

THE QUERY COMPLEXITY OF ESTIMATING ENTROPY

by

Jafar Jafarov

B.S., Computer Engineering, Boğaziçi University, 2014

Submitted to the Institute for Graduate Studies in
Science and Engineering in partial fulfillment of
the requirements for the degree of
Master of Science

Graduate Program in Computer Engineering
Boğaziçi University

2016

ACKNOWLEDGEMENTS

I am deeply indebted to my advisor A. C. Cem Say for introducing me to the area of theoretical computer science. It was his simple yet deep and poetic introduction that inspired me to pursue science. His endless support and guidance was truly invaluable.

I am obliged to Ryan O'Donnell not only for proposing the topic or contributing greatly to this thesis but also for being an inspiration both as a person and as a scientist.

I am grateful to Ali Taylan Cemgil and Wolfgang Hörmann, my thesis committee members, for devoting their time and energy to this work.

I would like to express my gratitude to Jayadev Acharya for his helpful answers to my questions.

But above all, I would like to thank my family for their continual support and my friends for making my years at Boğaziçi University so extraordinary.

ABSTRACT

THE QUERY COMPLEXITY OF ESTIMATING ENTROPY

We investigate the query complexity of additively estimating entropy of a discrete probability distribution in two settings. Let \mathbf{p} be an unknown probability distribution on $[n] := \{1, 2, \dots, n\}$, and define two kinds of queries: A **SAMP** query takes no input and returns $x \in [n]$ with probability $\mathbf{p}[x]$; a **PMF** query takes as input $x \in [n]$ and returns the value $\mathbf{p}[x]$. In the **SAMP** model of query complexity, the only allowed interaction with \mathbf{p} is via **SAMP** queries. In the **SAMP+PMF** model, both **SAMP** and **PMF** queries are utilized to interact with \mathbf{p} .

In particular, we consider the task of estimating the entropy of \mathbf{p} to within $\pm\Delta$ (with high probability). For the usual Shannon entropy $H(\mathbf{p})$, we review the matching upper and lower bounds established by Valiant and Valiant in the **SAMP** model, and describe the algorithm constructed by Canonne and Rubinfeld in the **SAMP+PMF** model. For the Rényi entropy $H_\alpha(\mathbf{p})$, we analyze three different matching upper and lower bound pairs introduced by Acharya *et al.* in the **SAMP** model.

We show that $\Omega(\log^2 n/\Delta^2)$ queries are necessary to estimate the Shannon entropy $H(\mathbf{p})$ in the **SAMP+PMF** model, matching a recent upper bound of Canonne and Rubinfeld. In addition, we prove that $\Theta(n^{1-1/\alpha})$ queries are necessary and sufficient to estimate the Rényi entropy $H_\alpha(\mathbf{p})$ in the **SAMP+PMF** model, where $\alpha > 1$. This complements recent work of Acharya *et al.* in the **SAMP** model that showed $O(n^{1-1/\alpha})$ queries suffice when α is an integer, but roughly n queries are necessary when α is a noninteger. All of our lower bounds extend to the **SAMP+CDF** model, where **SAMP** and **CDF** queries (given x , return $\sum_{y \leq x} \mathbf{p}[y]$) are allowed. We give a matching lower bound on estimating the support size (the number of domain elements with nonzero probability) of an unknown distribution \mathbf{p} in the **SAMP+CDF** model. Lastly, we present an upper bound on additively estimating Tsallis entropy in the **SAMP+PMF** model.

ÖZET

ENTROPİ KESTİRİMİNİN SORGU KARMAŞIKLIĞI

Bu çalışmada, ayrık olasılık dağılımının entropisinin toplanır hata payıyla kestirimi iki farklı kurguda irdelenmektedir. Buna göre bilinmeyen bir olasılık dağılımı \mathbf{p} 'ye erişim iki farklı sorgu türüyle sağlanmaktadır. Herhangi bir girdisi olmayan SAMP sorgusu $\mathbf{p}[x]$ olasılığıyla $x \in [n]$ döndürmektedir. Girdi olarak $x \in [n]$ alan PMF sorgusunun ise çıktısı $\mathbf{p}[x]$ 'dir. SAMP modeli ismini verdiğimiz ilk kurguda \mathbf{p} ile sadece SAMP sorgusu vasıtasıyla iletişim sağlanmaktadır. SAMP+PMF modeli olarak adlandırdığımız ikinci kurgudaysa hem SAMP hem de PMF sorguları kullanılabilir.

Daha kesin bir ifadeyle, bu çalışmanın odak noktası olasılık dağılımı \mathbf{p} 'nin entropisinin yüksek ihtimalle $\pm\Delta$ toplanır hata payıyla kestirimi problemidir. Shannon entropisi $H(\mathbf{p})$ 'nin kestirimini incelediğimiz bölümde Valiant ve Valiant'ın SAMP modelinde göstermiş olduğu eşleşen alt ve üst sınırları ve Canonne ve Rubinfeld'in SAMP+PMF modelinde inşa etmiş olduğu algoritmayı tasvir ediyoruz. Rényi entropisi $H_\alpha(\mathbf{p})$ 'yi incelediğimiz bölümdeyse Acharya ve diğerleri tarafından sunulan üç farklı eşleşen alt ve üst sınır çiftini analiz ediyoruz.

Kendi katkımız olarak, önce SAMP+PMF modelinde Shannon entropisi $H(\mathbf{p})$ 'nin kestirimi probleminin $\Omega(\frac{\log^2 n}{\Delta^2})$ sayıda sorgu gerektirdiğini kanıtlayarak Canonne ve Rubinfeld'in sunduğu üst sınırın optimal olduğunu gösteriyoruz. İkinci olarak, yine SAMP+PMF modelinde Rényi entropisi $H_\alpha(\mathbf{p})$ 'yi $\alpha > 1$ değerlerinde kestirebilmek için $\Theta(n^{1-1/\alpha})$ sayıda sorgunun gerekli ve yeterli olduğunu ispatlıyoruz. Böylelikle Acharya ve diğerleri tarafından yakın zamanda elde edilen, SAMP modelinde Rényi entropisi $H_\alpha(\mathbf{p})$ 'yi $\alpha > 1$ tamsayı değerlerinde kestirebilmek için $O(n^{1-1/\alpha})$ sayıda sorgunun yeterli olduğu fakat $\alpha > 1$ tamsayı olmayan değerlerinde kestirebilmek için kabaca n sayıda sorgunun gerekli olduğu yönündeki sonuçları tamamlamış oluyoruz.

TABLE OF CONTENTS

ACKNOWLEDGEMENTS	iv
ABSTRACT	v
ÖZET	vi
LIST OF FIGURES	viii
LIST OF SYMBOLS	ix
LIST OF ACRONYMS/ABBREVIATIONS	x
1. INTRODUCTION	1
1.1. Organization	3
2. SAMP MODEL	4
2.1. Shannon Entropy	5
2.1.1. Upper Bound I	7
2.1.2. Lower Bound	16
2.1.3. Upper Bound II	20
2.2. Rényi Entropy	23
2.2.1. Upper Bounds	25
2.2.2. Lower Bounds	29
3. SAMP+PMF MODEL	32
4. OPTIMAL BOUNDS FOR ESTIMATING ENTROPY WITH PMF QUERIES	35
4.1. Our Results, and Comparison with Prior Work	35
4.2. First Main Theorem	36
4.3. Second Main Theorem	45
4.3.1. Lower Bound	45
4.3.2. Upper Bound	48
4.4. Support Size	50
4.5. Tsallis Entropy	52
5. CONCLUSION	54
REFERENCES	56
APPENDIX A: INEQUALITIES	60

LIST OF FIGURES

Figure 2.1.	Canonical Estimator of Shannon Entropy	13
Figure 2.2.	Linear Estimator of Shannon Entropy	22
Figure 2.3.	SAMP Estimator of Rényi Entropy (of integer degree $\alpha > 1$)	28
Figure 2.4.	SAMP Estimator of Rényi Entropy (of noninteger degree $\alpha > 1$) . .	29
Figure 3.1.	SAMP+PMF Estimator of Shannon Entropy	34
Figure 4.1.	SAMP+PMF Estimator of Rényi Entropy (of degree $\alpha > 1$)	50

LIST OF SYMBOLS

$\text{Bin}(\cdot, \cdot)$	Binomial distribution
$\mathbf{Cov}[\cdot]$	Covariance operator
$\exp(\cdot)$	Exponential operator
$\mathbf{E}[\cdot]$	Expected value operator
$\mathcal{F}_{\mathcal{X}}$	Fingerprint of a sample set \mathcal{X}
h	Histogram of a distribution \mathbf{p}
$H(\mathbf{p})$	Shannon entropy of a distribution \mathbf{p}
$H_{\alpha}(\mathbf{p})$	Rényi entropy of degree α of a distribution \mathbf{p}
$\mathbf{KL}(\cdot)$	Kullback-Leibler divergence
$\mathcal{M}_{\alpha}(\mathbf{p})$	Moment of degree α of a distribution \mathbf{p}
$\widehat{\mathcal{M}}_{\alpha}^e$	Empirical estimator of $\mathcal{M}_{\alpha}(\mathbf{p})$
$\widehat{\mathcal{M}}_{\alpha}^u$	Bias-corrected estimator of $\mathcal{M}_{\alpha}(\mathbf{p})$
$[n]$	Set of natural numbers $\{1, \dots, n\}$
N_i	Multiplicity of a domain element i
\mathbf{p}	Discrete probability distribution on domain $[n]$
$\text{Pois}(\cdot)$	Poisson distribution
$\text{supp}(\mathbf{p})$	Support size of a distribution \mathbf{p}
$\mathcal{U}([n])$	Uniform distribution on the domain $[n]$
$\mathbf{Var}[\cdot]$	Variance operator
\mathcal{X}	Set of m independent samples $\{X_1, \dots, X_m\}$
δ	Error probability
Δ	Additive accuracy
$\mathbb{1}_E$	Indicator function of an event E

LIST OF ACRONYMS/ABBREVIATIONS

CLT	Central Limit Theorem
CDF	Cumulative Distribution Function
LP	Linear Program
KL	Kullback-Leibler
PMF	Probability Mass Function
SAMP	Sampling

1. INTRODUCTION

The question of what to infer about an unknown probability distribution \mathbf{p} given samples from it is fundamental to the field of statistics and has been researched for decades. However, the practicality of traditional techniques has been shaken due to the rapidly growing size of data in scientific study, a phenomenon designated as big data. In the absence of simplifying assumptions about a probability distribution such as being of a specific type or possessing certain smoothness properties, the number of samples utilized by such techniques grows linearly in the size of the domain of a distribution which is huge in the realm of big data. Thus, the task of constructing algorithms with sublinear sample complexity has become all-important, and the aforementioned question has been recently investigated within the theoretical computer science framework of *property testing*. In this framework, as its name implies, a certain characteristic of a distribution is put under the microscope. In addition, the only assumption made about \mathbf{p} is that it is a discrete probability distribution on a finite domain $[n] := \{1, 2, \dots, n\}$ where $n \in \mathbb{N}$. For a detailed exploration of the field, see the surveys by Rubinfeld [1] and Canonne [2].

One of the most significant characteristics of a probability distribution is its Shannon entropy, $H(\mathbf{p}) = -\sum_{i=1}^n \mathbf{p}[i] \log \mathbf{p}[i]$,¹ which represents the “amount of randomness” a distribution possesses. The first focal point of this work is estimating Shannon entropy to within additive error Δ with probability at least $1 - \delta$. (In a typical scenario $\Delta = 1$ and $\delta = 1/3$.) The reason we limit ourselves to additive rather than multiplicative estimation is that it is directly related to the estimation of mutual information. That is, if \mathbf{p} is a joint probability distribution of two random variables X, Y , then additively estimating mutual information $I(X, Y) = H(X) + H(Y) - H(X, Y)$ is realized via additively estimating $H(\mathbf{p})$. For deeper analysis of Shannon entropy and mutual information see Paninski [3]. The second focal point of this work is estimating another popular type of entropy, Rényi entropy $H_\alpha(\mathbf{p}) = \frac{1}{1-\alpha} \log \left(\sum_{i=1}^n \mathbf{p}[i]^\alpha \right)$, to within additive error Δ with probability at least $1 - \delta$, where $\alpha \in [0, 1) \cup (1, \infty)$. Both tasks

¹In this work, \log denotes \log_2 .

are investigated in two different settings.

In the conventional setting, which we refer to as the **SAMP** model, the only allowed interaction with a probability distribution \mathbf{p} is via independent samples. As recently shown in [3–6], $\Theta\left(\frac{n}{\log n}\right)$ samples are necessary and sufficient to estimate Shannon entropy to within a constant additive error with high probability. The case for estimating Rényi entropy is more complicated; three different results for three classes of α are obtained in [7], the most efficient one being for the class of integer $\alpha > 1$ with $\Theta\left(n^{1-1/\alpha}\right)$ sample complexity. These quantities are not always convenient, considering the aforementioned tendency in science and technology.

In the “unconventional” setting referred to as the **SAMP+PMF** model, aside from drawing independent samples as in the **SAMP** model, querying a probability mass function (PMF)² of an arbitrary domain element, that is, learning $\mathbf{p}[i]$ of an element $i \in [n]$ is allowed. This extended version of the **SAMP** model, called the “Generation+Evaluation” model in [8] and the “combined model” in [9], is introduced to overcome the difficulty described above. The results achieved in [10] imply that estimating Shannon entropy is possible with $\text{polylog}(n)$ **SAMP+PMF** queries, exponentially better than the $\Omega\left(\frac{n}{\log n}\right)$ queries in the **SAMP** model.

Although described as unconventional, the **SAMP+PMF** model becomes practical in many applications. For instance, the number of occurrences, and therefore the probability of a certain element in a sorted database can be calculated via at most logarithmically many interactions with the database. For a concrete example, consider the Google n-gram database in which the frequency of each n-gram is published, and a random n-gram is easily obtained from the underlying text corpus. Another motivation for the **SAMP+PMF** model stems from its strong relation with the *streaming* model [11], where entropy estimation has been thoroughly studied [12–17]. To exhibit the relation between the two, note that any q -query estimation algorithm in the **SAMP+PMF** model can be converted to a $q \cdot \text{polylog}(n)$ -space streaming algorithm with one or two passes

²In this work, PMF, CDF and **SAMP** are abbreviations for probability mass function, cumulative distribution function and sampling, respectively.

(details of the conversion depend on the model for how the items in the stream are ordered). For more motivation and results for the **SAMP+PMF** model, see Canonne and Rubinfeld [10].

Our main contribution [18] is to establish a lower bound matching the upper bound obtained in [10] on additively estimating Shannon entropy, $\Omega(\log^2 n)$. In addition, we found upper and lower bounds matching in their dependence on n for additive estimation of Rényi entropy when $\alpha > 1$, $\Theta(n^{1-1/\alpha})$.

1.1. Organization

In Chapter 2, we focus on the task of additively estimating entropy in the **SAMP** model. Section 2.1 is devoted to rigorously analyzing optimal upper and lower bounds achieved in three successive works [19–21]. Section 2.2 explores three different matching upper and lower bound pairs on additive estimation of Rényi entropy obtained in [7].

In Chapter 3, we concentrate on additively approximating the Shannon entropy in the **SAMP+PMF** model. We describe the exponentially better algorithm constructed by Canonne and Rubinfeld [10]. In addition, we introduce the **SAMP+CDF** model, which is an extension of the **SAMP+PMF** model.

In Chapter 4, we demonstrate our contribution to the problem. Section 4.1 includes a comparison between our results and the prior work. In Section 4.2, we build an optimal lower bound on additively estimating Shannon entropy, and in Section 4.3, we present upper and lower bounds on additive estimation of Rényi entropy in the **SAMP+PMF** and **SAMP+CDF** models. In Section 4.4, we establish a lower bound on estimating the support size of a probability distribution in the **SAMP+CDF** model. Section 4.5 is devoted to constructing an algorithm for additively estimating Tsallis entropy in the **SAMP+PMF** model.

Finally, we state some open questions in the conclusion and give the inequalities used throughout this work in the appendix.

2. SAMP MODEL

We start with the formal definition of the conventional model.

Definition 2.1. (*sampling-only model*) Let \mathbf{p} be a probability distribution on $[n]$ and *SAMP* denote a type of query which takes no input and returns $i \in [n]$ with probability $\mathbf{p}[i]$ independently of all previous calls. The *SAMP* model is a model of query complexity in which the only allowed interaction with \mathbf{p} is via *SAMP* queries.

Although there are different metrics to measure the distance between two probability distributions $\mathbf{p}_1, \mathbf{p}_2$, the most commonly used metric is *total variation distance*, denoted by d_{TV} and defined as

$$d_{\text{TV}}(\mathbf{p}_1, \mathbf{p}_2) := \frac{1}{2} \|\mathbf{p}_1 - \mathbf{p}_2\|_1 = \frac{1}{2} \sum_{i=1}^n |\mathbf{p}_1[i] - \mathbf{p}_2[i]|. \quad (2.1)$$

The following identity unveils the ‘‘mystery’’ behind the constant factor $\frac{1}{2}$.

$$d_{\text{TV}}(\mathbf{p}_1, \mathbf{p}_2) = \max_{E \subseteq [n]} \{\mathbf{p}_1(E) - \mathbf{p}_2(E)\} \quad (2.2)$$

Proof. Let $A = \{i : \mathbf{p}_1[i] \geq \mathbf{p}_2[i]\}$. Then

$$d_{\text{TV}}(\mathbf{p}_1, \mathbf{p}_2) = \frac{1}{2} (\mathbf{p}_1(A) - \mathbf{p}_2(A)) + \frac{1}{2} (\mathbf{p}_2(\bar{A}) - \mathbf{p}_1(\bar{A})) = \mathbf{p}_1(A) - \mathbf{p}_2(A),$$

where $\bar{A} := [n] \setminus A$. The next step is to show that A is an event maximizing the right-hand side of Equation 2.2. If one adds another element $j \in [n]$ to A , the difference $\mathbf{p}_1(A) - \mathbf{p}_2(A)$ decreases, since by definition $\mathbf{p}_1[j] < \mathbf{p}_2[j]$. Similarly, removing an element j from A leads to a decrease in $\mathbf{p}_1(A) - \mathbf{p}_2(A)$, since $\mathbf{p}_1[j] \geq \mathbf{p}_2[j]$. \square

2.1. Shannon Entropy

Shannon entropy, named after Claude E. Shannon [22], represents the expected information a probability distribution contains, thus, measures the randomness in a distribution and the compressibility of the data produced by that distribution. Shannon entropy is defined as

$$H(\mathbf{p}) = - \sum_{i=1}^n \mathbf{p}[i] \log \mathbf{p}[i]. \quad (2.3)$$

By convention, the quantity $\mathbf{p}[i] \log \mathbf{p}[i]$ is set to 0 in the case of $\mathbf{p}[i] = 0$ for some i which is consistent with the following: $\lim_{x \rightarrow 0^+} x \log x = 0$. Note that

$$0 \leq H(\mathbf{p}) \leq \log n. \quad (2.4)$$

The left-hand side of Inequality 2.4 is trivial, since $-\mathbf{p}[i] \log \mathbf{p}[i] \geq 0$ for all i . The right-hand side of Inequality 2.4 follows from the fact that $H(\mathbf{p}) = \mathbf{E}_{i \sim \mathbf{p}} \left[\log \frac{1}{\mathbf{p}[i]} \right]$ and $\log x$ is a concave function. By applying Jensen's inequality,³

$$H(\mathbf{p}) = \mathbf{E}_{i \sim \mathbf{p}} \left[\log \frac{1}{\mathbf{p}[i]} \right] \leq \log \left(\mathbf{E}_{i \sim \mathbf{p}} \left[\frac{1}{\mathbf{p}[i]} \right] \right) = \log n. \quad (2.5)$$

Shannon entropy has many applications such as measuring genetic diversity [23], quantifying neural activity [3], and detecting network anomalies [13].

