# Dimension-Free Empirical Entropy Estimation

Doron Cohen
Department of Computer Science
Ben-Gurion University of the Negev
Beer-Sheva, Israel
doronv@post.bgu.ac.il

Aryeh Kontorovich
Department of Computer Science
Ben-Gurion University of the Negev
Beer-Sheva, Israel
karyeh@cs.bgu.ac.il

Aaron Koolyk
Department of Computer Science
Hebrew University, Jerusalem, Israel
aaron.koolyk@mail.huji.ac.il

Geoffrey Wolfer
RIKEN, Center for AI Project
Tokyo, Japan *
geo-wolfer@m2.tuat.ac.jp

December 27, 2022

## Abstract

We seek an entropy estimator for discrete distributions with fully empirical accuracy bounds. As stated, this goal is infeasible without some prior assumptions on the distribution. We discover that a certain information moment assumption renders the problem feasible. We argue that the moment assumption is natural and, in some sense, *minimalistic* — weaker than finite support or tail decay conditions. Under the moment assumption, we provide the first finite-sample entropy estimates for infinite alphabets, nearly recovering the known minimax rates. Moreover, we demonstrate that our empirical bounds are significantly sharper than the state-of-the-art bounds, for various natural distributions and non-trivial sample regimes. Along the way, we give a dimension-free analogue of the Cover-Thomas result on entropy continuity (with respect to total variation distance) for finite alphabets, which may be of independent interest. Additionally, we resolve all of the open problems posed by Jürgensen and Matthews, 2010.

# 1  Introduction

Estimating the entropy of a discrete distribution based on a finite iid sample is a classic problem with theoretical and practical ramifications. Considerable progress has been made in the case of a finite alphabet, and the countably infinite case has also attracted a fair amount of attention in recent years. (See Section 4 for a some background, motivation, and related work.) A less-addressed issue is one of *empirical* accuracy estimates: data-dependent bounds adaptive to the particular distribution being sampled.

Our point of departure is the simpler (to analyze) problem of estimating a discrete distribution $\boldsymbol{\mu}$ in total variation norm $\|\cdot\|_{\mathsf{TV}} = \frac{1}{2} \|\cdot\|_1$, where the most recent advance was made by Cohen et al. (2020); see therein for a literature review. If $\boldsymbol{\mu}$ is a distribution on $\mathbb{N}$ and $\hat{\boldsymbol{\mu}}_n$ is its empirical realization based on a sample of size $n$, then Theorem 2.1 of Cohen et al. states that with probability at least $1 - \delta$,

$$\|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_n\|_1 \quad \leq \quad \frac{2}{\sqrt{n}} \sum_{j \in \mathbb{N}} \sqrt{\hat{\boldsymbol{\mu}}_n(j)} + 6\sqrt{\frac{\log(2/\delta)}{2n}}. \tag{1}$$

This bound has the advantage of being valid for all distributions on $\mathbb{N}$, without any prior assumptions, and being fully empirical: it yields a risk estimate that is computable based on the observed sample,

---

not depending on any unknown quantities. (Additionally, Cohen et al. argue that (1) is near-optimal in a well-defined sense.) The question we set out to explore in this paper is: What analogues of (1) are possible for discrete entropy estimation?

When $\boldsymbol{\mu}$ has support size $d < \infty$, an answer to our question is readily provided by combining (1) with Cover and Thomas (2006, Theorem 17.3.3), which asserts that, for $\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_1 \leq 1/2$, we have

$$
|\mathrm{H}(\boldsymbol{\mu}) - \mathrm{H}(\boldsymbol{\nu})| \quad \leq \quad \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_1 \log \frac{d}{\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_1}, \tag{2}
$$

where $\mathrm{H}(\cdot)$ is the entropy functional defined in (3). Indeed, taking $\boldsymbol{\mu}$ as in (1) and $\boldsymbol{\nu}$ to be $\hat{\boldsymbol{\mu}}_n$ yields a fully empirical estimate on $|\mathrm{H}(\boldsymbol{\mu}) - \mathrm{H}(\hat{\boldsymbol{\mu}}_n)|$. For fixed $d < \infty$, no technique relying on the plug-in estimator can yield minimax rates (Wu and Yang, 2016). The plug-in is, however, asymptotically optimal for fixed $d < \infty$ (Paninski, 2003) as well as strongly universally consistent even for $d = \infty$ (Antos and Kontoyiannis, 2001a), and is among the few methods for which explicitly computable finite-sample risk bounds are known.

The thrust of this paper is to replace the restrictive finite-support assumption with considerably more general moment conditions. It is well-known that when estimating the mean of some random variable $X$, the first-moment assumption $\mathbb{E}|X| \leq M$ is not sufficient to yield any finite-sample information.[1] Strengthening the assumption to $\mathbb{E}|X|^\alpha \leq M$, for any $\alpha > 1$, immediately yields finite-sample empirical estimates on $\left|\mathbb{E}X - \frac{1}{n}\sum_{i=1}^n X_i\right|$ via the von Bahr and Esseen (1965) inequality.[2] In this sense, a bound on the $(1 + \varepsilon)$th moment is a *minimal* requirement for empirical mean estimation. However, it is not immediately obvious how to apply this insight to the entropy estimation problem: the corresponding random variable is $X = -\log \boldsymbol{\mu}(I)$, where $I \sim \boldsymbol{\mu}$, but rather than being given iid samples of $X$, we are only given draws of $I$. In Corollary 1, we provide an empirical entropy estimate under a $(1 + \varepsilon)$th moment assumption (for any $\varepsilon > 0$) on $X = -\log \boldsymbol{\mu}(I)$.

**Our contribution.** In Theorem 1, we provide a dimension-free analogue of (2), which, combined with (1), allows for empirical accuracy bounds on the plug-in entropy estimator under a minimalistic moment assumption. Moreover, for this rich class of distributions, the plug-in estimator turns out to be asymptotically optimal, as we show in Theorem 5. Our moment assumption is natural and considerably less restrictive than the finite-alphabet and tail conditions studied in previous works (see Sections 7 and .1). Moreover, as we argue in Theorem 4, without such a moment assumption, an empirical bound is not feasible. As we demonstrate in Section 6, the rates provided by our empirical bound compare favorably against the state of the art.

## 2 Definitions and notation

Our logarithms will always be base e by default. For discrete distributions, there is no loss of generality in taking the domain to be the natural numbers $\mathbb{N} = \{1, 2, 3, \ldots\}$. For $k \in \mathbb{N}$, we write $[k] := \{i \in \mathbb{N} : i \leq k\}$. The set of all probability distributions on $\mathbb{N}$ will be denoted by $\Delta_{\mathbb{N}}$. For $d \in \mathbb{N}$, we write $\Delta_d \subset \Delta_{\mathbb{N}}$ to denote those $\boldsymbol{\mu}$ whose support is contained in $[d]$.

We define the operator $(\cdot)^\downarrow$, which maps any $\boldsymbol{\mu} \in \Delta_{\mathbb{N}}$ to its non-increasing rearrangement $\boldsymbol{\mu}^\downarrow$. The set of all non-increasing distributions will be denoted by $\Delta_{\mathbb{N}}^\downarrow := \{\boldsymbol{\mu}^\downarrow : \boldsymbol{\mu} \in \Delta_{\mathbb{N}}\}$.

We write $\mathbb{R}_+ := [0, \infty)$. For any $\xi : \mathbb{N} \to \mathbb{R}_+$ and $\alpha \geq 0$, define

$$
\mathrm{H}^{(\alpha)}(\xi) := \sum_{j \in \mathbb{N}: \xi(j) > 0} \xi(j) \left|\log \xi(j)\right|^\alpha. \tag{3}
$$

For $\xi \in \mathbb{R}^{\mathbb{N}}$, denote by $|\xi| \in \mathbb{R}_+^{\mathbb{N}}$ the elementwise application of $|\cdot|$ to $\xi$. When $\xi \in \Delta_{\mathbb{N}}$ and $\alpha = 1$, (3) recovers the standard definition of entropy, which we denote by $\mathrm{H}(\xi) := \mathrm{H}^{(1)}(\xi)$. For general $\alpha > 0$,

---

[1]Even distinguishing, for $X \geq 0$, between $\mathbb{E}X = 0$ and $\mathbb{E}X = M$ based on a finite sample is impossible with any degree of confidence. Of course, $\frac{1}{n}\sum_{i=1}^n X_i \to \mathbb{E}X$ almost surely, by the strong law of large numbers.

[2]Put $Y = X - \mathbb{E}X$; then $\mathbb{E}|Y| \leq 2M$. For $1 < \alpha < 2$, a sharper version of the Bahr-Esseen inequality (Pinelis, 2015) states that $\mathbb{E}\left[\left|\sum_{i=1}^n Y_i\right|^\alpha\right] \leq 2n(2M)^\alpha$, which implies tail bounds via Markov's inequality. Better rates are available via the median-of-means estimator, see Lugosi and Mendelson (2019).

this quantity may be referred to as the $\alpha$th *moment of information*. For $h \geq 0$, define

$$\Delta_{\mathbb{N}}^{(\alpha)}[h] = \left\{ \boldsymbol{\mu} \in \Delta_{\mathbb{N}} : \mathrm{H}^{(\alpha)}(\boldsymbol{\mu}) \leq h \right\}$$

and also $\Delta_{\mathbb{N}}^{(\alpha)} := \bigcup_{h \geq 0} \Delta_{\mathbb{N}}^{(\alpha)}[h]$ and $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h] := \Delta_{\mathbb{N}}^{\downarrow} \cap \Delta_{\mathbb{N}}^{(\alpha)}[h]$.

For $n \in \mathbb{N}$ and $\boldsymbol{\mu} \in \Delta_{\mathbb{N}}$, we write $\boldsymbol{X} = (X_1, \ldots, X_n) \sim \boldsymbol{\mu}^n$ to mean that the components of the vector $\boldsymbol{X}$ are drawn iid from $\boldsymbol{\mu}$. The empirical measure $\hat{\boldsymbol{\mu}}_n \in \Delta_{\mathbb{N}}$ induced by the sample $\boldsymbol{X}$ is defined by $\hat{\boldsymbol{\mu}}_n(j) = \frac{1}{n} \sum_{i \in [n]} \mathbf{1}[X_i = j]$. For any $\xi \in \mathbb{R}^{\mathbb{N}}$ and $0 < p < \infty$, the $\ell_p$ (pseudo)norm is defined by $\|\xi\|_p^p = \sum_{j \in \mathbb{N}} |\xi(j)|^p$ and $\|\xi\|_\infty = \sup_{j \in \mathbb{N}} |\xi(j)|$.

For $\alpha, h > 0$, and $n \in \mathbb{N}$, define the $L_1$ *minimax risk* for the $\alpha$th moment by

$$\mathcal{R}_n^{(\alpha)}(h) := \inf_{\hat{H}} \sup_{\boldsymbol{\mu} \in \Delta_{\mathbb{N}}^{(\alpha)}[h]} \mathbb{E}|\hat{H}(X_1, \ldots, X_n) - \mathrm{H}(\boldsymbol{\mu})|, \tag{4}$$

where the infimum is over all mappings $\hat{H} : \mathbb{N}^n \to \mathbb{R}_+$.

# 3 Main results

Our first result is a dimension-free analogue of (2):

**Theorem 1.** *For all $\alpha > 1$, $\mathrm{H} : \Delta_{\mathbb{N}}^{(\alpha)} \to \mathbb{R}_+$ is uniformly continuous under $\ell_1$. In particular, for all $\boldsymbol{\mu}, \boldsymbol{\nu} \in \Delta_{\mathbb{N}}^{(\alpha)}$ satisfying $\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_\infty < \mathrm{e}^{-\alpha}$, we have*

$$|\mathrm{H}(\boldsymbol{\mu}) - \mathrm{H}(\boldsymbol{\nu})| \leq \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_1^{1-1/\alpha} \left( 2\mathrm{e}^\alpha \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_\infty \log^\alpha \frac{1}{\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_\infty} + \mathrm{H}^{(\alpha)}(\boldsymbol{\mu}) + \mathrm{H}^{(\alpha)}(\boldsymbol{\nu}) \right)^{1/\alpha} \tag{5}$$

$$\leq \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_1^{1-1/\alpha} \left( 2\mathrm{e} \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_\infty^{1/\alpha} \log \frac{1}{\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_\infty} + \mathrm{H}^{(\alpha)}(\boldsymbol{\mu})^{1/\alpha} + \mathrm{H}^{(\alpha)}(\boldsymbol{\nu})^{1/\alpha} \right). \tag{6}$$

*Moreover, a weaker form of (5) holds with $\alpha^\alpha$ in place of $\mathrm{e}^\alpha \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_\infty \log^\alpha \frac{1}{\|\boldsymbol{\mu}-\boldsymbol{\nu}\|_\infty}$ under the weaker condition $\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_\infty < 1/2$, from which the correspondingly weaker form of (6) follows as well.*

The requirement in Theorem 1 that $\alpha > 1$ cannot be dispensed with, as the function $\mathrm{H} : \Delta_{\mathbb{N}}^{(\alpha)}[h] \to \mathbb{R}_+$ is not continuous under $\ell_1$ for $\alpha = 1$ (see Remark following Lemma 4), and, a fortiori, is not uniformly continuous. Thus, there can be no function $F : \mathbb{R}_+^2 \to \mathbb{R}_+$ satisfying

$$|\mathrm{H}(\boldsymbol{\mu}) - \mathrm{H}(\boldsymbol{\nu})| \leq F(\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_1, h), \qquad h > 0, \boldsymbol{\mu}, \boldsymbol{\nu} \in \Delta_{\mathbb{N}}^{(1)}[h]$$

with the additional property that for any two sequences $\boldsymbol{\mu}_n, \boldsymbol{\nu}_n \in \Delta_{\mathbb{N}}$ satisfying $\varepsilon_n := \|\boldsymbol{\mu}_n - \boldsymbol{\nu}_n\|_1 \to 0$, it holds that $F(\varepsilon_n, h) \to 0$. Moreover, the upper bound in Theorem 1 is tight, up to a constant factor:

**Theorem 2.** *For every $0 < \varepsilon < 1/2$ and $\alpha \geq 1$, there are $\boldsymbol{\mu}, \boldsymbol{\nu} \in \Delta_{\mathbb{N}}$ such that $\varepsilon = \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_1 \geq \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_\infty$ and*

$$|\mathrm{H}(\boldsymbol{\mu}) - \mathrm{H}(\boldsymbol{\nu})| \geq c\varepsilon^{1-1/\alpha} \left( 2\mathrm{e}\varepsilon^{1/\alpha} \log \frac{1}{\varepsilon} + \mathrm{H}^{(\alpha)}(\boldsymbol{\mu})^{1/\alpha} + \mathrm{H}^{(\alpha)}(\boldsymbol{\nu})^{1/\alpha} \right), \tag{7}$$

*where $c \geq \frac{1}{2\mathrm{e}+2}$ is a universal constant.*

Perhaps surprisingly,[3] it turns out that $\mathrm{H} : \Delta_{\mathbb{N}}^{(\alpha)}[h] \to \mathbb{R}_+$ is uniformly continuous not only under $\ell_1$, but actually under all $\ell_p$ norms:

**Theorem 3.** *There is a function $F : \mathbb{R}_+^4 \to \mathbb{R}_+$ such that*

$$|\mathrm{H}(\boldsymbol{\mu}) - \mathrm{H}(\boldsymbol{\nu})| \leq F(\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_p, h, \alpha, p), \qquad h > 0, \alpha > 1, p \in [1, \infty], \boldsymbol{\mu}, \boldsymbol{\nu} \in \Delta_{\mathbb{N}}^{(\alpha)}[h]$$

*with the additional property that whenever $\varepsilon_n := \|\boldsymbol{\mu}_n - \boldsymbol{\nu}_n\|_p \to 0$, we have $F(\varepsilon_n, h, \alpha, p) \to 0$.*

---

[3]Since $\ell_1$ dominates all of the $\ell_p$ norms, continuity of a function under $\ell_p$ trivially implies continuity under $\ell_1$, but the reverse implication is generally not true.

