

Gradient of Mutual Information in Linear Vector Gaussian Channels

Daniel P. Palomar, *Member, IEEE*, and Sergio Verdú, *Fellow, IEEE*

Abstract—This paper considers a general linear vector Gaussian channel with arbitrary signaling and pursues two closely related goals: i) closed-form expressions for the gradient of the mutual information with respect to arbitrary parameters of the system, and ii) fundamental connections between information theory and estimation theory. Generalizing the fundamental relationship recently unveiled by Guo, Shamai, and Verdú, we show that the gradient of the mutual information with respect to the channel matrix is equal to the product of the channel matrix and the error covariance matrix of the best estimate of the input given the output. Gradients and derivatives with respect to other parameters are then found via the differentiation chain rule.

Index Terms—De Bruijn’s identity, divergence, Gaussian noise, minimum mean-square error (MMSE), multiple-input multiple-output (MIMO) channels, mutual information, nonlinear estimation, precoder optimization.

I. INTRODUCTION AND MOTIVATION

THIS paper considers general linear vector channels with Gaussian noise and arbitrary input distributions. Our purpose is twofold: i) to find closed-form expressions for the gradient of the mutual information with respect to arbitrary parameters of the system, and ii) to explore the fundamental connections between information theory and estimation theory. In fact, both goals are achieved simultaneously since the gradient of the mutual information happens to be directly related to the performance of the conditional mean estimator.

Closed-form expressions for the gradient of the mutual information with respect to arbitrary parameters are useful in both analysis and design. The most direct application is to analyze and understand the sensitivity and robustness of a system to variations in certain parameters. Learning the weaknesses of a system may teach us how to make it more robust. Indeed, an engineer may have the freedom to modify or even to design a specific part of the system. In such a case, the availability of expressions for the gradient of an objective function with respect to

the design parameters is of paramount importance to optimize the system. One common example is the design of a transmit precoder for a given communication system; in particular, the precoder can be flexible and adapt to the channel realization to increase the system performance [1].

Early connections between fundamental quantities in information theory and estimation theory were De Bruijn’s identity, used by Stam [2] to prove the entropy-power inequality, the representation of mutual information as a function of causal filtering error by Duncan [3] and by Kadota, Zakai, and Ziv [4], and the rate-distortion bounds of Ziv and Zakai [5]. Recently, a fundamental relation between the mutual information and the minimum mean-square error (MMSE) was unveiled in [6] for discrete-time and continuous-time, scalar and vector channels with Gaussian noise. In the special case of the scalar Gaussian channel $y = \sqrt{\text{snr}}x + n$ (complex-valued inputs/outputs are considered throughout this paper) and regardless of the input distribution:

$$\frac{d}{d\text{snr}}I(x; \sqrt{\text{snr}}x + n) = \text{mmse}(\text{snr}) \quad (1)$$

where mutual information is in nats and $\text{mmse}(\text{snr})$ is the MMSE corresponding to the best estimation of x upon the observation y for a given signal-to-noise ratio (SNR) snr , i.e.,

$$\text{mmse}(\text{snr}) = \mathbb{E}[|x - \mathbb{E}[x | \sqrt{\text{snr}}x + n]|^2]. \quad (2)$$

An extension of (1) to the vector case was also given in [6] for the linear vector Gaussian channel $\mathbf{y} = \sqrt{\text{snr}}\mathbf{H}\mathbf{x} + \mathbf{n}$ as

$$\begin{aligned} \frac{d}{d\text{snr}}I(\mathbf{x}; \sqrt{\text{snr}}\mathbf{H}\mathbf{x} + \mathbf{n}) \\ = \mathbb{E}[||\mathbf{H}\mathbf{x} - \mathbb{E}[\mathbf{H}\mathbf{x} | \sqrt{\text{snr}}\mathbf{H}\mathbf{x} + \mathbf{n}]||^2] \end{aligned} \quad (3)$$

where the right-hand side is the expected squared Euclidean norm of the error in the estimation of $\mathbf{H}\mathbf{x}$ (rather than \mathbf{x}).

The representation of mutual information in (1) as an integral of $\text{mmse}(\text{snr})$ has found various interesting applications (cf. [6]). The derivative of the mutual information with respect to the signal-to-noise ratio (SNR) in a communication system is clearly a useful quantity for an engineer. In addition, it may also be of interest to generalize such sensitivity analysis to arbitrary parameters of the system rather than just the SNR. In particular, to obtain a full description of the sensitivity of the mutual information, one should obtain the partial derivatives with respect to arbitrary parameters affecting each combination of transmit–receive dimension. A compact way to describe the sensitivity of mutual information for the linear vector channel with independent Gaussian noise $\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{n}$ is via the gradient with respect

Manuscript received April 29, 2005; revised August 4, 2005. This work was supported in part by the Fulbright Program and the Ministry of Education and Science of Spain; the U.S. National Science Foundation under Grant NCR-0074277; and through collaborative participation in the Communications and Networks Consortium sponsored by the U.S. Army Research Laboratory under the Collaborative Technology Alliance Program, Cooperative Agreement DAAD19-01-2-0011. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation thereon.

The authors are with the Department of Electrical Engineering, Princeton University, Engineering Quadrangle, Princeton, NJ 08544 USA (e-mail: danielp@princeton.edu; verdu@princeton.edu).

Communicated by Y. Steinberg, Associate Editor for Shannon Theory.
Digital Object Identifier 10.1109/TIT.2005.860424

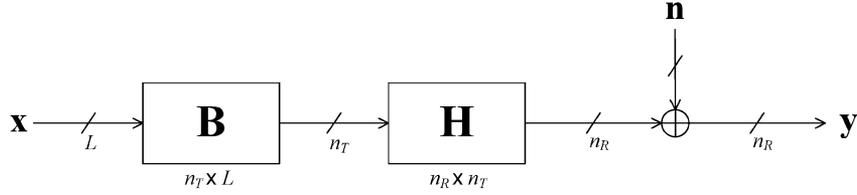


Fig. 1. Linear vector channel including linear precoding.

to the deterministic matrix \mathbf{H} . The main result of this paper is the following formula for the gradient:

$$\nabla_{\mathbf{H}} I(\mathbf{x}; \mathbf{H}\mathbf{x} + \mathbf{n}) = \mathbf{H}\mathbf{E} \quad (4)$$

where \mathbf{E} is the covariance matrix of the estimation error vector, sometimes referred to as the MMSE matrix [7]. The MMSE matrix is the full generalization of the scalar MMSE in (1) to the vector case. As we will see, applying the chain rule for differentiation to the basic gradient in (4), one can obtain the sensitivity of the mutual information with respect to any arbitrary parameter such as the SNR, the transmit covariance matrix, or a precoding matrix.

The paper is organized as follows. Section II describes the signal model for the general linear vector Gaussian channel, along with several particularizations. The main results of the paper are given in Section III: the relation between various gradients of the mutual information and the MMSE matrix. Section IV contains a number of gradients of information measures (namely, particularization of the general result to SNR gradients, generalization of De Bruijn's identity, gradient of the divergence between the conditional and unconditional outputs, gradient of the divergence between the signal-plus-noise and noise distributions, and gradient of the non-Gaussianness). Section V illustrates the potential of the theoretical results with a practical application consisting of the numerical optimization of a linear precoder. Section VI draws the final conclusions and summarizes the paper.

Notation: Boldface upper case letters denote matrices, boldface lower case letters denote column vectors, and italics denote scalars. The superscripts $(\cdot)^T$, $(\cdot)^*$, and $(\cdot)^\dagger$ denote transpose, complex conjugate, and Hermitian operations, respectively. $[\mathbf{X}]_{ij}$ denotes the (i th, j th) element of matrix \mathbf{X} and x_i denotes the i th element of vector \mathbf{x} . $\text{Tr}(\cdot)$ and $\det(\cdot)$ denote the trace and determinant of a matrix, respectively. $\|\mathbf{X}\|$ denotes the Frobenius norm of matrix \mathbf{X} and is given by $\|\mathbf{X}\| = \sqrt{\text{Tr}(\mathbf{X}\mathbf{X}^\dagger)}$.

II. SIGNAL MODEL

Consider a general discrete-time linear vector Gaussian channel represented by the following vector signal model with n_T transmit dimensions and n_R receive dimensions

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{n} \quad (5)$$

where all quantities are complex-valued, \mathbf{x} is the n_T -dimensional transmitted vector, \mathbf{H} is the $n_R \times n_T$ matrix that denotes the linear transformation undergone by the signal, \mathbf{y} is the

n_R -dimensional received vector, and \mathbf{n} is an n_R -dimensional proper complex Gaussian noise vector independent of \mathbf{x} . The input and the noise covariance matrices are

$$\Sigma_{\mathbf{x}} = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^\dagger] \quad (6)$$

and

$$\Sigma_{\mathbf{n}} = \mathbb{E}[(\mathbf{n} - \mathbb{E}[\mathbf{n}])(\mathbf{n} - \mathbb{E}[\mathbf{n}])^\dagger]. \quad (7)$$

The general channel model in (5) describes many different communication systems, e.g., wireless multiple-antenna systems, code-division multiple-access (CDMA) systems, wireline digital subscriber line systems, or even single-antenna frequency-selective wideband channels.

It is interesting to further generalize the model in (5) to include an additional linear transformation represented by the $n_T \times L$ matrix \mathbf{B} , where L is now the dimension of the data vector \mathbf{x} (Fig. 1)

$$\mathbf{y} = \mathbf{H}\mathbf{B}\mathbf{x} + \mathbf{n}. \quad (8)$$

The additional matrix \mathbf{B} can play different roles: i) in a wireless multiple-antenna system it may represent a *beamforming matrix* that uses some knowledge about the physical channel \mathbf{H} to properly steer the transmitted signal through the best channel eigenmodes [8], [9]; ii) \mathbf{B} can denote a *linear precoding matrix* or *shaping matrix* that adapts or shapes the transmitted signal to the channel realization [1], [10]–[12], [9]; iii) the overall input-output linear transformation may factor into two matrices \mathbf{H} and \mathbf{B} (one of which may be controllable by the designer of the system).

A special case of (8) by setting $\mathbf{B} = \sqrt{\text{snr}}\mathbf{I}$ yields the model used in (3)

$$\mathbf{y} = \sqrt{\text{snr}}\mathbf{H}\mathbf{x} + \mathbf{n}. \quad (9)$$

The output conditional probability density function (pdf) corresponding to the linear vector Gaussian model in (8) (assuming a zero-mean noise without loss of generality) is

$$p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) = \frac{1}{\det(\pi\Sigma_{\mathbf{n}})} \exp\left(-(\mathbf{y} - \mathbf{H}\mathbf{B}\mathbf{x})^\dagger \Sigma_{\mathbf{n}}^{-1} (\mathbf{y} - \mathbf{H}\mathbf{B}\mathbf{x})\right). \quad (10)$$

Consider the estimation of the input signal \mathbf{x} based on the observation of the output \mathbf{y} . In the scalar case, the mean-square error (MSE) of an estimate $\hat{x}(y)$ of the input x based on the observation y is defined as $\mathbb{E}[|x - \hat{x}(y)|^2]$, which is the variance of the estimator error $x - \hat{x}(y)$ provided that the estimator is unbiased. The conditional mean $\hat{x}(y) = \mathbb{E}[x|y]$ achieves the

MMSE and is referred to as *MMSE estimator*. In the more general vector setup, the MMSE estimator $\hat{\mathbf{x}}(\mathbf{y})$ achieves simultaneously the MMSE for all components of the estimation error vector and is again given by the conditional mean estimator

$$\hat{\mathbf{x}}(\mathbf{y}) = \mathbb{E}[\mathbf{x} | \mathbf{y}]. \quad (11)$$

The full description of the performance of the vector MMSE estimator is given by the *MMSE matrix*, i.e., the covariance of the estimation error vector

$$\mathbf{E} \triangleq \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x} | \mathbf{y}])(\mathbf{x} - \mathbb{E}[\mathbf{x} | \mathbf{y}])^\dagger]. \quad (12)$$

An alternative limited description of the performance, which is sometimes useful, is given by the mean-squared norm of the estimation error $\mathbb{E}[\|\mathbf{x} - \mathbb{E}[\mathbf{x} | \mathbf{y}]\|^2]$, which is equal to the trace of the MMSE matrix.

III. MAIN RESULTS

A. Basic Result: Gradient With Respect to \mathbf{H}

To start with, consider the signal model in (5) with *Gaussian signaling*, where the covariance matrix of the transmitted signal \mathbf{x} is denoted by Σ_x and the covariance matrix of the noise is the identity matrix. The mutual information, $I(\mathbf{x}; \mathbf{y})$, is¹ (e.g., [13])

$$I = \log \det(\mathbf{I} + \mathbf{H}\Sigma_x\mathbf{H}^\dagger). \quad (13)$$

The MMSE estimator is

$$\hat{\mathbf{x}} = \Sigma_x\mathbf{H}^\dagger(\mathbf{I} + \mathbf{H}\Sigma_x\mathbf{H}^\dagger)^{-1}\mathbf{y} \quad (14)$$

(assuming \mathbf{x} and \mathbf{n} with zero mean), which is linear since the input is Gaussian, and the MMSE matrix is (e.g., [7, Theorem 11.1])

$$\mathbf{E} = (\Sigma_x^{-1} + \mathbf{H}^\dagger\mathbf{H})^{-1}. \quad (15)$$

Taking the gradient² of (13) we obtain

$$\begin{aligned} \nabla_{\mathbf{H}}I &= \nabla_{\mathbf{H}} \log \det(\mathbf{I} + \mathbf{H}^\dagger\mathbf{H}\Sigma_x) \\ &= \mathbf{H}\Sigma_x(\mathbf{I} + \mathbf{H}^\dagger\mathbf{H}\Sigma_x)^{-1} = \mathbf{H}\mathbf{E}. \end{aligned} \quad (16)$$

The main result of the paper shows that the Gaussian assumption on the input is unnecessary for (4) to hold.

Theorem 1: Consider the signal model in (5), where \mathbf{H} is an arbitrary deterministic matrix, the signaling \mathbf{x} is arbitrarily distributed (with finite second-order moments), and the noise \mathbf{n} is Gaussian, independent of the input \mathbf{x} , with normalized covariance matrix $\Sigma_n = \mathbf{I}$. Then, the mutual information I and the MMSE matrix \mathbf{E} satisfy

$$\nabla_{\mathbf{H}}I(\mathbf{x}; \mathbf{H}\mathbf{x} + \mathbf{n}) = \mathbf{H}\mathbf{E}. \quad (17)$$

Proof: See Appendix B. \square

¹Throughout this paper nats are used as the information units and log denotes natural logarithm.

²See Appendix A for some observations on the definition of complex gradient.

B. Gradient With Respect to Arbitrary Parameters

With the basic gradient result of Theorem 1, one can easily find the gradient with respect to arbitrary parameters through a gradient chain rule as stated next.

