

Capacity-Achieving Probability Measure for Conditionally Gaussian Channels With Bounded Inputs

Terence H. Chan, *Member, IEEE*, Steve Hranilovic, *Member, IEEE*, and Frank R. Kschischang, *Senior Member, IEEE*

Abstract—A conditionally Gaussian channel is a vector channel in which the channel output, given the channel input, has a Gaussian distribution with (well-behaved) input-dependent mean and covariance. We study the capacity-achieving probability measure for conditionally Gaussian channels subject to bounded-input constraints and average cost constraints. Many practical communication systems, including additive Gaussian noise channels, certain optical channels, fading channels, and interference channels fall within this framework. Subject to bounded-input constraint (and average cost constraints), we show that the channel capacity is achievable and we derive a necessary and sufficient condition for a probability measure to be capacity achieving. Under certain conditions, the capacity-achieving measure is proved to be discrete.

Index Terms—Bounded-input constraint, capacity-achieving measure, conditionally Gaussian channel, optical channel, Rayleigh-fading channel.

I. INTRODUCTION

WE STUDY the capacity-achieving probability measure under boundedness constraint for conditionally Gaussian channels, a class of vector channels whose conditional output distribution, given the channel input, is Gaussian with input-dependent mean and, in general, input-dependent covariance. Classical additive white Gaussian noise channels, multiple-input multiple-output (MIMO) Rayleigh-fading channels, certain interference channels, and certain optical channels with signal-dependent noise are examples of communication channel models that fall within this framework.

Of course, determining the capacity of a channel subject to various input constraints is a classical problem of information theory. Shannon demonstrated that in the case of a scalar additive Gaussian noise channel subject to an average power constraint, the capacity-achieving distribution is Gaussian. In practice, however, this source distribution is not realizable due to its

unbounded amplitude. A limit on the peak amplitude, as well as on the average power, is necessary to more accurately reflect the physical limitations present in most practical communication systems.

In [1], Smith studied the capacity of a scalar Gaussian channel subject both to an average power constraint and to a peak power constraint. The capacity-achieving distribution was shown to be discrete, with a finite number of probability mass points. Using a similar approach, Shamai and his colleagues [2]–[5] showed that the capacity-achieving measures are also discrete for many other channels, including Poisson channels, quadrature Gaussian channels, Rayleigh-fading channels and so on. In particular, since the publication of [4], there has been a wide interest in this area, and many channels have been found to have discrete capacity-achieving measures. In Table I, we list some channel models that we are aware of having been shown to have discrete capacity-achieving measures.

In this paper, we extend the work of Smith and Shamai to study the capacity-achieving measure for conditionally Gaussian channels. We have organized the paper as follows. In Section II, we present the channel model and we specify the channel input constraints. In Section III, we provide the main results of this paper. Subject to a bounded-input constraint and average cost constraints, we show that the channel capacity is achievable, and derive a necessary and sufficient condition for a probability measure to be capacity achieving. Using this necessary and sufficient condition, we propose an algorithm to find the capacity-achieving measure of a signal-dependent optical channel, which is traditionally difficult to analyze. Using an approach similar to that of [1], we prove that, under suitable criteria, the capacity-achieving measure is discrete. In the special case when the conditionally Gaussian channel is constant, we further prove that the capacity-achieving measures approach a Gaussian distribution if the bounded-input constraint is relaxed. The proofs for our main results rely on a collection of intermediate propositions and lemmas, which are stated and proved in Appendix B. In Section IV, we apply our results to analyze several practical channels, including certain optical channels, fading channels, interference channels, and parallel Gaussian channels. Finally, we provide some conclusions in Section V.

II. SYSTEM MODEL

A. Notation

In this paper, we adopt the following notation. For positive integers N , M , and K , let $\mathcal{N} = \{1, \dots, N\}$, $\mathcal{M} = \{1, \dots, M\}$,

Manuscript received September 17, 2003; revised December 7, 2004. The work of T. H. Chan was supported in part by a Croucher Foundation Fellowship. The material in this paper was presented at the IEEE International Symposium on Information Theory, Chicago, IL, June/July 2004.

T. H. Chan is with the Department of Computer Science, University of Regina, Regina, SK S4S 0A2, Canada (e-mail: terence@cs.uregina.ca).

S. Hranilovic is with the Department of Electrical and Computer Engineering, McMaster University, Hamilton, ON L8S 4K1, Canada (e-mail: hranilovic@mcmaster.ca).

F. R. Kschischang is with the Department of Electrical and Computer Engineering, University of Toronto, Toronto, ON M5S 3G4, Canada (e-mail: frank@comm.utoronto.ca).

Communicated by A. Lapidoth, Associate Editor for Shannon Theory.
Digital Object Identifier 10.1109/TIT.2005.847707

TABLE I
A LIST OF CHANNELS WITH DISCRETE CAPACITY-ACHIEVING MEASURES

| Channel model | Constraints | References |
|--|---------------------------------------|------------------------------------|
| Scalar additive Gaussian | PP , AP ¹ | Smith [1] 1971 |
| Additive noise with piecewise constant noise density | PP | Oettli [6] 1974 |
| Poisson | PP , AP | Shamai [2] 1990 |
| Quadrature Gaussian | PP , AP | Shamai <i>et al.</i> [3] 1995 |
| Scalar additive heavy-tailed noise | AP | Das [7] 2000 |
| Rayleigh fading | AP | Abou-Faycal <i>et al.</i> [4] 2001 |
| Non-coherent scalar AWGN | AP | Katz <i>et al.</i> [5] 2002 |
| Additive vector Gaussian | PP , AP | Palanki [8] 2002 |
| Single antenna Rayleigh block fading | AP | Palanki [8] 2002 |
| General scalar fading | higher-moment constraint ² | Palanki [8] 2002 |
| Non-coherent Rician fading | AP | Gursoy <i>et al.</i> [9] 2003 |
| Non-coherent block independent AWGN | AP | Nuriyev <i>et al.</i> [10] 2003 |
| Partially coherent AWGN | AP | Hou <i>et al.</i> [11] 2003 |
| UIUO ³ | PP | Huang <i>et al.</i> [12] 2003 |
| General scalar additive | PP | Tchamkerten [13] 2004 |

¹ **PP** and **AP** stand for “peak power” and “average power” respectively.

² A higher moment constraint is an average cost constraint whose cost function is in the form: $x^\alpha - \gamma$ where $\alpha > 2$.

³ This is a name we gave to the channel described in [12]. Here, UIUO stands for “Unbounded-Input Unbounded-Output.” A UIUO channel is a scalar channel such that for any $c > 0$, $\lim_{|x| \rightarrow \infty} \Pr(|Y| < c | X = x) = 0$ where X and Y are channel input and output, respectively. In other words, if the channel input goes to infinity, the probability that the channel output is contained in a given bounded subset of \mathbb{R} goes to zero.

and $K = \{1, \dots, K\}$ be index sets. The transpose, the determinant and the trace of a matrix \mathbf{A} are denoted by \mathbf{A}^T , $\det \mathbf{A}$, and $\text{tr}(\mathbf{A})$, respectively. An identity matrix of size n is denoted by \mathbf{I}_n . The real part and imaginary part of a complex matrix \mathbf{A} are denoted by $\text{re}(\mathbf{A})$ and $\text{im}(\mathbf{A})$, respectively. Let \mathbf{j} be the imaginary number $\sqrt{-1}$. Then for any complex matrix \mathbf{A} , we have $\mathbf{A} = \text{re}(\mathbf{A}) + \mathbf{j} \text{im}(\mathbf{A})$.

For any $\mathbf{x} = [x_1, \dots, x_N]^T \in \mathbb{R}^N$, its amplitude $\|\mathbf{x}\|$ is defined as $\sqrt{\mathbf{x}^T \mathbf{x}}$ and we say its **phase**¹ is the normalized vector $\mathbf{x}/\|\mathbf{x}\|$. Let $\mathbb{B}_N(s)$ be the ball in \mathbb{R}^N of radius \sqrt{s} , i.e.,

$$\mathbb{B}_N(s) = \{\mathbf{x} \in \mathbb{R}^N : \mathbf{x}^T \mathbf{x} \leq s\}.$$

Then, the phase of a vector \mathbf{x} is defined on the “surface” of the unit ball $\mathbb{B}_N(1)$. Clearly, \mathbf{x} is uniquely determined by its

¹To follow the similar terminology defined in previous work, the way we use the term “phase” in this paper is slightly different from the traditional usage in which it refers to the special case when $N = 2$ and the normalized vector is represented by an angle.

amplitude and its phase. If $N = 1$, then \mathbf{x} is a scalar and is simply denoted by x .

Let \mathbf{X} be a random variable defined on \mathbb{R}^N and μ be its probability measure. Then μ is called **discrete in amplitude and uniform in phase (DAUP)** [11] if the probability distribution of $\|\mathbf{X}\|$ is discrete with a finite number of probability mass points, and the phase of \mathbf{X} is uniformly distributed on the surface of $\mathbb{B}_N(1)$.

B. Channel Model

Consider a discrete-time memoryless channel whose input is a real-valued N -tuple $\mathbf{X} = [X_1, \dots, X_N]^T$ subject to input constraints to be defined and whose output is a real-valued M -tuple $\mathbf{Y} = [Y_1, \dots, Y_M]^T$. In a practical communication system, due to various physical restrictions imposed on the system, not every channel input $\mathbf{x} \in \mathbb{R}^N$ can be generated by the transmitter. For example, all amplifiers have a limitation on their maximum allowable input and output amplitudes. A channel input \mathbf{x} is said to be **admissible** if it can be produced by the transmitter. We will denote the set of admissible channel input vectors by \mathbb{S} , and we will assume that \mathbb{S} is a closed and bounded subset of \mathbb{R}^N , unless specified explicitly. We refer to this channel input constraint as the **bounded-input constraint**. If the probability measure of \mathbf{X} is μ , then the bounded-input constraint means that $\mu(\mathbb{S}) = 1$. For example, under a peak total power constraint in N dimensions of s^2 , we have $\mathbb{S} = \mathbb{B}_N(s^2)$.

We call a channel **conditionally Gaussian (CG)** if the conditional probability distribution of the channel output is Gaussian distributed. In particular, for a given input $\mathbf{x} \in \mathbb{S}$, the conditional expectation vector and the conditional covariance matrix of the channel output \mathbf{Y} are denoted by $\mathbf{b}_\mathbf{x}$ and $\Delta_\mathbf{x}$. Thus, the conditional probability density function $P_{Y|X}(\mathbf{y}|\mathbf{x})$ of \mathbf{Y} given $\mathbf{X} = \mathbf{x}$ is given by

$$P_{Y|X}(\mathbf{y}|\mathbf{x}) = \frac{\exp(-\frac{1}{2}(\mathbf{y} - \mathbf{b}_\mathbf{x})^T \nabla_\mathbf{x}(\mathbf{y} - \mathbf{b}_\mathbf{x}))}{\sqrt{(2\pi)^M \det \Delta_\mathbf{x}}} \quad (1)$$

where $\nabla_\mathbf{x}$ is the matrix inverse of $\Delta_\mathbf{x}$.

In this paper, the following assumptions are made: 1) for all $\mathbf{x} \in \mathbb{S}$, the covariance matrix $\Delta_\mathbf{x}$ is positive-definite (hence, $\det \Delta_\mathbf{x} > 0$), and 2) the entries in $\mathbf{b}_\mathbf{x}$ and $\Delta_\mathbf{x}$ are **well behaved**, in the sense that they can be extended holomorphically over \mathbb{C}^N and are real over \mathbb{R}^N . For clarity, when the functions are extended over \mathbb{C}^N , a complex variable \mathbf{w} is used instead of \mathbf{x} to denote the function argument. Since \mathbb{S} is assumed to be closed and bounded, the continuity of $\mathbf{b}_\mathbf{x}$ and $\Delta_\mathbf{x}$ implies that there exist $\kappa_l, \kappa_h, \vartheta > 0$ such that for all $\mathbf{x} \in \mathbb{S}$, all eigenvalues of $\nabla_\mathbf{x}$ are within the interval (κ_l, κ_h) and $\|\mathbf{b}_\mathbf{x}\| < \vartheta$. Hence, $1/\kappa_h^M \leq \det \Delta_\mathbf{x} \leq 1/\kappa_l^M$.

We also introduce cost constraints, as some channel inputs require more “effort” by the transmitter to generate than others. In this paper, we consider a set of K cost measures $g_k(\mathbf{x})$ indexed by K . The k th average cost constraint is characterized by the inequality $G_k(\mu) \leq 0$, where $G_k(\mu) \triangleq \mathbf{E}_\mu [g_k(\mathbf{X})]$ and the expectation is taken with respect to the input probability measure μ . For instance, the commonly used constraint that the average signal power should not exceed γ is expressed

as $\mathbf{E}_\mu [\|\mathbf{X}\|^2 - \gamma] \leq 0$. Again, we assume that the cost functions $g_k(\mathbf{x})$ are all well behaved.

For simplicity, we use the tuple $(\Delta_{\mathbf{x}}, \mathbf{b}_{\mathbf{x}}, \{g_k(\mathbf{x}) : k \in \mathcal{K}\}, \mathbb{S})$ to denote a CG channel subject to the above bounded-input constraint and average cost constraints.

We call a CG channel **quadratic** if its input and output are related by the equation $\mathbf{Y} = \mathbf{H}\mathbf{X} + \mathbf{Z}$, where 1) \mathbf{H} is a random $M \times N$ channel matrix whose entries are Gaussian distributed, and 2) $\mathbf{Z} = [Z_1, Z_2, \dots, Z_M]^\top$ is an additive Gaussian noise that is independent of \mathbf{H} and \mathbf{X} . For the special case that \mathbf{H} is deterministic, we call the channel **constant quadratic**, or simply **constant**.

Clearly, for a quadratic CG channel, $\mathbf{b}_{\mathbf{x}} = \mathbf{E}[\mathbf{H}]\mathbf{x} + \mathbf{E}[\mathbf{Z}]$, and

$$\Delta_{\mathbf{x}} = \mathbf{E}[\mathbf{H}\mathbf{x}\mathbf{x}^\top \mathbf{H}^\top] - \mathbf{E}[\mathbf{H}]\mathbf{x}\mathbf{x}^\top \mathbf{E}[\mathbf{H}^\top] + \mathbf{E}[\mathbf{Z}\mathbf{Z}^\top] - \mathbf{E}[\mathbf{Z}]\mathbf{E}[\mathbf{Z}^\top].$$

If the channel is constant, then $\mathbf{b}_{\mathbf{x}} = \mathbf{H}\mathbf{x} + \mathbf{E}[\mathbf{Z}]$ and

$$\Delta_{\mathbf{x}} = \mathbf{E}[\mathbf{Z}\mathbf{Z}^\top] - \mathbf{E}[\mathbf{Z}]\mathbf{E}[\mathbf{Z}^\top].$$

Thus, $\Delta_{\mathbf{x}}$ is constant for all \mathbf{x} and will simply be denoted by Δ .

The following are some communication system channel models that fall within our framework.

1) *A Scalar Additive White Gaussian Channel*: Subject to peak and average power constraints, the channel is characterized by the tuple

$$(\sigma^2, x, \{x^2 - \gamma\}, [-s, s]).$$

The discreteness of the capacity-achieving measure for this channel was first shown in [1]. It is obvious that a scalar additive Gaussian channel is a constant CG channel whose channel ‘‘matrix’’ \mathbf{H} is equal to 1.

2) *A System of N Parallel Independent Gaussian Channels*: Subject to peak and average total power constraints, the channel is characterized by the tuple

$$(\sigma^2 \mathbf{I}_N, \mathbf{x}, \{\mathbf{x}^\top \mathbf{x} - \gamma\}, \mathbb{B}_N(s^2)).$$

This is a generalization of the scalar-Gaussian channel. The special case when $N = 2$ was studied in [3], and it was shown that the capacity-achieving measure is DAUP.

3) *An Optical Intensity Modulated/Direct Detected Channel (Using Raised Quadrature Amplitude Modulation (QAM) Basis Functions [14])*: Under the assumption that the noise is dominated by the background illumination, subject to peak and average power constraints, the channel is a constant CG channel and is characterized by the tuple

$$(\sigma^2 \mathbf{I}_3, \mathbf{x}, \{x_1 - \gamma\}, \mathbb{S})$$

where

$$\mathbb{S} = \left\{ \mathbf{x} \in \mathbb{R}^3 : \sqrt{2x_2^2 + 2x_3^2} \leq x_1 \leq s - \sqrt{2x_2^2 + 2x_3^2} \right\}. \quad (2)$$

Here, for any input \mathbf{x} , the optical power is given by x_1 . See [14] for more details concerning this channel model.

4) *A Signal-Dependent Noise Optical Pulse Amplitude Modulation (PAM) Channel*: Subject to peak and average power constraints, the channel is characterized by the tuple

$$(\sigma_1^2 x + \sigma_0^2, x, \{x - \gamma\}, [0, s]).$$

This example will be elaborated in detail in Section IV-A.

5) *A MIMO Rayleigh-Fading Channel*: Assuming that there are $N/2$ and $M/2$ transmitting and receiving antennae, subject to peak and average power constraints, the channel is characterized by the tuple

$$\left(\frac{1}{2} (\chi^2 \mathbf{x}^\top \mathbf{x} + \sigma^2) \mathbf{I}_M, \mathbf{0}, \{\mathbf{x}^\top \mathbf{x} - \gamma\}, \mathbb{B}_N(s^2) \right).$$

It was proved in [4] that the capacity-achieving measure for this channel is discrete in the sense that the distribution of $\|\mathbf{X}\|$ is discrete with a finite number of probability mass points. Furthermore, the discreteness of the capacity-achieving measure remains valid even when $s = \infty$, i.e., $\mathbb{S} = \mathbb{R}^N$.

