

Efficient Approximation of Channel Capacities

Tobias Sutter,^{1,*} David Sutter,^{2,†} Peyman Mohajerin Esfahani,^{1,*} and John Lygeros^{1,*}

¹*Automatic Control Laboratory, ETH Zurich, Switzerland*

²*Institute for Theoretical Physics, ETH Zurich, Switzerland*

We propose an iterative method for approximately computing the capacity of discrete memoryless channels, possibly under additional constraints on the input distribution. Based on duality of convex programming, we derive explicit upper and lower bounds for the capacity. The presented method requires $O(M^2N\sqrt{\log N}/\varepsilon)$ to provide an estimate of the capacity to within ε , where N and M denote the input and output alphabet size; a single iteration has a complexity $O(MN)$. We also show how to approximately compute the capacity of memoryless channels having a bounded continuous input alphabet and a countable output alphabet under some mild assumptions on the decay rate of the channel's tail. It is shown that discrete-time Poisson channels fall into this problem class. As an example, we compute sharp upper and lower bounds for the capacity of a discrete-time Poisson channel with a peak-power input constraint.

1. INTRODUCTION

A discrete memoryless channel (DMC) comprises a finite input alphabet $\mathcal{X} = \{1, 2, \dots, N\}$, a finite output alphabet $\mathcal{Y} = \{1, 2, \dots, M\}$, and a conditional probability mass function expressing the probability of observing the output symbol y given the input symbol x , denoted by $W(y|x)$. In his seminal 1948 paper [1], Shannon proved that the channel capacity for a DMC is

$$C(W) = \max_{p \in \Delta_N} I(p, W), \quad (1)$$

where $\Delta_N := \{x \in \mathbb{R}^N : x \geq 0, \sum_{i=1}^N x_i = 1\}$ denotes the N -simplex and $I(p, W) := \sum_{x \in \mathcal{X}} p(x) D(W(\cdot|x) || (pW)(\cdot))$ the mutual information. $W(y|x) = \mathbb{P}[Y = y | X = x]$ describes the channel law and $(pW)(\cdot)$ is the probability distribution of the channel output induced by p and W , i.e., $(pW)(y) := \sum_{x \in \mathcal{X}} p(x)W(y|x)$. $D(\cdot || \cdot)$ denotes the relative entropy that is defined as $D(W(\cdot|x) || (pW)(\cdot)) := \sum_{y \in \mathcal{Y}} W(y|x) \log \left(\frac{W(y|x)}{(pW)(y)} \right)$. Shannon also showed that in case of an additional average cost constraint on the input distribution of the form $\mathbb{E}[s(X)] \leq S$, where $s : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ denotes a cost function and $S \geq 0$, the capacity is given by

$$C_S(W) = \begin{cases} \max_p I(p, W) \\ \text{s. t. } \mathbb{E}[s(X)] \leq S \\ p \in \Delta_N. \end{cases} \quad (2)$$

For a few DMCs it is known that the capacity can be computed analytically, however in general there is no closed-form solution. It is therefore of interest to have an algorithm that solves (2) in a reasonable amount of time. Since for a fixed channel the mutual information is a concave function in p , the optimization problem (2) is a finite dimensional convex optimization problem. Solving

The material in this paper was presented in part at the IEEE International Symposium on Information Theory, June 2014.

* {sutter, mohajerin, lygeros}@control.ee.ethz.ch

† suttetdav@phys.ethz.ch

(2) with convex programming solvers, however, turned out to be computationally inefficient even for small alphabet sizes [2].

Shannon's formula for the capacity of a DMC generalizes to the case of memoryless channels with continuous input and output alphabets, i.e. $\mathcal{X} = \mathcal{Y} = \mathbb{R}$. However, when considering such channels, it is essential to introduce additional constraints on the channel input to obtain physically meaningful results, more details can be found in [3, Chapter 7]. In addition to average cost type constraints, peak-power constraints are also often considered. A peak-power constraint demands that $X \in \mathbb{A}$ for some compact set $\mathbb{A} \subset \mathcal{X}$ with probability one. For such a setup, i.e., having average and peak-power constraints, the capacity is given by

$$C_{\mathbb{A},S}(W) = \begin{cases} \sup_p I(p, W) \\ \text{s. t. } \mathbb{E}[s(X)] \leq S \\ p \in \mathcal{P}(\mathbb{A}), \end{cases} \quad (3)$$

where $\mathcal{P}(\mathbb{A})$ denotes the set of all probability distributions on the Borel σ -algebra $\mathcal{B}(\mathbb{A})$ and the mutual information is defined as $I(p, W) := \int_{\mathbb{A}} D(W(\cdot|x) \| (pW)(\cdot)) p(dx)$. The channel is described by a transition density defined by $\mathbb{P}[Y \in dy | X = x] = W(y|x)dy$ and $(pW)(\cdot)$ is the probability distribution of the channel output induced by p and W which is given by $(pW)(y) := \int_{\mathbb{A}} W(y|x)p(dx)$ and the relative entropy that is defined as $D(W(\cdot|x) \| (pW)(\cdot)) := \int_{\mathcal{Y}} W(y|x) \log \left(\frac{W(y|x)}{(pW)(y)} \right) dy$. The optimization problem (3) is an infinite dimensional convex optimization problem and as such in general computationally intractable (NP-hard).

Previous Work and Contributions.— Historically one of the first attempts to numerically solve (2) is the so-called *Blahut-Arimoto algorithm* [2, 4], that exploits the special structure of the mutual information and approximates iteratively the capacity of any DMC. Each iteration step has a computational complexity $O(MN)$. It was shown that this algorithm, in case of no additional input constraints has an *a priori* error bound of the form $|C(W) - C_{\text{approx}}^{(n)}(W)| \leq O\left(\frac{\log(N)}{n}\right)$, where n denotes the number of iterations [4, Corollary 1]. Hence, the overall computational complexity of finding an additive ε -solution is given by $O\left(\frac{MN \log(N)}{\varepsilon}\right)$. As such the computational cost required for an acceptable accuracy for channels with large input alphabets can be considerable. This undesirable property together with the complexity per iteration prevents the algorithm from being useful for a large class of channels, e.g., a Rayleigh channel with a discrete input alphabet [5]. There have been several improvements of the Blahut-Arimoto algorithm [6–8], which achieve a better convergence for certain channels. However, since they all rely on the original Blahut-Arimoto algorithm they inherit its overall computational complexity as well as its complexity per iteration step. Therefore, even with improved Blahut-Arimoto algorithms, approximating the capacity for channels having large input alphabets remains computationally expensive. Based on sequential Monte-Carlo integration methods (a.k.a. particle filters), the Blahut-Arimoto algorithm has been extended to memoryless channels with continuous input and output alphabets [9–12]. As shown in several examples, this approach seems to be powerful in practice, however a rate of convergence has not been proven.

Another recent approach towards approximating (2) is presented in [13] by Mung and Boyd, where they introduce an efficient method to derive upper bounds on the channel capacity problem, based on geometric programming. Huang and Meyn [14] developed a different approach based on cutting plane methods, where the mutual information is iteratively approximated by linear functionals and in each iteration step, a finite dimensional linear program is solved. It has been shown that this method converges to the optimal value, however no rate of convergence is provided.

In this article, we present a new approach to solve (2) that is based on its dual formulation. It turns out that the dual problem of (2) has a particular structure that allows us to apply Nesterov's

smoothing method [15]. In the absence of input cost constraints, this leads to an a priori error bound of the order $|C(W) - C_{\text{approx}}^{(n)}(W)| \leq O\left(\frac{M\sqrt{\log(N)}}{n}\right)$, where n denotes the number of iterations and each iteration step has a computational complexity of $O(NM)$. Thus, the overall computational complexity of finding an ε -solution is given by $O\left(\frac{M^2N\sqrt{\log(N)}}{\varepsilon}\right)$. In particular for large input alphabets our method has a computational advantage over the Blahut-Arimoto algorithm. In addition the novel method provides primal and dual optimizers leading to an *a posteriori* error which is often much smaller than the a priori error.

Due to the favorable structure of the capacity problem and its dual formulation, the presented method can be extended to approximate the capacity of memoryless channels having a bounded continuous input alphabet and a countable output alphabet, under some assumptions on the tail of $W(\cdot|x)$, i.e., problem (3) is addressed for a countable output alphabet. As a concrete example, this is demonstrated on the discrete-time Poisson channel with a peak-power constraint. To the best of our knowledge, for this scenario up to now only lower bounds exist [16].

Structure.— Section 2 introduces our method for approximating the channel capacity for DMCs. We provide a priori and a posteriori bounds for the approximation error and present two numerical examples that illustrate its computational performance compared to the Blahut-Arimoto algorithm. In Section 3, we generalize the approximation scheme to channels having bounded continuous input alphabets and countable output alphabets. We then show how the presented results can be used to compute the capacity of discrete-time Poisson channels under a peak-power constraint and possibly average-power constraints on the input. We conclude in Section 4 with a summary and potential subjects of further research. In the interest of readability, some of the technical proofs and details are given in the appendices.

Notation.— The logarithm with basis 2 is denoted by $\log(\cdot)$ and the natural logarithm by $\ln(\cdot)$. In Section 2 we consider DMCs with a finite input alphabet $\mathcal{X} = \{1, 2, \dots, N\}$ and a finite output alphabet $\mathcal{Y} = \{1, 2, \dots, M\}$. The channel law is summarized in a matrix $W \in \mathbb{R}^{N \times M}$, where $W_{ij} := \mathbb{P}[Y = j|X = i] = W(j|i)$. We define the standard n -simplex as $\Delta_d := \left\{x \in \mathbb{R}^d : x \geq 0, \sum_{i=1}^d x_i = 1\right\}$. The input and output probability mass functions are denoted by the vectors $p \in \Delta_N$ and $q \in \Delta_M$. The input cost constraint can be written as $\mathbb{E}[s(X)] = p^\top s \leq S$, where $s \in \mathbb{R}_{\geq 0}^N$ denotes the cost vector and $S \in \mathbb{R}_{\geq 0}$ is the given total cost. The binary entropy function is denoted by $H_b(\alpha) := -\alpha \log(\alpha) - (1 - \alpha) \log(1 - \alpha)$, for $\alpha \in [0, 1]$. For a probability mass function $p \in \Delta_N$ we denote the entropy by $H(p) := \sum_{i=1}^N -p_i \log(p_i)$. It is convenient to introduce an additional variable for the conditional entropy of Y given $\{X = i\}$ as $r \in \mathbb{R}^N$, where $r_i = -\sum_{j=1}^M W_{ij} \log(W_{ij})$. For a probability density p supported at a measurable set $B \subset \mathbb{R}$ we denote the differential entropy by $h(p) = -\int_B p(x) \log(p(x)) dx$. For two vectors $x, y \in \mathbb{R}^n$, we denote the canonical inner product by $\langle x, y \rangle := x^\top y$. We denote the maximum (resp. minimum) between a and b by $a \vee b$ (resp. $a \wedge b$). For $\mathbb{A} \subset \mathbb{R}$ and $1 \leq p \leq \infty$, let $L^p(\mathbb{A})$ denote the space of L^p -functions on the measure space $(\mathbb{A}, \mathcal{B}(\mathbb{A}), dx)$, where $\mathcal{B}(\mathbb{A})$ denotes the Borel σ -algebra and dx the Lebesgue measure. The capacity of a channel W is denoted by $C(W)$. For the channel law matrix $W \in \mathbb{R}^{N \times M}$ we consider the norm $\|W\| := \max_{\lambda \in \mathbb{R}^M, p \in \mathbb{R}^N} \{\langle W\lambda, p \rangle : \|\lambda\|_2 = 1, \|p\|_1 = 1\}$, and note that an upper bound is given by

$$\|W\| = \max_{\|p\|_1=1} \max_{\|\lambda\|_2=1} \lambda^\top W^\top p \leq \max_{\|p\|_1=1} \|W^\top p\|_2 \leq \max_{\|p\|_1=1} \|W^\top p\|_1 = \max_{\|p\|_1=1} \|p\|_1 = 1. \quad (4)$$

2. DISCRETE MEMORYLESS CHANNEL

To keep notation simple we consider a single average-input cost constraint as the extension to multiple average-input cost constraints is straightforward. In a first step, we introduce the output distribution $q \in \Delta_M$ as an additional decision variable, as done in [13, 17, 18] and note that the mutual information $I(X; Y)$ is equal to $H(Y) - H(Y|X)$.

Lemma 2.1. *Let $\mathcal{F} := \arg \max_{p \in \Delta_N} I(p, W)$ and $S_{\max} := \min_{p \in \mathcal{F}} s^\top p$. If $S \geq S_{\max}$ the optimization problem (2) has the same optimal value as*

$$\text{P : } \begin{cases} \max_{p, q} & -r^\top p + H(q) \\ \text{s. t.} & W^\top p = q \\ & p \in \Delta_N, q \in \Delta_M. \end{cases} \quad (5)$$

If $S < S_{\max}$ the optimization problem (2) has the same optimal value as

$$\text{P : } \begin{cases} \max_{p, q} & -r^\top p + H(q) \\ \text{s. t.} & W^\top p = q \\ & s^\top p = S \\ & p \in \Delta_N, q \in \Delta_M. \end{cases} \quad (6)$$

Proof. The proof can be found in Appendix A1. □

Note that we later add an assumption on our channel (Assumption 2.3) that guarantees uniqueness of the optimizer maximizing the mutual information, i.e., \mathcal{F} is a singleton. In this case the optimizer to (6) (resp. (5)) is also feasible for the original problem (2). Computing S_{\max} is straightforward once \mathcal{F} is known. The singleton \mathcal{F} can be seen as the maximizer of a channel capacity problem with no additional input cost constraint and can as such be computed with the scheme we present in this article.

For the rest of the section we restrict attention to (6), since the less constrained problem (5) can be solved in a similar, more direct way. We tackle this optimization problem through its Lagrangian dual problem. The dual function turns out to be a non-smooth function. As such, it is known that the efficiency estimate of a black-box first-order method is of the order $O(\frac{1}{\varepsilon^2})$ if no specific problem structure is used, where ε is the desired absolute accuracy of the approximate solution in function value [19]. We show, however, that P has a certain structure that allows us to use Nesterov's approach for approximating non-smooth problems with smooth ones [15] leading to an efficiency estimate of the order $O(\frac{1}{\varepsilon})$. This, together with the low complexity of each iteration step in the approximation scheme leads to a numerical method for the channel capacity problem that has a very attractive computational complexity.

A. Preliminaries

Some preliminaries are needed in order to present our capacity approximation method. We begin by recalling Nesterov's seminal work [15] in the context of structural convex optimization, which is our main tool in the proposed capacity approximation scheme.

