On the Efficient Implementation of High Accuracy Optimality of Profile Maximum Likelihood

Moses Charikar Stanford University moses@cs.stanford.edu Zhihao Jiang Stanford University faebdc@stanford.edu Kirankumar Shiragur Stanford University shiragur@stanford.edu

Aaron Sidford Stanford University sidford@stanford.edu

Abstract

We provide an efficient unified plug-in approach for estimating symmetric properties of distributions given n independent samples. Our estimator is based on profile-maximum-likelihood (PML) and is sample optimal for estimating various symmetric properties when the estimation error $\epsilon\gg n^{-1/3}$. This result improves upon the previous best accuracy threshold of $\epsilon\gg n^{-1/4}$ achievable by polynomial time computable PML-based universal estimators [ACSS21] [ACSS20]. Our estimator reaches a theoretical limit for universal symmetric property estimation as [Han21] shows that a broad class of universal estimators (containing many well known approaches including ours) cannot be sample optimal for every 1-Lipschitz property when $\epsilon\ll n^{-1/3}$.

1 Introduction

Given n independent samples $y_1, ..., y_n \in \mathcal{D}$ from an unknown discrete distribution $\mathbf{p} \in \Delta^{\mathcal{D}}$ the problem of estimating properties of \mathbf{p} , e.g. entropy, distance to uniformity, support size and coverage are among the most fundamental in statistics and learning. Further, the problem of estimating *symmetric properties* of distributions \mathbf{p} (i.e. properties invariant to label permutations) are well studied and have numerous applications [Cha84, BF93, CCG+12, TE87, Für05, KLR99, PBG+01] DS13, [RCS+09], GTPB07, HHRB01].

Over the past decade, symmetric property estimation has been studied extensively and there have been many improvements to the time and sample complexity for estimating different properties, e.g. support [VV11b] WY15], coverage [ZVV+16] OSW16], entropy [VV11b] WY16 [VVHV15], and distance to uniformity [VV11a] [HW16]. Towards unifying the attainment of computationally-efficient, sample-optimal estimators a striking work of [ADOS17] provided a universal plug-in approach based on a (approximate) profile maximum likelihood (PML) distribution, that (approximately) maximizes the likelihood of the observed profile (i.e. multiset of observed frequencies).

Formally, [ADOS17] showed that given $y_1, ..., y_n$ if there exists an estimator for a symmetric property f achieving accuracy ϵ and failure probability δ , then this PML-based plug-in approach achieves error 2ϵ with failure probability $\delta \exp{(3\sqrt{n})}$. As the failure probability δ for many estimators for well-known properties (e.g. support size and coverage, entropy, and distance to uniformity) is roughly $\exp{(-\epsilon^2 n)}$, this result implied a sample optimal unified approach for estimating these properties when the estimation error $\epsilon \gg n^{-1/4}$.

This result of [ADOS17] laid the groundwork for a line of work on the study of computational and statistical aspects of PML-based approaches to symmetric property estimation. For example, follow up work of [HS21] improved the analysis of [ADOS17] and showed that the failure probability of PML is at most $\delta^{1-c} \exp(-n^{1/3+c})$, for any constant c>0 and therefore it is sample optimal in the regime $\epsilon\gg n^{-1/3}$. The condition $\epsilon\gg n^{-1/3}$ on the optimality of PML is tight [Han21], in the sense that, PML is known to be not sample optimal in the regime $\epsilon\ll n^{-1/3}$. In fact, no estimator (that obeys some mild conditions), is sample optimal for estimating all symmetric properties in the regime $\epsilon\ll n^{-1/3}$; see Section 2 after Theorem 2.6 for details.

We also remark that the statistical guarantees in [ADOS17, [HS21]] hold for any β -approximate PML for suitable values of β . In particular, [HS21]] showed that any β -approximate PML for $\beta > \exp(-n^{1-c'})$ and any constant c'>0, has a failure probability of $\delta^{1-c}\exp(n^{1/3+c}+n^{1-c'})$ for any constant c>0. These results further imply a sample optimal estimator in the regime $\epsilon \gg n^{-\min(1/3,c'/2)}$ for properties with failure probability less than $\exp(-\epsilon^2 n)$. Note that better approximation leads to a larger range of ϵ for which the estimator is sample optimal.