Remark. Although Theorem 3 establishes uniform continuity, it gives no hint as to the functional dependence of the modulus of continuity $F$ on $\alpha$, $p$, $h$, and $\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_p$. We leave this as a fascinating open problem — even though the practical applications are likely to be limited: it follows from Wyner and Foster (2003) and Theorem 5 that for $p = \alpha = 2$ and fixed $h$, $F(\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_2, h, 2, 2)$ cannot decay at a faster rate than $1/\log(1/\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_2)$.

Combining Theorem 1 with (1) yields an empirical (under moment assumptions) bound for the plug-in entropy estimator:

**Corollary 1.** *For all $\alpha > 1$, $h > 0$, $\delta \in (0, 1)$, $n \geq 2\log\frac{4}{\delta}$, and $\boldsymbol{\mu} \in \Delta_{\mathbb{N}}^{(\alpha)}[h]$, we have that*

$$|\mathrm{H}(\boldsymbol{\mu}) - \mathrm{H}(\hat{\boldsymbol{\mu}}_n)| \;\leq\; \left(2\alpha^\alpha + h + \mathrm{H}^{(\alpha)}(\hat{\boldsymbol{\mu}}_n)\right)^{1/\alpha} \left(\frac{2\|\hat{\boldsymbol{\mu}}_n\|_{1/2}^{1/2}}{\sqrt{n}} + 6\sqrt{\frac{\log(4/\delta)}{2n}}\right)^{1-1/\alpha} \tag{8}$$

*holds with probability at least $1 - \delta$. For all $\alpha > 1$, $h > 0$, $0 < \varepsilon < \mathrm{e}^{-\alpha}$, $\delta \in (0, 1)$, $n \geq \frac{2}{\varepsilon^2}\log\frac{4}{\delta}$, and $\boldsymbol{\mu} \in \Delta_{\mathbb{N}}^{(\alpha)}[h]$, we have that*

$$\begin{aligned}
|\mathrm{H}(\boldsymbol{\mu}) - \mathrm{H}(\hat{\boldsymbol{\mu}}_n)| \;\leq\; & \left(2\mathrm{e}^\alpha \varepsilon \log^{1/\alpha}\frac{1}{\varepsilon}\right. \\
& \left. + \; h + \mathrm{H}^{(\alpha)}(\hat{\boldsymbol{\mu}}_n)\right)^{1/\alpha} \left(\frac{2\|\hat{\boldsymbol{\mu}}_n\|_{1/2}^{1/2}}{\sqrt{n}} + 6\sqrt{\frac{\log(4/\delta)}{2n}}\right)^{1-1/\alpha}
\end{aligned} \tag{9}$$

*holds with probability at least $1 - \delta$. For all $\alpha > 1$, $h > 0$, $\delta \in (0, 1)$, $n \geq 2e^{2\alpha}\log\frac{4}{\delta}$, and $\boldsymbol{\mu} \in \Delta_{\mathbb{N}}^{(\alpha)}[h]$, we have that*

$$\begin{aligned}
|\mathrm{H}(\boldsymbol{\mu}) - \mathrm{H}(\hat{\boldsymbol{\mu}}_n)| \;\leq\; & \left(2\left(\frac{\mathrm{e}}{2}\right)^\alpha \sqrt{\frac{2}{n}\log\frac{4}{\delta}}\left|\log\left(\frac{2}{n}\log\frac{4}{\delta}\right)\right|^\alpha\right. \\
& \left. + \; h + \mathrm{H}^{(\alpha)}(\hat{\boldsymbol{\mu}}_n)\right)^{1/\alpha} \left(\frac{2\|\hat{\boldsymbol{\mu}}_n\|_{1/2}^{1/2}}{\sqrt{n}} + 6\sqrt{\frac{\log(4/\delta)}{2n}}\right)^{1-1/\alpha}
\end{aligned} \tag{10}$$

*holds with probability at least $1 - \delta$.*

Remark. Since the estimates in Corollary 1 involve the random quantity $\|\hat{\boldsymbol{\mu}}_n\|_{1/2}$, it is natural to inquire as to the behavior of the latter. It follows from Cohen et al. (2020, Proposition C.1) that $n^{-1/2}\|\hat{\boldsymbol{\mu}}_n\|_{1/2}^{1/2} \to 0$ in almost surely. The rate of convergence must necessarily depend on $\boldsymbol{\mu}$ itself (cf. Berend and Kontorovich (2013, Remark 9)).

In Section 6, we compare the rates implied by Corollary 1 to the state of the art on various distributions.

Next, we examine the optimality of the plug-in estimate by analyzing the minimax risk, defined in (4). It was known (Silva, 2018, Appendix A) that assuming $\mathrm{H}(\boldsymbol{\mu}) < \infty$ does not suffice to yield a minimax rate for the $L_2$ risk:

$$\inf_{\hat{H}:\mathbb{N}^n \to \mathbb{R}_+} \sup_{\boldsymbol{\mu} \in \Delta_{\mathbb{N}}^{(1)}} \mathbb{E}\left(\hat{H}(X_1, \ldots, X_n) - \mathrm{H}(\boldsymbol{\mu})\right)^2 = \infty.$$

This technique yields an analogous result for the $L_1$ risk as well. We strengthen these results in two ways: (i) by lower-bounding the $L_1$ risk (rather than $L_2$, which is never smaller), and (ii) by restricting $\boldsymbol{\mu}$ to $\Delta_{\mathbb{N}}^{(1)}[h]$ and obtaining a finitary, quantitative lower bound:

**Theorem 4.** *For $\alpha = 1$, there is a universal constant $C > 0$ such that for all $h > 1$ and $n \in \mathbb{N}, n \geq 2$, we have $\mathcal{R}_n^{(1)}(h) \geq Ch$.*

Remark. The above result complements — but is not directly comparable to — Antos and Kontoyiannis (2001a, Theorem 4). Ours gives a quantitative dependence on $h$ but constructs an adversarial distribution for each sample size $n$; theirs is asymptotic only but a single adversarial distribution suffices for all $n$.

Remark. Our technique immediately yields a lower bound of $Ch^2$ on the $L_2$ minimax risk.

In contradistinction to the $\alpha = 1$ case, where no minimax rate exists, we show that the plug-in estimator is minimax for all $\alpha > 1$:

**Theorem 5.** *The following bounds hold for the $L_1$ minimax risk:*

*(a) Upper bound: for all $h > 0, \alpha > 1$,*

$$\mathcal{R}_n^{(\alpha)}(h) \leq \frac{1 + \log n}{\sqrt{n}} + \frac{2^{\alpha-1}h}{\log^{\alpha-1} n}, \qquad n \in \mathbb{N};$$

*further, this bound is achieved by the plug-in estimate $\mathrm{H}(\hat{\boldsymbol{\mu}}_n)$.*

*(b) Lower bound: for each $\alpha > 0$, $n \in \mathbb{N}$ there is an $h > 0$ such that*

$$\mathcal{R}_n^{(\alpha)}(h) \geq \frac{h}{4 \cdot 3^\alpha \log^{\alpha-1} n}.$$

**Open problem.** Close the gap in the dependence on $\alpha$ in the upper and lower bounds.


**Open problem.** Another gap between the upper and lower bounds is the quantified on $h$: in the upper bound, it is "for all", while in the lower bound, it is "exists". Closing this gap is also of interest.

Finally, in Section 7.3, we resolve most of the conjectures posed by Jürgensen and Matthews (2010).


## 4 Related work

**Continuity, convergence, moments of information** Zhang (2007) gave a sharpened version of (2) and Ho and Yeung (2010) presented analogous bounds; Audenaert (2007) proved a non-commutative generalization. Sason (2013, Theorem 5) upper-bounds $|\mathrm{H}(\boldsymbol{\mu}) - \mathrm{H}(\boldsymbol{\nu})|$ in terms of quantities related to $\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_1$, where (at most) one of them is allowed to have infinite support. Even though $\mathrm{H}(\cdot)$ is not continuous on $\Delta_{\mathbb{N}}$, the plug-in estimate $\mathrm{H}(\hat{\boldsymbol{\mu}}_n)$ converges to $\mathrm{H}(\boldsymbol{\mu})$ almost surely and in $L_2$ (Antos and Kontoyiannis, 2001a). Silva (2018) studied a variety of restrictions on distributions over infinite alphabets to derive strong consistency results and rates of convergence. Moments of information were apparently first defined in Golomb (1966).


**Entropy estimation** Recent surveys of entropy estimation results may be found in Jiao et al. (2015); Verdú (2019). The finite-alphabet case is particularly well-understood. For fixed alphabet size $d < \infty$, the plug-in estimate is asymptotically minimax optimal (Paninski, 2003). Paninski (2004) non-constructively established the existence of a sublinear (in $d$) entropy estimator. The optimal dependence on $d$ (at fixed accuracy) was settled by Valiant and Valiant (2011a, 2017) as being $\Theta(d/\log d)$.

The $\Theta(d/\log d)$ dependence on the alphabet size is also relevant in the so-called *high dimensional* asymptotic regime, where $d$ grows with $n$. Here, the plug-in estimate is no longer optimal, and more sophisticated techniques are called for (Valiant and Valiant, 2011a,b, 2017). The works of Wu and Yang (2016); Jiao et al. (2015); Han et al. (2015); Jiao et al. (2017) characterized the minimax rates for the high-dimensional regime: a small additive error of $\varepsilon$ requires $\Theta(d/\varepsilon \log d)$ samples. Building off of these polynomial-approximation based constructions, Acharya et al. (2017) design an additional optimal estimator, this one based on a profile maximum likelihood approach that can also estimate a variety of other important statistics. Fukuchi and Sakuma (2017, 2018) generalize the optimal estimators to estimate any additive functional, recovering in particular the optimal rates for entropy. Acharya et al. (2019) modify these optimal estimators with the added goal of low space complexity.

Finally, there is the infinite-alphabet case. Although here the plug-in estimate is again universally strongly consistent, control of the convergence rate requires some assumption on the sampling distribution — and Antos and Kontoyiannis (2001a) compellingly argue that moment assumptions are natural and minimalistic. Absent any prior assumptions, the $L_1$ (and hence $L_2$) convergence rate of *any* estimator can be made arbitrarily slow (Theorem 4 ibid.). The present paper proves a variant of this result (see Theorem 4 and the Remark following it). Antos and Kontoyiannis (2001a) further show that even under moment assumptions, there is no polynomial rate of convergence for the plug-in estimate: there is no $\beta > 0$ such that its risk decays as $O(n^{-\beta})$. Wyner and Foster (2003) showed that the plug-in estimate achieves a rate of $O(\frac{1}{\log n})$ for bounded second moment, and this is minimax optimal. Brautbar and Samorodnitsky (2007) exhibited a function of the higher moments that can be used in place of alphabet size to give a multiplicative approximation to the entropy.

The empirical nature of Corollary 1 can be seen as a distribution-dependent improvement over otherwise worst-case minimax guarantees. It can be compared, in this light, to the "instance-optimality" program of Hao et al. (2018); Hao and Orlitsky (2020) and the adaptive guarantees of Han et al. (2015).

## 5 Proofs

### 5.1 Proof of Theorem 1

We begin with a subadditivity result for the $\alpha$th moment of information (which we state for $\alpha > 0$, even though only the range $\alpha > 1$ will be needed).

**Lemma 1.** *For $\alpha > 0$ and $\boldsymbol{\mu}, \boldsymbol{\nu} \in \Delta_{\mathbb{N}}^{(\alpha)}$, we have*

$$\mathrm{H}^{(\alpha)}(|\boldsymbol{\mu} - \boldsymbol{\nu}|) \leq 2\alpha^{\alpha} + \mathrm{H}^{(\alpha)}(\boldsymbol{\mu}) + \mathrm{H}^{(\alpha)}(\boldsymbol{\nu}). \tag{11}$$

*If, additionally, $\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_{\infty} \leq \mathrm{e}^{-\alpha}$, then*

$$\mathrm{H}^{(\alpha)}(|\boldsymbol{\mu} - \boldsymbol{\nu}|) \leq 2\mathrm{e}^{\alpha} \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_{\infty} \log^{\alpha} \frac{1}{\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_{\infty}} + \mathrm{H}^{(\alpha)}(\boldsymbol{\mu}) + \mathrm{H}^{(\alpha)}(\boldsymbol{\nu}). \tag{12}$$

*Proof.* Define $\mathrm{h}^{(\alpha)} \colon [0,1] \to \mathbb{R}_+$ by $z \mapsto z \log^{\alpha}(1/z)$, where $\mathrm{h}^{(\alpha)}(0) = 0$. The function $\mathrm{h}^{(\alpha)}$ is increasing on $[0, \mathrm{e}^{-\alpha}]$ and decreasing on $[\mathrm{e}^{-\alpha}, 1]$. The maximum is therefore achieved at $z = \mathrm{e}^{-\alpha}$, and

$$\max_{z \in [0,1]} \mathrm{h}^{(\alpha)}(z) = \mathrm{h}^{(\alpha)}(\mathrm{e}^{-\alpha}) = \mathrm{e}^{-\alpha} \alpha^{\alpha}. \tag{13}$$

Now decompose $\mathrm{H}^{(\alpha)}$:

$$\mathrm{H}^{(\alpha)}(|\boldsymbol{\mu} - \boldsymbol{\nu}|) = \sum_{i: \boldsymbol{\mu}(i) \vee \boldsymbol{\nu}(i) > \mathrm{e}^{-\alpha}} \mathrm{h}^{(\alpha)}(|\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|) + \sum_{i: \boldsymbol{\mu}(i) \vee \boldsymbol{\nu}(i) \leq \mathrm{e}^{-\alpha}} \mathrm{h}^{(\alpha)}(|\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|). \tag{14}$$