6) *A MIMO Rician Fading Channel*: Subject to a peak power constraint, the channel is characterized by the tuple

$$\left(\frac{1}{2} (\chi^2 \mathbf{x}^\top \mathbf{x} + \sigma^2) \mathbf{I}_M, \mathbf{E}[\mathbf{H}]\mathbf{x}, \{\}, \mathbb{B}_N(s^2) \right).$$

7) *A MIMO Rayleigh Fading Channel With Receiver-Side Channel Information*: Assuming that there are $N/2$ transmitting and one receiving antennae, subject to peak power constraint, the channel is characterized by the tuple

$$(\Delta_{\mathbf{x}}, \mathbf{0}, \{\}, \mathbb{B}_N(s^2))$$

where $\mathbf{x}^\top = [\mathbf{x}_{\text{re}}^\top, \mathbf{x}_{\text{im}}^\top]$ and $\Delta_{\mathbf{x}}$ is defined as

$$\frac{1}{2} \begin{bmatrix} \chi^2 \mathbf{I}_N & \mathbf{0} & \chi^2 \mathbf{x}_{\text{re}} & \chi^2 \mathbf{x}_{\text{im}} \\ \mathbf{0} & \chi^2 \mathbf{I}_N & -\chi^2 \mathbf{x}_{\text{im}} & \chi^2 \mathbf{x}_{\text{re}} \\ \chi^2 \mathbf{x}_{\text{re}}^\top & -\chi^2 \mathbf{x}_{\text{im}}^\top & \chi^2 \mathbf{x}^\top \mathbf{x} + \sigma^2 & 0 \\ \chi^2 \mathbf{x}_{\text{im}}^\top & \chi^2 \mathbf{x}_{\text{re}}^\top & 0 & \chi^2 \mathbf{x}^\top \mathbf{x} + \sigma^2 \end{bmatrix}. \quad (3)$$

8) *An Interference Channel*: Subject to peak power constraint, the channel is characterized by the tuple

$$((\chi^2 \mathbf{x}^\top \mathbf{x} + \sigma^2) \mathbf{I}_N, \mathbf{x}, \{\}, \mathbb{B}_N(s^2)).$$

This channel is a special case of MIMO Rician fading channel where $\mathbf{E}[\mathbf{H}]$ is equal to \mathbf{I}_N . We will elaborate more in Section IV-C.

III. MAIN RESULTS

Consider a CG channel characterized by the tuple $(\Delta_{\mathbf{x}}, \mathbf{b}_{\mathbf{x}}, \{g_k(\mathbf{x}) : k \in \mathcal{K}\}, \mathbb{S})$. We are interested in the probability measure which satisfies the channel input constraints and maximizes the mutual information between the channel input and output. Let $\Lambda_{\mathbb{S}}$ be the set of all input probability measures μ satisfying the bounded-input constraint, and Λ be the subset of $\Lambda_{\mathbb{S}}$ which also satisfies the average cost constraints (if any). For any input probability measure μ of \mathbf{X} , we denote the μ -induced probability density function of \mathbf{Y} by $P_{\mathbf{Y}}(\mathbf{y}; \mu)$, the differential entropy of \mathbf{Y} by $H_{\mathbf{Y}}(\mu)$, the conditional differential entropy of

\mathbf{Y} given \mathbf{X} by $H_{\mathbf{Y}|\mathbf{X}}(\mu)$, and the mutual information between \mathbf{X} and \mathbf{Y} by $I(\mu)$. Thus,

$$\begin{aligned} P_{\mathbf{Y}}(\mathbf{y}; \mu) &= \int P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}) d\mu \\ H_{\mathbf{Y}}(\mu) &= - \int P_{\mathbf{Y}}(\mathbf{y}; \mu) \log P_{\mathbf{Y}}(\mathbf{y}; \mu) d\mathbf{y} \\ H_{\mathbf{Y}|\mathbf{X}}(\mu) &= \frac{1}{2} \int \log[(2\pi e)^M \det \Delta_{\mathbf{x}}] d\mu \\ I(\mu) &= H_{\mathbf{Y}}(\mu) - H_{\mathbf{Y}|\mathbf{X}}(\mu). \end{aligned} \quad (4)$$

As a result, the **channel capacity problem** is the following optimization problem:

$$\begin{cases} \text{Maximize} & I(\mu) \\ \text{subject to} & \mu \in \Lambda_{\mathcal{S}}, \text{ and for all } k \in \mathbf{K}, G_k(\mu) \leq 0. \end{cases} \quad (5)$$

Our first result shows that a capacity-achieving measure exists.

Theorem 1 (Existence and Uniqueness): There exists μ_o in Λ which solves the channel capacity problem. Furthermore, the capacity-achieving output distribution is unique. In other words, if ν is also a capacity-achieving input measure, then $P_{\mathbf{Y}}(\mathbf{y}; \mu_o) = P_{\mathbf{Y}}(\mathbf{y}; \nu)$.

Suppose, in addition, that the CG channel is constant and the channel matrix \mathbf{H} is left-invertible. Then the capacity-achieving input measure is unique.

Sketch of Proof: The proof of Theorem 1 follows the same approach used in [1]. First, we will prove that Λ is convex and sequentially compact² (see Proposition 2). As the mutual information function $I(\mu)$ is continuous over Λ (see Proposition 3), it achieves its maximum in Λ .

For any $0 < \theta < 1$ and $\mu_o, \mu_1 \in \Lambda_{\mathcal{S}}$, let $\mu_{\theta} = (1 - \theta)\mu_o + \theta\mu_1$. It is well known that $I(\mu)$ and $H_{\mathbf{Y}}(\mu)$ are concave, i.e., $I(\mu_{\theta}) \geq (1 - \theta)I(\mu_o) + \theta I(\mu_1)$. By the strict concavity of the function $f(c) = -c \log c$, we can show that equality holds if and only if $P_{\mathbf{Y}}(\mathbf{y}; \mu_o) = P_{\mathbf{Y}}(\mathbf{y}; \mu_1)$.

Furthermore, if the CG channel is constant with a left-invertible channel matrix \mathbf{H} , then by Lemma 4, there is a one-to-one correspondence between the input and output probability measures, i.e., $P_{\mathbf{Y}}(\mathbf{y}; \mu_o) = P_{\mathbf{Y}}(\mathbf{y}; \mu_1)$ if and only if $\mu_o = \mu_1$. As a result, $I(\mu)$ is a strictly concave function of μ . The uniqueness of the capacity-achieving input measure then follows. \square

Definition 1: A point $\mathbf{x} \in \mathbb{R}^N$ is said to be a **point of increase** of μ if for any open subset \mathcal{O} of \mathbb{R}^N containing \mathbf{x} , $\mu(\mathcal{O}) > 0$.

Let \mathbb{E}_{μ} be the set of points of increase of μ . Then $\mu(\mathbb{E}_{\mu}) = 1$. In fact, \mathbb{E}_{μ} is the minimal closed subset of \mathbb{R}^N whose probability is 1.

Theorem 2 (Necessity and Sufficiency): Let μ_o be an admissible input probability measure, i.e., $\mu_o \in \Lambda$. Let $C(\mathbf{x})$ be $\frac{1}{2} \log[(2\pi e)^M \det \Delta_{\mathbf{x}}]$. Then μ_o is capacity achieving if and only if there exists $\{\lambda_k \geq 0 : k \in \mathbf{K}\}$ such that for all $\mathbf{x} \in \mathcal{S}$

$$Q(\mathbf{x}; \mu_o) - C(\mathbf{x}) - I(\mu_o) - \sum_{k \in \mathbf{K}} \lambda_k g_k(\mathbf{x}) \leq 0 \quad (6)$$

²A space is sequentially compact if every infinite sequence in the space has a convergent subsequence.

where $Q(\mathbf{x}; \mu_o)$ is defined as $-\int P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_o) d\mathbf{y}$. Furthermore, if \mathbf{x} belongs to \mathbb{E}_{μ_o} , the inequality (6) is satisfied with equality.

Sketch of Proof: The proof of this theorem also follows a similar approach as in [1]. By the method of Lagrange multipliers [16], μ_o is capacity achieving if and only if 1) there exists $\{\lambda_k \geq 0 : k \in \mathbf{K}\}$ such that $\lambda_k G_k(\mu_o) = 0$ for all $k \in \mathbf{K}$ and 2) for all $\mu \in \Lambda_{\mathcal{S}}$, $J(\mu_o) \geq J(\mu)$ where $J(\mu) = I(\mu) - \sum_{k \in \mathbf{K}} \lambda_k G_k(\mu)$. Since the function $J(\mu)$ is concave, the condition 2) is equivalent to that $J'_{\mu_o}(\mu) \leq 0$ for all $\mu \in \Lambda_{\mathcal{S}}$ where $J'_{\mu_o}(\mu)$ is the weak derivative of $J(\mu)$ at μ_o (see Definition 6). Since $G_k(\mu)$ is a linear function of μ , we have

$$J'_{\mu_o}(\mu) = I'_{\mu_o}(\mu) - \sum_{k \in \mathbf{K}} \lambda_k (G_k(\mu) - G_k(\mu_o)).$$

Therefore, given that the condition 1) is true, the condition 2) is equivalent to

$$I'_{\mu_o}(\mu) - \sum_{k \in \mathbf{K}} \lambda_k G_k(\mu) \leq 0 \quad (7)$$

for all $\mu \in \Lambda_{\mathcal{S}}$. In Proposition 4, we prove that

$$I'_{\mu_o}(\mu) = \int Q(\mathbf{x}; \mu_o) d\mu - H_{\mathbf{Y}|\mathbf{X}}(\mu) - I(\mu_o).$$

For simplicity, the function $a(\mathbf{x}, \mu_o, \{\lambda_k : k \in \mathbf{K}\})$ is used to denote the following formula:

$$Q(\mathbf{x}; \mu_o) - \frac{1}{2} \log[(2\pi e)^M \det \Delta_{\mathbf{x}}] - I(\mu_o) - \sum_{k \in \mathbf{K}} \lambda_k g_k(\mathbf{x}). \quad (8)$$

(Proof of Necessity) Suppose μ_o is capacity achieving. Let $\{\lambda_k \geq 0 : k \in \mathbf{K}\}$ be selected such that the inequality (7) is satisfied. For any $\mathbf{x}^* \in \mathcal{S}$, let μ be the probability measure such that $\mu(\{\mathbf{x}^*\}) = 1$, i.e., μ is the measure with a single probability mass at \mathbf{x}^* . Obviously, $\mu \in \Lambda_{\mathcal{S}}$. Substituting μ into (7), we have $a(\mathbf{x}^*, \mu_o, \{\lambda_k : k \in \mathbf{K}\}) \leq 0$. The necessity is thus proved.

(Proof of Sufficiency) Suppose the inequality (6) is satisfied for selected $\{\lambda_k \geq 0 : k \in \mathbf{K}\}$. Since $\int Q(\mathbf{x}; \mu_o) d\mu_o = H_{\mathbf{Y}}(\mu_o)$ and $\int C(\mathbf{x}) d\mu_o = H_{\mathbf{Y}|\mathbf{X}}(\mu_o)$, if we integrate both sides of (6) with respect to μ_o , we have

$$\begin{aligned} 0 &\geq \int Q(\mathbf{x}; \mu_o) d\mu_o - \int C(\mathbf{x}) d\mu_o - I(\mu_o) \\ &\quad - \sum_{k \in \mathbf{K}} \lambda_k \int g_k(\mathbf{x}) d\mu_o \\ &= - \sum_{k \in \mathbf{K}} \lambda_k G_k(\mu_o) \\ &\geq 0 \end{aligned}$$

where the last inequality follows from that $G_k(\mu_o) \leq 0$. Hence, $\lambda_k G_k(\mu_o) = 0$ for all $k \in \mathbf{K}$. On the other hand, if we integrate (6) with respect to any $\mu \in \Lambda_{\mathcal{S}}$, we have

$$\begin{aligned} 0 &\geq \int Q(\mathbf{x}; \mu_o) d\mu - \int C(\mathbf{x}) d\mu - I(\mu_o) \\ &\quad - \sum_{k \in \mathbf{K}} \lambda_k \int g_k(\mathbf{x}) d\mu \\ &= I'_{\mu_o}(\mu) - \sum_{k \in \mathbf{K}} \lambda_k G_k(\mu). \end{aligned} \quad (9)$$

$$= I'_{\mu_o}(\mu) - \sum_{k \in \mathbf{K}} \lambda_k G_k(\mu). \quad (10)$$

Hence, μ_o is capacity-achieving.

Finally, it remains to prove that the inequality (6) is satisfied with equality at \mathbb{E}_{μ_o} . Let $\mathbf{x}^* \in \mathbb{E}_{\mu_o}$, and suppose, to the contrary, that $a(\mathbf{x}^*, \mu_o, \{\lambda_k : k \in \mathbb{K}\}) = -\epsilon < 0$. By the continuity of $Q(\mathbf{x}; \mu_o)$ (see Lemma 5), the function $a(\mathbf{x}, \mu_o, \{\lambda_k : k \in \mathbb{K}\})$ is also continuous on \mathbb{S} . Hence, there is an open subset \mathbb{O} containing \mathbf{x}^* such that $a(\mathbf{x}, \mu_o, \{\lambda_k : k \in \mathbb{K}\}) \leq -\epsilon/2$ for all $\mathbf{x} \in \mathbb{O}$. Consequently, if we integrate both sides of (6) with respect to μ_o , we have $-\epsilon\mu_o(\mathbb{O})/2 \geq -\sum_{k \in \mathbb{K}} \lambda_k G_k(\mu_o) = 0$. A contradiction occurs, and the theorem is proved. \square

By Theorem 1, if μ_o and ν are both capacity achieving, then $P_{\mathbf{Y}}(\mathbf{y}; \mu_o) = P_{\mathbf{Y}}(\mathbf{y}; \nu)$, which further implies that $Q(\mathbf{x}; \mu_o) = Q(\mathbf{x}; \nu)$. Hence,

$$a(\mathbf{x}, \mu_o, \{\lambda_k : k \in \mathbb{K}\}) = a(\mathbf{x}, \nu, \{\lambda_k : k \in \mathbb{K}\})$$

for all $\mathbf{x} \in \mathbb{S}$. According to Theorem 2, the inequality (6) is thus tight at any points of increase of a capacity-achieving measure.

Definition 2: A complex-valued function $f(\mathbf{w})$ defined on an open subset \mathbb{U} of \mathbb{C}^N is called **holomorphic** [17] on \mathbb{U} if it is analytic in each individual variable on \mathbb{U} .

Definition 3: A subset \mathbb{F} of \mathbb{R}^N is defined to be **sparse** (in \mathbb{R}^N) if there exists a nonzero holomorphic function f defined on a connected open subset \mathbb{U} of \mathbb{C}^N containing the closure of \mathbb{F} such that $f(\mathbf{w}) = 0$ for all $\mathbf{w} \in \mathbb{F}$.

Suppose \mathbb{F} is not sparse and $f(\mathbf{w})$ is holomorphic on a connected open subset \mathbb{U} of \mathbb{C}^N containing the closure of \mathbb{F} . If f is zero on \mathbb{F} , then it is also zero on \mathbb{U} .

By the identity theorem, a bounded subset \mathbb{F} of \mathbb{R} is sparse in \mathbb{R} if and only if it is finite. However, for a high-dimensional space \mathbb{R}^N , it is difficult in general to determine whether a bounded subset of \mathbb{R}^N is sparse or not. In the following lemma, we show that when \mathbb{F} is a collection of ‘‘concentric shells’’ and is bounded, it is sparse if and only if its number of shells is finite.

Lemma 1: Suppose $\{p_i : i \in \mathbb{I}\}$ is a set of distinct non-negative real numbers which are bounded above by a positive real number. Let \mathbb{F} be a collection of ‘‘concentric shells’’ of radii $\{p_i : i \in \mathbb{I}\}$, i.e.,

$$\mathbb{F} = \bigcup_{i \in \mathbb{I}} \{\mathbf{x} \in \mathbb{R}^N : \mathbf{x}^\top \mathbf{x} = p_i^2\}.$$

Then \mathbb{F} is sparse if and only if \mathbb{I} is finite.

Proof: See Appendix A. \square

Definition 4: A probability measure μ is said to be **discrete** if its set of points of increase \mathbb{E}_μ is sparse.

Definition 5: Let $\Delta_{\mathbf{w}}$ be a holomorphic extension of $\Delta_{\mathbf{x}}$ to \mathbb{C}^N . The **well-behaved region** of $\Delta_{\mathbf{w}}$, denoted by $\mathbb{A}(\Delta_{\mathbf{w}})$, is the subset of \mathbb{C}^N such that for all its elements \mathbf{w} , $\text{re}(\det \Delta_{\mathbf{w}}) > 0$ and $\text{re}(\nabla_{\mathbf{w}})$ is positive definite.

Note that $\text{re}(\det \Delta_{\mathbf{w}}) > 0$ implies that the logarithm and the square root of $\det \Delta_{\mathbf{w}}$ are well defined and holomorphic over $\mathbb{A}(\Delta_{\mathbf{w}})$. Furthermore, as $\text{re}(\nabla_{\mathbf{w}})$ is a symmetric matrix, the condition that $\text{re}(\nabla_{\mathbf{w}})$ is positive definite is equivalent to the condition that the eigenvalues of $\text{re}(\nabla_{\mathbf{w}})$ are positive. As the eigenvalues of $\text{re}(\nabla_{\mathbf{w}})$ are continuous functions of \mathbf{w} , the

well-behaved region $\mathbb{A}(\Delta_{\mathbf{w}})$ is an open subset in \mathbb{C}^N . In the special case when the channel is constant, $\Delta_{\mathbf{w}}$ is real, constant, and positive definite for all $\mathbf{w} \in \mathbb{C}^N$. Hence, $\mathbb{A}(\Delta) = \mathbb{C}^N$.

Theorem 3 (Discreteness): Suppose μ_o is capacity achieving and $\{\lambda_k \geq 0 : k \in \mathbb{K}\}$ are chosen such that the inequality (6) is satisfied. Let \mathbb{E} be the subset of \mathbb{S} at which the inequality (6) is tight. In the following three cases, the set \mathbb{E} is sparse, and hence, μ_o is discrete.

[Case A] There is no average cost constraint (i.e., \mathbb{K} is empty), and there exists a connected open subset \mathbb{W} of $\mathbb{A}(\Delta_{\mathbf{w}})$ containing \mathbb{S} and a sequence $\{\mathbf{w}^{(i)}\}_{i=1}^\infty$ in \mathbb{W} such that 1) $\mathbf{b}_{\mathbf{w}^{(i)}}$ and $\Delta_{\mathbf{w}^{(i)}}$ are real for all positive integers i , and 2) $\lim_{i \rightarrow \infty} \text{tr}(\Delta_{\mathbf{w}^{(i)}}) = +\infty$.

[Case B] There exists a connected open subset \mathbb{W} of $\mathbb{A}(\Delta_{\mathbf{w}})$ containing \mathbb{S} and a convergent sequence $\{\mathbf{w}^{(i)}\}_{i=1}^\infty$ in \mathbb{W} with a limit $\mathbf{w}^{(0)}$ (not necessary in \mathbb{W}), such that 1) $\mathbf{b}_{\mathbf{w}^{(i)}}$ and $\Delta_{\mathbf{w}^{(i)}}$ are real for all positive integers i , and 2) $\lim_{i \rightarrow \infty} \det \Delta_{\mathbf{w}^{(i)}} = 0$.

