# Optimal Quantization for Distribution Synthesis

Georg Böcherer and Bernhard C. Geiger

## Abstract

Finite precision approximations of discrete probability distributions are considered, applicable for distribution synthesis, e.g., probabilistic shaping. Two algorithms are presented that find the optimal $M$-type approximation $Q$ of a distribution $P$ in terms of the variational distance $\|Q - P\|_1$ and the informational divergence $D(Q\|P)$. Bounds on the approximation errors are derived and shown to be asymptotically tight. Several examples illustrate that the variational distance optimal approximation can be quite different from the informational divergence optimal approximation.

## Index Terms

distribution synthesis, distribution quantization, $M$-type approximation, variational distance, informational divergence, Kullback-Leibler divergence.

## I. INTRODUCTION

Probabilistic models are often used for information processing. In practice, such models are represented with finite precision, e.g., discrete probabilities are represented by rational numbers with finitely many digits. If each probability can be written as an integer multiple of $1/M$ for some integer $M$, then the resulting distribution is called an *M-type distribution*. The integer $M$ characterizes the precision by which the rational distribution approximates the true distribution. Additionally, $M$ influences the space needed to store the rational distribution and the complexity to process it. This work studies approximating target distributions $\boldsymbol{t} = (t_1, t_2, \dots)$ by $M$-type distributions.

### A. Quality-of-Synthesis Criteria

One way to measure how good $\hat{\boldsymbol{t}}$ approximates $\boldsymbol{t}$ is the variational distance

$$\|\boldsymbol{t} - \hat{\boldsymbol{t}}\|_1 = \sum_i |t_i - \hat{t}_i| \tag{1}$$

which is symmetric in its arguments $\boldsymbol{t}, \hat{\boldsymbol{t}}$. Another criterion is the informational divergence

$$\mathbb{D}(\hat{\boldsymbol{t}}\|\boldsymbol{t}) = \sum_{i\,:\,\hat{t}_i > 0} \hat{t}_i \log \frac{\hat{t}_i}{t_i} \tag{2}$$

where the expectation is taken w.r.t. the approximating distribution $\hat{\boldsymbol{t}}$. The informational divergence with exchanged order of arguments is

$$\mathbb{D}(\boldsymbol{t}\|\hat{\boldsymbol{t}}) = \sum_{i\,:\,t_i > 0} t_i \log \frac{t_i}{\hat{t}_i}. \tag{3}$$

Note that the expectation in (3) is taken with respect to the target distribution $\boldsymbol{t}$. The informational divergence is asymmetric, i.e., (2) and (3) are different in general.

In this work we are interested in the scenario where the approximating distribution $\hat{\boldsymbol{t}}$ *synthesizes* the distribution $\boldsymbol{t}$, i.e, we take expectation with respect to the approximating distribution $\hat{\boldsymbol{t}}$. We will therefore consider (1) and (2) as quality-of-synthesis criteria. Several rationales for this choice are as follows:

*1) Empirical Probability:* In distribution synthesis, the approximation $\hat{\boldsymbol{t}}$ is the "true" distribution and describes a random experiment where the random variable $I$ takes on the integer values $1, 2, 3, \ldots$ according to $\hat{\boldsymbol{t}}$, i.e., $\Pr(I = i) = \hat{t}_i$. Denote by $i_1, i_2, \ldots, i_m$ the outcomes of performing the random experiment $m$ times. By the law of large numbers,

$$\frac{\sum_{j=1}^{m} \log \frac{\hat{t}_{i_j}}{t_{i_j}}}{m} \approx \mathbb{D}(\hat{\boldsymbol{t}} \| \boldsymbol{t}). \tag{4}$$

There is no such interpretation for the measures (1) and (3).

*2) Infinite Support:* Many important probability distributions have infinite support, e.g., Poisson, Boltzmann, Borel, and Yule-Simon distributions. $M$-type distributions have finite support, and if the target distribution $\boldsymbol{t}$ has infinite support, then the measure (3) is infinity. The measures (1) and (2) do not have this issue.

*3) Probabilistic Shaping:* Suppose the target distribution $\boldsymbol{t}$ is the capacity-achieving input distribution of some communication channel and suppose further that $\hat{\boldsymbol{t}}$ is the actual input distribution generated by a communication system. Denote by $W$ the transition probability matrix of the channel. The mutual information $\mathbb{I}(\hat{\boldsymbol{t}}, W)$ that results from using the approximation $\hat{\boldsymbol{t}}$ at the channel input is bounded as

$$\mathsf{C} \geq \mathbb{I}(\hat{\boldsymbol{t}}, W) \overset{\text{(a)}}{=} \mathsf{C} - \mathbb{D}(\hat{\boldsymbol{t}}W \| \boldsymbol{t}W)$$
$$\overset{\text{(b)}}{\geq} \mathsf{C} - \mathbb{D}(\hat{\boldsymbol{t}} \| \boldsymbol{t}) \tag{5}$$

where $\hat{\boldsymbol{t}}W$ and $\boldsymbol{t}W$ are the output distributions that result from the input distributions $\hat{\boldsymbol{t}}$ and $\boldsymbol{t}$, respectively, and where $\mathsf{C}$ is the capacity of the channel. The equality in (a) follows by [2, Sec. III],[3, Proposition 3.11] and (b) by the data processing inequality [4, Lemma 3.11]. The bound (5) shows that as (2) approaches zero, the mutual information $\mathbb{I}(\hat{\boldsymbol{t}}, W)$ approaches capacity.

## B. Related Work

For probabilistic shaping, Gallager suggested in [5, p. 208] to choose $\hat{\boldsymbol{t}}$ as an $M$-type approximation of the capacity-achieving distribution $\boldsymbol{t}$. Several works propose to use *dyadic distributions* in Gallager's scheme, which are $M$-type distributions where $M$ is an integer power of two and where every probability can be written as $2^k/M$ for some integer $k$. The authors in [6] calculate a dyadic approximation by rounding the entries of $\boldsymbol{t}$, which minimizes the variational distance (1). The authors in [7] calculate the dyadic approximation of $\boldsymbol{t}$ that minimizes (2) by Geometric Huffman Coding [2]. Gallager's scheme also works for $M$-type distributions that are not dyadic. In [8], the authors calculate an $M$-type distribution by a sub-optimal algorithm that aims at minimizing (3). In [1], we proposed to use $M$-type distributions that minimize (2) in Gallager's scheme.

The authors in [9], [10] propose a quantization algorithm that minimizes the variational distance (1), the Euclidean distance, and the $L_\infty$ norm. The authors also use a Taylor series approximation to analyze their algorithm in terms of the informational divergence (3) for $M$ significantly larger than the support size of the distribution.

Resolution coding uses an $M$-type input distribution to approximate a target output distribution [11]. For the identity channel, [11, Sec. III.B] constructs an $M$-type approximation that is asymptotically optimal for the variational distance (1). In [12, Sec. VI.A], informational divergence (2) optimal $M$-type approximations are constructed. The authors in [13] derive fundamental limits of resolution coding for the identity channel with respect to various approximation measures including (1) and a normalized version

of (2). For noisy channels, resolution coding with respect to variational distance (1) is considered in [11], informational divergence (2) is considered in [14] and a normalized version of (2) is considered in [11], [15]. Most of the work presented in [11], [13]–[15] focuses on fundamental limits, i.e., the existence of asymptotically optimal $M$-type approximations is shown but no practical algorithms to construct them are provided.

### C. Contributions and Outline

We propose two simple algorithms that find the $M$-type approximations $\boldsymbol{t}^{\mathrm{id}}$ and $\boldsymbol{t}^{\mathrm{vd}}$ minimizing the informational divergence (2) and the variational distance (1), respectively. We provide bounds on the approximation errors for target distributions with finite and countably infinite supports. The bounds are asymptotically tight, i.e., any target distribution can be approximated arbitrarily well by an $M$-type approximation with sufficiently large $M$. In Sec. V, we show that variational distance (1) and informational divergence (2) lead to fundamentally different $M$-type approximations. In particular, we provide an example where the variational distance optimal approximation $\boldsymbol{t}^{\mathrm{vd}}$ results in an informational divergence equal to one for arbitrarily large $M$. Furthermore, we show that the informational divergence minimizing approximation $\boldsymbol{t}^{\mathrm{id}}$ can have significantly smaller support size than the variational distance minimizing approximation $\boldsymbol{t}^{\mathrm{vd}}$.

## II. PRELIMINARIES

Let $\boldsymbol{t}$ be a target probability distribution with a finite or countably infinite support. We denote by $n$ the support size of $\boldsymbol{t}$. If the support is infinite, then $n = \infty$. Without loss of generality, we assume that $\boldsymbol{t}$ is ordered so that $t_1 \geq t_2 \geq \cdots$. We define the complement of the cumulative distribution function as

$$T_k := \sum_{i > k} t_i. \tag{6}$$

Let $M$ be a positive integer. A distribution $\boldsymbol{p}$ is $M$-type, if each entry can be written as $p_i = c_i/M$ for some non-negative integer $c_i \leq M$. We want to determine the $M$-type distribution $\boldsymbol{p}$ that best approximates the target distribution $\boldsymbol{t}$. Two quality measures for approximation are considered, namely, the informational divergence and the variational distance as defined in (2) and (1), respectively. Pinsker's inequality [16, Lem. 11.6.1] bounds the informational divergence from below in terms of the variational distance:

$$\|\boldsymbol{p} - \boldsymbol{t}\|_1 \leq \sqrt{2\,\mathbb{D}(\boldsymbol{p}\|\boldsymbol{t})}. \tag{7}$$

There have been several works on bounding the informational divergence from above in terms of the variational distance; see [17] for a recent improvement and an overview over available bounds. The most useful for our purposes is adapted from [18]:

**Lemma 1** ([18, Thm. 7]). *For two probability distributions $\boldsymbol{p}$ and $\boldsymbol{t}$,*

$$\mathbb{D}(\boldsymbol{p}\|\boldsymbol{t}) \leq \frac{1}{2} \frac{r \log r}{r - 1} \|\boldsymbol{p} - \boldsymbol{t}\|_1 \tag{8}$$

*where* $r := \sup\limits_{i \,:\, p_i > 0} \dfrac{p_i}{t_i} \geq 1$.