Regarding computational aspects of PML, [CSS19a] provided the first efficient algorithm with a non-trivial approximation guarantee of $\exp(-n^{2/3}\log n)$, which further implied a sample optimal universal estimator for $\epsilon\gg n^{-1/6}$. This was then improved by [ACSS21] which showed how to efficiently compute PML to higher accuracy of $\exp(-\sqrt{n}\log n)$ thereby achieving a sample optimal universal estimator in the regime $\epsilon\gg n^{-1/4}$. The current best polynomial time approximate PML algorithm by [ACSS20] achieves an accuracy of $\exp(-k\log n)$, where k is the number of distinct observed frequencies. Although this result achieves better instance based statistical guarantees, in the worst case it still only implies a sample optimal universal estimator in the regime $\epsilon\gg n^{-1/4}$.

In light of these results, a key open problem is to close the gap between the regimes $\epsilon\gg n^{-1/3}$ and $\epsilon\gg n^{-1/4}$, where the former is the regime in which PML based estimators are statistically optimal and the later is the regime where efficient PML based estimators exist. In this work we ask:

Is there an efficient approximate PML-based estimator that is sample optimal for $\epsilon \gg n^{-1/3}$.

In this paper, we answer this question in the affirmative. In particular, we give an efficient PML-based estimator that has failure probability at most $\delta^{1-c} \exp(n^{1/3+c}+n^{1-c'})$, and consequently is sample optimal in the regime $\epsilon\gg n^{-1/3}$. As remarked, this result is tight in the sense that PML and a broad class of estimators are known to be not optimal in the regime $\epsilon\ll n^{-1/3}$.

To obtain this result we depart slightly from the previous approaches in [ADOS17] [CSS19a] [ACSS21]. Rather than directly compute an approximate PML distribution we compute a weaker notion of approximation which we show suffices to get us the desired universal estimator. We propose a notion of a β -weak approximate PML distribution inspired by [HS21] and show that an $\exp(-n^{1/3}\log n)$ -weak approximate PML achieves the desired failure probability of $\delta^{1-c}\exp(n^{1/3+c})$ for any constant c>0. Further, we provide an efficient algorithm to compute an $\exp(-n^{1/3}\log n)$ -weak approximate PML distribution. Our paper can be viewed as an efficient algorithmic instantiation of [HS21].

Ultimately, our algorithms use the convex relaxation presented in [CSS19a ACSS21] and provide a new rounding algorithm. We differ from the previous best $\exp(-k \log n)$ approximate PML algorithm [ACSS20] only in the matrix rounding procedure which controls the approximation guarantee. At a high level, the approximation guarantee for the rounding procedure in [ACSS20] is exponential in the sum of matrix dimensions. In the present work, we need to round a rectangular matrix with an approximation exponential in the smaller dimension, which may be infeasible for arbitrary matrices. Our key technical innovation is to introduce a *swap* operation (see Section [4.1]) which facilitates such an approximation guarantee. In addition to a better approximation guarantee than [ACSS20], our algorithm also exhibits better run times (see Section [2]).

Organization: We introduce preliminaries in Section [1.1] In Section 2, we state our main results and also cover related work. In Section 3, we provide the convex relaxation to PML studied in [CSS19a, ACSS21]. Finally, in Section 4, we provide a proof sketch of our main computational result. Many proofs are then differed to the appendix.

 $^{^{1}\}beta$ -approximate PML is a distribution that achieves a multiplicative β -approximation to the PML objective.

1.1 Preliminaries

General notation: For matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{s \times t}$, we use $\mathbf{A} \leq \mathbf{B}$ to denote that $\mathbf{A}_{ij} \leq \mathbf{B}_{ij}$ for all $i \in [s]$ and $j \in [t]$. We let [a,b] and $[a,b]_{\mathbb{R}}$ denote the interval $\geq a$ and $\leq b$ of integers and reals respectively. We use $\widetilde{O}(\cdot)$, $\widetilde{\Omega}(\cdot)$ notation to hide all polylogarithmic factors in n and N. We let $a_n \gg b_n$ to denote that $a_n \in \Omega(b_n n^c)$ or $b_n \in O(n^{-c}a_n)$, for some small constant c > 0.