To prove the lemma, we bound the two terms of (14) separately. The second term can be bound in two ways, yielding (11) and (12), respectively. To bound the first term of (14), notice that $\boldsymbol{\mu} \in \Delta_{\mathbb{N}}$ implies that $|\{i \in \mathbb{N} \colon \boldsymbol{\mu}(i) > \mathrm{e}^{-\alpha}\}| \leq \mathrm{e}^{\alpha}$, and similarly for $\boldsymbol{\nu}$. Thus,

$$\sum_{i: \boldsymbol{\mu}(i) \vee \boldsymbol{\nu}(i) > \mathrm{e}^{-\alpha}} \mathrm{h}^{(\alpha)}(|\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|) \leq \left( |\{i \colon \boldsymbol{\mu}(i) > \mathrm{e}^{-\alpha}\}| + |\{i \colon \boldsymbol{\nu}(i) > \mathrm{e}^{-\alpha}\}| \right) \max_{z \in [0,1]} \mathrm{h}^{(\alpha)}(z) \tag{15}$$

$$\leq 2\mathrm{e}^{\alpha} \mathrm{e}^{-\alpha} \alpha^{\alpha} = 2\alpha^{\alpha}.$$

For the second term of (14), notice that when $\boldsymbol{\mu}(i) \vee \boldsymbol{\nu}(i) \leq \mathrm{e}^{-\alpha}$, the monotonicity of $\mathrm{h}^{(\alpha)}$ implies

$$\mathrm{h}^{(\alpha)}(|\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|) \leq \mathrm{h}^{(\alpha)}(\boldsymbol{\mu}(i) \vee \boldsymbol{\nu}(i)),$$

and hence

$$\sum_{i \in \mathbb{N}: \boldsymbol{\mu}(i) \vee \boldsymbol{\nu}(i) \leq \mathrm{e}^{-\alpha}} \mathrm{h}^{(\alpha)}(|\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|) \leq \sum_{i: \boldsymbol{\mu}(i) \vee \boldsymbol{\nu}(i) \leq \mathrm{e}^{-\alpha}} \mathrm{h}^{(\alpha)}(\boldsymbol{\mu}(i) \vee \boldsymbol{\nu}(i))$$

$$\leq \sum_{i: \boldsymbol{\mu}(i) \vee \boldsymbol{\nu}(i) \leq \mathrm{e}^{-\alpha}} \mathrm{h}^{(\alpha)}(\boldsymbol{\mu}(i)) + \mathrm{h}^{(\alpha)}(\boldsymbol{\nu}(i))$$

$$\leq \mathrm{H}^{(\alpha)}(\boldsymbol{\mu}) + \mathrm{H}^{(\alpha)}(\boldsymbol{\nu});$$

6

this proves (11). Given the additional condition $\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_\infty \leq \mathrm{e}^{-\alpha}$, to prove (12), put $\varepsilon = \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_\infty$ and modify (15) as follows:

$$\sum_{i:\boldsymbol{\mu}(i)\vee\boldsymbol{\nu}(i)>\mathrm{e}^{-\alpha}} \mathrm{h}^{(\alpha)}(|\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|) \leq \left(\left|\{i\colon \boldsymbol{\mu}(i) > \mathrm{e}^{-\alpha}\}\right| + \left|\{i\colon \boldsymbol{\nu}(i) > \mathrm{e}^{-\alpha}\}\right|\right)\mathrm{h}^{(\alpha)}(\varepsilon)$$

$$\leq 2\mathrm{e}^{\alpha}\mathrm{h}^{(\alpha)}(\varepsilon).$$

The latter holds since $|\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)| \leq \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_\infty$, and so $\mathrm{h}^{(\alpha)}(|\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|) \leq \mathrm{h}^{(\alpha)}(\varepsilon)$, by $\mathrm{h}^{(\alpha)}$'s monotonicity on $[0, \mathrm{e}^{-\alpha}]$. $\qquad\square$

*Proof of Theorem 1.* The concavity argument in the proof of Cover and Thomas (2006, Theorem 17.3.3), immediately implies

$$|\mathrm{H}(\boldsymbol{\mu}) - \mathrm{H}(\boldsymbol{\nu})| \leq \mathrm{H}(|\boldsymbol{\mu} - \boldsymbol{\nu}|).$$

Then, via an application of Hölder's inequality,

$$
\begin{aligned}
\mathrm{H}(|\boldsymbol{\mu} - \boldsymbol{\nu}|) &= \sum_{i\in\mathbb{N}} |\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)| \log \frac{1}{|\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|} \\
&= \sum_{i\in\mathbb{N}} |\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|^{1-1/\alpha} \cdot |\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|^{1/\alpha} \log \frac{1}{|\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|} \\
&\leq \left(\sum_{i\in\mathbb{N}} \left(|\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|^{1-1/\alpha}\right)^{1/(1-1/\alpha)}\right)^{1-1/\alpha} \left(\sum_{i\in\mathbb{N}} \left(|\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|^{1/\alpha} \log \frac{1}{|\boldsymbol{\mu}(i) - \boldsymbol{\nu}(i)|}\right)^{\alpha}\right)^{1/\alpha} \\
&= \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_1^{1-1/\alpha} \mathrm{H}^{(\alpha)}(|\boldsymbol{\mu} - \boldsymbol{\nu}|)^{1/\alpha}.
\end{aligned}
$$

The claim follows by invoking Lemma 1 and the subadditivity of $t \mapsto t^{1/\alpha}$ for $t \geq 0$ and $\alpha > 1$. $\qquad\square$

## 5.2 Proof of Theorem 2

Let $0 < \varepsilon \leq 1/2$ be given and choose $\boldsymbol{\mu}, \boldsymbol{\nu} \in \Delta_\mathbb{N}$ as follows: $\boldsymbol{\mu} = (1, 0, 0, \ldots)$ and $\boldsymbol{\nu} = (1 - \varepsilon, \varepsilon, 0, 0, \ldots)$. Then left-hand side of (7) is $L(\varepsilon) = \mathrm{H}(\boldsymbol{\nu})$:

$$L(\varepsilon) = (1 - \varepsilon)\log\frac{1}{1-\varepsilon} + \varepsilon\log\frac{1}{\varepsilon} =: L_1(\varepsilon) + L_2(\varepsilon),$$

and right-hand side of (7), without the constant $c$, is

$$
\begin{aligned}
R(\varepsilon) &= \varepsilon^{1-1/\alpha}\left(2\mathrm{e}\varepsilon^{1/\alpha}\log\frac{1}{\varepsilon} + \left((1-\varepsilon)\left|\log\frac{1}{1-\varepsilon}\right|^\alpha + \varepsilon\left|\log\frac{1}{\varepsilon}\right|^\alpha\right)^{1/\alpha}\right) \\
&\leq \varepsilon^{1-1/\alpha}\left(2\mathrm{e}\varepsilon^{1/\alpha}\log\frac{1}{\varepsilon} + 2^{1/\alpha}\max\left\{(1-\varepsilon)^{1/\alpha}\log\frac{1}{1-\varepsilon}, \varepsilon^{1/\alpha}\log\frac{1}{\varepsilon}\right\}\right) \\
&\leq 2\varepsilon^{1-1/\alpha}(1-\varepsilon)^{1/\alpha}\log\frac{1}{1-\varepsilon} + (2\mathrm{e}+2)\varepsilon\log\frac{1}{\varepsilon} \\
&=: R_1(\varepsilon) + R_2(\varepsilon).
\end{aligned}
$$

Now $L_1(\varepsilon) \geq R_1(\varepsilon)/2$ for $\varepsilon \in (0, \frac{1}{2}]$ and $L_2(\varepsilon) = \frac{R_2(\varepsilon)}{2\mathrm{e}+2}$, and therefore $L_1(\varepsilon) + L_2(\varepsilon) \geq R_1(\varepsilon)/2 + \frac{1}{2\mathrm{e}+2}R_2(\varepsilon) \geq \frac{1}{2\mathrm{e}+2}R(\varepsilon)$. $\qquad\square$

## 5.3 Proof of Theorem 3

The following fact (Lieb and Loss, 2001, Theorem 3.5 and Eq. (5) on p. 83) will be useful[4]:

$$\|\boldsymbol{\mu}^\downarrow - \boldsymbol{\nu}^\downarrow\|_p \leq \|\boldsymbol{\mu} - \boldsymbol{\nu}\|_p, \qquad p \in [1, \infty], \; \boldsymbol{\mu}, \boldsymbol{\nu} \in \Delta_\mathbb{N}. \tag{16}$$

---

[4]The result is stated for functions in $f \in L_2(\mathbb{R}^n)$ and their symmetric-decreasing rearrangements $f^*$, but the specialization to discrete distributions is straightforward. We convert $\boldsymbol{\mu}$ to a function $f\colon \mathbb{R}_+ \to \mathbb{R}_+$ via $f(x) = \mu(\lceil x\rceil)$ and $\boldsymbol{\nu}$ to $g(x)$ analogously. A direct calculation then shows that $\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_p = \|f - g\|_p$ and $\|\boldsymbol{\mu}^\downarrow - \boldsymbol{\nu}^\downarrow\|_p = \|f^* - g^*\|_p$, to which the result from Lieb and Loss (2001) applies to yield (16).

A result of Scheffé (1947) (more accurately credited to Riesz, 1928 (Kusolitsch, 2010)) implies that a sequence $\{\xi_{n\in\mathbb{N}}\} \subset \ell_1(\mathbb{N})$ converging pointwise to some $\xi \in \ell_1(\mathbb{N})$ also converges in $\ell_1$ iff $\|\xi_n\|_1 \to \|\xi\|_1$. This immediately implies

**Lemma 2.** *If $\{\boldsymbol{\mu}_{n\in\mathbb{N}}\} \subset \Delta_{\mathbb{N}}$ converges pointwise to some $\boldsymbol{\mu} \in \Delta_{\mathbb{N}}$, then it also converges in $\ell_1$.*

Berend et al. (2017, Lemma 1) showed that $\Delta_{\mathbb{N}}^{\downarrow(1)}[h]$ is compact under $\ell_1$. We begin by extending this result to general $\alpha, p$.

**Lemma 3.** *For all $\alpha \geq 1$, $p \in [1, \infty]$, and $h > 0$, the set $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ is compact under $\ell_p$.*

Remark. This is quite false if either the non-increasing or the bounded-entropy condition is omitted. For a counterexample to the former, consider the sequence $\boldsymbol{\mu}_n \in \Delta_{\mathbb{N}}$ defined by $\boldsymbol{\mu}_n(i) = \mathbf{1}[i = n]$. For a counterexample to the latter, consider the sequence $\boldsymbol{\mu}_n \in \Delta_{\mathbb{N}}$, where $\boldsymbol{\mu}_n$ is uniform on $[n]$.

*Proof.* We closely follow the proof strategy of Berend et al. (2017, Lemma 1). In a metric space, compactness and sequential compactness are equivalent. Let $\boldsymbol{\mu}_{n\in\mathbb{N}}$ be a sequence in $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$. Since $[0, 1]$ is compact, every $\{\boldsymbol{\mu}_n(i) : n \in \mathbb{N}\}$ has a convergent subsequence, and hence $\boldsymbol{\mu}_{n\in\mathbb{N}}$ has a pointwise convergent subsequence. There is thus no loss of generality in assuming that $\boldsymbol{\mu}_n \to \boldsymbol{\mu}$ pointwise. Obviously, $\boldsymbol{\mu}$ is non-negative and non-increasing. It remains to show that

(a) $\sum_{i\in\mathbb{N}} \boldsymbol{\mu}(i) = 1$,

(b) $\mathrm{H}^{(\alpha)}(\boldsymbol{\mu}) \leq h$,

(c) $\|\boldsymbol{\mu}_n - \boldsymbol{\mu}\|_p \to 0$.

To show (a), assume, for a contradiction, that $\sum_{i\in\mathbb{N}} \boldsymbol{\mu}(i) > 1$. Then there must be an $i_0 \in \mathbb{N}$ such that $\sum_{i=1}^{i_0} \boldsymbol{\mu}(i) > 1$. But the latter must then hold for all $\boldsymbol{\mu}_n$ with $n$ sufficiently large, which contradicts $\boldsymbol{\mu}_n \in \Delta_{\mathbb{N}}$. Now assume $\varepsilon := 1 - \sum_{i\in\mathbb{N}} \boldsymbol{\mu}(i) > 0$. For any $i_0 \in \mathbb{N}$, we have $\sum_{i=1}^{i_0} \boldsymbol{\mu}_n(i) < 1 - \varepsilon/2$ for all sufficiently large $n$. Now every $\boldsymbol{\nu} \in \Delta_{\mathbb{N}}^{\downarrow}$ satisfies $\boldsymbol{\nu}(i) \leq \frac{1}{i}(\boldsymbol{\nu}(1) + \boldsymbol{\nu}(2) + \ldots + \boldsymbol{\nu}(i)) \leq \frac{1}{i}$. Hence,

$$\sum_{i=i_0+1}^{\infty} \boldsymbol{\mu}_n(i) \left|\log \boldsymbol{\mu}_n(i)\right|^\alpha \geq \sum_{i=i_0+1}^{\infty} \boldsymbol{\mu}_n(i)(\log i_0)^\alpha > \frac{\varepsilon}{2}(\log i_0)^\alpha.$$

Choosing $i_0$ sufficiently large makes the latter expression exceed $h$, violating the assumption $\boldsymbol{\mu}_n \in \Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$. Thus (a) holds.

To show (b), assume for a contradiction that $\mathrm{H}^{(\alpha)}(\boldsymbol{\mu}) > h$ — and, in particular, $\sum_{i=1}^{i_0} \boldsymbol{\mu}(i) \left|\log \boldsymbol{\mu}(i)\right|^\alpha > h$ for some $i_0 \in \mathbb{N}$. But the latter must hold for all $\boldsymbol{\mu}_n$ with $n$ sufficiently large, a contradiction.

Finally, to show (c), we invoke Lemma 2: if $\{\boldsymbol{\mu}_{n\in\mathbb{N}}\} \subset \Delta_{\mathbb{N}}$ converges pointwise to some $\boldsymbol{\mu} \in \Delta_{\mathbb{N}}$, then it also converges in $\ell_1$. Since $\ell_1$ dominates every $\ell_p$, $p > 1$, this proves (c). $\square$

Next, we examine the continuity of $\mathrm{H}(\cdot)$ on $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ under $\ell_p$.

**Lemma 4.** *Fix $h > 0$, $\alpha > 1$, and $p \in [1, \infty]$. If $\{\boldsymbol{\mu}_{n\in\mathbb{N}}\} \subset \Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ converges in $\ell_p$, then its limit is some $\boldsymbol{\mu} \in \Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ and furthermore, $\mathrm{H}(\boldsymbol{\mu}_n) \to \mathrm{H}(\boldsymbol{\mu})$. In other words, $\mathrm{H}(\cdot)$ is continuous on $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ under $\ell_p$.*

Remark. We note that $\mathrm{H}(\cdot)$ is not continuous on $\Delta_{\mathbb{N}}^{\downarrow(1)}[h]$ under $\ell_p$, $p \in [1, \infty]$, as evidenced by the sequence $\boldsymbol{\mu}_n = (1 - \varepsilon_n, \varepsilon_n/n, \ldots, \varepsilon/n, 0, 0, \ldots)$, with support size $n + 1$. We can choose $\varepsilon_n$ so that $\mathrm{H}(\boldsymbol{\mu}_n) = h$, but of course the limiting $\boldsymbol{\mu}$ has $\mathrm{H}(\boldsymbol{\mu}) = 0$ (see Example 1 in Berend et al. (2017)).