In Lemma 1 and throughout the remainder of this work, $\log$ denotes the natural logarithm.

Note that the upper bound (8) depends on the distributions not only via the variational distance $\|\boldsymbol{p} - \boldsymbol{t}\|_1$, but also via $r$. We therefore call (8) *distribution dependent*. Any reverse Pinsker's inequality must be distribution dependent, see [19, Sec. I.A]. Note further that Lemma 1 was refined in [17, Thm. 1].

**Algorithm 1.** Variational distance optimal approximation.

---

Initialize $\boldsymbol{t}^{\mathrm{vd}} = \boldsymbol{0}$
Compute $t_i^{\mathrm{vd}} \leftarrow \frac{\lfloor M t_i \rfloor}{M}$, $i = 1, \ldots, \min\{n, M\}$.
Compute $e_i \leftarrow t_i - t_i^{\mathrm{vd}}$, $i = 1, \ldots, \min\{n, M\}$.
Compute $L \leftarrow M - M \cdot \sum_{i=1}^{\min\{n, M\}} t_i^{\mathrm{vd}}$.
**repeat** $L$ times
   Choose $j = \min\{\underset{i}{\operatorname{argmax}}\, e_i\}$. //*choose the smallest index first.*
   Update $t_j^{\mathrm{vd}} \leftarrow t_j^{\mathrm{vd}} + \frac{1}{M}$.
   Update $e_j \leftarrow t_j - t_j^{\mathrm{vd}}$.
**end repeat**
Return $\boldsymbol{t}^{\mathrm{vd}}$.

---

## III. VARIATIONAL DISTANCE OPTIMAL QUANTIZATION

### A. Algorithm 1

An $M$-type approximation of a target distribution $\boldsymbol{t}$ can be calculated as follows. First, round off the entries of $\boldsymbol{t}$ and then distribute the remaining mass among the entries with the largest error. We call this method *Algorithm 1*, see the top of Page 4.

Formally, we first calculate the pre-approximation

$$\tilde{t}_i^{\mathrm{vd}} = \frac{\lfloor M t_i \rfloor}{M}, \quad i = 1, \ldots, n. \tag{9}$$

Note that in Algorithm 1 we can restrict this computation to the first $\min\{n, M\}$ indices of $\boldsymbol{t}$ since, by assumption, $\boldsymbol{t}$ is ordered, and since not more than $M$ masses can be distributed. Thus, if $n > M$, we can be sure that $t_i^{\mathrm{vd}} = 0$ for $i > M$.

In general, the entries of $\tilde{\boldsymbol{t}}^{\mathrm{vd}}$ do not sum to one. The pre-approximation gives rise to the non-negative errors

$$e_i := t_i - \tilde{t}_i^{\mathrm{vd}} \geq 0, \quad i = 1, \ldots, n \tag{10}$$

which sum to the *rest mass*

$$\sum_{i=1}^n e_i = \frac{L}{M} \tag{11}$$

for some integer $L$. Note that the rest mass is bounded as $0 \leq L \leq M$, and it is equal to zero if and only if the target distribution $\boldsymbol{t}$ is itself $M$-type.

**Example 1.** For the 2-type target distribution $\boldsymbol{t} = (\frac{1}{2}, \frac{1}{2})$ and $M = 2$, we have $\tilde{\boldsymbol{t}}^{\mathrm{vd}} = \boldsymbol{t}$ and rest mass 0, i.e., $L = 0$. For the 3-type target distribution $\boldsymbol{t} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ and $M = 2$, we have $\tilde{\boldsymbol{t}}^{\mathrm{vd}} = (0, 0, 0)$ and rest mass 1, i.e., $L = M$.

Let $\mathcal{L}$ be a set of the indices with the $|\mathcal{L}| = L$ largest error terms, i.e., we have

$$i \in \mathcal{L}, j \notin \mathcal{L} \Rightarrow e_i \geq e_j. \tag{12}$$

We distribute the remaining $L$ unit masses to the indices in $\mathcal{L}$, i.e., we choose

$$t_i^{\mathrm{vd}} = \begin{cases} \tilde{t}_i^{\mathrm{vd}} + \frac{1}{M}, & i \in \mathcal{L} \\ \tilde{t}_i^{\mathrm{vd}}, & \text{otherwise.} \end{cases} \tag{13}$$

Note that the set $\mathcal{L}$ is not unique, and consequently the approximation $\boldsymbol{t}^{\mathrm{vd}}$ is not unique either. We illustrate this by an example.

**Example 2.** Let $\boldsymbol{t} = (\frac{3}{4}, \frac{1}{4})$ and suppose $M = 2$. Then

$$\tilde{t}_1^{\mathrm{vd}} = \frac{1}{2}, \quad \tilde{t}_2^{\mathrm{vd}} = 0 \tag{14}$$

and

$$e_1 = e_2 = \frac{1}{4}. \tag{15}$$

Thus, either $\mathcal{L} = \{1\}$ or $\mathcal{L} = \{2\}$. The corresponding approximations $\boldsymbol{t}^{\mathrm{vd}} = (1, 0)$ and $\boldsymbol{t}^{\mathrm{vd}} = (1/2, 1/2)$ both lead to the same approximation error, namely $\|\boldsymbol{t}^{\mathrm{vd}} - \boldsymbol{t}\|_1 = \frac{1}{2}$.

Algorithm 1 resolves this ambiguity by taking entries with lower indices first. From now on, $\boldsymbol{t}^{\mathrm{vd}}$ denotes the unique $M$-type approximation of $\boldsymbol{t}$ that is calculated by Algorithm 1.

### B. Elementwise Properties

From (10) and (13), we see that for each index $i$, we have

$$|t_i - t_i^{\mathrm{vd}}| < \frac{1}{M} \tag{16}$$

and $\boldsymbol{t}^{\mathrm{vd}}$ is a *uniform approximation* of $\boldsymbol{t}$. Also by (10) and (13), it follows that the approximation $\boldsymbol{t}^{\mathrm{vd}}$ assigns no mass to entries of $\boldsymbol{t}$ that are equal to zero, i.e., we have

$$t_i = 0 \quad \Rightarrow \quad t_i^{\mathrm{vd}} = 0. \tag{17}$$

Furthermore, if $\boldsymbol{t}^{\mathrm{vd}}$ assigns zero mass to some entry $t_i$, then it also assigns zero mass to all entries smaller than $t_i$:

**Lemma 2.** $t_j < t_i$ *and* $t_i^{\mathrm{vd}} = 0 \Rightarrow t_j^{\mathrm{vd}} = 0$.

*Proof:* Assume $t_j < t_i$. In the pre-approximation step, Algorithm 1 ensures that $t_i^{\mathrm{vd}} \geq \frac{\lfloor Mt_i \rfloor}{M}$, hence $t_i^{\mathrm{vd}} = 0$ implies $1/M > t_i > t_j$. Thus, the errors after pre-approximation satisfy $e_i = t_i$, $e_j = t_j$, and $e_j < e_i$. Algorithm 1 can only assign a remaining unit mass to $t_j$ and not to $t_i$ if $e_j \geq e_i$. Whence, $t_j^{\mathrm{vd}} = 0$. ∎

To prove the optimality of Algorithm 1, we make use of the following lemma.

**Lemma 3.** *Let $\boldsymbol{t}$ be a target distribution with finite or countably infinite support and let $M$ be a positive integer. Every $M$-type approximation $\boldsymbol{p}$ of $\boldsymbol{t}$ that is optimal w.r.t. the variational distance satisfies (16).*

*Proof:* See Section VI-A. ∎

### C. Optimality of Algorithm 1 and Performance Bounds

**Proposition 1.** *Let $\boldsymbol{t}$ be an ordered target distribution with finite or countably infinite support and let $M$ be a positive integer. Among all $M$-type distributions $\boldsymbol{p}$, $\boldsymbol{p} = \boldsymbol{t}^{\mathrm{vd}}$ minimizes $\|\boldsymbol{p} - \boldsymbol{t}\|_1$.*

*Proof:* According to Lemma 3, any optimal approximation satisfies (16). Hence, any optimal approximation $\boldsymbol{p}^*$ can be written as

$$p_i^* = \begin{cases} \tilde{t}_i^{\mathrm{vd}} + \frac{1}{M}, & i \in \mathcal{L}' \\ \tilde{t}_i^{\mathrm{vd}}, & \text{otherwise} \end{cases} \tag{18}$$

where $\tilde{\boldsymbol{t}}^{\text{vd}}$ is the pre-approximation (9) and where $\mathcal{L}'$ is some set of indices with $|\mathcal{L}'| = L$, where $L$ is given by (11). We have

$$\|\boldsymbol{p}^* - \boldsymbol{t}\|_1 = \sum_{i \in \mathcal{L}'} \left( \frac{1}{M} - e_i \right) + \sum_{i \notin \mathcal{L}'} e_i \tag{19}$$

where the error terms $e_i$ are defined in (10). The residual (19) is minimized if $\mathcal{L}'$ consists of the indices of the $L$ largest error terms $e_i$. According to (12), the approximation calculated by Algorithm 1 has this property. ∎

We next bound the variational distance in terms of $M$. If the target distribution $\boldsymbol{t}$ has finite support of cardinality $n$, then

$$\begin{aligned}
\|\boldsymbol{t}^{\text{vd}} - \boldsymbol{t}\|_1 &= \sum_{i=1}^{n} |t_i^{\text{vd}} - t_i| \\
&\overset{(a)}{\leq} \sum_{i=1}^{n} \frac{1}{M} \\
&= \frac{n}{M}
\end{aligned} \tag{20}$$

where (a) follows by (16). For $n = \infty$, the bound (20) is infinity for any finite $M$. Thus, we need a different approach to derive a useful bound for the case of infinite support. The next lemma lets us tighten bound (20) if $M \geq n$ and it will also lead to a useful bound for $n = \infty$. The underlying observation is that we can apply Algorithm 1 also to a *sub-probability* distribution, i.e., a target vector whose entries are positive and sum to a value less than or equal to one.