Lemma 1: Let f be a real-valued function which depends on θ through \mathbf{H} . The following chain rules hold.³

- If θ is complex, then

$$\nabla_{\theta}f = \text{Tr}(\nabla_{\mathbf{H}}f \times \nabla_{\theta}\mathbf{H}^\dagger) + \text{Tr}(\nabla_{\mathbf{H}}^\dagger f \times \nabla_{\theta}\mathbf{H}). \quad (18)$$

- If θ is real, then

$$\nabla_{\theta}f = 2\text{ReTr}(\nabla_{\mathbf{H}}f \times \nabla_{\theta}\mathbf{H}^\dagger). \quad (19)$$

For example, the gradient of the mutual information along a specific direction given by \mathbf{D} is

$$\left. \frac{d}{dt}I(\mathbf{x}; (\mathbf{H} + t\mathbf{D})\mathbf{x} + \mathbf{n}) \right|_{t=0} = 2\text{ReTr}(\mathbf{D}^\dagger\mathbf{H}\mathbf{E}). \quad (20)$$

We now give the gradient of the mutual information with respect to several quantities of interest.

Theorem 2: Consider the general signal model in (8) including a linear precoder \mathbf{B} at the transmitter, where \mathbf{H} is an arbitrary deterministic matrix, the signaling \mathbf{x} is arbitrarily distributed with covariance matrix Σ_x , the noise \mathbf{n} is Gaussian, independent of the input \mathbf{x} , and has positive definite covariance matrix Σ_n , the transmit covariance matrix is $\mathbf{Q} = \mathbf{B}\Sigma_x\mathbf{B}^\dagger$ (which includes the precoding), and the squared linear precoder is $\mathbf{Q}_B = \mathbf{B}\mathbf{B}^\dagger$. Then, the mutual information I and the MMSE matrix \mathbf{E} satisfy

$$\nabla_{\mathbf{H}}I(\mathbf{x}; \mathbf{H}\mathbf{B}\mathbf{x} + \mathbf{n}) = \Sigma_n^{-1}\mathbf{H}\mathbf{B}\mathbf{E}\mathbf{B}^\dagger \quad (21)$$

$$\nabla_{\mathbf{B}}I(\mathbf{x}; \mathbf{H}\mathbf{B}\mathbf{x} + \mathbf{n}) = \mathbf{H}^\dagger\Sigma_n^{-1}\mathbf{H}\mathbf{B}\mathbf{E} \quad (22)$$

$$\nabla_{\mathbf{Q}}I(\mathbf{x}; \mathbf{H}\mathbf{B}\mathbf{x} + \mathbf{n})\mathbf{B}\Sigma_x = \mathbf{H}^\dagger\Sigma_n^{-1}\mathbf{H}\mathbf{B}\mathbf{E} \quad (23)$$

$$\nabla_{\mathbf{Q}_B}I(\mathbf{x}; \mathbf{H}\mathbf{B}\mathbf{x} + \mathbf{n})\mathbf{B} = \mathbf{H}^\dagger\Sigma_n^{-1}\mathbf{H}\mathbf{B}\mathbf{E} \quad (24)$$

$$\nabla_{\Sigma_x}I(\mathbf{x}; \mathbf{H}\mathbf{B}\mathbf{x} + \mathbf{n})\Sigma_x = \mathbf{B}^\dagger\mathbf{H}^\dagger\Sigma_n^{-1}\mathbf{H}\mathbf{B}\mathbf{E} \quad (25)$$

$$\nabla_{\Sigma_n^{-1}}I(\mathbf{x}; \mathbf{H}\mathbf{B}\mathbf{x} + \mathbf{n}) = \mathbf{H}\mathbf{B}\mathbf{E}\mathbf{B}^\dagger\mathbf{H}^\dagger \quad (26)$$

$$\nabla_{\Sigma_n}I(\mathbf{x}; \mathbf{H}\mathbf{B}\mathbf{x} + \mathbf{n}) = -\Sigma_n^{-1}\mathbf{H}\mathbf{B}\mathbf{E}\mathbf{B}^\dagger\mathbf{H}^\dagger\Sigma_n^{-1}. \quad (27)$$

Proof: See Appendix C. \square

C. Random Channels

We now consider the interesting scenario in which the channel matrix \mathbf{H} is random. In such a case, we can still obtain expressions for gradients with respect to different parameters of the system (with the exception of the channel matrix \mathbf{H} , of course). It will be instrumental to make the distinction according to whether the channel is known or unknown at the receiver. For concreteness, we consider the signal model

$$\mathbf{y} = \mathbf{C}\mathbf{H}\mathbf{x} + \mathbf{n} \quad (28)$$

³For the sake of clarity, in (18)–(19) we indicate the product between matrices by \times .

where \mathbf{H} is a random matrix with finite second-order moments, and compute the gradient with respect to \mathbf{C} . Note that for a deterministic channel matrix the gradient is

$$\nabla_{\mathbf{C}} I(\mathbf{x}; \mathbf{C}\mathbf{H}\mathbf{x} + \mathbf{n}) = \Sigma_n^{-1} \mathbf{C} \mathbf{H} \mathbf{E} \mathbf{H} \mathbf{E} \mathbf{H}^\dagger. \quad (29)$$

The case of channel known at the receiver reduces to taking expectations with respect to the channel matrix in (29) with the following result (cf. Appendix H)

$$\begin{aligned} \nabla_{\mathbf{C}} I(\mathbf{x}; \mathbf{C}\mathbf{H}\mathbf{x} + \mathbf{n} | \mathbf{H}) \\ = \Sigma_n^{-1} \mathbf{C} \mathbb{E}[(\mathbf{H}\mathbf{x} - \mathbb{E}[\mathbf{H}\mathbf{x} | \mathbf{y}, \mathbf{H}])(\mathbf{H}\mathbf{x} - \mathbb{E}[\mathbf{H}\mathbf{x} | \mathbf{y}, \mathbf{H}])^\dagger] \end{aligned} \quad (30)$$

where the outer expectation is with respect to \mathbf{x}, \mathbf{y} , and \mathbf{H} . Gradients with respect to other parameters of the system can be similarly derived (except with respect to the random channel matrix) similarly to Theorem 2. As an illustrative example, we can consider the signal model $\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{n}$ with $\Sigma_n = N_0 \mathbf{I}$ and obtain the first-order expansion of the mutual information with respect to N_0 as

$$I(\mathbf{x}; \mathbf{H}\mathbf{x} + \mathbf{n} | \mathbf{H}) = N_0^{-1} \text{Tr}(\mathbb{E}[\mathbf{H}\Sigma_x \mathbf{H}^\dagger]) + o(N_0^{-1}) \quad (31)$$

which agrees with [14, Theorem 3].

The case of channel unknown at the receiver is significantly more complicated and can be successfully tackled using the following result for an arbitrary random transformation of the input, which generalizes Theorem 1 and [6, Theorem 10].

Theorem 3: Consider the following signal model including an arbitrary random transformation $\mathbf{z} = f(\mathbf{x})$ on the input (independent of the noise \mathbf{n}):

$$\mathbf{y} = \mathbf{H}f(\mathbf{x}) + \mathbf{n} \quad (32)$$

where all the terms are defined as in Theorem 1. Then, the gradient of the mutual information is

$$\begin{aligned} \nabla_{\mathbf{H}} I(\mathbf{x}; \mathbf{y}) = \mathbf{H} \mathbb{E}[(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}])(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}])^\dagger] \\ - \mathbf{H} \mathbb{E}[(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}, \mathbf{x}])(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}, \mathbf{x}])^\dagger]. \end{aligned} \quad (33)$$

In addition, if the transformation f is deterministic, then the second term in the right-hand side of (33) vanishes, obtaining

$$\begin{aligned} \nabla_{\mathbf{H}} I(\mathbf{x}; \mathbf{H}f(\mathbf{x}) + \mathbf{n}) \\ = \mathbf{H} \mathbb{E}[(f(\mathbf{x}) - \mathbb{E}[f(\mathbf{x}) | \mathbf{y}])(f(\mathbf{x}) - \mathbb{E}[f(\mathbf{x}) | \mathbf{y}])^\dagger] \end{aligned} \quad (34)$$

which happens to be equal to $\nabla_{\mathbf{H}} I(f(\mathbf{x}); \mathbf{H}f(\mathbf{x}) + \mathbf{n})$.

Proof: The proof hinges on the fact that $\mathbf{x} \rightarrow \mathbf{z} \rightarrow \mathbf{y}$ is a Markov chain and then $I(\mathbf{x}; \mathbf{y}) = I(\mathbf{z}; \mathbf{y}) - I(\mathbf{z}; \mathbf{y} | \mathbf{x})$ to which Theorem 1 can be applied. If f is deterministic, then $I(\mathbf{z}; \mathbf{y} | \mathbf{x}) = 0$. \square

Now, using Theorem 3, we can easily obtain the gradient for the case of channel unknown at the receiver (similarly to [6, Theorem 10] for the scalar case) by modeling the system as the Markov chain $\mathbf{x} \rightarrow \mathbf{z} \rightarrow \mathbf{y}$, where $\mathbf{z} = \mathbf{H}\mathbf{x}$ is the output of a random transformation applied to \mathbf{x} that describes the effect to the random channel. The gradient is found as (cf. Appendix H)

$$\begin{aligned} \nabla_{\mathbf{C}} I(\mathbf{x}; \mathbf{C}\mathbf{H}\mathbf{x} + \mathbf{n}) \\ = \Sigma_n^{-1} \mathbf{C} \mathbb{E}[(\mathbf{H}\mathbf{x} - \mathbb{E}[\mathbf{H}\mathbf{x} | \mathbf{y}])(\mathbf{H}\mathbf{x} - \mathbb{E}[\mathbf{H}\mathbf{x} | \mathbf{y}])^\dagger] \\ - \Sigma_n^{-1} \mathbf{C} \mathbb{E}[(\mathbf{H}\mathbf{x} - \mathbb{E}[\mathbf{H}\mathbf{x} | \mathbf{y}, \mathbf{x}])(\mathbf{H}\mathbf{x} - \mathbb{E}[\mathbf{H}\mathbf{x} | \mathbf{y}, \mathbf{x}])^\dagger] \end{aligned} \quad (35)$$

where the outer expectations are with respect to \mathbf{x}, \mathbf{y} , and \mathbf{H} . Gradients with respect to other parameters of the system can be similarly derived as long as they appear after the random transformation $\mathbf{H}\mathbf{x}$, i.e., as long as they occur in the last step $\mathbf{z} \rightarrow \mathbf{y}$ of the Markov chain $\mathbf{x} \rightarrow \mathbf{z} \rightarrow \mathbf{y}$. As an illustrative example, we can again consider the signal model $\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{n}$ with $\Sigma_n = N_0 \mathbf{I}$ and obtain the first-order expansion of the mutual information with respect to N_0 as

$$I(\mathbf{x}; \mathbf{H}\mathbf{x} + \mathbf{n}) = N_0^{-1} \mathbb{E}[\|\mathbb{E}[\mathbf{H}](\mathbf{x} - \mathbb{E}[\mathbf{x}])\|^2] + o(N_0^{-1}) \quad (36)$$

which agrees with [14, Theorem 2].

D. First-Order Approximation of Mutual Information

In this subsection, we obtain first-order approximations of the mutual information as a function of different parameters using the gradients in Theorems 1 and 2.

Consider the signal model $\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{n}$ (with $\Sigma_n = \mathbf{I}$). The first-order expansion of the mutual information as a function of the transmit (positive definite) covariance matrix Σ_x , denoted by $I(\Sigma_x)$, is obtained from the gradient $\nabla_{\Sigma_x} I = \mathbf{H}^\dagger \mathbf{H} \mathbf{E} \Sigma_x^{-1}$ (Theorem 2) as [15]

$$\begin{aligned} I(\Sigma_{x,0} + \Delta) &= I(\Sigma_{x,0}) + \text{Tr}(\nabla_{\Sigma_x}^\dagger I(\Sigma_{x,0}) \Delta) + o(\|\Delta\|) \\ &= I(\Sigma_{x,0}) + \text{Tr}(\Sigma_{x,0}^{-1} \mathbf{E} \mathbf{H}^\dagger \mathbf{H} \Delta) + o(\|\Delta\|). \end{aligned} \quad (37)$$

This expansion can be easily particularized around $\Sigma_{x,0} = \mathbf{0}$ (low-power regime) as follows. The gradient $\nabla_{\Sigma_x} I$ in Theorem 2 is not fully characterized when Σ_x is singular, so we need the following more general result (see Appendix D)

$$\nabla_{\Sigma_x} I = \mathbf{H}^\dagger \mathbf{H} \lim_{\epsilon \rightarrow 0} \mathbf{E}(\Sigma_x + \epsilon \mathbf{I})^{-1} \quad (38)$$

where the MMSE matrix \mathbf{E} implicitly depends on ϵ , and the limit is equal to $\mathbf{E}\Sigma_x^{-1}$, when Σ_x is nonsingular, and to the identity matrix \mathbf{I} , when $\Sigma_x = \mathbf{0}$ (in which case $\nabla_{\Sigma_x} I = \mathbf{H}^\dagger \mathbf{H}$). More generally, one can work with the subspaces of Σ_x associated to the null eigenvalues and the positive eigenvalues to obtain a closed-form expression for the limit. For example, if $\mathbf{H} = \text{diag}(\{\mathbf{H}_1, \mathbf{H}_2\})$ and $\Sigma_x = \text{diag}(\{\Sigma_{x,1}, \mathbf{0}\})$ with nonsingular $\Sigma_{x,1}$, then the limit still admits a simple expression

$$\lim_{\epsilon \rightarrow 0} \mathbf{E}(\Sigma_x + \epsilon \mathbf{I})^{-1} = \begin{bmatrix} \mathbf{E}_1 \Sigma_{x,1}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \quad (39)$$

where \mathbf{E}_1 is the MMSE matrix corresponding to the components of \mathbf{x} associated to $\Sigma_{x,1}$, i.e., the subspace of Σ_x with positive eigenvalues.

Using the result in (38), we can immediately obtain the following first-order expansion of the mutual information around $\Sigma_{x,0} = \mathbf{0}$:

$$I(\Sigma_x) = \text{Tr}(\mathbf{H}\Sigma_x \mathbf{H}^\dagger) + o(\|\Sigma_x\|). \quad (40)$$

Note that this expression only depends on the signaling through the transmit covariance matrix Σ_x , which implies the well-known result (e.g., [16]) that Gaussian signaling is not required to achieve $E_b/N_0|_{\min}$.

For illustrative purposes, we now particularize the signal model with $\Sigma_n = N_0 \mathbf{I}$ and obtain the first-order expansion of

the mutual information as a function of the inverse noise power N_0^{-1} as

$$I(N_0^{-1}) = N_0^{-1} \text{Tr}(\mathbf{H}\Sigma_x \mathbf{H}^\dagger) + o(N_0^{-1}) \quad (41)$$

which agrees with [14, Corollary 2].