Nesterov's smoothing approach [15]

Consider finite-dimensional real vector spaces E_i endowed with a norm $\|\cdot\|_i$ and denote its dual space by E_i^* for $i = 1, 2$. Each dual pair of vector spaces comes with a bilinear form $\langle \cdot, \cdot \rangle_i : E_i^* \times E_i \rightarrow \mathbb{R}$. For a linear operator $A : E_1 \rightarrow E_2^*$ the operator norm is defined as $\|A\|_{1,2} = \max_{x,u} \{\langle Ax, u \rangle_2 : \|x\|_1 = 1, \|u\|_2 = 1\}$. We are interested in the following optimization problem

$$\min_x \{f(x) : x \in Q_1\}, \quad (7)$$

where $Q_1 \subset E_1$ is a compact convex set and f is a continuous convex function on Q_1 . We assume that the objective function has the following structure

$$f(x) = \hat{f}(x) + \max_u \{\langle Ax, u \rangle_2 - \hat{\phi}(u) : u \in Q_2\}, \quad (8)$$

where $Q_2 \subset E_2$ is a compact convex set, \hat{f} is a continuously differentiable convex function whose gradient is Lipschitz continuous with constant L on Q_1 and $\hat{\phi}$ is a continuous convex function on Q_2 . It is assumed that $\hat{\phi}$ and Q_2 are simple enough such that the maximization in (8) is available in closed form. The dual program to (7) can be given as

$$\max_u \left\{ -\hat{\phi}(u) + \min_x \{\langle Ax, u \rangle_2 + \hat{f}(x) : x \in Q_1\} : u \in Q_2 \right\}. \quad (9)$$

The main difficulty in solving (7) efficiently is its non-smooth objective function. Without using any specific problem structure the complexity for subgradient-type methods is $O\left(\frac{1}{\varepsilon^2}\right)$, where ε is the desired absolute accuracy of the approximate solution in function value. Nesterov's work suggests that when approximating problems with the particular structure (8) by smooth ones, a solution to the non-smooth problem can be constructed with complexity in order of $O\left(\frac{1}{\varepsilon}\right)$. In addition, Nesterov shows that when solving the smooth problem, a solution to the dual problem (9) can be obtained, and as such an a posteriori statement about the duality gap is available that often is significantly tighter than the $O\left(\frac{1}{\varepsilon}\right)$ complexity bound. Consider the the smooth approximation to problem (7) given by

$$\min_x \{f_\nu(x) : x \in Q_1\}, \quad (10)$$

where $\nu > 0$ and the objective function is given by

$$f_\nu(x) = \hat{f}(x) + \max_u \{\langle Ax, u \rangle_2 - \hat{\phi}(u) - \nu d(u) : u \in Q_2\}, \quad (11)$$

where $d : Q_2 \rightarrow \mathbb{R}$ is continuous and strongly convex with convexity parameter σ . It can be shown that f_ν has a Lipschitz continuous gradient with Lipschitz constant $L + \frac{\|A\|_{1,2}^2}{\nu\sigma}$ [15, Theorem 1]. In this light, the optimization problem (10) belongs to a class of problems that can be solved in $O\left(\frac{1}{\sqrt{\varepsilon}}\right)$ using a fast gradient method. The result [15, Theorem 3] explicitly details how, having solved the smooth problem (10), primal and dual solutions to the non-smooth problems (7) and (9) can be obtained and how good they are.

Entropy maximization

As a second preliminary result for some $c \in \mathbb{R}^N$ we consider the following optimization problem, that, if feasible, has an analytical solution

$$\begin{cases} \max_p & H(p) - c^\top p \\ \text{s.t.} & s^\top p = S \\ & p \in \Delta_N. \end{cases} \quad (12)$$

Lemma 2.2. Let $p^* = [p_1^*, \dots, p_N^*]$ with $p_i^* = 2^{\mu_1 - c_i + \mu_2 s_i}$, where μ_1 and μ_2 are chosen such that p^* satisfies the constraints in (12). Then p^* uniquely solves (12).

Proof. See Appendix A 2. □

B. Capacity Approximation Scheme

In the following we focus on the input constrained channel capacity problem (6) and the scenario of no input constraints (5) is discussed as a special case within this section. Consider the convex optimization problem (6), whose optimal value, according to Lemma 2.1 is the capacity $C_S(W)$. The Lagrange dual program to (6) is

$$D : \begin{cases} \min_{\lambda} & G(\lambda) + F(\lambda) \\ \text{s.t.} & \lambda \in \mathbb{R}^M, \end{cases} \quad (13)$$

where $F, G : \mathbb{R}^M \rightarrow \mathbb{R}$ are given by

$$G(\lambda) = \begin{cases} \max_p & -r^\top p + \lambda^\top W^\top p \\ \text{s.t.} & s^\top p = S \\ & p \in \Delta_N \end{cases} \quad \text{and} \quad F(\lambda) = \begin{cases} \max_q & H(q) - \lambda^\top q \\ \text{s.t.} & q \in \Delta_M. \end{cases} \quad (14)$$

Note that since the coupling constraint $W^\top p = q$ in the primal program (6) is affine, the set of optimal solutions to the dual program (13) is nonempty [20, Proposition 5.3.1] and as such the optimum is attained. It can be seen that the dual program (13) structurally resembles the problem (7) with (8), without a bounded feasible set, however. To ensure that the set of dual optimizers is compact, we need to impose the following assumption on the channel matrix W , that we will maintain for the remainder of Section 2.

Assumption 2.3. $\gamma := \min_{i,j} W_{ij} > 0$

Assumption 2.3 excludes situations where the channel matrix has zero entries. Even though this may seem restrictive at first glance, it holds for a large class of channels. Moreover, in a finite dimensional setting, for a fixed input distribution, the mutual information is well known to be continuous in the channel matrix entries. Therefore, singular cases where the channel matrix contains zero entries can be avoided by slight perturbations of those entries. (This is discussed in more detail in Remark 2.13.) Under Assumption 2.3 for a fixed channel, the mutual information can be seen to be a strictly concave function in the input distribution. Therefore, the capacity achieving input distribution is unique. With Assumption 2.3 one can derive an explicit bound on the norm of the dual optimizers, which is crucial in the subsequent derivation of the main result in this section, namely Theorem 2.9.

Lemma 2.4. Under Assumption 2.3, the dual program (13) is equivalent to

$$\begin{cases} \min_{\lambda} & G(\lambda) + F(\lambda) \\ \text{s.t.} & \lambda \in Q, \end{cases} \quad (15)$$

where $Q := \{\lambda \in \mathbb{R}^M : \|\lambda\|_2 \leq M (\log(\gamma^{-1}) \vee \frac{1}{\ln 2})\}$.

Proof. See Appendix A 3. □

Lemma 2.5. *Strong duality holds between (6) and (13).*

Proof. The proof follows by a standard strong duality result of convex optimization, see [20, Proposition 5.3.1, p. 169]. \square

Note that the optimization problem defining $F(\lambda)$ is of the form given in (12). Hence, according to Lemma 2.2, $F(\lambda)$ has a unique optimizer q^* with components $q_j^* = 2^{\mu - \lambda_j}$, where $\mu \in \mathbb{R}$ needs to be chosen such that $q^* \in \Delta_M$, i.e.,

$$\mu = -\log \left(\sum_{j=1}^M 2^{-\lambda_j} \right).$$

Therefore,

$$F(\lambda) = \sum_{j=1}^M (-q_j^* \log(q_j^*) - \lambda_j q_j^*) = -\sum_{j=1}^M \mu 2^{\mu - \lambda_j} = -\mu 2^\mu \sum_{j=1}^M 2^{-\lambda_j} = \log \left(\sum_{j=1}^M 2^{-\lambda_j} \right). \quad (16)$$

$F(\lambda)$ is a smooth function with gradient

$$(\nabla F(\lambda))_i = \frac{-2^{-\lambda_i}}{\sum_{j=1}^M 2^{-\lambda_j}}. \quad (17)$$

According to [15, Theorem 1] and the fact that the negative entropy is strongly convex with convexity parameter 1 [15, Lemma 3], $\nabla F(\lambda)$ is Lipschitz continuous with Lipschitz constant 1. The main difficulty in solving (15) efficiently is that $G(\cdot)$ is non-smooth. Following Nesterov's smoothing technique [15], we alleviate this difficulty by approximating $G(\cdot)$ by a function with a Lipschitz continuous gradient. This smoothing step is efficient in our case because of the particular structure of (15). Following [15] and (11), consider

$$G_\nu(\lambda) = \begin{cases} \max_p & \lambda^\top W^\top p - r^\top p + \nu H(p) - \nu \log(N) \\ \text{s.t.} & s^\top p = S \\ & p \in \Delta_N, \end{cases} \quad (18)$$

with smoothing parameter $\nu \in \mathbb{R}_{>0}$ and denote by $p_\nu(\lambda)$ the optimizer to (18), which is unique because the objective function is strictly concave. Clearly for any $\lambda \in Q$, $G_\nu(\lambda)$ is a uniform approximation of the non-smooth function $G(\lambda)$, since $G_\nu(\lambda) \leq G(\lambda) \leq G_\nu(\lambda) + \nu \log(N)$. Using Lemma 2.2, the optimizer $p_\nu(\lambda)$ to (18) is analytically given by

$$p_\nu(\lambda, \mu)_i = 2^{\mu_1 + \frac{1}{\nu}(W\lambda - r)_i + \mu_2 s_i}, \quad (19)$$

where $\mu_1, \mu_2 \in \mathbb{R}$ have to be chosen so that $s^\top p_\nu(\lambda, \mu) = S$ and $p_\nu(\lambda, \mu) \in \Delta_N$; for this choice of μ_1, μ_2 we denote the solution by $p_\nu(\lambda)$.

Remark 2.6. In case of no input constraints, the unique optimizer to (18) is given by

$$p_\nu(\lambda)_i = \frac{2^{\frac{1}{\nu}(W\lambda - r)_i}}{\sum_{i=1}^N 2^{\frac{1}{\nu}(W\lambda - r)_i}} \quad \text{for } i = 1, \dots, N,$$

whose straightforward evaluation is numerically difficult for small ν . One can circumvent this problem, however, by following the numerically stable technique that we present in Remark 2.11.

By Dubin's theorem it can be shown that the capacity of a memoryless channel with a discrete output alphabet of size M and input alphabet size $N \geq M$, is achieved by a discrete input distribution with M mass points [3, 21]. Computing the exact positions and weights of this optimal input distribution may be difficult, though it is worth noting that our analytical solution in (19) converges to this optimal input distribution as ν tends to 0.

Remark 2.7 (Additional input constraints). In case of additional input constraints, we need an efficient method to find the coefficients μ_1 and μ_2 in (19). In particular if there are multiple input constraints (leading to multiple μ_i) the efficiency of the method computing them becomes important. Instead of solving a system of nonlinear equations, one can show ([22, Theorem 4.8], [23, p. 257 ff.]) that the coefficients μ_i are the unique maximizers to the following convex optimization problem

$$\max_{\mu \in \mathbb{R}^2} \left\{ y^\top \mu - \sum_{i=1}^N p_\nu(\lambda, \mu)_i \right\}, \quad (20)$$

where $y := (1, S)$. Notice that (20) is an unconstrained maximization of a strictly concave function, whose gradient and Hessian can be directly computed as

$$\begin{pmatrix} y_1 - \ln 2 \sum_{i=1}^N p_\nu(\lambda, \mu)_i \\ y_2 - \ln 2 \sum_{i=1}^N s_i p_\nu(\lambda, \mu)_i \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} -(\ln 2)^2 \sum_{i=1}^N p_\nu(\lambda, \mu)_i & -(\ln 2)^2 \sum_{i=1}^N s_i p_\nu(\lambda, \mu)_i \\ -(\ln 2)^2 \sum_{i=1}^N s_i p_\nu(\lambda, \mu)_i & -(\ln 2)^2 \sum_{i=1}^N s_i^2 p_\nu(\lambda, \mu)_i \end{pmatrix},$$

which allows the use of efficient second-order methods such as Newton's method. This method directly extends to multiple input constraints. Let us point out that Theorem 2.9, quantifying the approximation error of the presented algorithm, is based on the assumption that the maximum entropy solution (19) is available, meaning that one can solve (20) for optimality. In the case of a finite input alphabet this assumption is not restrictive as we have argued that (20) is easy to solve. For a continuous input alphabet, that we shall discuss in the subsequent section, however, finding the maximum entropy solution is numerically difficult as it involves integration problems. Therefore, in Remark 3.12, we comment on how the presented channel capacity algorithm behaves, when having access only to an approximate solution to the mentioned maximum entropy problem.

Finally, we can show that the uniform approximation $G_\nu(\lambda)$ is smooth and has a Lipschitz continuous gradient, with known Lipschitz constant.

Proposition 2.8. $G_\nu(\lambda)$ is well defined and continuously differentiable at any $\lambda \in Q$. Moreover, it is convex and its gradient $\nabla G_\nu(\lambda) = W^\top p_\nu(\lambda)$ is Lipschitz continuous with Lipschitz constant $\frac{1}{\nu}$.

Proof. The proof follows directly from the proof of Theorem 1 and Lemma 3 in [15] together with (4). \square

We consider the smooth, convex optimization problem

$$D_\nu : \begin{cases} \min_{\lambda} & F(\lambda) + G_\nu(\lambda) \\ \text{s.t.} & \lambda \in Q, \end{cases} \quad (21)$$

whose objective function has a Lipschitz continuous gradient with Lipschitz constant $1 + \frac{1}{\nu}$. As such D_ν can be approximated with Nesterov's optimal scheme for smooth optimization [15], which is summarized in Algorithm 1, where $\pi_Q(x)$ denotes the projection operator of the set Q , defined in Lemma 2.4, with $R := M (\log(\gamma^{-1}) \vee \frac{1}{\ln 2})$

$$\pi_Q(x) := \begin{cases} R \frac{x}{\|x\|_2}, & \|x\|_2 > R \\ x, & \text{otherwise.} \end{cases}$$

Algorithm 1: Optimal scheme for smooth optimization

Choose some $\lambda_0 \in Q$

For $k \geq 0$ **do***

Step 1: Compute $\nabla F(x_k) + \nabla G_\nu(x_k)$

Step 2: $y_k = \pi_Q \left(-\frac{1}{L_\nu} (\nabla F(x_k) + \nabla G_\nu(x_k)) + x_k \right)$

Step 3: $z_k = \pi_Q \left(-\frac{1}{L_\nu} \sum_{i=0}^k \frac{i+1}{2} (\nabla F(x_i) + \nabla G_\nu(x_i)) \right)$

Step 4: $x_{k+1} = \frac{2}{k+3} z_k + \frac{k+1}{k+3} y_k$

[*The stopping criterion is explained in Remark 2.10]

The following theorem provides explicit error bounds for the solution provided by Algorithm 1 after n iterations. Define the constants $D_1 := \frac{1}{2}(M \log(\gamma^{-1}) \vee \frac{1}{\ln 2})^2$ and $D_2 := \log(N)$.

Theorem 2.9 ([15]). *Under Assumption 2.3, for $n \in \mathbb{N}$ consider a smoothing parameter*

$$\nu = \nu(n) = \frac{2}{n+1} \sqrt{\frac{D_1}{D_2}}.$$

Then after n iterations of Algorithm 1 we can generate the approximate solutions to the problems (13) and (2), namely,

$$\hat{\lambda} = y_n \in Q \quad \text{and} \quad \hat{p} = \sum_{k=0}^n \frac{2(k+1)}{(n+1)(n+2)} p_\nu(x_k) \in \Delta_N, \quad (22)$$

which satisfy

$$0 \leq F(\hat{\lambda}) + G(\hat{\lambda}) - I(\hat{p}, W) \leq \frac{4}{n+1} \sqrt{D_1 D_2} + \frac{4D_1}{(n+1)^2}. \quad (23)$$

Thus, the complexity of finding an ε -solution to the problems (13) and (2) does not exceed

$$4\sqrt{D_1 D_2} \frac{1}{\varepsilon} + 2\sqrt{\frac{D_1}{\varepsilon}}. \quad (24)$$

Proof. The proof follows along the lines of [15, Theorem 3] and in particular requires Lemma 2.4, Lemma 2.5 and Proposition 2.8. \square

Note that Theorem 2.9 provides an explicit error bound (23), also called a *a priori error*. In addition this theorem gives an approximation to the optimal input distribution (22), i.e., the optimizer of the primal problem. Thus, by comparing the values of the primal and the dual optimization problem, one can also compute an *a posteriori error* which is the difference of the dual and the primal problem, namely $F(\hat{\lambda}) + G(\hat{\lambda}) - I(\hat{p}, W)$.