This work concentrates on additive approximation to the entropy, though we also state the results regarding its multiplicative counterpart. Batu, Dasgupta, Kumar and Rubinfeld [9] construct an algorithm approximating $H(\mathbf{p})$ of a distribution \mathbf{p} within a multiplicative factor of γ using $\tilde{O} \left(n^{(1+o(1))/\gamma^2} \right)$ samples given $H(\mathbf{p}) = \Omega(\gamma)$ for any $\gamma > 1$. They also show that no algorithm exists which γ -approximates the entropy of every distribution. Furthermore, Valiant [5] proves that $\Omega \left(n^{1/\gamma^2} \right)$ samples are necessary for the task given $H(\mathbf{p}) = \Omega \left(\frac{\log n}{\gamma^2} \right)$.

³For the definitions of the inequalities used throughout this work, see Appendix A.

Initially, we introduce a trivial upper bound on additively estimating Shannon entropy.

Fact 2.2. *There exists an algorithm estimating (with high probability) the Shannon entropy to within arbitrarily small constant Δ using $O\left(\frac{n \log^2 n}{\Delta^2}\right)$ samples.*

Proof. The overall idea is to build a distribution \mathbf{p}' which is close to \mathbf{p} in total variation distance and output $H(\mathbf{p}')$ as an approximation of $H(\mathbf{p})$. We utilize the “plug-in” distribution for \mathbf{p}' . Namely, let X_1, \dots, X_m be m independent samples drawn from \mathbf{p} and define $\mathbf{p}'[i] = \frac{1}{m} \sum_{j=1}^m \mathbb{1}_{\{X_j=i\}}$ for each $i \in [n]$ where $\mathbb{1}_E$ is the indicator function for an event E . Observe that $\mathbf{p}'[i]$ is an unbiased estimator of $\mathbf{p}[i]$ since $\mathbf{E}[\mathbf{p}'[i]] = \mathbf{p}[i]$ for each i . To show that $|H(\mathbf{p}') - H(\mathbf{p})| \leq \Delta$ we use the following fact:

Fact 2.3. ([24], Lemma 8) *Let $\mathbf{p}_1, \mathbf{p}_2$ be two arbitrary probability distributions such that $d_{\text{TV}}(\mathbf{p}_1, \mathbf{p}_2) \leq \frac{\Delta}{4 \log n}$, then $|H(\mathbf{p}_1) - H(\mathbf{p}_2)| \leq \Delta$ where Δ is an arbitrarily small constant.*

The final step is to prove that $d_{\text{TV}}(\mathbf{p}', \mathbf{p}) \leq \frac{\Delta}{4 \log n}$ with high probability when $m = O\left(\frac{n \log^2 n}{\Delta^2}\right)$. Similarly, $\mathbf{p}'[E] := \frac{1}{m} \sum_{j=1}^m \mathbb{1}_{\{X_j \in E\}}$ is an unbiased estimator of $\mathbf{p}[E]$ for any event $E \subseteq [n]$. Observe that $\Pr\left[|\mathbf{p}'[E] - \mathbf{p}[E]| > \frac{\Delta}{4 \log n}\right] \leq 2e^{-\frac{\Delta^2}{8 \log^2 n} m} = e^{-O(n)}$ for such m , which is easily derived via the Hoeffding bound. Then,

$$\begin{aligned} \Pr\left[d_{\text{TV}}(\mathbf{p}', \mathbf{p}) > \frac{\Delta}{4 \log n}\right] &= \Pr\left[\max_{E \subseteq [n]} \{\mathbf{p}'(E) - \mathbf{p}(E)\} > \frac{\Delta}{4 \log n}\right] \\ &\leq \Pr\left[\bigcup_{E \subseteq [n]} \left\{\mathbf{p}'(E) - \mathbf{p}(E) > \frac{\Delta}{4 \log n}\right\}\right] \\ &\leq \sum_{E \subseteq [n]} \Pr\left[\mathbf{p}'(E) - \mathbf{p}(E) > \frac{\Delta}{4 \log n}\right] \\ &\leq 2^n \cdot e^{-O(n)} = o(1). \quad \square \end{aligned}$$

The task of achieving an additive estimation of entropy is fulfilled in its entirety in three successive works [19], [20] and [21]. Valiant and Valiant establish matching upper

and lower bounds on additively estimating a major class of symmetric properties (a property is symmetric if it is immune to any permutation of domain elements) including entropy. It is proven that $\Theta\left(\frac{n}{\log n}\right)$ samples are necessary and sufficient for additive estimation of entropy of a probability distribution with support size at most n . We use a rather technical narrative for two reasons. First, each work has an intricate structure requiring an attentive analysis. Second, the techniques deployed in the process utilize a wide range of mathematical notions that may be of independent interest.

2.1.1. Upper Bound I

In this part, we introduce the algorithm constructed in [20] estimating entropy up to an arbitrarily small constant using $O\left(\frac{n}{\log n}\right)$ independent samples. As the title of the paper (*Estimating the unseen: A sublinear-sample canonical estimator of distributions*) suggests, Valiant and Valiant employ a canonical approach to delicately approximate an unobserved portion of a probability distribution rather than directly estimating entropy. In other words, based on the samples drawn from an unknown probability distribution \mathbf{p} , the estimator builds a probability distribution \mathbf{p}' such that with high probability the two are “close”, and returns the entropy of \mathbf{p}' as an approximation of the entropy of \mathbf{p} . The success of obtaining sublinear-sample complexity is due to exploiting features of symmetry and using a different distance metric to better capture the “closeness” between \mathbf{p} and \mathbf{p}' . Before going into the details, we introduce some key definitions.

Definition 2.4. *A property of a distribution is a function $\pi : \mathbf{p}^n \rightarrow \mathbb{R}$, where \mathbf{p}^n is the set of distributions on domain $[n]$. A property π is called a **symmetric property** if for all distributions \mathbf{p} , and all permutations σ , $\pi(\mathbf{p}) = \pi(\mathbf{p} \circ \sigma)$.*

Note that entropy is a symmetric property.

Definition 2.5. *Given a sequence of samples $\mathcal{X} = \{X_1, \dots, X_m\}$, let the associated **fingerprint**, denoted by $\mathcal{F}_{\mathcal{X}}$, be the vector whose i^{th} component, $\mathcal{F}_{\mathcal{X}}(i)$ is the number of domain elements that occur exactly $i \geq 1$ times in sample \mathcal{X} .*

Intuitively, the fingerprint of a sample should hold all necessary information about a sample for the task of estimating a symmetric property. We formalize this intuition via the following fact.

Fact 2.6. ([9], Lemma 8) *For any algorithm \mathcal{A} that approximates the entropy of a distribution to within additive Δ from samples, there exists an algorithm \mathcal{A}' which gets as input only the fingerprint of the generated sample and has the error probability upper bounded by that of \mathcal{A} .*

Proof. Let \mathbf{p} be an unknown distribution and $\mathcal{F}_{\mathcal{X}}$ denote the fingerprint of the set of samples $\mathcal{X} = \{X_1, \dots, X_m\}$ drawn from \mathbf{p} . Algorithm \mathcal{A}' is constructed as follows:

- Choose $\mathcal{F}_{\mathcal{X}}(i)$ elements at random from $[n]$ without replacement for each i ,⁴
- Build \mathcal{X}' so that an element chosen in step i occurs exactly i times in \mathcal{X}' ,
- Output the value that \mathcal{A} outputs on \mathcal{X}' .

The next step is to prove the correctness of \mathcal{A}' . Let π be a permutation on $[n]$ and define a permuted distribution $\pi(\mathbf{p})$ such that $\pi(\mathbf{p})[i] = \mathbf{p}[\pi(i)]$. Let $\pi(\mathcal{X})$ be a set of samples by relabeling the members of \mathcal{X} according to π . Observe that the set \mathcal{X}' generated by \mathcal{A}' is $\pi(\mathcal{X})$ for some random permutation π . Lastly, let $\mathcal{A}(\mathcal{X})$ denote the output of \mathcal{A} on the sample set \mathcal{X} . Then,

$$\begin{aligned}
& \Pr[\mathcal{A}' \text{ estimates } H(\mathbf{p}) \text{ to within } \Delta] \\
&= \sum_{\mathcal{X}} \Pr[\mathbf{p} \text{ generates } \mathcal{X}] \cdot \mathbf{E}_{\pi} [\Pr[\mathcal{A}(\pi(\mathcal{X})) \text{ estimates } H(\mathbf{p}) \text{ to within } \Delta]] \\
&= \mathbf{E}_{\pi} \left[\sum_{\mathcal{X}} \Pr[\mathbf{p} \text{ generates } \mathcal{X}] \cdot \Pr[\mathcal{A}(\pi(\mathcal{X})) \text{ estimates } H(\mathbf{p}) \text{ to within } \Delta] \right] \\
&= \mathbf{E}_{\pi} \left[\sum_{\mathcal{X}} \Pr[\pi(\mathbf{p}) \text{ generates } \pi(\mathcal{X})] \cdot \Pr[\mathcal{A}(\pi(\mathcal{X})) \text{ estimates } H(\mathbf{p}) \text{ to within } \Delta] \right] \\
&= \mathbf{E}_{\pi} [\Pr[\mathcal{A} \text{ estimates } H(\pi(\mathbf{p})) \text{ to within } \Delta]] \\
&\geq \min_{\pi} \Pr[\mathcal{A} \text{ estimates } H(\pi(\mathbf{p})) \text{ to within } \Delta],
\end{aligned}$$

⁴Note that $n - \|\mathcal{F}_{\mathcal{X}}\|_1$ is the number of elements not seen in the sample set \mathcal{X} .

which is the correctness probability of \mathcal{A} . \square

Similarly we define a histogram of the distribution which categorizes the domain elements according to their probability values.

Definition 2.7. The *histogram* of a distribution \mathbf{p} is a mapping $h : (0, 1] \rightarrow \mathbb{Z}$, where $h(x) = |\{i : \mathbf{p}[i] = x\}|$. Additionally, generalized histograms are allowed which do not necessarily take integral values.

Observe that, a symmetric property is a function of the histogram of a distribution. For instance, Shannon entropy can be written as

$$H(\mathbf{p}) = - \sum_{i=1}^n \mathbf{p}[i] \log \mathbf{p}[i] = \sum_{x:h(x) \neq 0} h(x) x \log \frac{1}{x}. \quad (2.6)$$

We now define a new distance metric to obtain a better measure for proximity between distributions.

Definition 2.8. For two histograms (or generalized histograms) h_1, h_2 , let the *relative earthmover distance* between them, $R(h_1, h_2)$, be the minimum over all schemes of moving the probability mass of the first histogram to yield the second histogram, of the cost of moving that mass, where the per-unit cost of moving mass from probability x to y is $|\log(x/y)|$.

Distributions which are close to each other according to this new metric have similar entropies.

Definition 2.9. A symmetric property π is (Δ, ϑ) -continuous if for all distributions $\mathbf{p}_1, \mathbf{p}_2$ with respective histograms h_1, h_2 satisfying $R(h_1, h_2) \leq \vartheta$ it follows that

$$|\pi(\mathbf{p}_1) - \pi(\mathbf{p}_2)| \leq \Delta. \quad (2.7)$$

Fact 2.10. ([20], Fact 9) For a probability distribution \mathbf{p} , and $\Delta > 0$ the Shannon entropy, $H(\mathbf{p}) = - \sum_{i=1}^n \mathbf{p}[i] \log \mathbf{p}[i]$ is (Δ, Δ) -continuous, with respect to the relative earthmover distance.

We describe a well-known sampling technique known as *Poisson sampling*. Recall that in the standard approach of estimating a certain $\mathbf{p}[i]$, one needs to draw m independent samples X_1, \dots, X_m from \mathbf{p} and calculate the multiplicity of a domain element i , that is, the number of occurrences of i among m samples. Observe that the multiplicities of any two elements are not independent, complicating the overall analysis: for a start, the sum of the multiplicities of all domain elements must be equal to m . To overcome this difficulty, instead of drawing exactly m independent samples from \mathbf{p} we draw $M \sim \text{Pois}(m)$ samples, where $\text{Pois}(m)$ is the Poisson distribution with parameter m . Let X_1, \dots, X_M be independent samples drawn from \mathbf{p} , then the multiplicity of a domain element i is defined as

$$N_i = |\{1 \leq j \leq M : X_j = i\}|. \quad (2.8)$$

Fact 2.11. *The multiplicities $\{N_i\}$ are independent random variables and distributed as $\text{Pois}(m \cdot \mathbf{p}[i])$.*

Proof. Let $\text{Bin}(y, z)$ be a Binomial distribution with parameters $y \in \mathbb{N}$ and $z \in [0, 1]$. Then

$$\begin{aligned} \Pr[N_i = j] &= \sum_{M=0}^{\infty} (\Pr[\text{Pois}(m) = M] \cdot \Pr[\text{Bin}(M, \mathbf{p}[i]) = j]) \\ &= \sum_{M=j}^{\infty} (\Pr[\text{Pois}(m) = M] \cdot \Pr[\text{Bin}(M, \mathbf{p}[i]) = j]) \\ &= \sum_{M=j}^{\infty} \left(\frac{e^{-m} m^M}{M!} \cdot \frac{M!}{j!(M-j)!} \cdot \mathbf{p}[i]^j (1 - \mathbf{p}[i])^{M-j} \right) \\ &= \frac{e^{-m} \mathbf{p}[i]^j}{j!} \sum_{M'=0}^{\infty} \left(\frac{m^{M'+j}}{M'!} \cdot (1 - \mathbf{p}[i])^{M'} \right) \\ &= \frac{e^{-m} \cdot (m\mathbf{p}[i])^j}{j!} \sum_{M'=0}^{\infty} \left(\frac{(m(1 - \mathbf{p}[i]))^{M'}}{M'!} \right) \\ &= \frac{e^{-m} \cdot (m\mathbf{p}[i])^j \cdot e^{m - m\mathbf{p}[i]}}{j!} \\ &= \frac{e^{-m\mathbf{p}[i]} \cdot (m\mathbf{p}[i])^j}{j!} = \text{Pois}(m\mathbf{p}[i], j). \end{aligned}$$

Now we show that the N_i 's are independent.

$$\begin{aligned}
\Pr [N_1 = M_1 \ \& \ \dots \ \& \ N_n = M_n] &= \Pr [M] \cdot \Pr [N_1 = M_1 \ \& \ \dots \ \& \ N_n = M_n \mid M] \\
&= \frac{e^{-m} m^M}{M!} \cdot \frac{M!}{M_1! \dots M_n!} \cdot \mathbf{p}[1]^{M_1} \dots \mathbf{p}[n]^{M_n} \\
&= \frac{e^{-m\mathbf{p}[1]} (m\mathbf{p}[1])^{M_1}}{M_1!} \dots \frac{e^{-m\mathbf{p}[n]} (m\mathbf{p}[n])^{M_n}}{M_n!} \\
&= \Pr [N_1 = M_1] \dots \Pr [N_n = M_n]. \quad \square
\end{aligned}$$

Observe that N_i/m is an unbiased estimator for $\mathbf{p}[i]$ because $\mathbf{E}[\frac{N_i}{m}] = \frac{\mathbf{E}[N_i]}{m} = \mathbf{p}[i]$. Consider the distribution of the j^{th} entry of a $\text{Pois}(m)$ -sample fingerprint,

$$\mathcal{F}(j) = \sum_{i=1}^n \mathbb{1}\{N_i = j\} \Rightarrow \mathbf{E}[\mathcal{F}(j)] = \sum_{i=1}^n \text{Pois}(m\mathbf{p}[i], j) = \sum_{x:h(x) \neq 0} h(x) \text{Pois}(mx, j). \quad (2.9)$$

The direct consequence of N_i 's being independent is that for each j , $\mathcal{F}(j)$ is closely concentrated around its expectation, having an easy proof by a direct application of the Chernoff bound. The other advantage of Poisson sampling is that its sample complexity is comparable to the sample complexity of usual sampling, because Poisson distribution also has a concentration around its expectation.

Proposition 2.12. ([20], Proposition 21) *Given $m > 30$, and any set of fingerprints A , let \bar{A} be the set of fingerprints that can be obtained by adding or removing at most $m^{0.6}$ samples from some fingerprint in set A . Let \mathcal{F} denote a random m -sample fingerprint, and let \mathcal{F}' denote a fingerprint obtained from choosing $M \sim \text{Pois}(m)$, random samples. Then*

$$\Pr [\mathcal{F} \in A] \leq \Pr [\mathcal{F}' \in \bar{A}] + e^{-m^{0.1/2}}.$$

The construction of the estimator starts with building a linear program that focuses on the low probability portion of a distribution. Based on a fingerprint \mathcal{F} of an unknown histogram h , the objective is to derive a histogram h' such that for

each domain element i and for a fingerprint $\mathcal{F}_{\mathcal{X}}(i)$ obtained from a sample \mathcal{X} of h' , $\mathbf{E}[\mathcal{F}_{\mathcal{X}}(i)] \approx \mathcal{F}(i)$. However, for the high probability portion of a distribution, that is, for elements with probability at least $m^{-1+\alpha}$ for some small constant $\alpha \in (0, 1)$, the histogram is set as $h'(\frac{j}{m}) = \mathcal{F}(j)$.

Definition 2.13 (The Linear Program LP). *Given an m -sample fingerprint \mathcal{F} and $\alpha = 1/50, c \in [1, 2]$, bounds $A := cm^{-1+\alpha}, B := 4m^{-1+0.6\alpha}$, and a real number $\gamma := m^{-3/2}$, the linear program consists of variables $v_x \geq 0$ for all $x \leq A + B/2$ in the set $X := \{\gamma, 2^2\gamma, 3^2\gamma, \dots, A + B/2\}$, subject to the following condition:*

1. $\sum_{x \in X: x \geq A} xv_x \leq 16m^{-0.4\alpha}$
2. $\sum_{x \in X} xv_x + \sum_{j \geq m(A+B)} \frac{j}{m} \mathcal{F}_j = 1$
3. For all integers $i \leq m(A + B/4)$,

$$\sum_{x \in X} v_x \text{Pois}(mx, i) \in [\mathcal{F}(i) - 4m^{0.6+\alpha}, \mathcal{F}(i) + 4m^{0.6+\alpha}].$$

The set X is chosen carefully to adjust the time complexity of LP to be linear in the number of samples. The first constraint is to guarantee that the probability mass residing in the neighborhood of the threshold probability, $A \approx m^{-1+\alpha}$, is small. The second constraint is to guarantee that the sum of the total probability mass of rarely occurring elements and the total probability mass of frequently occurring elements is 1. The third constraint is to guarantee that for rarely occurring domain elements the expectation of a fingerprint distribution \mathcal{F}' of a histogram h' is close to a fingerprint \mathcal{F} of a histogram h .

Observe that a solution $\mathbf{v} = \{v_x\}$ does not necessarily yield a proper histogram h' since v_x 's can be nonintegers. The following definition is to construct a proper histogram which is referred as the histogram associated to a solution \mathbf{v} .