IV. GRADIENTS OF INFORMATION MEASURES

In this section, we generalize some of the results of [6] by particularizing Theorem 1, and also deal with gradients of differential entropy (showing a new generalized version of De Bruijn's identity) and divergence.

A. Particularizations to SNR Gradients

The fundamental relation between the mutual information and the MMSE was thoroughly explored in [6] for the scalar Gaussian channel $y = \sqrt{\text{snr}}x + n$ and, to some extent, also for the vector Gaussian channel $\mathbf{y} = \sqrt{\text{snr}}\mathbf{x} + \mathbf{n}$ considering the trace of the MSE matrix as the estimation performance. We now show how the main results of this paper, namely Theorems 1 and 2, can be readily particularized to extend the results of [6]. The special case of the following result for $\Sigma_n = \mathbf{I}$ is given in [6, Theorem 2].

Corollary 1: Consider the signal model in (9), where all the terms are defined as in Theorem 2. Then

$$\frac{dI}{d\text{snr}} = \text{Tr}(\mathbf{H}^\dagger \Sigma_n^{-1} \mathbf{H} \mathbf{E}) \quad (42)$$

$$= \mathbb{E} \left[\left\| \Sigma_n^{-1/2} \mathbf{H} \mathbf{x} - \mathbb{E} \left[\Sigma_n^{-1/2} \mathbf{H} \mathbf{x} | \mathbf{y} \right] \right\|^2 \right]. \quad (43)$$

Proof: This is just a particular case of (8) with $\mathbf{B} = \sqrt{\text{snr}} \mathbf{I}$ (notice that $\mathbf{Q}_B = \mathbf{B} \mathbf{B}^\dagger$ in Theorem 2 plays the role of snr here). The result follows simply from the chain rule (note that $\mathbf{Q}_B = \text{snr} \mathbf{I}$) and Theorem 2 ($\nabla_{\mathbf{Q}_B} \mathbf{I} \mathbf{B} = \mathbf{H}^\dagger \Sigma_n^{-1} \mathbf{H} \mathbf{B} \mathbf{E}$)

$$\begin{aligned} \frac{dI}{d\text{snr}} &= \text{Tr} \left(\frac{\partial I}{\partial \mathbf{Q}_B^T} \frac{\partial \mathbf{Q}_B}{\partial \text{snr}} \right) \\ &= \text{Tr} \left(\frac{\partial I}{\partial \mathbf{Q}_B^*} \right) = \text{Tr}(\mathbf{H}^\dagger \Sigma_n^{-1} \mathbf{H} \mathbf{E}). \end{aligned} \quad (44)$$

□

Observe that (43) is decreasing in snr (since decreasing the noise variance improves the MMSE), which implies that $I(\text{snr})$ is an (increasing) concave function.

Now we can use Corollary 1 to generalize in several ways the result (1) obtained in [6, Theorem 1] for the scalar case. One obvious generalization is to include a scalar channel h in the signal model: $y = \sqrt{\text{snr}}hx + n$, obtaining

$$\frac{dI}{d\text{snr}} = |h|^2 \text{mmse}(\text{snr}) \quad (45)$$

where $\text{mmse}(\text{snr}) \triangleq \mathbb{E}[|x - \mathbb{E}[x | y]|^2]$. However, this signal model assumes a memoryless channel and also a memoryless source. Another interesting generalization is to allow a correlated input source (still with a memoryless channel) possibly originated by a code. Consider a block transmission

corresponding to N channel uses where the inputs in different blocks are independent

$$y_k = \sqrt{\text{snr}} h_k x_k + n_k, \quad 1 \leq k \leq N. \quad (46)$$

A direct application of Corollary 1 (after rewriting the signal model in vector form $\mathbf{y} = \sqrt{\text{snr}} \mathbf{H} \mathbf{x} + \mathbf{n}$ with $\mathbf{H} = \text{diag}(\{h_k\})$) gives (see also [6, Corollary 3])

$$\frac{dI}{d\text{snr}} = \sum_{k=1}^N |h_k|^2 \text{mmse}_k(\text{snr}) \quad (47)$$

where $\text{mmse}_k(\text{snr}) \triangleq \mathbb{E}[|x_k - \mathbb{E}[x_k | \mathbf{y}]|^2]$.

Next, we state a result for a multiuser scenario, which generalizes [6, Theorem 4].

Corollary 2: Consider the following signal model:

$$\mathbf{y} = \mathbf{H} \mathbf{\Gamma} \mathbf{x} + \mathbf{n} \quad (48)$$

where $\mathbf{\Gamma}$ is a square invertible matrix and the rest of the terms are defined as in Theorem 2. Then

$$\nabla_{\mathbf{\Gamma}^\dagger} I = \mathbf{H}^\dagger \Sigma_n^{-1} \mathbf{H} \mathbf{\Gamma} \mathbf{E} \mathbf{\Gamma}^{-1}. \quad (49)$$

Proof: This result follows from Theorem 2 with $\mathbf{B} = \mathbf{\Gamma}$. □

Observe that, when (48) models a multiuser system with the symbols transmitted by all K users stacked in \mathbf{x} and $\mathbf{\Gamma} = \text{diag}(\{\sqrt{\text{snr}_k}\}_{k=1}^K)$, then (49) is particularly useful since $\partial I / \partial \text{snr}_k = [\nabla_{\mathbf{\Gamma}^\dagger} I]_{kk}$.

B. Matrix Generalization of De Bruijn's Identity

For a multivariate density function $p_{\mathbf{y}}(\mathbf{y})$, De Bruijn's identity [2], [17]–[19] relates the derivative of the differential entropy $h(\cdot)$ to the Fisher information matrix defined as

$$\mathbf{J}(\mathbf{y}) \triangleq \mathbb{E}_{\mathbf{y}}[\nabla_{\mathbf{y}} \log p_{\mathbf{y}}(\mathbf{y}) \nabla_{\mathbf{y}}^\dagger \log p_{\mathbf{y}}(\mathbf{y})]. \quad (50)$$

Note that this is a special form of the Fisher information matrix (with respect to a translation parameter) which does not involve an explicit parameter as in its most general definition [19].⁴ The scalar version is defined as

$$J(\mathbf{y}) \triangleq \text{Tr}(\mathbf{J}(\mathbf{y})). \quad (51)$$

The original De Bruijn's identity was obtained for the scalar case [2], [17] as

$$\frac{d}{dt} h(z + \sqrt{t}n) = J(z + \sqrt{t}n) \quad (52)$$

where the noise n is standard Gaussian. The scalar version of the multivariate De Bruijn's identity (assuming a normalized covariance matrix for the noise $\Sigma_n = \mathbf{I}$) is [18], [19, Theorem 14], [6, eq. (51)]

$$\frac{d}{dt} h(\mathbf{z} + \sqrt{t}\mathbf{n}) = J(\mathbf{z} + \sqrt{t}\mathbf{n}). \quad (53)$$

⁴The Fisher information matrix with respect to a parameter $\boldsymbol{\theta}$ is defined as $\mathbf{J}(\mathbf{y}) \triangleq \mathbb{E}_{\mathbf{y}}[\nabla_{\boldsymbol{\theta}} \log p_{\mathbf{y}}(\mathbf{y}; \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}}^\dagger \log p_{\mathbf{y}}(\mathbf{y}; \boldsymbol{\theta})]$ [7].

A more general multivariate De Bruijn's identity is⁵

$$\frac{d}{dt}h(\mathbf{z} + \sqrt{t}\mathbf{A}\mathbf{n}) = \text{Tr}(\mathbf{A}^\dagger \mathbf{J}(\mathbf{z} + \sqrt{t}\mathbf{A}\mathbf{n})\mathbf{A}) \quad (54)$$

where again the noise has a normalized noise covariance matrix and \mathbf{A} is an arbitrary square invertible matrix. Note that in the context of the signal model in (5) we can set $\mathbf{z} = \mathbf{H}\mathbf{x}$ to particularize the results.

As can be observed from (52)–(54), the existing versions of De Bruijn's identity take a derivative with respect to a scalar parameter which takes the role of a noise-to-signal ratio. We now generalize these results by considering the gradient with respect to the noise covariance matrix and, more explicitly, with respect to an arbitrary linear transformation of the noise $\mathbf{T}^{1/2}\mathbf{n}$, where the linear transformation $\mathbf{T}^{1/2}$ plays the role of \sqrt{t} in the scalar version of the identity.

Theorem 4 (Matrix Version of the Multivariate De Bruijn's Identity): Consider an arbitrary random variable (with finite second-order moments) \mathbf{z} contaminated with a Gaussian noise \mathbf{n} independent of \mathbf{z} and with positive definite covariance matrix Σ_n . Then

$$\nabla_{\Sigma_n} h(\mathbf{z} + \mathbf{n}) = \mathbf{J}(\mathbf{z} + \mathbf{n}). \quad (55)$$

If \mathbf{z} is contaminated with the Gaussian random vector $\mathbf{T}^{1/2}\mathbf{n}$, where $\mathbf{T}^{1/2}$ is an arbitrary square invertible matrix, then

$$\nabla_{\mathbf{T}} h(\mathbf{z} + \mathbf{T}^{1/2}\mathbf{n}) = \mathbf{J}(\mathbf{z} + \mathbf{T}^{1/2}\mathbf{n}) \times \mathbf{T}^{1/2}\Sigma_n\mathbf{T}^{-1/2}. \quad (56)$$

Proof: See Appendix E. \square

Observe that when the noise has a normalized covariance matrix $\Sigma_n = \mathbf{I}$, (56) simplifies to

$$\nabla_{\mathbf{T}} h(\mathbf{z} + \mathbf{T}^{1/2}\mathbf{n}) = \mathbf{J}(\mathbf{z} + \mathbf{T}^{1/2}\mathbf{n}) \quad (57)$$

which is the natural generalization of the scalar identities (52)–(53).

The scalar expression of the multivariate De Bruijn's identity in (53) can be obtained from a simple application of the chain rule to Theorem 4 particularized for $\mathbf{T} = t\mathbf{I}$ as follows:

$$\frac{dh}{dt} = \text{Tr} \left(\frac{\partial h}{\partial \mathbf{T}^T} \frac{\partial \mathbf{T}}{\partial t} \right) = \text{Tr} \left(\frac{\partial h}{\partial \mathbf{T}^*} \right) = \text{Tr}(\mathbf{J}(\mathbf{z} + \sqrt{t}\mathbf{n})). \quad (58)$$

The more general expression for the the multivariate De Bruijn's identity in (54) also follows directly from Theorem 4 by particularizing to $\mathbf{T} = t\mathbf{A}\mathbf{A}^\dagger$ ($\mathbf{T}^{1/2} = \sqrt{t}\mathbf{A}$) and using the chain rule

$$\frac{dh}{dt} = \text{Tr} \left(\frac{\partial h}{\partial \mathbf{T}^T} \frac{\partial \mathbf{T}}{\partial t} \right) = \text{Tr} \left(\frac{\partial h}{\partial \mathbf{T}^*} \mathbf{A}\mathbf{A}^\dagger \right). \quad (59)$$

The multidimensional Cramer–Rao bound states that the Fisher information matrix is lower-bounded (i.e., the difference is nonnegative definite) by the inverse of its covariance matrix (e.g., [19]) (setting now $\mathbf{z} = \mathbf{H}\mathbf{x}$)

$$\begin{aligned} \mathbf{J}(\mathbf{y}) &\geq (\mathbb{E}[(\mathbf{y} - \mathbb{E}[\mathbf{y}])(\mathbf{y} - \mathbb{E}[\mathbf{y}])^\dagger])^{-1} \\ &= (\mathbf{H}\Sigma_x\mathbf{H}^\dagger + \Sigma_n)^{-1} \end{aligned} \quad (60)$$

⁵The special case of (54) for $t = 0$ is given in [20]. An alternative expression for the multivariate De Bruijn's identity was obtained in [21].

with equality if and only if \mathbf{y} is a Gaussian vector. Therefore, for Gaussian signaling, the following relation holds between the Fisher information matrix $\mathbf{J}(\mathbf{y})$ and the MMSE matrix \mathbf{E} :

$$\Sigma_x\mathbf{H}^\dagger\mathbf{J}(\mathbf{H}\mathbf{x} + \mathbf{n})\mathbf{H}\Sigma_x = \Sigma_x - \mathbf{E} \quad (61)$$

where we have used (15) and the matrix inversion lemma.⁶

C. Divergence Gradients

1) *Divergence for a Given Input:* The mutual information $I(\mathbf{x}; \mathbf{y})$ is equal to the conditional divergence $D(p_{\mathbf{y}|\mathbf{x}} \parallel p_{\mathbf{y}} | p_{\mathbf{x}})$. Since this quantity is an average of the divergence between conditional and unconditional output distributions, the following result is a refinement of Theorem 1.

Theorem 5: Consider the signal model in (5), where all the terms are defined as in Theorem 1. Then, the gradient of the divergence between the output conditioned on a given fixed input \mathbf{x}_0 and the unconditional output is

$$\begin{aligned} \nabla_{\mathbf{H}} D(p_{\mathbf{y}|\mathbf{x}=\mathbf{x}_0} \parallel p_{\mathbf{y}}) \\ = \mathbf{H}\mathbb{E}[(\mathbf{x}_0 - \mathbb{E}[\mathbf{x}|\mathbf{y}])(\mathbf{x}_0 - \mathbb{E}[\mathbf{x}|\mathbf{y}])^\dagger | \mathbf{x} = \mathbf{x}_0] \end{aligned} \quad (62)$$

where the expectation is with respect to the conditional output distribution $p_{\mathbf{y}|\mathbf{x}=\mathbf{x}_0}$.

Proof: See Appendix F. \square

Taking the expectation over \mathbf{x}_0 in (62) with respect to $p_{\mathbf{x}}$, we get

$$\mathbb{E}[\nabla_{\mathbf{H}} D(p_{\mathbf{y}|\mathbf{x}=\mathbf{x}_0} \parallel p_{\mathbf{y}})] = \mathbf{H}\mathbf{E}. \quad (63)$$

Interchanging gradient and expectation in (63), we obtain Theorem 1.⁷ We consider now the gradient of the conditional divergence but with the arguments swapped, i.e., of $D(p_{\mathbf{y}} \parallel p_{\mathbf{y}|\mathbf{x}})$.

Theorem 6: Consider the signal model in (5), where all the terms are defined as in Theorem 1. Then, the gradient of the divergence between the unconditional output and the output conditioned on a given fixed input \mathbf{x}_0 is

$$\begin{aligned} \nabla_{\mathbf{H}} D(p_{\mathbf{y}} \parallel p_{\mathbf{y}|\mathbf{x}=\mathbf{x}_0}) &= \mathbf{H}(\mathbb{E}[(\mathbf{x} - \mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0)^\dagger] - \mathbf{E}) \\ &= \mathbf{H}\mathbb{E}[(\mathbf{x}_0 - \mathbb{E}[\mathbf{x}|\mathbf{y}])(\mathbf{x}_0 - \mathbb{E}[\mathbf{x}|\mathbf{y}])^\dagger] \end{aligned} \quad (64)$$

where the expectation in (65) is with respect to the unconditional output distribution $p_{\mathbf{y}}$.