Throughout this paper, we assume we receive a sequence of n independent samples from a distribution $\mathbf{p} \in \Delta^{\mathcal{D}}$, where $\Delta^{\mathcal{D}} \stackrel{\mathrm{def}}{=} \{ \mathbf{q} \in [0,1]^{\mathcal{D}}_{\mathbb{R}} | \| q \|_1 = 1 \}$ is the set of all discrete distributions supported on domain \mathcal{D} . Let \mathcal{D}^n be the set of all length n sequences of elements of \mathcal{D} and for $y^n \in \mathcal{D}^n$ let y^n_i denoting its ith element. Let $\mathbf{f}(y^n,x) \stackrel{\mathrm{def}}{=} | \{i \in [n] \mid y^n_i = x\} |$ and \mathbf{p}_x be the frequency and probability of $x \in \mathcal{D}$ respectively. For a sequence $y^n \in \mathcal{D}^n$, let $\mathbf{M} = \{\mathbf{f}(y^n,x)\}_{x \in \mathcal{D}} \setminus \{0\}$ be the set of all its non-zero distinct frequencies and $\mathbf{m}_1, \mathbf{m}_2, \ldots, \mathbf{m}_{|\mathbf{M}|}$ be these distinct frequencies. The profile of a sequence y^n denoted $\phi = \Phi(y^n)$ is a vector in $\mathbb{Z}^{|\mathbf{M}|}$, where $\phi_j \stackrel{\mathrm{def}}{=} | \{x \in \mathcal{D} \mid \mathbf{f}(y^n,x) = \mathbf{m}_j\} |$ is the number of domain elements with frequency \mathbf{m}_j . We call n the length of profile ϕ and let Φ^n denote the set of all profiles of length n. The probability of observing sequence y^n and profile ϕ for distribution \mathbf{p} are $\mathbb{P}(\mathbf{p},y^n) = \prod_{x \in \mathcal{D}} \mathbf{p}_x^{\mathbf{f}(y^n,x)}$ and $\mathbb{P}(\mathbf{p},\phi) = \sum_{\{y^n \in \mathcal{D}^n \mid \Phi(y^n) = \phi\}} \mathbb{P}(\mathbf{p},y^n)$.

Profile maximum likelihood: A distribution $\mathbf{p}_{\phi} \in \Delta^{\mathcal{D}}$ is a *profile maximum likelihood* (PML) distribution for profile $\phi \in \Phi^n$ if $\mathbf{p}_{\phi} \in \operatorname{argmax}_{\mathbf{p} \in \Delta^{\mathcal{D}}} \mathbb{P}(\mathbf{p}, \phi)$. Further, a distribution $\mathbf{p}_{\phi}^{\beta}$ is a β -approximate PML distribution if $\mathbb{P}(\mathbf{p}_{\phi}^{\beta}, \phi) \geq \beta \cdot \mathbb{P}(\mathbf{p}_{\phi}, \phi)$. For a distribution \mathbf{p} and a length n, let \mathbf{X} be a random variable that takes value $\phi \in \Phi^n$ with probability $\mathbb{P}(\mathbf{p}, \phi)$. We call $H(\mathbf{X})$ (entropy of \mathbf{X}) the *profile entropy* with respect to (\mathbf{p}, n) and denote it by $H(\Phi^n, \mathbf{p})$.

Probability discretization: Let $\mathbf{R} \stackrel{\mathrm{def}}{=} \{\mathbf{r}_i\}_{i \in [1,\ell]}$ be a finite discretization of the probability space, where $\mathbf{r}_i \in [0,1]_{\mathbb{R}}$ and $\ell \stackrel{\mathrm{def}}{=} |\mathbf{R}|$. We call $\mathbf{q} \in [0,1]_{\mathbb{R}}^{\mathcal{D}}$ a pseudo-distribution if $\|\mathbf{q}\|_1 \leq 1$ and a discrete pseudo-distribution with respect to \mathbf{R} if all its entries are in \mathbf{R} as well. We use $\Delta_{\mathrm{pseudo}}^{\mathcal{D}}$ and $\Delta_{\mathbf{R}}^{\mathcal{D}}$ to denote the set of all pseudo-distributions and discrete pseudo-distributions with respect to \mathbf{R} respectively. In our work, we use the following most commonly used [CSS19a] [ACSS21] [ACSS20] probability discretization set. For any $\alpha > 0$,