*Proof.* It follows from Lemma 3 that the limiting $\boldsymbol{\mu}$ belongs to $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$. Further, Lemma 2 implies that $\boldsymbol{\mu}_n \to \boldsymbol{\mu}$ in $\ell_1$. Invoking the continuity result in Theorem 1 proves the claim. $\square$

*Proof of Theorem 3.* It follows from Lemma 4 that $H(\cdot)$ is continuous on $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ under $\ell_p$. Since, by Lemma 3, $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ is compact under $\ell_p$, it follows that $H(\cdot)$ is uniformly continuous on $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$: there is a function $F$ such that

$$|H(\boldsymbol{\mu}) - H(\boldsymbol{\nu})| \leq F(\|\boldsymbol{\mu} - \boldsymbol{\nu}\|_p, h, \alpha, p), \qquad \boldsymbol{\mu}, \boldsymbol{\nu} \in \Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$$

and $\varepsilon_n := \|\boldsymbol{\mu}_n - \boldsymbol{\nu}_n\|_p \to 0 \implies F(\varepsilon_n, h, \alpha, p) \to 0$. Now, for all $\boldsymbol{\mu}, \boldsymbol{\nu} \in \Delta_{\mathbb{N}}^{(\alpha)}[h]$ we have

$$|H(\boldsymbol{\mu}) - H(\boldsymbol{\nu})| = \left|H(\boldsymbol{\mu}^{\downarrow}) - H(\boldsymbol{\nu}^{\downarrow})\right| \leq F(\|\boldsymbol{\mu}^{\downarrow} - \boldsymbol{\nu}^{\downarrow}\|_p, h, \alpha, p).$$

It follows from (16) that $\|\boldsymbol{\mu}_n - \boldsymbol{\nu}_n\|_p \to 0 \implies \|\boldsymbol{\mu}_n^{\downarrow} - \boldsymbol{\nu}_n^{\downarrow}\|_p \to 0$, which concludes the proof. $\square$

### 5.4 Proof of Corollary 1

*Proof.* Fix $0 < \varepsilon < e^{-\alpha}$. Consider two potential "bad" events: $B_1$, where $\|\hat{\boldsymbol{\mu}}_n - \boldsymbol{\mu}\|_\infty > \varepsilon$, and $B_2$, where $\|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_n\|_1 > \frac{2\|\hat{\boldsymbol{\mu}}_n\|_{1/2}^{1/2}}{\sqrt{n}} + 6\sqrt{\frac{\log(4/\delta)}{2n}}$. Our assumption on the sample size $n$, together with the Dvoretzky-Kiefer-Wolfowitz inequality (Massart, 1990), implies that $\mathbb{P}(B_1) \leq \delta/2$ and (1) implies that $\mathbb{P}(B_2) \leq \delta/2$. Thus, with probability at least $1 - \delta$, neither of $B_1$ or $B_2$ occurs, and on the event where $\|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_n\|_\infty < \varepsilon < e^{-\alpha}$, we may invoke Theorem 1, from which the claims immediately follow. $\square$

### 5.5 Proof of Theorem 4

For $h > 1$ and $n \in \mathbb{N}, n \geq 2$, put $a_n = (1 - 1/(2n))\log(1 - 1/(2n))$ and define the support size $S = S(h, n)$ by $S = \lfloor(1/2n)\exp(2n(h + a_n))\rfloor$. Consider the distributions $\boldsymbol{\mu}_0 = (1, 0, 0, \dots)$ and $\boldsymbol{\mu}_n$ defined by $\boldsymbol{\mu}_n(1) = 1 - 1/(2n)$, and

$$\boldsymbol{\mu}_n(i) = \frac{1}{2nS}, \qquad 2 \leq i \leq 1 + S(h, n).$$

We compute the Kullback-Leibler divergence and entropy:

$$\begin{aligned}
D_{\mathrm{KL}}(\boldsymbol{\mu}_0 \| \boldsymbol{\mu}_n) &= \log\frac{1}{1 - 1/(2n)} \leq \frac{1}{1 - 1/(2n)} - 1 \leq \frac{1}{n} \qquad (17) \\
H(\boldsymbol{\mu}_0) &= 0 \leq h.
\end{aligned}$$

For $x \geq 2$, always $\lfloor x \rfloor \geq x/2$. Additionally, from $2na_n \geq -1$, and $\frac{1}{2n}\exp(2nh - 1) > 2$, we obtain that $S > (1/4n)\exp(2n(h + a_n))$, hence we also have that $h \geq H(\boldsymbol{\mu}_n) > h - \frac{1}{2n}\log 2$. Since $\frac{1}{2x}\log 2 \leq 1/2$ on $[1, \infty)$ and $h > 1$, it follows that $H(\boldsymbol{\mu}_n) \geq \frac{h}{2}$, whence $|H(\boldsymbol{\mu}_0) - H(\boldsymbol{\mu}_n)| \geq h/2$. To bound the $L_1$ minimax risk (defined in (4)), we invoke Markov's inequality:

$$\mathbb{E}|\hat{H}(X_1, \dots, X_n) - H(\boldsymbol{\mu})| \geq \frac{h}{4}\mathbb{P}\left(|\hat{H}(X_1, \dots, X_n) - H(\boldsymbol{\mu})| > \frac{h}{4}\right).$$

It follows via Le Cam's two point method (Tsybakov, 2008, Section 2.4.2) that

$$\mathcal{R}_n^{(1)}(h) \geq \frac{h}{8}e^{-nD_{\mathrm{KL}}(\boldsymbol{\mu}_0\|\boldsymbol{\mu})} \geq \frac{h}{8e},$$

where the second inequality stems from (17).

$\square$

### 5.6 Proof of Theorem 5

We begin with an auxiliary lemma, of possible independent interest.

**Lemma 5.** *For all $\boldsymbol{\mu} \in \Delta_{\mathbb{N}}$ and $n \in \mathbb{N}$, we have*

$$H(\boldsymbol{\mu}) \geq \mathbb{E}H(\hat{\boldsymbol{\mu}}_n) \geq H(\boldsymbol{\mu}) - \inf_{0 < \varepsilon < 1}\left[\sum_{i \in \mathbb{N}: \boldsymbol{\mu}(i) < \varepsilon} \boldsymbol{\mu}(i)\log\frac{1}{\boldsymbol{\mu}(i)} + \log\left(1 + \frac{1}{\varepsilon n}\right)\right].$$

*Proof.* The first inequality follows from Jensen's, since $H(\cdot)$ is concave and $\mathbb{E}\hat{\boldsymbol{\mu}}_n = \boldsymbol{\mu}$. To prove the second inequality, choose $\varepsilon > 0$, put $J := \{i \in \mathbb{N} : \boldsymbol{\mu}(i) < \varepsilon\}$, and compute

$$
\begin{aligned}
\mathbb{E}H(\hat{\boldsymbol{\mu}}_n) &= \mathbb{E}\left[\sum_{i \in \mathbb{N}\backslash J} \hat{\boldsymbol{\mu}}_n(i) \log \frac{1}{\hat{\boldsymbol{\mu}}_n(i)} + \sum_{i \in J} \hat{\boldsymbol{\mu}}_n(i) \log \frac{1}{\hat{\boldsymbol{\mu}}_n(i)}\right] \\
&\geq \mathbb{E}\left[\sum_{i \in \mathbb{N}\backslash J} \hat{\boldsymbol{\mu}}_n(i) \log \frac{1}{\hat{\boldsymbol{\mu}}_n(i)} + \left(\sum_{i \in J} \hat{\boldsymbol{\mu}}_n(i)\right) \log \frac{1}{\sum_{i \in J} \hat{\boldsymbol{\mu}}_n(i)}\right] \\
&=: \mathbb{E}H(\tilde{\boldsymbol{\mu}}_n),
\end{aligned}
$$

where $\tilde{\boldsymbol{\mu}}_n$ is the "collapsed" version of $\hat{\boldsymbol{\mu}}_n$, where all of the masses in $J$ have been replaced by a single mass equal to their sum, and the inequality holds because conditioning reduces entropy (Cover and Thomas, 2006, Eq.(2.157)). We observe that $\tilde{\boldsymbol{\mu}}_n$ has support size at most $1 + 1/\varepsilon$ and invoke Paninski (2003, Proposition 1):

$$
\mathbb{E}H(\tilde{\boldsymbol{\mu}}_n) \geq H(\tilde{\boldsymbol{\mu}}) - \log\left(1 + \frac{1}{\varepsilon n}\right), \tag{18}
$$

where $\tilde{\boldsymbol{\mu}}$ is the "collapsed" version of $\boldsymbol{\mu}$. Now

$$
\begin{aligned}
H(\tilde{\boldsymbol{\mu}}) &= H(\boldsymbol{\mu}) + \left(\sum_{i \in j} \boldsymbol{\mu}(i)\right) \log \frac{1}{\sum_{i \in J} \boldsymbol{\mu}(i)} - \sum_{i \in J} \boldsymbol{\mu}(i) \log \frac{1}{\boldsymbol{\mu}(i)} \\
&\geq H(\boldsymbol{\mu}) - \sum_{i \in J} \boldsymbol{\mu}(i) \log \frac{1}{\boldsymbol{\mu}(i)},
\end{aligned}
$$

which concludes the proof. $\qquad\square$

The first part of the theorem will follow from the following proposition.

**Proposition 1.** *For $\alpha \geq 1$, $h > 0$, $n \in \mathbb{N}$ and $\boldsymbol{\mu} \in \Delta_{\mathbb{N}}^{(\alpha)}[h]$, we have*

$$
\mathbb{E}|H(\boldsymbol{\mu}) - H(\hat{\boldsymbol{\mu}}_n)| \leq \frac{\log n}{\sqrt{n}} + \inf_{0 < \varepsilon < 1}\left[\left(\log \frac{1}{\varepsilon}\right)^{1-\alpha} h + \log\left(1 + \frac{1}{\varepsilon n}\right)\right].
$$

*Proof.* Since by Lemma 5, $|H(\boldsymbol{\mu}) - \mathbb{E}H(\hat{\boldsymbol{\mu}}_n)| = H(\boldsymbol{\mu}) - \mathbb{E}H(\hat{\boldsymbol{\mu}}_n)$, it follows from the triangle and Jensen inequalities that

$$
\begin{aligned}
\mathbb{E}|H(\boldsymbol{\mu}) - H(\hat{\boldsymbol{\mu}}_n)| &\leq \mathbb{E}|H(\hat{\boldsymbol{\mu}}_n) - \mathbb{E}H(\hat{\boldsymbol{\mu}}_n)| + H(\boldsymbol{\mu}) - \mathbb{E}H(\hat{\boldsymbol{\mu}}_n) \\
&\leq \sqrt{\mathbb{V}\mathrm{ar}\left[H(\hat{\boldsymbol{\mu}}_n)\right]} + H(\boldsymbol{\mu}) - \mathbb{E}H(\hat{\boldsymbol{\mu}}_n) \\
&\leq \frac{\log n}{\sqrt{n}} + H(\boldsymbol{\mu}) - \mathbb{E}H(\hat{\boldsymbol{\mu}}_n), \tag{19}
\end{aligned}
$$

where the variance bound is from Antos and Kontoyiannis (2001b, Proposition 1(iv)).

For any $\varepsilon > 0$, Lemma 5 implies

$$
\begin{aligned}
\mathbb{E}H(\hat{\boldsymbol{\mu}}_n) &\geq H(\boldsymbol{\mu}) - \sum_{i \in \mathbb{N}:\boldsymbol{\mu}(i)<\varepsilon} \boldsymbol{\mu}(i) \log \frac{1}{\boldsymbol{\mu}(i)} - \log\left(1 + \frac{1}{\varepsilon n}\right) \\
&\geq H(\boldsymbol{\mu}) - \left(\log \frac{1}{\varepsilon}\right)^{1-\alpha} \sum_{i \in \mathbb{N}:\boldsymbol{\mu}(i)<\varepsilon} \boldsymbol{\mu}(i) \left(\log \frac{1}{\boldsymbol{\mu}(i)}\right)^{\alpha} - \log\left(1 + \frac{1}{\varepsilon n}\right) \\
&\geq H(\boldsymbol{\mu}) - \left(\log \frac{1}{\varepsilon}\right)^{1-\alpha} H^{(\alpha)}(\boldsymbol{\mu}) - \log\left(1 + \frac{1}{\varepsilon n}\right), \tag{20}
\end{aligned}
$$

where the second and third inequalities follow from the obvious relations

$$\sum_{i:\boldsymbol{\mu}(i)<\varepsilon}\boldsymbol{\mu}(i)\log\frac{1}{\boldsymbol{\mu}(i)} \leq \left(\log\frac{1}{\varepsilon}\right)^{1-\alpha}\sum_{i:\boldsymbol{\mu}(i)<\varepsilon}\boldsymbol{\mu}(i)\left(\log\frac{1}{\boldsymbol{\mu}(i)}\right)^{\alpha}$$

$$\leq \left(\log\frac{1}{\varepsilon}\right)^{1-\alpha}\mathrm{H}^{(\alpha)}(\boldsymbol{\mu}).$$

The claim follows by combining (19) with (20). □

*Proof of Theorem 5(a).* Use the fact that $\mathcal{R}_n^{(\alpha)}(h) \leq \mathbb{E}|\mathrm{H}(\boldsymbol{\mu}) - \mathrm{H}(\hat{\boldsymbol{\mu}}_n)|$, invoke Proposition 1 with $\varepsilon = \frac{1}{\sqrt{n}}$ and use $\log(1+x) \leq x$. □

We now prove the second half of the theorem.