**Lemma 4.** *Let $\boldsymbol{t}$ be an ordered sub-probability distribution with $k \leq M$ entries and total mass $1 - T_k$, and let $M$ be a positive integer. Then we have*

$$\|\boldsymbol{t}^{\text{vd}} - \boldsymbol{t}\|_1 \begin{cases} \leq \frac{k}{2M} + \frac{MT_k^2}{2k}, & \text{always} \\ = T_k, & \text{if } T_k \geq \frac{k}{M} \end{cases} \tag{21}$$

$$\leq \frac{k}{2M} + T_k. \tag{22}$$

Note that for $T_k = k/M$ both cases in (21) coincide.

*Proof:* The proof is given in Sec. VI-B. ∎

A distribution can be split into two sub-probability distributions, one containing the first $k$ indices, and one containing the tail of the distribution. More specifically, we can split $\boldsymbol{t}$ into two vectors $\boldsymbol{t}_{1:k}$ and $\boldsymbol{t}_{\text{tail}}$ with the same length but disjoint support sets: The entries of $\boldsymbol{t}_{1:k} := (t_1, \ldots, t_k, 0, 0, 0, \ldots)$ are zero for indices larger than $k$, while for $\boldsymbol{t}_{\text{tail}} := (0, 0, \ldots, t_{k+1}, \ldots, t_n)$ the first $k$ entries are zero. Let $\boldsymbol{t}_{1:k}^{\text{vd}}$ denote the approximation that results from applying Algorithm 1 to $\boldsymbol{t}_{1:k}$. We have

$$\|\boldsymbol{t}_{1:k}^{\text{vd}} - \boldsymbol{t}\|_1 = \|\boldsymbol{t}_{1:k}^{\text{vd}} - \boldsymbol{t}_{1:k}\|_1 + T_k \tag{23}$$

where $\|\boldsymbol{t}_{1:k}^{\text{vd}} - \boldsymbol{t}_{1:k}\|_1$ can be bounded by Lemma 4. This divide-and-conquer approach is useful when the number of entries of the target distribution exceeds the type $M$ of the approximating distribution. Approach (23) is also used in the proof of the following proposition, which states various bounds on the approximation error of $\boldsymbol{t}^{\text{vd}}$.

**Proposition 2.** *Let $\boldsymbol{t}$ be an ordered target distribution and let $M$ be a positive integer.*

1) *If $\boldsymbol{t}$ has finite support of cardinality $n \leq M$, then*

$$\|\boldsymbol{t}^{\text{vd}} - \boldsymbol{t}\|_1 \leq \frac{n}{2M}. \tag{24}$$

2) *If $t$ has finite or countably infinite support of cardinality $n > M$, then*

$$\|t^{\text{vd}} - t\|_1 \leq \frac{k}{2M}\left(1 + \frac{MT_k}{k}\right)^2 \leq \frac{2k}{M} \tag{25}$$

*where $k$ is the support size of $t^{\text{vd}}$.*

3) *For $n = \infty$, the support size $k$ of $t^{\text{vd}}$ satisfies $k \overset{M\to\infty}{\longrightarrow} \infty$ and $k/M \overset{M\to\infty}{\longrightarrow} 0$.*

*Proof:* The proof is given in Sec. VI-C. ∎

We next give examples that illustrate the tightness of the bounds.

**Example 3.** For $n < \infty$, the bound (24) is tight for a uniform target distribution and $M = 3n/2$. For $M < n$, the bound (25) is tight for, e.g., $M = 5$ and $t_1 = t_2 = t_3 = 4/15$ and $t_i < 1/15$ for all $i > 3$ ($n$ arbitrary).

*D. Asymptotic Optimality*

For target vectors with finitely many entries, the bound (24) guarantees that the approximation error of $t^{\text{vd}}$ can be made arbitrarily small by choosing $M$ large enough. The same is true for infinitely many entries. This follows by bound (25) together with Statement 3) of Proposition 2. Furthermore, by (16) the $M$-type approximation converges uniformly to the target distribution. We summarize these observations as a corollary to Proposition 2.

**Corollary 1.** *Let $t$ be an ordered target distribution with finite or countably infinite support. For $M \to \infty$, the approximation $t^{\text{vd}}$ converges uniformly to the target distribution $t$.*

For $M \geq n$ the variational distance decreases with $\mathcal{O}(1/M)$. For $M < n$ no such convergence guarantee can be given. This is illustrated in the next example.

**Example 4.** Consider the Yule-Simon distribution [20] with $t_i = \rho B(i, \rho + 1)$, where $\rho > 0$ and where $B(\cdot, \cdot)$ is the beta-function. Lemma 2 ensures that Algorithm 1 assigns unit masses to at most the first $M$ indices. For $M > 1$, we have

$$\|t^{\text{vd}} - t\|_1 = \sum_{i=1}^{\infty} |t_i^{\text{vd}} - t_i| \geq T_M$$

$$= MB(M, \rho + 1) \tag{26}$$

$$\geq \frac{K(\rho)}{(M + \rho + 1)^\rho} \tag{27}$$

where $K(\rho)$ is a positive constant that does not depend on $M$, see Sec. VI-D for the derivation. Thus, the convergence of Algorithm 1 is at best $\mathcal{O}(1/M^\rho)$.

## IV. INFORMATIONAL DIVERGENCE OPTIMAL QUANTIZATION

We now consider $M$-type quantization with respect to the informational divergence, i.e., we want to solve the problem

$$\begin{aligned} \underset{p}{\text{minimize}} \quad & \mathbb{D}(p\|t) \\ \text{subject to} \quad & p \text{ is } M\text{-type}. \end{aligned} \tag{28}$$

**Algorithm 2.** Informational divergence optimal quantization.

---

Initialize $c_i \leftarrow 0$, $i = 1, \ldots, n$.
**for** $m = 1, 2, \ldots, M$
   Choose $j = \min\{\arg\min_i \Delta_i(c_i + 1)\}$. //*choose the smallest index first.*
   Update $c_j \leftarrow c_j + 1$.
**end for**
Return $\boldsymbol{c}$.

---

### A. Equivalent Problem

Recall that each entry $p_i$ of an $M$-type distribution can be written as $p_i = c_i/M$ for some non-negative integer $c_i$. We have

$$
\mathbb{D}(\boldsymbol{p}\|\boldsymbol{t}) = \sum_{i:\, c_i > 0} \frac{c_i}{M} \log \frac{\frac{c_i}{M}}{t_i}
$$

$$
= \frac{1}{M}\left( \sum_{i:\, c_i > 0} c_i \log \frac{c_i}{t_i} \right) - \log M \tag{29}
$$

so that Problem (28) is equivalent to

$$
\begin{aligned}
\underset{c_1, \ldots, c_n}{\text{minimize}} \quad & \sum_{i:\, c_i > 0} c_i \log \frac{c_i}{t_i} \\
\text{subject to} \quad & c_i \in \{0, 1, 2, \ldots, M\}, \quad i = 1, \ldots, n \\
& \sum_{i=1}^{n} c_i = M.
\end{aligned} \tag{30}
$$

If $\boldsymbol{c}^*$ is a solution of Problem (30), then $\boldsymbol{p}^* = \boldsymbol{c}^*/M$ is a solution of Problem (28).

### B. Algorithm 2

To solve problem (30), we write the objective function as a telescoping sum

$$
\sum_{i:\, c_i > 0} c_i \log \frac{c_i}{t_i} = \sum_{i=1}^{n} \sum_{k=1}^{c_i} \left[ k \log \frac{k}{t_i} - (k-1) \log \frac{k-1}{t_i} \right]
$$

$$
= \sum_{i=1}^{n} \sum_{k=1}^{c_i} \Delta_i(k) \tag{31}
$$

where the increment function is

$$
\Delta_i(k) = k \log k - (k-1) \log(k-1) + \log \frac{1}{t_i}. \tag{32}
$$

Evaluating $\Delta_i(x)$ as a function of a real number $x$ and taking the derivative,

$$
\frac{\partial}{\partial x} \Delta_i(x) = \log \frac{x}{x-1}, \tag{33}
$$

we conclude that $\Delta_i(k)$ is strictly monotonically increasing in $k$. Moreover, rewriting (32) as

$$
\Delta_i(k) = k \log \frac{k}{k-1} + \log(k-1) + \log \frac{1}{t_i}
$$

$$
\geq \log(k-1). \tag{34}
$$

(which holds trivially for $k = 1$) shows that the increment function grows without bound with $k$. The following lemma summarizes the properties of the increment function.