Proof: See Appendix G. \square

Observe the similarity between (62), where the expectation over \mathbf{y} is with respect to the conditional $p_{\mathbf{y}|\mathbf{x}=\mathbf{x}_0}$, and (65), where the expectation over \mathbf{y} is with respect to the unconditional $p_{\mathbf{y}}$.

Taking the expectation over \mathbf{x}_0 in (64) with respect to $p_{\mathbf{x}}$, we get

$$\mathbb{E}[\nabla_{\mathbf{H}} D(p_{\mathbf{y}} \parallel p_{\mathbf{y}|\mathbf{x}=\mathbf{x}_0})] = \mathbf{H}(2\Sigma_x - \mathbf{E}). \quad (66)$$

⁶Matrix inversion lemma:

$$(\mathbf{A} + \mathbf{B}\mathbf{C}\mathbf{D})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D}\mathbf{A}^{-1}\mathbf{B} + \mathbf{C}^{-1})^{-1}\mathbf{D}\mathbf{A}^{-1}.$$

⁷The validity of the interchange between gradient and expectation follows from the boundedness of (62) and the Lebesgue Dominated Convergence Theorem [22].

2) *Divergence of Signal-Plus-Noise Versus Noise*: The following result relates two fundamental quantities: the likelihood ratio (between the hypotheses that the random signal \mathbf{x} has been transmitted and nothing has been transmitted) and the MMSE estimation, which play central roles in detection and estimation theory, respectively, [7], [23] (see [6, Theorem 5] for the scalar version).

Theorem 7: Consider the signal model in (5), where all the terms are defined as in Theorem 1. Then, the gradient of the divergence between the signal-plus-noise and noise distributions is

$$\nabla_{\mathbf{H}} D(p_{\mathbf{y}} \parallel p_{\mathbf{n}}) = \mathbf{H}(\mathbb{E}[\mathbf{x}\mathbf{x}^\dagger] - \mathbf{E}) \quad (67)$$

$$= \mathbf{H}\mathbb{E}[\mathbb{E}[\mathbf{x} \mid \mathbf{y}]\mathbb{E}[\mathbf{x}^\dagger \mid \mathbf{y}]]. \quad (68)$$

Proof: Set $\mathbf{x}_0 = 0$ in Theorem 6. \square

3) *Non-Gaussianness*: The non-Gaussianness of a random variable is defined in [24] as the divergence between its distribution and a Gaussian distribution with the same first- and second-order moments. The gradient of the non-Gaussianness can be directly obtained from Theorem 1 (similarly, Theorem 2) by noting that the mutual information admits the following decomposition (see, for example, [24]):

$$I(\mathbf{x}; \mathbf{y}) = I(\mathbf{x}_G; \mathbf{y}_G) - D(p_{\mathbf{y}} \parallel p_{\mathbf{y}_G}) \quad (69)$$

where $(\mathbf{x}_G, \mathbf{y}_G)$ denote jointly Gaussian-distributed random vectors with the same mean and covariance matrix as the joint distribution of (\mathbf{x}, \mathbf{y}) and $D(p_{\mathbf{y}} \parallel p_{\mathbf{y}_G})$ is the non-Gaussianness of \mathbf{y} .

Theorem 8: Consider the signal model in (5), where all the terms are defined as in Theorem 2. Then, the gradient of the non-Gaussianness of \mathbf{y} is

$$\nabla_{\mathbf{H}} D(p_{\mathbf{y}} \parallel p_{\mathbf{y}_G}) = \Sigma_n^{-1} \mathbf{H}(\mathbf{E}_{\text{lin}} - \mathbf{E}) \quad (70)$$

where $\mathbf{E}_{\text{lin}} = (\Sigma_x^{-1} + \mathbf{H}^\dagger \Sigma_n^{-1} \mathbf{H})^{-1}$ is the linear MMSE matrix, i.e., the MMSE matrix among the class of *linear* estimators.

Using (70) and the fact that

$$D(p_{\mathbf{H}\mathbf{x}} \parallel p_{\mathbf{H}\mathbf{x}_G}) = \int_0^\infty (dD(p_{\mathbf{y}} \parallel p_{\mathbf{y}_G})/d\text{snr}) d\text{snr}$$

for the signal model in (9) (since $D(p_{\mathbf{H}\mathbf{x}} \parallel p_{\mathbf{H}\mathbf{x}_G}) = \lim_{\text{snr} \rightarrow \infty} D(p_{\mathbf{y}} \parallel p_{\mathbf{y}_G})$ [6, Lemma 7]), we can write the following alternative expression for the differential entropy of a random vector (setting $\mathbf{H} = \mathbf{I}$ and $\Sigma_n = \mathbf{I}$) which generalizes the corresponding scalar version in [6, Theorem 14]:

$$\begin{aligned} h(\mathbf{x}) &= h(\mathbf{x}_G) - D(p_{\mathbf{x}} \parallel p_{\mathbf{x}_G}) \\ &= \log \det(\pi e \Sigma_x) - \int_0^\infty \text{Tr}(\mathbf{E}_{\text{lin}}(\text{snr}) - \mathbf{E}(\text{snr})) d\text{snr} \end{aligned} \quad (71)$$

where in this case $\mathbf{E}_{\text{lin}}(\text{snr}) = (\Sigma_x^{-1} + \text{snr} \mathbf{I})^{-1}$. Note that the integral is always a nonnegative quantity since $\mathbf{E}_{\text{lin}}(\text{snr}) \geq \mathbf{E}(\text{snr})$, which agrees with the fact that the Gaussian distribution has the highest differential entropy for given second-order moments. Recall that $\text{Tr}(\mathbf{E}(\text{snr})) = \mathbb{E}[\|\mathbf{x} - \mathbb{E}[\mathbf{x} \mid \mathbf{y}]\|^2]$ is the mean-squared norm of the estimation error vector.

V. PRACTICAL APPLICATION: PRECODER OPTIMIZATION

The availability of the mutual information gradient enables us to optimize with respect to parameters of interest. One can always compute the gradient numerically and use it in a gradient-based algorithm to search the maximizing solution (which may be only locally optimum if the problem is not convex). In some special cases, it is even possible to derive optimality conditions in closed form. Recall that for a Gaussian input, the maximizing solution is the well-known channel diagonalization plus water-filling power allocation (e.g., [13]). In the following, we focus on a numerical optimization (for an arbitrary input) based on a simple gradient projection method.

Consider the optimization of the linear precoder \mathbf{B} to maximize the mutual information of the general signal model $\mathbf{y} = \mathbf{H}\mathbf{B}\mathbf{x} + \mathbf{n}$ subject to the power constraint $\text{Tr}(\mathbf{B}\mathbf{B}^\dagger) \leq P_T$. To deal with this constrained optimization problem, we can use a simple gradient update with an additional projection to guarantee the feasibility of the new solution [25]–[27]

$$\mathbf{B}(k+1) = [\mathbf{B}(k) + \mu \nabla_{\mathbf{B}} I(k)]_{\text{Tr}(\mathbf{B}\mathbf{B}^\dagger) \leq P_T}^+ \quad (72)$$

where k denotes the iteration index, μ is the stepsize of the update, the gradient of the mutual information at the k th iteration is given by $\nabla_{\mathbf{B}} I(k) = \mathbf{H}^\dagger \Sigma_n^{-1} \mathbf{H}\mathbf{B}(k)\mathbf{E}(k)$ (from Theorem 2), and $[\cdot]_{\text{Tr}(\mathbf{B}\mathbf{B}^\dagger) \leq P_T}^+$ denotes the projection onto the feasible set $\text{Tr}(\mathbf{B}\mathbf{B}^\dagger) \leq P_T$. For simplicity, we choose the stepsize μ as a fixed small number; however, other more sophisticated choices can be made with better convergence properties [26].

For the numerical computation of the MMSE matrix, which is required to evaluate the gradient $\nabla_{\mathbf{B}} I(k)$ in (72), closed-form expressions may be used in case they are available; otherwise, one can resort to a Monte Carlo method to approximate the MMSE matrix as follows:

$$\begin{aligned} \mathbf{E} &= \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x} \mid \mathbf{y}])(\mathbf{x} - \mathbb{E}[\mathbf{x} \mid \mathbf{y}])^\dagger] \\ &\approx \frac{1}{S} \sum_{\mathbf{x}_0, \mathbf{y}_0} (\mathbf{x}_0 - \mathbb{E}[\mathbf{x} \mid \mathbf{y} = \mathbf{y}_0])(\mathbf{x}_0 - \mathbb{E}[\mathbf{x} \mid \mathbf{y} = \mathbf{y}_0])^\dagger \end{aligned} \quad (73)$$

where the S random samples $\mathbf{x}_0, \mathbf{y}_0$ are generated from $p_{\mathbf{x}} p_{\mathbf{y} \mid \mathbf{x}}$ and the conditional estimate is computed as

$$\mathbb{E}[\mathbf{x} \mid \mathbf{y} = \mathbf{y}_0] = \sum_{\mathbf{x}_0} \mathbf{x}_0 p_{\mathbf{x} \mid \mathbf{y}}(\mathbf{x}_0 \mid \mathbf{y}_0) \quad (74)$$

where the summation is over all possible input vectors.

The projection of a point \mathbf{B}_0 onto the feasible set $\text{Tr}(\mathbf{B}\mathbf{B}^\dagger) \leq P_T$ (as required in (72)) is in turn the solution to the following convex optimization problem:

$$\begin{aligned} &\underset{\mathbf{B}}{\text{minimize}} && \|\mathbf{B} - \mathbf{B}_0\|^2 \\ &\text{subject to} && \|\mathbf{B}\|^2 \leq P_T. \end{aligned} \quad (75)$$

By forming the Lagrangian and deriving the Karush–Kuhn–Tucker (KKT) optimality conditions [28], one can easily find that the optimal solution to (75) is proportional to \mathbf{B}_0 up to a scaling factor chosen to satisfy the power constraint.

Summarizing, we propose the following iterative numerical algorithm to optimize the linear precoder \mathbf{B} .

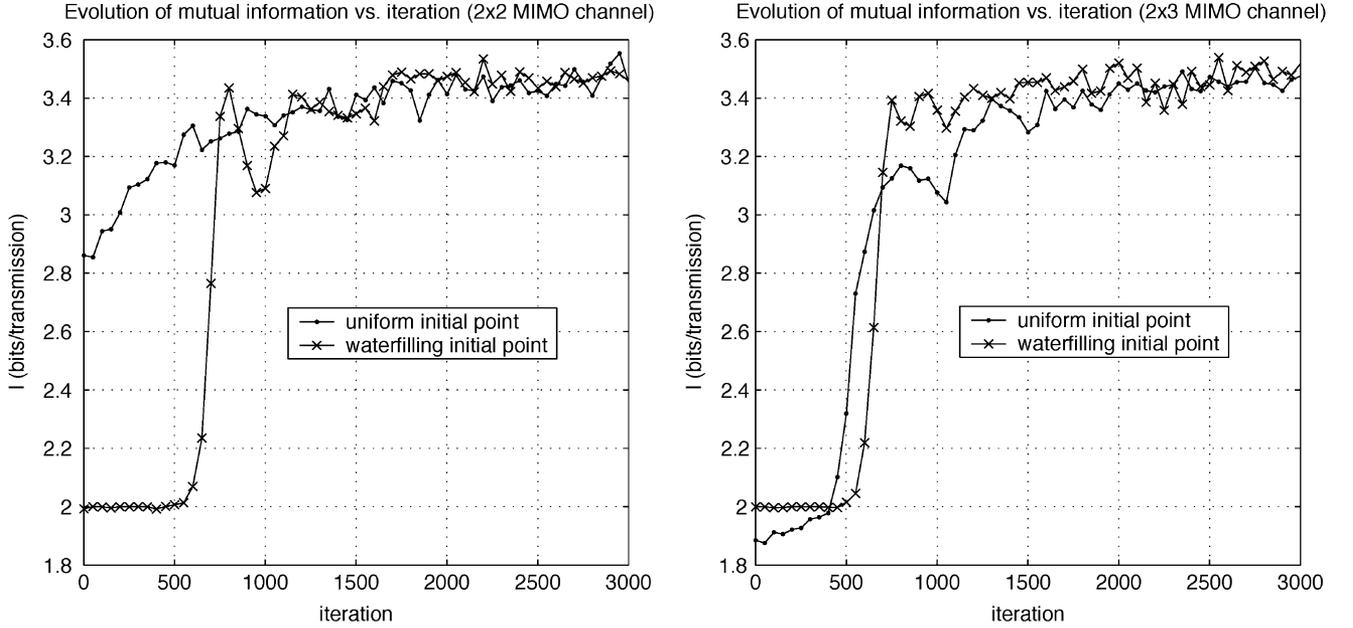


Fig. 2. Evolution of the mutual information (for a 2×2 channel matrix and for a 2×3 channel matrix) as the linear precoder is iteratively optimized with the proposed gradient algorithm.

Algorithm:

0. Initialization: set $k = 0$ and initialize the precoder $\mathbf{B}(0)$ to some sensible choice such as
 - uniform power allocation (equivalent to no precoder): $\mathbf{B}(0) = P_T/L[\mathbf{I} \ \mathbf{0}]^T$;
 - channel-diagonalization plus water filling (optimum for Gaussian inputs):

$$\mathbf{B}(0) = \mathbf{U} \text{diag}(\{\sqrt{p_i}\}),$$

where \mathbf{U} contains as columns the right singular vectors of the channel matrix corresponding to the L largest squared singular values λ_i and $p_i = (\mu - \lambda_i^{-1})^+$ is the water-filling power allocation.

- 1 Update the precoder with the gradient (after updating $\mathbf{E}(k)$)

$$\mathbf{B}(k+1) \leftarrow \mathbf{B}(k) + \mu \mathbf{H}^\dagger \Sigma_n^{-1} \mathbf{H} \mathbf{B}(k) \mathbf{E}(k). \quad (76)$$

- 2 If $\|\mathbf{B}(k+1)\|^2 > P_T$, then normalize

$$\mathbf{B}(k+1) \leftarrow \mathbf{B}(k+1) \sqrt{P_T / \|\mathbf{B}(k+1)\|^2}. \quad (77)$$

- 3 Set $k \leftarrow k+1$ and go to Step 1 (until some termination criterion is activated).