*Proof of Theorem 5(b).* Let $\alpha > 0$, $n \in \mathbb{N}$ and define two families of distributions:

$$\mathcal{U}_1 := \left\{\boldsymbol{\mu}_1 = \mathrm{Uniform}([n^3])\right\},$$
$$\mathcal{U}_2 := \left\{\boldsymbol{\mu}_2 = \mathrm{Uniform}(A) : A \subset [n^3], |A| = n^2\right\}.$$

Let $h := 3^\alpha \log^\alpha n$ and note that $\mathcal{U}_1 \cup \mathcal{U}_2 \subseteq \Delta_{\mathbb{N}}^{(\alpha)}[h]$. Let $E$ be the event that $\boldsymbol{X} = (X_1, \ldots, X_n)$ has no repeating elements, i.e $|\{X_1, X_2, \ldots, X_n\}| = n$. Let $\boldsymbol{\mu}_1 \in \mathcal{U}_1, \boldsymbol{\mu}_2 \in \mathcal{U}_2$ and consider the values $\mathbb{P}_{\boldsymbol{X}\sim\boldsymbol{\mu}_1^n}(E)$ and $\mathbb{P}_{\boldsymbol{X}\sim\boldsymbol{\mu}_2^n}(E)$. For $m \in \mathbb{N}$, define $K(m)$ to be the smallest $k$ such that when uniformly throwing $m$ balls into $k$ buckets, the probability of collision is at least $1/2$. Since $K(m)$ is known[5] to be at least $\sqrt{m}$ (and hence $K(n^2) > n$) we have a lower bound of $\frac{1}{2}$ on both $\mathbb{P}_{\boldsymbol{X}\sim\boldsymbol{\mu}_1^n}(E)$ and $\mathbb{P}_{\boldsymbol{X}\sim\boldsymbol{\mu}_2^n}(E)$. Define $\boldsymbol{\mu}_1^n|E$ as the distribution on $\mathbb{N}^n$ induced by conditioning the product $\boldsymbol{\mu}_1^n$ on the event $E$, and define $\boldsymbol{\mu}_2^n|E$ analogously. Our key observation is that conditional on $E$, $\boldsymbol{\mu}_1^n$ is uniform on $([n^3])_n$ whereas $\boldsymbol{\mu}_2^n = \mathrm{Uniform}(A)^n$ is uniform on $(A)_n$, where $(J)_k := \left\{(x_1, \ldots, x_k) \in J^k : |\{x_1, \ldots, x_k\}| = k\right\}$ is the set of all possible ordered samples of size $k$ from a distribution supported on $J$. To verify this observation, take $x = (x_1, \ldots, x_n) \in [n^3]$ and note that $(\boldsymbol{\mu}_1^n|E)(x) = \mathbb{P}_{\boldsymbol{X}\sim\boldsymbol{\mu}_1^n}(\{X = x\} \cap E) / \mathbb{P}_{\boldsymbol{X}\sim\boldsymbol{\mu}_1^n}(E) = \mathbb{P}_{\boldsymbol{X}\sim\boldsymbol{\mu}_1^n}(\{X = x\}) / \mathbb{P}_{\boldsymbol{X}\sim\boldsymbol{\mu}_1^n}(E)$ if $x$ has no repeating elements (meaning $\{X = x\} \subseteq E$) and $(\boldsymbol{\mu}_1^n|E)(x) = 0$ otherwise. Then

$$\mathcal{R}_n^{(\alpha)}(h) \geq \inf_{\hat{H}}\sup_{\boldsymbol{\mu}\in\mathcal{U}_1\cup\mathcal{U}_2}\mathbb{E}_{\boldsymbol{X}\sim\boldsymbol{\mu}^n}\left[|\hat{H}(\boldsymbol{X}) - \mathrm{H}(\boldsymbol{\mu})|\right]$$

$$\overset{(a)}{\geq} \inf_{\hat{H}}\sup_{\boldsymbol{\mu}\in\mathcal{U}_1\cup\mathcal{U}_2}\mathbb{E}_{\boldsymbol{X}\sim\boldsymbol{\mu}^n|E}\left[|\hat{H}(\boldsymbol{X}) - \mathrm{H}(\boldsymbol{\mu})|\right]\mathbb{P}_{\boldsymbol{X}\sim\boldsymbol{\mu}^n}(E)$$

$$\geq \inf_{\hat{H}}\frac{1}{2}\sup_{\boldsymbol{\mu}\in\mathcal{U}_1\cup\mathcal{U}_2}\mathbb{E}_{\boldsymbol{X}\sim\boldsymbol{\mu}^n|E}\left[|\hat{H}(\boldsymbol{X}) - \mathrm{H}(\boldsymbol{\mu})|\right]$$

$$\overset{(b)}{\geq} \inf_{\hat{H}}\frac{1}{4}\left(\mathbb{E}_{\boldsymbol{X}\sim\boldsymbol{\mu}_1^n|E}\left[|\hat{H}(\boldsymbol{X}) - \mathrm{H}(\boldsymbol{\mu}_1)|\right] + \sup_{\boldsymbol{\mu}_2\in\mathcal{U}_2}\mathbb{E}_{\boldsymbol{X}\sim\boldsymbol{\mu}_2^n|E}\left[|\hat{H}(\boldsymbol{X}) - \mathrm{H}(\boldsymbol{\mu}_2)|\right]\right)$$

$$\overset{(c)}{\geq} \inf_{\hat{H}}\frac{1}{4}\left(\mathbb{E}_{\boldsymbol{X}\sim\boldsymbol{\mu}_1^n|E}\left[|\hat{H}(\boldsymbol{X}) - \mathrm{H}(\boldsymbol{\mu}_1)|\right] + \mathbb{E}_{\boldsymbol{\mu}_2\sim\mathrm{Uniform}(\mathcal{U}_2)}\left[\mathbb{E}_{\boldsymbol{X}\sim\boldsymbol{\mu}_2^n|E}\left[|\hat{H}(\boldsymbol{X}) - \mathrm{H}(\boldsymbol{\mu}_2)|\right]\right]\right)$$

$$\overset{(d)}{=} \inf_{\hat{H}}\frac{1}{4}\left(\mathbb{E}_{\boldsymbol{X}\sim\boldsymbol{\mu}_1^n|E}\left[|\hat{H}(\boldsymbol{X}) - \mathrm{H}(\boldsymbol{\mu}_1)|\right] + \mathbb{E}_{\boldsymbol{X}\sim\boldsymbol{\mu}_1^n|E}\left[|\hat{H}(\boldsymbol{X}) - \mathrm{H}(\boldsymbol{\mu}_2)|\right]\right)$$

$$= \inf_{\hat{H}}\frac{1}{4}\left(\mathbb{E}_{\boldsymbol{X}\sim\boldsymbol{\mu}_1^n|E}\left[|\hat{H}(\boldsymbol{X}) - \mathrm{H}(\boldsymbol{\mu}_1)| + |\hat{H}(\boldsymbol{X}) - \mathrm{H}(\boldsymbol{\mu}_2)|\right]\right)$$

$$\overset{(e)}{\geq} \frac{1}{4}|\mathrm{H}(\boldsymbol{\mu}_1) - \mathrm{H}(\boldsymbol{\mu}_2)| = \frac{1}{4}\log n = \frac{1}{4}\frac{h}{3^\alpha \log^{\alpha-1} n},$$

where (a) is from the law of total expectation (the complement of $E$ is discarded), (b) and (c) are bounding a supremum by an average, (e) is from the triangle inequality, and (d) is by observing that,

---

[5]Better bounds exist (Brink, 2012).

by symmetry, the operators $\mathbb{E}_{\boldsymbol{\mu}_2 \sim \text{Uniform}(\mathcal{U}_2)} \left[ \mathbb{E}_{\boldsymbol{X} \sim \boldsymbol{\mu}_2^n | E} [\cdot] \right] = \mathbb{E}_{A \sim \text{Uniform}([n^2])} \left[ \mathbb{E}_{\boldsymbol{X} \sim \text{Uniform}((A)_n)} [\cdot] \right]$ and $\mathbb{E}_{\boldsymbol{X} \sim \boldsymbol{\mu}_1^n | E} [\cdot] = \mathbb{E}_{\boldsymbol{X} \sim \text{Uniform}([n^3])} [\cdot]$ are equivalent. (There is a minor abuse of notation in transitions after (c), since we write $\boldsymbol{\mu}_2$ without specifying a *particular* member of $\mathcal{U}_2$. However, $\boldsymbol{\mu}_2$ only occurs therein as $\text{H}(\boldsymbol{\mu}_2)$, and this value is identical for all $\boldsymbol{\mu}_2 \in \mathcal{U}_2$.) $\qquad\square$

# 6 Comparative Rates

Our bounds have the crucial characteristic of being empirical (under moment assumptions). When we *observe* favorable distributions (even without a priori knowledge of the fact), we will benefit from tighter bounds. This entails some cost, and in the worst case our bounds will be sub-optimal. In this section, we illustrate these trade-offs for various natural classes of distributions.

For the class of all finite alphabet distributions, our bound is sub-optimal. The MLE (plug-in estimator) is competitive with the optimal estimator up to logarithmic factors in $d$, but our bounds on the MLE are loose nearly quadratically in $d/n$, in the worst case. The convergence of the empirical distribution on a finite alphabet in $\ell_1$ occurs at rate $\Theta(\sqrt{d/n})$, whereas the MLE entropy estimator converges at rate $O\left( \sqrt{\left(\frac{d}{n}\right)^2 + \frac{\log^2 d}{n}} \right)$, as follows from Wu and Yang (2016, Proposition 1). So any approach that upper bounds the entropy risk via $\ell_1$ (as our Theorem 1 or Section 4 of Ho and Yeung (2010)) will be worst-case suboptimal for this class of distributions.

Nevertheless, for certain classes of distributions our bounds (Theorem 1 and Corollary 1) can significantly outperform the state of the art, for small and moderate-sized samples. To calculate the expected rate of our approach, we apply Hölder's inequality, as in the proof of Theorem 1:

$$\mathbb{E}|\text{H}(\hat{\boldsymbol{\mu}}_n) - \text{H}(\boldsymbol{\mu})| \quad \leq \quad \left( \mathbb{E}\left[ 2\alpha^\alpha + \text{H}^{(\alpha)}(\boldsymbol{\mu}) + \text{H}^{(\alpha)}(\hat{\boldsymbol{\mu}}_n) \right] \right)^{1/\alpha} \left( \mathbb{E}\|\hat{\boldsymbol{\mu}}_n - \boldsymbol{\mu}\|_1 \right)^{1-1/\alpha}.$$

Now, as in the proof of Lemma 1 (recall that $\text{h}^{(\alpha)}(z) := z \log^\alpha(1/z)$),

$$
\begin{aligned}
\mathbb{E}\text{H}^{(\alpha)}(\hat{\boldsymbol{\mu}}_n) \;=\; & \mathbb{E} \sum_{\substack{i \in [d] \\ \boldsymbol{\mu}(i) \vee \hat{\boldsymbol{\mu}}_n(i) \geq \mathrm{e}^{-\alpha}}} \text{h}^{(\alpha)}(\hat{\boldsymbol{\mu}}_n(i)) + \sum_{\substack{i \in [d] \\ \boldsymbol{\mu}(i) < \mathrm{e}^{-\alpha}}} \text{h}^{(\alpha)}(\hat{\boldsymbol{\mu}}_n(i))\mathbf{1}[\hat{\boldsymbol{\mu}}_n(i) < \mathrm{e}^{-\alpha}] \\
\leq \;& 2\mathrm{e}^\alpha \max_{z \in [\mathrm{e}^{-\alpha}, 1]} \text{h}^{(\alpha)}(z) + \sum_{\substack{i \in [d] \\ \boldsymbol{\mu}(i) < \mathrm{e}^{-\alpha}}} \mathbb{E}\text{h}^{(\alpha)}(\hat{\boldsymbol{\mu}}_n(i))\mathbf{1}[\hat{\boldsymbol{\mu}}_n(i) < \mathrm{e}^{-\alpha}] \\
= \;& 2\mathrm{e}^\alpha \max_{z \in [\mathrm{e}^{-\alpha}, 1]} \text{h}^{(\alpha)}(z) + \sum_{\substack{i \in [d] \\ \boldsymbol{\mu}(i) < \mathrm{e}^{-\alpha}}} \mathbb{E}\text{h}^{(\alpha)}(\hat{\boldsymbol{\mu}}_n(i)\mathbf{1}[\hat{\boldsymbol{\mu}}_n(i) < \mathrm{e}^{-\alpha}]) \\
\overset{(i)}{\leq} \;& 2\mathrm{e}^\alpha \max_{z \in [\mathrm{e}^{-\alpha}, 1]} \text{h}^{(\alpha)}(z) + \sum_{\substack{i \in [d] \\ \boldsymbol{\mu}(i) < \mathrm{e}^{-\alpha}}} \text{h}^{(\alpha)}(\mathbb{E}\hat{\boldsymbol{\mu}}_n(i)\mathbf{1}[\hat{\boldsymbol{\mu}}_n(i) < \mathrm{e}^{-\alpha}]) \\
\overset{(ii)}{\leq} \;& 2\mathrm{e}^\alpha \max_{z \in [\mathrm{e}^{-\alpha}, 1]} \text{h}^{(\alpha)}(z) + \sum_{\substack{i \in [d] \\ \boldsymbol{\mu}(i) < \mathrm{e}^{-\alpha}}} \text{h}^{(\alpha)}(\mathbb{E}\hat{\boldsymbol{\mu}}_n(i)) \\
\leq \;& 2\mathrm{e}^\alpha \max_{z \in [\mathrm{e}^{-\alpha}, 1]} \text{h}^{(\alpha)}(z) + \text{H}^{(\alpha)}(\boldsymbol{\mu}) \overset{(iii)}{\leq} 2\alpha^\alpha + \text{H}^{(\alpha)}(\boldsymbol{\mu}),
\end{aligned}
$$

where $(i)$ follows from Jensen's inequality, $(ii)$ is because $\text{h}^{(\alpha)}(z)$ is increasing on $z \in [0, \mathrm{e}^{-\alpha}]$ and $(iii)$ from (13).

By Berend and Kontorovich (2013, Lemma 6), we have $\mathbb{E}\|\hat{\boldsymbol{\mu}}_n - \boldsymbol{\mu}\|_1 \leq \Lambda_n(\boldsymbol{\mu})$, where

$$\Lambda_n(\boldsymbol{\mu}) := 2 \sum_{\boldsymbol{\mu}(j) < 1/n} \boldsymbol{\mu}(j) + \frac{1}{\sqrt{n}} \sum_{\boldsymbol{\mu}(j) \geq 1/n} \sqrt{\boldsymbol{\mu}(j)}.$$
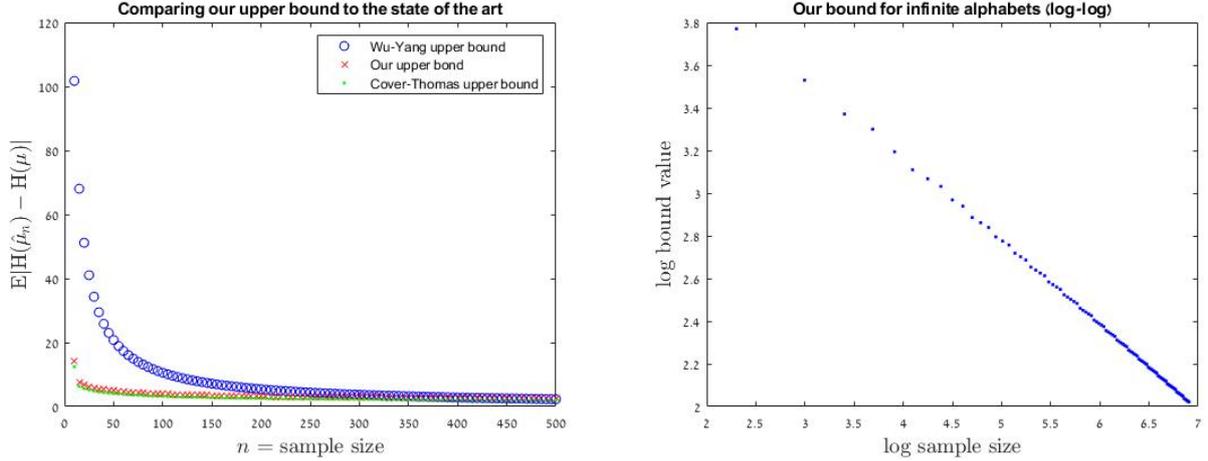
Figure 1: Left: A comparison of the three bounds for $d = 10$, $D = 1000$, $p = 0.95$. Our bound considerably outperforms Wu and Yang (2016) on small samples, and performs nearly as well as the finite-dimensional Cover-Thomas bound. Right: for our value of $q = 2$, the log-log plot shows roughly the correct slope of $-1/2$.