**Lemma 5.** *For all $m > 0$, the increment function $\Delta_i(k)$ grows without bound with $k$ and satisfies*

$$\ell > m \Rightarrow \Delta_i(\ell) > \Delta_i(m) \tag{35}$$

$$t_i > t_j \Rightarrow \Delta_i(m) < \Delta_j(m). \tag{36}$$

An allocation $\boldsymbol{c}$ can be obtained by initially assigning the zero vector $\boldsymbol{0}$ to a pre-allocation $\tilde{\boldsymbol{c}}$ and successively incrementing the entry of $\tilde{\boldsymbol{c}}$ by one for which the corresponding increment cost $\Delta(\tilde{c}_i + 1)$ is smallest. After $M$ iterations, the constraint $\sum_i \tilde{c}_i = M$ is fulfilled and $\boldsymbol{c} = \tilde{\boldsymbol{c}}$ is a valid allocation. If more than one entry of $\tilde{\boldsymbol{c}}$ has the smallest increment cost in some step, then either of them can be chosen, so the allocation obtained by this strategy is not unique. We illustrate this by the following example.

**Example 5.** Suppose $\boldsymbol{t} = (\frac{4}{5}, \frac{1}{5})$ and $M = 2$. We have $\Delta_1(1) = \log \frac{5}{4}$ and $\Delta_2(1) = \log 5$, so after the first step, $\tilde{\boldsymbol{c}} = (1, 0)$. In the second step, we have

$$\Delta_1(2) = 2 \log(2) + \log \frac{5}{4} = \log 5, \quad \Delta_2(1) = \log 5, \tag{37}$$

so the final allocation is either $\boldsymbol{c}_1 = (2, 0)$ or $\boldsymbol{c}_2 = (1, 1)$. The corresponding approximations are $\boldsymbol{p}_1 = (1, 0)$ and $\boldsymbol{p}_2 = (\frac{1}{2}, \frac{1}{2})$. Both approximations lead to the same informational divergence, namely

$$\mathbb{D}(\boldsymbol{p}_1 \| \boldsymbol{t}) = \mathbb{D}(\boldsymbol{p}_2 \| \boldsymbol{t}) = \log \frac{5}{4}. \tag{38}$$

Algorithm 2 resolves this ambiguity by incrementing entries with lower index first. From now on, we denote by $\boldsymbol{t}^{\mathrm{id}}$ the unique $M$-type approximation of $\boldsymbol{t}$ that is calculated by Algorithm 2.

### C. Elementwise Properties

The informational divergence is a weighted sum of $\log \frac{t_i^{\mathrm{id}}}{t_i}$. We therefore expect that for a good approximation $\boldsymbol{t}^{\mathrm{id}}$, the ratio $t_i^{\mathrm{id}}/t_i$ is close to one. The next lemma states this property.

**Lemma 6.** *Let $\boldsymbol{t}$ be a target distribution with finite or countably infinite support and let $M$ be a positive integer. Every $M$-type approximation $\boldsymbol{p}$ of $\boldsymbol{t}$ that is optimal w.r.t. the informational divergence satisfies*

$$\frac{p_i}{t_i} < \frac{e}{t_1}, \quad \forall i \leq k \tag{39}$$

*where $k$ is the support size of $\boldsymbol{p}$. In particular*

$$\frac{1}{M t_k} \leq \frac{e}{t_1}. \tag{40}$$

*Proof:* See Section VI-E. ∎

Lemma 6 directly implies

$$t_i = 0 \Rightarrow t_i^{\mathrm{id}} = 0. \tag{41}$$

Furthermore, if $\boldsymbol{t}^{\mathrm{id}}$ assigns zero mass to some entry $t_i$, then it also assigns zero mass to all entries smaller than $t_i$:

**Lemma 7.** $t_j < t_i$ *and* $t_i^{\mathrm{id}} = 0 \Rightarrow t_j^{\mathrm{id}} = 0$.

*Proof:* The statement follows by (36) for $m = 1$. ∎

### D. Optimality and Performance Bounds

**Proposition 3.** *Let $\boldsymbol{t}$ be an ordered target distribution with finite or countably infinite support and let $M$ be a positive integer. Among all $M$-type distributions $\boldsymbol{p}$, $\boldsymbol{p} = \boldsymbol{t}^{\mathrm{id}}$ minimizes $\mathbb{D}(\boldsymbol{p}\|\boldsymbol{t})$.*

*Proof:* See Section VI-F. ∎

The increment in the $m$-th iteration of Algorithm 2 does not depend on $M$. This means that the algorithm not only calculates the optimal $M$-type quantization, but actually *all* optimal $m$-type quantizations for $m = 1, 2, \ldots, M$. We state this property as a corollary of Proposition 3.

**Corollary 2.** *Let $\boldsymbol{c}$ be the pre-allocation calculated by Algorithm 2 in the $m$-th iteration and define*

$$\boldsymbol{t}_m^{\mathrm{id}} := \left(\frac{c_1}{m}, \ldots, \frac{c_n}{m}\right).$$

*Among all $m$-type distributions $\boldsymbol{p}$, $\boldsymbol{p} = \boldsymbol{t}_m^{\mathrm{id}}$ minimizes $\mathbb{D}(\boldsymbol{p}\|\boldsymbol{t})$.*

We next bound the informational divergence in terms of $M$. We start with the case when the support size of the target distribution is finite $(n < \infty)$. We have

$$
\begin{aligned}
\mathbb{D}(\boldsymbol{t}^{\mathrm{id}}\|\boldsymbol{t}) &\overset{(a)}{\leq} \mathbb{D}(\boldsymbol{t}^{\mathrm{vd}}\|\boldsymbol{t}) \\
&\overset{(b)}{\leq} \sum_{i:\, t_i^{\mathrm{vd}}>0} t_i^{\mathrm{vd}} \left(\frac{t_i^{\mathrm{vd}}}{t_i} - 1\right) \\
&\overset{(c)}{\leq} \sum_{i:\, t_i^{\mathrm{vd}}>0} t_i^{\mathrm{vd}} \left(\frac{t_i + \frac{1}{M}}{t_i} - 1\right) \\
&\leq \frac{1}{t_n M}
\end{aligned}
\tag{42}
$$

where (a) follows by the optimality of $\boldsymbol{t}^{\mathrm{id}}$, (b) by $\log(x) \leq x - 1$, and (c) by (16). For $n = \infty$, we have $t_i \overset{i \to \infty}{\to} 0$, so bound (42) becomes useless. The next proposition tightens (42) for $n < \infty$ and $M \geq n$ and it provides a bound for $M < n$, which is important when the support of $\boldsymbol{t}$ is infinite.

**Proposition 4.** *Let $\boldsymbol{t}$ be an ordered target distribution and let $M$ be a positive integer.*

1) *If $\boldsymbol{t}$ has finite support of cardinality $n \leq M$, then*

$$\mathbb{D}(\boldsymbol{t}^{\mathrm{id}}\|\boldsymbol{t}) < \log\left(1 + \frac{n}{2 t_n M^2}\right). \tag{43}$$

2) *If $\boldsymbol{t}$ has finite or countably infinite support of cardinality $n > M$, then*

$$\mathbb{D}(\boldsymbol{t}^{\mathrm{id}}\|\boldsymbol{t}) < \frac{1}{2} \frac{r \log r}{r - 1} \left(\frac{k}{2M} + 2 T_k\right) \tag{44}$$

*with $r = \frac{1}{1 - T_k} + \frac{e}{t_1}$.*

3) *For $n = \infty$, the support size $k$ of $\boldsymbol{t}^{\mathrm{id}}$ satisfies $k \overset{M \to \infty}{\longrightarrow} \infty$ and $k/M \overset{M \to \infty}{\longrightarrow} 0$.*

*Proof:* See the Section VI-G. ∎

We briefly discuss the intuition behind the bounds in Proposition 4. The bound (43) follows by evaluating the informational divergence of the variational distance optimal approximation $\boldsymbol{t}^{\mathrm{vd}}$. To derive bound (44), we apply Lemma 1. First, we determine the support size $k$ of $\boldsymbol{t}^{\mathrm{id}}$. Then, we use Algorithm 1 to approximate the sub-probability distribution $\boldsymbol{t}_{1:k}$. This lets us bound both the ratio $r$ and the variational distance in Lemma 1. Note that (43) and (44) are not tight for finite $M$.

## E. *Asymptotic Optimality*

For target distributions with finite support, bound (43) guarantees that the informational divergence can be made arbitrarily small by choosing $M$ large enough. This result is also valid for target distributions with infinite support by using Statement 3) of Proposition 4 in (44). We summarize these observations as a corollary to Proposition 4.

**Corollary 3.** *Let $\boldsymbol{t}$ be an ordered target distribution with finite or countably infinite support. For $M \to \infty$, the informational divergence of $\boldsymbol{t}^{\mathrm{id}}$ and $\boldsymbol{t}$ approaches zero.*

For $M \geq n$, the informational divergence approaches zero as $\mathcal{O}(1/M^2)$ by bound (43). For $M < n$, no such speed of convergence guarantee can be stated. We illustrate this by the following example.

**Example 6.** By Lemma 7, $\boldsymbol{t}^{\mathrm{id}}$ assigns mass only to at most the first (largest) $M$ indices. As in Example 4, we consider the Yule-Simon distribution. By Pinsker's inequality (7) and Example 4, the convergence of Algorithm 2 is at best $\mathcal{O}(1/M^{2\rho})$.