Observe that, in principle, at each iteration of the algorithm one has to estimate the MMSE matrix for the given linear precoder at the current iteration. This can be too time-consuming and a convenient solution is to update the MMSE matrix recursively. The rationale underneath this practical trick is that if the stepsize μ is sufficiently small, then the linear precoder will change very slowly, which implies that we can take advantage of the previously computed MMSE matrix. In particular, we have updated the MMSE matrix as follows:

$$\mathbf{E}(k+1) = \alpha \mathbf{E}(k) + (1 - \alpha) \mathbf{E}^{\text{new}}(k+1) \quad (78)$$

where α is the forgetting factor which measures how fast old data are forgotten (with a typical value around 0.9–0.99) and $\mathbf{E}^{\text{new}}(k+1)$ is the MMSE matrix estimated with new samples (typically just a few samples) based on the new precoder $\mathbf{B}(k+1)$.

Fig. 2 illustrates the convergence of the algorithm for the 2×2 channel matrix

$$\mathbf{H} = \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix}$$

over which two independent and equiprobable quaternary phase-shift keying (QPSK) symbols (drawn from $\{\pm 1, \pm j\}$) are transmitted. The precoder is first initialized to a uniform power allocation $\mathbf{B}(0) = \mathbf{I}$ and then to the channel-diagonalizing plus water-filling solution, which happens to be a bad initialization in this case because it only allocates power to one channel eigenmode which translates into a maximum of 2 bits/transmission. Each new MMSE matrix \mathbf{E}^{new} is computed with just $S = 10$ random samples and the averaged matrix \mathbf{E} is obtained with a forgetting factor of $\alpha = 0.99$.

Fig. 2 also illustrates the convergence for a 2×3 channel matrix obtained by appending a zero column on the left of the previous channel matrix. The 3×2 precoder matrix is again initialized to both a uniform power allocation $\mathbf{B}(0) = [\mathbf{I} \ \mathbf{0}]^T$ and to the channel-diagonalizing plus water-filling solution. As expected, the precoder matrix evolves to have a zero top row (to avoid the zero column of the channel matrix) and the two bottom rows are as in the previous 2×2 case, obtaining therefore the same final value of mutual information.

VI. CONCLUSION

Building upon the setting in [6], this paper has obtained the gradient of the mutual information of a general linear vector

Gaussian channel with arbitrary signaling with respect to several key system parameters. The basic expression of the gradient with respect to the channel matrix happens to involve the MMSE matrix, connecting the two key quantities from information theory and estimation theory. Gradients and derivatives with respect to arbitrary parameters are then easily obtained with the application of the chain rule for differentiation. Several interesting formulas have been obtained by applying the main result: first-order approximations of the mutual information, matrix generalization of De Bruijn's identity, gradient of the divergence between the conditional and unconditional outputs, gradient of the divergence between signal-plus-noise and noise, and gradient of the non-Gaussianness. We have illustrated the practical application of our results with an optimization problem that arises in the design of a linear precoder that maximizes mutual information.

APPENDIX

A. A Note on Complex Derivatives and Gradients

In general, commonly encountered functions of complex variables are not analytic and a more generalized definition of complex derivative needs to be used [29]. In particular, we use the well-known definition of the complex derivative of a real-valued scalar function f

$$\frac{df}{dx^*} \triangleq \frac{1}{2} \left(\frac{\partial f}{\partial \text{Re}\{x\}} + j \frac{\partial f}{\partial \text{Im}\{x\}} \right) \quad (79)$$

as in [29], so that, for example

$$\frac{\partial}{\partial \mathbf{x}^*} (\mathbf{x}^\dagger \mathbf{R} \mathbf{x}) = \mathbf{R} \mathbf{x} \quad (80)$$

(for \mathbf{R} independent of \mathbf{x}). The derivative of a complex-valued scalar function is defined componentwise on the real and imaginary part, i.e., $\partial f / \partial \mathbf{x}^* = \partial \text{Re}\{f\} / \partial \mathbf{x}^* + j \partial \text{Im}\{f\} / \partial \mathbf{x}^*$.

For the sake of notation, we define the complex gradient vector as $\nabla_{\mathbf{x}} f \triangleq \partial f / \partial \mathbf{x}^*$ as in [27].⁸ The complex gradient matrix is similarly defined as $\nabla_{\mathbf{X}} f \triangleq \partial f / \partial \mathbf{X}^*$, where $[\nabla_{\mathbf{X}} f]_{ij} = \partial f / \partial [\mathbf{X}^*]_{ij}$, so that, for example, [15]

$$\nabla_{\mathbf{X}} \text{Tr}(\mathbf{X}^\dagger \mathbf{R}) = \mathbf{R} \quad (81)$$

and

$$\nabla_{\mathbf{X}} \log \det(\mathbf{I} + \mathbf{X}^\dagger \mathbf{R} \mathbf{X}) = \mathbf{R} \mathbf{X} (\mathbf{I} + \mathbf{X}^\dagger \mathbf{R} \mathbf{X})^{-1}. \quad (82)$$

The chain rule for a nonanalytic real-valued scalar function f of a real vector \mathbf{r} and a complex vector \mathbf{c} , denoted explicitly by $f(\mathbf{r}, \mathbf{c}, \mathbf{c}^*)$, is

$$\begin{aligned} \frac{\partial f(\mathbf{r}, \mathbf{c}, \mathbf{c}^*)}{\partial x} &= \left(\frac{\partial f(\mathbf{r}, \mathbf{c}, \mathbf{c}^*)}{\partial \mathbf{r}} \right)^T \frac{\partial \mathbf{r}}{\partial x} \\ &+ \left(\frac{\partial f(\mathbf{r}, \mathbf{c}, \mathbf{c}^*)}{\partial \mathbf{c}} \right)^T \frac{\partial \mathbf{c}}{\partial x} + \left(\frac{\partial f(\mathbf{r}, \mathbf{c}, \mathbf{c}^*)}{\partial \mathbf{c}^*} \right)^T \frac{\partial \mathbf{c}^*}{\partial x} \end{aligned}$$

where the last term vanishes for analytic functions, obtaining then the classical chain rule.

⁸Note, however, that in [27] there is an additional factor of 2.

B. Proof of Theorem 1

Observe that the means of the transmitted signal and noise are completely arbitrary in the statement of the theorem. In fact, the result is not affected by the means since both the mutual information and the MMSE matrix are insensitive to the means. For simplicity, we will assume the noise to be zero mean without loss of generality. The justification of the interchange of the order of integrals and derivatives is left for the end of the proof. In the scalar case, this proof reduces to the proof in [6, Appendix III].

Since the noise is Gaussian with zero mean and normalized covariance matrix, the conditional output pdf is

$$p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) = \frac{1}{\pi^{n_R}} \exp(-\|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2) \quad (83)$$

and the unconditional output pdf is $p_{\mathbf{y}}(\mathbf{y}) = \mathbb{E}_{\mathbf{x}}[p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})]$. The mutual information can be written as

$$\begin{aligned} I(\mathbf{x}; \mathbf{y}) &= \mathbb{E} \left[\log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})} \right] \\ &= -n_R \log(\pi e) - \mathbb{E}[\log p_{\mathbf{y}}(\mathbf{y})] \\ &= -n_R \log(\pi e) - \int p_{\mathbf{y}}(\mathbf{y}) \log p_{\mathbf{y}}(\mathbf{y}) d\mathbf{y}. \quad (84) \end{aligned}$$

Then

$$\frac{\partial I}{\partial \mathbf{H}^*} = - \int (1 + \log p_{\mathbf{y}}(\mathbf{y})) \frac{\partial p_{\mathbf{y}}(\mathbf{y})}{\partial \mathbf{H}^*} d\mathbf{y} \quad (85)$$

$$= - \int (1 + \log p_{\mathbf{y}}(\mathbf{y})) \mathbb{E}_{\mathbf{x}} \left[\frac{\partial p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})}{\partial \mathbf{H}^*} \right] d\mathbf{y}. \quad (86)$$

The derivative of the conditional output can be written as

$$\begin{aligned} \frac{\partial p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})}{\partial \mathbf{H}^*} &= -p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) \frac{\partial}{\partial \mathbf{H}^*} ((\mathbf{y} - \mathbf{H}\mathbf{x})^\dagger (\mathbf{y} - \mathbf{H}\mathbf{x})) \\ &= p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) (\mathbf{y} - \mathbf{H}\mathbf{x}) \mathbf{x}^\dagger \\ &= -\nabla_{\mathbf{y}} p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) \mathbf{x}^\dagger. \quad (87) \end{aligned}$$

Therefore, we have

$$\begin{aligned} \frac{\partial I}{\partial \mathbf{H}^*} &= \int (1 + \log p_{\mathbf{y}}(\mathbf{y})) \mathbb{E}_{\mathbf{x}} [\nabla_{\mathbf{y}} p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) \mathbf{x}^\dagger] d\mathbf{y} \\ &= \mathbb{E} \left[\left(\int (1 + \log p_{\mathbf{y}}(\mathbf{y})) \nabla_{\mathbf{y}} p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) d\mathbf{y} \right) \mathbf{x}^\dagger \right] \\ &= \mathbb{E} \left[\left(- \int \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})} \nabla_{\mathbf{y}} p_{\mathbf{y}}(\mathbf{y}) d\mathbf{y} \right) \mathbf{x}^\dagger \right] \quad (88) \end{aligned}$$

where the last equality follows from integrating by parts [22]⁹ and noting that $p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})(1 + \log p_{\mathbf{y}}(\mathbf{y})) \rightarrow 0$ as $\|\mathbf{y}\| \rightarrow \infty$.

⁹The following real integration by parts should be applied to the real and imaginary parts of each element of \mathbf{y} (noting that integration on \mathbb{C}^n can be interpreted as integration on \mathbb{R}^{2n} [7])

$$\begin{aligned} &\int_{-\infty}^{+\infty} (1 + \log_e p_{\mathbf{y}}) (\partial p_{\mathbf{y}|\mathbf{x}} / \partial t) dt \\ &= [(1 + \log_e p_{\mathbf{y}}) p_{\mathbf{y}|\mathbf{x}}]_{t=-\infty}^{t=+\infty} - \int_{-\infty}^{+\infty} (1/p_{\mathbf{y}}) (\partial p_{\mathbf{y}} / \partial t) p_{\mathbf{y}|\mathbf{x}} dt \end{aligned}$$

where t represents a real variable [22].

Noting that, for almost every \mathbf{y} , $p_{\mathbf{x}|\mathbf{y}} \ll p_{\mathbf{x}}$ and the Radon–Nikodym derivative is $\frac{dp_{\mathbf{x}|\mathbf{y}}}{dp_{\mathbf{x}}}(\mathbf{x}, \mathbf{y}) = p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})/p_{\mathbf{y}}(\mathbf{y})$, we can write

$$\begin{aligned} \frac{\partial I}{\partial \mathbf{H}^*} &= - \int \nabla_{\mathbf{y}} p_{\mathbf{y}}(\mathbf{y}) \mathbb{E}_{\mathbf{x}} \left[\frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})} \mathbf{x}^\dagger \right] d\mathbf{y} \\ &= - \int \nabla_{\mathbf{y}} p_{\mathbf{y}}(\mathbf{y}) \mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}] d\mathbf{y}. \end{aligned} \quad (89)$$

Now, using

$$\begin{aligned} \nabla_{\mathbf{y}} p_{\mathbf{y}}(\mathbf{y}) &= \nabla_{\mathbf{y}} \mathbb{E}_{\mathbf{x}} [p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})] = \mathbb{E}_{\mathbf{x}} [\nabla_{\mathbf{y}} p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})] \\ &= -\mathbb{E}_{\mathbf{x}} [p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})(\mathbf{y} - \mathbf{H}\mathbf{x})] \\ &= -\mathbb{E} [p_{\mathbf{y}}(\mathbf{y})(\mathbf{y} - \mathbf{H}\mathbf{x}) | \mathbf{y}] \\ &= -p_{\mathbf{y}}(\mathbf{y})(\mathbf{y} - \mathbf{H}\mathbb{E}[\mathbf{x} | \mathbf{y}]) \end{aligned} \quad (90)$$

we can write

$$\begin{aligned} \frac{\partial I}{\partial \mathbf{H}^*} &= \int p_{\mathbf{y}}(\mathbf{y})(\mathbf{y} - \mathbf{H}\mathbb{E}[\mathbf{x} | \mathbf{y}]) \mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}] d\mathbf{y} \\ &= \mathbb{E}[\mathbf{y}\mathbf{x}^\dagger] - \mathbb{E}[\mathbf{H}\mathbb{E}[\mathbf{x} | \mathbf{y}]\mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}]] \\ &= \mathbf{H}(\mathbb{E}[\mathbf{x}\mathbf{x}^\dagger] - \mathbb{E}[\mathbb{E}[\mathbf{x} | \mathbf{y}]\mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}]]). \end{aligned} \quad (91)$$

Finally, noting that

$$\begin{aligned} \mathbf{E} &= \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x} | \mathbf{y}])(\mathbf{x} - \mathbb{E}[\mathbf{x} | \mathbf{y}])^\dagger] \\ &= \mathbb{E}[\mathbf{x}\mathbf{x}^\dagger + \mathbb{E}[\mathbf{x} | \mathbf{y}]\mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}] - \mathbb{E}[\mathbf{x} | \mathbf{y}]\mathbf{x}^\dagger - \mathbf{x}\mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}]] \\ &= \mathbb{E}[\mathbf{x}\mathbf{x}^\dagger] - \mathbb{E}[\mathbb{E}[\mathbf{x} | \mathbf{y}]\mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}]] \end{aligned} \quad (92)$$

we obtain the desired result $\partial I / \partial \mathbf{H}^* = \mathbf{H}\mathbf{E}$.

Interchange of Order of Integrals and Derivatives: Consider first the interchange of order in (86) given by

$$\partial \mathbb{E}_{\mathbf{x}} [p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})] / \partial [\mathbf{H}^*]_{ij} = \mathbb{E}_{\mathbf{x}} [\partial p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) / \partial [\mathbf{H}^*]_{ij}].$$

Simply note from (87) that

$$\partial p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) / \partial [\mathbf{H}^*]_{ij} = p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) [(\mathbf{y} - \mathbf{H}\mathbf{x})\mathbf{x}^\dagger]_{ij}$$

and then

$$\begin{aligned} \left| \frac{\partial p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})}{\partial [\mathbf{H}^*]_{ij}} \right| &\leq |[(\mathbf{y} - \mathbf{H}\mathbf{x})\mathbf{x}^\dagger]_{ij}| \\ &\leq |y_i| |x_j| + \sum_k |[\mathbf{H}]_{ik}| |x_k| |x_j| \end{aligned}$$

which is an upper bound valid for the given channel \mathbf{H} . To make the upper bound valid for a neighborhood around \mathbf{H} , just add some ϵ to each term $|[\mathbf{H}]_{ik}|$. The expectation over \mathbf{x} of the right-hand side is clearly finite since the second-order moments are finite¹⁰ and Lemma 2 below can be invoked.