Definition 2.14. *Let $X := \{\gamma, 2^2\gamma, 3^2\gamma, \dots, A + B/2\}$ be the set of probabilities for which LP solves. Given a m -fingerprint \mathcal{F} and a solution \mathbf{v} to the associated LP, the corresponding histogram $h^{\mathbf{v}}$ is derived from \mathbf{v} according to the following process.*

1. set $h^{\mathbf{v}}(*) = 0$.
2. for all $x \in X$ let $h^{\mathbf{v}}(x) = v_x$.
3. for all integers $j \geq m(A+B)$, let $h^{\mathbf{v}}(\frac{j}{m}) = \mathcal{F}(j)$.
4. for all x such that $h^{\mathbf{v}}(x) \neq 0$, set $h^{\mathbf{v}}((1+\epsilon)x) = \lfloor h^{\mathbf{v}}(x) \rfloor$ where $\epsilon = \frac{\sum_{x \in X} x(v_x - \lfloor v_x \rfloor)}{1 - \sum_{x \in X} x(v_x - \lfloor v_x \rfloor)}$.

The first and second steps assign \mathbf{v} to $h^{\mathbf{v}}$. The third step makes the histogram $h^{\mathbf{v}}$ for frequently occurring elements agree with the \mathcal{F} of a histogram h . The last step converts the histogram values to integers while compensating the resulting loss in total probability mass by renormalizing the distribution.

Lastly, we establish the connection between LP and the aforementioned $O\left(\frac{n}{\log n}\right)$ sample complexity.

Algorithm: ESTIMATOR I

- Fix $m = O\left(\frac{n}{\log n}\right)$. (For details see [20].)
- Draw $M \sim \text{Pois}(m)$ independent samples X_1, \dots, X_M from \mathbf{p} .
- Construct LP corresponding to the fingerprint \mathcal{F} obtained from X_1, \dots, X_M .
- Find a solution \mathbf{v} to LP.
- Compute histogram $h^{\mathbf{v}}$ associated to solution \mathbf{v} , as defined in Definition 2.14.
- Output $\sum_{x: h^{\mathbf{v}}(x) \neq 0} h^{\mathbf{v}}(x) x \log \frac{1}{x}$.

Figure 2.1. Canonical Estimator of Shannon Entropy.

Theorem 2.15. ([20], Theorem 2) *For a constant $\epsilon \in (0, 1]$, consider a sample consisting of m independent samples from a histogram h of support size at most $\epsilon m \log m$. With probability at least $1 - e^{-m^{0.04}}$, LP has a solution and furthermore, for any solution to LP, \mathbf{v} , the histogram $h^{\mathbf{v}}$ associated to \mathbf{v} as in Definition 2.14 satisfies*

$$R(h, h^{\mathbf{v}}) = O\left(\sqrt{\epsilon} \cdot \max\{1, |\log \epsilon|\}\right).$$

Theorem 2.15 combining with Fact 2.10 imply that entropy of the resulting histogram h^v is close to the entropy of an unknown histogram h . The proof of Theorem 2.15 consists of two parts; in the first part it is shown that with the claimed probability there exists a feasible point \hat{v} such that the associated histogram $h^{\hat{v}}$ is close to h . We informally describe how the existence of a feasible point \hat{v} is proven.

A feasible point \hat{v} is manually constructed via discretizing the low probability portion of an unknown histogram h . For each probability $y \leq A + \frac{B}{2}$ let x_i, x_{i+1} be consecutive members of the set X such that $x_i \leq y \leq x_{i+1}$, then $h(y) > 0$ is distributed between \hat{v}_{x_i} and $\hat{v}_{x_{i+1}}$ as follows:

All initially being equal to 0, set $\hat{v}_{x_i} := \hat{v}_{x_i} + h(y) \frac{x_{i+1}-y}{x_{i+1}-x_i}$, and $\hat{v}_{x_{i+1}} := \hat{v}_{x_{i+1}} + h(y) \frac{y-x_i}{x_{i+1}-x_i}$. Observe that such interpolation preserves both the probability mass residing on y and the quantity $h(y)$. Then it is proven that for each y ,

$$h(y) \left| \left(\frac{x_{i+1}-y}{x_{i+1}-x_i} \text{Pois}(mx_i, j) + \frac{y-x_i}{x_{i+1}-x_i} \text{Pois}(mx_{i+1}, j) \right) - \text{Pois}(my, j) \right| \quad (2.10)$$

is bounded. Since fingerprint entries are closely concentrated around their expectations, bounding Expression 2.10 implies that the third condition of LP is satisfied. Intuitively, Expression 2.10 is bounded due to the sufficiently dense structure of the set X . In the next step of the construction of v' , a normalization is realized to ensure that the second constraint of LP is satisfied, that is, the probability mass shared among the members of the set X and the probability mass in the empirical distribution derived from the fingerprint of frequently occurring elements add up to 1. Then it is shown that such normalization has a small effect on Expression 2.10. Therefore, \hat{v} is in the feasible region with high probability. Intuitively, the associated histogram $h^{\hat{v}}$ is close to h with respect to relative earthmover distance, since for the low probability portion of the distribution, the two histograms are highly similar by construction, and for the high probability portion of the distribution, $h^{\hat{v}}\left(\frac{i}{m}\right) = \mathcal{F}_i$ is close to $h\left(\frac{i}{m}\right)$ because each element of the fingerprint, with high probability, has true probability close to its observed probability.

In the second part of the proof of Theorem 2.15, it is shown that for any two solutions \mathbf{v}, \mathbf{w} the associated histograms $h^{\mathbf{v}}$ and $h^{\mathbf{w}}$ are close. We start with key definitions.

Definition 2.16. For a given m , a β -bump earthmoving scheme is defined by a sequence of positive real numbers $\{c_i\}$, the bump centers, and a sequence of functions $\{f_i\} : (0, 1] \rightarrow \mathbb{R}$ such that $\sum_{i=0}^{\infty} f_i(x) = 1$ and for each x , and each function f_i may be expressed as a linear combination of Poisson functions $f_i(x) = \sum_{j=0}^{\infty} a_{ij} \text{Pois}(mx, j)$ such that $\sum_{j=0}^{\infty} |a_{ij}| \leq \beta$. Given a generalized histogram h , the scheme works as follows: for each x such that $h(x) \neq 0$, and each integer $i > 0$, move $xh(x) \cdot f_i(x)$ probability mass from x to c_i . We denote the histogram resulting from this scheme by $(c, f)(h)$.

Definition 2.17. For given n, m , a bump earthmoving scheme (c, f) is ϵ -good if for any generalized histogram h , the relative earthmover distance between h and $(c, f)(h)$ is at most ϵ .

As these definitions hint, it is sufficient to construct a “good” bump earthmoving scheme such that, given two solutions \mathbf{v}, \mathbf{w} , when the scheme is applied to the associated histograms $h^{\mathbf{v}}, h^{\mathbf{w}}$, the resulting histograms $(c, f)(h^{\mathbf{v}})$ and $(c, f)(h^{\mathbf{w}})$ have small distance. Then it immediately follows that $(c, f)(h^{\mathbf{v}})$ and $(c, f)(h^{\mathbf{w}})$ have similar entropies, leading to the same conclusion about the entropies of $h^{\mathbf{v}}$ and $h^{\mathbf{w}}$ which ends the proof.

Lemma 2.18. ([20], Lemma 16) For $n > m$, letting ς be such that $n = \varsigma m \log m$, there exists an $O(\sqrt{\varsigma} \cdot \max\{1, |\log \varsigma|\})$ -good $m^{0.3}$ -bump earthmoving scheme.

We only present the general idea of the proof due to its laborious details which may be tiresome for a reader to follow. The building block of the construction of a “good” bump earthmoving scheme is to use two different classes of functions for $\{f_i\}$. For $i \geq \log m$ Poisson functions $\text{Pois}(mx, i)$ are utilized. For $i < \log m$ Chebyshev polynomials $T_i(x)$ are employed where the i^{th} Chebyshev polynomial is the polynomial of degree i such that $T_i(\cos y) = \cos(i \cdot y)$. The sequence $\{c_i\}$ is assigned to be $\{\frac{i}{m}\}$ for

$i \geq \log m$, and $\left\{ \frac{2 \log m}{5m} \left(1 - \cos \left(\frac{5(i+1)\pi}{\log m} \right) \right) \right\}$ for $i < \log m$. Then by leveraging the fact that any two solutions \mathbf{v}, \mathbf{w} must have close fingerprint expectations, it is shown that the resulting histograms $(c, f)(h^{\mathbf{v}})$ and $(c, f)(h^{\mathbf{w}})$ have small earthmoving distance.

2.1.2. Lower Bound

In the second part, we introduce the lower bound established in [19] on estimating entropy. Valiant and Valiant prove that the task of estimating entropy up to an additive error has sample complexity $\Omega\left(\frac{n}{\log n}\right)$ by constructing two probability distributions \mathbf{p}^+ and \mathbf{p}^- with large relative earthmover distance yet having close fingerprint distributions. In other words, \mathbf{p}^+ and \mathbf{p}^- are built such that the difference, $|H(\mathbf{p}^+) - H(\mathbf{p}^-)|$, is big enough to distinguish \mathbf{p}^+ from \mathbf{p}^- , however, no algorithm can differentiate between the fingerprint obtained from \mathbf{p}^+ and the fingerprint obtained from \mathbf{p}^- by drawing $o\left(\frac{n}{\log n}\right)$ samples. In addition, Valiant and Valiant provide two new central limit theorems (CLT) for establishing the lower bound. We choose to avoid the details of the proof of the central limit theorems in order to preserve the smoothness of the explanation.

Definition 2.19. *The generalized multinomial distribution parameterized by a non-negative matrix ρ each of whose rows sum to at most 1, is denoted M^ρ , and is defined by the following random process: for each row $\rho(i, \cdot)$ of matrix ρ , interpret it as a probability distribution over the columns of ρ — including, if $\sum_{j=1}^m \rho(i, j) < 1$, an “invisible” column 0 — and draw a column index from this distribution; return a row vector recording the total number of samples falling into each column (the histogram of the samples).*

The generalized multinomial distribution is employed to capture the fingerprint distribution of a probability distribution. We introduce the first central limit theorem that relates the sum of independent distributions to a Gaussian distribution with respect to earthmover distance.

Theorem 2.20. ([19], Theorem 2) *Given n independent distributions $\{Z_i\}$ of mean 0 in \mathbb{R}^m and a bound β such that $\|Z_i\| < \beta$ for any i and any sample, the earthmover*

distance between $\sum_{i=1}^n Z_i$ and the normal distribution of corresponding mean (0) and covariance is at most $\beta m(2.7 + 0.83 \log n)$.

The second central limit theorem approximates a generalized multinomial distribution by a Gaussian distribution with respect to statistical distance. Note that the statistical distance between a generalized multinomial distribution and a Gaussian distribution is 1 since the first is discrete but the second is continuous. Therefore, Gaussian distribution is discretized by rounding to the nearest lattice points.

Definition 2.21. *The m -dimensional discretized Gaussian distribution, with mean μ and covariance matrix Σ , denoted $\mathcal{N}^{disc}(\mu, \Sigma)$, is the distribution with support \mathbb{Z}^m obtained by picking a sample according to the Gaussian $\mathcal{N}(\mu, \Sigma)$, then rounding each coordinate to the nearest integer.*

Theorem 2.22. ([19], Theorem 4) *Given a generalized multinomial distribution M^ρ , with m dimensions and n rows, let μ denote its mean and Σ denote its covariance matrix, then*

$$d_{\text{TV}}(M_\rho, \mathcal{N}^{disc}(\mu, \Sigma)) \leq \frac{m^{4/3}}{\sigma^{1/3}} \cdot 2.2 \cdot (3.1 + 0.83 \log n)^{2/3},$$

where σ^2 is the minimum eigenvalue of Σ .

We present how distributions \mathbf{p}^+ and \mathbf{p}^- are constructed.

Definition 2.23. *Given a real number $\phi \in (0, \frac{1}{4})$, consider the degree $\log m + 2$ polynomial $M_{\log m, \phi}(x) := -\left(x - \phi \frac{1}{\log m}\right) \left(x - 2\phi \frac{1}{\log m}\right) L_{\log m}(x)$ such that $L_j(x) = \frac{e^x}{j!} \frac{d^j}{dx^j} (e^{-x} x^j)$ is the j^{th} Laguerre polynomial. Let $v(x)$ be the function that takes value $1/M'_{\log m, \phi}(x)$ for every x where $M_{\log m, \phi}(x) = 0$, and is 0 otherwise, where M' is the derivative of M . Define the distributions $\mathbf{p}_{\log m, \phi}^+, \mathbf{p}_{\log m, \phi}^-$ such that for each x where $v(x) > 0$, the distribution $\mathbf{p}_{\log m, \phi}^+$ contains $v(x)e^{x/32}$ probability mass at probability $\frac{1}{32m}x$, and for each x where $v(x) < 0$ the distribution $\mathbf{p}_{\log m, \phi}^-$ contains $|v(x)|e^{x/32}$ probability mass at probability $\frac{1}{32m}x$, where each distribution is then normalized to have total probability mass 1.*

Observe that since the probability of each element of either $\mathbf{p}_{\log m, \phi}^+$ or $\mathbf{p}_{\log m, \phi}^-$ is defined to be at least $\frac{\phi}{32m \log m}$, both distributions have support at most $\frac{32}{\phi}m \log m$. Thus, the connection between the sample and domain sizes necessary for the lower bound is formed. The second condition on $\mathbf{p}_{\log m, \phi}^+, \mathbf{p}_{\log m, \phi}^-$ is that the difference $|H(\mathbf{p}_{\log m, \phi}^+) - H(\mathbf{p}_{\log m, \phi}^-)|$ is sufficiently big. This can be achieved by proving that $\mathbf{p}_{\log m, \phi}^+, \mathbf{p}_{\log m, \phi}^-$ are close in the relative earthmover distance to two different distributions $\mathbf{q}^+, \mathbf{q}^-$, respectively, such that $H(\mathbf{q}^+)$ is distant from $H(\mathbf{q}^-)$.

Lemma 2.24. ([19], Lemma 13) *Distributions $\mathbf{p}_{\log m, \phi}^+, \mathbf{p}_{\log m, \phi}^-$ are $O(\phi |\log \phi|)$ -close, respectively, in the relative earthmover distance to the uniform distributions on $\frac{32}{\phi}m \log m$ and $\frac{16}{\phi}m \log m$ elements.*

Note that the relative earthmover distance, therefore by Fact 2.10, the difference of the entropies is $H(\mathcal{U}(\lceil \frac{32}{\phi}m \log m \rceil)) - H(\mathcal{U}(\lceil \frac{16}{\phi}m \log m \rceil)) = 1$, where $\mathcal{U}([n])$ denotes the uniform distribution with support size equal to n . That is, $\mathbf{p}_{\log m, \phi}^+$ and $\mathbf{p}_{\log m, \phi}^-$ are distinguishable. The third condition is that $\mathbf{p}_{\log m, \phi}^+, \mathbf{p}_{\log m, \phi}^-$ have fingerprint distributions $\mathcal{F}_{\mathbf{p}^+}, \mathcal{F}_{\mathbf{p}^-}$, respectively, such that the statistical distance between $\mathcal{F}_{\mathbf{p}^+}$ and $\mathcal{F}_{\mathbf{p}^-}$ is small. Similarly, it is shown that $\mathcal{F}_{\mathbf{p}^+}, \mathcal{F}_{\mathbf{p}^-}$ are approximated by two statistically close, discretized Gaussian distributions $\mathcal{N}_+^{disc}, \mathcal{N}_-^{disc}$, respectively. First, we state a weaker result.

Lemma 2.25. ([19], Lemma 16) *For any i , the i^{th} fingerprint expectations for distributions $\mathbf{p}_{\log m, \phi}^+, \mathbf{p}_{\log m, \phi}^-$ are equal to within $o(1)$.*

It is necessary for $\mathbf{p}_{\log m, \phi}^+, \mathbf{p}_{\log m, \phi}^-$ to have bigger variance in both directions due to obtaining better bounds when aforementioned central limit theorems are applied. Therefore, $\mathbf{p}_{\log m, \phi}^+, \mathbf{p}_{\log m, \phi}^-$ are modified to get “fat” distributions $\mathbf{p}_{\log m, \phi}^{F+}, \mathbf{p}_{\log m, \phi}^{F-}$ such that both are statistically close to their “thin” counterparts, respectively. A “fat” distribution is constructed as follows:

Definition 2.26. *Define the fattening operator F that, given a histogram h , constructs a new histogram h^F as follows:*

- Provisionally set $h^F(x) = \left(1 - \frac{\log m - 1}{2 \log^2 m}\right) h(x)$ for each x ;
- For each integer $i \in \{1, \dots, \log m\}$, increment $h^F\left(\frac{i}{m}\right) \leftarrow h^F\left(\frac{i}{m}\right) + \frac{m}{\log^3 m}$.

Note that operator F returns a proper probability distribution and preserves the previous upper bound on support size since no element with probability less than $1/m$ is added to the support. Intuitively, it also “fattens” a distribution because the number of low probability elements is increased substantially. Moreover, operator F does not negatively affect the bounds of Lemma 2.25 since both distributions $\mathbf{p}_{\log m, \phi}^+$, $\mathbf{p}_{\log m, \phi}^-$ are modified identically. We are ready to state the main result.

Proposition 2.27. ([19], Proposition 21) *For a positive constant $\phi < 1/4$, the statistical distance between the distribution of $\text{Pois}(m)$ –sample fingerprints from $\mathbf{p}_{\log m, \phi}^+$ and $\mathbf{p}_{\log m, \phi}^-$ goes to 0 as m goes to infinity.*

Recall that Theorem 2.22 is utilized to approximate the fingerprint distributions of $\mathbf{p}_{\log m, \phi}^+$, $\mathbf{p}_{\log m, \phi}^-$ by two statistically close, discretized Gaussian distributions, respectively. We only explain the necessary conditions for Theorem 2.22, skipping the details of its application.

The condition to be satisfied is that fingerprint distributions $\mathcal{F}_{\mathbf{p}_{F^+}}$, $\mathcal{F}_{\mathbf{p}_{F^-}}$ obtained from distributions $\mathbf{p}_{\log m, \phi}^+$, $\mathbf{p}_{\log m, \phi}^-$, respectively, have close variance and covariance. This is proven by using the result in Lemma 2.25, that is, fingerprint distributions having close expectations. Recall that for a histogram h , expectation of the i^{th} fingerprint entry is $\mathbf{E}[\mathcal{F}_i] = \sum_{x:h(x) \neq 0} h(x) \cdot \text{Pois}(mx, i)$, and covariance of two random variables X, Y is defined as $\mathbf{Cov}[X, Y] = \mathbf{E}[XY] - \mathbf{E}[X]\mathbf{E}[Y]$. Combining with Equation 2.9, it follows that covariance of the i^{th} and j^{th} fingerprint entries, for $i \neq j$, equals $\mathbf{Cov}[\mathcal{F}_i, \mathcal{F}_j] = \sum_{x:h(x) \neq 0} -h(x) \cdot \text{Pois}(xm, i) \text{Pois}(xm, j)$. After the simplification,

$$\text{Pois}(xm, i) \text{Pois}(xm, j) = \frac{(xm)^{i+j} e^{-2xm}}{i!j!} = 2^{-(i+j)} \binom{i+j}{i} \text{Pois}(2xm, i+j).$$

Recall that variance of a random variable X is $\mathbf{Var}[X] = \mathbf{E}[X^2] - \mathbf{E}^2[X]$. Then for the i^{th} fingerprint entry, $\mathbf{Var}[\mathcal{F}_i] = \sum_{x:h(x) \neq 0} h(x) \cdot (\text{Pois}(mx, i) - \text{Pois}(mx, i)^2)$. Note that $\text{Pois}(mx, i)^2 = 2^{-2i} \binom{2i}{i} \cdot \text{Pois}(2mx, 2i)$. The following relates Poisson functions with different parameters.