## V. COMPARISON OF INFORMATIONAL DIVERGENCE AND VARIATIONAL DISTANCE

### A. *Elementwise Properties*

The variational distance optimal approximation $\boldsymbol{t}^{\mathrm{vd}}$ guarantees a bounded per-entry approximation error $|t_i - t_i^{\mathrm{vd}}|$ by (16). Correspondingly, the informational divergence optimal approximation $\boldsymbol{t}^{\mathrm{id}}$ guarantees a bounded per-entry ratio $t_i^{\mathrm{id}}/t_i$ by (39). The approximations $\boldsymbol{t}^{\mathrm{vd}}$ and $\boldsymbol{t}^{\mathrm{id}}$ can violate the per-entry bounds of the other. We illustrate this by the following two examples.

**Example 7.** Let $t_1 = 1/M$ and $t_2 = \cdots = t_n = \frac{M-1}{(n-1)M}$, for $n > M$. Hence $\boldsymbol{t}^{\mathrm{vd}} = (\frac{1}{M}, \ldots, \frac{1}{M})$, and

$$\frac{t_2^{\mathrm{vd}}}{t_2} = \frac{(n-1)M}{M(M-1)} = \frac{n-1}{M-1} \tag{45}$$

can be arbitrarily large. The approximation $\boldsymbol{t}^{\mathrm{id}}$ guarantees that, by (39), we have

$$\frac{t_2^{\mathrm{id}}}{t_2} \leq \frac{e}{t_1} = eM \tag{46}$$

independent of $n$.

**Example 8.** Let $\boldsymbol{t} = (0.97, 0.01, 0.01, 0.01)$ and $M = 256$. It follows that $L = 2$ and we obtain $\boldsymbol{t}^{\mathrm{vd}} = (248, 3, 3, 2)/256$ from Algorithm 1. Algorithm 2, however, yields $\boldsymbol{t}^{\mathrm{id}} = (247, 3, 3, 3)/256$, where

$$t_1 - t_1^{\mathrm{id}} = \frac{1.32}{M} \tag{47}$$

violates (16).

Let $\boldsymbol{t} = (0.4, \varepsilon, \varepsilon, \ldots, \varepsilon)^T$ and $M = 2$. It follows that $L = 2$ and we obtain $\boldsymbol{t}^{\mathrm{vd}} = (1/2, 1/2, \ldots, 0, 0)^T$ from Algorithm 1. However, if $n$ is sufficiently large such that $\varepsilon < 0.1$, it can be shown that Algorithm 2 yields $\boldsymbol{t}^{\mathrm{id}} = (1, 0, \ldots, 0, 0)^T$, where

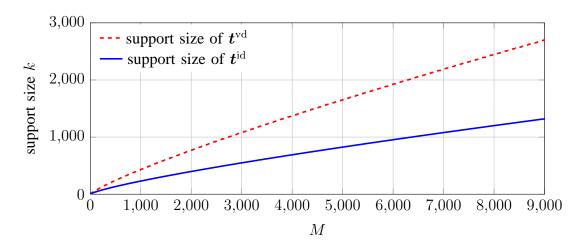$$t_1^{\mathrm{id}} - t_1 = \frac{1.2}{M} \tag{48}$$

violates (16).

Fig. 1. Support sizes of $\boldsymbol{t}^{\mathrm{vd}}$ and $\boldsymbol{t}^{\mathrm{id}}$ for the Yule-Simon distribution with $\rho = 0.2$.

## B. Support

Suppose the target distribution $\boldsymbol{t}$ has infinite support. By Statement 3) in Proposition 2 and Statement 3) in Proposition 4, the supports of the approximations $\boldsymbol{t}^{\mathrm{vd}}$ and $\boldsymbol{t}^{\mathrm{id}}$ both increase without bound and sublinearly with $M$. However, the following example shows that the support of $\boldsymbol{t}^{\mathrm{vd}}$ can grow much faster than the support of $\boldsymbol{t}^{\mathrm{id}}$. The reason is that assigning probability masses to indices with small target probabilities has a much higher cost in terms of informational divergence than in terms of variational distance. We illustrate this phenomenon by the following example.

**Example 9.** Consider the Yule-Simon distribution (see Example 4) with $\rho = 0.2$ and let $M$ take values from 1 to 10000 in steps of 10. The resulting support sizes of $\boldsymbol{t}^{\mathrm{vd}}$ and $\boldsymbol{t}^{\mathrm{id}}$ are displayed in Fig. 1. The support size of $\boldsymbol{t}^{\mathrm{vd}}$ is around twice the support size of $\boldsymbol{t}^{\mathrm{id}}$. The considered Yule-Simon distribution has a heavy tail with $T_{10000} \approx 0.15$. In other words, the first 10000 entries of $\boldsymbol{t}$ contain only 85% of the total probability mass.

The next example shows that the support of $\boldsymbol{t}^{\mathrm{vd}}$ is not always larger than the support of $\boldsymbol{t}^{\mathrm{id}}$.

**Example 10.** In Example 5 we showed that for $\boldsymbol{t} = (4/5, 1/5)$ and $M = 2$ both $\hat{\boldsymbol{t}}_1 = (1, 0)$ and $\hat{\boldsymbol{t}}_2 = (1/2, 1/2)$ are optimal in terms of the informational divergence. As it can be easily shown, $\hat{\boldsymbol{t}}_1$ is the unique approximation that is optimal in terms of the variational distance. We now modify the target distribution to $\boldsymbol{t} = (4/5 - \epsilon, 1/5 + \epsilon)$ with $0 < \epsilon < 1/20$. The vector $\hat{\boldsymbol{t}}_1$ remains the unique variational distance optimal approximation and $\hat{\boldsymbol{t}}_2$ is now the unique informational divergence optimal approximation. The support of $\hat{\boldsymbol{t}}_2$ is strictly larger than the support of $\hat{\boldsymbol{t}}_1$.

## C. Asymptotic Optimality

Corollaries 1 and 3 state that $\boldsymbol{t}^{\mathrm{vd}}$ and $\boldsymbol{t}^{\mathrm{id}}$ are asymptotically optimal w.r.t. variational distance and informational divergence, respectively. By Pinsker's inequality (7), $\boldsymbol{t}^{\mathrm{id}}$ is also asymptotically optimal w.r.t. the variational distance. In contrast, the variational distance optimal approximation $\boldsymbol{t}^{\mathrm{vd}}$ is in general not asymptotically optimal w.r.t. the informational divergence. This is illustrated by the following example.

**Example 11.** Consider the distribution $\boldsymbol{t}$ that is constructed from the geometric distribution $\tilde{t}_i = 2^{-i}$ as follows: First, $t_1 = \tilde{t}_1$. Then, the next probability mass $\tilde{t}_2$ is split into so many pieces that for $M = 2$ the informational divergence equals $\log 2$. For $M = 2$, Algorithm 1 yields $\boldsymbol{t}_2^{\mathrm{vd}} = (\frac{1}{2}, \frac{1}{2})$, where the first entry is approximated perfectly. The informational divergence of $\boldsymbol{t}_2^{\mathrm{vd}}$ and $\boldsymbol{t}$ evaluates to

$$\mathbb{D}(\boldsymbol{t}_2^{\mathrm{vd}} \| \boldsymbol{t}) = \frac{1}{2} \log \frac{1}{2 t_2} \stackrel{!}{=} \log 2 \tag{49}$$

from which $t_2 = 1/8$ follows. Thus, $t_2 = t_3 = 1/8$, which sums to $1/4$. Repeating the procedure for $M = 8$, the first three indices are approximated without error, and the two remaining masses are placed on the following indices, such that $\boldsymbol{t}_8^{\mathrm{vd}} = (\frac{1}{2}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8})$. To ensure that the informational divergence remains equal to $\log 2$, one again must split the next probability mass $\tilde{t}_3$ into sufficiently many pieces. It follows that $t_4 = \cdots = t_{19} = 1/128$, which sum to $1/8$. Repeating this procedure yields $\boldsymbol{t}$ satisfying

$$t_j = 2^{1-2^i},$$

$$\text{if } \sum_{k=0}^{i-1} 2^{2^k - k - 1} \leq j \leq \sum_{k=0}^{i} 2^{2^k - k - 1} - 1,$$

$$i \in \mathbb{N}. \quad (50)$$

For this, the subsequence $\{M_i\}_{i \in \mathbb{N}} = \{2^{2^i - 1}\}$ yields an informational divergence equal to $\log 2$, while the variational distance is bounded by $2/2^i$, i.e., twice the remaining mass of the geometric distribution. Hence, by Corollary 1, $\|\boldsymbol{t}^{\mathrm{vd}} - \boldsymbol{t}\|_1 \overset{M \to \infty}{\to} 0$, while $\limsup_{M \to \infty} \mathbb{D}(\boldsymbol{t}^{\mathrm{vd}} \| \boldsymbol{t}) = \log 2$.