Consider now the interchange of order in (90) given by $\nabla_{\mathbf{y}} \mathbb{E}_{\mathbf{x}} [p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})] = \mathbb{E}_{\mathbf{x}} [\nabla_{\mathbf{y}} p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})]$. Simply note from (87) that $\partial p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) / \partial y_i^* = -p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) [(\mathbf{y} - \mathbf{H}\mathbf{x})_i]$ and then

$$\left| \frac{\partial p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})}{\partial y_i^*} \right| = |(\mathbf{y} - \mathbf{H}\mathbf{x})_i| \leq |y_i| + \sum_k |[\mathbf{H}]_{ik}| |x_k|.$$

As before, to guarantee that the upper bound holds on some neighborhood of \mathbf{y} just add some small ϵ to the term $|y_i|$. The

¹⁰Observe that if $\int |x|^2 f(x) dx$ is finite, so is $\int |x| f(x) dx$ [30]. Also, by Hölder's inequality (for $p = q = 2$) $\mathbb{E}[|x||y|] \leq \mathbb{E}^{1/2}[|x|^2] \mathbb{E}^{1/2}[|y|^2]$ [30], if $\mathbb{E}[|x|^2]$ and $\mathbb{E}[|y|^2]$ are finite, so is $\mathbb{E}[|x||y|]$.

expectation over \mathbf{x} of the right-hand side is clearly finite since the second-order moments are finite and Lemma 2 below can be invoked. \square

The following result was used in the foregoing proof.

Lemma 2: Consider a function $f(\mathbf{x}; \theta)$ and a nonnegative function $g(\mathbf{x})$. The relation

$$\frac{\partial}{\partial \theta} \int g(\mathbf{x}) f(\mathbf{x}; \theta) d\mathbf{x} = \int g(\mathbf{x}) \frac{\partial}{\partial \theta} f(\mathbf{x}; \theta) d\mathbf{x} \quad (93)$$

holds¹¹ if for each θ_0 there exists a neighborhood of $\theta_0, \mathcal{N}_{\theta_0}$, and a function $M(\mathbf{x}; \theta_0)$ such that

$$\sup_{\theta \in \mathcal{N}_{\theta_0}} \left| \frac{\partial}{\partial \theta} f(\mathbf{x}; \theta) \right| \leq M(\mathbf{x}; \theta_0) \quad \text{a.e.}$$

with $\int g(\mathbf{x}) M(\mathbf{x}; \theta) d\mathbf{x} < \infty$.

Proof: For any θ_0 , the mean value theorem [22] says that

$$f_\delta(\mathbf{x}; \theta_0) \triangleq \frac{f(\mathbf{x}; \theta_0 + \delta) - f(\mathbf{x}; \theta_0)}{\delta} = \frac{\partial}{\partial \theta} f(\mathbf{x}; \tilde{\theta})$$

where $\tilde{\theta}$ is in a neighborhood of θ_0 . Therefore, $f_\delta(\mathbf{x}; \theta_0)$ can be upper-bounded as

$$|f_\delta(\mathbf{x}; \theta_0)| \leq \sup_{\theta \in \mathcal{N}_{\theta_0}} \left| \frac{\partial}{\partial \theta} f(\mathbf{x}; \theta) \right| \leq M(\mathbf{x}; \theta_0) \quad \text{a.e.}$$

Now, if $\int g(\mathbf{x}) M(\mathbf{x}; \theta_0) d\mathbf{x} < \infty$, we can invoke the Lebesgue Dominated Convergence Theorem [22], [30] to show that

$$\lim_{\delta \rightarrow 0} \int g(\mathbf{x}) f_\delta(\mathbf{x}; \theta_0) d\mathbf{x} = \int \lim_{\delta \rightarrow 0} g(\mathbf{x}) f_\delta(\mathbf{x}; \theta_0) d\mathbf{x}. \quad \square$$

C. Proof of Theorem 2

We first give the following two chain rules which will be instrumental in proving Theorem 2.

Lemma 3:

- Let f be a scalar real-valued function which depends on \mathbf{B} through $\mathbf{H} = \mathbf{A}\mathbf{B}\mathbf{C}$, where \mathbf{A} and \mathbf{C} are arbitrary fixed matrices. Then

$$\nabla_{\mathbf{B}} f = \mathbf{A}^\dagger \nabla_{\mathbf{H}} f \mathbf{C}^\dagger. \quad (94)$$

- Let f be a scalar real-valued function which depends on \mathbf{B} through $\mathbf{R} \triangleq \mathbf{H}\mathbf{H}^\dagger = \mathbf{A}\mathbf{B}\mathbf{C}\mathbf{C}^\dagger\mathbf{B}^\dagger\mathbf{A}^\dagger$, where \mathbf{A} and \mathbf{C} are arbitrary fixed matrices. Then

$$\nabla_{\mathbf{B}} f = \mathbf{A}^\dagger \nabla_{\mathbf{R}} f \mathbf{A}\mathbf{B}\mathbf{C}\mathbf{C}^\dagger. \quad (95)$$

Proof: Consider the first chain rule. Given $[\mathbf{H}^\dagger]_{ij} = \text{Tr}(\mathbf{C}^\dagger \mathbf{B}^\dagger \mathbf{A}^\dagger \mathbf{e}_j \mathbf{e}_i^\dagger)$, it follows that

$$\begin{aligned} \partial [\mathbf{H}^\dagger]_{ij} / \partial \mathbf{B}^* &= \mathbf{A}^\dagger \mathbf{e}_j \mathbf{e}_i^\dagger \mathbf{C}^\dagger, \partial [\mathbf{H}^\dagger]_{ij} / \partial [\mathbf{B}^*]_{kl} \\ &= \mathbf{e}_k^\dagger \mathbf{A}^\dagger \mathbf{e}_j \mathbf{e}_i^\dagger \mathbf{C}^\dagger \mathbf{e}_l \\ &= \mathbf{e}_i^\dagger \mathbf{C}^\dagger \mathbf{e}_l \mathbf{e}_k^\dagger \mathbf{A}^\dagger \mathbf{e}_j \end{aligned}$$

¹¹It is implicitly assumed in (93) that $\int g(\mathbf{x}) f(\mathbf{x}; \theta) d\mathbf{x}$ and $\partial f(\mathbf{x}; \theta) / \partial \theta$ exist.

and, more compactly, $\partial \mathbf{H}^\dagger / \partial [\mathbf{B}^*]_{kl} = \mathbf{C}^\dagger \mathbf{e}_l \mathbf{e}_k^\dagger \mathbf{A}^\dagger$ (note also that $\partial \mathbf{H} / \partial [\mathbf{B}^*]_{kl} = \mathbf{0}$). Now, we can use the chain rule

$$\begin{aligned} \frac{\partial f}{\partial [\mathbf{B}^*]_{kl}} &= \text{Tr} \left(\left(\frac{\partial f}{\partial \mathbf{H}^\dagger} \right)^T \frac{\partial \mathbf{H}^\dagger}{\partial [\mathbf{B}^*]_{kl}} \right) \\ &= \text{Tr} \left(\mathbf{e}_k^\dagger \mathbf{A}^\dagger \frac{\partial f}{\partial \mathbf{H}^\dagger} \mathbf{C}^\dagger \mathbf{e}_l \right) \end{aligned} \quad (96)$$

to obtain the desired result $\partial f / \partial \mathbf{B}^* = \mathbf{A}^\dagger \partial f / \partial \mathbf{H}^\dagger \mathbf{C}^\dagger$. Consider now the second chain rule. Given

$$[\mathbf{R}]_{ij} = \text{Tr}(\mathbf{A} \mathbf{B} \mathbf{C} \mathbf{C}^\dagger \mathbf{B}^\dagger \mathbf{A}^\dagger \mathbf{e}_j \mathbf{e}_i^\dagger)$$

it follows that

$$\begin{aligned} \partial [\mathbf{R}]_{ij} / \partial \mathbf{B}^* &= \mathbf{A}^\dagger \mathbf{e}_j \mathbf{e}_i^\dagger \mathbf{A} \mathbf{B} \mathbf{C} \mathbf{C}^\dagger \\ \partial [\mathbf{R}]_{ij} / \partial [\mathbf{B}^*]_{kl} &= \mathbf{e}_k^\dagger \mathbf{A}^\dagger \mathbf{e}_j \mathbf{e}_i^\dagger \mathbf{A} \mathbf{B} \mathbf{C} \mathbf{C}^\dagger \mathbf{e}_l \\ &= \mathbf{e}_i^\dagger \mathbf{A} \mathbf{B} \mathbf{C} \mathbf{C}^\dagger \mathbf{e}_l \mathbf{e}_k^\dagger \mathbf{A}^\dagger \mathbf{e}_j \end{aligned}$$

and, more compactly, $\partial \mathbf{R} / \partial [\mathbf{B}^*]_{kl} = \mathbf{A} \mathbf{B} \mathbf{C} \mathbf{C}^\dagger \mathbf{e}_l \mathbf{e}_k^\dagger \mathbf{A}^\dagger$. Now, we can use the chain rule (noticing that $\mathbf{R}^\dagger = \mathbf{R}$):

$$\begin{aligned} \frac{\partial f}{\partial [\mathbf{B}^*]_{kl}} &= \text{Tr} \left(\left(\frac{\partial f}{\partial \mathbf{R}} \right)^T \frac{\partial \mathbf{R}}{\partial [\mathbf{B}^*]_{kl}} \right) \\ &= \text{Tr} \left(\frac{\partial f}{\partial \mathbf{R}^*} \mathbf{A} \mathbf{B} \mathbf{C} \mathbf{C}^\dagger \mathbf{e}_l \mathbf{e}_k^\dagger \mathbf{A}^\dagger \right) \\ &= \mathbf{e}_k^\dagger \mathbf{A}^\dagger \frac{\partial f}{\partial \mathbf{R}^*} \mathbf{A} \mathbf{B} \mathbf{C} \mathbf{C}^\dagger \mathbf{e}_l, \end{aligned} \quad (97)$$

to obtain the desired result

$$\frac{\partial f}{\partial \mathbf{B}^*} = \mathbf{A}^\dagger \frac{\partial f}{\partial \mathbf{R}^*} \mathbf{A} \mathbf{B} \mathbf{C} \mathbf{C}^\dagger. \quad (98)$$

□

Proof of Theorem 2: First, write the equivalent whitened received signal $\tilde{\mathbf{y}} = \Sigma_n^{-1/2} \mathbf{y}$ (note that both the mutual information and the MMSE remain unchanged by this invertible transformation) as

$$\tilde{\mathbf{y}} = \Sigma_n^{-1/2} \mathbf{H} \mathbf{B} \mathbf{x} + \tilde{\mathbf{n}} \quad (99)$$

where $\tilde{\mathbf{n}} = \Sigma_n^{-1/2} \mathbf{n}$ is the normalized version of \mathbf{n} (assuming that $\Sigma_n = \Sigma_n^{1/2} \Sigma_n^{\dagger/2}$). Now, we can define $\tilde{\mathbf{H}} = \Sigma_n^{-1/2} \mathbf{H} \mathbf{B}$ and invoke Theorem 1 to obtain

$$\nabla_{\tilde{\mathbf{H}}} I = \tilde{\mathbf{H}} \mathbf{E} = \Sigma_n^{-1/2} \mathbf{H} \mathbf{B} \mathbf{E}. \quad (100)$$

By Lemma 3 a), we have $\nabla_{\mathbf{H}} I = \Sigma_n^{-1/2} \nabla_{\tilde{\mathbf{H}}} I \mathbf{B}^\dagger$ and $\nabla_{\mathbf{B}} I = \mathbf{H}^\dagger \Sigma_n^{-1/2} \nabla_{\tilde{\mathbf{H}}} I$, from which (21) and (22) are obtained, respectively.

Also, by Lemma 3 b), we have $\nabla_{\mathbf{B}} f = \nabla_{\mathbf{Q}} f \mathbf{B} \Sigma_x$ and $\nabla_{\mathbf{B}} f = \nabla_{\mathbf{Q}_B} f \mathbf{B}$, from which (23) and (24) are obtained, respectively.

To obtain (25), simply invoke Lemma 3 a) from which $\nabla_{\Sigma_x} I = \mathbf{B}^\dagger \nabla_{\mathbf{Q}} I \mathbf{B}$.

To derive (26), first invoke Lemma 3 a) to obtain

$$\nabla_{\Sigma_n^{-1/2}} I = \nabla_{\tilde{\mathbf{H}}} I \mathbf{B}^\dagger \mathbf{H}^\dagger = \Sigma_n^{-1/2} \mathbf{H} \mathbf{B} \mathbf{E} \mathbf{B}^\dagger \mathbf{H}^\dagger$$

and then invoke Lemma 3 b) (noting that $\Sigma_n = \Sigma_n^{1/2} \Sigma_n^{\dagger/2}$ and $\Sigma_n^{-1} = \Sigma_n^{-\dagger/2} \Sigma_n^{-1/2}$) to obtain

$$\nabla_{\Sigma_n^{-1/2}} I = \Sigma_n^{-1/2} \nabla_{\Sigma_n^{-1}} I$$

which gives $\nabla_{\Sigma_n^{-1}} I = \mathbf{H} \mathbf{B} \mathbf{E} \mathbf{B}^\dagger \mathbf{H}^\dagger$.

Finally, to obtain (27), simply use

$$\partial f / \partial \mathbf{X}^* = -\mathbf{X}^{-1} (\partial f / \partial \mathbf{X}^{-*}) \mathbf{X}^{-1}$$

for a positive definite (Hermitian) matrix \mathbf{X} [15] to get $\nabla_{\Sigma_n} I = -\Sigma_n^{-1} \nabla_{\Sigma_n^{-1}} I \Sigma_n^{-1}$.

D. Gradient of Mutual Information for Arbitrary Σ_x

Define $f_n(\Sigma_x) = I(\Sigma_x + \epsilon_n \mathbf{I})$ where $\epsilon_n \rightarrow 0$ as $n \rightarrow \infty$. Clearly, $f_n(\Sigma_x) \rightarrow I(\Sigma_x)$ and (from Theorem 2)

$$\nabla_{\Sigma_x} f_n(\Sigma_x) = \nabla_{\Sigma_x} I(\Sigma_x + \epsilon_n \mathbf{I}) = \mathbf{H}^\dagger \mathbf{H} \mathbf{E}(\Sigma_x + \epsilon_n \mathbf{I})^{-1} \quad (101)$$

where \mathbf{E} is the MMSE matrix when the transmit covariance matrix is $\Sigma_x + \epsilon_n \mathbf{I}$. To prove (38), it suffices to show that $\nabla_{\Sigma_x} \lim_{n \rightarrow \infty} f_n(\Sigma_x) = \lim_{n \rightarrow \infty} \nabla_{\Sigma_x} f_n(\Sigma_x)$, which holds if $\mathbf{E}(\Sigma_x + \epsilon_n \mathbf{I})^{-1}$ is uniformly continuous on a set containing Σ_x [22]. Since we are only interested in the point Σ_x , we can choose a convenient compact set containing Σ_x so that uniform continuity is implied by pointwise continuity [22]. Therefore, we just need to show that $\mathbf{E}(\Sigma_x + \epsilon_n \mathbf{I})^{-1}$ has a limit.