[Case C] The channel is constant with a nonzero channel matrix \mathbf{H} , and the cost functions are of second order, i.e., $g_k(\mathbf{x}) \triangleq \mathbf{x}^\top \mathbf{Q}_k \mathbf{x} + \mathbf{L}_k \mathbf{x} - \gamma_k$ for some $N \times N$ real matrices \mathbf{Q}_k , $1 \times N$ real row vectors \mathbf{L}_k and real scalars γ_k .

Sketch of Proof: Let μ_o be the capacity-achieving measure for the channel $(\Delta_{\mathbf{x}}, \mathbf{b}_{\mathbf{x}}, \{g_k(\mathbf{x}) : k \in \mathbb{K}\}, \mathbb{S})$. For any $\mathbf{w} \in \mathbb{A}(\Delta_{\mathbf{w}})$, define $d(\mathbf{w})$ as

$$- \int \frac{\exp(-\frac{1}{2}(\mathbf{y} - \mathbf{b}_{\mathbf{w}})^\top \nabla_{\mathbf{w}}(\mathbf{y} - \mathbf{b}_{\mathbf{w}}))}{\sqrt{(2\pi)^M \det \Delta_{\mathbf{w}}}} \log P_{\mathbf{Y}}(\mathbf{y}; \mu_o) d\mathbf{y}. \quad (11)$$

It is clear from the definition that for all $\mathbf{x} \in \mathbb{S}$, $d(\mathbf{x}) = Q(\mathbf{x}; \mu_o)$. Let

$$a(\mathbf{w}) \triangleq d(\mathbf{w}) - \frac{1}{2} \log(2\pi e)^M \det \Delta_{\mathbf{w}} - I(\mu_o) - \sum_{k \in \mathbb{K}} \lambda_k g_k(\mathbf{w}).$$

Then, by definition, $a(\mathbf{x}) = 0$ for all $\mathbf{x} \in \mathbb{E}$.

In the above three cases, there exists a connected open subset \mathbb{W} of $\mathbb{A}(\Delta_{\mathbf{w}})$ containing \mathbb{S} . In particular, \mathbb{W} is equal to \mathbb{C}^N in Case C. Since \mathbb{S} is a closed and bounded subset containing \mathbb{E} , \mathbb{W} also contains the closure of \mathbb{E} . As $d(\mathbf{w})$ is holomorphic on $\mathbb{A}(\Delta_{\mathbf{w}})$ (see Proposition 5), $a(\mathbf{w})$ is also holomorphic on $\mathbb{A}(\Delta_{\mathbf{w}})$. Hence, if \mathbb{E} is not sparse, then for all $\mathbf{w} \in \mathbb{W}$, $a(\mathbf{w}) = 0$, or equivalently

$$d(\mathbf{w}) = \frac{1}{2} \log[(2\pi e)^M \det \Delta_{\mathbf{w}}] + I(\mu_o) + \sum_{k \in \mathbb{K}} \lambda_k g_k(\mathbf{w}).$$

In the following, we will show that if \mathbb{E} is not sparse, then a contradiction occurs in each of the three cases and hence Theorem 3 is proved.

Contradiction in Case A:

Suppose the requirements in Case A are satisfied, and \mathbb{E} is not sparse. By Lemma 6, there exist $\alpha, \beta > 0$ such that $d(\mathbf{w}^{(i)}) \geq -\alpha + \beta \text{tr}[\Delta_{\mathbf{w}^{(i)}}]$. Without loss of generality, assume that $\xi_1(\mathbf{w}^{(i)}) \geq \dots \geq \xi_M(\mathbf{w}^{(i)}) > 0$ are the eigenvalues of $\Delta_{\mathbf{w}^{(i)}}$. Then

$$\text{tr}[\Delta_{\mathbf{w}^{(i)}}] = \sum_{m \in \mathbb{M}} \xi_m(\mathbf{w}^{(i)})$$

and

$$\det \Delta_{\mathbf{w}^{(i)}} = \prod_{m \in \mathbb{M}} \xi_m(\mathbf{w}^{(i)}).$$

As a result

$$\begin{aligned} \frac{1}{2} \log(2\pi e)^M + \frac{1}{2} \sum_{m \in \mathbb{M}} \log \xi_m(\mathbf{w}^{(i)}) + I(\mu_o) \\ \geq -\alpha + \beta \sum_{m \in \mathbb{M}} \xi_m(\mathbf{w}^{(i)}) \end{aligned}$$

and hence,

$$\frac{1}{2} \log(2\pi e)^M + I(\mu_o) + \alpha \geq \beta \xi_1(\mathbf{w}^{(i)}) - \frac{M}{2} \log \xi_1(\mathbf{w}^{(i)}).$$

Since

$$\lim_{i \rightarrow \infty} \overline{\text{tr}}(\Delta_{\mathbf{w}^{(i)}}) = +\infty$$

$\xi_1(\mathbf{w}^{(i)})$ and $\beta \xi_1(\mathbf{w}^{(i)}) - \frac{M}{2} \log \xi_1(\mathbf{w}^{(i)})$ are not bounded above. A contradiction thus occurs.

Contradiction in Case B:

Suppose the requirements in Case B are satisfied, and \mathbb{E} is not sparse. Again, by Lemma 6, there exist $\alpha, \beta > 0$ such that $d(\mathbf{w}^{(i)}) \geq -\alpha + \beta \text{tr}[\Delta_{\mathbf{w}^{(i)}}]$. Therefore,

$$\begin{aligned} \frac{1}{2} \log[(2\pi e)^M \det \Delta_{\mathbf{w}^{(i)}}] \\ \geq -\alpha + \beta \text{tr}[\Delta_{\mathbf{w}^{(i)}}] - I(\mu_o) - \sum_{k \in \mathbb{K}} \lambda_k g_k(\mathbf{w}^{(i)}) \end{aligned}$$

and consequently

$$\begin{aligned} \lim_{i \rightarrow \infty} \frac{1}{2} \log[(2\pi e)^M \det \Delta_{\mathbf{w}^{(i)}}] \\ \geq -\alpha + \beta \text{tr}[\Delta_{\mathbf{w}^{(0)}}] - I(\mu_o) - \sum_{k \in \mathbb{K}} \lambda_k g_k(\mathbf{w}^{(0)}). \end{aligned}$$

Since $\lim_{i \rightarrow \infty} \log \det \Delta_{\mathbf{w}^{(i)}} = -\infty$, a contradiction thus occurs.

Contradiction in Case C:

Suppose the requirements in Case C are satisfied, and \mathbb{E} is not sparse. Then, for all $\mathbf{x} \in \mathbb{R}^N$, we have

$$d(\mathbf{x}) = H_Y(\mu_o) + \sum_{k \in \mathbb{K}} \lambda_k (\mathbf{x}^\top \mathbf{Q}_k \mathbf{x} + \mathbf{L}_k \mathbf{x} - \gamma_k). \quad (12)$$

Following the same approach as in [1], we will show that the equality (12) implies that the capacity-achieving output distribution is Gaussian. Since $\mathbf{Y} = \mathbf{H}\mathbf{X} + \mathbf{Z}$ and \mathbf{Z} is Gaussian distributed, $\mathbf{H}\mathbf{X}$ is also Gaussian distributed. As the set of points of increase of a Gaussian distribution is not bounded, the Gaussianity of $\mathbf{H}\mathbf{X}$ contradicts the assumption that μ_o satisfies the bounded-input constraint.

Now, it remains to prove that \mathbf{Y} is Gaussian distributed. First, assume that the rank of \mathbf{H} is equal to M . Without loss of generality, we can also assume that \mathbf{Z} is zero mean. Let $P_{\mathbf{Z}}(\mathbf{z})$ be the probability density function of the additive Gaussian noise. Then, it is clear that for $\mathbf{x} \in \mathbb{R}^N$

$$d(\mathbf{x}) = - \int P_{\mathbf{Z}}(\mathbf{H}\mathbf{x} - \mathbf{y}) \log P_Y(\mathbf{y}; \mu_o) d\mathbf{Y} \quad (13)$$

$$= - P_{\mathbf{Z}}(\mathbf{H}\mathbf{x}) * \log P_Y(\mathbf{H}\mathbf{x}; \mu_o) \quad (14)$$

where $*$ denotes convolution. Let $\mathbf{r} = \mathbf{H}\mathbf{x}$. The right-hand side of (14) depends only on \mathbf{r} . Hence,

$$H_Y(\mu_o) + \sum_{k \in \mathbb{K}} \lambda_k (\mathbf{x}^\top \mathbf{Q}_k \mathbf{x} + \mathbf{L}_k \mathbf{x} - \gamma_k)$$

also depends on \mathbf{r} only and we can thus construct matrices $\hat{\mathbf{Q}}_k$ and $\hat{\mathbf{L}}_k$ such that

$$\begin{aligned} H_Y(\mu_o) + \sum_{k \in \mathbb{K}} \lambda_k (\mathbf{x}^\top \mathbf{Q}_k \mathbf{x} + \mathbf{L}_k \mathbf{x} - \gamma_k) \\ = H_Y(\mu_o) + \sum_{k \in \mathbb{K}} \lambda_k (\mathbf{r}^\top \hat{\mathbf{Q}}_k \mathbf{r} + \hat{\mathbf{L}}_k \mathbf{r} - \gamma_k). \end{aligned}$$

Thus, we have

$$\begin{aligned} -P_{\mathbf{Z}}(\mathbf{r}) * \log P_Y(\mathbf{r}; \mu_o) \\ = H_Y(\mu_o) + \sum_{k \in \mathbb{K}} \lambda_k (\mathbf{r}^\top \hat{\mathbf{Q}}_k \mathbf{r} + \hat{\mathbf{L}}_k \mathbf{r} - \gamma_k). \quad (15) \end{aligned}$$

By direct integration, we can immediately verify that $\hat{\mathbf{L}}_k \mathbf{r} = P_{\mathbf{Z}}(\mathbf{r}) * \hat{\mathbf{L}}_k \mathbf{r}$ and

$$\mathbf{r}^\top \hat{\mathbf{Q}}_k \mathbf{r} = P_{\mathbf{Z}}(\mathbf{r}) * \left(\mathbf{r}^\top \hat{\mathbf{Q}}_k \mathbf{r} - \sum_{i,j \in \mathbb{M}} \hat{\mathbf{Q}}_k^{(i,j)} \Delta^{(i,j)} \right)$$

where $\hat{\mathbf{Q}}_k^{(i,j)}$ and $\Delta^{(i,j)}$ are the (i,j) th entries of $\hat{\mathbf{Q}}_k$ and Δ , respectively. Hence, (15) implies that for all $\mathbf{r} \in \mathbb{R}^M$, we have

$$P_{\mathbf{Z}}(\mathbf{r}) * \left(\log P_Y(\mathbf{r}; \mu_o) + H_Y(\mu_o) + \sum_{k \in \mathbb{K}} \lambda_k \left(\mathbf{r}^\top \hat{\mathbf{Q}}_k \mathbf{r} + \hat{\mathbf{L}}_k \mathbf{r} - \gamma_k - \sum_{i,j \in \mathbb{M}} \hat{\mathbf{Q}}_k^{(i,j)} \Delta^{(i,j)} \right) \right) = \mathbf{0}.$$

As the Fourier transform of $P_{\mathbf{Z}}$ is nonzero everywhere, by Corollary 9, the function

$$\begin{aligned} \log P_Y(\mathbf{y}; \mu_o) + H_Y(\mu_o) \\ + \sum_{k \in \mathbb{K}} \lambda_k \left(\mathbf{y}^\top \hat{\mathbf{Q}}_k \mathbf{y} - \sum_{i,j \in \mathbb{M}} \hat{\mathbf{Q}}_k^{(i,j)} \Delta^{(i,j)} + \hat{\mathbf{L}}_k \mathbf{y} - \gamma_k \right) \end{aligned}$$

is equal to 0, which further implies that \mathbf{Y} is Gaussian distributed.

Now, assume that the rank of \mathbf{H} is equal to $M' < M$. As Δ is real, symmetric, and positive definite, there exists an invertible matrix \mathbf{B} such that $\mathbf{B}\Delta\mathbf{B}^\top = \mathbf{I}_M$. Using singular-value decomposition on $\mathbf{B}\mathbf{H}$, we can construct an $M' \times M$ matrix \mathbf{V} such that 1) $I(\mathbf{X}; \mathbf{H}\mathbf{X} + \mathbf{Z}) = I(\mathbf{X}; \mathbf{V}\mathbf{B}\mathbf{H}\mathbf{X} + \mathbf{V}\mathbf{B}\mathbf{Z})$, and 2) the rank of $\mathbf{V}\mathbf{B}\mathbf{H}$ is M' . Hence, the original constant CG channel is equivalent to another constant CG channel where the channel matrix is $\mathbf{V}\mathbf{B}\mathbf{H}$ and the additive noise is $\mathbf{V}\mathbf{B}\mathbf{Z}$. Using the same argument as above, we can conclude that if the requirements in Case C are satisfied, then \mathbb{E} must be sparse. \square

As a remark, subject to an average power constraint ($\mathbf{E}_\mu[\mathbf{x}^\top \mathbf{x} - \gamma] \leq 0$), the capacity-achieving measure for a constant CG channel whose channel gain matrix is an identity matrix has also been studied separately in [8], in which it was proved that the capacity-achieving measure is discrete if $\mathbb{R}^N \setminus \mathbb{S}$ has a positive Lebesgue measure.

Corollary 1: Let $(\Delta_{\mathbf{x}}, \mathbf{b}_{\mathbf{x}}, \{\}, \mathbb{S})$ be a quadratic CG channel subject only to a bounded-input constraint. The ca-

capacity-achieving measure is discrete except in the trivial case that the channel is constant with a zero channel matrix \mathbf{H} .

Proof: Let \mathbf{H} be the channel matrix for the quadratic CG channel. Then

$$\Delta_{\mathbf{x}} = \mathbf{E}[\mathbf{H}\mathbf{x}\mathbf{x}^T\mathbf{H}^T] - \mathbf{E}[\mathbf{H}]\mathbf{x}\mathbf{x}^T\mathbf{E}[\mathbf{H}^T] + \mathbf{E}[\mathbf{Z}\mathbf{Z}^T] - \mathbf{E}[\mathbf{Z}]\mathbf{E}[\mathbf{Z}^T].$$

Clearly, $\Delta_{\mathbf{x}}$ and $\mathbf{b}_{\mathbf{x}}$ are real over \mathbb{R}^N , and \mathbb{R}^N is a subset of $\mathbb{A}(\Delta_{\mathbf{w}})$.

Suppose the quadratic CG channel is not constant. Let \mathbb{W} be an open subset of $\mathbb{A}(\Delta_{\mathbf{w}})$ containing \mathbb{R}^N . Then $\text{tr}[\Delta_{\mathbf{x}}]$ is unbounded over \mathbb{W} . Hence, the requirements in Case A of Theorem 3 are satisfied, and thus the capacity-achieving measure is discrete.

If the quadratic CG channel is constant, then the requirements in Case C of Theorem 3 are also satisfied. Hence, the capacity-achieving measure is still discrete. \square

It is well known that for a scalar additive Gaussian channel subject only to an average power constraint, the capacity-achieving input distribution is Gaussian. When an additional bounded-input constraint $\mathbb{S} = [-s, s]$ is imposed, the capacity-achieving measure becomes discrete. An interesting question then arises: if we relax the bounded-input constraint by increasing s , then what is its effect on the capacity-achieving measure? In the next theorem, we will prove that as s increases, the capacity-achieving measure approaches the capacity-achieving distribution when there is no bounded-input constraint.

Consider a sequence of CG channels

$$\{(\Delta, \mathbf{H}\mathbf{x}, \{\mathbf{x}^T\mathbf{x} - \gamma\}, \mathbb{S}_j)\}_{j=1}^{\infty}$$

where \mathbf{H} is a nonzero real matrix and \mathbb{S}_j is a closed and bounded subset of \mathbb{R}^N . Let μ_j and Φ be the capacity-achieving measures for the channels $(\Delta, \mathbf{H}\mathbf{x}, \{\mathbf{x}^T\mathbf{x} - \gamma\}, \mathbb{S}_j)$ and $(\Delta, \mathbf{H}\mathbf{x}, \{\mathbf{x}^T\mathbf{x} - \gamma\}, \mathbb{R}^N)$, respectively. Note that μ_j is discrete while Φ is Gaussian distributed.

Theorem 4: Suppose for any closed and bounded set \mathbb{V} of \mathbb{R}^N , we have $\mathbb{V} \subseteq \mathbb{S}_j$ for all j sufficiently large. Then $\{\mu_j\}_{j=1}^{\infty}$ converges to Φ in Lévy metric.

Sketch of Proof: First, consider the special case when \mathbf{H} is an $N \times N$ diagonal matrix and $\Delta = \mathbf{I}_N$. For clarity, we use \mathbf{D} to denote \mathbf{H} and χ_n to denote its diagonal entries $\mathbf{D}^{(n,n)}$ for $n \in \mathbb{N}$. Using the water-filling algorithm [18], [19], there exist $P_1, \dots, P_N \geq 0$ and $\vartheta > 0$ such that 1) $\sum_{n \in \mathbb{N}} P_n = \gamma$, 2) $P_n = \max(0, \vartheta - 1/\chi_n^2)$, and 3) for all $n \in \mathbb{N}$, we have $\chi_n^2/(\chi_n^2 P_n + 1) \leq 1/\vartheta$ and $P_n \chi_n^2/(\chi_n^2 P_n + 1) = P_n/\vartheta$. Let Φ be the Gaussian input distribution of \mathbf{X} such that the components of \mathbf{X} are zero mean and independent with variances P_1, \dots, P_N , respectively. Then Φ is capacity achieving for the channel $(\mathbf{I}_N, \mathbf{D}\mathbf{x}, \{\mathbf{x}^T\mathbf{x} - \gamma\}, \mathbb{R}^N)$.