This quantity is always finite and $\Lambda_n(\boldsymbol{\mu}) \xrightarrow[n\to\infty]{} 0$ for all $\boldsymbol{\mu} \in \Delta_{\mathbb{N}}$ (ibid). Thus, we obtain the bound

$$\mathbb{E}|\mathrm{H}(\hat{\boldsymbol{\mu}}_n) - \mathrm{H}(\boldsymbol{\mu})| \leq \left(4\alpha^\alpha + 2\mathrm{H}^{(\alpha)}(\boldsymbol{\mu})\right)^{1/\alpha} \Lambda_n(\boldsymbol{\mu})^{1-1/\alpha}. \tag{21}$$

**Finite support**   For distributions with a large support but concentrated mass, the bound in (21) compares favorably to the state of the art, especially for smaller sample sizes. To illustrate this, consider a mixture of two distributions with support sizes $d$ and $D$: $\boldsymbol{\mu}'$ is uniform over $[d]$, $\boldsymbol{\mu}''$ is uniform over $[d + D]$, and $\boldsymbol{\mu} := p\boldsymbol{\mu}' + (1 - p)\boldsymbol{\mu}''$, for some $p \in [0, 1]$.

The state-of-the-art upper bound for the plug-in estimator can be inferred from Wu and Yang (2016, Appendix D), and has the form

$$\mathbb{E}|\mathrm{H}(\hat{\boldsymbol{\mu}}_n) - \mathrm{H}(\boldsymbol{\mu})| \leq \mathrm{WY}(d, D, p, n) := \frac{d + D}{n} + \min\left(C\frac{\log(d + D)}{\sqrt{n}}, \frac{\log n}{\sqrt{n}}\right)$$

for some $C > 1$; notice that it is insensitive to $p$. For a fair comparison to (21), our estimator's only a priori knowledge of $\boldsymbol{\mu}$ is that its support is of size at most $d + D$. By Proposition 2, we have $\max_{\boldsymbol{\mu} \in \Delta_K} \mathrm{H}^{(\alpha)}(\boldsymbol{\mu}) \leq \max\{\alpha, \log K\}^\alpha + (\alpha/\mathrm{e})^\alpha$. This allows us to optimize over $\alpha$ for each $n$:

$$\mathrm{OUR}(d, D, p, n) := \inf_{\alpha > 1} \left(4\alpha^\alpha + 2\max\{\alpha, \log(d + D)\}^\alpha + 2(\alpha/\mathrm{e})^\alpha\right)^{1/\alpha} \Lambda_n(\boldsymbol{\mu})^{1-1/\alpha}.$$

Since $\boldsymbol{\mu}$ has finite support, the Cover-Thomas inequality (2) also applies to yield an adaptive estimate when combined with (1). As $t\log(1/t)$ is concave, the latter has the form

$$\begin{aligned} \mathbb{E}|\mathrm{H}(\hat{\boldsymbol{\mu}}_n) - \mathrm{H}(\boldsymbol{\mu})| &\leq \mathbb{E}\left[\|\hat{\boldsymbol{\mu}}_n - \boldsymbol{\mu}\|_1 \log\frac{d + D}{\|\hat{\boldsymbol{\mu}}_n - \boldsymbol{\mu}\|_1}\right] \\ &\leq \Lambda_n(\boldsymbol{\mu})\log\frac{d + D}{\Lambda_n(\boldsymbol{\mu})} =: \mathrm{CT}(d, D, p, n). \end{aligned}$$

The comparisons are plotted in Figure 1 (Left).

**Infinite support**   In some cases our bound is nearly tight (at least for the plug-in estimate), such as for the family of zeta distributions $\boldsymbol{\mu}_q(i) \sim 1/i^q$ with parameter $q > 1$. For this family, Antos and Kontoyiannis (2001a, Theorem 7) establish a lower bound of order $n^{\frac{1-q}{q}}$ on $\mathbb{E}|\mathrm{H}(\hat{\boldsymbol{\mu}}_n) - \mathrm{H}(\boldsymbol{\mu}_q)|$. It is straightforward to verify[6] that $\boldsymbol{\mu}_q \in \Delta_{\mathbb{N}}^{(\alpha)}$ for all $q, \alpha > 1$. Thus, we can optimize our bound in (21)

---

[6]One can, for example, apply Cauchy's condensation test, followed up by the ratio test.

over all $\alpha > 1$; the results are presented in Figure 1 (Right).

# 7  Moments of Information

We motivate our bounded-moment assumption as being *considerably* less restrictive than the finite-alphabet assumptions and tail conditions studied in previous works (Wu and Yang, 2016; Silva, 2018); see Section .1 for a detailed comparison. Obtaining moment-based results is essentially a desideratum laid out by Antos and Kontoyiannis (2001a), in which it is hypothesized that — in parallel to the asymptotic distribution for the finite alphabet case — moment conditions are the correct notion to achieve finite-sample estimates in the infinite alphabet case. Our Theorem 5(a) shows that, under these assumptions, there is an inverse logarithmic convergence rate (similar to, though distinct from, the results of Wyner and Foster (2003)) and, furthermore, using empirical quantities, this rate can be very much accelerated, as demonstrated in Corollary 1.

In this section, we study some of the mathematical properties of moments of information.

## 7.1  Maximizing the $\alpha$-moment over a fixed support size

**Proposition 2.** *For $K \geq 2$ and $\alpha \geq 1$,*

$$\max\{\log K, (\alpha/\mathrm{e})\}^\alpha \ \leq \ \max_{\boldsymbol{\mu} \in \Delta_K} \mathrm{H}^{(\alpha)}(\boldsymbol{\mu}) \ \leq \ \max\{\log K, \alpha\}^\alpha + (\alpha/\mathrm{e})^\alpha.$$

We will need the following useful (and likely known) result.

**Lemma 6** (folklore). *Suppose that $0 < a < 1$ and $f : [0,1] \to \mathbb{R}$ is strictly concave on $[0,a]$ and strictly convex on $[a,1]$. Define the function $F : \Delta_K \to \mathbb{R}$ by*

$$F(\boldsymbol{\mu}) = \sum_{i=1}^{K} f(\boldsymbol{\mu}(i)).$$

*Then any maximizer $\boldsymbol{\mu}^\star$ of $F$ is either the uniform distribution or else has exactly 1 "heavy" mass $v \in [a,1]$ and $K-1$ identical "light" masses $(1-v)/(K-1)$.*

*Proof.* A standard "smoothing" argument (Loh, 2013) shows that if two masses $u \leq v$ occur in the interval $(a,1)$, there is an $\varepsilon > 0$ such that $f(u-\varepsilon) + f(v+\varepsilon) > f(u) + f(v)$. In other words, such masses can be pushed apart (keeping their sum fixed) to increase the value of $F$, until one of them reaches the boundary of $[a,1]$. Furthermore, since $0 < a < u < v$ and $u + v \leq 1$, repeated iteration of the "pushing apart" operation will hit the left endpoint (i.e., $a$) rather than the right one (i.e., 1). Having exhausted the "pushing apart" process, we are left with one "heavy" mass $v \in [a,1]$ and $K-1$ "lighter" ones in $[0,a]$. But concavity implies that $F$ will be maximized by pulling the lighter masses in (as opposed to pushing them apart), which amounts to replacing each of them by the average of the $K-1$ values. $\qquad\square$

*Proof of Proposition 2.* Choosing $\boldsymbol{\mu}$ to be the uniform distribution yields $\mathrm{H}^{(\alpha)}(\boldsymbol{\mu}) = \log^\alpha K$, and choosing $\boldsymbol{\mu}$ such that $v := \boldsymbol{\mu}(1) = \mathrm{e}^{-\alpha}$ yields $\mathrm{H}^{(\alpha)}(\boldsymbol{\mu}) \geq v \log(1/v)^\alpha = (\alpha/\mathrm{e})^\alpha$. Thus, the lower bound is proven and it only remains to prove the upper bound.

Let $\boldsymbol{\mu}^\star$ be a maximizer for given $\alpha, K$. Recall the function $\mathrm{h}^{(\alpha)}(z) = z \log^\alpha(1/z)$ and note that it is strictly concave on $[0, \mathrm{e}^{-(\alpha-1)}]$ and strictly convex on $[\mathrm{e}^{-(\alpha-1)}, 1]$. Then Lemma 6 shows that $\boldsymbol{\mu}^\star$ will either be uniform or else attains at most one value $v \in [\mathrm{e}^{-(\alpha-1)}, 1]$ in the convex interval, with the remaining values equal to $\frac{1-v}{K-1} \in [0, \mathrm{e}^{-(\alpha-1)}]$ in the concave interval. Only the latter case is non-trivial:

$$\mathrm{H}^{(\alpha)}(\boldsymbol{\mu}^\star) \ = \ v \left( \log \frac{1}{v} \right)^\alpha + (1-v) \left( \log \frac{K-1}{1-v} \right)^\alpha$$

for some $v$ satisfying

$$0 < \frac{1-v}{K-1} \leq \mathrm{e}^{-(\alpha-1)} \leq v < 1. \tag{22}$$

14

Now $v \left( \log \frac{1}{v} \right)^{\alpha}$ is maximized over $[0,1]$ by $v = e^{-\alpha}$, which yields the value $(\alpha/e)^{\alpha}$.

To bound the second term, $g(v) := (1-v) \left( \log \frac{K-1}{1-v} \right)^{\alpha}$, we consider two cases: (i) $K - 1 < e^{\alpha}$ and (ii) $K - 1 \geq e^{\alpha}$. In case (i), $g$ is maximized by $v^{\star} = 1 - (K-1)/e^{\alpha}$ and

$$g(v^{\star}) = (1 - v^{\star}) \left( \log \frac{K-1}{1-v^{\star}} \right)^{\alpha} \leq \left( \log \frac{K-1}{1-v^{\star}} \right)^{\alpha} = \alpha^{\alpha}.$$

In case (ii), $g$ is monotonically decreasing in $v$. The constraint $\frac{1-v}{K-1} \leq v$ from (22) implies $v \geq 1/K$, so in this case,

$$g(v) \leq \left( \log \frac{K-1}{1 - 1/K} \right)^{\alpha} = \log^{\alpha} K.$$

This proves the upper bound. $\qquad\square$

## 7.2  Moments of Information vs. Moments of Distributions

Since for all $r \geq 1$ and $\boldsymbol{\mu} \in \Delta_{\mathbb{N}}$, we trivially have $\|\boldsymbol{\mu}\|_r \leq 1$, it is only for $r < 1$ that $\|\boldsymbol{\mu}\|_r$ conveys nontrivial tail information. However, as a measure of tail decay, the latter is rather crude: $\|\boldsymbol{\mu}\|_r < \infty$ for *any* $r < 1$ implies $H^{(\alpha)}(\boldsymbol{\mu}) < \infty$ for *all* $\alpha > 0$:

**Proposition 3.** *For all $\alpha > 0$ and $r \in (0,1)$, we have*

$$H^{(\alpha)}(\boldsymbol{\mu}) \;\; \leq \;\; \left( \frac{\alpha}{e(1-r)} \right)^{2\alpha} \|\boldsymbol{\mu}\|_r \,.$$

Remark. The bound above is quite loose. For example, for $\alpha = 1$ the AM-GM inequality readily yields, for any $r \in (0,1)$,

$$H(\boldsymbol{\mu}) \leq \frac{r}{1-r} \ln \|\boldsymbol{\mu}\|_r \,.$$

It may be of interest to investigate bounds of the form

$$H^{(\alpha)}(\boldsymbol{\mu}) \leq a(\alpha) \log^{b(\alpha)} \|\boldsymbol{\mu}\|_{c(\alpha)} \,,$$

for some functions $a, b, c : \alpha \mapsto (0, \infty)$.

*Proof.* We first claim that for $\alpha > 0$ and $r \in (0,1)$,

$$x \log^{\alpha} \frac{1}{x} < x^r, \qquad x \in [0,1].$$

Indeed, the function $x \mapsto x^{1-r} \log^{\alpha}(\frac{1}{x})$ is maximized at $x = e^{\frac{\alpha}{r-1}}$, attaining the maximum value of $e^{-\alpha} (\frac{\alpha}{1-r})^{\alpha}$. The latter is less than 1 whenever $\alpha < e(1-r)$. Likewise, whenever $\alpha < ce(1-r)$, we have $x \log^{\alpha}(\frac{1}{x}) < c^{\alpha} x^r$, and so $H^{(\alpha)}(\boldsymbol{\mu}) < c^{2\alpha} \|\boldsymbol{\mu}\|_r$. Choosing $c = \frac{\alpha}{e(1-r)}$ proves the claim. $\qquad\square$

## 7.3  Resolution of Jürgensen and Matthews Conjectures

In this section, we give a complete resolution of the conjectures posed by Jürgensen and Matthews (2010).

**Conjecture 10.1** Jürgensen and Matthews (2010, Conjecture 10.1) posits that for $d = 2$, $\max_{\boldsymbol{\mu} \in \Delta_2} \mathrm{H}^{(\alpha)}(\boldsymbol{\mu})$ has two maximizers $\pi_1^{(\alpha)} = \left(\frac{1}{2} + x^{(\alpha)}, \frac{1}{2} - x^{(\alpha)}\right)$ and $\pi_2^{(\alpha)} = \left(\frac{1}{2} - x^{(\alpha)}, \frac{1}{2} + x^{(\alpha)}\right)$ for some value $x^{(\alpha)}$ such that $x^{(2)} = \frac{1}{2e}\sqrt{e^2 - 4}$ and $x^{(\alpha)}$ is strictly increasing as $\alpha \to \infty$ and $\lim_{\alpha \to \infty} x^{(\alpha)} = \frac{1}{2}$.