## VI. PROOFS

### A. Proof of Lemma 3

We prove that every optimal $\boldsymbol{p}$ satisfies (16) by contradiction: Suppose that $p_i \leq t_i - \frac{1}{M}$ for some $i$. Since both $t_i$ and $p_i$ must sum to one, there must be a $j$ for which $p_j > t_j$. Define $\boldsymbol{p}^\circ$ by $p_i^\circ = p_i + \frac{1}{M}$, $p_j^\circ = p_j - \frac{1}{M}$, and $p_\ell^\circ = p_\ell$ for all $\ell \neq i, j$. We calculate

$$\|\boldsymbol{p} - \boldsymbol{t}\|_1 - \|\boldsymbol{p}^\circ - \boldsymbol{t}\|_1 = t_i - p_i - t_i + p_i^\circ + |p_j - t_j| - |p_j^\circ - t_j|$$

$$= \frac{1}{M} + |p_j - t_j| - |p_j - \frac{1}{M} - t_j| \quad (51)$$

$$= \frac{1}{M} + |p_j - t_j| - \left| |p_j - t_j| - \frac{1}{M} \right| \quad (52)$$

$$> 0. \quad (53)$$

where (53) follows because $p_j > t_j$. We conclude that an optimal algorithm cannot lead to $p_i \leq t_i - \frac{1}{M}$. That $p_i \geq t_i + \frac{1}{M}$ is sub-optimal follows along the same lines. $\square$

### B. Proof of Lemma 4

We claim that the two bounds in (21) relate as

$$T_k \leq \frac{k}{2M} + \frac{M T_k^2}{2k}. \quad (54)$$

This can be seen from

$$\left( \frac{k}{2M} + \frac{M T_k^2}{2k} \right) - T_k = \frac{M}{2k} \left( \frac{k^2}{M^2} - 2\frac{k}{M} T_k + T_k^2 \right)$$

$$= \frac{M}{2k} \left( \frac{k}{M} - T_k \right)^2 \geq 0. \quad (55)$$

The general bound in (22) follows by loosening the right-hand side (left-hand side) of (54) if $T_k \leq k/M$ (if $T_k \geq k/M$).

We next consider the two cases $T_k \geq k/M$ and $T_k \leq k/M$ separately.

**Case $T_k \geq k/M$:** We show that $\|\boldsymbol{t}^{\mathrm{vd}} - \boldsymbol{t}\|_1 = T_k$ and the general bound follows by (54). We have

$$\frac{k}{M} \leq T_k = 1 - \sum_{i=1}^{k} t_i = \sum_{i=1}^{k} (t_i^{\mathrm{vd}} - t_i). \quad (56)$$

In Algorithm 1, the rest mass $L/M$ after the initialization step cannot be smaller than $T_k$. Thus

$$\frac{L}{M} \geq T_k \geq \frac{k}{M} \tag{57}$$

which implies $L \geq k$. Thus, in the finalization step of Algorithm 1, each of the entries $j = 1, \ldots, k$ will get assigned at least one more mass $1/M$, so

$$\text{for each } j = 1, \ldots, k \colon (t_i^{\text{vd}} - t_i) \geq 0. \tag{58}$$

Altogether, we have

$$\|\boldsymbol{t}^{\text{vd}} - \boldsymbol{t}\|_1 = \sum_{i=1}^{k} |t_i^{\text{vd}} - t_i| \overset{(a)}{=} \sum_{i=1}^{k} (t_i^{\text{vd}} - t_i) \overset{(b)}{=} T_k \tag{59}$$

where (a) follows by (58) and where (b) follows by (56).

**Case** $T_k \leq k/M$**:** If $t_i^{\text{vd}} - t_i \geq 0$ for all $i = 1, \ldots, k$, then $\|\boldsymbol{t}^{\text{vd}} - \boldsymbol{t}\|_1 = T_k$ by (59) and (54) implies that the general bound claimed by the lemma holds. It remains to show that the general bound also holds when

$$t_j^{\text{vd}} - t_j < 0 \text{ for some } j \tag{60}$$

which implies

$$t_i^{\text{vd}} - t_i < \frac{1}{M}, \quad i = 1, \ldots, k. \tag{61}$$

In particular, (60) implies $L < k$ for the rest mass after the initialization step in Algorithm 1, which implies further that in the finalization step, each entry $i = 1, \ldots, k$ gets assigned at most one additional mass $1/M$. The error mass after the initialization step is

$$\sum_{i=1}^{k} e_i = \sum_{i=1}^{k} t_i - \sum_{i=1}^{k} \frac{\lfloor Mt_i \rfloor}{M}$$
$$= \frac{L}{M} - T_k. \tag{62}$$

Now reorder the $k$ errors such that $\tilde{e}_i \geq \tilde{e}_{i+1}$. We bound the mean error from below and above by

$$\frac{1}{L} \sum_{i=1}^{L} \tilde{e}_i \geq \frac{L}{Mk} - \frac{T_k}{k} \geq \frac{1}{k-L} \sum_{i=L+1}^{k} \tilde{e}_i. \tag{63}$$

Equality holds if $\tilde{e}_i = \frac{L}{Mk} - \frac{T_k}{k}$ for all $i = 1, \ldots, k$. After the update step in Algorithm 1, the $L$ largest errors $\tilde{e}_i$ are replaced by the final errors $1/M - \tilde{e}_i$. The other errors remain unchanged. We bound

$$\sum_{i=1}^{k} |t_i^{\text{vd}} - t_i| = \sum_{i=1}^{L} \left( \frac{1}{M} - \tilde{e}_i \right) + \sum_{i=L+1}^{k} \tilde{e}_i$$
$$\overset{(a)}{\leq} \frac{L}{M} + \left( \frac{L}{Mk} - \frac{T_k}{k} \right) (k - 2L) \tag{64}$$

where (a) follows by (63). The maximum is achieved for $L = (k + MT_k)/2$, which yields

$$\|\boldsymbol{t}^{\text{vd}} - \boldsymbol{t}\|_1 \leq \frac{k}{2M} + \frac{MT_k^2}{2k}. \tag{65}$$

$\square$

## C. Proof of Proposition 2

*1) :* The proof follows from Lemma 4 for $k = n$ and $T_k = T_n \equiv 0$.

*2) :* Let $k$ be the support size of $\boldsymbol{t}^{\mathrm{vd}}$, and let $\boldsymbol{t}_{1:k}$ be the sub-probability distribution obtained by taking the first $k$ indices of $\boldsymbol{t}$. Then, we have

$$\|\boldsymbol{t}^{\mathrm{vd}} - \boldsymbol{t}\|_1 = \|\boldsymbol{t}^{\mathrm{vd}} - \boldsymbol{t}_{1:k}\|_1 + T_k. \tag{66}$$

If $k$ is the support size, then by Lemma 2 the first $k$ indices get $M$ masses. Since the algorithm satisfies (16), we have

$$T_k = 1 - \sum_{i=1}^{k} t_i = \sum_{i=1}^{k} (t_i^{\mathrm{vd}} - t_i) \leq \frac{k}{M}. \tag{67}$$

Thus we can bound $\|\boldsymbol{t}^{\mathrm{vd}} - \boldsymbol{t}_{1:k}\|_1$ by Lemma 4 and get

$$
\begin{aligned}
\|\boldsymbol{t}^{\mathrm{vd}} - \boldsymbol{t}\|_1 &\leq \frac{k}{2M} + \frac{M T_k^2}{2k} + T_k \\
&= \frac{k}{2M} \left( 1 + \frac{2M T_k}{k} + \frac{M^2 T_k^2}{k^2} \right) \\
&= \frac{k}{2M} \left( 1 + \frac{M T_k}{k} \right)^2.
\end{aligned} \tag{68}
$$

*3) :* The support size $k$ of $\boldsymbol{t}^{\mathrm{vd}}$ grows without bound with $M$ because for every $l$ there exists an $M$ such that $t_l > 1/M$, hence this index gets probability mass already in the initialization step of Algorithm 1.

We show that the support size $k \equiv k(M)$ grows sublinearly with $M$ by contradiction. Suppose there exists a $0 < c \leq 1$ such that

$$\limsup_{M \to \infty} \frac{k(M)}{M} = c. \tag{69}$$

Thus, for each $\epsilon > 0$, there exists a sequence $\{M_i\}_{i \in \mathbb{N}}$, $M_1 < M_2 < M_3 < \cdots$, such that

$$(c - \epsilon) M_i < k(M_i) < (c + \epsilon) M_i, \quad i \in \mathbb{N}. \tag{70}$$

Now choose $i < j \in \mathbb{N}$. Applying the algorithm for $M_i$ and $M_j$ increases the support size from $k(M_i)$ to $k(M_j)$. In total, the algorithm has $M_j$ masses to distribute, some of which are distributed to the first $k(M_i)$ indices. In particular, in the first step the algorithm assigns

$$\sum_{l=1}^{k(M_i)} \lfloor M_j t_l \rfloor \tag{71}$$

masses to the first $k(M_i)$ indices. The difference in support sizes is thus bounded from above by

$$
\begin{aligned}
k(M_j) - k(M_i) &\leq M_j - \sum_{l=1}^{k(M_i)} \lfloor M_j t_l \rfloor \\
&< M_j - \sum_{l=1}^{\lfloor (c-\epsilon)M_i \rfloor} \lfloor M_j t_l \rfloor \\
&= M_j \left( 1 - \sum_{l=1}^{\lfloor (c-\epsilon)M_i \rfloor} \frac{\lfloor M_j t_l \rfloor}{M_j} \right) \\
&= M_j \left( T_{\lfloor (c-\epsilon)M_i \rfloor} + \sum_{l=1}^{\lfloor (c-\epsilon)M_i \rfloor} \left( t_l - \frac{\lfloor M_j t_l \rfloor}{M_j} \right) \right) \\
&< M_j \left( T_{\lfloor (c-\epsilon)M_i \rfloor} + \frac{(c-\epsilon)M_i}{M_j} \right).
\end{aligned}
\tag{72}
$$

Now choose $i$ large enough such that $T_{\lfloor (c-\epsilon)M_i \rfloor} < \epsilon$ and choose $j$ large enough such that $M_i/M_j < 1/4$. We have

$$
\frac{k(M_j) - k(M_i)}{M_j} < \epsilon + \frac{c-\epsilon}{4}.
\tag{73}
$$

A lower bound on the support size difference is obtained from (70):

$$
\frac{k(M_j) - k(M_i)}{M_j} > (c-\epsilon) - (c+\epsilon)\frac{M_i}{M_j} > \frac{3c}{4} - \frac{5\epsilon}{4}.
\tag{74}
$$