By Proposition 6, we show that $\lim_{j \rightarrow \infty} I(\mu_j) = I(\Phi)$, or equivalently, $\lim_{j \rightarrow \infty} H_{\mathbf{Y}}(\mu_j) = H_{\mathbf{Y}}(\Phi)$. Let Γ be the set of all input probability measures μ of \mathbf{X} such that $\mathbf{E}_{\mu}[\mathbf{X}^T\mathbf{X}] \leq \gamma$. In Lemma 10, we prove that Γ is sequentially compact. Hence, $\{\mu_j\}_{j=1}^{\infty}$ converges to Φ if all convergent subsequences of $\{\mu_j\}_{j=1}^{\infty}$ converge to Φ . Assume without loss of generality that

$\{\mu_j\}_{j=1}^{\infty}$ converges to an input probability measure μ_0 in Γ . As $\lim_{j \rightarrow \infty} H_{\mathbf{Y}}(\mu_j) = H_{\mathbf{Y}}(\Phi)$, by Lemma 9, we have

$$\lim_{j \rightarrow \infty} \int |P_{\mathbf{Y}}(\mathbf{y}; \mu_j) - P_{\mathbf{Y}}(\mathbf{y}; \Phi)| d\mathbf{y} = 0.$$

On the other hand, as $\{\mu_j\}_{j=1}^{\infty}$ converges to μ_0

$$\lim_{j \rightarrow \infty} P_{\mathbf{Y}}(\mathbf{y}; \mu_j) = P_{\mathbf{Y}}(\mathbf{y}; \mu_0), \quad \text{for all } \mathbf{y} \in \mathbb{R}^M.$$

Therefore, for any $L > 0$

$$\begin{aligned} & \int_{\|\mathbf{y}\| \leq L} |P_{\mathbf{Y}}(\mathbf{y}; \mu_0) - P_{\mathbf{Y}}(\mathbf{y}; \Phi)| d\mathbf{y} \\ & \leq \lim_{j \rightarrow \infty} \left(\int |P_{\mathbf{Y}}(\mathbf{y}; \mu_j) - P_{\mathbf{Y}}(\mathbf{y}; \Phi)| d\mathbf{y} \right. \\ & \quad \left. + \int_{\|\mathbf{y}\| \leq L} |P_{\mathbf{Y}}(\mathbf{y}; \mu_j) - P_{\mathbf{Y}}(\mathbf{y}; \mu_0)| d\mathbf{y} \right) \\ & = 0. \end{aligned}$$

As L is arbitrary, by the continuity of $P_{\mathbf{Y}}(\mathbf{y}; \mu_0)$ and $P_{\mathbf{Y}}(\mathbf{y}; \Phi)$, we have $P_{\mathbf{Y}}(\mathbf{y}; \mu_0) = P_{\mathbf{Y}}(\mathbf{y}; \Phi)$ for all $\mathbf{y} \in \mathbb{R}^N$. Hence, $\mathbf{E}_{\mu_0}[Y_n^2] = \mathbf{E}_{\Phi}[Y_n^2]$ for all $n \in \mathbb{N}$. Let \mathbb{L} be the subset $\{n \in \mathbb{N} : \chi_n \neq 0\}$. As $\mathbf{E}_{\mu_0}[Y_n^2] = \chi_n^2 \mathbf{E}_{\mu_0}[X_n^2] + 1$ and $\mathbf{E}_{\Phi}[Y_n^2] = \chi_n^2 \mathbf{E}_{\Phi}[X_n^2] + 1$, we have $\mathbf{E}_{\Phi}[X_n^2] = \mathbf{E}_{\mu_0}[X_n^2]$ for all $n \in \mathbb{L}$. Consequently

$$\gamma = \sum_{n \in \mathbb{L}} \mathbf{E}_{\Phi}[X_n^2] = \sum_{n \in \mathbb{L}} \mathbf{E}_{\mu_0}[X_n^2]$$

which implies that $\mathbf{E}_{\mu_0}[X_n^2] = 0$ for $n \notin \mathbb{L}$. In other words, with respect to μ_0 and Φ , X_n is deterministic and equal to zero for $n \notin \mathbb{L}$.

Let $P_{\mathbf{D}\mathbf{X}}(\mu_0)$ and $P_{\mathbf{D}\mathbf{X}}(\Phi)$ be the probability measures of $\mathbf{D}\mathbf{X}$ when the input probability measures of \mathbf{X} are μ_0 and Φ , respectively. Since $P_{\mathbf{Y}}(\mathbf{y}; \mu_0) = P_{\mathbf{Y}}(\mathbf{y}; \Phi)$, by Lemma 4, $P_{\mathbf{D}\mathbf{X}}(\mu_0) = P_{\mathbf{D}\mathbf{X}}(\Phi)$. As X_n is deterministic with respect to μ_0 and Φ for those n such that $\chi_n = 0$, $P_{\mathbf{D}\mathbf{X}}(\mu_0) = P_{\mathbf{D}\mathbf{X}}(\Phi)$ implies that $\mu_0 = \Phi$. Hence, Theorem 4 holds for this simple case.

Now, consider a general channel $(\Delta, \mathbf{H}\mathbf{x}, \{\mathbf{x}^T\mathbf{x} - \gamma\}, \mathbb{S})$. Let M' be the rank of \mathbf{H} . Using singular-value decomposition, we can construct an $N \times N$ orthonormal matrix \mathbf{U} , an $M' \times N$ diagonal matrix \mathbf{D}^* , an $M' \times M$ matrix \mathbf{V} , and an $M \times M$ matrix \mathbf{B} such that 1) $I(\mathbf{X}; \mathbf{H}\mathbf{X} + \mathbf{Z}) = I(\mathbf{X}; \mathbf{V}\mathbf{B}\mathbf{H}\mathbf{X} + \mathbf{V}\mathbf{B}\mathbf{Z})$, 2) the rank of $\mathbf{V}\mathbf{B}\mathbf{H}$ is M' , 3) $\mathbf{V}\mathbf{B}\mathbf{H} = \mathbf{D}^*\mathbf{U}$, and 4) $\mathbf{E}[\mathbf{V}\mathbf{B}\mathbf{Z}\mathbf{Z}^T\mathbf{B}^T\mathbf{V}^T] = \mathbf{I}_{M'}$. Then the channel is equivalent to $\mathbf{Y}^* = \mathbf{D}^*\mathbf{X}^* + \mathbf{Z}^*$ where $\mathbf{X}^* = \mathbf{U}\mathbf{X}$, $\mathbf{Z}^* = \mathbf{V}\mathbf{B}\mathbf{Z}$, and is denoted by the tuple $(\mathbf{I}_{M'}, \mathbf{D}^*\mathbf{x}^*, \{\|\mathbf{x}^*\|^2 - \gamma\}, \mathbb{S}^*)$ where $\mathbb{S}^* = \{\mathbf{U}\mathbf{x} : \mathbf{x} \in \mathbb{S}\}$.

By adding $N - M'$ “dummy” output variables, the channel $(\mathbf{I}_{M'}, \mathbf{D}^*\mathbf{x}^*, \{\|\mathbf{x}^*\|^2 - \gamma\}, \mathbb{S}^*)$ is equivalent to $(\mathbf{I}_N, \hat{\mathbf{D}}\mathbf{x}^*, \{\|\mathbf{x}^*\|^2 - \gamma\}, \mathbb{S}^*)$, where $\hat{\mathbf{D}}$ is obtained by appending $N - M'$ rows of zeros in \mathbf{D}^* . It is obvious that μ_j and Φ are still capacity-achieving measures for the channels $(\mathbf{I}_N, \hat{\mathbf{D}}\mathbf{x}^*, \{\|\mathbf{x}^*\|^2 - \gamma\}, \mathbb{S}_j^*)$ and $(\mathbf{I}_N, \hat{\mathbf{D}}\mathbf{x}^*, \{\|\mathbf{x}^*\|^2 - \gamma\}, \mathbb{R}^N)$ where $\mathbb{S}_j^* = \{\mathbf{U}\mathbf{x} : \mathbf{x} \in \mathbb{S}_j\}$. Since $\hat{\mathbf{D}}$ is a diagonal square matrix, using the same argument as before, we prove that $\{\mu_j\}_{j=1}^{\infty}$ converges to Φ . \square

IV. EXAMPLES

A. Optical Channel—PAM

Consider an optical channel using intensity modulation and direct detection, in which an electrical signal is directly con-

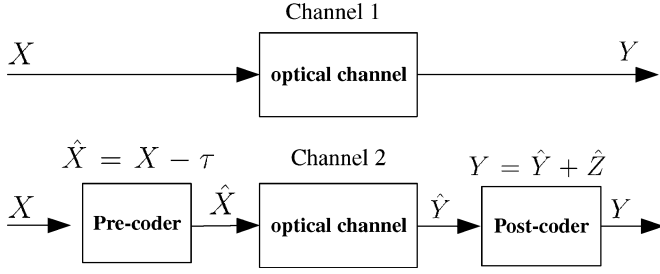


Fig. 1. Two optical channels.

verted to an optical intensity signal (e.g., by a laser diode). A discrete-time model for this channel is $Y = X + Z$ where Z is a zero-mean additive Gaussian noise, whose power may depend on the channel input. Specifically, we assume that $\Delta_x = \sigma_0^2 + \sigma_1^2 x$ for $\sigma_0^2 > 0$ and $\sigma_1^2 \geq 0$. This system is subject to three channel input constraints: a) the average power constraint, b) the peak power constraint, and c) the nonnegativity constraint (due to the nonnegativity of intensity). Such a channel can be characterized by the tuple $(\sigma_0^2 + \sigma_1^2 x, x, \{x - \gamma\}, [0, s])$.

By Theorem 1, there exists a probability measure μ_o in Λ which solves the channel capacity problem. In fact, by Theorem 2, there exists $\lambda \geq 0$ such that for all $x \in [0, s]$

$$Q(x; \mu_o) - \frac{1}{2} \log[2\pi e(\sigma_0^2 + \sigma_1^2 x)] - I(\mu_o) - \lambda(x - \gamma) \leq 0. \quad (16)$$

Furthermore, the inequality (16) is satisfied with equality for all $x \in \mathbb{E}_{\mu_o}$.

Corollary 2: Let μ_o be the capacity-achieving measure for the channel $(\sigma_0^2 + \sigma_1^2 x, x, \{x - \gamma\}, [0, s])$. Then μ_o is a discrete distribution with a finite number of probability mass points.

Proof: Suppose $\sigma_1^2 > 0$. Since \mathbb{S} is connected, it is contained in one of the maximal connected open subsets of $\mathbb{A}(\Delta_w)$. Denote this connected open subset by \mathbb{W} . As $\Delta_w = \sigma_0^2 + \sigma_1^2 w$ and $\mathbf{b}_w = w$, the interval $(-\sigma_0^2/\sigma_1^2, +\infty)$ is a subset of \mathbb{W} . Let $\{w^{(i)}\}_{i=1}^\infty$ be a convergent sequence in the interval $(-\sigma_0^2/\sigma_1^2, +\infty)$ with limit $-\sigma_0^2/\sigma_1^2$. Then $\lim_{i \rightarrow \infty} \det \Delta_{w^{(i)}} = 0$ and for all positive integers i , $\mathbf{b}_{w^{(i)}}$ and $\Delta_{w^{(i)}}$ are real. Clearly, the requirements in Case B of Theorem 3 are satisfied, and thus, μ_o is discrete.

In the case $\sigma_1^2 = 0$, the channel is constant and is subject to a linear average cost constraint. Hence, the requirements in Case C of Theorem 3 are satisfied, and the capacity-achieving probability measure μ_o is thus discrete.

In both cases, \mathbb{E}_{μ_o} , the set of points of increase of μ_o , is a sparse set. Since \mathbb{E}_{μ_o} is a bounded subset of \mathbb{R} , it is a finite set, i.e., μ_o has a finite number of probability mass points. \square

Proposition 1: Suppose μ_o is the capacity-achieving measure for the channel $(\sigma_0^2 + \sigma_1^2 x, x, \{x - \gamma\}, [0, s])$. Then $x = 0$ is a point of increase of μ_o .

Proof: Suppose, to the contrary, that $x = 0$ is not a point of increase of μ_o . Let τ be the minimal element in \mathbb{E}_{μ_o} , and hence $\tau > 0$. Consider the two channels depicted in Fig. 1. Channel 1 is the original optical channel, and channel 2 is obtained from channel 1 by appending a “pre-coder” and a “post-coder” before and after the inner optical channel. Specifically, $\hat{X} = X - \tau$ and

$Y = \hat{Y} + \hat{Z}$ where \hat{Z} is an independent additive Gaussian noise \hat{Z} with mean τ and covariance $\tau\sigma_1^2$.

For any $x \geq \tau$, the conditional probability density functions of Y given $X = x$ is the same in both channels. Hence, if the input probability measure of the two channels is μ_o , then the joint probability measures of X and Y in the two channels are still the same, and consequently, so is the mutual information between the channel input and output.

In the second channel, as X , \hat{X} , \hat{Y} , and Y form a Markov chain $X \rightarrow \hat{X} \rightarrow \hat{Y} \rightarrow Y$, we have $I(\hat{X}; \hat{Y}) \geq I(X; Y)$ by data processing inequality. Let ν be the corresponding probability measure of \hat{X} when the probability measure of X is μ_o . Clearly, ν satisfies the bounded-input constraint and the average power constraint. Thus, ν is also capacity-achieving for channel 1, and hence, $P_Y(y; \mu_o) = P_Y(y; \nu)$ by Theorem 1. As a result, for channel 2, given the input probability measure of X is μ_o , the probability density function for Y and \hat{Y} are the same, which is not possible since $\mathbf{E}[Y] = \mathbf{E}[\hat{Y}] + \tau$. Hence, the proposition follows. \square

Corollary 3: If μ_o is capacity achieving for the channel $(\sigma_0^2 + \sigma_1^2 x, x, \{x - \gamma\}, [0, s])$, and it satisfies (16), then

$$\lambda = \left[I(\mu_o) - Q(0; \mu_o) + \frac{1}{2} \log 2\pi e \sigma_0^2 \right] / \gamma.$$

Proof: Since 0 is a point of increase of μ_o , the inequality (16) is tight at 0. Therefore, the Lagrange multiplier λ can be obtained by putting $x = 0$ in (16). \square

For each $b \in \{2, 3, \dots\}$, let $\tau^{(b)}$ be an input probability measure in Λ that maximizes $I(\mu)$ and has b or fewer points of increase. Using an approach similar to that of [1], the capacity-achieving measure of this optical channel can be found via the following search algorithm.

Search Algorithm for Capacity-Achieving Measures

- Step 1: Set $b = 2$.
- Step 2: Solve for $\tau^{(b)}$.
- Step 3: Let $\lambda^{(b)} = [I(\tau^{(b)}) - Q(0; \tau^{(b)}) + \log 2\pi e \sigma_0^2] / \gamma$. If $\lambda^{(b)} < 0$, increase b by 1 and go back to step 2.
- Step 4: Verify whether the inequality

$$Q(x; \tau^{(b)}) - I(\tau^{(b)}) - \frac{1}{2} \log[2\pi e(\sigma_0^2 + \sigma_1^2 x)] - \lambda^{(b)}(x - \gamma) \leq 0$$

holds for all $x \in [0, s]$. If so, then $\tau^{(b)}$ is capacity achieving. If otherwise, increase b by 1 and go back to step 2.

Numerical Results

The parameter σ_1^2 determines the degree of signal dependency of the optical channel. It is interesting to know how the parameter σ_1^2 affects the capacity-achieving measure. In the above optical channel, we fix the following parameters: $\sigma_0^2 = 0.1$, $\gamma = 2$, and $s = 6$. Let $\mu[\sigma_1^2]$ be the capacity-achieving measure for the channel $(\sigma_0^2 + \sigma_1^2 x, x, \{x - \gamma\}, [0, s])$. With respect to different specified values of σ_1^2 , via the above search algorithm, $\mu[\sigma_1^2]$ was found and its set of points of increases are plotted in Fig. 2. Specifically, the horizontal axis denotes the values of σ_1^2 and the vertical axis denotes the input signal. If a dot is indicated in the

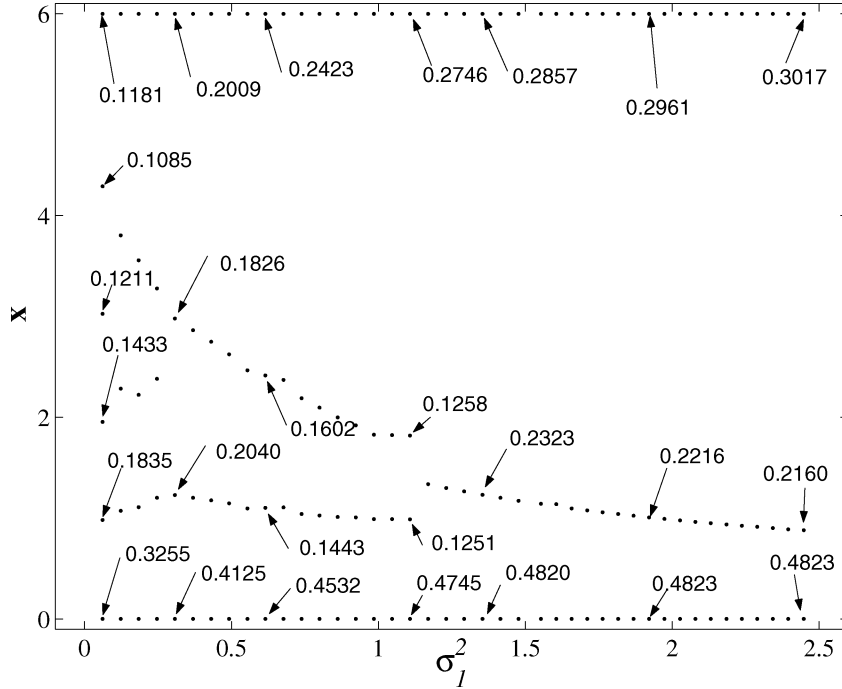


Fig. 2. Capacity-achieving measures.

position (σ_1^2, x^*) , then x^* is a point of increase of $\mu[\sigma_1^2]$. Moreover, the probabilities at the points of increase of $\mu[\sigma_1^2]$ are also shown in the figure for selected values of σ_1^2 .

From the figure, we observe the following.

1. 0 and s are always points of increase of $\mu[\sigma_1^2]$. It is proved in Proposition 1 that 0 is always a point of increase of $\mu[\sigma_1^2]$. However, it is not known whether s is also a point of increase of $\mu[\sigma_1^2]$ in general.
2. The distance between two neighboring points of increases of $\mu[\sigma_1^2]$ varies significantly. In particular, the separation is smaller if the interval is “closer” to 0. This phenomenon is especially apparent when the channel is highly signal dependent (i.e., when σ_1^2 is large). We believe that this phenomenon is due to the fact that when σ_1^2 is large, input signals of larger power require greater separation than those of smaller power.
3. As σ_1^2 increases, the size of $\mathbb{E}_{\mu[\sigma_1^2]}$ decreases. This happens because when σ_1^2 increases, the average noise power also increases; hence, a signaling scheme with a smaller constellation is more favorable.
4. The probabilities of “ $x = 0$ ” and “ $x = s$ ” increase, as σ_1^2 increases. For a fixed channel input, when σ_1^2 increases, the induced noise power also increases. Hence, for large σ_1^2 , it is more favorable to use inputs of smaller power. This explains why the probability of “ $x = 0$ ” increases. An unexpected observation from our numerical results is that the probability of “ $x = s$ ” also increases. We believe that this is because the increase in probabilities of “ $x = s$ ” and “ $x = 0$ ” also increases the “average distances” between channel codewords. According to our numerical results, this advantage seems to outweigh the disadvantage of the increase in average noise power.