Indeed, if Σ_x is nonsingular, then $\mathbf{E}(\Sigma_x + \epsilon_n \mathbf{I})^{-1} \rightarrow \mathbf{E} \Sigma_x^{-1}$; if $\Sigma_x = \mathbf{0}$, then $\mathbf{E}(\Sigma_x + \epsilon_n \mathbf{I})^{-1} \rightarrow \mathbf{I}$ (as shown next); and in general we can work with the subspaces of Σ_x associated to null eigenvalues and positive eigenvalues to obtain a limit. To deal with the case $\Sigma_x = \mathbf{0}$, define the normalized vector $\tilde{\mathbf{x}} = \mathbf{x} / \sqrt{\epsilon_n}$, such that its covariance matrix is equal to the identity matrix, and define the MMSE of $\tilde{\mathbf{x}}$ as $\tilde{\mathbf{E}}$, which satisfies $\tilde{\mathbf{E}} = \mathbf{E} / \epsilon_n$ (observe that $\tilde{\mathbf{E}}$ can be interpreted as the MMSE matrix of a communication with noise $\mathbf{n} / \sqrt{\epsilon_n}$). Then, $\mathbf{E}(\epsilon_n \mathbf{I})^{-1} = \tilde{\mathbf{E}}$ and the desired results follows from $\tilde{\mathbf{E}} \rightarrow \mathbf{I}$, since $\mathbb{E}[\tilde{\mathbf{x}} | \mathbf{y}] \rightarrow \mathbb{E}[\tilde{\mathbf{x}}]$.

E. Proof of Theorem 4

Consider the signal model $\mathbf{y} = \mathbf{z} + \mathbf{T}^{1/2} \mathbf{n}$, where $\mathbf{T} = \mathbf{T}^{1/2} \mathbf{T}^{\dagger/2}$ ($\mathbf{T}^{\dagger/2} \triangleq (\mathbf{T}^{1/2})^\dagger$). Equivalently, we can write $\mathbf{y} = \mathbf{z} + \tilde{\mathbf{n}}$, where we have defined the equivalent noise vector as $\tilde{\mathbf{n}} = \mathbf{T}^{1/2} \mathbf{n}$ with covariance matrix $\tilde{\Sigma}_n = \mathbf{T}^{1/2} \Sigma_n \mathbf{T}^{\dagger/2}$. Noting that [13]

$$\begin{aligned} h(\mathbf{z} + \tilde{\mathbf{n}}) &= I(\mathbf{z}; \mathbf{y}) + h(\mathbf{z} + \tilde{\mathbf{n}} | \mathbf{z}) \\ &= I(\mathbf{z}; \mathbf{y}) + \log \det(\pi e \tilde{\Sigma}_n) \end{aligned} \quad (102)$$

we have that

$$\begin{aligned} \frac{\partial}{\partial \tilde{\Sigma}_n} h(\mathbf{z} + \tilde{\mathbf{n}}) &= \frac{\partial}{\partial \tilde{\Sigma}_n} I(\mathbf{z}; \mathbf{y}) + \frac{\partial}{\partial \tilde{\Sigma}_n} \log \det(\pi e \tilde{\Sigma}_n) \\ &= -\tilde{\Sigma}_n^{-1} \mathbf{E} \tilde{\Sigma}_n^{-1} + \tilde{\Sigma}_n^{-1} \\ &= \tilde{\Sigma}_n^{-1} (\mathbf{I} - \mathbf{E} \tilde{\Sigma}_n^{-1}) \end{aligned} \quad (103)$$

where we have used (27). Now, invoking Lemma 3 b) twice we get $\nabla_{\mathbf{T}^{1/2}} f = \nabla_{\tilde{\Sigma}_n} f \mathbf{T}^{1/2} \Sigma_n$ and $\nabla_{\mathbf{T}^{1/2}} f = \nabla_{\mathbf{T}} f \mathbf{T}^{1/2}$, which implies $\nabla_{\tilde{\Sigma}_n} f \tilde{\Sigma}_n = \nabla_{\mathbf{T}} f \mathbf{T}$. Hence, we can finally write

$$\begin{aligned} \nabla_{\mathbf{T}} h(\mathbf{z} + \tilde{\mathbf{n}}) &= \nabla_{\tilde{\Sigma}_n} h(\mathbf{z} + \tilde{\mathbf{n}}) \tilde{\Sigma}_n \mathbf{T}^{-1} \\ &= \tilde{\Sigma}_n^{-1} (\tilde{\Sigma}_n - \mathbf{E}) \mathbf{T}^{-1}. \end{aligned} \quad (104)$$

Now, we compute $\mathbf{J}(\mathbf{y}) = \mathbb{E}_{\mathbf{y}}[\nabla \log p_{\mathbf{y}}(\mathbf{y}) \nabla^\dagger \log p_{\mathbf{y}}(\mathbf{y})]$ as follows (for simplicity of notation and without affecting the result we now consider that the noise is zero mean). First, obtain the inner gradient (from (90)) as

$$\nabla \log p_{\mathbf{y}}(\mathbf{y}) = \frac{1}{p_{\mathbf{y}}(\mathbf{y})} \nabla p_{\mathbf{y}}(\mathbf{y}) = -\tilde{\Sigma}_n^{-1} (\mathbf{y} - \mathbb{E}[\mathbf{z} | \mathbf{y}]) \quad (105)$$

and then

$$\begin{aligned} \mathbf{J}(\mathbf{y}) &= \tilde{\Sigma}_n^{-1} \mathbb{E}[(\mathbf{y} - \mathbb{E}[\mathbf{z} | \mathbf{y}])(\mathbf{y} - \mathbb{E}[\mathbf{z} | \mathbf{y}])^\dagger] \tilde{\Sigma}_n^{-1} \\ &= \tilde{\Sigma}_n^{-1} \mathbb{E}[(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}]) + \tilde{\mathbf{n}})((\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}]) + \tilde{\mathbf{n}})^\dagger] \tilde{\Sigma}_n^{-1} \\ &= \tilde{\Sigma}_n^{-1} \left(\mathbb{E}[(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}])(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}])^\dagger] + \tilde{\Sigma}_n \right. \\ &\quad \left. + \mathbb{E}[(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}])\tilde{\mathbf{n}}^\dagger] + \mathbb{E}[\tilde{\mathbf{n}}(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}])^\dagger] \right) \tilde{\Sigma}_n^{-1} \\ &= \tilde{\Sigma}_n^{-1} (\tilde{\Sigma}_n - \mathbf{E}) \tilde{\Sigma}_n^{-1} \end{aligned} \quad (106)$$

where we have used

$$\begin{aligned} \mathbb{E}[\tilde{\mathbf{n}}(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}])^\dagger] &= \mathbb{E}[(\mathbf{y} - \mathbf{z})(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}])^\dagger] \\ &= -\mathbb{E}[(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}])^\dagger] \\ &= -\mathbb{E}[(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}])(\mathbf{z} - \mathbb{E}[\mathbf{z} | \mathbf{y}])^\dagger] \\ &= -\mathbf{E}. \end{aligned} \quad (107)$$

Combining both results, we can finally write

$$\frac{\partial}{\partial \mathbf{T}^*} h(\mathbf{z} + \mathbf{T}^{1/2} \mathbf{n}) = \mathbf{J}(\mathbf{y}) \tilde{\Sigma}_n \mathbf{T}^{-1} = \mathbf{J}(\mathbf{y}) \mathbf{T}^{1/2} \Sigma_n \mathbf{T}^{-1/2} \quad (108)$$

which is the desired result in (56). Now, (57) is readily obtained particularizing for $\Sigma_n = \mathbf{I}$, which is equivalent to the expression in (55).

F. Proof of Theorem 5

This proof is a refinement of the proof of Theorem 1 in Appendix B, so we will skip many details. We will assume the noise to be zero mean without loss of generality.

The divergence of interest is

$$D(p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0} \| p_{\mathbf{y}}) = -n_R \log(\pi e) - \mathbb{E}[\log p_{\mathbf{y}}(\mathbf{y}) | \mathbf{x} = \mathbf{x}_0]. \quad (109)$$

Then

$$\begin{aligned} \nabla_{\mathbf{H}} D(p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0} \| p_{\mathbf{y}}) &= -\frac{\partial}{\partial \mathbf{H}^*} \int p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) \log p_{\mathbf{y}}(\mathbf{y}) d\mathbf{y} \\ &= -\int \left(\frac{\partial p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y})}{\partial \mathbf{H}^*} \log p_{\mathbf{y}}(\mathbf{y}) \right. \\ &\quad \left. + p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) \frac{\partial \log p_{\mathbf{y}}(\mathbf{y})}{\partial \mathbf{H}^*} \right) d\mathbf{y} \end{aligned} \quad (110)$$

where the interchange of the derivative and the integral follows from Lebesgue Dominated Convergence Theorem and the finiteness of \mathbf{x}_0 .

Now, expanding the derivatives as in Appendix B by using

$$\nabla_{\mathbf{H}} p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) = p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) (\mathbf{y} - \mathbf{H} \mathbf{x}_0) \mathbf{x}_0^\dagger \quad (111)$$

and

$$\nabla_{\mathbf{y}} p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) = -p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) (\mathbf{y} - \mathbf{H} \mathbf{x}_0) \quad (112)$$

we obtain

$$\begin{aligned} \nabla_{\mathbf{H}} D(p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0} \| p_{\mathbf{y}}) &= \int \nabla_{\mathbf{y}} p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) \mathbf{x}_0^\dagger \log p_{\mathbf{y}}(\mathbf{y}) d\mathbf{y} \\ &\quad + \int \frac{p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})} \mathbb{E}_{\mathbf{x}} [\nabla_{\mathbf{y}} p_{\mathbf{y} | \mathbf{x}=\mathbf{x}}(\mathbf{y}) \mathbf{x}^\dagger] d\mathbf{y} \\ &= \int \log p_{\mathbf{y}}(\mathbf{y}) \nabla_{\mathbf{y}} p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) d\mathbf{y} \mathbf{x}_0^\dagger \\ &\quad + \mathbb{E}_{\mathbf{x}} \left[\int \frac{p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})} \nabla_{\mathbf{y}} p_{\mathbf{y} | \mathbf{x}=\mathbf{x}}(\mathbf{y}) d\mathbf{y} \mathbf{x}^\dagger \right] \\ &= -\int p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) \nabla_{\mathbf{y}} \log p_{\mathbf{y}}(\mathbf{y}) d\mathbf{y} \mathbf{x}_0^\dagger \\ &\quad - \mathbb{E}_{\mathbf{x}} \left[\int p_{\mathbf{y} | \mathbf{x}=\mathbf{x}}(\mathbf{y}) \nabla_{\mathbf{y}} \left(\frac{p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})} \right) d\mathbf{y} \mathbf{x}^\dagger \right] \end{aligned} \quad (113)$$

where the last equality follows from integrating by parts [22] as in Appendix B.

Proceeding, we get

$$\begin{aligned} \nabla_{\mathbf{H}} D(p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0} \| p_{\mathbf{y}}) &= -\int \frac{p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})} \nabla_{\mathbf{y}} p_{\mathbf{y}}(\mathbf{y}) \mathbf{x}_0^\dagger d\mathbf{y} \\ &\quad - \int \nabla_{\mathbf{y}} p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) \mathbb{E}_{\mathbf{x}} \left[\frac{p_{\mathbf{y} | \mathbf{x}=\mathbf{x}}(\mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})} \mathbf{x}^\dagger \right] d\mathbf{y} \\ &\quad + \int \frac{\nabla_{\mathbf{y}} p_{\mathbf{y}}(\mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})} p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) \mathbb{E}_{\mathbf{x}} \left[\frac{p_{\mathbf{y} | \mathbf{x}=\mathbf{x}}(\mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})} \mathbf{x}^\dagger \right] d\mathbf{y} \\ &= -\int \frac{\nabla_{\mathbf{y}} p_{\mathbf{y}}(\mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})} p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) \mathbf{x}_0^\dagger d\mathbf{y} \\ &\quad - \int \nabla_{\mathbf{y}} p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) \mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}] d\mathbf{y} \\ &\quad + \int \frac{\nabla_{\mathbf{y}} p_{\mathbf{y}}(\mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})} p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) \mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}] d\mathbf{y} \\ &= -\int \frac{\nabla_{\mathbf{y}} p_{\mathbf{y}}(\mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})} p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) (\mathbf{x}_0 - \mathbb{E}[\mathbf{x} | \mathbf{y}])^\dagger d\mathbf{y} \\ &\quad - \int \nabla_{\mathbf{y}} p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) \mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}] d\mathbf{y} \\ &= \int \mathbb{E}[\mathbf{y} - \mathbf{H} \mathbf{x} | \mathbf{y}] p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) (\mathbf{x}_0 - \mathbb{E}[\mathbf{x} | \mathbf{y}])^\dagger d\mathbf{y} \\ &\quad + \int p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0}(\mathbf{y}) (\mathbf{y} - \mathbf{H} \mathbf{x}_0) \mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}] d\mathbf{y} \end{aligned} \quad (114)$$

where we have used $\nabla_{\mathbf{y}} p_{\mathbf{y}}(\mathbf{y}) = -p_{\mathbf{y}}(\mathbf{y}) \mathbb{E}[\mathbf{y} - \mathbf{H} \mathbf{x} | \mathbf{y}]$ in the last equality.

Finally, expanding the products, we obtain the desired result

$$\begin{aligned} \nabla_{\mathbf{H}} D(p_{\mathbf{y} | \mathbf{x}=\mathbf{x}_0} \| p_{\mathbf{y}}) &= \mathbb{E}[(\mathbf{y} - \mathbf{H} \mathbb{E}[\mathbf{x} | \mathbf{y}])(\mathbf{x}_0 - \mathbb{E}[\mathbf{x} | \mathbf{y}])^\dagger] \\ &\quad + (\mathbf{y} - \mathbf{H} \mathbf{x}_0) \mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}] | \mathbf{x} = \mathbf{x}_0] \\ &= \mathbb{E}[\mathbf{y} \mathbf{x}_0^\dagger + \mathbf{H} \mathbb{E}[\mathbf{x} | \mathbf{y}] \mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}] \\ &\quad - \mathbf{H} \mathbb{E}[\mathbf{x} | \mathbf{y}] \mathbf{x}_0^\dagger - \mathbf{H} \mathbf{x}_0 \mathbb{E}[\mathbf{x}^\dagger | \mathbf{y}] | \mathbf{x} = \mathbf{x}_0] \\ &= \mathbf{H} \mathbb{E}[(\mathbf{x}_0 - \mathbb{E}[\mathbf{x} | \mathbf{y}])(\mathbf{x}_0 - \mathbb{E}[\mathbf{x} | \mathbf{y}])^\dagger | \mathbf{x} = \mathbf{x}_0]. \end{aligned} \quad (115)$$

□

G. Proof of Theorem 6

The proof is a direct consequence of Theorem 1 as we now show. Alternatively, it can also be directly proved without resorting to Theorem 1 (as in [6, Theorem 5]), but the proof becomes tedious and is omitted for lack of additional interest.