Combining (73) and (74) yields an upper bound on $c$:

$$
\frac{3c}{4} - \frac{5\epsilon}{4} < \frac{c-\epsilon}{4} + \epsilon.
\tag{75}
$$

After rearranging we have $c < 4\epsilon$ for any $\epsilon > 0$, and thus

$$
\limsup_{M \to \infty} \frac{k(M)}{M} = 0.
\tag{76}
$$

$\square$

*D. Proof of* (27)

We make use of the following lower bound on the beta function [21, eq. (2)]

$$
B(x,y) \geq \frac{x^{x-1} y^{y-1}}{(x+y)^{x+y-1}}
\tag{77}
$$

which in our case gives

$$
\begin{aligned}
M \cdot B(M, \rho+1) &\geq \frac{M^M (\rho+1)^\rho}{(M+\rho+1)^{M+\rho}} \\
&= \frac{M^M}{(M+\rho+1)^M} \frac{(\rho+1)^\rho}{(M+\rho+1)^\rho} \\
&= \frac{(\rho+1)^\rho}{(1+\frac{\rho+1}{M})^M} \frac{1}{(M+\rho+1)^\rho} \\
&\geq \frac{(\rho+1)^\rho}{e^{\rho+1}} \frac{1}{(M+\rho+1)^\rho}
\end{aligned}
\tag{78}
$$

where (78) follows because $(1 + \frac{\rho+1}{M})^M$ approaches $e^{\rho+1}$ from below. This shows the existence of the constant $K(\rho)$ in (27).

## E. Proof of Lemma 6

The case $M = 1$ (hence $k = 1$) is trivial; we focus on $M \geq 2$. Suppose that $\boldsymbol{p}$ is an $M$-type distribution (not necessarily optimal) and that $\boldsymbol{p}^\circ$ is such that $p_i^\circ = p_i + \frac{1}{M} \leq 1$, $p_j^\circ = p_j - \frac{1}{M} \geq 0$ and $p_\ell = p_\ell^\circ$ for all $\ell \neq i, j$. We now show that $\mathbb{D}(\boldsymbol{p}\|\boldsymbol{t}) > \mathbb{D}(\boldsymbol{p}^\circ\|\boldsymbol{t})$ holds if $\boldsymbol{p}$ violates the statement of Lemma 6, i.e., that $\boldsymbol{p}$ is not optimal is not optimal in this case. To this end, notice that

$$
\begin{aligned}
\mathbb{D}(\boldsymbol{p}\|\boldsymbol{t}) - \mathbb{D}(\boldsymbol{p}^\circ\|\boldsymbol{t}) &= p_i \log \frac{p_i}{t_i} + p_j \log \frac{p_j}{t_j} - \left(p_i + \frac{1}{M}\right) \log \frac{p_i + \frac{1}{M}}{t_i} - \left(p_j - \frac{1}{M}\right) \log \frac{p_j - \frac{1}{M}}{t_j} \\
&\overset{(a)}{=} p_j \log \frac{p_j}{t_j} - \left(p_j - \frac{1}{M}\right) \log \frac{p_j - \frac{1}{M}}{t_j} - \frac{1}{M}\left(\Delta_i(Mp_i + 1) - \log M\right) \\
&\overset{(b)}{>} \frac{1}{M} \log \frac{p_j}{t_j} + \underbrace{\left(p_j - \frac{1}{M}\right) \log \frac{p_j}{p_j - \frac{1}{M}}}_{>0} - \frac{1}{M}\left(\Delta_i(M) - \log M\right) \\
&> \frac{1}{M} \log \frac{p_j}{t_j} + \underbrace{\frac{M-1}{M} \log \frac{M-1}{M}}_{\geq -\frac{1}{M}} + \frac{1}{M} \log t_i \\
&\geq \frac{1}{M} \log \frac{p_j}{t_j} - \frac{1}{M} \log \frac{e}{t_i}
\end{aligned}
$$

where $(a)$ is due to (32) and $(b)$ follows by (35). Hence, if

$$
\frac{p_j}{t_j} \geq \frac{e}{t_i} \tag{79}
$$

for any pair of indices $i$ and $j$, then above difference of informational divergences is positive as well. Thus, an optimal $\boldsymbol{p}$ may not fulfill (79) for any such pair of indices. The best bound is obtained for $i = 1$, hence Lemma 6 follows. The result for index $k$ results from $p_k \geq 1/M$. $\qquad\square$

## F. Proof of Proposition 3

To prove optimality, we need the following lemma.

**Lemma 8.** *Let $\boldsymbol{c}^*$ be an optimal allocation. Let $\boldsymbol{c}$ be a pre-allocation with $\sum_i c_i < M$ and $c_i \leq c_i^*$ for $i = 1, \ldots, n$. Define*

$$
j = \underset{i}{\operatorname{argmin}} \, \Delta_i(c_i + 1). \tag{80}
$$

*Then there exists an optimal allocation $\tilde{\boldsymbol{c}}$ with*

$$
c_j + 1 \leq \tilde{c}_j \tag{81}
$$
$$
c_i \leq \tilde{c}_i, \quad i = 1, \ldots, n. \tag{82}
$$

*Proof:* Suppose we have

$$
c_j + 1 > c_j^*. \tag{83}
$$

Since $c_j \leq c_j^*$ by assumption, (83) implies

$$
c_j + 1 = c_j^* + 1. \tag{84}
$$

Since $\sum_i c_i < M$ and $\sum_i c_i^* = M$, there must be at least one $\ell \neq j$ with

$$
c_\ell^* \geq c_\ell + 1. \tag{85}
$$

By decreasing $c_\ell^*$ by one and increasing $c_j^*$ by one, the change of the objective function is $\Delta_j(c_j^* + 1) - \Delta_\ell(c_\ell^*)$. We bound this change as follows:

$$\Delta_j(c_j^* + 1) - \Delta_\ell(c_\ell^*) \overset{(a)}{\leq} \Delta_j(c_j^* + 1) - \Delta_\ell(c_\ell + 1) \tag{86}$$
$$\overset{(b)}{=} \Delta_j(c_j + 1) - \Delta_\ell(c_\ell + 1)$$
$$\overset{(c)}{\leq} 0 \tag{87}$$

where (a) follows by (85) and Lemma 5, (b) follows by (84), and (c) follows by the definition of $j$ in (80).

We must consider two cases. First, suppose we have strict inequality in either (86) or (87). Then the objective function is decreased, which contradicts the assumption that $\boldsymbol{c}^*$ is optimal. Thus, the supposition (83) is false and the statements of the lemma hold for $\tilde{\boldsymbol{c}} = \boldsymbol{c}^*$. Second, suppose we have equality both in (86) and (87). In this case, define the allocation

$$\tilde{c}_\ell = c_\ell^* - 1, \quad \tilde{c}_j = c_j^* + 1, \quad \tilde{c}_i = c_i^* \text{ for } i \neq j, \ell. \tag{88}$$

Equality in (86)–(87) implies optimality of $\tilde{\boldsymbol{c}}$. By (84) and (85), we can verify that $\tilde{\boldsymbol{c}}$ fulfills the statements of the lemma. This concludes the proof of Lemma 8. ∎

We are now ready to prove Proposition 3. By Lemma 8, there is an optimal allocation $\tilde{\boldsymbol{c}}$ such that in each iteration of Algorithm 2 we have

$$c_i \leq \tilde{c}_i, \qquad i = 1, \ldots, n. \tag{89}$$

After Algorithm 2 terminates, we have

$$M = \sum_i c_i \leq \sum_i \tilde{c}_i = M. \tag{90}$$

Statements (89) and (90) can be true simultaneously only if $c_i = \tilde{c}_i$ for all $i = 1, \ldots, n$. Consequently, the constructed allocation $\boldsymbol{c}$ is optimal. □

### G. Proof of Proposition 4

*1) Case $M \geq n$:* By Proposition 3, $\boldsymbol{t}^{\text{id}}$ is optimal w.r.t. the informational divergence and

$$\mathbb{D}(\boldsymbol{t}^{\text{id}}\|\boldsymbol{t}) \leq \mathbb{D}(\boldsymbol{t}^{\text{vd}}\|\boldsymbol{t}). \tag{91}$$

Moreover,

$$\mathbb{D}(\boldsymbol{t}^{\text{vd}}\|\boldsymbol{t}) = \sum_{i=1}^n t_i^{\text{vd}} \log \frac{t_i^{\text{vd}}}{t_i}$$
$$\overset{(a)}{\leq} \log \left( \sum_{i=1}^n \frac{(t_i^{\text{vd}})^2}{t_i} \right)$$
$$= \log \left( 1 + \sum_{i=1}^n \frac{(t_i^{\text{vd}} - t_i)^2}{t_i} \right) \tag{92}$$

where $(a)$ is Jensen's inequality (see also the proof of [17, Thm. 3]) and where the sum inside the logarithm is Pearson's $\chi^2$-distance $\chi^2(\boldsymbol{t}^{\text{vd}}\|\boldsymbol{t})$. Note that (92) equals $\mathbb{D}_2(\boldsymbol{t}^{\text{vd}}\|\boldsymbol{t})$, the Rényi divergence of second order. The inequality in $(a)$ is then a direct consequence of the fact that Rényi divergence is non-decreasing in the order [22, Thm. 3].