B. Fading Channel

Consider a general multiple-antenna system operated in a fading environment. The channel input and output are complex-valued \tilde{N} -tuple $\tilde{\mathbf{X}}$ and \tilde{M} -tuple $\tilde{\mathbf{Y}}$, respectively, such that $\tilde{\mathbf{X}}$ and $\tilde{\mathbf{Y}}$ are related as $\tilde{\mathbf{Y}} = \tilde{\mathbf{H}}\tilde{\mathbf{X}} + \tilde{\mathbf{Z}}$. Here, we do not make any assumption on the complex channel gain matrix $\tilde{\mathbf{H}}$ and the additive noise $\tilde{\mathbf{Z}}$ except that 1) they are complex Gaussian distributed and 2) $\tilde{\mathbf{H}}$ and $\tilde{\mathbf{Z}}$ are independent to each other. Such a channel can also be formulated as a quadratic CG channel $\mathbf{Y} = \mathbf{H}\mathbf{X} + \mathbf{Z}$, where

$$\mathbf{Y} = \begin{bmatrix} \text{re}(\tilde{\mathbf{Y}}) \\ \text{im}(\tilde{\mathbf{Y}}) \end{bmatrix}, \quad \mathbf{H} = \begin{bmatrix} \text{re}(\tilde{\mathbf{H}}) & -\text{im}(\tilde{\mathbf{H}}) \\ \text{im}(\tilde{\mathbf{H}}) & \text{re}(\tilde{\mathbf{H}}) \end{bmatrix}$$

$$\mathbf{X} = \begin{bmatrix} \text{re}(\tilde{\mathbf{X}}) \\ \text{im}(\tilde{\mathbf{X}}) \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} \text{re}(\tilde{\mathbf{Z}}) \\ \text{im}(\tilde{\mathbf{Z}}) \end{bmatrix}. \quad (17)$$

Example: Rayleigh-Fading Channels: Consider the Rayleigh-fading channel, in which all the entries in $\tilde{\mathbf{H}}$ and $\tilde{\mathbf{Z}}$ are independent, zero-mean, complex Gaussian distributed with variance χ^2 and σ^2 , respectively. Hence, for any input $\mathbf{x} \in \mathbb{R}^N$, $\mathbf{b}_x = \mathbf{0}$ and

$$\Delta_x = \frac{1}{2} \left(\sum_{n \in \mathbb{N}} \chi^2 x_n^2 + \sigma^2 \right) \mathbf{I}_M.$$

Subject to peak and average total power constraints, the channel is denoted by the tuple

$$\left(\frac{1}{2} \left(\sum_{n \in \mathbb{N}} \chi^2 x_n^2 + \sigma^2 \right) \mathbf{I}_M, \mathbf{0}, \{\mathbf{x}^\top \mathbf{x} - \gamma\}, \mathbb{B}_N(s^2) \right)$$

where s^2 and γ are the maximal peak and average total power respectively.

Corollary 4: Let μ_o be the capacity-achieving measure for the above Rayleigh-fading channel. Then μ_o is discrete. In particular, its set of points of increase is contained in a finite number

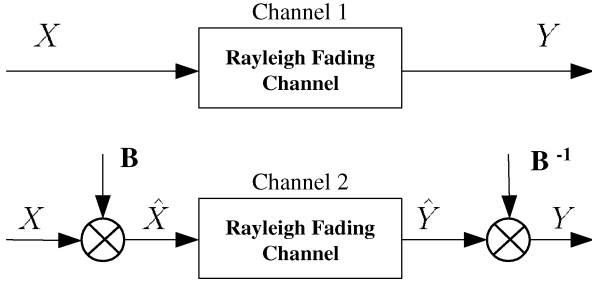


Fig. 3. Two equivalent Rayleigh-fading channels.

of “concentric shells” centered at $\mathbf{0}$. In other words, the probability distribution of $\|\mathbf{X}\|$ is discrete with a finite number of probability mass points.

Proof: Let \mathbb{W} be the maximal connected open subset of $\mathbb{A}(\Delta_{\mathbf{w}})$ containing $\mathbb{B}_N(s^2)$. Then, for any $0 \leq \epsilon < \sigma/\chi$, the vector $(\mathbf{j}\epsilon, 0, \dots, 0) \in \mathbb{W}$. Let $\{\mathbf{w}^{(i)} \triangleq (\mathbf{j}\epsilon_i, 0, \dots, 0)\}_{i=1}^{\infty}$ be a sequence in \mathbb{W} such that $\lim_{i \rightarrow \infty} \epsilon_i = \sigma/\chi$ and $0 \leq \epsilon_i < \sigma/\chi$ for all positive integers i . Obviously, $\lim_{i \rightarrow \infty} \det \Delta_{\mathbf{w}^{(i)}} = 0$ and for all positive integers i , $\mathbf{b}_{\mathbf{w}^{(i)}}$ and $\Delta_{\mathbf{w}^{(i)}}$ are real. The requirements in Case B of Theorem 3 are thus satisfied, and hence, μ_o is discrete. It remains to prove that \mathbb{E}_{μ_o} is a subset of a finite number of concentric shells centered at $\mathbf{0}$.

Consider the two channels depicted in Fig. 3. Channel 1 is the Rayleigh-fading channel, and channel 2 is obtained from channel 1 by multiplying a real orthonormal matrix \mathbf{B} and its inverse \mathbf{B}^{-1} before and after the Rayleigh-fading channel. It is easy to see that for any $\mathbf{x} \in \mathbb{R}^N$, the conditional probability density functions of \mathbf{Y} given $\mathbf{X} = \mathbf{x}$ is the same for both channels. Hence, if the input probability measures of the two channels are equal to μ_o , then the mutual information $I(\mathbf{X}; \mathbf{Y})$ in both channels is the same.

In the second channel, as $\mathbf{X}, \hat{\mathbf{X}}, \hat{\mathbf{Y}}$ and \mathbf{Y} form a Markov chain $\mathbf{X} \rightarrow \hat{\mathbf{X}} \rightarrow \hat{\mathbf{Y}} \rightarrow \mathbf{Y}$, we have $I(\hat{\mathbf{X}}; \hat{\mathbf{Y}}) \geq I(\mathbf{X}; \mathbf{Y})$. Let ν be the probability measure of $\hat{\mathbf{X}}$. It is obvious that ν satisfies all channel input constraints. Thus, ν is also capacity achieving, and $I(\mu_o) = I(\nu)$.

By Theorem 2, there exists $\lambda \geq 0$ such that for all $\mathbf{x} \in \mathbb{B}_N(s^2)$

$$Q(\mathbf{x}; \mu_o) - \frac{1}{2} \log[(2\pi e)^M \det \Delta_{\mathbf{x}}] - I(\mu_o) - \lambda(\mathbf{x}^T \mathbf{x} - \gamma) \leq 0. \quad (18)$$

Let \mathbb{E} be the subset of $\mathbb{B}_N(s^2)$ such that (18) is tight. If $\mathbf{x} \in \mathbb{E}_{\mu_o}$, then $\mathbf{B}\mathbf{x} \in \mathbb{E}_{\nu}$. Hence, \mathbf{x} and $\mathbf{B}\mathbf{x}$ are both in \mathbb{E} . Since \mathbf{B} is an arbitrary orthonormal matrix, \mathbb{E} contains the set

$$\{\mathbf{x} \in \mathbb{R}^N : \exists \mathbf{w} \in \mathbb{E}_{\mu_o} \text{ such that } \|\mathbf{x}\|^2 = \|\mathbf{w}\|^2\}$$

which is a collection of concentric shells centered at $\mathbf{0}$. By Lemma 1, the sparseness of \mathbb{E} implies that the number of concentric shells is finite and the corollary then follows. \square

As a remark, the capacity-achieving measure for the Rayleigh-fading channel has been studied separately in [4], in which it is proved that the capacity-achieving measure is discrete even when there is no peak power constraint.

Suppose the receiver can estimate the channel correctly. In other words, suppose the channel matrix \mathbf{H} is known to the receiver. Let $\hat{\mathbf{H}}$ be an $MN \times 1$ column vector such that it contains all the entries of \mathbf{H} . Such a fading channel with receiver-side channel information can still be formulated as a CG channel, where the output of the channel is $[\hat{\mathbf{H}}^T, \mathbf{Y}^T]^T$. It is easy to prove that for all $\mathbf{x} \in \mathbb{R}^N$, $\mathbf{b}_{\mathbf{x}}$ and $\Delta_{\mathbf{x}}$ are real,

$$\text{tr}(\Delta_{\mathbf{x}}) = MN\chi^2/2 + M\sigma^2/2 + M\|\mathbf{x}\|^2\chi^2/2,$$

and

$$\log \det \Delta_{\mathbf{x}} = MN \log(\chi^2/2) + M \log(\sigma^2/2), \quad \text{for all } \mathbf{x} \in \mathbb{R}^N.$$

Corollary 5: For the above Rayleigh-fading channel with channel side information available at the receiver, subject only to the peak total power constraint,³ there exists a capacity-achieving measure μ_o which is DAUP.

Proof: Let μ^* be a capacity-achieving measure and \mathbb{W} be the maximal connected open subset of $\mathbb{A}(\Delta_{\mathbf{w}})$ containing \mathbb{R}^N . As $\text{tr}(\Delta_{\mathbf{x}})$ is unbounded over \mathbb{R}^N , the requirements of Case A in Theorem 3 are satisfied. Hence, μ^* is discrete.

For any orthonormal matrix \mathbf{B} , let $\hat{\mathbf{X}} = \mathbf{B}\mathbf{X}$ and ν be its probability measure. Applying a similar technique when we proved Corollary 4, we can show that ν is also capacity-achieving. Since $I(\mu)$ is concave, any linear combination of μ^* and ν is still capacity achieving. By “averaging” over all possible orthonormal matrices \mathbf{B} , we can construct from μ^* a capacity-achieving measure μ_o which is uniform in phase. The phase uniformity of μ_o implies that \mathbb{E}_{μ_o} consists of a number of concentric shells centered at $\mathbf{0}$. Hence, the sparseness of \mathbb{E}_{μ_o} implies that the number of concentric shells in \mathbb{E}_{μ_o} is finite. In other words, μ_o is DAUP. \square

Example: MIMO Rayleigh Block-Fading Channels: Consider a MIMO Rayleigh block-fading channel [15], whose channel gain matrix remains constant for T symbol periods, after which it changes to a new set of values independently, and again, maintain for another T symbol periods, and so on. Let \tilde{N} and \tilde{M} be the number of transmitting and receiving antennas, respectively. Within a block of T channel input vectors, for $l = 1, \dots, T$, let

$$\tilde{\mathbf{X}}^{(l)} = [\tilde{X}_1^{(l)}, \dots, \tilde{X}_{\tilde{N}}^{(l)}]^T$$

be the l th channel input vector and

$$\tilde{\mathbf{Y}}^{(l)} = [\tilde{Y}_1^{(l)}, \dots, \tilde{Y}_{\tilde{M}}^{(l)}]^T$$

be the corresponding channel output vector. Then $\tilde{\mathbf{X}}^{(l)}$ and $\tilde{\mathbf{Y}}^{(l)}$ are related by the equation $\tilde{\mathbf{Y}}^{(l)} = \tilde{\mathbf{H}}\tilde{\mathbf{X}}^{(l)} + \tilde{\mathbf{Z}}^{(l)}$ where $\tilde{\mathbf{H}}$ is the complex channel gain matrix for the Rayleigh channel as defined in Section IV-B and $\tilde{\mathbf{Z}}^{(l)}$ is the corresponding additive complex Gaussian noise. For simplicity, let

$$\mathbf{X}^{(l)} = \begin{bmatrix} \text{re}(\tilde{\mathbf{X}}^{(l)}) \\ \text{im}(\tilde{\mathbf{X}}^{(l)}) \end{bmatrix}, \quad \mathbf{Y}^{(l)} = \begin{bmatrix} \text{re}(\tilde{\mathbf{Y}}^{(l)}) \\ \text{im}(\tilde{\mathbf{Y}}^{(l)}) \end{bmatrix}$$

³The peak total power is the sum of the powers of all transmitting antennae

$$\begin{aligned} \mathbf{Z}^{(l)} &= \begin{bmatrix} \text{re}(\tilde{\mathbf{Z}}^{(l)}) \\ \text{im}(\tilde{\mathbf{Z}}^{(l)}) \end{bmatrix}, \mathbf{H}^* = \begin{bmatrix} \text{re}(\tilde{\mathbf{H}}) & -\text{im}(\tilde{\mathbf{H}}) \\ \text{im}(\tilde{\mathbf{H}}) & \text{re}(\tilde{\mathbf{H}}) \end{bmatrix} \\ \mathbf{X} &= \left[(\mathbf{X}^{(1)})^\top, \dots, (\mathbf{X}^{(T)})^\top \right]^\top \\ \mathbf{Y} &= \left[(\mathbf{Y}^{(1)})^\top, \dots, (\mathbf{Y}^{(T)})^\top \right]^\top \\ \mathbf{Z} &= \left[(\mathbf{Z}^{(1)})^\top, \dots, (\mathbf{Z}^{(T)})^\top \right]^\top. \end{aligned} \quad (19)$$

It is easy to show that the MIMO Rayleigh block-fading channel can be rewritten as a quadratic CG channel $\mathbf{Y} = \mathbf{H}\mathbf{X} + \mathbf{Z}$ where \mathbf{H} is a $2\tilde{M}T \times 2\tilde{N}T$ “block-diagonal matrix” whose diagonal blocks are all equal to \mathbf{H}^* . Since the channel is quadratic, the well-behaved region $\mathbb{A}(\Delta_{\mathbf{w}})$ contains $\mathbb{R}^{2T\tilde{N}}$. As \mathbf{H} is zero mean, $\mathbf{b}_{\mathbf{x}} = \mathbf{0}$ for all $\mathbf{x} \in \mathbb{R}^{2\tilde{N}T}$.

Corollary 6: Let μ_o be the capacity-achieving measure for the above Rayleigh block-fading channel subject to peak and average power constraint. Then μ_o is discrete.

Proof: Let \mathbb{W} be the maximal connected open subset of $\mathbb{A}(\Delta_{\mathbf{w}})$ that contains $\mathbb{R}^{2\tilde{N}T}$. Let $\mathbf{w}^* = (j\epsilon, 0, \dots, 0)$ be a vector in $\mathbb{C}^{2\tilde{N}T}$ where $0 \leq \epsilon < \sigma/\chi$. By direct substitution, it can be verified easily that $\mathbf{w}^* \in \mathbb{W}$. Furthermore, $\Delta_{\mathbf{w}^*}$ are real and

$$\det \Delta_{\mathbf{w}^*} = \left(\frac{\sigma^2 - \chi^2 \epsilon^2}{2} \right)^{2\tilde{M}} \left(\frac{\sigma^2}{2} \right)^{2(T-1)\tilde{M}}. \quad (20)$$

Let $\{\mathbf{w}^{(i)} \triangleq (j\epsilon_i, 0, \dots, 0)\}_{i=1}^\infty$ be a sequence in \mathbb{W} such that $\lim_{i \rightarrow \infty} \epsilon_i = \sigma/\chi$ and $0 \leq \epsilon_i < \sigma/\chi$ for all positive integers i . Obviously, $\lim_{i \rightarrow \infty} \det \Delta_{\mathbf{w}^{(i)}} = 0$ and for all positive integers i , $\mathbf{b}_{\mathbf{w}^{(i)}}$ and $\Delta_{\mathbf{w}^{(i)}}$ are real. The requirements in Case B of Theorem 3 are thus satisfied, and hence, μ_o is discrete. \square

As a remark, the capacity-achieving measure for the single-antenna Rayleigh block-fading channel has been studied separately in [8], in which it is proved that the capacity-achieving measure is discrete even when there is no peak power constraint.

Example: Rician Fading Channels: The Rician fading channel is a generalization of the Rayleigh-fading channel in the sense that $\mathbf{E}[\mathbf{H}]$ is not necessarily a zero matrix as in the case of Rayleigh-fading channel. Specifically, given the channel input is \mathbf{x} ,

$$\Delta_{\mathbf{x}} = \frac{1}{2} \left(\sum_{n \in \mathbb{N}} \chi^2 x_n^2 + \sigma^2 \right)$$

and $\mathbf{b}_{\mathbf{x}} = \mathbf{E}[\mathbf{H}]\mathbf{x}$ which is nonzero in general.

Corollary 7: The capacity-achieving probability measure for the Rician fading channel

$$\left(\frac{1}{2} \left(\sum_{n \in \mathbb{N}} \chi^2 x_n^2 + \sigma^2 \right) \mathbf{I}_M, \mathbf{E}[\mathbf{H}]\mathbf{x}, \{\}, \mathbb{B}_N(s^2) \right)$$

is discrete.

Proof: Since a Rician fading channel is a quadratic CG channel, by Corollary 1, the result follows. \square

As a remark, many fading channels can be formulated as quadratic CG channels. Therefore, according to Corollary 1, their capacity-achieving measures are discrete if there is no average cost constraint.

C. Example: Interference Channels

Consider a communication system consisting of N pairs of transmitters and receivers. Each pair is connected by a point-to-point subchannel, and the subchannels interfere with each other. Specifically, the channel input \mathbf{X} and output \mathbf{Y} are related as $\mathbf{Y} = \mathbf{H}\mathbf{X} + \mathbf{Z}$ such that 1) the entries of \mathbf{H} and \mathbf{Z} are Gaussian distributed independently with variances χ^2 and σ^2 , respectively, 2) \mathbf{Z} is zero mean and $\mathbf{E}[\mathbf{H}]$ is the identity matrix. Therefore, for any input \mathbf{x} , $\mathbf{b}_{\mathbf{x}} = \mathbf{x}$, and

$$\Delta_{\mathbf{x}} = \left(\sum_{n \in \mathbb{N}} \chi^2 x_n^2 + \sigma^2 \right) \mathbf{I}_N.$$

Corollary 8: There exists a DAUP capacity-achieving measure μ_o for the interference channel

$$\left(\left(\sum_{n \in \mathbb{N}} \chi^2 x_n^2 + \sigma^2 \right) \mathbf{I}_N, \mathbf{x}, \{\}, \mathbb{B}_N(s^2) \right).$$

Proof: Since the channel is quadratic and there is no average cost constraint, the capacity-achieving measure is discrete. The proof of that μ_o is DAUP is the same as the one given in Corollary 5. \square

D. Example: Parallel Gaussian Channels

Consider a system of N independent and identically scalar additive Gaussian channels in parallel subject to an average total power constraint $\mathbf{E}_\mu[\mathbf{X}^\top \mathbf{X} - \gamma] \leq 0$ and a bounded-input constraint. In other words, the channel is characterized by the tuple $(\sigma^2 \mathbf{I}_N, \mathbf{x}, \{\mathbf{x}^\top \mathbf{x} - \gamma\}, \mathbb{S})$. By Theorems 1 and 3, the capacity-achieving measure is unique and discrete. We will consider two cases of bounded-input constraints: cubic and spherical constraints.