Remark 2.10 (Stopping criterion of Algorithm 1). There are two immediate approaches to define a stopping criterion for Algorithm 1.

- (i) *A priori stopping criterion:* Choose an a priori error $\varepsilon > 0$. Setting the right hand side of (23) equal to ε defines a number of iterations n_ε required to ensure an ε -close solution.

- (ii) *A posteriori stopping criterion*: Choose an a posteriori error $\varepsilon > 0$. Choose the smoothing parameter $\nu(n_\varepsilon)$ for n_ε as defined above in the a priori stopping criterion. Fix a (small) number of iterations ℓ that are run using Algorithm 1. Compute the a posteriori error $e_\ell := F(\hat{\lambda}) + G(\hat{\lambda}) - I(\hat{p}, \rho)$ according to Theorem 2.9. If $e_\ell \leq \varepsilon$ terminate the algorithm otherwise continue with another ℓ iterations. Continue until the a posteriori error is below ε .

Remark 2.11 (Computational stability). In the special case of no input cost constraints, one can derive an analytical expression for $G_\nu(\lambda)$ and its gradient as

$$\begin{aligned} G_\nu(\lambda) &= \nu \log \left(\sum_{i=1}^N 2^{\frac{1}{\nu}(\mathbf{W}\lambda - r)_i} \right) - \nu \log(N) \\ \nabla G_\nu(\lambda) &= \frac{1}{S(\lambda)} \sum_{i=1}^N 2^{\frac{1}{\nu}(\mathbf{W}\lambda - r)_i} \mathbf{W}_{i,\cdot}, \end{aligned} \quad (25)$$

where $S(\lambda) := \sum_{i=1}^N 2^{\frac{1}{\nu}(\mathbf{W}\lambda - r)_i}$. In order to achieve an ε -precise solution the smoothing factor ν has to be chosen in the order of ε , according to Theorem 2.9. A straightforward computation of $\nabla G_\nu(\lambda)$ via (25) for a small enough ν is numerically difficult. In the light of [15, p. 148], we present a numerically stable technique for computing $\nabla G_\nu(\lambda)$. By considering the functions $\mathbb{R}^M \ni \lambda \mapsto f(\lambda) = \mathbf{W}\lambda - r \in \mathbb{R}^N$ and $\mathbb{R}^N \ni x \mapsto R_\nu(x) = \nu \log \left(\sum_{i=1}^N 2^{\frac{x_i}{\nu}} \right) \in \mathbb{R}$ it is clear that $\nabla_\lambda R_\nu(f(\lambda)) = \nabla G_\nu(\lambda)$. The basic idea is to define $\bar{f}(\lambda) := \max_{1 \leq i \leq N} f_i(\lambda)$ and then consider a function $g : \mathbb{R}^M \rightarrow \mathbb{R}^N$ given by $g_i(\lambda) = f_i(\lambda) - \bar{f}(\lambda)$, such that all components of $g(\lambda)$ are non-positive. One can show that

$$\nabla_\lambda R_\nu(f(\lambda)) = \nabla_\lambda R_\nu(g(\lambda)) + \nabla \bar{f}(\lambda),$$

where the term on the right-hand side can be computed with a small numerical error.

Remark 2.12 (Computational complexity). In case of no input cost constraint, one can see by (25) that the computational complexity of a single iteration step of Algorithm 1 is $O(MN)$. Furthermore, according to (24), the complexity in terms of number of iterations to achieve an ε -precise solution is $O\left(\frac{M\sqrt{\log N}}{\varepsilon}\right)$. This finally gives a computational complexity for finding an additive ε -solution of $O\left(\frac{M^2 N \sqrt{\log N}}{\varepsilon}\right)$. Let us point out that the constants in the computational complexity, explicitly given in (24) and in particular the dependency on the parameter γ , can have a significant impact on the runtime of the proposed approximation method in practice. In the following remark, however, we present a way to circumvent ill-conditioned channels with very small (or even vanishing) γ parameter.

Remark 2.13 (Removing Assumption 2.3). The continuity of the channel capacity can be used to remove Assumption 2.3. Let $\mathbf{W}_1 \in \mathbb{R}^{N \times M}$ be a channel transition matrix that does not satisfy Assumption 2.3, i.e., that contains zero entries. Define a new channel matrix $\mathbf{W}_2 \in \mathbb{R}^{N \times M}$ by adding a perturbation $\varepsilon > 0$ to all zero entries of \mathbf{W}_1 and then normalizing the rows. According to [24]

$$|C(\mathbf{W}_1) - C(\mathbf{W}_2)| \leq 3 \|\mathbf{W}_1 - \mathbf{W}_2\|_{\triangleright} \log(M \vee N) + 2\eta(\|\mathbf{W}_1 - \mathbf{W}_2\|_{\triangleright}), \quad (26)$$

where $\eta(t) = -t \log t$ and the norm $\|\cdot\|_{\triangleright}$ on $\mathbb{R}^{N \times M}$ is defined as $\|A\|_{\triangleright} := \max_{b \in \Delta_N} \|bb^\top A\|_{\text{tr}}$. Since \mathbf{W}_2 by construction satisfies Assumption 2.3, we can run Algorithm 1 for channel \mathbf{W}_2 and as such get the following upper and lower bounds for the capacity of the singular channel \mathbf{W}_1

$$\begin{aligned} C_{\text{LB}}(\mathbf{W}_1) &:= C_{\text{LB}}(\mathbf{W}_2) - 3 \|\mathbf{W}_1 - \mathbf{W}_2\|_{\triangleright} \log(M \vee N) - 2\eta(\|\mathbf{W}_1 - \mathbf{W}_2\|_{\triangleright}) \\ C_{\text{UB}}(\mathbf{W}_1) &:= C_{\text{UB}}(\mathbf{W}_2) + 3 \|\mathbf{W}_1 - \mathbf{W}_2\|_{\triangleright} \log(M \vee N) + 2\eta(\|\mathbf{W}_1 - \mathbf{W}_2\|_{\triangleright}). \end{aligned}$$

See in Example 2.15 how this perturbation method behaves numerically.

C. Simulation Results

This section presents two examples to illustrate the theoretical results developed in the preceding sections and their performance. All the simulations in this section are performed on a 2.3 GHz Intel Core i7 processor with 8 GB RAM.

Example 2.14. Consider a DMC W having a channel matrix $W \in \mathbb{R}^{N \times M}$ with $N = 10000$ and $M = 100$, such that $W_{ij} = \frac{V_{ij}}{\sum_{j=1}^M V_{ij}}$, where V_{ij} is chosen i.i.d. uniformly distributed in $[0, 1]$ for all $1 \leq i \leq N$ and $1 \leq j \leq M$. The parameter γ happens to be $1.0742 \cdot 10^{-8}$. Figure 2.14 and Table I compare the performance of the Blahut-Arimoto algorithm with that of Algorithm 1, which has the a priori error bound predicted by Theorem 2.9, namely

$$C_{\text{UB}}(W) - C_{\text{LB}}(W) \leq \frac{2M\sqrt{2\log(N)}}{n+1} (\log(\gamma^{-1}) \vee \frac{1}{\ln 2}) + \frac{2M^2}{(n+1)^2} (\log(\gamma^{-1}) \vee \frac{1}{\ln 2})^2,$$

where n denotes the number of iterations and γ is equal to the smallest entry in the channel matrix W . Recall that the Blahut-Arimoto algorithm has an a priori error bound of the form $C(W) - C_{\text{LB}}(W) \leq \frac{\log(N)}{n}$ [4, Corollary 1]. Moreover, the new method provides us with an a posteriori error, which the Blahut-Arimoto algorithm does not.

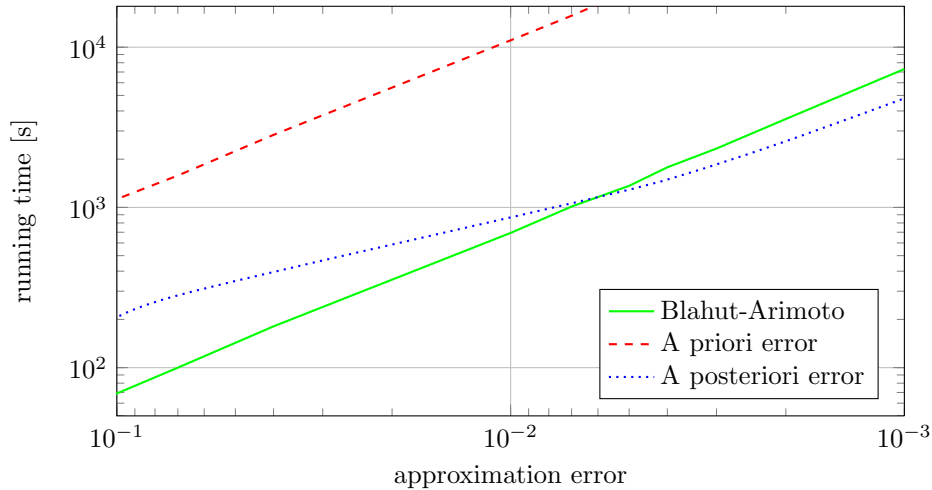


FIG. 1. For Example 2.14, this plot depicts the runtime of Algorithm 1 with respect to the a priori and a posteriori stopping criterion, as explained in Remark 2.10. As a reference, the runtime of the Blahut-Arimoto algorithm is shown.

TABLE I. Some specific simulation points of Example 2.14.

	Blahut-Arimoto Algorithm				Algorithm 1			
A priori error	1	0.1	0.01	0.001	1	0.1	0.01	0.001
$C_{\text{UB}}(W)$	—	—	—	—	0.4419	0.4131	0.4092	0.4088
$C_{\text{LB}}(W)$	0.2930	0.4008	0.4088	0.4088	0.3094	0.4069	0.4088	0.4088
A posteriori error	—	—	—	—	0.1325	0.0063	$4.0 \cdot 10^{-4}$	$3.7 \cdot 10^{-5}$
Time [s]	7.4	69	693	7306	114	1127	11 036	110 987
Iterations	14	133	1329	13 288	27 797	273 447	2 729 860	27 294 000

Example 2.15. Consider a binary erasure channel with erasure probability α whose channel transition matrix is given by $W = \begin{pmatrix} 1-\alpha & \alpha & 0 \\ 0 & \alpha & 1-\alpha \end{pmatrix}$ and as such does not satisfy Assumption 2.3. We use the perturbation method introduced in Remark 2.13 to approximate its capacity that is analytically known to be $1 - \alpha$ [25, p. 189]. Table II shows the performance of this perturbation method and Algorithm 1.

TABLE II. Some specific simulation points of Example 2.15 for $\alpha = 0.4$

Perturbation ε	10^{-4}	10^{-5}	10^{-6}	10^{-7}
A priori error	0.01	0.01	0.01	0.01
$C_{\text{UB}}(W)$	0.6024	0.6003	0.6000	0.6000
$C_{\text{LB}}(W)$	0.5949	0.5994	0.5999	0.6000
A posteriori error	0.0075	$9.2 \cdot 10^{-4}$	$1.1 \cdot 10^{-4}$	$1.2 \cdot 10^{-5}$
Time [s]	0.70	0.54	0.66	0.78
Iterations	9056	7402	8523	9896

3. CHANNELS WITH CONTINUOUS INPUT AND COUNTABLE OUTPUT ALPHABETS

In this section we generalize the approximation scheme introduced in Section 2 to memoryless channels with continuous input and countable output alphabets. The class of discrete-time Poisson channels is an example of such channels with particular interest in applications, for example to model direct detection optical communication systems [10, 26, 27]. Consider $\mathcal{X} \subseteq \mathbb{R}$ as the input alphabet set and $\mathcal{Y} = \mathbb{N}_0$ as the output alphabet set. The channel is described by the conditional probability $W(i|x) := \mathbb{P}[Y = i | X = x]$. Given a channel W and an integer M , we introduce an M -truncated version of the channel by

$$W_M(i|x) := \begin{cases} W(i|x) + \frac{1}{M} \sum_{j \geq M} W(j|x), & i \in \{0, 1, \dots, M-1\} \\ 0, & i \geq M. \end{cases} \quad (27)$$

W_M can be seen as a channel with input alphabet \mathcal{X} and output alphabet $\{0, 1, \dots, M-1\}$. Figure 2 shows a pictorial representation of a channel and its M -truncated counterpart. The finiteness of the output alphabet of W_M allows us to deploy an approximation scheme similar to the one developed in Section 2 to numerically approximate $C(W_M)$.

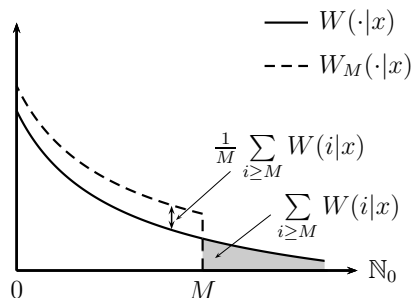


FIG. 2. Pictorial representation of the M -truncated channel counterpart.

The following definition is a key feature of the channel required for the theoretical results developed in this section which, roughly speaking, imposes a certain decay rate for the output distribution uniformly in the input alphabet.

Definition 3.1 (Polynomial tail). The channel W features a k -ordered polynomial tail if for $M \in \mathbb{N}_0$ and $k \in \mathbb{R}_{\geq 0}$

$$R_k(M) := \sum_{i \geq M} \left(\sup_{x \in \mathcal{X}} W(i|x) \right)^k < \infty. \quad (28)$$

The following assumptions hold throughout this section.

Assumption 3.2.

(i) The channel W has a k -ordered polynomial tail for some $k \in (0, 1)$ in the sense of Definition 3.1.

(ii) The mapping $x \mapsto W(i|x)$ is Lipschitz continuous for any $i \in \mathbb{N}_0$ with Lipschitz constant L .

Assumption 3.2 allows us to relate the capacity of the original channel to that of its truncated counterpart.

Theorem 3.3. Suppose channel W satisfies Assumption 3.2(i) with the order $k \in (0, 1)$. Then, for any $M \in \mathbb{N}_0$ and for any probability distribution $p \in \mathcal{P}(\mathcal{X})$ we have

$$|I(p, W) - I(p, W_M)| \leq \frac{2 \log(e)}{e(1-k)} \left[M^{1-k} (R_1(M))^k + R_k(M) \right],$$

where $R_k(M)$ is as defined in (28).

Proof. See Appendix A 4. □

Note that Theorem 3.3 directly implies an upper bound to the capacity since

$$|C(W) - C(W_M)| = \left| \sup_{p \in \mathcal{P}(\mathcal{X})} I(p, W) - \sup_{p \in \mathcal{P}(\mathcal{X})} I(p, W_M) \right| \leq \sup_{p \in \mathcal{P}(\mathcal{X})} |I(p, W) - I(p, W_M)|.$$

We consider two types of input cost constraints: a peak-power constraint $\mathbb{P}[X \in \mathbb{A}] = 1$ for a compact set $\mathbb{A} \subseteq \mathcal{X}$ and an average-power constraint $\mathbb{E}[s(X)] \leq S$ for $S \in \mathbb{R}_{\geq 0}$ and a continuous function s on \mathcal{X} . The primal capacity problem for the channel W_M is given by

$$C_{\mathbb{A}, S}(W_M) = \begin{cases} \sup_p I(p, W_M) \\ \text{s. t. } \mathbb{E}[s(X)] \leq S \\ p \in \mathcal{P}(\mathbb{A}), \end{cases} \quad (29)$$

where $\mathcal{P}(\mathbb{A})$ denotes the space of all probability distributions supported on \mathbb{A} , cf. (3). Our method always requires a peak-power constraint, whereas the average-power constraint is optimal. The following proposition allows us to restrict the optimization variables from probability distributions to probability densities.