Since for an arbitrary $p_{\mathbf{y}'}$ such that $p_{\mathbf{y}} \ll p_{\mathbf{y}'}$

$$\begin{aligned} I(\mathbf{x}; \mathbf{y}) &= D(p_{\mathbf{y}|\mathbf{x}} \| p_{\mathbf{y}} | p_{\mathbf{x}}) \\ &= D(p_{\mathbf{y}|\mathbf{x}} \| p_{\mathbf{y}'|\mathbf{x}} | p_{\mathbf{x}}) - D(p_{\mathbf{y}} \| p_{\mathbf{y}'}) \end{aligned}$$

we can choose $p_{\mathbf{y}'} = p_{\mathbf{y}|\mathbf{x}=\mathbf{x}_0}$ to express

$$\begin{aligned} D(p_{\mathbf{y}} \| p_{\mathbf{y}|\mathbf{x}=\mathbf{x}_0}) &= D(p_{\mathbf{y}|\mathbf{x}} \| p_{\mathbf{y}|\mathbf{x}=\mathbf{x}_0} | p_{\mathbf{x}}) - I(\mathbf{x}; \mathbf{y}) \\ &= \mathbb{E}[(\mathbf{x} - \mathbf{x}_0)^\dagger \mathbf{H}^\dagger \mathbf{H}(\mathbf{x} - \mathbf{x}_0)] - I(\mathbf{x}; \mathbf{y}) \end{aligned} \quad (116)$$

where the second expression follows simply from the explicit evaluation of the divergence between two proper complex Gaussian distributions (with arbitrary units) [16, eq. (59)]

$$\begin{aligned} D(\mathcal{CN}(\mathbf{m}_1, \mathbf{\Sigma}_1), \mathcal{CN}(\mathbf{m}_0, \mathbf{\Sigma}_0)) &= \log \det(\mathbf{\Sigma}_0) - \log \det(\mathbf{\Sigma}_1) \\ &\quad + (\mathbf{m}_1 - \mathbf{m}_0)^\dagger \mathbf{\Sigma}_0^{-1} (\mathbf{m}_1 - \mathbf{m}_0) \log e \\ &\quad + \text{Tr}(\mathbf{\Sigma}_0^{-1} \mathbf{\Sigma}_1 - \mathbf{I}) \log e. \end{aligned} \quad (117)$$

Therefore, invoking Theorem 1, we obtain the desired result

$$\begin{aligned} \nabla_{\mathbf{H}} D(p_{\mathbf{y}} \| p_{\mathbf{y}|\mathbf{x}=\mathbf{x}_0}) &= \nabla_{\mathbf{H}} \mathbb{E}[(\mathbf{x} - \mathbf{x}_0)^\dagger \mathbf{H}^\dagger \mathbf{H}(\mathbf{x} - \mathbf{x}_0)] - \nabla_{\mathbf{H}} I(\mathbf{x}; \mathbf{y}) \\ &= \mathbf{H}(\mathbb{E}[(\mathbf{x} - \mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0)^\dagger] - \mathbf{E}). \end{aligned} \quad (118)$$

The second part of the result follows from the definition of the MMSE matrix \mathbf{E} in (12).

H. Proof of Derivatives for Random \mathbf{H}

For the case of the channel known at the receiver, note that

$$I(\mathbf{x}; \mathbf{C}\mathbf{H}\mathbf{x} + \mathbf{n} | \mathbf{H}) = \mathbb{E}[I(\mathbf{x}; \mathbf{C}\mathbf{H}\mathbf{x} + \mathbf{n} | \mathbf{H} = \mathbf{H})]. \quad (119)$$

To prove (30), it suffices to show that

$$\begin{aligned} \nabla_{\mathbf{C}} \mathbb{E}[I(\mathbf{x}; \mathbf{C}\mathbf{H}\mathbf{x} + \mathbf{n} | \mathbf{H} = \mathbf{H})] &= \mathbb{E}[\nabla_{\mathbf{C}} I(\mathbf{x}; \mathbf{C}\mathbf{H}\mathbf{x} + \mathbf{n} | \mathbf{H} = \mathbf{H})] \end{aligned} \quad (120)$$

and then use (29). The validity of the interchange of the order of the derivative and expectation can be shown by the Lebesgue Dominated Convergence Theorem: simply invoke Lemma 2 observing that

$$\begin{aligned} \left| \left[\mathbf{\Sigma}_n^{-1} \mathbf{C} \mathbf{H} \mathbf{E} \mathbf{H}^\dagger \right]_{ij} \right| &= |\mathbf{a}^\dagger \mathbf{H} \mathbf{E} \mathbf{H}^\dagger \mathbf{b}| \\ &\leq \|\mathbf{a}\| \|\mathbf{b}\| \lambda_{\max}(\mathbf{E}) \lambda_{\max}(\mathbf{H} \mathbf{H}^\dagger) \\ &\leq \|\mathbf{a}\| \|\mathbf{b}\| \text{Tr}(\mathbf{\Sigma}_x) \text{Tr}(\mathbf{H} \mathbf{H}^\dagger) \end{aligned} \quad (121)$$

where $\mathbf{a}^\dagger = \mathbf{e}_i^\dagger \mathbf{\Sigma}_n^{-1} \mathbf{C}$ (\mathbf{e}_i is the canonical i th vector, i.e., the all-zero vector except a one in the i th position), $\mathbf{b} = \mathbf{e}_j$, and we have used the Cauchy–Schwarz inequality in the first inequality

and the fact that the maximum eigenvalue of a matrix is upper-bounded by the trace. Observe now that

$$\mathbb{E}[\text{Tr}(\mathbf{\Sigma}_x) \text{Tr}(\mathbf{H} \mathbf{H}^\dagger)] = \text{Tr}(\mathbf{\Sigma}_x) \text{Tr}(\mathbb{E}[\mathbf{H} \mathbf{H}^\dagger]) < \infty$$

due to the finite second-order moments of \mathbf{H} and \mathbf{x} .

For the case of a channel unknown at the receiver, rewrite the signal model as $\mathbf{y} = \mathbf{C}\mathbf{z} + \mathbf{n}$, where $\mathbf{z} = \mathbf{H}\mathbf{x}$ is a random variable. Since $\mathbf{x} \rightarrow \mathbf{z} \rightarrow \mathbf{y}$ is a Markov chain, then $I(\mathbf{x}; \mathbf{y}) = I(\mathbf{z}; \mathbf{y}) - I(\mathbf{z}; \mathbf{y} | \mathbf{x})$ [13] and using (29), we get (35) provided that

$$\nabla_{\mathbf{C}} \mathbb{E}[I(\mathbf{z}; \mathbf{y} | \mathbf{x} = \mathbf{x})] = \mathbb{E}[\nabla_{\mathbf{C}} I(\mathbf{z}; \mathbf{y} | \mathbf{x} = \mathbf{x})] \quad (122)$$

where the validity of the interchange of the order of the derivative and expectation follows again from the Lebesgue Dominated Convergence Theorem by invoking Lemma 2 and observing that

$$\begin{aligned} &\left| \left[\mathbf{\Sigma}_n^{-1} \mathbf{C} \mathbb{E}[(\mathbf{H}\mathbf{x} - \mathbb{E}[\mathbf{H}\mathbf{x} | \mathbf{y}, \mathbf{x}])(\mathbf{H}\mathbf{x} - \mathbb{E}[\mathbf{H}\mathbf{x} | \mathbf{y}, \mathbf{x}])^\dagger | \mathbf{x}] \right]_{ij} \right| \\ &= |\mathbf{a}^\dagger \mathbb{E}[(\mathbf{H}\mathbf{x} - \mathbb{E}[\mathbf{H}\mathbf{x} | \mathbf{y}, \mathbf{x}])(\mathbf{H}\mathbf{x} - \mathbb{E}[\mathbf{H}\mathbf{x} | \mathbf{y}, \mathbf{x}])^\dagger | \mathbf{x}] \mathbf{b}| \\ &\leq \|\mathbf{a}\| \|\mathbf{b}\| \lambda_{\max}(\mathbb{E}[(\mathbf{H}\mathbf{x} - \mathbb{E}[\mathbf{H}\mathbf{x} | \mathbf{y}, \mathbf{x}])(\mathbf{H}\mathbf{x} \\ &\quad - \mathbb{E}[\mathbf{H}\mathbf{x} | \mathbf{y}, \mathbf{x}])^\dagger | \mathbf{x}]) \\ &\leq \|\mathbf{a}\| \|\mathbf{b}\| \text{Tr}(\mathbb{E}[(\mathbf{H} - \mathbb{E}[\mathbf{H}])\mathbf{x}\mathbf{x}^\dagger(\mathbf{H} - \mathbb{E}[\mathbf{H}])^\dagger]). \end{aligned} \quad (123)$$

Observe now that

$$\begin{aligned} \mathbb{E}[\text{Tr}(\mathbb{E}[(\mathbf{H} - \mathbb{E}[\mathbf{H}])\mathbf{x}\mathbf{x}^\dagger(\mathbf{H} - \mathbb{E}[\mathbf{H}])^\dagger])] &= \text{Tr}(\mathbb{E}[(\mathbf{H} - \mathbb{E}[\mathbf{H}])^\dagger(\mathbf{H} - \mathbb{E}[\mathbf{H}])\mathbb{E}[\mathbf{x}\mathbf{x}^\dagger]]) < \infty \end{aligned} \quad (124)$$

due to the finite second-order moments of \mathbf{H} and \mathbf{x} .

REFERENCES

- [1] K. H. Lee and D. P. Petersen, "Optimal linear coding for vector channels," *IEEE Trans. Commun.*, vol. COM-24, no. 12, pp. 1283–1290, Dec. 1976.
- [2] A. J. Stam, "Some inequalities satisfied by the quantities of information of Fisher and Shannon," *Inf. Contr.*, vol. 2, pp. 101–112, Jun. 1959.
- [3] T. E. Duncan, "On the calculation of mutual information," *SIAM J. Applied Math.*, vol. 19, pp. 215–220, Jul. 1970.
- [4] T. T. Kadota, M. Zakai, and J. Ziv, "Mutual information of the white Gaussian channel with and without feedback," *IEEE Trans. Inf. Theory*, vol. IT-17, no. 4, pp. 368–371, Jul. 1971.
- [5] J. Ziv and M. Zakai, "Some lower bounds on signal parameter estimation," *IEEE Trans. Inf. Theory*, vol. IT-15, no. 3, pp. 386–391, May 1969.
- [6] D. Guo, S. Shamai (Shitz), and S. Verdú, "Mutual information and minimum mean-square error in Gaussian channels," *IEEE Trans. Inf. Theory*, vol. 51, no. 4, pp. 1261–1282, Apr. 2005.
- [7] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*. Englewood Cliffs, NJ: Prentice-Hall, 1993.
- [8] G. Jöngren, M. Skoglund, and B. Ottersen, "Combining beamforming and orthogonal space-time block coding," *IEEE Trans. Inf. Theory*, vol. 48, no. 3, pp. 611–627, Mar. 2002.
- [9] D. P. Palomar, J. M. Cioffi, and M. A. Lagunas, "Joint Tx-Rx beamforming design for multicarrier MIMO channels: A unified framework for convex optimization," *IEEE Trans. Signal Process.*, vol. 51, no. 9, pp. 2381–2401, Sep. 2003.
- [10] J. Salz, "Digital transmission over cross-coupled linear channels," *AT&T Tech. J.*, vol. 64, no. 6, pp. 1147–1159, Jul.–Aug. 1985.
- [11] J. Yang and S. Roy, "On joint transmitter and receiver optimization for multiple-input-multiple-output (MIMO) transmission systems," *IEEE Trans. Commun.*, vol. 42, no. 12, pp. 3221–3231, Dec. 1994.

- [12] A. Scaglione, G. B. Giannakis, and S. Barbarossa, "Redundant filterbank precoders and equalizers. Part I: Unification and optimal designs," *IEEE Trans. Signal Process.*, vol. 47, no. 7, pp. 1988–2006, Jul. 1999.
- [13] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York: Wiley, 1991.
- [14] V. V. Prelov and S. Verdú, "Second-order asymptotics of mutual information," *IEEE Trans. Inf. Theory*, vol. 50, no. 8, pp. 1567–1580, Aug. 2004.
- [15] J. R. Magnus and H. Neudecker, *Matrix Differential Calculus With Applications in Statistics and Econometrics*. New York: Wiley, 1999.
- [16] S. Verdú, "Spectral efficiency in the wideband regime," *IEEE Trans. Inf. Theory*, vol. 48, no. 6, pp. 1319–1343, Jun. 2002.
- [17] N. M. Blachman, "The convolution inequality for entropy powers," *IEEE Trans. Inf. Theory*, vol. IT-11, no. 2, pp. 267–211, Apr. 1965.
- [18] M. H. M. Costa, "A new entropy power inequality," *IEEE Trans. Inf. Theory*, vol. IT-31, no. 6, pp. 751–760, Nov. 1985.
- [19] A. Dembo, T. M. Cover, and J. A. Thomas, "Information theoretic inequalities," *IEEE Trans. Inf. Theory*, vol. 37, no. 6, pp. 1501–1518, Nov. 1991.
- [20] A. Dembo, "Information Inequalities and Uncertainty Principles," Dept. Statistics, Stanford Univ., Stanford, CA, Tech. Rep. 75, Jul. 1990.
- [21] O. Johnson and Y. Suhov, "Entropy and random vectors," *J. Statist. Phys.*, vol. 104, no. 1–2, pp. 145–165, Jul. 2001.
- [22] W. Rudin, *Principles of Mathematical Analysis*, 3rd ed. New York: McGraw-Hill, 1976.
- [23] H. V. Poor, *An Introduction to Signal Detection and Estimation*, 2nd ed. New York: Springer-Verlag, 1994.
- [24] M. S. Pinsker, V. V. Prelov, and S. Verdú, "Sensitivity of channel capacity," *IEEE Trans. Inf. Theory*, vol. 41, no. 6, pp. 1877–1888, Nov. 1995.
- [25] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and Distributed Computation: Numerical Methods*. Englewood Cliffs, NJ: Prentice-Hall, 1989.
- [26] D. P. Bertsekas, *Nonlinear Programming*, 2nd ed. Belmont, MA: Athena Scientific, 1999.
- [27] S. Haykin, *Adaptive Filter Theory*, 4th ed. Upper Saddle River, NJ: Prentice-Hall, 2002.
- [28] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [29] D. H. Brandwood, "A complex gradient operator and its application in adaptive array theory," *Proc. Inst. Elec. Eng.*, vol. 130, no. 1, pp. 11–16, Feb. 1983.
- [30] P. Billingsley, *Probability and Measure*, 3rd ed. New York: Wiley, 1995.