Lemma 2.28. ([19], Lemma 20) *For any $\epsilon > 0$ and integer $i \geq 0$, one may approximate $\text{Pois}(2x, i)$ as a linear combination $\sum_{j=0}^{\infty} \alpha(j) \text{Pois}(x, j)$ such that*

1. *For all $x \geq 0$, $|\text{Pois}(2x, i) - \sum_{j=0}^{\infty} \alpha(j) \text{Pois}(x, j)| \leq \epsilon$; and*
2. *$\sum_{j=0}^{\infty} |\alpha(j)| \leq \frac{1}{\epsilon} \cdot 200 \max\left\{\sqrt[4]{i}, 24 \log^{3/2} \frac{1}{\epsilon}\right\}$.*

Therefore, both variance and covariance of fingerprint distributions of $\mathbf{p}_{\log m, \phi}^{F+}$ and $\mathbf{p}_{\log m, \phi}^{F-}$ can be expressed as linear combinations of Poisson functions. It is already known that fingerprint distributions of $\mathbf{p}_{\log m, \phi}^{F+}$, $\mathbf{p}_{\log m, \phi}^{F-}$ have close expectations which themselves are expressed as linear combinations of Poisson functions. Then, at least intuitively, fingerprint distributions of $\mathbf{p}_{\log m, \phi}^{F+}$, $\mathbf{p}_{\log m, \phi}^{F-}$ have similar variance and covariance.

Theorem 2.29. ([19], Theorem 1) *For any positive constant $\phi < 1/4$ there exists a pair of distributions $\mathbf{p}^+, \mathbf{p}^-$ that are $O(\phi |\log \phi|)$ – close in the relative earthmover distance, respectively, to the uniform distributions on n and $n/2$ elements, but which are indistinguishable to $m = \frac{\phi}{32} \cdot \frac{n}{\log n}$ – sample testers.*

2.1.3. Upper Bound II

In the final part of this section, we present an improved upper bound on estimating entropy up to an additive Δ . The bounds given in [19, 20] are matching in their dependence on n , whereas for the dependence on Δ this is not the case. The estimator constructed in [20] has sample complexity $O\left(\frac{n}{\Delta^2 \log n}\right)$, while the lower bound established in [19] is $\Omega\left(\frac{n}{\Delta \log n}\right)$, which leaves open the question of error decrease rate. The problem is resolved in [21] by constructing an optimal estimator which estimates the entropy of a distribution to within additive accuracy Δ , with probability at least $1 - o(\text{poly}(n))$, given $O\left(\frac{n}{\Delta \log n}\right)$ independent samples from a distribution with sup-

port size at most n . Differently than the previous one which is based on a canonical approach, Valiant and Valiant construct an estimator focusing directly on entropy.

Definition 2.30. A symmetric property π is **linear** if there exists some function $f_\pi : [0, 1] \rightarrow \mathbb{R}$, denoted as **characteristic function** of π , such that for any distribution \mathbf{p} with histogram $h_{\mathbf{p}}$, $\pi(\mathbf{p}) = \sum_{x:h_{\mathbf{p}}(x) \neq 0} h_{\mathbf{p}}(x) f_\pi(x)$.

Observe that Shannon entropy is a linear property, and its characteristic function is $f(x) = x|\log x|$. The new estimator is based on approximating the characteristic function of entropy as a linear combination of Poisson functions. The following clarifies the motivation behind this approach. Assume there exists a sequence of coefficients $\{\beta_i\}$ such that for all $x \in (0, 1]$, $\sum_{i=1}^{\infty} \beta_i \text{Pois}(mx, i) = x|\log x| = f(x)$. Then,

$$\begin{aligned} \sum_{x:h(x) \neq 0} h(x)f(x) &= \sum_{x:h(x) \neq 0} h(x) \sum_{i \geq 1} \beta_i \text{Pois}(mx, i) \\ &= \sum_{i \geq 1} \beta_i \sum_{x:h(x) \neq 0} h(x) \text{Pois}(xm, i) \\ &= \sum_{i \geq 1} \beta_i \mathbf{E}[\mathcal{F}_i] = \mathbf{E} \left[\sum_{i \geq 1} \beta_i \mathcal{F}_i \right]. \end{aligned} \tag{2.11}$$

That is, the quantity $\sum_{i \geq 1} \beta_i \mathcal{F}_i$ is an unbiased estimator for entropy. Recall that for each i , a fingerprint entry \mathcal{F}_i is concentrated around its expectation. Then, roughly, for $\sum_{i \geq 1} \beta_i \mathcal{F}_i$ having relatively small variance one needs the coefficients $\{\beta_i\}$ to be small comparing to $1/\sqrt{m}$. However, instead of approximating the characteristic function $f(x) = x|\log x|$ directly, the function $\frac{f(x)}{x} = |\log x|$ is expressed as a linear combination of Poisson functions, $\sum_{i=0}^{\infty} z_i \text{Pois}(mx, i)$. Observe that these approaches are equivalent in the sense that $\beta_i = \frac{i}{m} \cdot z_{i-1}$, since $x \text{Pois}(mx, i) = \text{Pois}(mx, i+1) \frac{i+1}{m}$. The following formalizes the relationship between the magnitudes of coefficients, error in approximating $|\log x|$ and the estimator defined above.

Proposition 2.31. ([21], Proposition 17) *Given integers m, n , and a set of coefficients z_0, z_1, \dots such that if for positive real numbers a, b, c the following conditions hold:*

1. $\left| |\log x| - \sum_{i=0}^{\infty} z_i \text{Pois}(mx, i) \right| < a + \frac{b}{x},$

2. for all $j \geq 1$ let $\beta_j = \frac{j}{m} z_{j-1}$ with $\beta_0 = 0$, then for any j, l such that $|j - l| \leq \sqrt{j} \log m$ we have $|\beta_j - \beta_l| \leq c \sqrt{\frac{j}{m}}$.

Then the estimator described in Equation 2.11 estimates entropy with error at most $a + bn + c \log m$, with probability at least $1 - o(1/\text{poly}(m))$ when given a fingerprint derived from a set of m independent samples chosen from a distribution with support size at most n .

The task of finding such coefficients $\{z_i\}$ is realized via linear programming. A linear program is constructed with constraints describing the conditions of Proposition 2.31 and with the objective function minimizing error in the estimation.

Algorithm: ESTIMATOR II

- Fix $m = O\left(\frac{n}{\Delta \log n}\right)$. (For details see [21].)
- Draw $M \sim \text{Pois}(m)$ independent samples X_1, \dots, X_M from \mathbf{p} .
- Construct the linear program as in Definition 18 [21], corresponding to \mathcal{F} .
- Find a solution $\{z_i\}$ to the the linear program.
- Calculate the coefficients $\{\beta_i\}$.
- Output $\sum_{i \geq 1} \beta_i \mathcal{F}_i$.

Figure 2.2. Linear Estimator of Shannon Entropy.

Aside from employing linear programming to find convenient coefficients $\{z_i\}$, Valiant and Valiant explicitly construct an optimal estimator such that given $O\left(\frac{n}{\Delta \log n}\right)$ independent samples from a distribution with support size at most n , it estimates entropy of a distribution to within additive accuracy Δ , with probability at least $1 - o(\text{poly}(n))$. We briefly describe how an optimal estimator is constructed. Recall that the objective is to approximate the function $\log x$ as a linear combination of Poisson functions, $\sum_{i=0}^{\infty} z_i \text{Pois}(mx, i)$. A straightforward choice for coefficients $\{z_i\}$ is the sequence $\{\log \frac{i}{m}\}$. Note that $\log \frac{i}{m}$ is a “plug-in” estimator for entropy. The following lemma bounds the precision of any “plug-in” estimator.

Lemma 2.32. ([21], Lemma 19) *Given a function $f : \mathbb{R} \rightarrow \mathbb{R}$ whose fourth derivative at x is bounded in magnitude by $\frac{\alpha}{x^4}$ for $x \geq 1$ and by α for $x \leq 1$, and whose third derivative at x is bounded by $\frac{\alpha}{x^3}$, then for any real x , $\sum_{i=0}^{\infty} f(i) \text{Pois}(x, i)$ is within $O\left(\frac{\alpha}{x^2}\right)$ of $f(x) + \frac{1}{2}xf''(x)$.*

For a “plug-in” estimator $\log \frac{i}{m}$, this lemma suggests that

$$\log x - \sum_{i=0}^{\infty} \log(i/m) \text{Pois}(mx, i) = \frac{-1}{2 \ln 2 \cdot mx} + O\left(\frac{1}{m^2 x^2}\right). \quad (2.12)$$

Observe that for a high probability, error of approximation is small, whereas for a low probability such as $x \leq \frac{1}{m}$, the right-hand side of Equation 2.12 becomes unbounded. In other words, “plug-in” estimator is satisfactory for high probability portion of a distribution, however, it behaves poorly for low probability portion of a distribution. Similar techniques are used as in bump earthmoving scheme described in Subsection 2.1.1 to resolve the issue. We only give an outline of the proof due to its laborious details.

Two different functions are used for approximating $\log x$ as a linear combination of Poisson functions. For probabilities $x \geq O(\log m)$ “plug-in” estimator $\log \frac{i}{m}$ and for probabilities $x < O(\log m)$ a function of Chebyshev polynomials referred as the Chebyshev bumps are utilized. The next step is to establish a Chebyshev bump version of Lemma 2.32 and to show that the Chebyshev bumps can be expressed as a linear combination of Poisson functions with relatively small coefficients. The proof ends by applying Proposition 2.31.

2.2. Rényi Entropy

We investigate Rényi entropy, first introduced by Alfréd Rényi [25] which is a popular generalization of Shannon entropy. It is defined as follows:

Definition 2.33. *Let $\alpha \geq 0$ be a real number. The Rényi entropy of order α of a distribution \mathbf{p} , denoted by $H_{\alpha}(\mathbf{p})$, is*

for $\alpha \neq 1$

$$H_\alpha(\mathbf{p}) = \frac{1}{1-\alpha} \log \left(\sum_x \mathbf{p}[x]^\alpha \right), \quad (2.13)$$

and for $\alpha = 1$

$$H_1(\mathbf{p}) = \lim_{\alpha \rightarrow 1} H_\alpha(\mathbf{p}) \quad (2.14)$$

For $\alpha = 0$, $H_0(\mathbf{p}) = \log \text{supp}(\mathbf{p})$, where $\text{supp}(\mathbf{p})$ denotes the support size of the distribution \mathbf{p} . For $\alpha = 1$, $H_1(\mathbf{p}) = H(\mathbf{p})$, that is, Rényi entropy becomes Shannon entropy, as easily derived via L'Hôpital's rule. For $\alpha = \infty$, $H_\alpha(\mathbf{p})$ is the min-entropy $H_\infty(\mathbf{p})$, where by definition $H_\infty(\mathbf{p}) = -\log \max_i \mathbf{p}[i]$.

Rényi entropy has many applications. Particularly, $H_2(\mathbf{p})$ is used for measuring the quality of random number generators [26], for testing the closeness of probability distributions [27, 28], for characterizing the number of reads needed to reconstruct a DNA sequence [29], etc.

Recall that $\Theta\left(\frac{n}{\log n}\right)$ samples are necessary and sufficient for estimating Shannon entropy, which are only better by a polylogarithmic factor than $\Theta(n \log^2 n)$, a trivial upper bound for this task. Thus, being a generalization of Shannon entropy, determining the complexity of estimating Rényi entropy becomes additionally intriguing. Acharya, Orlitsky, Suresh, and Tyagi [7] provide near-optimal upper and lower bounds for three different cases of α ; it is shown that to estimate Rényi entropy to within an additive error one requires (i) for $\alpha < 1$, super-linear, roughly $n^{1/\alpha}$ samples (ii) for noninteger $\alpha > 1$, near-linear, roughly n samples (iii) for integer $\alpha > 1$, sub-linear, $\Theta(n^{1-1/\alpha})$ samples. Note that in the case of $\alpha > 1$ being integer, estimating Rényi entropy becomes substantially easier than estimating Shannon entropy.

2.2.1. Upper Bounds

In the first part we illustrate the upper bounds established in [7]. Defining the α^{th} moment of \mathbf{p} as $\mathcal{M}_\alpha(\mathbf{p}) = \sum_{i=1}^n (\mathbf{p}[i])^\alpha$, Rényi entropy can be expressed as $H_\alpha(\mathbf{p}) = \frac{1}{1-\alpha} \log \mathcal{M}_\alpha(\mathbf{p})$. Observe that estimating $H_\alpha(\mathbf{p})$ to an additive accuracy of $\pm\Delta$ is equivalent to estimating $\mathcal{M}_\alpha(\mathbf{p})$ to a multiplicative accuracy of $2^{\pm\Delta(1-\alpha)}$. Acharya *et al.* construct two different multiplicative estimators, denoted by $\widehat{\mathcal{M}}_\alpha^e$ and $\widehat{\mathcal{M}}_\alpha^u$, the first for $\alpha \notin \mathbb{Z}$ and the latter for $\alpha \in \mathbb{Z}$.

To simplify the analysis, Poisson sampling technique is utilized as in the case of estimating Shannon entropy whose explicit description is given in Section 2.1. Recall that in the $\text{Pois}(m)$ -sampling scheme N_i/m is an unbiased estimator for $\mathbf{p}[i]$ where N_i denotes the multiplicity of an element i . We define the aforementioned estimators $\widehat{\mathcal{M}}_\alpha^e, \widehat{\mathcal{M}}_\alpha^u$.

Definition 2.34. For $\alpha \notin \mathbb{Z}$, let the empirical estimator for $\mathcal{M}_\alpha(\mathbf{p})$, denoted by $\widehat{\mathcal{M}}_\alpha^e$, be

$$\widehat{\mathcal{M}}_\alpha^e = \sum_{i \in [n]} \left(\frac{N_i}{m} \right)^\alpha. \quad (2.15)$$

Observe that $\widehat{\mathcal{M}}_\alpha^e$ is biased. For $\alpha \in \mathbb{Z}^+$, let $n^\alpha = n \cdot (n-1) \cdots (n-\alpha+1)$ denote the α^{th} falling power of n .

Definition 2.35. For integer $\alpha > 1$, let the bias-corrected estimator for $\mathcal{M}_\alpha(\mathbf{p})$, denoted by $\widehat{\mathcal{M}}_\alpha^u$, be

$$\widehat{\mathcal{M}}_\alpha^u = \sum_{i \in [n]} \frac{N_i^\alpha}{m^\alpha}. \quad (2.16)$$

We demonstrate the central lemma for establishing upper bounds, which essentially constructs an error reduction algorithm to increase the accuracy of an estimator given certain bounds on its bias and variance.

Lemma 2.36. ([7], Lemma 6) *For $M \sim \text{Pois}(m)$, let the estimator $\widehat{\mathcal{M}}_\alpha$ have bias and variance satisfying*

$$\begin{aligned} \left| \mathbf{E} \left[\widehat{\mathcal{M}}_\alpha \right] - \mathcal{M}_\alpha(\mathbf{p}) \right| &\leq \frac{\gamma}{2} \mathcal{M}_\alpha(\mathbf{p}), \\ \mathbf{Var} \left[\widehat{\mathcal{M}}_\alpha \right] &\leq \frac{\gamma^2}{12} \mathcal{M}_\alpha(\mathbf{p})^2. \end{aligned}$$

Then, there exists an estimator $\widehat{\mathcal{M}}'_\alpha$ that uses $O(m \log(1/\delta))$ samples, and ensures

$$\Pr \left[\left| \widehat{\mathcal{M}}'_\alpha - \mathcal{M}_\alpha(\mathbf{p}) \right| > \gamma \mathcal{M}_\alpha(\mathbf{p}) \right] \leq \delta.$$

Proof. By Chebyshev's inequality,

$$\begin{aligned} &\Pr \left[\left| \widehat{\mathcal{M}}_\alpha - \mathcal{M}_\alpha(\mathbf{p}) \right| > \gamma \mathcal{M}_\alpha(\mathbf{p}) \right] \\ &\leq \Pr \left[\left| \widehat{\mathcal{M}}_\alpha - \mathbf{E} \left[\widehat{\mathcal{M}}_\alpha \right] \right| + \left| \mathbf{E} \left[\widehat{\mathcal{M}}_\alpha \right] - \mathcal{M}_\alpha(\mathbf{p}) \right| > \gamma \mathcal{M}_\alpha(\mathbf{p}) \right] \\ &\leq \Pr \left[\left| \widehat{\mathcal{M}}_\alpha - \mathbf{E} \left[\widehat{\mathcal{M}}_\alpha \right] \right| > \frac{\gamma}{2} \mathcal{M}_\alpha(\mathbf{p}) \right] \\ &\leq \frac{4 \mathbf{Var} \left[\widehat{\mathcal{M}}_\alpha \right]}{\gamma^2 \widehat{\mathcal{M}}_\alpha^2} \leq \frac{1}{3}. \end{aligned}$$

$\mathcal{M}_\alpha(\mathbf{p})$ is estimated in t independent rounds, and $\widehat{\mathcal{M}}'_\alpha$ is assigned to be the sample median of these rounds. More specifically, let $\widehat{\mathcal{M}}_i$ denote the result of round i , and let $\mathbb{1}_{E_i}$ be the indicator function of the event $E_i = \left\{ \left| \widehat{\mathcal{M}}_i - \mathcal{M}_\alpha(\mathbf{p}) \right| > \gamma \mathcal{M}_\alpha(\mathbf{p}) \right\}$. Then the expectation satisfies $\mathbf{E}[\mathbb{1}_{E_i}] \leq 1/3$, and by the Hoeffding bound,

$$\Pr \left[\sum_{i=1}^t \mathbb{1}_{E_i} > \frac{t}{2} \right] \leq e^{-t/18}.$$

To reduce the probability of error to δ , set $t = 18 \log(1/\delta)$. What this means is that with probability at least $1 - \delta$, the majority of the rounds, therefore, the sample median $\widehat{\mathcal{M}}'_\alpha$ satisfies the condition $\left| \widehat{\mathcal{M}}'_\alpha - \mathcal{M}_\alpha(\mathbf{p}) \right| \leq \gamma \mathcal{M}_\alpha(\mathbf{p})$. \square

The next step is to bound the bias and variance of both estimators, $\widehat{\mathcal{M}}_\alpha^e$, $\widehat{\mathcal{M}}_\alpha^u$, in order to apply Lemma 2.36, which finalizes the proof. Note that independence gained due to Poisson sampling simplifies the analysis in bounding variance. We only state related results regarding the bias-corrected estimator.

Lemma 2.37. ([7], Lemma 2) *Let $X \sim \text{Pois}(\lambda)$. Then, for all $r \in \mathbb{N}$*

$$\begin{aligned}\mathbf{E}[X^r] &= \lambda^r, \\ \mathbf{Var}[X^r] &\leq \lambda^r ((\lambda + r)^r - \lambda^r).\end{aligned}$$

Proof. The expected value is

$$\mathbf{E}[X^r] = \sum_{i=0}^{\infty} \text{Pois}(\lambda, i) \cdot i^r = \sum_{i=r}^{\infty} e^{-\lambda} \frac{\lambda^i}{i!} \cdot \frac{i!}{(i-r)!} = \lambda^r \sum_{i=0}^{\infty} e^{-\lambda} \frac{\lambda^i}{i!} = \lambda^r.$$

We bound the following

$$\begin{aligned}\mathbf{E}[(X^r)^2] &= \sum_{i=0}^{\infty} \text{Pois}(\lambda, i) \cdot (i^r)^2 \\ &= \sum_{i=r}^{\infty} e^{-\lambda} \frac{\lambda^i}{i!} \cdot \frac{(i!)^2}{(i-r)!^2} \\ &= \lambda^r \sum_{i=0}^{\infty} e^{-\lambda} \frac{\lambda^i}{i!} (i+r)^r \\ &= \lambda^r \mathbf{E}[(X+r)^r] \\ &= \lambda^r \mathbf{E} \left[\prod_{j=1}^r [(X+1-j) + r] \right] \\ &\leq \lambda^r \mathbf{E} \left[\sum_{j=0}^r \binom{r}{j} X^j \cdot r^{r-j} \right] \\ &= \lambda^r \sum_{j=0}^r \binom{r}{j} \lambda^j \cdot r^{r-j} \\ &= \lambda^r (\lambda + r)^r.\end{aligned}$$

Hence,

$$\mathbf{Var}[X^r] = \mathbf{E}[(X^r)^2] - \mathbf{E}[X^r]^2 \leq \lambda^r ((\lambda + r)^r - \lambda^r). \quad \square$$

Finally, we state the upper bound results for estimating Rényi entropy.