We now bound (92) by

$$\sum_{i=1}^{n} \frac{(t_i^{\mathrm{vd}} - t_i)^2}{t_i} \leq \frac{1}{t_n} \sum_{i=1}^{n} (t_i^{\mathrm{vd}} - t_i)^2$$

$$= \frac{1}{t_n} \sum_{i=1}^{n} |t_i^{\mathrm{vd}} - t_i| \underbrace{|t_i^{\mathrm{vd}} - t_i|}_{< \frac{1}{M} \text{ by (16)}}$$

$$< \frac{1}{t_n M} \|\boldsymbol{t}^{\mathrm{vd}} - \boldsymbol{t}\|_1$$

$$\overset{(a)}{\leq} \frac{n}{2 t_n M^2} \tag{93}$$

where (a) follows by Statement 1) in Proposition 2.

*2) Case $M < n$:* Let $k$ be the support size of $\boldsymbol{t}^{\mathrm{id}}$. Define the auxiliary distribution $\tilde{\boldsymbol{t}} := \boldsymbol{t}_{1:k}/(1 - T_k)$. Because of the normalization by $1 - T_k$, the entries of $\tilde{\boldsymbol{t}}$ sum to one and $\tilde{\boldsymbol{t}}$ is a distribution. Denote by $\tilde{\boldsymbol{t}}^{\mathrm{vd}}$ the approximation that results from applying Algorithm 1 to $\tilde{\boldsymbol{t}}$. We have

$$\mathbb{D}(\boldsymbol{t}^{\mathrm{id}} \| \boldsymbol{t}) \leq \mathbb{D}(\tilde{\boldsymbol{t}}^{\mathrm{vd}} \| \boldsymbol{t}) \overset{(a)}{\leq} \frac{1}{2} \frac{r \log r}{r - 1} \|\tilde{\boldsymbol{t}}^{\mathrm{vd}} - \boldsymbol{t}\|_1$$

$$\text{with } r = \max_{i \leq k} \frac{\tilde{t}_i^{\mathrm{vd}}}{t_i} \tag{94}$$

where (a) follows by Lemma 1. It remains to bound the ratio $r$ and the variational distance $\|\tilde{\boldsymbol{t}}^{\mathrm{vd}} - \boldsymbol{t}\|_1$.

*Bounding $r$:* By (16), we have

$$\tilde{t}_i^{\mathrm{vd}} < \tilde{t}_i + \frac{1}{M} = \frac{t_i}{1 - T_k} + \frac{1}{M}. \tag{95}$$

Thus, for each $i \leq k$, we have

$$\frac{\tilde{t}_i^{\mathrm{vd}}}{t_i} < \frac{1}{1 - T_k} + \frac{1}{t_i M} \leq \frac{1}{1 - T_k} + \frac{1}{t_k M} \tag{96}$$

which implies

$$r < \frac{1}{1 - T_k} + \frac{1}{t_k M}$$

$$\overset{(a)}{\leq} \frac{1}{1 - T_k} + \frac{e}{t_1} \tag{97}$$

where (a) follows by (40) in Lemma 6.

*Bounding $\|\tilde{\boldsymbol{t}}^{\mathrm{vd}} - \boldsymbol{t}\|_1$:* We bound

$$\|\tilde{\boldsymbol{t}}^{\mathrm{vd}} - \boldsymbol{t}\|_1 = \|\tilde{\boldsymbol{t}}^{\mathrm{vd}} - \boldsymbol{t}_{1:k}\|_1 + T_k$$

$$= \|\tilde{\boldsymbol{t}}^{\mathrm{vd}} - \tilde{\boldsymbol{t}}(1 - T_k)\|_1 + T_k$$

$$\overset{(a)}{\leq} \|\tilde{\boldsymbol{t}}^{\mathrm{vd}} - \tilde{\boldsymbol{t}}\|_1 + \|\tilde{\boldsymbol{t}} T_k\|_1 + T_k$$

$$= \|\tilde{\boldsymbol{t}}^{\mathrm{vd}} - \tilde{\boldsymbol{t}}\|_1 + 2 T_k$$

$$\leq \frac{k}{2M} + 2 T_k \tag{98}$$

where (a) follows by the triangle inequality. Using (97) and (98) in (94) completes the proof.

*3) :* The support $k$ grows without bound because the increment functions $\Delta_i$ grow without bound by (34), i.e., for every positive integer $\ell$ there exists an $M$ large enough such that, for all $i = 1, \ldots, \ell - 1$,

$$\log \frac{1}{t_\ell} < \Delta_i(c_i + 1) \tag{99}$$

where the sum over all $c_i$ is less than $M$. In other words, after assigning a specific number of masses to indices 1 to $\ell - 1$, assigning a mass to index $\ell$ must have lower cost than assigning additional masses to the first $\ell - 1$ indices.

The result $k(M)/M \overset{M \to \infty}{\to} 0$ can be seen as follows. Increasing $M$ by one increases the support size $k$ at most by one. This is a consequence of the update rule in Algorithm 2. Thus, the sequence $k \equiv k(M)$ contains each integer $1, 2, 3, \ldots$ at least once and we can define a sequence $M(k)$, $k = 1, 2, 3, \ldots$. Note that some integers may not occur in the sequence $M(k)$. By (40) in Lemma 6, we can bound the $k$-th probability $t_k$ by

$$t_k > \frac{t_1}{eM(k)} \tag{100}$$

and we have

$$1 = \sum_{k=1}^{\infty} t_k > \sum_{k=1}^{\infty} \frac{t_1}{eM(k)}. \tag{101}$$

If $M(k)$ grows only linearly with $k$, then the sum on the right-hand side diverges, which contradicts that the probabilities need to sum to one. Thus, $M(k)$ grows super-linearly with $k$ and equivalently, $k(M)$ grows sub-linearly with $M$. $\qquad\square$

## VII. ACKNOWLEDGMENT

## REFERENCES

[1] G. Böcherer, "Optimal non-uniform mapping for probabilistic shaping," in *Proc. Int. ITG Conf. Sys., Commun, Coding (SCC)*, Munich, Jan. 2013, pp. 1–6.

[2] G. Böcherer and R. Mathar, "Matching dyadic distributions to channels," in *Proc. Data Compression Conf. (DCC)*, Mar. 2011, pp. 23–32.

[3] G. Böcherer, "Capacity-achieving probabilistic shaping for noisy and noiseless channels," Ph.D. dissertation, RWTH Aachen University, 2012. [Online]. Available: http://www.georg-boecherer.de/capacityAchievingShaping.pdf

[4] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*. Cambridge University Press, 2011.

[5] R. G. Gallager, *Information Theory and Reliable Communication*. John Wiley & Sons, Inc., 1968.

[6] D. Raphaeli and A. Gurevitz, "Constellation shaping for pragmatic turbo-coded modulation with high spectral efficiency," *IEEE Trans. Commun.*, vol. 52, no. 3, pp. 341–345, 2004.

[7] M. Yankov, S. Forchhammer, K. J. Larsen, and L. P. Christensen, "Rate-adaptive constellation shaping for near-capacity achieving turbo coded BICM," in *Proc. IEEE Int. Conf. Commun. (ICC)*, 2014, pp. 2112–2117.

[8] F. Schreckenbach and P. Henkel, "Signal shaping using non-unique symbol mappings," in *Proc. Allerton Conf. Commun., Contr., Comput.*, Sep. 2005.

[9] Y. A. Reznik, "An algorithm for quantization of discrete probability distributions," in *Proc. Data Compression Conf. (DCC)*, Snowbird, UT, Mar. 2011, pp. 333–342.

[10] Y. A. Reznik, V. Chandrasekhar, G. Takacs, D. M. Chen, S. S. Tsai, and B. Girod, "Fast quantization and matching of histogram-based image features," in *Proc. SPIE 7798, App. Digital Image Proc. XXXIII*, San Diego, Aug. 2010.

[11] T. S. Han and S. Verdu, "Approximation theory of output statistics," *IEEE Trans. Inf. Theory*, vol. 39, no. 3, pp. 752–772, 1993.

[12] G. Böcherer and R. A. Amjad, "Fixed-to-variable length resolution coding for target distributions," in *Proc. IEEE Inf. Theory Workshop (ITW)*, Sep. 2013, pp. 1–5.

[13] Y. Steinberg and S. Verdu, "Simulation of random processes and rate-distortion theory," *IEEE Trans. Inf. Theory*, vol. 42, no. 1, pp. 63–86, Jan. 1996.

[14] J. Hou and G. Kramer, "Informational divergence approximations to product distributions," in *Proc. Canadian Workshop Inf. Theory (CWIT)*, Jun. 2013, pp. 76–81.

[15] A. D. Wyner, "The common information of two dependent random variables," *IEEE Trans. Inf. Theory*, vol. 21, no. 2, pp. 163–179, Mar. 1975.

[16] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed.   John Wiley & Sons, Inc., 2006.

[17] I. Sason, "On reverse Pinsker inequalities," Mar. 2015, `arXiv:1503.07118 [cs.IT]`.

[18] S. Verdú, "Total variation distance and the distribution of relative information," in *Proc. Inf. Theory and Applicat. Workshop (ITA)*, San Diego, CA, Feb. 2014, pp. 499–501.

[19] D. Berend, P. Harremoës, and A. Kontorovich, "Minimum KL-divergence on complements of $L_1$ balls," *IEEE Trans. Inf. Theory*, 2014.

[20] H. A. Simon, "On a class of skew distribution functions," *Biometrika*, vol. 42, no. 314, pp. 425–440, 1955.

[21] L. Grenié and G. Molteni, "Inequalities for the Beta function," *Mathematical Inequalities & Applications*, vol. 18, no. 4, pp. 1427–1442, 2015.

[22] T. van Erven and P. Harremoës, "Rényi divergence and Kullback-Leibler divergence," *IEEE Trans. Inf. Theory*, vol. 60, no. 7, pp. 3797–3820, Jul. 2014.