Proposition 3.4. The optimization problem (29) is equivalent to

$$C_{\mathbb{A}, S}(W_M) = \begin{cases} \sup_p I(p, W_M) \\ \text{s. t. } \mathbb{E}[s(X)] \leq S \\ p \in \mathcal{D}(\mathbb{A}), \end{cases}$$

where $\mathcal{D}(\mathbb{A})$ is the set of probability densities functions, i.e., $\mathcal{D}(\mathbb{A}) := \{f \in L^1(\mathbb{A}) : f \geq 0, \int_{\mathbb{A}} f(x) dx = 1\}$.

Proof. See Appendix A 5. □

We consider the pair of vector spaces $(L^1(\mathbb{A}), L^\infty(\mathbb{A}))$ together with the bilinear form

$$\langle f, g \rangle := \int_{\mathcal{X}} f(x)g(x)dx.$$

In the light of [28, Theorem 243G] this is a dual pair of vector spaces; we refer to [29, Section 3] for the details of the definition of dual pairs of vector spaces. Considering the standard inner product as a bilinear form on the dual pair $(\mathbb{R}^M, \mathbb{R}^M)$, we define the linear operator $\mathcal{W} : \mathbb{R}^M \rightarrow L^\infty(\mathbb{A})$ and its adjoint operator $\mathcal{W}^* : L^1(\mathbb{A}) \rightarrow \mathbb{R}^M$, given by

$$\mathcal{W}\lambda(x) := \sum_{i=1}^M W_M(i-1|x)\lambda_i, \quad (\mathcal{W}^*p)_i := \int_X W_M(i-1|x)p(x)dx.$$

Let $S_{\max} := \inf_{p \in \mathcal{D}(\mathbb{A})} \{ \langle p, s \rangle : I(p, W_M) = \sup_{q \in \mathcal{D}(\mathbb{A})} I(q, W_M) \}$. Following similar lines as in Lemma 2.1, one can deduce that in problem (29) the inequality input constraint can be replaced by equality (resp. removed) is $S < S_{\max}$ (resp. $S \geq S_{\max}$). That is, in view of Proposition 3.4, Lemma 2.1 and the discussion there, problem (29) (under Assumption 3.6, that we require later) is equivalent to

$$\mathbf{P} : \begin{cases} \sup_{p, q} -\langle p, r \rangle + H(q) \\ \text{s. t. } \mathcal{W}^*p = q \\ \langle p, s \rangle = S \\ p \in \mathcal{D}(\mathbb{A}), q \in \Delta_M, \end{cases} \quad (30)$$

where $r(\cdot) := -\sum_{j=0}^{M-1} W_M(j|\cdot) \log(W_M(j|\cdot))$ is an element in $L^\infty(\mathbb{A})$ by Assumption 3.2(ii). For the rest of the section we restrict attention to (30), since unconstrained problem can be solved in a similar way. We call (30) the primal program. Thanks to the dual vector space framework, the Lagrange dual program of P is given by

$$\mathbf{D} : \begin{cases} \inf_{\lambda} G(\lambda) + F(\lambda) \\ \text{s. t. } \lambda \in \mathbb{R}^M, \end{cases} \quad (31)$$

where

$$G(\lambda) = \begin{cases} \sup_p \langle p, \mathcal{W}\lambda \rangle - \langle p, r \rangle \\ \text{s. t. } \langle p, s \rangle = S \\ p \in \mathcal{D}(\mathbb{A}) \end{cases} \quad \text{and} \quad F(\lambda) = \begin{cases} \max_q H(q) - \lambda^\top q \\ \text{s. t. } q \in \Delta_M. \end{cases}$$

Lemma 3.5. *Strong duality holds between (30) and (31).*

Proof. Note that the dualized constraint is a linear equality constraint. Therefore the conditions of (1) in [30, Theorem 5] holds and as such strong duality follows by [30, Theorem 4]. □

In the remainder of this article we impose the following assumption on the channel.

Assumption 3.6. $\gamma_M := \min_{y \in \{0, 1, \dots, M-1\}} \min_{x \in \mathbb{A}} W_M(y|x) > 0$

In case $\sum_{j \geq M} W(j|x) > 0$ for all x , Assumption 3.6 holds according to (27) and a lower bound can be given by $\gamma_M \geq \frac{1}{M} \min_x \sum_{j \geq M} W(j|x)$. Under Assumption 3.6 we can show that we can again assume without loss of generality that λ takes values in a compact set.

Lemma 3.7. *Under Assumption 3.6, the dual program (31) is equivalent to*

$$\begin{cases} \min_{\lambda} & G(\lambda) + F(\lambda) \\ \text{s.t.} & \lambda \in Q, \end{cases}$$

where $Q := \{\lambda \in \mathbb{R}^M : \|\lambda\|_2 \leq M (\log(\gamma_M^{-1}) \vee \frac{1}{\ln 2})\}$.

Proof. The proof follows the same lines as in the proof of Lemma 2.4. \square

Note that $F(\lambda)$ is the same as in Section 2 and therefore given by (16) and its gradient by (17). As in Section 2, we consider the smooth approximation

$$G_\nu(\lambda) = \begin{cases} \sup_p & \langle p, \mathcal{W}\lambda \rangle - \langle p, r \rangle + \nu h(p) - \nu \log(\rho) \\ \text{s.t.} & \langle p, s \rangle = S \\ & p \in \mathcal{D}(\mathbb{A}), \end{cases} \quad (32)$$

with smoothing parameter $\nu \in \mathbb{R}_{>0}$ and ρ denoting the Lebesgue measure of \mathbb{A} . To analyze the properties of $G_\nu(\lambda)$ we need one more auxiliary lemma.

Lemma 3.8. *The function $\mathcal{D}(\mathbb{A}) \ni p \mapsto -h(p) + \log(\rho) \in \mathbb{R}_{\geq 0}$ is strongly convex with convexity parameter $\sigma = 1$.*

Proof. See Appendix A 6. \square

Furthermore, we can show that the uniform approximation $G_\nu(\lambda)$ is smooth and has a Lipschitz continuous gradient, with known constant. The following result is a generalization of Proposition 2.8.

Proposition 3.9. *$G_\nu(\lambda)$ is well defined and continuously differentiable at any $\lambda \in \mathbb{R}^M$. Moreover, this function is convex and its gradient $\nabla G_\nu(\lambda) = \mathcal{W}^* p_\nu^\lambda$ is Lipschitz continuous with constant $L_\nu = \frac{1}{\nu}$.*

Proof. See Appendix A 7. \square

We denote by p_ν^λ the optimizer to (32), that is unique since the objective function is strictly concave. To analyze the solution to (32) we consider the following optimization problem, that, if feasible, has a closed form solution

$$\begin{cases} \sup_p & h(p) + \langle p, c \rangle \\ \text{s.t.} & \langle p, s \rangle = S \\ & p \in \mathcal{D}(\mathbb{A}), \end{cases} \quad (33)$$

with $c, s \in L^\infty(\mathbb{A})$.

Lemma 3.10. *Let $p^*(x) = 2^{\mu_1 + c(x) + \mu_2 s(x)}$, where $\mu_1, \mu_2 \in \mathbb{R}$ are chosen such that p^* satisfies the constraints in (33). Then p^* uniquely solves (33).*

The proof directly follows from [25, p. 409] and the proof of Lemma 2.2. Hence, $G_\nu(\lambda)$ has a (unique) analytical optimizer

$$p_\nu^\lambda(x, \mu) = 2^{\mu_1 + \frac{1}{\nu}(\mathcal{W}\lambda(x) - r(x)) + \mu_2 s(x)}, \quad x \in \mathcal{X}, \quad (34)$$

where $\mu_1, \mu_2 \in \mathbb{R}$ have to be chosen such that $\langle p_\nu^\lambda(\cdot, \mu), s \rangle = S$ and $p_\nu^\lambda(\cdot, \mu) \in \mathcal{D}(\mathbb{A})$; for this choice of μ_1, μ_2 we denote the solution by $p_\nu^\lambda(\cdot)$.

Remark 3.11 (No input constraints). In case of no input constraints, the unique optimizer to (32) is given by

$$p_\nu^\lambda(x) = \frac{2^{\frac{1}{\nu}(\mathcal{W}\lambda(x)-r(x))}}{\int_{\mathbb{A}} 2^{\frac{1}{\nu}(\mathcal{W}\lambda(x)-r(x))} dx},$$

whose numerical evaluation can be done in a stable way by following Remark 2.11.

Remark 3.12 (Additional input constraints). As in Remark 2.7, in case of additional input constraints we need an efficient method to find the coefficients μ_i in (34). This problem can again be reduced to a finite dimensional convex optimization problem ([22, Theorem 4.8],[23, p. 257 ff.]), in the sense that the coefficients μ_i are the unique maximizers to

$$\max_{\mu \in \mathbb{R}^2} \left\{ y^\top \mu - \int_{\mathbb{A}} p_\nu^\lambda(x, \mu) dx \right\}, \quad (35)$$

where $y := (1, S)$. Note that (35) is an unconstrained maximization of a strictly concave function. The evaluation of the gradient and the Hessian of this objective function involves computing moments of the measure $p_\nu^\lambda(x, \mu) dx$, which unlike to the finite input alphabet case (Remark 2.7) is numerically difficult. In [23, p. 259 ff.], an efficient approximation of the mentioned gradient and Hessian in terms of two single semidefinite programs involving two linear matrix inequalities (LMI) is presented, where the desired accuracy is controlled by the size of the LMI constraints. As mentioned in Remark 2.7, this will provide a suboptimal solution to the maximum entropy problem (32) and as such the error bounds of Theorem 3.16 do not hold. By following [31], however, one can quantify the approximation error of Algorithm 1 in case of an inexact gradient. We also refer the interested reader to [32], for a related work on channel capacity approximation under inexact first-order information.

Note that the differential entropy $h(p) \leq \log(\rho)$ for all $p \in \mathcal{D}(\mathbb{A})$ and that there exists a function $\iota : \mathbb{R}_{>0} \rightarrow \mathbb{R}_{\geq 0}$ such that

$$G_\nu(\lambda) \leq G(\lambda) \leq G_\nu(\lambda) + \iota(\nu) \text{ for all } \lambda \in Q, \quad (36)$$

i.e., $G_\nu(\lambda)$ is a uniform approximation of the non-smooth function $G(\lambda)$. The following lemma, Lemma 3.15, provides an explicit expression for the function ι in (36) under some Lipschitz continuity assumptions, implying in particular that $\iota(\nu) \rightarrow 0$ as $\nu \rightarrow 0$.

Lemma 3.13. *Under Assumption 3.2(ii) and Assumption 3.6 the function $f_\lambda(\cdot) := \mathcal{W}\lambda(\cdot) - r(\cdot)$ is Lipschitz continuous uniformly in $\lambda \in Q$ with constant $L_f = LM^2(\log \frac{1}{\gamma_M} \vee \frac{1}{\ln 2}) + ML|\log \frac{1}{\gamma_M} - \frac{1}{\ln 2}|$.*

Proof. See Appendix A 8. □

Assumption 3.14 (Lipschitz continuity of the average-power constraint function). The average-power constraint function $s(\cdot)$ is Lipschitz continuous with constant L_s .

Lemma 3.15. *Under Assumptions 3.2(ii), 3.6 and 3.14 a possible choice of the function ι in (36) is given by*

$$\iota(\nu) = \begin{cases} \nu \left(\log \left(\frac{T_1}{\nu} + T_2 \right) + 1 \right), & \nu < \frac{T_1}{1-T_2} \text{ or } T_2 > 1 \\ \nu, & \text{otherwise,} \end{cases}$$

where $T_1 := L_f \rho + 2L_f L_s \rho^2 \left(\frac{1}{\underline{s}} \vee \frac{1}{\bar{s}} \right)$, $T_2 := L_s \rho (\underline{\mu} \vee \bar{\mu})$, $\underline{\mu} := \frac{2}{\underline{s}} \log \left(\frac{2L_s \rho}{\underline{s}} \vee 1 \right)$, $\bar{\mu} := \frac{2}{\bar{s}} \log \left(\frac{2L_s \rho}{\bar{s}} \vee 1 \right)$, $\rho := \int_{\mathbb{A}} dx$, $\underline{s} := -S + \min_{x \in \mathbb{A}} s(x)$ and $\bar{s} := -S + \max_{x \in \mathbb{A}} s(x)$.

Proof. See Appendix A 9. □

We consider the smooth, finite dimensional, convex optimization problem

$$D_\nu : \begin{cases} \inf_{\lambda} & F(\lambda) + G_\nu(\lambda) \\ \text{s.t.} & \lambda \in Q, \end{cases} \quad (37)$$

whose solution can be approximated with Algorithm 1 presented in Section 2, as follows. Define the constant $D_1 := \frac{1}{2}(M \log(\gamma^{-1}) \vee \frac{1}{\ln 2})^2$.

Theorem 3.16. *Under Assumptions 3.2(ii), 3.6 and 3.14, let $\alpha := 2(T_1 + T_2 + 1)$ where T_1 and T_2 are as defined in Lemma 3.15. Given precision $\varepsilon \in (0, \frac{\alpha}{4})$, we set the smoothing parameter $\nu = \frac{\varepsilon/\alpha}{\log(\alpha/\varepsilon)}$ and number of iterations $n \geq \frac{1}{\varepsilon} \sqrt{8D_1 \alpha} \sqrt{\log(\varepsilon^{-1}) + \log(\alpha) + \frac{1}{4}}$. Consider*

$$\hat{\lambda} = y_n \in Q \quad \text{and} \quad \hat{p} = \sum_{k=0}^n \frac{2(i+1)}{(n+1)(n+2)} p_\nu^{x_k} \in \mathcal{D}(\mathbb{A}), \quad (38)$$

where y_k computed at the k^{th} iteration of Algorithm 1 and $p_\nu^{x_k}$ is the analytical solution in (34). Then, $\hat{\lambda}$ and \hat{p} are the approximate solutions to the problems (31) and (29), i.e.,

$$0 \leq F(\hat{\lambda}) + G(\hat{\lambda}) - I(\hat{p}, W_M) \leq \varepsilon. \quad (39)$$

Therefore, Algorithm 1 requires $O\left(\frac{1}{\varepsilon} \sqrt{\log(\varepsilon^{-1})}\right)$ iterations to find an ε -solution to the problems (31) and (29).

Proof. See Appendix A 10. □

Hence, under Assumption 3.14 we can quantify the approximation error of the presented method to find the capacity of any channel W , satisfying Assumptions 3.2 and 3.6, by

$$\left| C(W) - C_{\text{approx}}^{(n)}(W_M) \right| \leq \underbrace{|C(W) - C(W_M)|}_{(\star)} + \underbrace{|C(W_M) - C_{\text{approx}}^{(n)}(W_M)|}_{(\star\star)},$$

where (\star) and $(\star\star)$ are addressed by Theorem 3.3 and Theorem 3.16, respectively. Let us highlight that for the term $(\star\star)$ we have two different quantitative bounds: First, the *a priori* bound ε for which Theorem 3.16 prescribes a lower bound for the required number of iterations; second, the *a posteriori* bound $F(\hat{\lambda}) + G(\hat{\lambda}) - I(\hat{p}, W_M)$ which can be computed after a number of iterations have been executed. In practice, the *a posteriori* bound often approaches ε much faster than the *a priori* bound. Note also that by (36) and Theorem 3.16

$$0 \leq F(\hat{\lambda}) + G_\nu(\hat{\lambda}) + \iota(\nu) - I(\hat{p}, W_M) \leq \iota(\nu) + \varepsilon,$$

which shows that $F(\hat{\lambda}) + G_\nu(\hat{\lambda}) + \iota(\nu)$ is an upper bound for the channel capacity with *a priori* error $\iota(\nu) + \varepsilon$. This bound can be particularly helpful in cases where an evaluation of $G(\lambda)$ for a given λ is hard.