Theorem 2.38. ([7], Theorem 9) *For an integer $\alpha > 1$, any $\Delta > 0$, and $0 < \delta < 1$, there exists an algorithm estimating with probability at least $1 - \delta$ the Rényi entropy $H_\alpha(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$ with $O\left(\frac{n^{1-1/\alpha}}{(1-2^{(1-\alpha)\Delta})^2} \log \frac{1}{\delta}\right)$ samples.*

Algorithm: ESTIMATOR III

- Fix $\gamma = 1 - 2^{(1-\alpha)\Delta}$ and $m = O\left(\frac{n^{1-1/\alpha}}{\gamma^2}\right)$. (For details see [7].)
- Repeat the following for $t = \lceil 18 \log \frac{1}{\delta} \rceil$ independent rounds.
 - Draw $M \sim \text{Pois}(m)$ independent samples X_1, \dots, X_M from \mathbf{p} .
 - Compute the multiplicity N_i based on the samples X_1, \dots, X_M for $1 \leq i \leq n$.
 - Set $\widehat{\mathcal{M}}_j = \sum_{i \in [n]} \frac{N_i^\alpha}{m^\alpha}$ for round $1 \leq j \leq t$.
- Output $\frac{1}{1-\alpha} \log \widehat{\mathcal{M}}_\alpha$ where $\widehat{\mathcal{M}}_\alpha$ is the median of the sequence $\{\widehat{\mathcal{M}}_j\}$.

Figure 2.3. SAMP Estimator of Rényi Entropy (of integer degree $\alpha > 1$)

Theorem 2.39. ([7], Theorem 7) *For $\alpha > 1$, $\Delta > 0$, and $0 < \delta < 1$, there exists an algorithm estimating with probability at least $1 - \delta$ the Rényi entropy $H_\alpha(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$ with $O\left(\frac{n}{\gamma^{\max\{4, 1/(\alpha-1)\}}} \log \frac{1}{\delta}\right)$ samples where $\gamma = 1 - 2^{(1-\alpha)\Delta}$.*

Theorem 2.40. ([7], Theorem 8) *For $\alpha < 1$, $\Delta > 0$, and $0 < \delta < 1$, there exists an algorithm estimating with probability at least $1 - \delta$ the Rényi entropy $H_\alpha(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$ with $O\left(\frac{n^{1/\alpha}}{\gamma^{\max\{4, 2/\alpha\}}} \log \frac{1}{\delta}\right)$ samples where $\gamma = 1 - 2^{(\alpha-1)\Delta}$.*

Algorithm: ESTIMATOR IV

- Fix $\gamma = 1 - 2^{(1-\alpha)\Delta}$ and $m = O\left(\frac{n}{\gamma^{\max\{4, 1/(\alpha-1)\}}}\right)$. (For details see [7].)
- Repeat the following for $t = \lceil 18 \log \frac{1}{\delta} \rceil$ independent rounds.
 - Draw $M \sim \text{Pois}(m)$ independent samples X_1, \dots, X_M from \mathbf{p} .
 - Compute the multiplicity N_i based on the samples X_1, \dots, X_M for $1 \leq i \leq n$.
 - Set $\widehat{\mathcal{M}}_j = \sum_{i \in [n]} \binom{N_i}{m}^\alpha$ for round $1 \leq j \leq t$.
- Output $\frac{1}{1-\alpha} \log \widehat{\mathcal{M}}_\alpha$ where $\widehat{\mathcal{M}}_\alpha$ is the median of the sequence $\{\widehat{\mathcal{M}}_j\}$.

Figure 2.4. SAMP Estimator of Rényi Entropy (of noninteger degree $\alpha > 1$)

We skip the algorithmic view of additively estimating the Rényi entropy for $\alpha < 1$, since it is easily obtained by setting $m = O\left(\frac{n^{1/\alpha}}{(1-2^{(\alpha-1)\Delta})^{\max\{4, 2/\alpha\}}}\right)$ in Figure 2.4.

Very recently, Acharya, Orlitsky, Suresh, and Tyagi [30] have improved the upper bound results for additively estimating the Rényi entropy for noninteger values of α . For $\alpha < 1$ and noninteger $\alpha > 1$, they construct algorithms additively estimating (with high probability) the Rényi entropy $H_\alpha(\mathbf{p})$ of an unknown distribution \mathbf{p} using $O\left(\frac{n^{1/\alpha}}{\log n}\right)$ and $O\left(\frac{n}{\log n}\right)$ samples, respectively. Both algorithms employ empirical and bias-corrected estimators, $\widehat{\mathcal{M}}_\alpha^e$ and $\widehat{\mathcal{M}}_\alpha^u$, and achieve a factor of $\log n$ improvement due to exploiting the polynomial approximation technique. In particular, for domain elements i with $N_i > \tau$, the empirical estimator $\widehat{\mathcal{M}}_\alpha^e$ is utilized, where $\tau = O(\log n)$. On the other hand, for domain elements i with $N_i \leq \tau$, the task of estimating $H_\alpha(\mathbf{p})$ is conducted in two steps. Firstly, the function x^α is approximated by an integer d -degree polynomial $p(x) = \sum_{j=0}^d c_j x^j$. Secondly, the quantity $(\mathbf{p}[i])^\alpha$ is estimated by utilizing the bias-corrected estimator $\widehat{\mathcal{M}}_\alpha^u$ for each term of $p(\mathbf{p}[i])$.

2.2.2. Lower Bounds

In the second part of this section, we exhibit lower bounds [7] on the task of estimating Rényi entropy to within an additive error. The techniques utilized for the task

resemble to the ones used for establishing a lower bound on estimating Shannon entropy. Unsurprisingly, since Rényi entropy is a symmetric property, the fingerprint of a sample is satisfactory for the analysis, as described in Section 2.1. Acharya *et al.* construct two probability distributions \mathbf{p} and \mathbf{q} such that the difference, $|H_\alpha(\mathbf{p}) - H_\alpha(\mathbf{q})|$, is sufficiently big to distinguish \mathbf{p} from \mathbf{q} . However, the fingerprint distributions, $\mathbf{p}_{\mathcal{F}}$, $\mathbf{q}_{\mathcal{F}}$, corresponding to \mathbf{p} , \mathbf{q} , respectively, have small total variation distance so that given a fingerprint it is impossible to decide which distribution it is obtained from when the number of samples is less than certain quantity. The following sets a bound on the total variation distance between the fingerprint distributions.

Theorem 2.41. ([7], Theorem 13) *Given distributions \mathbf{p} and \mathbf{q} such that $\max_x \max\{\mathbf{p}_x, \mathbf{q}_x\} \leq \frac{\epsilon}{40m}$, for Poisson sampling with $M \sim \text{Pois}(m)$, it holds that*

$$\|\mathbf{p}_{\mathcal{F}} - \mathbf{q}_{\mathcal{F}}\| \leq \frac{\epsilon}{2} + 5 \sum_{\alpha} m^{\alpha} |\mathcal{M}_{\alpha}(\mathbf{p}) - \mathcal{M}_{\alpha}(\mathbf{q})|.$$

Thus, it is sufficient to build \mathbf{p} , \mathbf{q} with distant Rényi entropies, yet having identical moments. For a probability distribution \mathbf{p} , let $\|\mathbf{p}\|_r = \left(\sum_{i=1}^n |\mathbf{p}[i]|^r \right)^{1/r}$, where r is a positive real number. We present the main ingredients of constructing such \mathbf{p} , \mathbf{q} pairs.

Lemma 2.42. ([7], Lemma 14) *For every $d \in \mathbb{N}$ and noninteger α , there exist positive vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ such that*

$$\|\mathbf{x}\|_r = \|\mathbf{y}\|_r, \quad 1 \leq r \leq d - 1$$

$$\|\mathbf{x}\|_d \neq \|\mathbf{y}\|_d,$$

$$\|\mathbf{x}\|_{\alpha} \neq \|\mathbf{y}\|_{\alpha}.$$

Definition 2.43. *For every positive integer d and every vector $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$, construct a distribution $\mathbf{p}^{\mathbf{x}}$ with support size dn as follows*

$$\mathbf{p}_{ij}^{\mathbf{x}} = \frac{|x_i|}{n\|\mathbf{x}\|_1}, \quad 1 \leq i \leq d, \quad 1 \leq j \leq n.^5$$

⁵Analyzing probability distributions with support size dn instead of the ones with support size n is only for practical reasons and does not affect the lower bound results stated above.

It only remains to construct distributions \mathbf{p}, \mathbf{q} based on Lemma 2.42 and Definition 2.43 for three different cases of α and apply Theorem 2.41 to finalize the proof. Before stating the lower bound results on estimating Rényi entropy to within additive error, we give one last definition.

Definition 2.44. Let $f(n) = \tilde{\Omega}(n^\beta)$ indicate that for all sufficiently large n and for all $\eta > 0$, $f(n) > n^{\beta-\eta}$ where $f: \mathbb{R} \rightarrow \mathbb{R}$ and $\beta \in \mathbb{R}$. That is, $f(n)$ grows polynomially in n with exponent not less than β .

Theorem 2.45. ([7], Theorem 15) For any integer $\alpha > 1$, $\Omega(n^{1-1/\alpha})$ samples are necessary to estimate (with high probability) the Rényi entropy $H_\alpha(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$.

Theorem 2.46. ([7], Theorem 16) For any noninteger $\alpha > 1$, $\tilde{\Omega}(n)$ samples are necessary to estimate (with high probability) the Rényi entropy $H_\alpha(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$.

Theorem 2.47. ([7], Theorem 17) For any $1 > \alpha > 0$, $\tilde{\Omega}(n^{1/\alpha})$ samples are necessary to estimate (with high probability) the Rényi entropy $H_\alpha(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$.

3. SAMP+PMF MODEL

We formally define the “unconventional” model.

Definition 3.1. Let \mathbf{p} be a probability distribution on $[n]$ and *PMF* denote a type of query which takes input $i \in [n]$ and returns its probability mass function $\mathbf{p}[i]$. The *SAMP+PMF* model is a model of query complexity in which both *SAMP* and *PMF* queries are utilized to interact with \mathbf{p} .

Note that a trivial upper bound on estimating entropy is n , since one can learn \mathbf{p} fully via n *PMF* queries. We start with the results regarding multiplicatively estimating entropy. As shown in [9, 12], estimating $H(\mathbf{p})$ (with high probability) to within a multiplicative factor of $1 + \gamma$ requires between

$$\Omega\left(\frac{\log n}{\max(\gamma, \gamma^2)}\right) \cdot \frac{1}{H(\mathbf{p})} \quad (3.1)$$

and

$$O\left(\frac{\log n}{\gamma^2}\right) \cdot \frac{1}{H(\mathbf{p})} \quad (3.2)$$

SAMP+PMF queries. One disadvantage of these bounds is them depending quantitatively on the entropy $H(\mathbf{p})$ itself.

The task of estimating (with high probability) Shannon entropy to within additive error is examined in [10]. Canonne and Rubinfeld construct an algorithm with $O(\log^2 n)$ query complexity, demonstrating superiority of *SAMP+PMF* model to *SAMP* model in which $\Omega\left(\frac{n}{\log n}\right)$ samples are necessary for the same task. They build the algorithm on the following:

$$H(\mathbf{p}) = \sum_{i \in [n]} \mathbf{p}[i] \log \frac{1}{\mathbf{p}[i]} = \mathbf{E}_{i \sim \mathbf{p}} \left[\log \frac{1}{\mathbf{p}[i]} \right] \quad (3.3)$$

Since it is difficult to give an upper bound on the quantity $\log \frac{1}{\mathbf{p}[i]}$, Identity 3.3 is reconstructed. Note that the function $f(x) = x \log \left(\frac{1}{x}\right)$ is increasing for $x \in \left(0, \frac{1}{e}\right)$ and $\lim_{x \rightarrow 0^+} f(x) = 0$. Then for any threshold $\tau \in \left(0, \frac{1}{e}\right)$,

$$\begin{aligned} H(\mathbf{p}) &= \sum_{i:\mathbf{p}[i] \geq \tau} \mathbf{p}[i] \log \frac{1}{\mathbf{p}[i]} + \sum_{i:\mathbf{p}[i] < \tau} \mathbf{p}[i] \log \frac{1}{\mathbf{p}[i]} \Rightarrow \\ H(\mathbf{p}) &\geq \sum_{i:\mathbf{p}[i] \geq \tau} \mathbf{p}[i] \log \frac{1}{\mathbf{p}[i]} = H(\mathbf{p}) - \sum_{i:\mathbf{p}[i] < \tau} \mathbf{p}[i] \log \frac{1}{\mathbf{p}[i]} \geq H(\mathbf{p}) - n \cdot \tau \log \frac{1}{\tau} \end{aligned} \quad (3.4)$$

Assume $\frac{\Delta}{n} < \frac{1}{2}$, and let $\tau = \frac{\frac{\Delta}{n}}{10 \log \frac{n}{\Delta}}$, so that $n \cdot \tau \log \frac{1}{\tau} \leq \frac{\Delta}{2}$. Define a function $\varphi(x) = \log \frac{1}{x} \mathbb{1}_{\{x \geq \tau\}}$. Then, Equation 3.4 implies that

$$H(\mathbf{p}) \geq \mathbf{E}_{i \sim \mathbf{p}} [\varphi(\mathbf{p}[i])] \geq H(\mathbf{p}) - \frac{\Delta}{2}.$$

Consequently, estimating $\mathbf{E}_{i \sim \mathbf{p}} [\varphi(\mathbf{p}[i])]$ to within additive $\frac{\Delta}{2}$ is sufficient for estimating $H(\mathbf{p})$ to within additive Δ . Observe that $0 \leq \varphi(\mathbf{p}[i]) \leq \log \frac{1}{\tau} \approx \log \frac{n}{\Delta}$. Let X_1, \dots, X_m be independent samples drawn from \mathbf{p} where $m = O\left(\frac{\log^2 \frac{n}{\Delta}}{\Delta^2}\right)$. If we compute the quantities $\varphi(\mathbf{p}[X_j])$ via PMF queries and apply the Hoeffding bound on random variables $Y_j = \frac{\varphi(\mathbf{p}[X_j])}{\log \frac{1}{\tau}}$ for $1 \leq j \leq m$,

$$\Pr \left[\left| \frac{1}{m} \sum_{j=1}^m \varphi(\mathbf{p}[X_j]) - \mathbf{E}_{i \sim \mathbf{p}} [\varphi(\mathbf{p}[i])] \right| \geq \frac{\Delta}{2} \right] \leq 2e^{-\frac{\Delta^2 m}{\log^2 \frac{1}{\tau}}} \leq \frac{1}{3}.$$

Theorem 3.2. ([10], Theorem 10) *In the SAMP+PMF model, there exists an algorithm estimating Shannon entropy (with high probability) to within $\pm\Delta$ with sample complexity $O\left(\frac{\log^2 \frac{n}{\Delta}}{\Delta^2}\right)$.*

Apart from achieving an exponentially better bound, the simplicity of the algorithm compared to the entangled nature of the estimators described in Section 2.1 is sufficient to illuminate the power a PMF query adds to the model. Yet, how helpful is it for estimating Rényi entropy? We will examine this question in Section 4.3.

Algorithm: ESTIMATOR V

- Fix $\tau = \frac{\Delta}{10 \log \frac{n}{\Delta}}$ and $m = \lceil \frac{\ln 6}{\Delta^2} \log^2 \frac{1}{\tau} \rceil$.
- Draw m independent samples X_1, \dots, X_m from \mathbf{p} .
- Compute $Y_j = \log \frac{1}{\mathbf{p}[X_j]} \mathbb{1}_{\{\mathbf{p}[X_j] \geq \tau\}}$ by evaluating PMF on X_j for $1 \leq j \leq m$.
- Output $\frac{1}{m} \sum_{j=1}^m Y_j$.

Figure 3.1. SAMP+PMF Estimator of Shannon Entropy

We can extend the SAMP+PMF model by making a slight modification on a PMF query.

Definition 3.3. Let \mathbf{p} be a probability distribution on $[n]$ and CDF denote a type of query which takes input $i \in [n]$ and returns its cumulative distribution function $\sum_{j=1}^i \mathbf{p}[j]$. The SAMP+CDF model is a model in which both SAMP and CDF queries are utilized to interact with \mathbf{p} .

Observe that since $\text{PMF}(i) = \text{CDF}(i) - \text{CDF}(i-1)$, a PMF query is simulated by at most two CDF queries. Canonne and Rubinfeld [10] show that in certain cases the SAMP+CDF model is more powerful than the SAMP+PMF model in estimating Shannon entropy. More specifically, if a probability distribution \mathbf{p} is known to be monotone, $\Omega(\log n)$ queries are necessary for estimating $H(\mathbf{p})$ in the SAMP+PMF model, however, there exists an algorithm in the SAMP+CDF model using $O(\log^2 \log n)$ queries for the same task. Does this hold true for estimating entropy of an arbitrary distribution? Answers to these and more questions are presented in the next chapter.

4. OPTIMAL BOUNDS FOR ESTIMATING ENTROPY WITH PMF QUERIES

4.1. Our Results, and Comparison with Prior Work

As described in Chapter 3, Canonne and Rubinfeld [10] build a SAMP+PMF algorithm estimating with high probability Shannon entropy to within ± 1 using $O(\log^2 n)$ queries. In addition, they prove that $\Omega(\log n)$ queries are necessary for the task. Our first main result is an improved, optimal lower bound:

First main theorem. *In the SAMP+PMF model, $\Omega(\log^2 n)$ queries are necessary to estimate (with high probability) the Shannon entropy $H(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within ± 1 .*

Remark 4.1. *Our lower bound and the lower bound for multiplicative estimation of Shannon entropy given in Expression 3.1 hold even under the promise that $H(\mathbf{p}) = \Theta(\log n)$. Note that Expression 3.1 yields a nonoptimal $\Omega(\log n)$ lower bound for additive estimation problem by taking $\gamma = \frac{1}{\log n}$.*

More precisely, Canonne and Rubinfeld show that $O(\frac{\log^2 n}{\Delta^2})$ queries are sufficient for estimating Shannon entropy to within $\pm \Delta$.⁶ However, this result is trivial once $\Delta \leq \frac{\log n}{\sqrt{n}}$ since \mathbf{p} can be learned exactly via n PMF queries. In fact, our first main theorem gives a matching lower bound for essentially the full range of Δ : we prove that $\Omega(\frac{\log^2 n}{\Delta^2})$ SAMP+PMF queries are necessary for any $\frac{1}{n^{0.4999}} \leq \Delta \leq \frac{\log n}{16 \cdot 10^6}$.

Our second main result is regarding the estimation of Rényi entropy $H_\alpha(\mathbf{p})$ for various parameters $0 \leq \alpha \leq \infty$. Recall that Acharya *et al.* [7] establish three different lower and upper bound pairs on additively estimating $H_\alpha(\mathbf{p})$ in the SAMP model. They prove that $\Theta(n^{1-1/\alpha})$ samples are necessary and sufficient when $\alpha > 1$ is an integer;

⁶They actually state $O(\frac{\log^2(n/\Delta)}{\Delta^2})$, but this is the same as $O(\frac{\log^2 n}{\Delta^2})$ because the range of interest is $\frac{\log n}{\sqrt{n}} \leq \Delta \leq \log n$.