Case A: Cubic Constraint: In the first case, the set of admissible channel input vectors is a hypercube, defined as $\mathbb{S} = \prod_{n \in \mathbb{N}} \mathbb{S}_n$ where $\mathbb{S}_n = \{x_n \in \mathbb{R} : x_n^2 \leq s^2\}$. In other words, the peak power of each component channel is s^2 . For any input probability measure μ of \mathbf{X} , we denote the corresponding marginal probability measure of X_n by $\mu^{(n)}$.

Let μ_o be the capacity-achieving measure for the channel $(\sigma^2 \mathbf{I}_N, \mathbf{x}, \{\mathbf{x}^\top \mathbf{x} - \gamma\}, \prod_{n \in \mathbb{N}} \mathbb{S}_n)$ and ν be the product measure of $\{\mu_o^{(1)}, \dots, \mu_o^{(N)}\}$. It can be proved easily that $I(\nu) \geq I(\mu_o)$ and $\nu \in \Lambda$. By the uniqueness of μ_o , we have $\mu_o = \nu$ and $\mu^{(1)} = \mu^{(2)} = \dots = \mu^{(N)}$. By direct substitution, we can verify immediately that $I(\nu) = NC$ where C is the channel capacity of the channel $(\sigma^2, x_n, \{x_n^2 - \gamma/N\}, [-s, s])$ and $\mu_o^{(n)}$ is the corresponding capacity-achieving measure. Again, by Theorem 3, $\mu_o^{(n)}$ is discrete and hence has a finite set of points of increase. As $\mathbb{E}_{\mu_o} = \prod_{n \in \mathbb{N}} \mathbb{E}_{\mu_o^{(n)}}$, \mathbb{E}_{μ_o} is also finite.

Case B: Spherical Constraint: In the second case, the set of admissible channel input vectors is $\mathbb{B}_N(s^2)$. Let μ_o be the capacity-achieving measure for the channel $(\sigma^2 \mathbf{I}_N, \mathbf{x}, \{\mathbf{x}^\top \mathbf{x} - \gamma\}, \mathbb{B}_N(s^2))$. For any orthonormal matrix \mathbf{B} , let $\tilde{\mathbf{X}} = \mathbf{B}\mathbf{X}$ and ν be its probability measure. Again, using the same technique as in proving Corollary 4, we can show that ν is also capacity achieving. The uniqueness of the capacity-achieving measures thus implies that $\mu_o = \nu$. Hence, μ_o is uniform in phase, and \mathbb{E}_{μ_o} consists of a number of concentric shells centered at $\mathbf{0}$. The sparseness of \mathbb{E}_{μ_o} thus implies that the number of concentric

shells in \mathbb{E}_{μ_o} is finite. In other words, μ_o is DAUP. The special case when $N = 2$ was proved by Shamai and Bar-David in [3].

V. CONCLUSION

In this paper, we have shown that, in many instances, the capacity-achieving probability measure for a conditionally Gaussian channel with bounded-input constraints is discrete. The criteria for discreteness given in Theorem 3 are not exhaustive, and hence, there may be many more examples of CG channels with a discrete capacity-achieving measure. In fact, according to Theorem 2, for a conditionally Gaussian channel subject to bounded-input constraints, if its capacity-achieving probability measure is not discrete, then the inequality in (6) must be tight for all admissible channel inputs \mathbf{x} . We believe that such a requirement is rather stringent in general. In other words, it is more natural to expect that most capacity-achieving measures are discrete.

There is much work yet to be done in the analysis of channels with input constraints. The tools we used allow us to verify whether a given input measure is capacity achieving or not. Once the capacity-achieving measure and the channel capacity are determined, they can serve as a benchmark to evaluate the efficiency of any practical scheme. Furthermore, the obtained capacity-achieving measure may lead to valuable insights in designing practical and efficient modulation schemes.

However, except in a few special cases, when we may guess the form (e.g., DAUP) of the capacity-achieving probability measure, these capacity-achieving probability measures are extremely difficult to obtain. More sophisticated techniques are yet to be discovered to overcome this problem.

APPENDIX I

PROOF OF LEMMA 1

Suppose \mathbb{I} is finite. Let $b(\mathbf{w}) = \prod_{i \in \mathbb{I}} (\mathbf{w}^T \mathbf{w} - p_i^2)$. Clearly, the function $b(\mathbf{w})$ is holomorphic over \mathbb{C}^N and is zero in \mathbb{F} . Therefore, \mathbb{F} is sparse. It remains to prove that if \mathbb{I} is not finite, then \mathbb{F} is not sparse.

Suppose \mathbb{I} is not finite. The Bolzano–Weierstrass theorem implies that $\{p_i : i \in \mathbb{I}\}$ has a limiting point p . Let $b(\mathbf{w})$ be a holomorphic function defined on a connected open subset \mathbb{U} of \mathbb{C}^N containing the closure of \mathbb{F} such that $b(\mathbf{w}) = 0$ for all $\mathbf{w} \in \mathbb{F}$.

First, consider the case when $p > 0$. Let \mathbf{v} be any point in \mathbb{R}^N such that $\mathbf{v}^T \mathbf{v} = p^2$. Then \mathbf{v} is a limiting point of \mathbb{F} , and is contained in \mathbb{U} . Let $\mathbb{W}(\epsilon, \delta)$ be as in the equation at the bottom of the page. By picking $\epsilon, \delta > 0$ small enough, $\mathbb{W}(\epsilon, \delta)$ can be made to be a subset of \mathbb{U} .

Let $\mathbf{x} \in \mathbb{W}(\epsilon, \delta)$ be a vector such that $\mathbf{x}^T \mathbf{x} = p^2$. Let $f(v) = b(v\mathbf{x})$ be the single-variable complex-valued function, and \mathbb{F} be the set

$$\{\zeta \in (1 - \delta, 1 + \delta) : \exists p_i \text{ such that } \zeta = p_i/p \text{ for some } i \in \mathbb{I}\}.$$

It is easy to see that $f(v)$ is an analytic function on an open subset of \mathbb{C} containing the interval $(1 - \delta, 1 + \delta)$, and $f(v)$ is zero on \mathbb{F} . As p is a limiting point of $\{p_i : i \in \mathbb{I}\}$, 1 is a limiting point of \mathbb{F} . By the identity theorem [20], $f(v) = 0$ for all $v \in (1 - \delta, 1 + \delta)$. In other words, $b(v\mathbf{x}) = 0$ for all $v \in (1 - \delta, 1 + \delta)$. As

$$\mathbb{W}(\epsilon, \delta) = \bigcup_{\mathbf{x} \in \mathbb{W}(\epsilon, \delta) : \mathbf{x}^T \mathbf{x} = p^2} \{v\mathbf{x} : 1 - \delta < v < 1 + \delta\} \quad (21)$$

$b(\mathbf{w}) = 0$ for all $\mathbf{w} \in \mathbb{W}(\epsilon, \delta)$. Since $\mathbb{W}(\epsilon, \delta)$ is open in \mathbb{R}^N , again by the identity theorem, $b(\mathbf{w}) = 0$ for all $\mathbf{w} \in \mathbb{U}$. This implies that \mathbb{F} is not sparse.

In the case when $p = 0$, then the origin $\mathbf{0}$ is a limiting point of \mathbb{F} . Let $\mathbb{W}(\delta) = \{\mathbf{u} \in \mathbb{R}^N : \|\mathbf{u}\| < \delta\}$. Picking $\delta > 0$ small enough, $\mathbb{W}(\delta)$ is a subset of \mathbb{U} . For any $\mathbf{x} \in \mathbb{R}^N$ such that $\mathbf{x}^T \mathbf{x} = 1$, let $f(v)$ be the single-variable complex-valued function $b(v\mathbf{x})$. Then $f(v)$ is analytic on an open subset of \mathbb{C} containing the interval $(-\delta, \delta)$. As $\{p_i \in (-\delta, \delta) : i \in \mathbb{I}\}$ has a limiting point at 0 and $f(v)$ is zero on this set, $f(v) = 0$ for $v \in (-\delta, \delta)$, and consequently, $b(\mathbf{w}) = 0$ for all $\mathbf{w} \in \mathbb{W}(\delta)$. Again, by the identity theorem, $b(\mathbf{w}) = 0$ for all $\mathbf{w} \in \mathbb{U}$. The lemma is thus proved.

APPENDIX II

SUPPLEMENTARY RESULTS

Let \mathcal{B}_N be the vector space of all bounded and continuous real-valued functions defined on \mathbb{R}^N and \mathcal{B}_N^* be its dual space [21]. The set of input probability measures, denoted by Ω , can be identified as a subset of \mathcal{B}_N^* . In this paper, we assume that the topology on Ω is induced by the weak* topology defined on \mathcal{B}_N^* . The induced topology on Ω is metrizable, whose metric is called the Lévy metric [23]. This topology is complete and satisfies the following properties.

Lemma 2: Let $\{\mu_i\}_{i=0}^\infty$ be a sequence in Ω .

1. The sequence converges to μ_0 in the Lévy metric, denoted by $\mu_i \implies \mu_0$, if and only if for every $f \in \mathcal{B}_N$

$$\lim_{i \rightarrow \infty} \int f(\mathbf{x}) d\mu_i = \int f(\mathbf{x}) d\mu_0.$$

In particular, if $\{\mu_i\}_{i=0}^\infty$ is a sequence in $\Lambda_{\mathbb{S}}$, then $\mu_i \implies \mu_0$ if and only if for all-continuous function f defined on \mathbb{S}

$$\lim_{i \rightarrow \infty} \int f(\mathbf{x}) d\mu_i = \int f(\mathbf{x}) d\mu_0.$$

2. Suppose $\mu_i \implies \mu_0$. For any closed subset \mathbb{K} of \mathbb{R}^N ,

$$\overline{\lim}_{i \rightarrow \infty} \mu_i(\mathbb{K}) \leq \mu_0(\mathbb{K}).$$

Moreover, if $f(\mathbf{x})$ is continuous on \mathbb{R}^N and is bounded below, then

$$\underline{\lim}_{i \rightarrow \infty} \int f d\mu_i \geq \int f d\mu_0.$$

3. A subset Υ of Ω is called **tight** if for any $\epsilon > 0$, there exists a closed and bounded subset \mathbb{K} of \mathbb{R}^N such that

$$\mathbb{W}(\epsilon, \delta) = \{\mathbf{u} \in \mathbb{R}^N : \|\mathbf{u}\| \|\mathbf{v}\| (1 - \epsilon) < \mathbf{u}^T \mathbf{v} \text{ and } (1 - \delta)p < \|\mathbf{u}\| < (1 + \delta)p\}.$$

$\mu(\mathbb{K}) \geq 1 - \epsilon$ for all $\mu \in \Upsilon$. Suppose Υ is closed and tight. Then it is sequentially compact.

Proof: See [23, Sec. 3.1]. \square

Proposition 2 (Convexity and Sequential Compactness): The subsets $\Lambda_{\mathcal{S}}$ and Λ of Ω are convex and sequentially compact.

Proof: The convexity of $\Lambda_{\mathcal{S}}$ and Λ are trivial. To prove that they are sequentially compact, it is sufficient to show that they are tight and closed. First, notice that \mathcal{S} is closed and bounded. As $\mu(\mathcal{S}) = 1$ for all $\mu \in \Lambda_{\mathcal{S}}$, $\Lambda_{\mathcal{S}}$ and its subset Λ are tight.

Let $\{\mu_i\}_{i=1}^{\infty}$ be a convergent sequence in $\Lambda_{\mathcal{S}}$ with limit μ_0 . Since \mathcal{S} is closed, $\overline{\lim}_{i \rightarrow \infty} \mu_i(\mathcal{S}) \leq \mu_0(\mathcal{S})$, which further implies that $\mu_0(\mathcal{S}) = 1$, or, equivalently, $\mu_0 \in \Lambda_{\mathcal{S}}$. Therefore, $\Lambda_{\mathcal{S}}$ is closed, and thus sequentially compact.

Similarly, let $\{\mu_i\}_{i=1}^{\infty}$ be a convergent sequence in Λ with limit μ_0 . The sequential compactness of $\Lambda_{\mathcal{S}}$ implies that $\mu_0 \in \Lambda_{\mathcal{S}}$. As $\mu_i \Rightarrow \mu_0$ and $g_k(\mathbf{x})$ is continuous on \mathbb{R}^N for all $k \in \mathbb{K}$, by Lemma 2, we have

$$0 \geq \lim_{i \rightarrow \infty} \int g_k(\mathbf{x}) d\mu_i = \int g_k(\mathbf{x}) d\mu_0.$$

Hence, $\mu_0 \in \Lambda$ and the sequential compactness of Λ then follows. \square

Lemma 3: There exist positive constants α, β, β' such that for all $\mathbf{x} \in \mathcal{S}$ and $\mathbf{y} \in \mathbb{R}^M$

$$\exp(-\alpha - \beta' \|\mathbf{y}\|^2) \leq P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}) \leq \exp(\alpha - \beta \|\mathbf{y}\|^2).$$

Hence, for all $\mu \in \Lambda_{\mathcal{S}}$

$$\exp(-\alpha - \beta' \|\mathbf{y}\|^2) \leq P_{\mathbf{Y}}(\mathbf{y}; \mu) \leq \exp(\alpha - \beta \|\mathbf{y}\|^2)$$

and $|\log P_{\mathbf{Y}}(\mathbf{y}; \mu)| \leq \alpha + \beta' \|\mathbf{y}\|^2$.

Proof: As stated in Section II, there exist $\kappa_l, \kappa_h, \vartheta > 0$ such that for all $\mathbf{x} \in \mathcal{S}$, all eigenvalues of $\nabla_{\mathbf{x}}$ are within the interval (κ_l, κ_h) and $\|\mathbf{b}_{\mathbf{x}}\| < \vartheta$. Hence, $\det \Delta_{\mathbf{x}}$ is within the interval $(1/\kappa_h^M, 1/\kappa_l^M)$. As a result

$$\begin{aligned} -\frac{\kappa_h}{2} \|\mathbf{y} - \mathbf{b}_{\mathbf{x}}\|^2 &\leq -\frac{1}{2} (\mathbf{y} - \mathbf{b}_{\mathbf{x}})^{\top} \nabla_{\mathbf{x}} (\mathbf{y} - \mathbf{b}_{\mathbf{x}}) \\ &\leq -\frac{\kappa_l}{2} \|\mathbf{y} - \mathbf{b}_{\mathbf{x}}\|^2. \end{aligned} \quad (22)$$

Hence, for any $0 < \epsilon < 1$ and $\|\mathbf{y}\| > \vartheta/\epsilon$, we have

$$\begin{aligned} -\frac{\kappa_h}{2} (1 + \epsilon)^2 \|\mathbf{y}\|^2 &\leq -\frac{1}{2} (\mathbf{y} - \mathbf{b}_{\mathbf{x}})^{\top} \nabla_{\mathbf{x}} (\mathbf{y} - \mathbf{b}_{\mathbf{x}}) \\ &\leq -\frac{\kappa_l}{2} (1 - \epsilon)^2 \|\mathbf{y}\|^2 \end{aligned} \quad (23)$$

and

$$\begin{aligned} \frac{\exp(-\frac{\kappa_h}{2} (1 + \epsilon)^2 \|\mathbf{y}\|^2)}{\sqrt{(2\pi/\kappa_l)^M}} &\leq P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}) \\ &\leq \frac{\exp(-\frac{\kappa_l}{2} (1 - \epsilon)^2 \|\mathbf{y}\|^2)}{\sqrt{(2\pi/\kappa_h)^M}}. \end{aligned} \quad (24)$$

Let $\beta' = \frac{\kappa_h}{2} (1 + \epsilon)^2$ and $\beta = \frac{\kappa_l}{2} (1 - \epsilon)^2$. Then $\beta', \beta > 0$. If α is large enough, then

$$\exp(-\alpha - \beta' \|\mathbf{y}\|^2) \leq P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}) \leq \exp(\alpha - \beta \|\mathbf{y}\|^2)$$

for all $\mathbf{y} \in \mathbb{R}^M$ and the lemma is proved. \square

Proposition 3 (Continuity): $H_{\mathbf{Y}}(\mu)$, $H_{\mathbf{Y}|\mathbf{X}}(\mu)$, and $I(\mu)$ are continuous over $\Lambda_{\mathcal{S}}$.

Proof: Let $\{\mu_i\}_{i=0}^{\infty}$ be a sequence in $\Lambda_{\mathcal{S}}$ such that $\mu_i \Rightarrow \mu_0$. For any fixed $\mathbf{y} \in \mathbb{R}^M$, $P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x})$ is a continuous function of \mathbf{x} over \mathcal{S} . Hence,

$$\lim_{i \rightarrow \infty} \int P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}) d\mu_i = \int P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}) d\mu_0$$

i.e., $P_{\mathbf{Y}}(\mathbf{y}; \mu_i)$ is pointwise convergent to $P_{\mathbf{Y}}(\mathbf{y}; \mu_0)$. By Lemma 3, there exist $\alpha, \beta, \beta' > 0$ such that

$$| -P_{\mathbf{Y}}(\mathbf{y}; \mu_i) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_i) | \leq (\alpha + \beta' \|\mathbf{y}\|^2) \exp(\alpha - \beta \|\mathbf{y}\|^2).$$

By the Lebesgue convergence theorem

$$\begin{aligned} \lim_{i \rightarrow \infty} - \int P_{\mathbf{Y}}(\mathbf{y}; \mu_i) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_i) d\mathbf{y} \\ = - \int P_{\mathbf{Y}}(\mathbf{y}; \mu_0) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_0) d\mathbf{y}. \end{aligned} \quad (25)$$

Thus, $H_{\mathbf{Y}}(\mu)$ is continuous over $\Lambda_{\mathcal{S}}$.