Remark 3.17 (Optimal tail truncation). Given a fixed number of iterations, the term $(\star\star)$ above is effected by the truncation level M for two reasons: the higher M the larger the size of the output as well as the lower the parameter γ_M . Therefore, term $(\star\star)$ increases as M increases, which can be quantified by (A17). On the other hand, term (\star) obviously has the opposite behavior. Namely,

the higher M leads to the better approximation of the channel W by the truncated version W_M as quantified in Theorem 3.3. Hence, given a channel W with the polynomial tail order k , there is an optimal value for the truncation parameter M , which thanks to the monotonicity explained above can be effectively computed in practice by techniques such as bisection.

Note, that this truncation procedure could also be applied to a finite output alphabet, given that the channel satisfies Assumption 3.2(i), and for example improve the performance of the method presented in Section 2.

Remark 3.18 (Without average-power constraint). In case of considering only a peak-power constraint and no average-power constraint, our proposed methodology allows us to access a closed form expression for $G_\nu(\lambda)$ and its gradient,

$$G_\nu(\lambda) = \nu \log \left(\int_{\mathbb{A}} 2^{\frac{1}{\nu}(\mathcal{W}\lambda(x)-r(x))} dx \right) - \nu \log(\rho) \quad (40)$$

$$\nabla G_\nu(\lambda) = \frac{\int_{\mathbb{A}} 2^{\frac{1}{\nu}(\mathcal{W}\lambda(x)-r(x))} W_M(\cdot|x) dx}{\int_{\mathbb{A}} 2^{\frac{1}{\nu}(\mathcal{W}\lambda(x)-r(x))} dx}.$$

Discrete-Time Poisson Channel

The discrete-time Poisson channel is a mapping from $\mathbb{R}_{\geq 0}$ to \mathbb{N}_0 , such that conditioned on the input $x \geq 0$ the output is Poisson distributed with mean $x + \eta$, i.e.,

$$W(y|x) = e^{-(x+\eta)} \frac{(x+\eta)^y}{y!}, \quad y \in \mathbb{N}_0, x \in \mathbb{R}_{\geq 0}, \quad (41)$$

where $\eta \geq 0$ denotes a constant sometimes referred to as *dark current*. A peak-power constraint on the transmitter is given by the peak-input constraint $X \leq A$ with probability one, i.e., $\mathbb{A} = [0, A]$ and an average-power constraint on the transmitter is considered by $\mathbb{E}[X] \leq S$.

Up to now, no analytic expression for the capacity of a discrete-time Poisson channel is known. However, for different scenarios lower and upper bounds exist. Brady and Verdú derived a lower and upper bound in the presence of only an average-power constraint [33]. Later, for $\eta = 0$ and only an average-power constraint, Martinez introduced better upper and lower bounds [34]. Lapidoth and Moser derived a lower bound and an asymptotic upper bound, which is valid only when the available peak and average power tend to infinity with their ratio held fixed, for the presence of a peak and average-power constraint [16]. Lapidoth *et al.* computed the asymptotic capacity of the discrete-time Poisson channel when the allowed average-input power tends to zero with the allowed peak power — if finite — held fixed and the dark current is constant or tends to zero proportionally to the average power [35].

In [11] a numerical algorithm is presented, where the Blahut-Arimoto algorithm is incorporated into the deterministic annealing method, that allows the computation of both the channel capacity under peak and average power constraints and its associated optimal input distribution. Furthermore, the works [11, 12] derive several fundamental properties of capacity achieving input distributions for the discrete-time Poisson channel.

Here, we numerically approximate the capacity of a discrete-time Poisson channel using the proposed algorithm. For simplicity, we consider the case where only a peak power constraint is imposed; the case where an additional average power constraint is present can be treated similarly. It was shown in [27] that in the case of a peak power constraint (with or without average power constraint), the capacity achieving input distribution is discrete. This, in the limit as the number of iterations in the proposed approximation method goes to infinity, is consistent with the optimal input distribution given in Remark 3.11.

The following proposition provides an upper bound for the k -polynomial tail for the Poisson channel W as defined in (41).

Proposition 3.19 (Poisson tail). *The Poisson channel (41) having a bounded input alphabet $\mathcal{X} = [0, A]$ and dark current parameter η has a k -polynomial tail for any $k \in (0, 1]$ in the sense of Definition 3.1, which is upper bounded for all $M \geq A + \eta$ by*

$$R_k(M) \leq \left(\alpha e^{(\alpha-1)(A+\eta)} \frac{(A+\eta)^M}{M!} \right)^k, \quad \alpha := 2^{(k^{-1}-1)}.$$

Proof. See Appendix A 11. □

In the following we present an example to illustrate the theoretical results developed in the preceding sections and their performance. Note that for the discrete-time Poisson channel Assumption 3.6 clearly holds.

Example 3.20. We consider a discrete-time Poisson channel W as defined in (41) with a peak-power constraint A and dark current $\eta = 1$. Up to now, the best known lower bound for the capacity is given by [16, Theorem 4]

$$C(W) \geq \frac{1}{\ln 2} \left(\frac{1}{2} \ln A + \left(\frac{A}{3} + 1 \right) \ln \left(1 + \frac{3}{A} \right) - 1 - \sqrt{\frac{\eta + \frac{1}{12}}{A}} \left(\frac{\pi}{4} + \frac{1}{2} \ln 2 \right) - \frac{1}{2} \ln \frac{\pi e}{2} \right). \quad (42)$$

To the best of our knowledge no upper bound for the capacity is known. In [16] an asymptotic upper bound is given which includes an unknown error term that is vanishing in the limit $A \rightarrow \infty$. According to Theorems 3.3 and 3.16, the algorithm introduced in this article leads to an approximation error after n iterations that is given by

$$\begin{aligned} \left| C_{\text{approx}}^{(n)}(W_M) - C(W) \right| &\leq \left| C_{\text{approx}}^{(n)}(W_M) - C(W_M) \right| + |C(W_M) - C(W)| \\ &\leq F(\hat{\lambda}) + G(\hat{\lambda}) - I(\hat{p}, W) + \mathcal{E}, \end{aligned}$$

where $\mathcal{E} = \frac{2 \log(e)}{e(1-k)} \left[M^{1-k} (R_1(M))^k + R_k(M) \right]$, $R_\ell(M) = \left(\alpha e^{(\alpha-1)(A+\eta)} \frac{(A+\eta)^M}{M!} \right)^\ell$ and $\alpha := 2^{(\ell^{-1}-1)}$ for any $k \in (0, 1)$ and $\ell \in (0, 1]$. The truncation parameter M was determined as described in Remark 3.17. This finally leads to the following upper and lower bounds on $C(W)$

$$2I(\hat{p}, W) - \left(F(\hat{\lambda}) + G(\hat{\lambda}) \right) - \mathcal{E} \leq C(W) \leq 2 \left(F(\hat{\lambda}) + G(\hat{\lambda}) \right) - I(\hat{p}, W) + \mathcal{E}. \quad (43)$$

Figure 3.20 compares the two bounds (42) and (43) for different values of A . Further details on the simulation can be found in Appendix B.

Remark 3.21 (AWGN channel with a quantized output). Another example of a channel that is well studied and can be treated by the proposed method is the discrete-time additive white Gaussian noise (AWGN) channel under output quantization. The output of the channel is described by

$$Y = Q(X + N),$$

where $X \in \mathbb{R}$ is the channel input, $N \sim \mathcal{N}(0, \sigma^2)$ for $\sigma^2 > 0$ is white Gaussian noise and $Q(\cdot)$ is a quantizer that maps the real valued input $X + N$ to one of M bins (where we assume $M < \infty$), which gives $Y \in \{y_1, \dots, y_M\}$. In addition an average and/or a peak power constraint at the input is considered. More information about this channel model and why it is of interest can be found in [36, 37]. By definition, the AWGN channel with a quantized output has a continuous input alphabet and a discrete output alphabet. Thus, the approximation method discussed in this section can be used to compute the capacity of such channels.

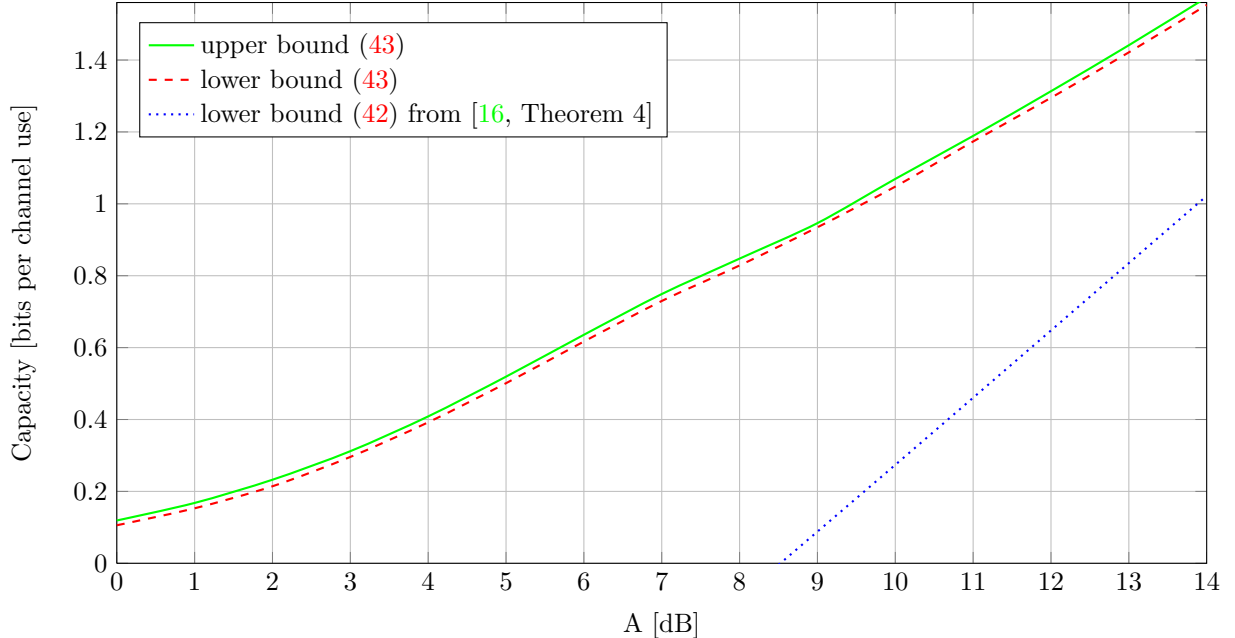


FIG. 3. This plot depicts the capacity of a discrete-time Poisson channel with dark current $\eta = 1$ as a function of the peak-power constraint parameter A . The red (resp. green) line shows the lower (resp. upper) bound (43) obtained for a moderate number of iterations, see Appendix B. As a comparison we plot the lower bound of [16], which to the best of our knowledge is the tightest lower bound available to date (blue line). The parameter A is given in decibels where $A[\text{dB}] = 10 \log_{10}(A)$.

4. CONCLUSION AND FUTURE WORK

We introduced a new approach to approximate the capacity of DMCs possibly having constraints on the input distribution. The dual problem of Shannon’s capacity formula turns out to have a particular structure such that the Lagrange dual function admits a closed form solution. Applying smoothing techniques to the non-smooth dual function enables us to solve the dual problem efficiently. This new approach, in the case of no constraints on the input distribution, has a computational complexity per iteration step of $O(MN)$, where N is the input alphabet size and M the size of the output alphabet. In comparison, the Blahut-Arimoto algorithm has the same computational cost of $O(MN)$ per iteration step. More precisely for no input power constraint, the total computational cost to find an ε -close solution is $O(\frac{M^2 N \sqrt{\log(N)}}{\varepsilon})$ for the algorithm developed in this article, whereas the Blahut-Arimoto algorithm requires $O(\frac{MN \log(N)}{\varepsilon})$. A strength of the new approach is that it provides an a posteriori error, i.e., after having run a certain number of iterations we can precisely estimate the actual error in the current approximation. This is computationally appealing as explicit (or a priori) error bounds often are conservative in practice. By exploiting this a posteriori bound we can stop the computation once the desired accuracy has been reached.

As a second contribution, we have shown how similar ideas can be used to approximate the capacity of memoryless channels with continuous bounded input alphabets and countable output alphabets under a mild assumption on the channels tail. This assumption holds, for example, for discrete-time Poisson channels, allowing us to efficiently approximate their capacity. As an example we derived upper and lower bounds for a discrete-time Poisson channel having a peak-power constraint at the input.

The presented optimization method highly depends on the Lipschitz constant estimate of the objective's gradient. The worse this estimate the more steps the method requires for an a priori ε -precision. For future work, we aim to study the derivation of local Lipschitz constants of the gradient. This technique has recently been shown to be very efficient in practice (up to three orders of magnitude reduction of computation time), while preserving the worst-case complexity [38].

In the case of a continuous input alphabet, the proposed method requires to evaluate the gradient $\nabla G_\nu(\cdot)$ in every step of Algorithm 1, that requires solving an integral over \mathbb{A} . As such the method used to compute those integrals has to be included to the complexity of the proposed algorithm. Therefore, it would be interesting to investigate under which structural properties on the channel the gradient $\nabla G_\nu(\cdot)$ can be evaluated efficiently.

The approach introduced in this article can be used to efficiently approximate the capacity of classical-quantum channels, i.e., channels that have classical input and quantum mechanical output, with a discrete or bounded continuous input alphabet. Using the idea of a universal encoder allows us to compute close upper and lower bounds for the Holevo capacity [32].

Appendix A: Proofs

This appendix collects the technical proofs omitted above.