$\tilde{\Omega}(n)$ samples are necessary when $\alpha > 1$ is a noninteger; $\tilde{\Omega}(n^{1/\alpha})$ samples are necessary when $1 > \alpha > 0$. We give matching upper and lower bounds on estimating $H_\alpha(\mathbf{p})$ for all $\alpha > 1$. Apparently, PMF queries provide no advantage in estimating Rényi entropy for integer α , whereas they *are* advantageous for noninteger α .

Second main theorem. *Let $\alpha > 1$ be a real number. In the SAMP+PMF model, $\Omega\left(\frac{n^{1-1/\alpha}}{2^{2\Delta}}\right)$ queries are necessary and $O\left(\frac{n^{1-1/\alpha}}{(1-2^{1-\alpha}\Delta)^2}\right)$ queries are sufficient to estimate (with high probability) the Rényi entropy $H_\alpha(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$.*

We show that our two bounds extend to the SAMP+CDF model, thus, answering the question of the previous chapter. In addition, we give a matching lower bound for estimating support size to within $\pm\epsilon n$ in the SAMP+CDF model. Lastly, we provide an upper bound on additively estimating Tsallis entropy in the SAMP+PMF model.

4.2. First Main Theorem

We present a well-known fact which is of key importance in establishing a lower bound on additively estimating Shannon entropy in the SAMP+PMF model.

Lemma 4.2. *For $\lambda \in (0, \frac{1}{4}]$, $\Theta\left(\frac{1}{\lambda^2}\right)$ samples are necessary and sufficient to distinguish between the uniform distribution $\mathbf{p}_1 = (\frac{1}{2}, \frac{1}{2})$ and the biased distribution $\mathbf{p}_2 = (\frac{1}{2} + \lambda, \frac{1}{2} - \lambda)$.⁷*

Proof. We start with proving the upper bound. Let X_1, \dots, X_m be m independent samples drawn from \mathbf{p} which is a probability distribution promised to be either a uniform or a biased distribution. Then $\frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{X_i=1\}}$ is an unbiased estimator for $\mathbf{p}[1]$, where $\mathbb{1}_{\{E\}}$ is an indicator function of an event E . We need to approximate (with high probability) $\mathbf{p}[1]$ to within some $\epsilon < \frac{\lambda}{2}$ to properly distinguish between two

⁷The range of λ can be easily extended to $(0, \frac{1}{2} - \xi]$, where ξ is an arbitrarily small constant.

distributions. By applying the Hoeffding bound,

$$\Pr \left[\left| \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{X_i=1\}} - \mathbf{p}[1] \right| > \epsilon \right] \leq 2e^{-2\epsilon^2 m}.$$

Obviously, to bound the right-hand side of the inequality $m = O\left(\frac{1}{\lambda^2}\right)$ samples are sufficient. Now we prove the lower bound.

Definition 4.3. Let $\mathbf{d}_1, \mathbf{d}_2$ be discrete probability distributions. The Kullback-Leibler (KL) divergence is defined as

$$\mathbf{KL}(\mathbf{d}_1 \parallel \mathbf{d}_2) = \sum_x \mathbf{d}_1[x] \log \frac{\mathbf{d}_1[x]}{\mathbf{d}_2[x]}. \quad (4.1)$$

By convention, $\mathbf{d}_1[x] \log \mathbf{d}_1[x]$ is set to 0 if $\mathbf{d}_1[x] = 0$. We exploit the KL-divergence for the particular class of probability distributions.

Fact 4.4. Let X_1, \dots, X_m be m i.i.d random variables drawn from \mathbf{d}_1 (respectively \mathbf{d}_2), where $m \in \mathbb{N}$. Let \mathbf{d}_1^m (respectively \mathbf{d}_2^m) denote a joint probability distribution of random variables X_1, \dots, X_m . Then

$$\mathbf{KL}(\mathbf{d}_1^m \parallel \mathbf{d}_2^m) = m \mathbf{KL}(\mathbf{d}_1 \parallel \mathbf{d}_2).$$

Proof. By definition $\mathbf{d}_1^m = (\mathbf{d}_1)^m$ and $\mathbf{d}_2^m = (\mathbf{d}_2)^m$. We prove the statement via strong induction.

Base case: For $m = 2$,

$$\begin{aligned} \mathbf{KL}(\mathbf{d}_1^2 \parallel \mathbf{d}_2^2) &= \sum_{X_1} \sum_{X_2} \mathbf{d}_1[X_1] \mathbf{d}_1[X_2] \cdot \log \left(\frac{\mathbf{d}_1[X_1] \mathbf{d}_1[X_2]}{\mathbf{d}_2[X_1] \mathbf{d}_2[X_2]} \right) \\ &= \sum_{X_1} \sum_{X_2} \mathbf{d}_1[X_1] \mathbf{d}_1[X_2] \cdot \left(\log \frac{\mathbf{d}_1[X_1]}{\mathbf{d}_2[X_1]} + \log \frac{\mathbf{d}_1[X_2]}{\mathbf{d}_2[X_2]} \right) \\ &= \sum_{X_1} \mathbf{d}_1[X_1] \log \frac{\mathbf{d}_1[X_1]}{\mathbf{d}_2[X_1]} \sum_{X_2} \mathbf{d}_1[X_2] + \sum_{X_1} \mathbf{d}_1[X_1] \sum_{X_2} \mathbf{d}_1[X_2] \log \frac{\mathbf{d}_1[X_2]}{\mathbf{d}_2[X_2]} \end{aligned}$$

$$\begin{aligned}
&= \sum_{X_1} \mathbf{d}_1[X_1] \log \frac{\mathbf{d}_1[X_1]}{\mathbf{d}_2[X_1]} + \sum_{X_2} \mathbf{d}_1[X_2] \log \frac{\mathbf{d}_1[X_2]}{\mathbf{d}_2[X_2]} \\
&= 2 \mathbf{KL}(\mathbf{d}_1 \| \mathbf{d}_2).
\end{aligned}$$

Induction step: Assume that the statement is true for all $i \leq m - 1$. Then

$$\begin{aligned}
\mathbf{KL}(\mathbf{d}_1^m \| \mathbf{d}_2^m) &= \sum_{X_1} \cdots \sum_{X_m} \left(\prod_{i=1}^m \mathbf{d}_1[X_i] \log \left(\prod_{j=1}^m \frac{\mathbf{d}_1[X_j]}{\mathbf{d}_2[X_j]} \right) \right) \\
&= \sum_{X_1} \cdots \sum_{X_m} \left(\prod_{i=1}^m \mathbf{d}_1[X_i] \left(\sum_{j=1}^m \log \frac{\mathbf{d}_1[X_j]}{\mathbf{d}_2[X_j]} \right) \right) \\
&= \sum_{X_1} \mathbf{d}_1[X_1] \log \frac{\mathbf{d}_1[X_1]}{\mathbf{d}_2[X_1]} \sum_{X_2} \mathbf{d}_1[X_2] \cdots \sum_{X_m} \mathbf{d}_1[X_m] \\
&\quad + \sum_{X_1} \mathbf{d}_1[X_1] \sum_{X_2} \cdots \sum_{X_m} \left(\prod_{i=2}^m \mathbf{d}_1[X_i] \left(\sum_{j=2}^m \log \frac{\mathbf{d}_1[X_j]}{\mathbf{d}_2[X_j]} \right) \right) \\
&= \sum_{X_1} \mathbf{d}_1[X_1] \log \frac{\mathbf{d}_1[X_1]}{\mathbf{d}_2[X_1]} + \sum_{X_2} \cdots \sum_{X_m} \left(\prod_{i=2}^m \mathbf{d}_1[X_i] \left(\sum_{j=2}^m \log \frac{\mathbf{d}_1[X_j]}{\mathbf{d}_2[X_j]} \right) \right) \\
&= \mathbf{KL}(\mathbf{d}_1 \| \mathbf{d}_2) + \mathbf{KL}(\mathbf{d}_1^{m-1} \| \mathbf{d}_2^{m-1}) \\
&= m \mathbf{KL}(\mathbf{d}_1 \| \mathbf{d}_2). \quad \square
\end{aligned}$$

Lemma 4.5. Let $\widehat{\mathbf{d}}_1, \widehat{\mathbf{d}}_2$ be discrete probability distributions on the universe U , and let $f : U \rightarrow [0, C]$ be a real-valued function for some positive constant C . Then

$$\left| \mathbf{E}_{\widehat{\mathbf{d}}_1}[f(x)] - \mathbf{E}_{\widehat{\mathbf{d}}_2}[f(x)] \right| \leq C \cdot \|\widehat{\mathbf{d}}_1 - \widehat{\mathbf{d}}_2\|_1.$$

Proof.

$$\begin{aligned}
\left| \mathbf{E}_{\widehat{\mathbf{d}}_1}[f(x)] - \mathbf{E}_{\widehat{\mathbf{d}}_2}[f(x)] \right| &= \left| \sum_{x \in U} \widehat{\mathbf{d}}_1[x] f(x) - \sum_{x \in U} \widehat{\mathbf{d}}_2[x] f(x) \right| \\
&= \left| \sum_{x \in U} (\widehat{\mathbf{d}}_1[x] - \widehat{\mathbf{d}}_2[x]) f(x) \right| \\
&\leq \sum_{x \in U} \left| (\widehat{\mathbf{d}}_1[x] - \widehat{\mathbf{d}}_2[x]) \right| f(x) \\
&\leq C \sum_{x \in U} \left| (\widehat{\mathbf{d}}_1[x] - \widehat{\mathbf{d}}_2[x]) \right|
\end{aligned}$$

$$= C \cdot \|\widehat{\mathbf{d}}_1 - \widehat{\mathbf{d}}_2\|_1. \quad \square$$

As a final ingredient we give an upper bound on the KL-divergence between the biased distribution and the uniform distribution.

$$\begin{aligned}
\mathbf{KL}(\mathbf{p}_2 \|\mathbf{p}_1) &= \left(\frac{1}{2} + \lambda\right) \log \frac{\frac{1}{2} + \lambda}{\frac{1}{2}} + \left(\frac{1}{2} - \lambda\right) \log \frac{\frac{1}{2} - \lambda}{\frac{1}{2}} \\
&= \frac{1}{2} \log((1 - 2\lambda)(1 + 2\lambda)) + \lambda \log \left(\frac{1 + 2\lambda}{1 - 2\lambda}\right) \\
&\leq \lambda \log \left(\frac{1 + 2\lambda}{1 - 2\lambda}\right) \\
&= \frac{\lambda}{\ln 2} \ln \left(1 + \frac{4\lambda}{1 - 2\lambda}\right) \\
&\leq \frac{4\lambda^2}{\ln 2} \cdot \frac{1}{1 - 2\lambda} \\
&\leq \frac{8\lambda^2}{\ln 2},
\end{aligned} \tag{4.2}$$

where the second inequality follows from the exponential inequality.

Suppose for the sake of contradiction that there exists a hypothetical algorithm \mathcal{D} such that given $m = o\left(\frac{1}{\lambda^2}\right)$ independent samples, decides whether $\mathbf{p} = \mathbf{p}_1$ or $\mathbf{p} = \mathbf{p}_2$ with probability of error at most $\frac{1}{3}$. Let \mathcal{D} output 1 to indicate that $\mathbf{p} = \mathbf{p}_1$ and output 0 to indicate that $\mathbf{p} = \mathbf{p}_2$. Denote by \mathcal{X} the set of m independent samples $\{X_1, \dots, X_m\}$. Then

$$\Pr[\mathcal{D}(\mathcal{X} \sim \mathbf{p}_1^m) = 1] \geq \frac{2}{3} \quad \text{and} \quad \Pr[\mathcal{D}(\mathcal{X} \sim \mathbf{p}_2^m) = 0] \geq \frac{2}{3}.$$

In other words,

$$\mathbf{E}_{\mathcal{X} \sim \mathbf{p}_1^m}[\mathcal{D}(\mathcal{X})] \geq \frac{2}{3} \quad \text{and} \quad \mathbf{E}_{\mathcal{X} \sim \mathbf{p}_2^m}[\mathcal{D}(\mathcal{X})] \leq \frac{1}{3},$$

which gives

$$\mathbf{E}_{x \sim \mathbf{p}_1^m} [\mathcal{D}(x)] - \mathbf{E}_{x \sim \mathbf{p}_2^m} [\mathcal{D}(x)] \geq \frac{1}{3}. \quad (4.3)$$

Letting $\widehat{\mathbf{d}}_1 = \mathbf{p}_1^m$, $\widehat{\mathbf{d}}_2 = \mathbf{p}_2^m$ and $f = \mathcal{D}$, Lemma 4.5 and Inequality 4.3 imply

$$\|\mathbf{p}_1^m - \mathbf{p}_2^m\|_1 \geq \frac{1}{3}. \quad (4.4)$$

Then

$$\begin{aligned} o\left(\frac{1}{\lambda^2}\right) = m &= \frac{\mathbf{KL}(\mathbf{p}_1^m \|\mathbf{p}_2^m)}{\mathbf{KL}(\mathbf{p}_1 \|\mathbf{p}_2)} && \text{(Fact 4.4)} \\ &\geq \frac{1}{2 \ln 2} \cdot \|\mathbf{p}_1^m - \mathbf{p}_2^m\|_1^2 \cdot \frac{1}{\mathbf{KL}(\mathbf{p}_1 \|\mathbf{p}_2)} && \text{(Pinsker's inequality)} \\ &\geq \frac{1}{18 \ln 2} \cdot \frac{1}{\mathbf{KL}(\mathbf{p}_1 \|\mathbf{p}_2)} && \text{(Inequality 4.4)} \\ &\geq \frac{1}{144 \lambda^2} = \Omega\left(\frac{1}{\lambda^2}\right), && \text{(Inequality 4.2)} \end{aligned}$$

which is a contradiction. \square

We establish a tight lower bound on estimating Shannon entropy in the following theorem.

Theorem 4.6. *In the SAMP+PMF model, $\Omega\left(\frac{\log^2 n}{\Delta^2}\right)$ queries are necessary to estimate (with high probability) the Shannon entropy $H(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$, where $\frac{1}{n^{0.4999}} \leq \Delta \leq \frac{\log n}{16 \cdot 10^6}$.*

Proof. We will show that a hypothetical SAMP+PMF algorithm \mathcal{E} that can estimate the entropy of an unknown distribution on $[n]$ to within $\pm\Delta$ using $o\left(\frac{\log^2 n}{\Delta^2}\right)$ queries would contradict Lemma 4.2 stating that $\Omega(1/\lambda^2)$ coin tosses are necessary to determine whether a given coin is fair, or comes up heads with probability $1/2 + \lambda$.

The idea is to use the given coin to realize the probability distribution that \mathcal{E} will work on. Let n be the smallest one millionth power of a natural number that satisfies $\frac{4 \cdot 10^6 \Delta}{\log n} \leq \lambda$.⁸ Partition the domain $[n]$ into $M = n^{0.999999}$ consecutive blocks I_1, \dots, I_M , each containing $K = \frac{n}{M} = n^{0.000001}$ elements. Each block will be labeled either as a tails or a heads block. The internal distribution of each heads block is uniform, i.e. each element has probability mass $\frac{1}{MK} = \frac{1}{n}$. In each tails block, the first element has probability mass $\frac{1}{n^{0.999999}}$, while the rest of the elements have probability mass 0. Note that the total probability mass of each block is $K \cdot \frac{1}{MK} = \frac{1}{M} = \frac{1}{n^{0.999999}}$, regardless of its label. We will now describe a costly method of constructing a probability distribution \mathbf{p} of this kind, using a coin that comes up heads with probability d :

- Throw the coin M times to obtain the outcomes X_1, \dots, X_M ,
- Set the label of block I_m to X_m , for all $m \in [M]$.

Let X be the number of heads blocks in \mathbf{p} . Then $\mu = \mathbf{E}[X] = Md$. Let $\bar{X} = \frac{X}{M}$ denote the proportion of heads blocks in \mathbf{p} . Then we can calculate the entropy $H(\mathbf{p})$ by calculating the individual entropies of the blocks. For a heads block, the entropy is $K \cdot \frac{1}{MK} \cdot \log(MK) = \frac{1}{M} \log n$. The entropy of a tails block is $\frac{1}{n^{0.999999}} \log(n^{0.999999}) = \frac{0.999999}{M} \log n$. Since there are $M\bar{X}$ heads blocks and $M(1 - \bar{X})$ tails blocks, the total entropy becomes

$$\begin{aligned} H(\mathbf{p}) &= M\bar{X} \cdot \frac{1}{M} \log n + M(1 - \bar{X}) \cdot \frac{0.999999}{M} \log n \\ &= \bar{X} \log n + 0.999999(1 - \bar{X}) \log n \\ &= (0.999999 + 0.000001\bar{X}) \log n. \end{aligned} \tag{4.5}$$

Note that this function is monotone with respect to \bar{X} . Define two families of distributions \mathcal{P}_1 and \mathcal{P}_2 constructed by the above process, taking d to be $p_1 = \frac{1}{2}$ and $p_2 = \frac{1}{2} + \lambda$, respectively. Let \mathbf{p}_1 (respectively \mathbf{p}_2) be a probability distribution randomly chosen from \mathcal{P}_1 (respectively \mathcal{P}_2).

⁸Since Lemma 4.2 holds for $\lambda \in (0, \frac{1}{2} - \xi]$, where ξ is an arbitrarily small constant, the upper bound on Δ can be extended to $\frac{(1-2\xi)\log n}{8 \cdot 10^6}$.

Proposition 4.7. \mathbf{p}_1 has entropy at most $0.9999995 \log n + \Delta$ with high probability.

Proof. We prove this by using the Chernoff bound on the number of heads blocks in the distribution.

$$\begin{aligned}
\Pr\left[X \geq \left(p_1 + \frac{10^6 \Delta}{\log n}\right)M\right] &= \Pr\left[X \geq \frac{M}{2} \left(1 + \frac{2 \cdot 10^6 \Delta}{\log n}\right)\right] \\
&\leq \exp\left(-\frac{\frac{4 \cdot 10^{12} \Delta^2}{\log^2 n} \frac{M}{2}}{2 + \frac{2 \cdot 10^6 \Delta}{\log n}}\right) \\
&= \exp\left(-\frac{10^{12} \Delta^2 M}{\log n (\log n + 10^6 \Delta)}\right) \\
&\leq \exp\left(-\frac{10^{12} \cdot n^{0.999999} / n^{0.999998}}{\log^2 n (1 + 10^6)}\right) \\
&= \exp\left(-\frac{10^{12} \cdot n^{0.000001}}{\log^2 n (1 + 10^6)}\right) \\
&= o(1).
\end{aligned}$$

The last term indicates that the number of heads blocks $X < \left(p_1 + \frac{10^6 \Delta}{\log n}\right)M$, and the proportion of the heads blocks $\bar{X} < \left(p_1 + \frac{10^6 \Delta}{\log n}\right)$ with high probability. Thus, with high probability

$$H[\mathbf{p}_1] = (0.999999 + 0.000001 \bar{X}) \log n < 0.9999995 \log n + \Delta. \quad \square$$

Proposition 4.8. \mathbf{p}_2 has entropy at least $0.9999995 \log n + 3\Delta$ with high probability.