$$\mathbf{R}_{n,\alpha} \stackrel{\text{def}}{=} \{1\} \cup \left\{ \frac{1}{2n^2} (1 + n^{-\alpha})^i \mid \text{for all } i \in \mathbb{Z}_{\geq 0} \text{ such that } \frac{1}{2n^2} (1 + n^{-\alpha})^i \leq 1 \right\} . \tag{1}$$

For all probability terms defined involving distributions \mathbf{p} , we extend those definitions to pseudo distributions \mathbf{q} by replacing \mathbf{p}_x with \mathbf{q}_x everywhere.

See ?? for the definition of an estimator and optimal sample complexity.

2 Results

Here we provide our main results. In our first result (Theorem 2.2), we show that a weaker notion of approximate PML suffices to obtain the desired universal estimator. Later we show that these weaker approximate PML distributions can be efficiently computed (Theorem 2.3).

Definition 2.1. Given a profile ϕ , we call a distribution $\mathbf{p}' \in \Delta^{\mathcal{D}}$ β -approximate PML distribution with respect to \mathbf{R} if $\mathbb{P}(\mathbf{p}', \phi) \geq \beta \cdot \max_{\mathbf{q} \in \Delta^{\mathcal{D}}_{\mathbf{R}}} \mathbb{P}\left(\frac{\mathbf{q}}{\|\mathbf{q}\|_{1}}, \phi\right)$.

The above definition generalizes β -approximate PML distributions which is simply the special case when $\mathbf{R} = [0,1]_{\mathbb{R}}$. Using our new definition, we show that for a specific choice of the discretization set $\mathbf{R}_{n,1/3}$, a distribution \mathbf{p}' that is an approximate PML with respect to $\mathbf{R}_{n,1/3}$ suffices to obtain a universal estimator; this result is formally stated below.

Theorem 2.2 (Competitiveness of an approximate PML w.r.t **R**). For symmetric property f, suppose there exists an estimator \hat{f} that takes input a profile $\phi \in \Phi^n$ drawn from $\mathbf{p} \in \Delta^{\mathcal{D}}$ and satisfies,

$$\mathbb{P}\left(|f(\mathbf{p}) - \widehat{f}(\phi)| \ge \epsilon\right) \le \delta ,$$

then for $\mathbf{R} = \mathbf{R}_{n,1/3}$ (See Equation (1)), a discrete pseudo distribution $\mathbf{q}' \in \Delta_{\mathbf{R}}^{\mathcal{D}}$ such that $\mathbf{q}'/\|\mathbf{q}'\|_1$ is an $\exp(-O(|\mathbf{R}|\log n))$ -approximate PML distribution with respect to the \mathbf{R} satisfies,

$$\mathbb{P}\left(\left|f\left(\frac{\boldsymbol{q}'}{\|\boldsymbol{q}'\|_1}\right) - f(\boldsymbol{p})\right| \ge 2\epsilon\right) \le \delta^{1-c} \exp(O(n^{1/3+c})), \quad \text{for any constant } c > 0. \tag{2}$$

The proof of the above theorem is implicit in the analysis of [HS21], however we provide a short simpler proof using their continuity lemma (Lemma 2 in [HS21]). Note that the bound on the failure probability we get is the same asymptotically as that of exact PML from [HS21], which is known to be tight [Han21]. Furthermore, to achieve such an improved failure probability bound all we need is an approximate PML distribution with respect to **R**, for some **R** which is of small size. Taking advantage of this fact and building upon [CSS19a, [ACSS21]], we provide a new rounding algorithm that outputs the desired approximate PML distribution with respect to **R**.