The continuity of $H_{\mathbf{Y}|\mathbf{X}}(\mu)$ follows from the fact that $\frac{1}{2} \log[(2\pi e)^M \det \Delta_{\mathbf{x}}]$ is continuous over \mathcal{S} and

$$H_{\mathbf{Y}|\mathbf{X}}(\mu) \triangleq \int \frac{1}{2} \log[(2\pi e)^M \det \Delta_{\mathbf{x}}] d\mu.$$

Finally, as $I(\mu) = H_{\mathbf{Y}}(\mu) - H_{\mathbf{Y}|\mathbf{X}}(\mu)$, the continuity of $I(\mu)$ is established. \square

Lemma 4: Consider a constant CG channel $\mathbf{Y} = \mathbf{H}\mathbf{X} + \mathbf{Z}$. Let $\mathbf{R} = \mathbf{H}\mathbf{X}$ and $P_{\mathbf{R}}(\mu)$ be the probability measure of \mathbf{R} where μ is the probability measure of \mathbf{X} . Then, $P_{\mathbf{Y}}(\mathbf{y}; \mu_1) = P_{\mathbf{Y}}(\mathbf{y}; \mu_2)$ if and only if $P_{\mathbf{R}}(\mu_1) = P_{\mathbf{R}}(\mu_2)$. Consequently, if \mathbf{H} is left-invertible, then $P_{\mathbf{Y}}(\mathbf{y}; \mu_1) = P_{\mathbf{Y}}(\mathbf{y}; \mu_2)$ if and only if $\mu_1 = \mu_2$.

Proof: It is obvious that $P_{\mathbf{Y}}(\mathbf{y}; \mu)$ is equal to the convolution of $P_{\mathbf{Z}}(\mathbf{z})$ and $P_{\mathbf{R}}(\mu)$, where $P_{\mathbf{Z}}(\mathbf{z})$ is the probability density function of the Gaussian noise \mathbf{Z} . Since the Fourier transform of $P_{\mathbf{Z}}(\mathbf{z})$ is nonzero everywhere, $P_{\mathbf{Y}}(\mathbf{y}; \mu)$ and $P_{\mathbf{R}}(\mu)$ are in one-to-one correspondence (see [1] and Lemma 7). In other words, for any input probability measures μ_1 and μ_2 , $P_{\mathbf{Y}}(\mathbf{y}; \mu_1) = P_{\mathbf{Y}}(\mathbf{y}; \mu_2)$ if and only if $P_{\mathbf{R}}(\mu_1) = P_{\mathbf{R}}(\mu_2)$. \square

Definition 6: Let $J(\mu)$ be a real-valued function defined on $\Lambda_{\mathcal{S}}$ and $\mu_0 \in \Lambda_{\mathcal{S}}$. If the limit

$$J'_{\mu_0}(\mu_1) = \lim_{\theta \rightarrow 0^+} \frac{J((1 - \theta)\mu_0 + \theta\mu_1) - J(\mu_0)}{\theta} \quad (26)$$

exists for all $\mu_1 \in \Lambda_{\mathcal{S}}$, then the function $J(\mu)$ is called weakly differentiable [1] at μ_0 , and the function $J'_{\mu_0}(\mu_1)$ is its weak derivative at μ_0 .

Note that, if $J(\mu)$ is concave and is weakly differentiable on $\Lambda_{\mathcal{S}}$, then $J(\mu_0)$ is maximized if and only if $J'_{\mu_0}(\mu_1) \leq 0$ for all $\mu_1 \in \Lambda_{\mathcal{S}}$.

Proposition 4 (Weak Differentiability): The mutual information function $I(\mu)$ is weakly differentiable in $\Lambda_{\mathcal{S}}$, and

$$I'_{\mu_0}(\mu_1) = \int Q(\mathbf{x}; \mu_0) d\mu_1 - H_{\mathbf{Y}|\mathbf{X}}(\mu_1) - I(\mu_0).$$

Proof: Let $0 \leq \theta \leq 1$ and $\mu_\theta = (1 - \theta)\mu_0 + \theta\mu_1$. It can be shown easily that $\frac{I(\mu_\theta) - I(\mu_0)}{\theta}$ is equal to

$$-I(\mu_0) - H_{\mathbf{Y}|\mathbf{X}}(\mu_1) - \int P_{\mathbf{Y}}(\mathbf{y}; \mu_1) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_\theta) d\mathbf{y} \\ + \frac{1 - \theta}{\theta} \int P_{\mathbf{Y}}(\mathbf{y}; \mu_0) \log \frac{P_{\mathbf{Y}}(\mathbf{y}; \mu_0)}{P_{\mathbf{Y}}(\mathbf{y}; \mu_\theta)} d\mathbf{y}. \quad (27)$$

Since $\int Q(\mathbf{x}; \mu_0) d\mu_1 = -\int P_{\mathbf{Y}}(\mathbf{y}; \mu_1) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_0) d\mathbf{y}$, to prove Proposition 4, it suffices to prove the following claims.

Claim 1:

$$\lim_{\theta \rightarrow 0^+} \int P_{\mathbf{Y}}(\mathbf{y}; \mu_1) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_\theta) d\mathbf{y} \\ = \int P_{\mathbf{Y}}(\mathbf{y}; \mu_1) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_0) d\mathbf{y}.$$

Claim 2:

$$\lim_{\theta \rightarrow 0^+} \frac{1}{\theta} \int P_{\mathbf{Y}}(\mathbf{y}; \mu_0) \log \frac{P_{\mathbf{Y}}(\mathbf{y}; \mu_\theta)}{P_{\mathbf{Y}}(\mathbf{y}; \mu_0)} d\mathbf{y} = 0.$$

As the limit $\lim_{\theta \rightarrow 0^+} P_{\mathbf{Y}}(\mathbf{y}; \mu_1) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_\theta)$ is equal to $P_{\mathbf{Y}}(\mathbf{y}; \mu_1) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_0)$ and there exists $\alpha, \beta, \beta' > 0$ such that

$$|P_{\mathbf{Y}}(\mathbf{y}; \mu_1) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_\theta)| \leq (\alpha + \beta' \|\mathbf{y}\|^2) \exp(\alpha - \beta \|\mathbf{y}\|^2)$$

the first claim thus follows from the Lebesgue convergence theorem.

To prove the second claim, first notice that for $c, c_1, c_2 > 0$, $\log(1 - \theta + \theta c) \geq \theta \log c$ and $c_1 \log(c_2/c_1) \leq c_2 - c_1$. Let

$$c = P_{\mathbf{Y}}(\mathbf{y}; \mu_1) / P_{\mathbf{Y}}(\mathbf{y}; \mu_0)$$

$$c_1 = P_{\mathbf{Y}}(\mathbf{y}; \mu_0)$$

and

$$c_2 = P_{\mathbf{Y}}(\mathbf{y}; \mu_\theta).$$

Then it is easy to show that

$$P_{\mathbf{Y}}(\mathbf{y}; \mu_0) \log \frac{P_{\mathbf{Y}}(\mathbf{y}; \mu_1)}{P_{\mathbf{Y}}(\mathbf{y}; \mu_0)} \leq \frac{P_{\mathbf{Y}}(\mathbf{y}; \mu_0)}{\theta} \log \frac{P_{\mathbf{Y}}(\mathbf{y}; \mu_\theta)}{P_{\mathbf{Y}}(\mathbf{y}; \mu_0)} \\ \leq P_{\mathbf{Y}}(\mathbf{y}; \mu_1) - P_{\mathbf{Y}}(\mathbf{y}; \mu_0).$$

By L'Hospital's rule, it can be shown that

$$\lim_{\theta \rightarrow 0^+} \frac{P_{\mathbf{Y}}(\mathbf{y}; \mu_0)}{\theta} \log \frac{P_{\mathbf{Y}}(\mathbf{y}; \mu_\theta)}{P_{\mathbf{Y}}(\mathbf{y}; \mu_0)} = P_{\mathbf{Y}}(\mathbf{y}; \mu_1) - P_{\mathbf{Y}}(\mathbf{y}; \mu_0).$$

Since $P_{\mathbf{Y}}(\mathbf{y}; \mu_0) \log \frac{P_{\mathbf{Y}}(\mathbf{y}; \mu_1)}{P_{\mathbf{Y}}(\mathbf{y}; \mu_0)}$ and $P_{\mathbf{Y}}(\mathbf{y}; \mu_1) - P_{\mathbf{Y}}(\mathbf{y}; \mu_0)$ are integrable, by the Lebesgue convergence theorem

$$\lim_{\theta \rightarrow 0^+} \int \frac{P_{\mathbf{Y}}(\mathbf{y}; \mu_0)}{\theta} \log \frac{P_{\mathbf{Y}}(\mathbf{y}; \mu_\theta)}{P_{\mathbf{Y}}(\mathbf{y}; \mu_0)} d\mathbf{y} \\ = \int P_{\mathbf{Y}}(\mathbf{y}; \mu_1) - P_{\mathbf{Y}}(\mathbf{y}; \mu_0) d\mathbf{y} \quad (28)$$

$$= 0. \quad (29)$$

The proposition is thus proved. \square

Lemma 5: $Q(\mathbf{x}; \mu_0)$ is a continuous function of \mathbf{x} over \mathcal{S} .

Proof: Suppose $\{\mathbf{x}^{(i)}\}_{i=1}^\infty$ is a convergent sequence in \mathcal{S} with limit $\mathbf{x}^{(0)}$. For any $\mathbf{y} \in \mathbb{R}^M$, the sequence

$$\{-P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}^{(i)}) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_0)\}_{i=1}^\infty$$

is pointwise convergent to $-P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}^{(0)}) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_0)$. Again, by the Lebesgue convergence theorem, we have

$$\lim_{i \rightarrow \infty} -\int P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}^{(i)}) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_0) d\mathbf{y} \\ = -\int P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}^{(0)}) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_0) d\mathbf{y} \quad (30)$$

and the result follows. \square

Proposition 5: Let $d(\mathbf{w})$ be defined as in (11). Then $d(\mathbf{w})$ is holomorphic over $\mathbb{A}(\Delta_{\mathbf{w}})$.

Proof: Suppose \mathbb{V} is a closed and bounded subset of $\mathbb{A}(\Delta_{\mathbf{w}})$. As the entries in $\mathbf{b}_{\mathbf{w}}$ and $\nabla_{\mathbf{w}}$ are continuous functions of \mathbf{w} , there exists $\vartheta, \tilde{\kappa}_h, \tilde{\kappa}_l, \tilde{\zeta}_h, \tilde{\zeta}_l, \tilde{\zeta} > 0$ such that for all $\mathbf{w} \in \mathbb{V}$, we have 1) $\|\mathbf{b}_{\mathbf{w}}\| \leq \vartheta$, 2) the eigenvalues of the matrix $\text{re} \nabla_{\mathbf{w}}$ are within the interval $(\tilde{\kappa}_l, \tilde{\kappa}_h)$, 3) $\|\text{im}(\nabla_{\mathbf{w}}) \text{im}(\mathbf{b}_{\mathbf{w}})\| < \tilde{\zeta}$, and 4) the magnitude of $\det \nabla_{\mathbf{w}}$ is within the interval $(\tilde{\zeta}_l, \tilde{\zeta}_h)$.

For any scalar ϵ and vector $\mathbf{y} \in \mathbb{R}^M$ such that $1 > \epsilon > 0$ and $\|\mathbf{y}\| > \max(\vartheta/\epsilon, (1 + \epsilon)\tilde{\zeta}/\epsilon)$, we have

$$\text{re}((\mathbf{y} - \mathbf{b}_{\mathbf{w}})^\top \nabla_{\mathbf{w}} (\mathbf{y} - \mathbf{b}_{\mathbf{w}})) \quad (31)$$

$$= \text{re}(\mathbf{y} - \mathbf{b}_{\mathbf{w}})^\top \text{re}(\nabla_{\mathbf{w}}) \text{re}(\mathbf{y} - \mathbf{b}_{\mathbf{w}}) \\ - \text{im}(\mathbf{b}_{\mathbf{w}})^\top \text{re}(\nabla_{\mathbf{w}}) \text{im}(\mathbf{b}_{\mathbf{w}}) \\ - 2\text{re}(\mathbf{y} - \mathbf{b}_{\mathbf{w}})^\top \text{im}(\nabla_{\mathbf{w}}) \text{im}(\mathbf{b}_{\mathbf{w}}) \quad (32)$$

$$\geq \tilde{\kappa}_l \|\text{re}(\mathbf{y} - \mathbf{b}_{\mathbf{w}})\|^2 - \tilde{\kappa}_h \|\text{im}(\mathbf{b}_{\mathbf{w}})\|^2 - 2\|\text{re}(\mathbf{y} - \mathbf{b}_{\mathbf{w}})\| \tilde{\zeta} \quad (33)$$

$$\geq \|\mathbf{y}\|^2 [\tilde{\kappa}_l(1 - \epsilon)^2 - \tilde{\kappa}_h \epsilon^2 - 2\epsilon]. \quad (34)$$

Let $\tilde{\beta}$ be $(\tilde{\kappa}_l(1 - \epsilon)^2 - \tilde{\kappa}_h \epsilon^2 - 2\epsilon)/2$. By choosing ϵ small enough, $\tilde{\beta}$ is positive. Let $q(\mathbf{y}; \mathbf{w})$ be

$$\exp(-\frac{1}{2}(\mathbf{y} - \mathbf{b}_{\mathbf{w}})^\top \nabla_{\mathbf{w}} (\mathbf{y} - \mathbf{b}_{\mathbf{w}})) / \sqrt{(2\pi)^M \det \Delta_{\mathbf{w}}}$$

and $f(\mathbf{y}, \mathbf{w})$ be $-q(\mathbf{y}, \mathbf{w}) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_0)$. Thus, for all

$$\|\mathbf{y}\| > \max(\vartheta/\epsilon, (1 + \epsilon)\tilde{\zeta}/\epsilon)$$

we have

$$|q(\mathbf{y}; \mathbf{w})| \leq \exp(-\tilde{\beta}\|\mathbf{y}\|^2) / \sqrt{(2\pi)^M / \tilde{\zeta}_h}.$$

As a result, by picking a positive number $\tilde{\alpha}$ large enough, we have

$$|q(\mathbf{y}; \mathbf{w})| \leq \exp(\tilde{\alpha} - \tilde{\beta}\|\mathbf{y}\|^2), \quad \text{for all } \mathbf{y} \in \mathbb{R}^M$$

and

$$|f(\mathbf{y}, \mathbf{w})| \leq (\alpha + \beta' \|\mathbf{y}\|^2) \exp(\tilde{\alpha} - \tilde{\beta}\|\mathbf{y}\|^2).$$

As a corollary, the integral $\int |f(\mathbf{y}, \mathbf{w})| d\mathbf{y}$ is uniformly convergent over \mathbb{V} .

On the other hand, since the entries in $\Delta_{\mathbf{w}}$ and $\mathbf{b}_{\mathbf{w}}$ are well behaved and $\text{re}(\det \Delta_{\mathbf{w}}) > 0$ for all $\mathbf{w} \in \mathbb{A}(\Delta_{\mathbf{w}})$, the function $f(\mathbf{y}; \mathbf{w})$ is a holomorphic function over $\mathbb{A}(\Delta_{\mathbf{w}})$ for any $\mathbf{y} \in \mathbb{R}^M$. Therefore, according to the differentiation lemma [24], $d(\mathbf{w})$ is holomorphic over $\mathbb{A}(\Delta_{\mathbf{w}})$. \square

Lemma 6: Let $d(\mathbf{w})$ be defined as in (11). If $\mathbf{w}^* \in \mathbb{A}(\Delta_{\mathbf{w}})$ such that $\nabla_{\mathbf{w}^*}$ and $\mathbf{b}_{\mathbf{w}^*}$ are real, then for some $\alpha \geq 0$ and $\beta \geq 0$, we have $d(\mathbf{w}^*) \geq -\alpha + \beta \text{tr}[\Delta_{\mathbf{w}^*}]$.

Proof: Since $\nabla_{\mathbf{w}^*}$ is real and $\mathbf{w}^* \in \mathbb{A}(\Delta_{\mathbf{w}})$, the matrix $\Delta_{\mathbf{w}^*}$ is symmetric and positive definite. By Lemma 3, we have $-\log P_{\mathbf{Y}}(\mathbf{y}; \mu_o) \geq -\alpha + \beta \|\mathbf{y}\|^2$. Let

$$C(\mathbf{y}, \mathbf{w}) = \frac{\exp(-\frac{1}{2}(\mathbf{y} - \mathbf{b}_{\mathbf{w}^*})^\top \nabla_{\mathbf{w}^*} (\mathbf{y} - \mathbf{b}_{\mathbf{w}^*}))}{\sqrt{(2\pi)^M \det \Delta_{\mathbf{w}^*}}}. \quad (35)$$

Then, we have

$$d(\mathbf{w}^*) = - \int C(\mathbf{y}, \mathbf{w}) \log P_{\mathbf{Y}}(\mathbf{y}; \mu_o) d\mathbf{y} \quad (36)$$

$$\geq \int C(\mathbf{y}, \mathbf{w}) (-\alpha + \beta \|\mathbf{y}\|^2) d\mathbf{y} \quad (37)$$

$$= -\alpha + \beta (\text{tr}[\Delta_{\mathbf{w}^*}] + \|\mathbf{b}_{\mathbf{w}^*}\|^2) \quad (38)$$

$$\geq -\alpha + \beta \text{tr}[\Delta_{\mathbf{w}^*}]. \quad (39)$$

□

Lemma 7: Let $\psi(\mathbf{y})$ be a Schwartz function [21] defined on \mathbb{R}^M such that its Fourier transform $\hat{\psi}(\mathbf{y})$ is nonzero for all $\mathbf{y} \in \mathbb{R}^M$. Suppose Ψ is a tempered distribution (i.e., a linear and continuous functional defined on the set of Schwartz functions) and $\Psi * \psi$, the convolution between Ψ and ψ , is the zero tempered distribution. Then Ψ and its Fourier transform $\hat{\Psi}$ are the zero tempered distribution. In addition, if there exist two probability measures μ_0 and μ_1 such that $\Psi = \mu_1 - \mu_0$, then $\mu_1 = \mu_0$.

Proof: Let \mathcal{D} be the set of Schwartz functions which have a compact support and ϕ_0 be a function in \mathcal{D} . As $\hat{\psi}(\mathbf{y})$ is nonzero for all $\mathbf{y} \in \mathbb{R}^M$, the function $\phi_1 \triangleq \phi_0 / \hat{\psi}$ (i.e., $\phi_1(\mathbf{y}) = \phi_0(\mathbf{y}) / \hat{\psi}(\mathbf{y})$ for all $\mathbf{y} \in \mathbb{R}^M$) is also in \mathcal{D} .

Since $(\Psi * \psi)$ is the zero tempered distribution, its Fourier transform $\hat{\Psi} \cdot \hat{\psi}$ is again the zero tempered distribution. Consequently, we have⁴

$$0 = \langle \phi_1, \hat{\Psi} \cdot \hat{\psi} \rangle \quad (40)$$

$$= \langle \phi_0 / \hat{\psi}, \hat{\Psi} \cdot \hat{\psi} \rangle \quad (41)$$

$$= \langle \phi_0, \hat{\Psi} \rangle. \quad (42)$$

As \mathcal{D} is dense in the set of Schwartz functions[22], $\langle \phi, \hat{\Psi} \rangle = 0$ for any Schwartz function ϕ . In other words, $\hat{\Psi}$, and consequently, Ψ , are the zero tempered distribution.