By Lemma 6, there are at most three maximizers. Since $\mathrm{H}^{(\alpha)}\left((e^{-\alpha}, 1 - e^{-\alpha})\right) > \left(\frac{\alpha}{e}\right)^\alpha > \log^\alpha(2)$, the uniform distribution is not a maximizer. So, including permutations, there are exactly two maximizers.

Let $(u_\alpha^\star, v_\alpha^\star)$ be the increasingly-ordered maximizing distribution. We cannot have $u_\alpha^\star < e^{-\alpha}$, because this would only decrease $\mathrm{H}^{(\alpha)}$ as compared to $\mathrm{H}^{(\alpha)}\left((e^{-\alpha}, 1 - e^{-\alpha})\right)$. By Lemma 6, $u_\alpha^\star \leq e^{-(\alpha-1)}$, and similarly $v_\alpha^\star \in [1 - e^{-(\alpha-1)}, 1 - e^{-\alpha}]$. By convexity and monotonicity of $\log \frac{1}{x}$ on $[0,1]$, the difference between $\left|\mathrm{h}^{(\alpha+1)}(e^{-\alpha}) - \mathrm{h}^{(\alpha+1)}(u_\alpha^\star)\right|$ and $\left|\mathrm{h}^{(\alpha)}(e^{-\alpha}) - \mathrm{h}^{(\alpha)}(u_\alpha^\star)\right|$ shrinks by more than the difference between $\left|\mathrm{h}^{(\alpha+1)}(v_\alpha^\star) - \mathrm{h}^{(\alpha+1)}(1 - e^{-\alpha})\right|$ and $\left|\mathrm{h}^{(\alpha)}(v_\alpha^\star) - \mathrm{h}^{(\alpha)}(1 - e^{-\alpha})\right|$ grows. So, for $\max_{\boldsymbol{\mu} \in \Delta_2} \mathrm{H}^{(\alpha)}(\boldsymbol{\mu})$ to be less than $\max_{\boldsymbol{\mu} \in \Delta_2} \mathrm{H}^{(\alpha+1)}(\boldsymbol{\mu})$, as occurs (for sufficiently large $\alpha$) by resolution of Conjecture 10.4 below, it must be that $u_{\alpha+1}^\star < u_\alpha^\star$ and $v_{\alpha+1}^\star > v_\alpha^\star$ (as $\alpha$ tends to infinity).

Furthermore, $e^{-(\alpha-1)} \to 0$ as $\alpha \to \infty$, and so $\lim_{\alpha \to \infty} x^{(\alpha)} = \frac{1}{2}$.

To find the value of $x^{(2)}$, set $x := x^{(2)}$ and find the critical points of

$$\mathrm{H}^{(2)}\left(x + \frac{1}{2}, x - \frac{1}{2}\right) := \left(\frac{1}{2} - x\right)\log^2\left(\frac{1}{2} - x\right) + \left(x + \frac{1}{2}\right)\log^2\left(x + \frac{1}{2}\right).$$

Differentiating and factoring, we get

$$\frac{\mathrm{d}}{\mathrm{d}x}\mathrm{H}^{(2)}\left(x + \frac{1}{2}, x - \frac{1}{2}\right) = -\left(\log\left(-x + \frac{1}{2}\right) - \log\left(x + \frac{1}{2}\right)\right)\left(2 + \log\left(-x + \frac{1}{2}\right) + \log\left(x + \frac{1}{2}\right)\right) = 0.$$

Now $x = 0$ is a solution which we know is not the maximum and we also get $x = \pm\frac{\sqrt{e^2-4}}{2e}$ which exactly what Jürgensen and Matthews (2010) conjectured.

**Conjecture 10.2** Jürgensen and Matthews (2010, Conjecture 10.2) posits that for $\pi_1^{(\alpha)}, \pi_2^{(\alpha)}$ as above and $\alpha \geq 2$, we have $\mathrm{H}^{(\alpha)}(\pi_1^{(\alpha)}) = \mathrm{H}^{(\alpha)}(\pi_2^{(\alpha)}) > (\log 2)^\alpha$ and moreover, this quantity is strictly increasing and unbounded as $\alpha \to \infty$.

Since $\log(2) < \frac{\alpha}{e}$, by Proposition 2, $\mathrm{H}^{(\alpha)}(\pi_1^{(\alpha)}) = \mathrm{H}^{(\alpha)}(\pi_2^{(\alpha)}) > (\log 2)^\alpha$ and unbounded.

**Conjecture 10.3** Jürgensen and Matthews (2010, Conjecture 10.3) posits that $\mathrm{H}^{(\alpha)}$ has $d$ local maxima for for $d > 2$ and $\alpha > 2$. By Lemma 6, the only maxima are the uniform distribution and the $d$ permutations of $\sup_{v \in [e^{-(\alpha-1)}, 1):(d-1)u+v=1}(d-1)u\log^\alpha(u) + v\log^\alpha(v)$, should the latter exist with $v$ in interior of interval. So there are either 1 (e.g. $e^\alpha = d$) or $d + 1$ local maxima.

**Conjecture 10.4** For $d, \alpha \in \mathbb{N}$, define $h_{d,\alpha}^\star := \max_{\boldsymbol{\mu} \in \Delta_d} \mathrm{H}^{(\alpha)}(\boldsymbol{\mu})$. Jürgensen and Matthews (2010, Conjecture 10.4) posits that $h_{d,\alpha+1}^\star > h_{d,\alpha}^\star$ and that $\lim_{\alpha \to \infty} h_{d,\alpha}^\star = \infty$. In light of the lower bound in Proposition 2, the latter claim (i.e., unboundedness) is immediate.

For $d > e^\alpha$, by Proposition 2, $\max_{\boldsymbol{\mu} \in \Delta_d} \mathrm{H}^{(\alpha)}(\boldsymbol{\mu}) \leq 2\log^\alpha d$ and $\log^{\alpha+1} d \leq \max_{\boldsymbol{\mu} \in \Delta_d} \mathrm{H}^{(\alpha+1)}(\boldsymbol{\mu})$. We find, therefore, that for $d > e^2$, $\max_{\boldsymbol{\mu} \in \Delta_d} \mathrm{H}^{(\alpha)}(\boldsymbol{\mu}) \leq \max_{\boldsymbol{\mu} \in \Delta_d} \mathrm{H}^{(\alpha+1)}(\boldsymbol{\mu})$.

But since the conjecture takes interest in the case of $\alpha$ tending to infinity, let us focus on $e^\alpha \geq d - 1$.

By Lemma 6, $\boldsymbol{\mu}_\alpha^\star := \arg\max_{\boldsymbol{\mu} \in \Delta_d} \mathrm{H}^{(\alpha)}(\boldsymbol{\mu})$, is either uniform or takes two distinct values $v \in [e^{-(\alpha-1)}, 1]$ and $u = \frac{1-v}{d-1} \in [0, e^{-(\alpha-1)}]$.

For $x \in [0, \frac{1}{e}]$, $\log(1/x) \geq 1$, so $\mathrm{h}^{(\alpha+1)}(x) \geq \mathrm{h}^{(\alpha)}(x)$. So if $u, v \in [0, \frac{1}{e}]$, then $\mathrm{H}^{(\alpha)}(\boldsymbol{\mu}_\alpha^\star) < \mathrm{H}^{(\alpha+1)}(\boldsymbol{\mu}_\alpha^\star) \leq \max_{\boldsymbol{\mu} \in \Delta_d} \mathrm{H}^{(\alpha+1)}(\boldsymbol{\mu})$

So assume instead that $v \in (\frac{1}{e}, 1]$. In this case, we can bound the difference $\mathrm{h}^{(\alpha)}(v) - \mathrm{h}^{(\alpha+1)}(v) \leq \mathrm{h}^{(\alpha)}(v) \leq \frac{1}{e}$.

Since $e^\alpha \geq d-1$, $u > e^{-\alpha}$ and must lie in $[e^{-\alpha}, e^{-(\alpha-1)}]$. But for this entire interval, $x \in [e^{-\alpha}, e^{-(\alpha-1)}]$ has $\mathrm{h}^{(\alpha+1)}(x) - \mathrm{h}^{(\alpha)}(x) \geq \frac{1}{e} \geq \mathrm{h}^{(\alpha)}(v)$. In order to see this, it suffices, since $\mathrm{h}^{(\alpha)}$ and $\mathrm{h}^{(\alpha+1)}$ are both decreasing on $[e^{-\alpha}, e^{-(\alpha-1)}]$, to show that $\mathrm{h}^{(\alpha+1)}(e^{-(\alpha-1)}) - \mathrm{h}^{(\alpha)}(e^{-\alpha}) > \frac{1}{e}$. This can be seen by observing $e(\alpha-1)^{\alpha+1} - \alpha^\alpha > e^{\alpha-1}$ for all $\alpha \geq 3$.

It follows, therefore, that for $\alpha \in \mathbb{N}$, when $d > e^{\alpha}$, for all $d \geq 8$, or for $e^{\alpha} \geq d - 1$, for all $\alpha \geq 3$, the former claim (monotonicity) holds, i.e.

$$\max_{\boldsymbol{\mu} \in \Delta_d} H^{(\alpha)}(\boldsymbol{\mu}) < \max_{\boldsymbol{\mu} \in \Delta_d} H^{(\alpha+1)}(\boldsymbol{\mu}).$$

(This can also be generalized to any $\beta(\neq \alpha + 1)$ for sufficiently large $d$ or $\alpha$ respectively, if one so desired).

**Our conjecture**   We close the section with a conjecture of our own.

**Conjecture 1.** *For $e^{\alpha} < d < \infty$, we have $\max_{\boldsymbol{\mu} \in \Delta_d} H_{\alpha}(\boldsymbol{\mu}) = \log^{\alpha} d$ and moreover, the maximum is achieved by the uniform distribution over $[d]$.*

# Acknowledgments

# References

Jayadev Acharya, Hirakendu Das, Alon Orlitsky, and Ananda Theertha Suresh. A unified maximum likelihood approach for estimating symmetric properties of discrete distributions. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 11–21. PMLR, 06–11 Aug 2017. URL http://proceedings.mlr.press/v70/acharya17a.html.

Jayadev Acharya, Sourbh Bhadane, Piotr Indyk, and Ziteng Sun. Estimating entropy of distributions in constant space. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2019. Curran Associates Inc.

András Antos and Ioannis Kontoyiannis. Convergence properties of functional estimates for discrete distributions. *Random Structures & Algorithms*, 19(3-4):163–193, 2001a.

András Antos and Ioannis Kontoyiannis. Estimating the entropy of discrete distributions. In *IEEE International Symposium on Information Theory*, pages 45–45, 2001b.

Koenraad M. R. Audenaert. A sharp continuity estimate for the von Neumann entropy. *J. Phys. A*, 40(28):8127–8136, 2007. ISSN 1751-8113. doi: 10.1088/1751-8113/40/28/S18. URL https://doi.org/10.1088/1751-8113/40/28/S18.

Daniel Berend and Aryeh Kontorovich. A sharp estimate of the binomial mean absolute deviation with applications. *Statistics & Probability Letters*, 83(4):1254–1259, 2013. ISSN 0167-7152. doi: https://doi.org/10.1016/j.spl.2013.01.023. URL https://www.sciencedirect.com/science/article/pii/S0167715213000242.

Daniel Berend, Aryeh Kontorovich, and Gil Zagdanski. The expected missing mass under an entropy constraint. *Entropy*, 19(7):315, 2017. doi: 10.3390/e19070315. URL https://doi.org/10.3390/e19070315.

Mickey Brautbar and Alex Samorodnitsky. Approximating entropy from sublinear samples. In Nikhil Bansal, Kirk Pruhs, and Clifford Stein, editors, *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2007, New Orleans, Louisiana, USA, January 7-9, 2007*, pages 366–375. SIAM, 2007. URL http://dl.acm.org/citation.cfm?id=1283383.1283422.

David Brink. A (probably) exact solution to the birthday problem. *The Ramanujan Journal*, 28, 06 2012. doi: 10.1007/s11139-011-9343-9.

Doron Cohen, Aryeh Kontorovich, and Geoffrey Wolfer. Learning discrete distributions with infinite support. In *Neural Information Processing Systems (NeurIPS)*, 2020.

Thomas M. Cover and Joy A. Thomas. *Elements of information theory*. Wiley-Interscience [John Wiley & Sons], Hoboken, NJ, second edition, 2006. ISBN 978-0-471-24195-9; 0-471-24195-4.

Kazuto Fukuchi and Jun Sakuma. Minimax optimal estimators for additive scalar functionals of discrete distributions. In *2017 IEEE International Symposium on Information Theory (ISIT)*, pages 2103–2107, 2017. doi: 10.1109/ISIT.2017.8006900.

Kazuto Fukuchi and Jun Sakuma. Minimax optimal additive functional estimation with discrete distribution: Slow divergence speed case. In *2018 IEEE International Symposium on Information Theory (ISIT)*, pages 1041–1045, 2018. doi: 10.1109/ISIT.2018.8437725.

S. Golomb. The information generating function of a probability distribution (corresp.). *IEEE Transactions on Information Theory*, 12(1):75–77, 1966. doi: 10.1109/TIT.1966.1053843.

Yanjun Han, Jiantao Jiao, and Tsachy Weissman. Adaptive estimation of shannon entropy. In *2015 IEEE International Symposium on Information Theory (ISIT)*, pages 1372–1376. IEEE, 2015.

Yi Hao and Alon Orlitsky. Data amplification: Instance-optimal property estimation. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 4049–4059. PMLR, 13–18 Jul 2020. URL http://proceedings.mlr.press/v119/hao20a.html.

Yi Hao, Alon Orlitsky, Ananda Theertha Suresh, and Yihong Wu. Data amplification: A unified and competitive approach to property estimation. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper/2018/file/a753a43564c29148df3150afb4475440-Paper.pdf.

Siu-Wai Ho and Raymond W Yeung. The interplay between entropy and variational distance. *IEEE Transactions on Information Theory*, 56(12):5906–5929, 2010.

J. Jiao, K. Venkat, Y. Han, and T. Weissman. Maximum likelihood estimation of functionals of discrete distributions. *IEEE Transactions on Information Theory*, 63(10):6774–6798, 2017. doi: 10.1109/TIT.2017.2733537.

Jiantao Jiao, Kartik Venkat, Yanjun Han, and Tsachy Weissman. Minimax estimation of functionals of discrete distributions. *IEEE Transactions on Information Theory*, 61(5):2835–2885, 2015.

Helmut Jürgensen and David E Matthews. Entropy and higher moments of information. *Journal of Universal Computer Science*, 16(5):749–794, 2010.

Norbert Kusolitsch. Why the theorem of Scheffé should be rather called a theorem of Riesz. *Period. Math. Hungar.*, 61(1-2):225–229, 2010. ISSN 0031-5303. doi: 10.1007/s10998-010-3225-6. URL https://doi.org/10.1007/s10998-010-3225-6.

Elliott H. Lieb and Michael Loss. *Analysis*, volume 14 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, second edition, 2001. ISBN 0-8218-2783-9. doi: 10.1090/gsm/014. URL https://doi.org/10.1090/gsm/014.