Proof. We find a similar bound by;

$$\begin{aligned}
\Pr\left[X \leq \left(p_2 - \frac{10^6 \Delta}{\log n}\right)M\right] &= \Pr\left[X \leq p_2 M \left(1 - \frac{10^6 \Delta}{p_2 \log n}\right)\right] \\
&\leq \exp\left(-\frac{\frac{10^{12} \Delta^2}{p_2^2 \log^2 n} p_2 M}{2}\right) \\
&= \exp\left(-\frac{10^{12} \Delta^2 M}{2 p_2 \log^2 n}\right) \\
&= \exp\left(-\frac{10^{12} \Delta^2 M}{2 \left(\frac{1}{2} + \lambda\right) \log^2 n}\right)
\end{aligned}$$

$$\begin{aligned}
&\leq \exp\left(-\frac{n^{0.000001}}{\log^2 n}\right) \\
&= o(1).
\end{aligned}$$

The last term indicates that the number of heads blocks $X > (p_2 - \frac{10^6 \Delta}{\log n})M$, and the proportion of the heads blocks $\bar{X} > (p_2 - \frac{10^6 \Delta}{\log n})$ with high probability. Thus, with high probability

$$\begin{aligned}
H[\mathbf{p}_2] &= (0.999999 + 0.000001\bar{X}) \log n \\
&> \left(0.999999 + 0.000001 \left(\frac{1}{2} + \lambda - \frac{10^6 \Delta}{\log n}\right)\right) \log n \\
&= 0.9999995 \log n + 0.000001 \left(\lambda - \frac{10^6 \Delta}{\log n}\right) \log n \tag{4.6} \\
&\geq 0.9999995 \log n + 0.000001 \left(\frac{3 \cdot 10^6 \Delta}{\log n}\right) \log n \\
&= 0.9999995 \log n + 3\Delta. \quad \square
\end{aligned}$$

Since the entropies of \mathbf{p}_1 and \mathbf{p}_2 are sufficiently far apart from each other, our hypothetical estimator \mathcal{E} can be used to determine whether the underlying coin has probability p_1 or p_2 associated with it. To arrive at the contradiction we want, we must ensure that the coin is not thrown too many times during this process. This is achieved by constructing the distribution “on-the-fly” [10] during the execution of \mathcal{E} , throwing the coin only when it is required to determine the label of a previously undefined block:

When \mathcal{E} makes a **SAMP** query, we choose a block I_m uniformly at random (since each block has probability mass $\frac{1}{M}$), and then flip the coin for I_m to decide its label if it is yet undetermined. We then draw a sample $i \sim \mathbf{d}_m$ from I_m , where \mathbf{d}_m is the normalized distribution of the m^{th} block.

When \mathcal{E} makes a **PMF** query on $i \in [n]$, we flip the coin to determine the label of the associated block I_m if it is yet undetermined. We then return the probability mass of i .

By this procedure, the queries of \mathcal{E} about the probability distribution \mathbf{p} (known to be either \mathbf{p}_1 or \mathbf{p}_2) can be answered by using at most one coin flip per query, i.e. $o\left(\frac{\log^2 n}{\Delta^2}\right)$ times in total.

Since we selected n so that $1/\lambda^2 = \Theta\left(\frac{\log^2 n}{\Delta^2}\right)$, this would mean that it is possible to distinguish between the two coins using only $o(1/\lambda^2)$ throws, which is a contradiction, letting us conclude that no algorithm can estimate the Shannon entropy $H(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$ with high probability making $o\left(\frac{\log^2 n}{\Delta^2}\right)$ queries. \square

We now give a similar lower bound for the SAMP+CDF model.

Corollary 4.9. *In the SAMP+CDF model, any algorithm estimating (with high probability) the Shannon entropy $H(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$ must make $\Omega\left(\frac{\log^2 n}{\Delta^2}\right)$ queries.*

Proof. The construction is identical to the one in the proof of Theorem 4.6, except that we now have to describe how the CDF queries of the estimation algorithm must be answered using the coin:

When \mathcal{E} makes a CDF query on $i \in [n]$, we flip the coin to determine the label of the associated block I_m if this is necessary. We then return the sum of the total probability mass of the blocks preceding I_m (which is $\frac{m-1}{M}$, since each block has a total probability mass of $\frac{1}{M}$ regardless of its label) and the probability masses of the elements from the beginning of I_m up to and including i itself. At most one coin flip per CDF query is therefore sufficient. \square

4.3. Second Main Theorem

4.3.1. Lower Bound

We present another well-known fact about distinguishing coins.

Fact 4.10. $\Omega\left(\frac{1}{\lambda}\right)$ samples are necessary to distinguish between the distribution $\mathbf{p}_1 = (1, 0)$ and the distribution $\mathbf{p}_2 = (1 - \lambda, \lambda)$.

Proof. Let X_1, \dots, X_m be m independent samples drawn from \mathbf{p} , which is a probability distribution promised to be either \mathbf{p}_1 or \mathbf{p}_2 . Note that observing the domain element 2 even once suffices to distinguish between two distributions, since $\mathbf{p}_1[2] = 0$. The probability of the most unfortunate case, not observing any 2 when $\mathbf{p} = \mathbf{p}_2$, is

$$\Pr\left[\sum_{i=1}^m \mathbb{1}_{\{X_i=1\}} = m\right] = (1 - \lambda)^m \leq e^{-\lambda m} \quad (4.7)$$

where the inequality follows from the exponential inequality. Obviously, to bound the right-hand side of Inequality 4.7, $m = \Omega\left(\frac{1}{\lambda}\right)$ samples are necessary. \square

We establish a lower bound on estimating Rényi entropy in the following theorem.

Theorem 4.11. For any $\alpha > 1$, $\Omega\left(\frac{n^{1-1/\alpha}}{2^{2\Delta}}\right)$ SAMP+PMF queries are necessary to estimate (with high probability) the Rényi entropy $H_\alpha(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$.

Proof. We will first prove the theorem for rational α , and show that it remains valid for irrationals at the end.

The proof has the same structure as that of Theorem 4.6. One difference is that we reduce from the problem of distinguishing a maximally biased coin that never comes up tails from a less biased one (instead of the problem of distinguishing a fair coin from a biased one).

Suppose that we are given a coin whose probability of coming up heads is promised to be either $p_1 = 1$ or $p_2 = 1 - \lambda$ for a specified number λ , and we must determine which is the case. As Fact 4.10 indicates, this task requires at least $\Omega(1/\lambda)$ coin throws. We will show that this fact is contradicted if one assumes that there exist natural numbers s and t , where $\alpha = \frac{s}{t} > 1$, such that it is possible to estimate (with high probability) the Rényi entropy $H_\alpha(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$ using an algorithm, say \mathcal{R} , that makes only $o\left(\frac{n^{1-1/\alpha}}{2^{2\Delta}}\right)$ SAMP+PMF queries.

Let n be the smallest number of the form $(\lceil 2^{2\Delta} \rceil j)^s$ that satisfies $\frac{5 \cdot \lceil 2^{2\Delta} \rceil}{n^{1-1/\alpha}} \leq \lambda$, where j is some natural number. Partition $[n]$ into $M = \frac{n^{1-1/\alpha}}{\lceil 2^{2\Delta} \rceil}$ consecutive blocks I_1, I_2, \dots, I_M , each of size $K = \lceil 2^{2\Delta} \rceil \cdot n^{1/\alpha}$. As in the proof of Theorem 4.6, a probability distribution \mathbf{p} can be realized by throwing a given coin M times to obtain the outcomes X_1, \dots, X_M , and setting the label of block I_m to X_m , for all $m \in [M]$, where each member of each heads block again has probability mass $1/n$. The first member of each tails block has probability mass $\frac{\lceil 2^{2\Delta} \rceil}{n^{1-1/\alpha}}$, and the remaining members have probability mass 0. We again have that each block has total probability mass $\frac{K}{n} = \frac{\lceil 2^{2\Delta} \rceil n^{1/\alpha}}{n} = \frac{1}{M}$ regardless of its label, so this process always results in a legal probability distribution.

If the coin is maximally biased, then \mathbf{p} becomes the uniform distribution, and $H_\alpha(\mathbf{p}) = \log n$. We will examine the probability of the same distribution being obtained using the less biased coin. Let \mathcal{P}_2 be the family of distributions constructed by the process described above, using a coin with probability p_2 of coming up heads. Let \mathbf{p}_2 be a probability distribution randomly chosen from \mathcal{P}_2 . The probability of the undesired case where \mathbf{p}_2 is the uniform distribution is

$$\Pr[\mathbf{p}_2 = \mathcal{U}([n])] = p_2^M = (1 - \lambda)^M \leq \left(1 - \frac{5 \cdot \lceil 2^{2\Delta} \rceil}{n^{1-1/\alpha}}\right)^M \leq e^{-\frac{5 \cdot \lceil 2^{2\Delta} \rceil}{n^{1-1/\alpha}} M} = e^{-5} \leq \frac{1}{1000} .$$

That is, with probability ≥ 0.999 , \mathbf{p}_2 has at least one element with probability mass $\frac{\lceil 2^{2\Delta} \rceil n^{1/\alpha}}{n}$. Let X be the number of heads outcomes, and let B and W denote the number of elements with probability mass $\frac{1}{n}$ and $\frac{\lceil 2^{2\Delta} \rceil n^{1/\alpha}}{n}$, respectively. It is not difficult to see that $B = K \cdot X$ and $W = M - X$. We just showed that $X < M$ with high

probability.

Then the Rényi entropy of the constructed distribution $\mathbf{p}_2 \in \mathcal{P}_2$ is, with high probability:

$$\begin{aligned}
H_\alpha(\mathbf{p}_2) &= \frac{1}{1-\alpha} \log \left(B \cdot \frac{1}{n^\alpha} + W \cdot \left(\frac{\lceil 2^{2\Delta} \rceil n^{\frac{1}{\alpha}}}{n} \right)^\alpha \right) \\
&= \frac{1}{1-\alpha} \log \left(\frac{K \cdot X/n + (M-X) \lceil 2^{2\Delta} \rceil^\alpha}{n^{\alpha-1}} \right) \\
&= \log n - \frac{1}{\alpha-1} \log (K \cdot X/n + (M-X) \lceil 2^{2\Delta} \rceil^\alpha) \\
&\leq \log n - \frac{1}{\alpha-1} \log (\lceil 2^{2\Delta} \rceil^\alpha) < \log n - 2\Delta.
\end{aligned}$$

Because $H_\alpha(\mathcal{U}([n])) - H_\alpha(\mathbf{p}_2) > 2\Delta$, \mathcal{R} has to be able to distinguish $\mathcal{U}([n])$ and \mathbf{p}_2 with high probability. We can then perform a simulation of \mathcal{R} involving an “on-the-fly” construction of distribution \mathbf{p} exactly as described in the proof of Theorem 4.6. As discussed in Section 4.2, this process requires no more coin throws than the number of SAMP+PMF queries made by \mathcal{R} , allowing us to determine the type of the coin using only $o\left(\frac{n^{1-1/\alpha}}{2^{2\Delta}}\right)$, that is, $o(1/\lambda)$ tosses with high probability, a contradiction.

Having thus proven the statement for rational α , it is straightforward to cover the case of irrational α : Note that $H_\alpha(\mathbf{p})$ is a continuous function of α for fixed \mathbf{p} . Given any \mathbf{p} and ε , for any irrational number α_i greater than 1, there exists a rational α_r which is so close to α_i such that $H_{\alpha_i}(\mathbf{p}) - H_{\alpha_r}(\mathbf{p}) < \varepsilon$. An efficient entropy estimation method for some irrational value of α would therefore imply the existence of an equally efficient method for some rational value, contradicting the result obtained above. \square

These results are generalized to the SAMP+CDF model in the same way as in Section 4.2.

Corollary 4.12. *For any $\alpha > 1$, $\Omega\left(\frac{n^{1-1/\alpha}}{2^{2\Delta}}\right)$ SAMP+PMF or SAMP+CDF queries are necessary to estimate (with high probability) the Rényi entropy $H_\alpha(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$.*

4.3.2. Upper Bound

We now show that PMF queries are useful for estimating H_α for noninteger α .

Theorem 4.13. *For any number $\alpha > 1$, there exists an algorithm estimating (with high probability) the Rényi entropy $H_\alpha(\mathbf{p})$ of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$ with $O\left(\frac{n^{1-1/\alpha}}{(1-2^{(1-\alpha)\Delta})^2}\right)$ SAMP+PMF queries.*

Proof. We will prove this statement for rational α . The generalization to irrational α discussed in the proof of Theorem 4.11 can be applied. Recall that Rényi entropy can be expressed as $H_\alpha(\mathbf{p}) = \frac{1}{1-\alpha} \log \mathcal{M}_\alpha(\mathbf{p})$, where $\mathcal{M}_\alpha(\mathbf{p}) = \sum_{i=1}^n (\mathbf{p}[i])^\alpha$ is the α^{th} moment of \mathbf{p} . Then, estimating $H_\alpha(\mathbf{p})$ to an additive accuracy of $\pm\Delta$ is equivalent to estimating $\mathcal{M}_\alpha(\mathbf{p})$ to a multiplicative accuracy of $2^{\pm\Delta(1-\alpha)}$. Therefore, we construct a multiplicative estimator for $\mathcal{M}_\alpha(\mathbf{p})$.

Let $\gamma = 1 - 2^{(1-\alpha)\Delta}$ and $m = \left\lceil \frac{100n^{1-1/\alpha}}{\gamma^2} \right\rceil$, and let X_1, \dots, X_m be i.i.d. random variables drawn from \mathbf{p} . Define $Y_i = (\mathbf{p}[X_i])^{\alpha-1}$, where $\mathbf{p}[X_i]$ can be calculated using a PMF query on X_i for $1 \leq i \leq m$. Note that

$$\mathbf{E}[Y_i] = \sum_{j=1}^n \mathbf{p}[j] (\mathbf{p}[j])^{\alpha-1} = \sum_{j=1}^n (\mathbf{p}[j])^\alpha = \mathcal{M}_\alpha(\mathbf{p}).$$

Then $\frac{1}{m} \sum_{i=1}^m Y_i$ is an unbiased estimator of $\mathcal{M}_\alpha(\mathbf{p})$, since

$$\mathbf{E}\left[\frac{1}{m} \sum_{i=1}^m Y_i\right] = \frac{1}{m} \sum_{i=1}^m \mathbf{E}[Y_i] = \mathcal{M}_\alpha(\mathbf{p}).$$

Moreover,

$$\mathbf{Var}[Y_i] = \mathbf{E}[Y_i^2] - \mathbf{E}[Y_i]^2 = \sum_{j=1}^n \mathbf{p}[j] (\mathbf{p}[j])^{2\alpha-2} - \mathbf{E}[Y_i]^2 = \mathcal{M}_{2\alpha-1}(\mathbf{p}) - \mathcal{M}_\alpha^2(\mathbf{p}).$$

Since the Y_i 's are also i.i.d. random variables,

$$\mathbf{Var} \left[\frac{1}{m} \sum_{i=1}^m Y_i \right] = \frac{1}{m^2} \sum_{i=1}^m \mathbf{Var} [Y_i] = \frac{m}{m^2} \mathbf{Var} [Y] = \frac{1}{m} (\mathcal{M}_{2\alpha-1}(\mathbf{p}) - \mathcal{M}_\alpha^2(\mathbf{p})) .$$

We use the following fact from [7] to find an upper bound for the variance of our empirical estimator.

Fact 4.14. ([7], Lemma 1) *For $\alpha > 1$ and $0 \leq \beta \leq \alpha$*

$$\mathcal{M}_{\alpha+\beta}(\mathbf{p}) \leq n^{(\alpha-1)(\alpha-\beta)/\alpha} \mathcal{M}_\alpha^2(\mathbf{p}) .$$

By taking $\beta = \alpha - 1$, we get

$$\begin{aligned} \sigma^2 &= \mathbf{Var} \left[\frac{1}{m} \sum_{i=1}^m Y_i \right] = \frac{1}{m} (\mathcal{M}_{2\alpha-1}(\mathbf{p}) - \mathcal{M}_\alpha^2(\mathbf{p})) \\ &\leq \frac{1}{m} (n^{(\alpha-1)/\alpha} \mathcal{M}_\alpha^2(\mathbf{p}) - \mathcal{M}_\alpha^2(\mathbf{p})) \\ &= \frac{1}{m} \mathcal{M}_\alpha^2(\mathbf{p}) (n^{1-1/\alpha} - 1) \\ &\leq \frac{\gamma^2}{100} \mathcal{M}_\alpha^2(\mathbf{p}) . \end{aligned}$$

We obtain a similar upper bound for the standard deviation of our empirical estimator,

$$\sigma = \sqrt{\mathbf{Var} \left[\frac{1}{m} \sum_{i=1}^m Y_i \right]} \leq \sqrt{\frac{\gamma^2}{100} \mathcal{M}_\alpha^2(\mathbf{p})} \leq \frac{\gamma}{10} \mathcal{M}_\alpha(\mathbf{p}) .$$

By Chebyshev's inequality we have

$$\begin{aligned} \Pr \left[\left| \frac{1}{m} \sum_{i=1}^m Y_i - \mathcal{M}_\alpha(\mathbf{p}) \right| > 10\sigma \right] &\leq \frac{1}{100} \Rightarrow \\ \Pr \left[\left| \frac{1}{m} \sum_{i=1}^m Y_i - \mathcal{M}_\alpha(\mathbf{p}) \right| \leq \gamma \mathcal{M}_\alpha(\mathbf{p}) \right] &\geq 0.99 \end{aligned}$$

Thus we can estimate $\mathcal{M}_\alpha(\mathbf{p})$ to a desired multiplicative accuracy with $O\left(\frac{n^{1-1/\alpha}}{(1-2^{(1-\alpha)\Delta})^2}\right)$ queries, which ends the proof. \square

Algorithm: ESTIMATOR VI

- Fix $\gamma = 1 - 2^{(1-\alpha)\Delta}$ and $m = \left\lceil \frac{100n^{1-1/\alpha}}{\gamma^2} \right\rceil$
- Draw m independent samples X_1, \dots, X_m from \mathbf{p} .
- Compute $Y_i = (\mathbf{p}[X_i])^{\alpha-1}$ by evaluating PMF on X_i for $1 \leq i \leq m$.
- Output $\frac{1}{1-\alpha} \log\left(\frac{1}{m} \sum_{i=1}^m Y_i\right)$.

Figure 4.1. SAMP+PMF Estimator of Rényi Entropy (of degree $\alpha > 1$)

4.4. Support Size

Definition 4.15. For a probability distribution \mathbf{p} , support size of \mathbf{p} is defined as $\text{supp}(\mathbf{p}) = |\{i : \mathbf{p}[i] \neq 0\}|$, the number of domain elements with nonzero probability.

Recall that $H_0(\mathbf{p}) = \log \text{supp}(\mathbf{p})$. The techniques described in Section 2.1 should be applicable to the distribution support size, since it is a symmetric property. In fact, Valiant and Valiant [19, 20] prove that for any positive constant $\epsilon < \frac{1}{4}$, estimating the support size of a distribution whose support members occur with probability at least $\frac{1}{n}$,⁹ to within $\pm\epsilon n$ requires $\Theta(\frac{n}{\log n})$ independent samples. In addition, Canonne and Rubinfeld [10] show that $\Theta(1/\epsilon^2)$ SAMP+PMF queries are necessary (and sufficient) for estimating $\text{supp}(\mathbf{p})$ with the guarantee that $\mathbf{p}[i] \geq \frac{1}{n}$ for all support members i , to within $\pm\epsilon n$. We modify their proof to establish a matching lower bound for this task in the SAMP+CDF model.

Theorem 4.16. $\Omega\left(\frac{1}{\epsilon^2}\right)$ SAMP+CDF queries are necessary to estimate (with high probability) the support size of an unknown distribution \mathbf{p} on domain $[n]$ to within $\pm\epsilon n$.

⁹It is typical to assume that all elements in the support occur with probability at least $\frac{1}{n}$; since without such a lower bound it is impossible to estimate support size.