1. Proof of Lemma 2.1

The mutual information $I(p, W)$ can be expressed as

$$\begin{aligned} I(p, W) &= \sum_{i=1}^N \sum_{j=1}^M W_{ij} p_i \log \left(\frac{W_{ij}}{\sum_{k=1}^N W_{kj} p_k} \right) \\ &= \sum_{i=1}^N \sum_{j=1}^M \left[p_i W_{ij} \log(W_{ij}) - p_i W_{ij} \log \left(\sum_{k=1}^N W_{kj} p_k \right) \right]. \end{aligned}$$

By adding the constraint $\sum_{i=1}^N p_i W_{ij} = q_j$ for all $j = 1, \dots, M$,

$$\begin{aligned} I(p, W) &= \sum_{i=1}^N \sum_{j=1}^M [p_i W_{ij} \log(W_{ij}) - p_i W_{ij} \log(q_j)] \\ &= \sum_{i=1}^N \sum_{j=1}^M p_i W_{ij} \log(W_{ij}) - \sum_{j=1}^M q_j \log(q_j) \\ &= -r^\top p + H(q), \end{aligned}$$

where $p \in \Delta_N$. Since $q = W^\top p$ and W^\top is a stochastic matrix, this implies $q \in \Delta_M$. By definition of S_{\max} it is obvious that the input cost constraint $s^\top p \leq S$ is inactive for $S \geq S_{\max}$, leading to the first optimization problem in Lemma 2.1. It remains to show that for $S < S_{\max}$, the input constraint can be written with equality, leading to the second optimization problem in Lemma 2.1. In order to keep the notation simple we define $C(S) := C_S(W)$ for a fixed channel W . We show that $C(S)$ is concave in S for $S \in [0, S_{\max}]$. Let $S^{(1)}, S^{(2)} \in [0, S_{\max}]$, $0 \leq \lambda \leq 1$ and $p^{(i)}$ probability mass functions that achieve $C(S^{(i)})$ for $i \in \{1, 2\}$. Consider the probability mass function $p^{(\lambda)} = \lambda p^{(1)} + (1 - \lambda) p^{(2)}$. We can write

$$s^\top p^{(\lambda)} = \lambda s^\top p^{(1)} + (1 - \lambda) s^\top p^{(2)}$$

$$\begin{aligned}
&\leq \lambda S^{(1)} + (1 - \lambda)S^{(2)} \\
&=: S^{(\lambda)} \in [0, S_{\max}].
\end{aligned} \tag{A1}$$

Using the concavity of the mutual information in the input distribution, we obtain

$$\begin{aligned}
\lambda C(S^{(1)}) + (1 - \lambda)C(S^{(2)}) &= \lambda I(p^{(1)}, W) + (1 - \lambda)I(p^{(2)}, W) \\
&\leq I(p^{(\lambda)}, W) \\
&\leq C(S^{(\lambda)}),
\end{aligned}$$

where the final inequality follows by Shannon's formula for the capacity given in (1). $C(S)$ clearly is non-decreasing in S since enlarging S relaxes the input cost constraint. Furthermore, we show that

$$C(S_{\max} - \varepsilon) < C(S_{\max}), \quad \text{for all } \varepsilon > 0. \tag{A2}$$

Suppose $C(S_{\max} - \varepsilon) = C(S_{\max})$ and denote $C^* := \max_{p \in \Delta_N} I(p, W)$. This then implies that there exists $\bar{p} \in \Delta_N$ such that $I(\bar{p}, W) = C^*$ and $s^\top \bar{p} \leq S_{\max} - \varepsilon$, which contradicts the definition of S_{\max} . Hence, the concavity of $C(S)$ together with the non-decreasing property and (A2) imply that $C(S)$ is strictly increasing in S . \square

2. Proof of Lemma 2.2

This proof is similar to the proof given in [25, Theorem 12.1.1]. Let q satisfy the constraints in (12). Then

$$\begin{aligned}
J(q) &= H(q) - c^\top q = -\sum_{i=1}^N q_i \log(q_i) - c^\top q \\
&= -\sum_{i=1}^N q_i \log\left(\frac{q_i}{p_i^*} p_i^*\right) - c^\top q = -D(q \| p^*) - \sum_{i=1}^N q_i \log(p_i^*) - c^\top q \\
&\leq -\sum_{i=1}^N q_i \log(p_i^*) - c^\top q
\end{aligned} \tag{A3a}$$

$$= -\sum_{i=1}^N q_i (\mu_1 + \mu_2 s_i) \tag{A3b}$$

$$= -\sum_{i=1}^N p_i^* (\mu_1 + \mu_2 s_i) - c^\top p^* + c^\top p^* \tag{A3c}$$

$$= -\sum_{i=1}^N p_i^* \log(p_i^*) - c^\top p^* = J(p^*).$$

The inequality follows from the non-negativity of the relative entropy. Equality (A3b) follows by the definition of p^* and (A3c) uses the fact that both p^* and q satisfy the constraints in (12). Note that equality holds in (A3a) if and only if $q = p^*$. This proves the uniqueness. \square

3. Proof of Lemma 2.4

Consider the following two convex optimization problems

$$P_\beta : \begin{cases} \max_{p,q,\varepsilon} & -r^\top p + H(q) - \beta\varepsilon \\ \text{s.t.} & \|W^\top p - q\|_\infty \leq \varepsilon \\ & s^\top p = S \\ & p \in \Delta_N, q \in \Delta_M, \varepsilon \in \mathbb{R}_{\geq 0} \end{cases} \quad \text{and} \quad D_\beta : \begin{cases} \min_{\lambda} & F(\lambda) + G(\lambda) \\ \text{s.t.} & \|\lambda\|_1 \leq \beta \\ & \lambda \in \mathbb{R}^M. \end{cases}$$

Claim A.1. *Strong duality holds between P_β and D_β .*

Proof. According to the identity $\|W^\top p - q\|_\infty = \max_{\|\lambda\|_1 \leq 1} \lambda^\top (W^\top p - q)$ [39, p. 7] the optimization problem P_β can be rewritten as

$$P_\beta : \begin{cases} \max_{p,q} & -r^\top p + H(q) + \min_{\|\lambda\|_1 \leq \beta} \lambda^\top (W^\top p - q) \\ \text{s.t.} & s^\top p = S \\ & p \in \Delta_N, q \in \Delta_M, \end{cases}$$

whose dual program, where strong duality holds according to [20, Proposition 5.3.1, p. 169] is given by

$$\begin{cases} \min_{\|\lambda\|_1 \leq \beta} & \max_{p,q} & -r^\top p + H(q) + \lambda^\top (W^\top p - q) \\ \text{s.t.} & & s^\top p = S \\ & & p \in \Delta_N, q \in \Delta_M, \end{cases}.$$

which clearly is equivalent to D_β with $F(\cdot)$ and $G(\cdot)$ as given in (14). \square

Denote by $\varepsilon^*(\beta)$ the optimizer of P_β with the respective optimal value J_β^* . We show that for a sufficiently large β the optimizer $\varepsilon^*(\beta)$ of P_β is equal to zero. Hence, in light of the duality relation, the constraint $\|\lambda\|_1 \leq \frac{\beta}{2}$ in D_β is inactive and as such D_β is equivalent to D in equation (13). Note that for

$$J(\varepsilon) := \begin{cases} \max_{p,q} & -r^\top p + H(q) \\ \text{s.t.} & \|W^\top p - q\|_\infty \leq \varepsilon \\ & s^\top p = S \\ & p \in \Delta_N, q \in \Delta_M \end{cases}, \quad (\text{A4})$$

the mapping $\varepsilon \mapsto J(\varepsilon)$, the so-called perturbation function, is concave [40, p. 268]. In the next step we write the optimization problem (A4) in another equivalent form

$$J(\varepsilon) = \begin{cases} \max_{p,v} & -r^\top p + H(W^\top p + \varepsilon v) \\ \text{s.t.} & \|v\|_\infty \leq 1 \\ & s^\top p = S \\ & p \in \Delta_N, v \in \text{Im}(W^\top) \subset \mathbb{R}^M \end{cases}. \quad (\text{A5})$$

By using Taylor's theorem, there exists $y_\varepsilon \in [0, \varepsilon]$ such that the entropy term in the objective function of (A5) can be bounded as

$$H(W^\top p + \varepsilon v) = H(W^\top p) - \left(\log(W^\top p) + \frac{1}{\ln 2} \mathbf{1}\right)^\top v \varepsilon - \sum_{j=1}^M \frac{v_j^2}{\sum_{i=1}^N W_{ij} p_i + y_\varepsilon v_j} \varepsilon^2 \frac{1}{\ln 2}$$

$$\leq H(\mathbf{W}^\top p) - \left(\log(\mathbf{W}^\top p) + \frac{1}{\ln 2} \mathbf{1}\right)^\top v \varepsilon + \frac{M}{\gamma \ln 2} \varepsilon^2. \quad (\text{A6})$$

Thus, the optimal value of problem \mathbf{P}_β can be expressed as

$$\begin{aligned} J_\beta^* &\leq \max_\varepsilon \{J(\varepsilon) - \beta \varepsilon\} \\ &\leq \max_\varepsilon \left\{ \max_{p,v} \left[-r^\top p + H(\mathbf{W}^\top p) - \left(\log(\mathbf{W}^\top p) + \frac{1}{\ln 2} \mathbf{1}\right)^\top v \varepsilon : s^\top p = S \right] + \frac{M}{\gamma \ln 2} \varepsilon^2 - \beta \varepsilon \right\} \end{aligned} \quad (\text{A7a})$$

$$\leq \max_\varepsilon \left\{ \max_{p,v} [-r^\top p + H(\mathbf{W}^\top p) : s^\top p = S] + (\rho - \beta) \varepsilon + \frac{M}{\gamma \ln 2} \varepsilon^2 \right\} \quad (\text{A7b})$$

$$= J(0) + \max_\varepsilon \left\{ (\rho - \beta) \varepsilon + \frac{M}{\gamma \ln 2} \varepsilon^2 \right\}, \quad (\text{A7c})$$

where $\rho = M \left(\log(\gamma^{-1}) \vee \frac{1}{\ln 2} \right)$. Note that (A7a) follows from (A5) and (A6). The equation (A7b) uses the fact that for $\|v\|_\infty \leq 1$, $-\left(\log(\mathbf{W}^\top p) + \frac{1}{\ln 2} \mathbf{1}\right)^\top v \leq M \left(\log(\gamma^{-1}) \vee \frac{1}{\ln 2} \right)$. Thus, for $\beta > \rho$ and $\varepsilon_1 = \frac{\gamma \ln 2}{M} (\beta - \rho)$, we have $\max_{\varepsilon \leq \varepsilon_1} \left\{ (\rho - \beta) \varepsilon + \frac{M}{\gamma \ln 2} \varepsilon^2 \right\} = 0$. Therefore, (A7c) together with the concavity of the mapping $\varepsilon \mapsto J(\varepsilon)$ imply that $J(0)$ is the global optimum of $J(\varepsilon)$ and as such $\varepsilon^*(\beta) = 0$ for $\beta > \rho$, indicating that \mathbf{P}_β is equivalent to \mathbf{P} in the sense that $J_\beta^* = J_0^*$. By strong duality this implies that the constraint $\|\lambda\|_1 \leq \beta$ in \mathbf{D}_β is inactive. Finally, $\|\lambda\|_2 \leq \|\lambda\|_1$ concludes the proof. \square

4. Proof of Theorem 3.3

To prove Theorem 3.3 we need a preliminary lemma.

Lemma A.2. *Given $k \in (0, 1)$ and $p \in [0, 1]$, we have for all $x \in [0, 1 - p]$*

$$|(p+x) \log(p+x) - p \log(p)| \leq \frac{\log(e)}{e(1-k)} x^k.$$

Proof. Note that for a fixed $x \in [0, 1]$, the mapping $p \mapsto (p+x) \log(p+x) - p \log(p)$ is non-decreasing; observe that the derivative of the mapping is non-negative for all $x \in [0, 1]$. Therefore, it suffices to verify the claim for $p \in \{0, 1\}$. For $p = 1$ and accordingly $x = 0$, Lemma A.2 holds trivially. Let $p = 0$ and $h(x) := \frac{\log(e)}{e(1-k)} x^{k-1} + \log(x)$. Note that $h(1) = \frac{\log(e)}{e(1-k)} > 0$ and $h(x) \rightarrow \infty$ as $x \rightarrow 0$. Hence, by setting $\frac{d}{dx} h(x^*) = 0$, it can be easily seen that

$$\min_{x \in (0, 1]} h(x) = h(x^*) = 0, \quad x^* := e^{\frac{1}{k-1}}.$$

Thus $h(x) \geq 0$, and consequently $xh(x) \geq 0$ for all $x \in (0, 1]$, which concludes the proof. \square

Proof of Theorem 3.3. We bound the mutual information difference uniformly in the input probability distribution $p \in \mathcal{P}(\mathcal{X})$. Observe that

$$\begin{aligned} &|I(p, W) - I(p, W_M)| \\ &= \left| \int_{\mathcal{X}} \left[-h(W(\cdot, x)) + h(W_M(\cdot, x)) \right] p(dx) + h\left(\int_{\mathcal{X}} W(\cdot, x) p(dx) \right) - h\left(\int_{\mathcal{X}} W_M(\cdot, x) p(dx) \right) \right| \\ &= \left| \int_{\mathcal{X}} \left[\sum_{i \in \mathbb{N}_0} W(i|x) \log(W(i|x)) - W_M(i|x) \log(W_M(i|x)) \right] p(dx) \right| \end{aligned}$$

$$\begin{aligned}
& + \sum_{i \in \mathbb{N}_0} - \left(\int_{\mathcal{X}} W(i|x)p(dx) \right) \log \left(\int_{\mathcal{X}} W(i|x)p(dx) \right) \\
& + \left(\int_{\mathcal{X}} W_M(i|x)p(dx) \right) \log \left(\int_{\mathcal{X}} W_M(i|x)p(dx) \right) \Big|.
\end{aligned}$$

By the definition of the truncated channel in (27) and applying Lemma A.2 to the above relation, we have

$$\begin{aligned}
|I(p, W) - I(p, W_M)| & \leq \frac{\log(e)}{e(1-k)} \left(\int_{\mathcal{X}} \left[\sum_{i < M} \left(\frac{1}{M} \sum_{j \geq M} W(j|x) \right)^k + \sum_{i \geq M} (W(i|x))^k \right] p(dx) \right. \\
& \quad \left. + \sum_{i < M} \left(\frac{1}{M} \sum_{j \geq M} \int_{\mathcal{X}} W(j|x)p(dx) \right)^k + \sum_{i \geq M} \left(\int_{\mathcal{X}} W(i|x)p(dx) \right)^k \right) \\
& \leq \frac{2 \log(e)}{e(1-k)} \left(M \left(\frac{R_1(M)}{M} \right)^k + R_k(M) \right),
\end{aligned}$$

which concludes the proof. \square

5. Proof of Proposition 3.4

We show that the optimization problem (29) is equivalent to

$$C_{\mathbb{A}, S}(W_M) = \sup_{p \in \mathfrak{D}(\mathbb{A})} \{I(p, W_M) : \mathbb{E}[s(X)] \leq S\},$$

where $\mathfrak{D}(\mathbb{A})$ is the space of probability measures that are absolutely continuous with respect to the Lebesgue measure. This completes the proof since optimizing over $\mathfrak{D}(\mathbb{A})$ is equivalent to optimizing over the space of probability densities $\mathcal{D}(\mathbb{A})$ according to the Radon-Nikodým Theorem [41, Theorem 3.8, p. 90].