Po-Shen Loh. Convexity. 2013. URL https://www.math.cmu.edu/~ploh/docs/math/mop2008/convexity-soln.pdf.

Gábor Lugosi and Shahar Mendelson. Mean estimation and regression under heavy-tailed distributions: A survey. *Found. Comput. Math.*, 19(5):1145–1190, 2019. doi: 10.1007/s10208-019-09427-x. URL https://doi.org/10.1007/s10208-019-09427-x.

P. Massart. The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality. *Ann. Probab.*, 18(3):1269–1283, 1990. ISSN 0091-1798. URL http://links.jstor.org/sici?sici=0091-1798(199007)18:3<1269:TTCITD>2.0.CO;2-Q&origin=MSN.

L. Paninski. Estimating entropy on $m$ bins given fewer than $m$ samples. *IEEE Transactions on Information Theory*, 50(9):2200–2203, 2004. doi: 10.1109/TIT.2004.833360.

Liam Paninski. Estimation of entropy and mutual information. *Neural computation*, 15(6):1191–1253, 2003. doi: 10.1162/089976603321780272.

Iosif Pinelis. Best possible bounds of the von Bahr–Esseen type. *Annals of Functional Analysis*, 6(4):1 – 29, 2015. doi: 10.15352/afa/06-4-1. URL https://doi.org/10.15352/afa/06-4-1.

Igal Sason. Entropy bounds for discrete random variables via maximal coupling. *IEEE Trans. Inf. Theory*, 59(11):7118–7131, 2013. doi: 10.1109/TIT.2013.2274515. URL https://doi.org/10.1109/TIT.2013.2274515.

Henry Scheffé. A useful convergence theorem for probability distributions. *Ann. Math. Statistics*, 18:434–438, 1947. ISSN 0003-4851. doi: 10.1214/aoms/1177730390. URL https://doi.org/10.1214/aoms/1177730390.

Jorge F Silva. Shannon entropy estimation in $\infty$-alphabets from convergence results: studying plug-in estimators. *Entropy*, 20(6):397, 2018.

Alexandre B Tsybakov. *Introduction to nonparametric estimation*. Springer Science & Business Media, 2008.

Gregory Valiant and Paul Valiant. Estimating the unseen: An n/log(n)-sample estimator for entropy and support size, shown optimal via new clts. In *Proceedings of the Forty-Third Annual ACM Symposium on Theory of Computing*, STOC '11, New York, NY, USA, 2011a. Association for Computing Machinery. ISBN 9781450306911. doi: 10.1145/1993636.1993727. URL https://doi.org/10.1145/1993636.1993727.

Gregory Valiant and Paul Valiant. The power of linear estimators. In *2011 IEEE 52nd Annual Symposium on Foundations of Computer Science*, pages 403–412, 2011b. doi: 10.1109/FOCS.2011.81.

Gregory Valiant and Paul Valiant. Estimating the unseen: Improved estimators for entropy and other properties. *J. ACM*, 64(6), October 2017. ISSN 0004-5411. doi: 10.1145/3125643. URL https://doi.org/10.1145/3125643.

Sergio Verdú. Empirical estimation of information measures: A literature guide. *Entropy*, 21(8):720, 2019.

Bengt von Bahr and Carl-Gustav Esseen. Inequalities for the $r$th Absolute Moment of a Sum of Random Variables, $1 \leqq r \leqq 2$. *The Annals of Mathematical Statistics*, 36(1):299 – 303, 1965. doi: 10.1214/aoms/1177700291. URL https://doi.org/10.1214/aoms/1177700291.

Yihong Wu and Pengkun Yang. Minimax rates of entropy estimation on large alphabets via best polynomial approximation. *IEEE Transactions on Information Theory*, 62(6):3702–3720, 2016.

Abraham J Wyner and Dean Foster. On the lower limits of entropy estimation. *IEEE Transactions on Information Theory, submitted for publication*, 2003.

Zhengmin Zhang. Estimating mutual information via kolmogorov distance. *IEEE Transactions on Information Theory*, 53(9):3280–3282, 2007. doi: 10.1109/TIT.2007.903122.

### .1 Comparison of tail-vs-moment assumptions

Expanding upon the observation of Antos and Kontoyiannis (2001b) that moment of information assumptions are "weaker (and somewhat more natural" than tail decay rates, we make some concrete comparisons between the two.

Let us list a number of conditions one might impose $\boldsymbol{\mu}$:

$A_1(\alpha)$: Finite $\alpha$-moment of information Antos and Kontoyiannis (2001b):
For some $\alpha > 1, \mathbb{E}\left[\log^\alpha \frac{1}{\boldsymbol{\mu}(X)}\right] < \infty$.

$B_1(\alpha, M_\alpha)$: Bounded $\alpha$-moment of information:
For some $\alpha > 1, \exists M_\alpha > 0, \mathbb{E}\left[\log^\alpha \frac{1}{\boldsymbol{\mu}(X)}\right] < M_\alpha$.

$A_2(\beta)$: Superlinear $\beta$ tail decay:
For some $\beta > 1, \boldsymbol{\mu}(i) = \mathcal{O}\left(\frac{1}{i^\beta}\right)$.

$B_2(\beta, \underline{c}_\beta, \overline{c}_\beta)$: Controlled superlinear $\beta$ tail decay (Antos and Kontoyiannis, 2001b):
For some $\beta > 1, \exists \underline{c}_\beta, \overline{c}_\beta > 0$ such that $\frac{\underline{c}_\beta}{i^\beta} \leq \boldsymbol{\mu}(i) \leq \frac{\overline{c}_\beta}{i^\beta}$.

$A_3(\gamma)$: Superlinearly $\gamma$-logarithmic tail decay:
For some $\gamma > 1, \boldsymbol{\mu}(i) = \mathcal{O}\left(\frac{1}{i \log^\gamma i}\right)$.

$B_3(\gamma, \underline{c}_\gamma, \overline{c}_\gamma)$: Controlled superlinearly $\gamma$-logarithmic tail decay (Antos and Kontoyiannis, 2001b):
For some $\gamma > 1, \exists \underline{c}_\gamma, \overline{c}_\gamma > 0$ such that $\frac{\underline{c}_\gamma}{i \log i^\gamma} \leq \boldsymbol{\mu}(i) \leq \frac{\overline{c}_\gamma}{i \log i^\gamma}$.

**Proposition 4.** *The following implications hold:*

(a) $A_3(\gamma), \gamma > 2 \implies A_1(\alpha), \forall \alpha < \gamma - 1$.

(b) $A_1(\alpha), \alpha > 1 \implies\!\!\!\!\!/ \;\; A_3(\gamma), \gamma < \alpha + 1$

(c) $B_1(\alpha, M_\alpha) \implies A_1(\alpha)$

(d) $A_2(\beta) \implies A_3(\gamma), \forall \gamma > 1$

(e) $A_2(\beta) \implies A_1(\alpha), \forall \alpha > 1$

(f) $B_2(\beta, \underline{c}_\beta, \overline{c}_\beta) \implies A_2(\beta)$

(g) $B_3(\gamma, \underline{c}_\gamma, \overline{c}_\gamma) \implies A_3(\gamma)$

(h) $B_2(\beta, \underline{c}_\beta, \overline{c}_\beta) \implies B_1(\alpha, M_\alpha)$ *with* $M_\alpha = M_\alpha(\alpha, \overline{c}_\beta, \beta)$

(i) $B_3(\gamma, \underline{c}_\gamma, \overline{c}_\gamma), \gamma > 2 \implies B_1(\alpha, M_\alpha), \forall \alpha < \gamma - 1$, *with* $M_\alpha = M_\alpha(\alpha, \overline{c}_\gamma, \gamma)$.

**Remark** We start by noticing that

$$\sum_{\substack{i \in \Omega \\ \boldsymbol{\mu}(i) \leq e^{-\alpha}}} \boldsymbol{\mu}(i) \log^\alpha \frac{1}{\boldsymbol{\mu}(i)} \leq \sum_{i \in \Omega} \boldsymbol{\mu}(i) \log^\alpha \frac{1}{\boldsymbol{\mu}(i)} \leq \alpha^\alpha + \sum_{\substack{i \in \Omega \\ \boldsymbol{\mu}(i) \leq e^{-\alpha}}} \boldsymbol{\mu}(i) \log^\alpha \frac{1}{\boldsymbol{\mu}(i)}, \tag{23}$$

and that on $[0, e^{-\alpha}]$, it holds that $x \to x \log^\alpha \frac{1}{x}$ is increasing. Since $\alpha^\alpha$ is finite, the convergence of the series is primarily governed by what happens or small probabilities.

**Proof** $(c), (d), (f), (g)$ Immediate.

**Proof** $(a)$   Suppose that assumption $A_3(\gamma)$ holds. Then $\exists N \in \mathbb{N}, C > 0$ such that for any $i \geq N$, $\boldsymbol{\mu}(i) \leq C\frac{1}{i \log^\gamma i}$. We focus on the rightmost term of (23).

$$
\sum_{\substack{i \in \Omega \\ \boldsymbol{\mu}(i) \leq \mathrm{e}^{-\alpha}}} \boldsymbol{\mu}(i) \log^\alpha \frac{1}{\boldsymbol{\mu}(i)} \leq N\mathrm{e}^{-\alpha}\alpha^\alpha + \sum_{\substack{i \in \Omega \\ \boldsymbol{\mu}(i) \leq \mathrm{e}^{-\alpha} \\ i > N}} \boldsymbol{\mu}(i) \log^\alpha \frac{1}{\boldsymbol{\mu}(i)}
$$

$$
\leq N\mathrm{e}^{-\alpha}\alpha^\alpha + C \sum_{\substack{i \in \Omega \\ \boldsymbol{\mu}(i) \leq \mathrm{e}^{-\alpha} \\ i > N}} \frac{1}{i \log^\gamma i} \log^\alpha \frac{i \log^\gamma i}{C}
$$

$$
\leq N\mathrm{e}^{-\alpha}\alpha^\alpha + C \sum_{\substack{i \in \Omega \\ \boldsymbol{\mu}(i) \leq \mathrm{e}^{-\alpha} \\ i > N}} \frac{(\log i + \gamma \log \log i + \log 1/C)^\alpha}{i \log^\gamma i}.
$$

Since $\log i$ dominates both $\log \log i$ and $\log 1/C$, the series converges whenever $\sum_{i \in \mathbb{N}} \frac{1}{i \log^{\gamma - \alpha} i}$ converges, which occurs exactly for $\gamma > \alpha + 1$.

**Proof** $(e)$   Let $\boldsymbol{\mu}$, and suppose that $A_2(\beta)$ holds for some $\beta$. Then $\exists N \in \mathbb{N}, C > 0$ such that $\forall i \in \mathbb{N}, i > N, \boldsymbol{\mu}(i) \leq C\frac{1}{i^\beta}$. Let $\alpha > 1$. We decompose the expression of $\mathbb{E}\left[\log^\alpha \frac{1}{\boldsymbol{\mu}(X)}\right]$:

$$
\sum_{i \in \Omega} \boldsymbol{\mu}(i) \log^\alpha \frac{1}{\boldsymbol{\mu}(i)} = \sum_{\substack{i \in \Omega \\ i \leq N}} \boldsymbol{\mu}(i) \log^\alpha \frac{1}{\boldsymbol{\mu}(i)} + \underbrace{\sum_{\substack{i \in \Omega \\ i > N \\ \boldsymbol{\mu}(i) \leq \mathrm{e}^{-\alpha}}} \boldsymbol{\mu}(i) \log^\alpha \frac{1}{\boldsymbol{\mu}(i)}}_{x \mapsto x \log^\alpha \frac{1}{x} \text{ increasing on } [0, \mathrm{e}^{-\alpha}]} + \sum_{\substack{i \in \Omega \\ i > N \\ \boldsymbol{\mu}(i) > \mathrm{e}^{-\alpha}}} \boldsymbol{\mu}(i) \log^\alpha \frac{1}{\boldsymbol{\mu}(i)}
$$

$$
\leq N \underbrace{\max_{x \in [0,1]}\left\{ x \log^\alpha \frac{1}{x} \right\}}_{= \mathrm{e}^{-\alpha}\alpha^\alpha} + C\beta^\alpha \underbrace{\sum_{\substack{i \in \Omega \\ i > N \\ \boldsymbol{\mu}(i) \leq \mathrm{e}^{-\alpha}}} \frac{1}{i^\beta} \log^\alpha i}_{\leq S_{\alpha, \beta} < \infty} + \underbrace{\sum_{\substack{i \in \Omega \\ i > N \\ \boldsymbol{\mu}(i) > \mathrm{e}^{-\alpha}}} \boldsymbol{\mu}(i) \log^\alpha \frac{1}{\boldsymbol{\mu}(i)}}_{\text{at most } \mathrm{e}^\alpha \text{ elements}}
$$

$$
\leq (N\mathrm{e}^{-\alpha} + 1)\alpha^\alpha + C\beta^\alpha S_{\alpha, \beta},
$$

such that for any $\alpha > 1$, there exists $M_\alpha < \infty$ that bounds the $\alpha$-moment of information. Notice that although existence is guaranteed, $N, C$ depend on the unknown $\boldsymbol{\mu}$. The asymptotic nature of assumption $A_2(\beta)$ is therefore not enough to specify what $M_\alpha$ is.

**Proof** $(h)$   Starting from (23),

$$
\sum_{i \in \Omega} \boldsymbol{\mu}(i) \log^\alpha \frac{1}{\boldsymbol{\mu}(i)} \leq \alpha^\alpha + \sum_{\substack{i \in \Omega \\ \boldsymbol{\mu}(i) \leq \mathrm{e}^{-\alpha}}} \boldsymbol{\mu}(i) \log^\alpha \frac{1}{\boldsymbol{\mu}(i)} \leq \alpha^\alpha + \sum_{\substack{i \in \Omega \\ \boldsymbol{\mu}(i) \leq \mathrm{e}^{-\alpha}}} \frac{\overline{c}_\beta}{i^\beta} \log^\alpha \frac{i^\beta}{\overline{c}_\beta}
$$

which is upper bounded by a converging series, whose value is entirely defined by $\alpha, \beta$ and $\overline{c}_\beta$.

**Proof** $(i)$   Follows the arguments of the proof for $(a)$ and the proof of $(h)$. The series converges exactly when $\alpha < \gamma - 1$, and if it does, the value of the converging series is a function of $\alpha, \gamma, \overline{c}_\gamma$.

**Proof** $(b)$   Assume that $\boldsymbol{\mu}$ satisfies $A_1(\alpha)$ for some $\alpha > 1$. Then, $\forall \gamma < \alpha + 1$, the collection of distributions such that $\boldsymbol{\mu}(i) \in \mathcal{O}\left(\frac{(\log \log i)^\delta}{i \log^\gamma i}\right), \delta > 1$ also verifies the hypothesis.