Proof. Assume that there exists a program \mathcal{S} which can accomplish the task specified in the theorem statement with only $o\left(\frac{1}{\epsilon^2}\right)$ queries. Let us show how \mathcal{S} can be used to determine whether a given a coin is fair, or comes up heads with probability $p_2 = \frac{1}{2} + \lambda$.

Set $\epsilon = \frac{\lambda}{6}$, and let n be the smallest even number satisfying $n \geq 10/\epsilon^2$. Partition the domain $[n]$ into $M = \frac{n}{2}$ blocks I_1, \dots, I_M where $I_m = \{2m - 1, 2m\}$ for all $m \in [M]$. The construction of a probability distribution \mathbf{p} based on coin flips is as follows:

- Throw the coin M times, with outcomes X_1, \dots, X_M ,
- Set $\mathbf{p}[2m - 1] = \frac{2}{n}$ and $\mathbf{p}[2m] = 0$ if X_m is heads,
- Set $\mathbf{p}[2m - 1] = \mathbf{p}[2m] = \frac{1}{n}$ if X_m is tails, for each $m \in [M]$.

Note that by construction $\mathbf{p}[2m - 1] + \mathbf{p}[2m] = \frac{2}{n}$ for all $m \in [M]$. Let \mathcal{P}_1 and \mathcal{P}_2 be the families of distributions constructed by the above process, using the fair and biased coin, respectively. Let \mathbf{p}_1 (respectively \mathbf{p}_2) be a probability distribution randomly chosen from \mathcal{P}_1 (respectively \mathcal{P}_2). Then

$$\begin{aligned} \mathbf{E}[\text{supp}(\mathbf{p}_1)] &= n - M \frac{1}{2} = n \left(1 - \frac{1/2}{2}\right) = \frac{3}{4}n, \\ \mathbf{E}[\text{supp}(\mathbf{p}_2)] &= n - Mp_2 = n \left(1 - \frac{p_2}{2}\right) = n \left(\frac{3}{4} - \frac{\lambda}{2}\right) = \left(\frac{3}{4} - 3\epsilon\right)n, \end{aligned}$$

and via the additive Chernoff bound,

$$\begin{aligned} \Pr \left[\text{supp}(\mathbf{p}_1) \leq \frac{3}{4}n - \frac{\epsilon}{2}n \right] &\leq e^{-\frac{\epsilon^2 n}{2}} \leq e^{-5} < \frac{1}{1000} \\ \Pr \left[\text{supp}(\mathbf{p}_2) \geq \frac{3}{4}n - \frac{5\epsilon}{2}n \right] &\leq e^{-\frac{\epsilon^2 n}{2}} \leq e^{-5} < \frac{1}{1000}. \end{aligned}$$

In other words, the resulting distributions will satisfy (with high probability) $\text{supp}(\mathbf{p}_1) - \text{supp}(\mathbf{p}_2) > 2\epsilon n$, distant enough for \mathcal{S} to distinguish between two families.

As in our previous proofs, we could use \mathcal{S} (if only it existed) to distinguish between the two possible coin types by using the coin for an on-the-fly construction of \mathbf{p} . As before, SAMP and CDF queries are answered by picking a block randomly, throwing

the coin if the type of this block has not been fixed before, and returning the answer depending on the type of the block. Since $o\left(\frac{1}{\epsilon^2}\right) = o\left(\frac{1}{\lambda^2}\right)$ coin tosses would suffice for this task, we have reached a contradiction based on Lemma 4.2. \square

4.5. Tsallis Entropy

Tsallis entropy [31], defined as

$$S_\alpha(\mathbf{p}) = \frac{\mathbf{k}_B}{\alpha - 1} \left(1 - \sum_{i=1}^n (\mathbf{p}[i])^\alpha \right), \quad (4.8)$$

is a generalization of Boltzmann-Gibbs entropy where $\alpha \in \mathbb{R}$ and \mathbf{k}_B is the Boltzmann constant. Harvey *et al.* [17] gives an algorithm to estimate the Tsallis entropy which is also used to approximate the Shannon entropy in the most general streaming model. Without loss of generality we focus on additively estimating the quantity $T_\alpha(\mathbf{p}) := \frac{S_\alpha(\mathbf{p})}{\mathbf{k}_B}$ which appears to be a generalization of Shannon entropy, easily derived via L'Hôpital's rule.

Lemma 4.17. *For any number $\alpha > 1$, there exists an algorithm estimating (with high probability) the Tsallis entropy of an unknown distribution \mathbf{p} on $[n]$ to within $\pm\Delta$ with $O\left(\frac{1}{(\alpha-1)^2\Delta^2}\right)$ SAMP+PMF queries.*

Proof. Observe that for $\alpha > 1$

$$n^{1-\alpha} \leq \mathcal{M}_\alpha(\mathbf{p}) \leq 1 \quad \Rightarrow \quad 0 \leq T_\alpha(\mathbf{p}) \leq \frac{1}{\alpha-1} - \frac{n^{1-\alpha}}{\alpha-1}.$$

To estimate $T_\alpha(\mathbf{p})$ to within additive error Δ , one needs to estimate $\mathcal{M}_\alpha(\mathbf{p})$ to within $\gamma = (\alpha - 1)\Delta$. Note that $\Delta < \frac{1}{\alpha-1}$, therefore, $\gamma < 1$ must satisfy for achieving nontrivial approximations of $T_\alpha(\mathbf{p})$ and $\mathcal{M}_\alpha(\mathbf{p})$, respectively. We use the estimator constructed for Rényi entropy in Theorem 4.13.

Let $m = \lceil \frac{3}{(\alpha-1)^2 \Delta^2} \rceil$ and draw m independent samples X_1, \dots, X_m from \mathbf{p} . Define $Y_i = (\mathbf{p}[X_i])^{\alpha-1}$ where $\mathbf{p}[X_i]$ is computed using a PMF query on X_i for $1 \leq i \leq m$. Recall that

$$\mathbf{E}[Y_i] = \sum_j (\mathbf{p}[j])^{\alpha-1} = \sum_j (\mathbf{p}[j])^\alpha = \mathcal{M}_\alpha(\mathbf{p}).$$

Obviously, $\frac{1}{m} \sum_{i=1}^m Y_i$ is an unbiased estimator of $\mathcal{M}_\alpha(\mathbf{p})$. Observe that $Y_i \in [0, 1]$ since $\alpha > 1$. By applying the Hoeffding bound, we get

$$\Pr \left[\left| \frac{1}{m} \sum_{i=1}^m Y_i - \mathcal{M}_\alpha(\mathbf{p}) \right| > \gamma \right] \leq 2e^{-2\gamma^2 m}.$$

It easily follows that $m = O\left(\frac{1}{(\alpha-1)^2 \Delta^2}\right) = O\left(\frac{1}{\gamma^2}\right)$ queries are sufficient to bound the right-hand side of the inequality. \square

Jiao, Venkat and Weissman [32] construct an algorithm approximating Tsallis entropy $S_\alpha(\mathbf{p})$ of a distribution \mathbf{p} to within additive error using $O(n^{2/\alpha-1})$ SAMP queries for $1 < \alpha < 2$. More interestingly, they also prove that one requires only constant number of SAMP queries to additively estimate $S_\alpha(\mathbf{p})$ for $\alpha \geq 2$.¹⁰

¹⁰Acharya *et al.* [30] improve this result by constructing an algorithm for this task using only $O(1)$ samples for all $\alpha > 1$.

5. CONCLUSION

In this work, we investigate the task of additively estimating entropy in two settings based on two types of queries. A **SAMP** query takes no input and returns $x \in [n]$ with probability $\mathbf{p}[x]$; a **PMF** query takes as input $x \in [n]$ and returns the value $\mathbf{p}[x]$. The first setting is the **SAMP** model, where only **SAMP** queries are allowed. The second setting is the **SAMP+PMF** model in which both **SAMP** and **PMF** queries are utilized.

The motivation behind this work has both practical and theoretical reasons. Firstly, the **SAMP+PMF** model can be practical in many applications. For a concrete example, consider the Google n-gram database in which the frequency of each n-gram is published, and a random n-gram is easily obtained from the underlying text corpus. Secondly, the **SAMP+PMF** model is strongly related to the streaming model of computation which is an important field of computer science. Moreover, the **SAMP+PMF** model may illuminate the limitations of estimating entropy in the **SAMP** model.

We thoroughly analyzed the optimal bounds for estimating the Shannon entropy and the near-optimal bounds for estimating the Rényi entropy in the **SAMP** model. We described the exponentially faster algorithm constructed for estimating the Shannon entropy in the **SAMP+PMF** model. We established a matching lower bound for the estimation of the Shannon entropy $H(\mathbf{p})$ in the **SAMP+PMF** model, $\Omega(\log^2 n)$. We gave optimal bounds for the estimation of the Rényi entropy $H_\alpha(\mathbf{p})$ in the **SAMP+PMF** model, $\Theta(n^{1-1/\alpha})$.

Apparently, **PMF** queries provided no advantage in estimating Rényi entropy for integer $\alpha > 1$, whereas they *were* advantageous for noninteger $\alpha > 1$. However, the proximity between $\Omega(n^{1-1/\alpha})$ and $O(n)$, the lower and upper bound results for estimating Rényi entropy $H_\alpha(\mathbf{p})$ for noninteger $\alpha > 1$ in the **SAMP+PMF** and **SAMP** models, respectively, indicated the amount of the advantage a **PMF** query added to the model. In addition, the exponential gap between $O(\log^2 n)$ and $\Omega(n^{1-1/\alpha})$, the

upper and lower bound results for estimating Shannon entropy and Rényi entropy in the SAMP+PMF model, respectively, implied the difficulty of the latter problem.

We proved that the bounds were easily extended to the SAMP+CDF model, where SAMP and CDF queries (given x , return $\sum_{y \leq x} \mathbf{p}[y]$) were allowed. We gave a matching lower bound for estimating support size to within $\pm \epsilon n$ in the SAMP+CDF model. Lastly, we constructed an algorithm for additively estimating Tsallis entropy $S_\alpha(\mathbf{p})$ using a constant number of SAMP+PMF queries.

One problem left open by our work is that of optimal bounds for estimating the Rényi entropy $H_\alpha(\mathbf{p})$ in the SAMP+PMF model for $\alpha < 1$. The work [7] shows that in the model where only SAMP are allowed, $\tilde{\Omega}(n^{1/\alpha})$ queries are necessary when $0 < \alpha < 1$. It is interesting to ask whether there exists a sublinear algorithm for this task in the SAMP+PMF model.

Moreover, it is obvious that the SAMP+PMF model is superior to the SAMP model. However, degree of the superiority of the SAMP+PMF model is an open problem. In other words, how is the SAMP+PMF model affected if one has restricted number of such as $O(1)$ or $o(\log^2 n)$ PMF queries?

REFERENCES

1. Rubinfeld, R., “Taming Big Probability Distributions”, *XRDS: Crossroads, The ACM Magazine for Students*, Vol. 19, No. 1, pp. 24–28, 2012.
2. Canonne, C., *A Survey on Distribution Testing: Your Data is Big. But is It Blue?*, Tech. Rep. TR15-063, ECCC, 2015.
3. Paninski, L., “Estimation of Entropy and Mutual Information”, *Neural Computation*, Vol. 15, No. 6, pp. 1191–1253, 2003.
4. Paninski, L., “Estimating Entropy on m Bins Given Fewer than m Samples”, *IEEE Transactions on Information Theory*, Vol. 50, No. 9, pp. 2200–2203, 2004.
5. Valiant, P., “Testing Symmetric Properties of Distributions”, *SIAM Journal on Computing*, Vol. 40, No. 6, pp. 1927–1968, 2011.
6. Valiant, G. and P. Valiant, “Estimating the Unseen: An $n/\log(n)$ -Sample Estimator for Entropy and Support Size, Shown Optimal via New CLTs”, *Proceedings of the 43rd Annual ACM Symposium on Theory of Computing*, pp. 685–694, 2011.
7. Acharya, J., A. Orlitsky, A. T. Suresh and H. Tyagi, “The Complexity of Estimating Rényi Entropy”, *Proceedings of the 26th Annual ACM-SIAM Symposium on Discrete Algorithms*, 2015.
8. Kearns, M., Y. Mansour, D. Ron, R. Rubinfeld, R. Schapire and L. Sellie, “On the Learnability of Discrete Distributions”, *Proceedings of the 26th Annual ACM Symposium on Theory of Computing*, pp. 273–282, 1994.
9. Batu, T., S. Dasgupta, R. Kumar and R. Rubinfeld, “The Complexity of Approximating the Entropy”, *SIAM Journal on Computing*, Vol. 35, No. 1, pp. 132–150, 2005.

10. Canonne, C. and R. Rubinfeld, *Testing Probability Distributions Underlying Aggregated Data*, Tech. Rep. TR14-021, Electronic Colloquium on Computational Complexity, 2014.
11. Alon, N., Y. Matias and M. Szegedy, “The Space Complexity of Approximating the Frequency Moments”, *Journal of Computer and System Sciences*, Vol. 58, No. 1, pp. 137–147, 1999.
12. Guha, S., A. McGregor and S. Venkatasubramanian, “Streaming and Sublinear Approximation of Entropy and Information Distances”, *Proceedings of the 17th Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 733–742, ACM, 2006.
13. Lall, A., V. Sekar, M. Ogihara, J. Xu and H. Zhang, “Data Streaming Algorithms for Estimating Entropy of Network Traffic”, *Proceedings of ACM SIGMETRICS*, pp. 145–156, 2006.
14. Chakrabarti, A., K. Do Ba and S. Muthukrishnan, “Estimating Entropy and Entropy Norm on Data Streams”, *Internet Mathematics*, Vol. 3, No. 1, pp. 63–78, 2006.
15. Bhuvanagiri, L. and S. Ganguly, “Estimating Entropy over Data Streams”, *Proceedings of the 14th Annual European Symposium on Algorithms*, pp. 148–159, 2006.
16. Chakrabarti, A., G. Cormode and A. McGregor, “A Near-Optimal Algorithm for Computing the Entropy of a Stream”, pp. 328–335, 2007.
17. Harvey, N., J. Nelson and K. Onak, “Sketching and Streaming Entropy via Approximation Theory”, *Proceedings of the 49th Annual IEEE Symposium on Foundations of Computer Science*, pp. 489–498, 2008.
18. Caferov, C., B. Kaya, R. O’Donnell and A. C. C. Say, “Optimal Bounds for Estimating Entropy with PMF Queries”, *Proceedings of the 40th International Sym-*

- posium on Mathematical Foundations of Computer Science*, pp. 187–198, MFCS, 2015.
19. Valiant, G. and P. Valiant, *A CLT and Tight Lower Bounds for Estimating Entropy*, Tech. Rep. TR10-179, Electronic Colloquium on Computational Complexity, 2011.
 20. Valiant, G. and P. Valiant, *Estimating the Unseen: A Sublinear-Sample Canonical Estimator of Distributions*, Tech. Rep. TR10-180, Electronic Colloquium on Computational Complexity, 2010.
 21. Valiant, G. and P. Valiant, “The Power of Linear Estimators”, *Proceedings of the 52nd Annual IEEE Symposium on Foundations of Computer Science*, pp. 403–412, 2011.
 22. Shannon, C. E., “A Mathematical Theory of Communication”, *Bell System Technical Journal*, Vol. 27, No. 3, pp. 379—423, 1948.
 23. Shenkin, P. S., B. Erman and L. D. Mastrandrea, “Information-Theoretical Entropy as a Measure of Sequence Variability”, *Proteins*, Vol. 11, No. 4, pp. 297–313, 1991.
 24. Valiant, P., *Testing Symmetric Properties of Distributions*, Tech. Rep. TR07-135, Electronic Colloquium on Computational Complexity, 2007.
 25. Rényi, A., “On Measures of Entropy and Information”, *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*, pp. 547–561, 1961.
 26. Oorschot, P. C. v. and M. J. Wiener, “Parallel Collision Search with Cryptanalytic Applications”, *Journal of Cryptology*, Vol. 12, No. 1, pp. 1–28, 1999.
 27. Batu, T., L. Fortnow, R. Rubinfeld, W. D. Smith and P. White, “Testing Closeness

- of Discrete Distributions”, *Journal of the ACM*, Vol. 60, No. 1, pp. 1927–1968, 2013.
28. Paninski, L., “A Coincidence-Based Test for Uniformity Given Very Sparsely Sampled Discrete Data”, *IEEE Transactions on Information Theory*, Vol. 54, No. 10, pp. 4750–4755, 2008.
 29. Motahari, A. S., G. Bresler and D. N. Tse, “Information Theory of DNA Shotgun Sequencing”, *IEEE Transactions on Information Theory*, Vol. 59, No. 10, pp. 6273–6289, 2013.
 30. Acharya, J., A. Orlitsky, A. T. Suresh and H. Tyagi, *Estimating Rényi Entropy of Discrete Distributions*, Tech. Rep. arXiv:1408.1000v3, 2016.
 31. Tsallis, C., *Possible Generalization of Boltzmann-Gibbs Statistics*, Tech. Rep. CBPF-NF-062/87, CBPF, 1987.
 32. Jiao, J., K. Venkat and T. Weissman, “Maximum Likelihood Estimation of Functionals of Discrete Distributions”, *CoRR*, Vol. abs/1406.6959, 2014.

APPENDIX A: INEQUALITIES

Theorem A.1. (Chebyshev's Inequality) *Let X be a real-valued random variable such that $\mathbf{Var}[X]$ is well-defined. Then, $\forall t > 0$,*

$$\Pr[|X - \mathbf{E}[X]| > t] \leq \frac{\mathbf{Var}[X]}{t^2}.$$

Theorem A.2. (Jensen's Inequality) *Let X be an integrable random variable and $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ be a convex function. Then,*

$$\mathbf{E}[\varphi(X)] \geq \varphi(\mathbf{E}[X]).$$

Theorem A.3. (Pinsker's Inequality) *Let \mathbf{p}_1 and \mathbf{p}_2 be two probability distributions on the universe U . Then*

$$\mathbf{KL}(\mathbf{p}_1 \parallel \mathbf{p}_2) \geq \frac{1}{2 \ln 2} \cdot \|\mathbf{p}_1 - \mathbf{p}_2\|_1^2$$

Theorem A.4. (Exponential Inequality) *Let x be a real number. Then,*

$$1 + x \leq \left(1 + \frac{x}{n}\right)^n \leq e^x \quad \text{for } n > 1, |x| \leq n.$$

Theorem A.5. (Chernoff Bound) *Let X_1, \dots, X_m be m independent random variables that take on values in $[0, 1]$, where $\mathbf{E}[X_i] = p_i$, and $\sum_{i=1}^m p_i = P$. For any $\gamma \in (0, 1]$ we have*

$$\Pr\left[\sum_{i=1}^m X_i > (1 + \gamma)P\right] \leq e^{-\gamma^2 P/3}, \quad \Pr\left[\sum_{i=1}^m X_i < (1 - \gamma)P\right] \leq e^{-\gamma^2 P/2}.$$

Theorem A.6. (Hoeffding Bound)¹¹ Let X_1, \dots, X_m be m independent random variables that take on values in $[0, 1]$, where $\mathbf{E}[X_i] = p_i$, and $\sum_{i=1}^m p_i = P$. For any $\gamma \in (0, 1]$ we have

$$\Pr \left[\sum_{i=1}^m X_i > P + \gamma m \right] \leq e^{-2\gamma^2 m}, \quad \Pr \left[\sum_{i=1}^m X_i < P - \gamma m \right] \leq e^{-2\gamma^2 m}.$$

¹¹Usually, the Hoeffding bound is referred to as the additive Chernoff bound, whereas Theorem A.5 is referred to as the multiplicative Chernoff bound.