It is known that the mapping $p \mapsto I(p, W_M)$ is weakly lower semicontinuous [42]. It then suffices to show that $\mathfrak{D}(\mathbb{A})$ is weakly dense in $\mathcal{P}(\mathbb{A})$. Let \mathbb{B} be a countable dense subset of \mathbb{A} , and $\Delta(\mathbb{B})$ be the family of probability measures whose supports are finite subsets of \mathbb{B} . It is well known that $\Delta(\mathbb{B})$ is weakly dense in $\mathcal{P}(\mathbb{A})$, i.e., $\mathcal{P}(\mathbb{A}) = \overline{\Delta(\mathbb{B})}$ [43, Theorem 4, p. 237], where $\overline{\Delta}$ is the weak closure of Δ . Moreover, thanks to the Lebesgue differentiation theorem [41, Theorem 3.21, p. 98], we know that for any $b \in \mathbb{B}$ the point measure $\delta_{\{b\}} \in \Delta(\mathbb{B})$ can be arbitrarily weakly approximated by measures in $\mathfrak{D}(\mathbb{A})$, i.e., $\delta_{\{b\}} \in \overline{\mathfrak{D}(\mathbb{A})}$. Hence, we have $\Delta(\mathbb{B}) = \overline{\mathfrak{D}(\mathbb{A})}$, which in light of the preceding assertion implies $\mathcal{P}(\mathbb{A}) = \overline{\mathfrak{D}(\mathbb{A})}$. \square

6. Proof of Lemma 3.8

The proof follows the ideas of [15]. It can easily be shown that for $d(p) := -h(p) + \log(\rho)$

$$\langle d''(p) \cdot g, g \rangle = \int_{\mathbb{A}} \frac{g(x)^2}{p(x)} dx.$$

Cauchy-Schwarz then implies

$$\langle d''(p) \cdot g, g \rangle \geq \frac{\left(\int_{\mathbb{A}} g(x) dx \right)^2}{\int_{\mathbb{A}} p(x) dx} = \|g\|^2.$$

\square

7. Proof of Proposition 3.9

It is known, according to Theorem 5.1 in [44], that $G_\nu(\lambda)$ is well defined and continuously differentiable at any $\lambda \in \mathbb{R}^M$ and that this function is convex and its gradient $\nabla G_\nu(\lambda) = \mathcal{W}^* p_\nu^\lambda$ is Lipschitz continuous with constant $L_\nu = \frac{1}{\nu} \|\mathcal{W}\|^2$, where we have also used Lemma 3.8. The operator norm can be simplified to

$$\begin{aligned}
\|\mathcal{W}\| &= \sup_{\lambda \in \mathbb{R}^M, p \in L^1(\mathbb{A})} \{ \langle p, \mathcal{W}\lambda \rangle : \|\lambda\|_2 = 1, \|p\|_1 = 1 \} \\
&\leq \sup_{\lambda \in \mathbb{R}^M, p \in L^1(\mathbb{A})} \{ \|\mathcal{W}^* p\|_2 \|\lambda\|_2 : \|\lambda\|_2 = 1, \|p\|_1 = 1 \} \\
&\leq \sup_{p \in L^1(\mathbb{A})} \{ \|\mathcal{W}^* p\|_1 : \|p\|_1 = 1 \} \\
&= \sup_{p \in L^1(\mathbb{A})} \left\{ \sum_{i=0}^{M-1} \int_{\mathcal{X}} W_M(i|x) p(x) dx : \|p\|_1 = 1 \right\} \\
&= \sup_{p \in L^1(\mathbb{A})} \left\{ \int_{\mathcal{X}} \|W_M(\cdot|x)\|_1 p(x) dx : \|p\|_1 = 1 \right\} \\
&\leq \sup_{x \in \mathbb{A}} \|W_M(\cdot|x)\|_1 \\
&\leq 1,
\end{aligned} \tag{A8}$$

where (A8) is due to Cauchy-Schwarz. \square

8. Proof of Lemma 3.13

Let $x_1, x_2 \in \mathcal{X}$, then by definition of $f_\lambda(\cdot)$ we obtain

$$\begin{aligned}
&|f_\lambda(x_1) - f_\lambda(x_2)| \\
&= \left| \sum_{i=1}^M W_M(i-1|x_1) \lambda_i + \sum_{j=1}^M W_M(j-1|x_1) \log W_M(j-1|x_1) \right. \\
&\quad \left. - \sum_{i=1}^M W_M(i-1|x_2) \lambda_i - \sum_{j=1}^M W_M(j-1|x_2) \log W_M(j-1|x_2) \right| \\
&\leq \left| \sum_{i=1}^M (W_M(i-1|x_1) - W_M(i-1|x_2)) \lambda_i \right| + |H(W_M(\cdot|x_1)) - H(W_M(\cdot|x_2))| \tag{A9a}
\end{aligned}$$

$$\leq \sum_{i=1}^M |(W_M(i-1|x_1) - W_M(i-1|x_2)) \lambda_i| + |H(W_M(\cdot|x_1)) - H(W_M(\cdot|x_2))| \tag{A9b}$$

$$\leq LM \|\lambda\|_1 |x_1 - x_2| + |H(W_M(\cdot|x_1)) - H(W_M(\cdot|x_2))| \tag{A9c}$$

$$\leq LM^2 \left(\log \frac{1}{\gamma_M} \vee \frac{1}{\ln 2} \right) |x_1 - x_2| + |H(W_M(\cdot|x_1)) - H(W_M(\cdot|x_2))| \tag{A9d}$$

$$\leq LM^2 \left(\log \frac{1}{\gamma_M} \vee \frac{1}{\ln 2} \right) |x_1 - x_2| + ML \left| \log \frac{1}{\gamma_M} - \frac{1}{\ln 2} \right| |x_1 - x_2|. \tag{A9e}$$

Inequalities (A9a) and (A9b) use the triangle inequality. Inequality (A9c) follows by Assumption 3.2(ii) and (A9d) can be derived by following the proof of Lemma 3.7, which is similar to the

one of Lemma 2.4. Finally, (A9e) follows from the fact that the function $\Delta_n \ni x^n \mapsto H(x^n) \in \mathbb{R}_{\geq 0}$ with $\min_{1 \leq i \leq n} x_i < c$ is Lipschitz continuous with constant $n \left| \log \frac{1}{c} - \frac{1}{\ln 2} \right|$ and from Assumption 3.2(ii). \square

9. Proof of Lemma 3.15

We start by the following definitions that simplify the proof below

$$\begin{aligned} f_{\lambda, \nu}(x) &:= \mathcal{W}\lambda(x) - r(x) + \nu\mu_\nu s(x), & \bar{f}_{\lambda, \nu} &:= \max_{x \in \mathbb{A}} f_{\lambda, \nu}(x) \\ B_{\lambda, \nu}(\varepsilon) &:= \{x \in \mathbb{A} \mid \bar{f}_{\lambda, \nu} - f_{\lambda, \nu}(x) < \varepsilon\}, & \eta_{\lambda, \nu}(\varepsilon) &:= \int_{B_{\lambda, \nu}(\varepsilon)} dx. \end{aligned}$$

By the Lipschitz continuity of $f_\lambda(\cdot)$ and $s(\cdot)$ we get the uniform lower bound

$$\eta_{\lambda, \nu}(\varepsilon) \geq \frac{\varepsilon}{L_f + |\nu\mu_\nu|L_s} \wedge \rho. \quad (\text{A10})$$

By using the solution to $G_\nu(\lambda)$, according to (34) we can write

$$G_\nu(\lambda) = -\nu \log(\rho) + \nu \log \left(\int_{\mathbb{A}} 2^{\frac{1}{\nu}} f_{\lambda, \nu}(x) dx \right) \quad (\text{A11a})$$

$$\leq \inf_{\ell \in \mathbb{R}} \max_{x \in \mathbb{A}} \{f_\lambda(x) + \ell s(x)\} \quad (\text{A11b})$$

$$= G(\lambda), \quad (\text{A11c})$$

where the equality (A11c) follows as (A11b) is the dual program to $G(\lambda)$ and strong duality holds. The inequality (A11b) then is due to $G_\nu(\lambda) \leq G(\lambda)$ for any λ , see (36). Therefore,

$$G(\lambda) - G_\nu(\lambda) \leq \bar{f}_{\lambda, \nu} - G_\nu(\lambda) \quad (\text{A12a})$$

$$= \nu \left(-\log \left(\int_{B_{\lambda, \nu}(\varepsilon)} 2^{\frac{1}{\nu}} (f_{\lambda, \nu}(x) - \bar{f}_{\lambda, \nu}) dx + \int_{B_{\lambda, \nu}^c(\varepsilon)} 2^{\frac{1}{\nu}} (f_{\lambda, \nu}(x) - \bar{f}_{\lambda, \nu}) dx \right) + \log(\rho) \right) \quad (\text{A12b})$$

$$\leq \nu \left(-\log \left(\int_{B_{\lambda, \nu}(\varepsilon)} 2^{\frac{1}{\nu}} (f_{\lambda, \nu}(x) - \bar{f}_{\lambda, \nu}) dx \right) + \log(\rho) \right)$$

$$\leq \nu \left(-\log \left(\eta_{\lambda, \nu}(\varepsilon) 2^{-\frac{\varepsilon}{\nu}} \right) + \log(\rho) \right) \quad (\text{A12c})$$

$$\leq \nu \left(-\log \left(\frac{\varepsilon}{L_f + |\nu\mu_\nu|L_s} \vee \rho \right) + \frac{\varepsilon}{\nu} + \log(\rho) \right) \quad (\text{A12d})$$

$$= \nu \log \left(\frac{(L_f + |\nu\mu_\nu|L_s)\rho}{\varepsilon} \vee 1 \right) + \varepsilon,$$

where (A12a) follows from (A11c) and (A12b) is due to (A11a). The inequality (A12c) results from the definitions of $B_{\lambda, \nu}(\varepsilon)$ and $\eta_{\lambda, \nu}(\varepsilon)$ above and (A12d) is implied by (A10). Finally, it can be seen that for $\nu < (L_f + |\nu\mu_\nu|L_s)\rho$, the optimal choice for ε is ν , which leads to

$$G(\lambda) - G_\nu(\lambda) \leq \nu \left(1 + \log \left(\frac{(L_f + |\nu\mu_\nu|L_s)\rho}{\nu} \vee 1 \right) \right). \quad (\text{A13})$$

It remains to upper bound the term $|\nu\mu_\nu|$. Define $\underline{f} := \min_{x, \lambda} f_\lambda(x)$, $\bar{f} := \max_{x, \lambda} f_\lambda(x)$, $\Delta_f := \bar{f} - \underline{f}$ and note that $\Delta_f \leq L_f \rho$. By (A11a), (36) and the fact that adding an additional constraint to a maximization problem cannot increase its objective value

$$G_\nu(\lambda) = \nu \log \left(\int_{\mathbb{A}} 2^{\frac{1}{\nu}} (f_\lambda(x) + \nu\mu_\nu s(x)) dx \right) - \nu \log(\rho) \leq \bar{f} = \nu \log \left(2^{\frac{1}{\nu}} \bar{f} \right),$$

which is equivalent to $\int_{\mathbb{A}} 2^{\frac{1}{\nu}(f_{\lambda}(x) - \bar{f} + \nu\mu_{\nu}s(x))} dx \leq \rho$ and implies

$$\int_{\mathbb{A}} 2^{\mu_{\nu}s(x)} dx \leq \rho 2^{\frac{\Delta_f}{\nu}}. \quad (\text{A14})$$

From (A14) two bounds can be derived. First, (A14) implies that $(\rho \wedge \frac{\varepsilon}{L_s}) 2^{\mu_{\nu}(\bar{s}-\varepsilon)} \leq \rho 2^{\frac{\Delta_f}{\nu}}$, which by choosing $\varepsilon = \frac{\bar{s}}{2}$ leads to $2^{\mu_{\nu}\frac{\bar{s}}{2}} \leq \left(\frac{2L_s\rho}{\bar{s}} \vee 1\right) 2^{\frac{\Delta_f}{\nu}}$ and finally,

$$\nu\mu_{\nu} \leq \frac{2}{\bar{s}} \log\left(\frac{2L_s\rho}{\bar{s}} \vee 1\right) \nu + \frac{2\Delta_f}{\bar{s}}. \quad (\text{A15})$$

Similarly one can derive a lower bound

$$\nu\mu_{\nu} \geq \frac{2}{\bar{s}} \log\left(\frac{2L_s\rho}{-\bar{s}} \vee 1\right) \nu + \frac{2\Delta_f}{\bar{s}}. \quad (\text{A16})$$

Equation (A13) together with (A15) and (A16) complete the proof. \square

10. Proof of Theorem 3.16

Following [15] and using Lemma 3.7, Lemma 3.5, Propostion 3.9 and Lemma 3.15, after n iterations of Algorithm 1 the following approximation error is obtained

$$0 \leq F(\hat{\lambda}) + G(\hat{\lambda}) - I(\hat{p}, W) \leq \iota(\nu) + \frac{4D_1}{\nu(n+1)^2} + \frac{4D_1}{(n+1)^2} =: \text{err}(\nu, n), \quad (\text{A17})$$

where for $\nu < \frac{T_1}{1-T_2}$ or $T_2 > 1$ we have $\iota(\nu) = \nu \left(\log\left(\frac{T_1}{\nu} + T_2\right) + 1\right)$, which is strictly increasing in ν . Let us redefine the smoothing term by $\nu := \frac{\delta}{\log(\delta^{-1})}$ for $\delta \in (0, 1)$ and define the function

$g(\delta) := \left(\frac{\log(T_1 \log(\delta^{-1}) + T_2 \delta) + 1}{\log(\delta^{-1})} + 1\right)$. One can see that $\iota(\nu) = \delta g(\delta)$ and that $\lim_{\delta \rightarrow 0} g(\delta) = 1$.

Furthermore $\delta \leq 2^{-1} \wedge 2^{-\frac{1}{T_1+T_2}}$ implies

$$g(\delta) - 1 \leq \frac{\log(2(T_1 + T_2) \log(\delta^{-1}))}{\log(\delta^{-1})} \leq T_1 + T_2, \quad (\text{A18})$$

where the first inequality is due to $\delta \leq 2^{-1}$ and the second follows from $\delta \leq 2^{-\frac{1}{T_1+T_2}}$. We seek for a lower bound of n and upper bound δ such that the error term (A17) is smaller than the preassigned $\varepsilon > 0$, i.e.,

$$\text{err}\left(\frac{\delta}{\log(\delta^{-1})}, n\right) = g(\delta)\delta + \frac{4D_1}{(n+1)^2} \left(\frac{\log(\delta^{-1})}{\delta} + 1\right) \leq \varepsilon \quad (\text{A19})$$

To this end, we introduce an auxiliary variable $\zeta \in (0, 1)$ such that $g(\delta)\delta = (1 - \zeta)\varepsilon$ and $\frac{4D_1}{(n+1)^2} \left(\frac{\log(\delta^{-1})}{\delta} + 1\right) \leq \zeta\varepsilon$, which implies (A19). Observe that $g(\delta)\delta = (1 - \zeta)\varepsilon$ is equivalent to $\delta = \frac{(1-\zeta)\varepsilon}{g(\delta)} =: \beta\varepsilon$. Hence $\zeta = 1 - \beta g(\delta)$ for $\beta \in [0, \frac{1}{g(\delta)}]$. Moreover,

$$\frac{4D_1}{(n+1)^2} \left(\frac{\log(\delta^{-1})}{\delta} + 1\right) = \frac{4D_1}{(n+1)^2} \left(\frac{\log((\beta\varepsilon)^{-1})}{\beta\varepsilon} + 1\right) \leq \zeta\varepsilon$$

is equivalent to

$$4D_1 \left(\frac{\log((\beta\varepsilon)^{-1}) + \beta\varepsilon}{\beta(1-g(\delta)\beta)\varepsilon^2} \right) = 4D_1 \left(\frac{\log(\varepsilon^{-1}) + \log 2g(\delta) + \frac{\varepsilon}{2g(\delta)}}{\frac{\varepsilon^2}{4g(\delta)}} \right) \leq (n+1)^2, \quad (\text{A20})$$

where we have chosen $\beta = \frac{1}{2g(\delta)}$ and as such is equivalent to

$$\frac{4}{\varepsilon} \sqrt{D_1 (g(\delta) \log(\varepsilon^{-1}) + g(\delta) \log(2g(\delta)) + \frac{\varepsilon}{2})} \leq n+1$$

Finally, using (A18) implies for $\nu = \frac{\varepsilon/\alpha}{\log(\alpha/\varepsilon)}$, where $\alpha := 2(T_1 + T_2 + 1)$

$$\text{err}(\nu, n) \leq \varepsilon \quad \text{for} \quad n \geq \frac{1}{\varepsilon} \sqrt{8D_1 \alpha} \sqrt{\log(\varepsilon^{-1}) + \log(\alpha) + \frac{1}{4}}.$$

□

11. Proof of Proposition 3.19

To prove Proposition 3.19, we need two lemmas.

