sample Conditioned MSE Performance of the Bayesian MMSE Estimator for Classification Error—Part I: Representation,” *IEEE Transactions on Signal Processing*, vol. 60, no. 5, pp. 2575–2587, 2012.
- [7] L. A. Dalton and E. R. Dougherty, “Exact Sample Conditioned MSE Performance of the Bayesian MMSE Estimator for Classification Error—Part II: Consistency and Performance Analysis,” *IEEE Transactions on Signal Processing*, vol. 60, no. 5, pp. 2588–2603, 2012.
- [8] D. Guo, S. Shamai (Shitz), and S. Verdú, “Mutual Information and Minimum Mean-Square Error in Gaussian Channels,” *IEEE Transactions on Information Theory*, vol. 51, no. 4, pp. 1261–1282, 2005.
- [9] S. Verdú and D. Guo, “A Simple Proof of the Entropy-Power Inequality,” *IEEE Transactions on Information Theory*, vol. 52, no. 5, pp. 2165–2166, 2006.
- [10] Y. Wu and S. Verdú, “Functional Properties of Minimum Mean-Square Error and Mutual Information,” *IEEE Transactions on Information Theory*, vol. 58, no. 3, pp. 1289–1301, 2012.
- [11] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, Wiley-Interscience, Hoboken, New Jersey, USA, 2 edition, 2006.
- [12] R. M. Corless, G. H. Gonnet, D. E. G. Hare, D. J. Jeffrey, and D. E. Knuth, “On the Lambert W Function,” *Advances in Computational Mathematics*, vol. 5, no. 1, pp. 329–359, 1996.
- [13] J. Bell, “Fréchet Derivatives and Gâteaux Derivatives,” 2014, available online: <http://individual.utoronto.ca/jordanbell/notes/frechetderivatives.pdf>.
- [14] F. Clarke, *Optimization and Nonsmooth Analysis*, Society for Industrial and Applied Mathematics, 1990.
- [15] P. Embrechts and M. Hofert, “A Note on Generalized Inverses,” *Mathematical Methods of Operations Research*, vol. 77, no. 3, pp. 423–432, 2013.
- [16] O. Brezhneva and A. A. Tret’yakov, “An Elementary Proof of the Karush–Kuhn–Tucker Theorem in Normed Linear Spaces for Problems With a Finite Number of Inequality Constraints,” *Optimization*, vol. 60, no. 5, pp. 613–618, 2011.
- [17] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, Cambridge, UK, 2004.
- [18] W. D. Penny, “Signal Processing Course,” 2000, available online: <http://www.fil.ion.ucl.ac.uk/wpenny/course/course.pdf>.
- [19] P. J. Davis, “Leonhard Euler’s Integral: A Historical Profile of the Gamma Function,” *The American Mathematical Monthly*, vol. 66, no. 10, pp. 849–869, 1959.
- [20] A. Dytso, R. Bustin, H. V. Poor, and S. Shlomo (Shitz), “On Additive Channels With Generalized Gaussian Noise,” in *Proc. of the IEEE International Symposium on Information Theory*, June 2017, pp. 426–430.
- [21] R. D. Gill and B. Y. Levit, “Applications of the van Trees Inequality: A Bayesian Cramér–Rao Bound,” *Bernoulli*, vol. 1, no. 1/2, pp. 59–79, 1995.
- [22] S. A. Kassam and J. B. Thomas, *Signal Detection in Non-Gaussian Noise*, Springer Texts in Electrical Engineering, Springer, New York City, New York, USA, 2012.
- [23] M. N. Do and M. Vetterli, “Wavelet-Based Texture Retrieval Using Generalized Gaussian Density and Kullback-Leibler Distance,” *IEEE Transactions on Image Processing*, vol. 11, no. 2, pp. 146–158, 2002.