Finally, if $\Psi = \mu_1 - \mu_0$ for some probability measures μ_0 and μ_1 , then $\langle \phi, \mu_1 \rangle = \langle \phi, \mu_0 \rangle$ for any $\phi \in \mathcal{D}$. By Riesz representation theorem [25], $\mu_1 = \mu_0$. □

Corollary 9: Let \mathbf{Z} be a Gaussian random variable and $P_{\mathbf{Z}}(\mathbf{z})$ be the corresponding probability density function. Suppose $g(\mathbf{z})$ is a continuous function such that $|g(\mathbf{z})| \leq \alpha + \beta \|\mathbf{z}\|^2$ for some $\alpha, \beta > 0$. If $P_{\mathbf{Z}}(\mathbf{z}) * g(\mathbf{z})$ is the zero function, then $g(\mathbf{z})$ is also the zero function.

Proof: It can be proved easily that 1) $P_{\mathbf{Z}}(\mathbf{z})$ is a Schwartz function such that its Fourier transform is nonzero everywhere, and 2) $g(\mathbf{z})$ is a tempered distribution. Therefore, by Lemma 7, $g(\mathbf{z})$ is the zero tempered distribution, which further implies that $g(\mathbf{z})$ is the zero function. □

Lemma 8: Consider a constant CG channel $\mathbf{Y} = \mathbf{H}\mathbf{X} + \mathbf{Z}$. Let \mathbf{X} be Gaussian distributed with probability measure

⁴For any Schwartz function ψ and tempered distribution Ψ , we denote the function value of Ψ at ψ by $\langle \phi, \Psi \rangle$.

ϕ . Then we can construct a sequence of probability measures $\{\nu_i\}_{i=1}^{\infty}$ such that 1) \mathbb{E}_{ν_i} is closed and bounded, 2) $\mathbf{E}_{\nu_i}[\mathbf{X}^\top \mathbf{X}] \leq \mathbf{E}_{\phi}[\mathbf{X}^\top \mathbf{X}]$, 3) $\nu_i \implies \phi$ in the Lévy metric, and 4) $\lim_{i \rightarrow \infty} H_{\mathbf{Y}}(\nu_i) = H_{\mathbf{Y}}(\phi)$, or equivalently, $\lim_{i \rightarrow \infty} I(\nu_i) = I(\phi)$.

Proof: Assume without loss of generality that $\mathbf{E}_{\phi}[\mathbf{X}\mathbf{X}^\top]$ is positive definite. Hence, ϕ is indeed characterized by a probability density function, which is denoted by $\phi(\mathbf{x})$ for simplicity. For any positive integer i , let $\mathbb{V}_i = \{\mathbf{x} \in \mathbb{R}^N : \mathbf{x}^\top \mathbf{x} \leq i\}$, and $\nu_i(\mathbf{x})$ be the following ‘‘truncated’’ Gaussian probability density function

$$\nu_i(\mathbf{x}) = \begin{cases} \alpha_i \phi(\mathbf{x}), & \text{if } \mathbf{x} \in \mathbb{V}_i \\ 0, & \text{otherwise} \end{cases} \quad (43)$$

where $\alpha_i > 1$ is a normalizing constant. Clearly, conditions 1) and 2) are satisfied, and $\lim_{i \rightarrow \infty} \alpha_i = 1$. For any continuous and bounded function $f(\mathbf{x})$ defined on \mathbb{R}^N , it is apparent that

$$\lim_{i \rightarrow \infty} \int f(\mathbf{x}) \nu_i(\mathbf{x}) d\mathbf{x} = \int f(\mathbf{x}) \phi(\mathbf{x}) d\mathbf{x}.$$

Hence, by Lemma 2, $\nu_i \implies \phi$. It remains to prove that $\lim_{i \rightarrow \infty} H_{\mathbf{Y}}(\nu_i) = H_{\mathbf{Y}}(\phi)$.

Let $b_i(\mathbf{y}) = P_{\mathbf{Y}}(\mathbf{y}; \nu_i) / \alpha_i$. Since $\nu_i \implies \phi$

$$\lim_{i \rightarrow \infty} b_i(\mathbf{y}) = P_{\mathbf{Y}}(\mathbf{y}; \phi).$$

Also, as $\nu_i(\mathbf{x}) / \alpha_i \leq \phi(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{R}^N$, $b_i(\mathbf{y}) \leq P_{\mathbf{Y}}(\mathbf{y}; \phi)$ for all $\mathbf{y} \in \mathbb{R}^M$.

Let $\vartheta > 0$ such that $P_{\mathbf{Y}}(\mathbf{y}; \phi) < 1/e$ for all $\|\mathbf{y}\| \geq \vartheta$. As $-t \log t$ is increasing and nonnegative on the interval $[0, 1/e]$, $|-b_i(\mathbf{y}) \log b_i(\mathbf{y})| \leq |-P_{\mathbf{Y}}(\mathbf{y}; \phi) \log P_{\mathbf{Y}}(\mathbf{y}; \phi)|$ for all $\|\mathbf{y}\| \geq \vartheta$. By the Lebesgue convergence theorem

$$\begin{aligned} \lim_{i \rightarrow \infty} \int_{\|\mathbf{y}\| \geq \vartheta} -b_i(\mathbf{y}) \log b_i(\mathbf{y}) d\mathbf{y} \\ = \int_{\|\mathbf{y}\| \geq \vartheta} -P_{\mathbf{Y}}(\mathbf{y}; \phi) \log P_{\mathbf{Y}}(\mathbf{y}; \phi) d\mathbf{y}. \end{aligned}$$

Similarly, as the function $|-b_i(\mathbf{y}) \log b_i(\mathbf{y})|$ is bounded for $\|\mathbf{y}\| \leq \vartheta$, the limit

$$\lim_{i \rightarrow \infty} \int_{\|\mathbf{y}\| \leq \vartheta} -b_i(\mathbf{y}) \log b_i(\mathbf{y}) d\mathbf{y}$$

is equal to $\int_{\|\mathbf{y}\| \leq \vartheta} -P_{\mathbf{Y}}(\mathbf{y}; \phi) \log P_{\mathbf{Y}}(\mathbf{y}; \phi) d\mathbf{y}$. Consequently, as i goes to infinity, $-\int b_i(\mathbf{y}) \log b_i(\mathbf{y}) d\mathbf{y}$ converges to

$$-\int P_{\mathbf{Y}}(\mathbf{y}; \phi) \log P_{\mathbf{Y}}(\mathbf{y}; \phi) d\mathbf{y}$$

and, as a result, the limit

$$\lim_{i \rightarrow \infty} - \int P_{\mathbf{Y}}(\mathbf{y}; \nu_i) \log P_{\mathbf{Y}}(\mathbf{y}; \nu_i) d\mathbf{y}$$

is also equal to

$$-\int P_{\mathbf{Y}}(\mathbf{y}; \phi) \log P_{\mathbf{Y}}(\mathbf{y}; \phi) d\mathbf{y}$$

i.e., $\lim_{i \rightarrow \infty} H_{\mathbf{Y}}(\nu_i) = H_{\mathbf{Y}}(\phi)$ and the lemma is proved. □

Proposition 6: Consider a sequence of constant CG channels $\{(\Delta, \mathbf{H}\mathbf{x}, \{\mathbf{x}^\top \mathbf{x} - \gamma\}, \mathbb{S}_j)\}_{j=1}^{\infty}$ such that for any closed and bounded subset \mathbb{V} of \mathbb{R}^N , $\mathbb{V} \subseteq \mathbb{S}_j$ for sufficiently large j .

Let μ_j and Φ be the capacity-achieving measures for the channels $(\Delta, \mathbf{H}\mathbf{x}, \{\mathbf{x}^\top \mathbf{x} - \gamma\}, \mathbb{S}_j)$ and $(\Delta, \mathbf{H}\mathbf{x}, \{\mathbf{x}^\top \mathbf{x} - \gamma\}, \mathbb{R}^N)$, respectively. Then $\lim_{j \rightarrow \infty} I(\mu_j) = I(\Phi)$, or equivalently, $\lim_{j \rightarrow \infty} H_{\mathbf{Y}}(\mu_j) = H_{\mathbf{Y}}(\Phi)$.

Proof: By Lemma 8, we can construct a sequence $\{\nu_i\}_{i=1}^{\infty}$ such that 1) \mathbb{E}_{ν_i} is closed and bounded, 2) $\mathbf{E}_{\nu_i}[\mathbf{X}^\top \mathbf{X}] \leq \mathbf{E}_{\Phi}[\mathbf{X}^\top \mathbf{X}]$, and 3) $\lim_{i \rightarrow \infty} I(\nu_i) = I(\Phi)$. It is obvious that for sufficiently large j , $\mathbb{E}_{\nu_i} \subseteq \mathbb{S}_j$. Hence, $I(\nu_i) \leq I(\mu_j) \leq I(\Phi)$, and thus, the proposition follows. \square

Consider a constant CG channel with a diagonal square channel matrix \mathbf{D} with diagonal entries χ_n . Let Φ be the capacity-achieving probability measure for the channel $(\mathbf{I}_N, \mathbf{D}, \{\mathbf{x}^\top \mathbf{x} - \gamma\}, \mathbb{R}^N)$, and $P_n \triangleq \mathbf{E}_{\Phi}[X_n^2]$ for $n \in \mathbb{N}$.

Lemma 9: Let μ be an input probability measure such that $\mathbf{E}_{\mu}[\mathbf{X}^\top \mathbf{X}] \leq \gamma$. For the above channel $(\mathbf{I}_N, \mathbf{D}, \{\mathbf{x}^\top \mathbf{x} - \gamma\}, \mathbb{R}^N)$, we have

$$D(P_{\mathbf{Y}}(\mathbf{y}; \mu) \| P_{\mathbf{Y}}(\mathbf{y}; \Phi)) \leq H_{\mathbf{Y}}(\Phi) - H_{\mathbf{Y}}(\mu)$$

where $D(P_{\mathbf{Y}}(\mathbf{y}; \mu) \| P_{\mathbf{Y}}(\mathbf{y}; \Phi))$ is the Kullback–Leibler distance between $P_{\mathbf{Y}}(\mathbf{y}; \mu)$ and $P_{\mathbf{Y}}(\mathbf{y}; \Phi)$. Consequently

$$\int |P_{\mathbf{Y}}(\mathbf{y}; \mu) - P_{\mathbf{Y}}(\mathbf{y}; \Phi)| d\mathbf{y} \leq \sqrt{2H_{\mathbf{Y}}(\Phi) - 2H_{\mathbf{Y}}(\mu)}$$

according to Pinsker's inequality [19].

Proof: According to the water-filling algorithm, there exists $\vartheta > 0$ such that $\chi_n^2/(\chi_n^2 P_n + 1) \leq 1/\vartheta$ and $P_n \chi_n^2/(\chi_n^2 P_n + 1) = P_n/\vartheta$ for all $n \in \mathbb{N}$. Let Δ^* be the covariance matrix of \mathbf{Y} when the input probability measure is Φ . Let $C = -H_{\mathbf{Y}}(\mu) + \frac{1}{2} \log[(2\pi)^N \det \Delta^*]$. We have

$$D(P_{\mathbf{Y}}(\mathbf{y}; \mu) \| P_{\mathbf{Y}}(\mathbf{y}; \Phi)) = C + \frac{1}{2} \sum_{n \in \mathbb{N}} \frac{\chi_n^2 \mathbf{E}_{\mu}[X_n^2] + 1}{\chi_n^2 P_n + 1} \quad (44)$$

$$\leq C + \frac{1}{2} \sum_{n \in \mathbb{N}} \frac{\mathbf{E}_{\mu}[X_n^2]}{\vartheta} + \frac{1}{2} \sum_{n \in \mathbb{N}} \frac{1}{\chi_n^2 P_n + 1} \quad (45)$$

$$\leq C + \frac{\gamma}{2\vartheta} + \frac{1}{2} \sum_{n \in \mathbb{N}} \frac{1}{\chi_n^2 P_n + 1} \quad (46)$$

$$= H_{\mathbf{Y}}(\Phi) - H_{\mathbf{Y}}(\mu). \quad (47)$$

\square

Lemma 10: Let Γ be the set of all input probability measures μ of \mathbf{X} such that $\mathbf{E}_{\mu}[\mathbf{X}^\top \mathbf{X}] \leq \gamma$. Then Γ is sequentially compact.

Proof: For any $\epsilon > 0$, let \mathbb{K} be the closed and bounded set $\{\mathbf{x} \in \mathbb{R}^N : \mathbf{x}^\top \mathbf{x} \leq \gamma/\epsilon\}$. Hence, if $\mu \in \Gamma$, then $\mu(\mathbb{K}) \geq 1 - \epsilon$. Therefore, Γ is tight. Let $\{\mu_i\}_{i=1}^{\infty}$ be a convergent sequence in Γ such that it converges to μ_0 . By Lemma 2, we have

$$\int \mathbf{x}^\top \mathbf{x} d\mu_0 \leq \liminf_{i \rightarrow \infty} \int \mathbf{x}^\top \mathbf{x} d\mu_i \leq \gamma.$$

Hence, $\mu_0 \in \Gamma$, and thus, Γ is closed. The sequential compactness of Γ then follows from Lemma 2. \square

ACKNOWLEDGMENT

The authors would like to thank the referees and the Associate Editor for their detailed comments and suggestions, which have significantly improved the readability of the paper.

REFERENCES

- [1] J. G. Smith, "The information capacity of amplitude and variance-constrained scalar Gaussian channels," *Inf. Contr.*, vol. 18, pp. 203–219, 1971.
- [2] S. Shamai (Shitz), "Capacity of a pulse amplitude modulated direct detection photon channel," *Proc. Inst. Elec. Eng.*, pt. 1, vol. 137, no. 6, pp. 424–430, Dec. 1990.
- [3] S. Shamai (Shitz) and I. Bar-David, "The capacity of average and peak-power-limited quadrature Gaussian channels," *IEEE Trans. Inf. Theory*, vol. 41, no. 4, pp. 1060–1071, Jul. 1995.
- [4] I. C. Abou-Faycal, M. D. Trott, and S. Shamai (Shitz), "The capacity of discrete-time memoryless Rayleigh-fading channels," *IEEE Trans. Inf. Theory*, vol. 47, no. 4, pp. 1290–1301, May 2001.
- [5] M. Katz and S. Shamai (Shitz), "On the capacity-achieving distribution of the discrete-time noncoherent additive white noise channel," in *Proc. 2002 IEEE Int. Symp. Information Theory*, Lausanne, Switzerland, Jun./Jul. 2002, p. 165.
- [6] W. Oetli, "The capacity-achieving input distribution for some amplitude limited channels with additive noise," *IEEE Trans. Inf. Theory*, vol. IT-20, no. 3, pp. 372–374, May 1974.
- [7] A. Das, "Capacity-achieving distributions for non-Gaussian additive noise channels," in *Proc. 2000 IEEE Int. Symp. Information Theory*, Sorrento, Italy, Jun. 2000, p. 432.
- [8] R. Palanki, "On the capacity achieving distributions of some fading channels," in *Proc. 40th Annu. Allerton Conf. Communication, Control, and Computing*, Monticello, IL, Oct. 2002, pp. 337–346.
- [9] M. C. Gursoy, H. V. Poor, and S. Verdú, "On the capacity-achieving distribution of the noncoherent Rician fading channel," in *Proc. 2003 Canadian Workshop on Information Theory*, Waterloo, ON, Canada, May 2003, pp. 24–27.
- [10] R. Nuriyev and A. Anastasopoulos, "Capacity characterization for the noncoherent block independent AWGN channel," in *Proc. 2003 IEEE Int. Symp. Information Theory*, Yokohama, Japan, Jun./Jul. 2003, p. 373.
- [11] P. Hou, B. J. Belzer, and T. R. Fischer, "On the capacity of the partially coherent additive white Gaussian noise channel," in *Proc. 2003 IEEE Int. Symp. Information Theory*, Yokohama, Japan, Jun./Jul. 2003, p. 372.
- [12] J. Huang and S. Meyn, "Characterization and computation of optimal distributions for channel coding," *IEEE Trans. Inf. Theory*, submitted for publication.
- [13] A. Tchamkerten, "On the discreteness of capacity-achieving distributions," *IEEE Trans. Inf. Theory*, vol. 50, no. 11, pp. 2773–2778, Nov. 2004.
- [14] S. Hranilovic and F. R. Kschischang, "Optical intensity-modulated direct detection channels: Signal space and Lattice codes," *IEEE Trans. Inf. Theory*, vol. 49, no. 6, pp. 1385–1399, Jun. 2003.
- [15] T. L. Marzetta and B. M. Hochwald, "Capacity of a mobile multiple-antenna communication link in Rayleigh fading," *IEEE Trans. Inf. Theory*, vol. 45, no. 1, pp. 139–157, Jan. 1999.
- [16] D. G. Luenberger, *Optimization by Vector Space Methods*. New York: Wiley, 1969.
- [17] R. C. Gunning, *Introduction to Holomorphic Functions of Several Variables*. Pacific Grove, CA: Wadsworth and Brooks Cole, 1990.
- [18] C. E. Shannon, "Communication in the presence of noise," *Proc. IRE*, vol. 37, no. 1, pp. 10–21, Jan. 1949.
- [19] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York: Wiley Interscience, 1991.
- [20] R. M. Range, *Holomorphic Functions and Integral Representations in Several Complex Variables*. New York: Springer-Verlag, 1986.
- [21] W. Rudin, *Functional Analysis*. New York: McGraw-Hill, 1973.
- [22] R. S. Pathak, *A Course in Distribution Theory and Applications*. Boca Raton, FL: CRC Press/Narosa, 2000.
- [23] D. W. Stroock, *Probability Theory, an Analytic View*. New York: Cambridge Univ. Press, 1993.
- [24] S. Lang, *Complex Analysis*. New York: Springer-Verlag, 1999.
- [25] G. B. Folland, *Real Analysis: Modern Techniques and Their Applications*. New York: Wiley, 1984.
- [26] G. N. Iyengar, "Voice channel," in *Proc. 31st Asilomar Conf. Signals, Systems and Computers*, vol. 2, Pacific Grove, CA, 1997, pp. 1354–1358.
- [27] A. Lapidoth and S. Shamai (Shitz), "The Poisson multiple-access channel," *IEEE Trans. Inf. Theory*, vol. 44, no. 2, pp. 488–501, Mar. 1998.
- [28] G. B. Folland, *Fourier Analysis and its Applications*. Pacific Grove, CA: Wadsworth and Brooks Cole Advanced Books and Software, 1992.