Lemma A.3. *For any $k \in (0, 1]$ and $a, b \geq 0$*

$$a^k + b^k \leq 2^{1-k}(a+b)^k.$$

Proof. Let $g(x) := 2^{1-k}(1+x)^k - x^k$. By setting $\frac{d}{dx}g(x^*) = 0$, one can easily see that $x^* = 1$ is the minimizer of function g over the interval $[0, 1]$, i.e., $g(x) \geq g(1) = 1$ for all $x \in [0, 1]$. Suppose, without loss of generality, that $a \geq b$. By virtue of the preceding result of function g , we know that

$$1 \leq g\left(\frac{b}{a}\right) = 2^{1-k} \left(1 + \frac{b}{a}\right)^k - \left(\frac{b}{a}\right)^k,$$

where by multiplying a^k it readily leads to the desired assertion. □

Lemma A.4. *Let $(a_i)_{i \in \mathbb{N}}$ be a non-negative sequence of real numbers. For any $k \in (0, 1]$*

$$\sum_{i \in \mathbb{N}} a_i^k \leq \left(\sum_{i \in \mathbb{N}} \alpha^i a_i \right)^k, \quad \alpha := 2^{(k^{-1}-1)}.$$

Proof. For the proof we make use of an induction argument. Note that for any $a_1 \geq 0$ it trivially holds that $a_1^k \leq 2^{1-k} a_1^k$. We now assume that for any sequence $(a_i)_{i=1}^N \subset \mathbb{R}_{\geq 0}$ we have

$$\sum_{i=1}^N a_i^k \leq \left(\sum_{i=1}^N 2^{(k^{-1}-1)i} a_i \right)^k. \quad (\text{A21})$$

Let $(a_i)_{i=1}^{N+1} \subset \mathbb{R}_{\geq 0}$. Then,

$$\begin{aligned} \sum_{i=1}^{N+1} a_i^k &= a_1^k + \sum_{i=2}^{N+1} a_i^k \leq a_1^k + \left(\sum_{i=2}^{N+1} 2^{(k^{-1}-1)(i-1)} a_i \right)^k \leq 2^{1-k} \left(a_1 + \sum_{i=2}^{N+1} 2^{(k^{-1}-1)(i-1)} a_i \right)^k \\ &= \left(2^{(k^{-1}-1)} a_1 + \sum_{i=2}^{N+1} 2^{(k^{-1}-1)i} a_i \right)^k = \left(\sum_{i=1}^{N+1} 2^{(k^{-1}-1)i} a_i \right)^k, \end{aligned} \quad (\text{A22})$$

where the first (resp. second) inequality in (A22) follows from (A21) (resp. Lemma A.3). □

Proof of Proposition 3.19. It is straightforward to see that

$$\max_{x \in [0, A]} e^{-x} x^i = e^{-\min\{A, i\}} (\min\{A, i\})^i. \quad (\text{A23})$$

Moreover, based on a Taylor series expansion, it is well known that for all $M \in \mathbb{N}$ and $x \in \mathbb{R}_{\geq 0}$

$$\sum_{i \geq M} \frac{x^i}{i!} \leq \frac{e^x}{M!} x^M. \quad (\text{A24})$$

Therefore, it follows that

$$R_k(M) := \sum_{i \geq M} \left(\sup_{x \in [0, A]} e^{-(x+\eta)} \frac{(x+\eta)^i}{i!} \right)^k \leq \sum_{i \geq M} \left(e^{-(A+\eta)} \frac{(A+\eta)^i}{i!} \right)^k \quad (\text{A25a})$$

$$\leq e^{-k(A+\eta)} \left(\sum_{i \geq M} \alpha^{(i-M+1)} \frac{(A+\eta)^i}{i!} \right)^k = \frac{e^{-k(A+\eta)}}{\alpha^{k(M-1)}} \left(\sum_{i \geq M} \frac{(\alpha(A+\eta))^i}{i!} \right)^k \quad (\text{A25b})$$

$$\leq \frac{e^{-k(A+\eta)}}{\alpha^{k(M-1)}} \left(\frac{e^{\alpha(A+\eta)}}{M!} \alpha^M (A+\eta)^M \right)^k = \left(\alpha e^{(\alpha-1)(A+\eta)} \frac{(A+\eta)^M}{M!} \right)^k, \quad (\text{A25c})$$

where (A25a) results from (A23) and the assumption $M \geq A + \eta$, and (A25b) (resp. (A25c)) follows from Lemma A.4 (resp. (A24)). \square

Appendix B: Simulation Details

This section provides some further details on the simulation in Example 3.20. The parameters considered are $k = \frac{1}{2}$, $L_f = 0$ and M is chosen according to Table III. All the simulations in this section are performed on a 2.3 GHz Intel Core i7 processor with 8 GB RAM with Matlab.

TABLE III. Simulation details to Example 3.20

A [dB]	0	1	2	3	4	5	6	7
M	16	17	19	20	22	25	28	31
Iterations n	$4 \cdot 10^4$	$4 \cdot 10^4$	$4 \cdot 10^4$	$5 \cdot 10^4$	$6 \cdot 10^4$	$7 \cdot 10^4$	$9 \cdot 10^4$	$1.2 \cdot 10^5$
ν	0.0026	0.0029	0.0036	0.0029	0.0027	0.0029	0.0026	0.0022
$F(\hat{\lambda}) + G(\hat{\lambda})$	0.1144	0.1626	0.2263	0.3063	0.4029	0.5129	0.6293	0.7423
$I(\hat{p}, W)$	0.1105	0.1583	0.2206	0.3015	0.3979	0.5072	0.6234	0.7365
\mathcal{E}	$9.3 \cdot 10^{-4}$	$9.7 \cdot 10^{-4}$	$4.8 \cdot 10^{-4}$	$8.5 \cdot 10^{-4}$	$8.2 \cdot 10^{-4}$	$4.9 \cdot 10^{-4}$	$5.0 \cdot 10^{-4}$	$9.5 \cdot 10^{-4}$

A [dB]	8	9	10	11	12	13	14
M	36	42	49	59	71	85	104
Iterations n	$2 \cdot 10^5$	$5 \cdot 10^5$	$2 \cdot 10^6$	$3 \cdot 10^6$	$4 \cdot 10^6$	$9 \cdot 10^6$	$1.5 \cdot 10^7$
ν	0.0016	$7.1 \cdot 10^{-5}$	$8.0 \cdot 10^{-4}$	$8.3 \cdot 10^{-4}$	$9.7 \cdot 10^{-4}$	$6.2 \cdot 10^{-4}$	$5.8 \cdot 10^{-4}$
$F(\hat{\lambda}) + G(\hat{\lambda})$	0.8410	0.9422	1.0591	1.1835	1.3070	1.4343	1.5671
$I(\hat{p}, W)$	0.8351	0.9388	1.0547	1.1788	1.3013	1.4219	1.5605
\mathcal{E}	$7.5 \cdot 10^{-4}$	$7.1 \cdot 10^{-4}$	$8.0 \cdot 10^{-4}$	$6.2 \cdot 10^{-4}$	$5.2 \cdot 10^{-4}$	$9.0 \cdot 10^{-4}$	$6.7 \cdot 10^{-4}$

ACKNOWLEDGMENT

The authors thank Yurii Nesterov, Renato Renner and Stefan Richter for helpful discussions and pointers to references.

-
- [1] Claude E. Shannon, “A mathematical theory of communication,” *Bell System Technical Journal* **27**, 379–423 (1948).
 - [2] Richard E. Blahut, “Computation of channel capacity and rate-distortion functions,” *IEEE Transactions on Information Theory* **18**, 460–473 (1972).
 - [3] Robert G. Gallager, *Information Theory and Reliable Communication* (John Wiley & Sons, 1968).
 - [4] Suguru Arimoto, “An algorithm for computing the capacity of arbitrary discrete memoryless channels,” *IEEE Transactions on Information Theory* **18**, 14–20 (1972).
 - [5] Ibrahim C. Abou-Faycal, Mitchell D. Trott, and Shlomo Shamai, “The capacity of discrete-time memoryless Rayleigh-fading channels,” *IEEE Transactions on Information Theory* **47**, 1290–1301 (2001).
 - [6] Jossy Sayir, “Iterating the Arimoto-Blahut algorithm for faster convergence,” *Proceedings IEEE International Symposium on Information Theory (ISIT)*, 235 (2000).
 - [7] Gerald Matz and Pierre Duhamel, “Information geometric formulation and interpretation of accelerated Blahut-Arimoto-type algorithms,” *Proceedings Information Theory Workshop (ITW)*, 66–70 (2004).
 - [8] Yaming Yu, “Squeezing the Arimoto-Blahut algorithm for faster convergence,” *IEEE Transactions on Information Theory* **56**, 3149–3157 (2010).
 - [9] Justin Dauwels, “Numerical computation of the capacity of continuous memoryless channels,” *Proceedings of the 26th Symposium on Information Theory in the BENELUX*, 221–228 (2005).
 - [10] Jihai Cao, S. Hranilovic, and Jun Chen, “Capacity and nonuniform signaling for discrete-time poisson channels,” *Optical Communications and Networking, IEEE/OSA Journal of* **5**, 329–337 (2013).
 - [11] Jihai Cao, S. Hranilovic, and Jun Chen, “Capacity-achieving distributions for the discrete-time poisson channel - part 1: General properties and numerical techniques,” *Communications, IEEE Transactions on* **62**, 194–202 (2014).
 - [12] Jihai Cao, S. Hranilovic, and Jun Chen, “Capacity-achieving distributions for the discrete-time poisson channel - part 2: Binary inputs,” *Communications, IEEE Transactions on* **62**, 203–213 (2014).
 - [13] Chiang Mung and Stephen Boyd, “Geometric programming duals of channel capacity and rate distortion,” *IEEE Transactions on Information Theory* **50**, 245–258 (2004).
 - [14] Jianyi Huang and Sean P. Meyn, “Characterization and computation of optimal distributions for channel coding,” *IEEE Transactions on Information Theory* **51**, 2336–2351 (2005).
 - [15] Yurii Nesterov, “Smooth minimization of non-smooth functions,” *Mathematical Programming* **103**, 127–152 (2005).
 - [16] Amos Lapidoth and Stefan M. Moser, “On the capacity of the discrete-time Poisson channel,” *IEEE Transactions on Information Theory* **55**, 303–322 (2009).
 - [17] Aharon Ben-Tal and Marc Teboulle, “Extension of some results for channel capacity using a generalized information measure,” *Applied Mathematics and Optimization* **17**, 121–132 (1988).
 - [18] Chiang Mung, “Geometric programming for communication systems,” *Foundations and Trends in Communications and Information Theory* **2**, 1–154 (2005).
 - [19] Yurii Nesterov, *Introductory Lectures on Convex Optimization: A Basic Course*, Applied Optimization (Springer, 2004).
 - [20] Dimitri P. Bertsekas, *Convex Optimization Theory*, Athena Scientific optimization and computation series (Athena Scientific, 2009).
 - [21] H.S. Witsenhausen, “Some aspects of convexity useful in information theory,” *IEEE Transactions on Information Theory* **26**, 265–271 (1980).
 - [22] J. M. Borwein and A. S. Lewis, “Duality relationships for entropy-like minimization problems,” *SIAM J. Control Optim.* **29**, 325–338 (1991).
 - [23] Jean B. Lasserre, *Moments, Positive Polynomials and Their Applications*, Imperial College Press optimization series (Imperial College Press, 2009).

- [24] Debbie Leung and Graeme Smith, “Continuity of quantum channel capacities,” *Communications in Mathematical Physics* **292**, 201–215 (2009).
- [25] Thomas M. Cover and Joy A. Thomas, *Elements of Information Theory* (Wiley Interscience, 2006).
- [26] Stefan M. Moser, “Duality-based bounds on channel capacity,” PhD thesis, ETH Zurich (2005).
- [27] Shlomo Shamai, “Capacity of a pulse amplitude modulated direct detection photon channel,” IEE Proceedings on Communications, Speech and Vision **137**, 424–430 (1990).
- [28] D. H. Fremlin, *Measure theory. Vol. 2* (Torres Fremlin, Colchester, 2010) pp. 563+12 pp. (errata), broad foundations, Second edition January 2010.
- [29] Edward J. Anderson and Peter Nash, *Linear programming in infinite-dimensional spaces: theory and applications*, Wiley-Interscience Series in Discrete Mathematics and Optimization (Wiley, 1987).
- [30] Sanjoy K. Mitter, “Convex optimization in infinite dimensional spaces,” in *Recent advances in learning and control*, Lecture Notes in Control and Inform. Sci., Vol. 371 (Springer, London, 2008) pp. 161–179.
- [31] Olivier Devolder, François Glineur, and Yurii Nesterov, “First-order methods of smooth convex optimization with inexact oracle,” *Mathematical Programming*, 1–39 (2013).
- [32] David Sutter, Tobias Sutter, Peyman Mohajerin Esfahani, and Renato Renner, “Efficient approximation of quantum channel capacities,” (2014), available at [arXiv:1407.8202](https://arxiv.org/abs/1407.8202).
- [33] David Brady and Sergio Verdú, “The asymptotic capacity of the direct detection photon channel with a bandwidth constraint,” in *28th Annual Allerton Conference on Communication, Control, and Computing* (1990) pp. 691–700.
- [34] Alfonso Martinez, “Spectral efficiency of optical direct detection,” *J. Opt. Soc. Am. B* **24**, 739–749 (2007).
- [35] Amos Lapidoth, Jeffrey H. Shapiro, Vinodh Venkatesan, and Ligong Wang, “The discrete-time Poisson channel at low input powers,” *IEEE Transactions on Information Theory* **57**, 3260–3272 (2011).
- [36] J. Singh, O. Dabeer, and U. Madhow, “On the limits of communication with low-precision analog-to-digital conversion at the receiver,” *IEEE Transactions on Communications* **57**, 3629–3639 (2009).
- [37] Tobias Koch and Amos Lapidoth, “At low SNR, asymmetric quantizers are better,” *IEEE Transactions on Information Theory* **59**, 5421–5445 (2013).
- [38] Michel Baes and Michael Bürgisser, “An acceleration procedure for optimal first-order methods,” *Optimization Methods and Software* **29**, 610–628 (2014).
- [39] Alexander S. Holevo, *Quantum Systems, Channels, Information* (De Gruyter Studies in Mathematical Physics 16, 2012).
- [40] Stephen Boyd and Lieven Vandenberghe, *Convex Optimization* (Cambridge University Press, Cambridge, 2004) pp. xiv+716, sixth printing with corrections, 2008.
- [41] Gerald B. Folland, *Real analysis: modern techniques and their applications*, Pure and applied mathematics (Wiley, 1999).
- [42] Yihong Wu and Sergio Verdú, “Functional properties of minimum mean-square error and mutual information,” *IEEE Transactions on Information Theory* **58**, 1289–1301 (2012).
- [43] Patrick Billingsley, *Convergence of probability measures*, Wiley Series in probability and Mathematical Statistics: Tracts on probability and statistics (Wiley, 1968).
- [44] Olivier Devolder, François Glineur, and Yurii Nesterov, “Double smoothing technique for large-scale linearly constrained convex optimization,” *SIAM Journal on Optimization* **22**, 702–727 (2012).