Theorem 2.3 (Computation of an approximate PML w.r.t **R**). We provide an algorithm that given a probability discretization set $\mathbf{R} = \mathbf{R}_{n,\alpha}$ for $\alpha > 0$ (See Equation (1)) and a profile ϕ with k distinct frequencies, runs in time $\widetilde{O}\left(|\mathbf{R}| + \frac{n}{\min(k,|\mathbf{R}|)}\left(\min(|\mathbf{R}|,n/k)k^{\omega} + \min(|\mathbf{R}|,k)k^2\right)\right)$, where $\omega < 2.373$ is the current matrix multiplication constant [Will2] Gal14] AW21] and returns a pseudo distribution $\mathbf{q}' \in \Delta^{\mathbf{D}}_{\mathbf{R}}$ such that,

$$\mathbb{P}\left(\frac{\pmb{q}'}{\|\pmb{q}'\|_1},\phi\right) \geq \exp\left(-O(\min(k,|\pmb{R}|)\log n)\right) \cdot \max_{\pmb{q} \in \Delta_{\pmb{R}}^{\mathcal{D}}} \mathbb{P}\left(\frac{\pmb{q}}{\|\pmb{q}\|_1},\phi\right) \; .$$

When $\mathbf{R} = \mathbf{R}_{n,1}$, our algorithm computes an $\exp(-O(k \log n))$ approximate PML distribution, therefore our result is at least as good as the previous best known approximate PML algorithm due to [ACSS20]. In comparison to [ACSS20], our rounding algorithm is simpler and we suspect, more practical. We provide a more detailed comparison to it later in this section.

Applications: Our main results have several applications which we discuss here. First note that, combining Theorem 2.2 and 2.3 immediately yields the following corollary.

Corollary 2.4 (Efficient unified estimator). Given a profile $\phi \in \Phi^n$ with k distinct frequencies, we can compute an approximate PML distribution \mathbf{q}' that satisfies Equation (2) in Theorem 2.2 in time $\widetilde{O}\left(\frac{n}{\min(k,n^{1/3})}\left(\min(n^{1/3},n/k)k^\omega+\min(n^{1/3},k)k^2\right)\right)$.

For many symmetric properties the failure probability is exponentially small as stated below.

Lemma 2.5 (Lemma 2 in [ADOS17], Theorem 3 in [HS21]). For distance to uniformity, entropy, support size and coverage, and sorted ℓ_1 distance there exists an estimator that is sample optimal and the failure probability is at most $\exp(-\epsilon^2 n^{1-\alpha})$ for any constant $\alpha > 0$.

The above result combined with Corollary 2.4, immediately yields the following theorem.

Theorem 2.6 (Efficient sample optimal unified estimator). There exists an efficient approximate PML-based estimator that for $\epsilon \gg n^{-1/3}$ and symmetric properties such as, distance to uniformity, entropy, support size and coverage, and sorted ℓ_1 distance achieves optimal sample complexity and has failure probability upper bounded by $\exp(-n^{1/3})$.

As our work computes an $\exp(-O(k \log n))$ approximate PML, we recover efficient version of Lemma 2.3 and Theorem 2.4 from [ACSS20]. The first result uses $\exp(-O(k \log n))$ approximate PML algorithm to efficiently implement an estimator that has better statistical guarantees based on profile entropy [HO20] (See Section 1.1). The second result provides an efficient implementation of the PseudoPML estimators [CSS19b] [HO19]. Please refer to the respective papers for further details.

Tightness of our result: Recall that [HS21] showed that the failure probability of an (approximate) PML based estimator is upper bounded by $\delta^{1-c} \exp(-n^{1/3+c})$, for any constant c>0. This result further implied a sample optimal universal estimator in the regime $\epsilon\gg n^{-1/3}$ for various symmetric properties (Theorem 2.6). In our work, we efficiently recover these results and a natural question to ask here is if these results can be improved.

As remarked earlier, [Han21] showed that the condition for optimality of PML ($\epsilon \gg n^{-1/3}$) is in some sense tight. More formally, they showed that PML is not sample optimal in estimating every

1-Lipschitz property in the regime $\epsilon \ll n^{-1/3}$. In fact, the results in [Han21] hold more broadly for any universal plug-in based estimator that outputs a distribution $\hat{\bf p}$ satisfying.

$$\max_{\mathbf{p} \in \Delta^{\mathcal{D}}} \mathbb{E} \|\mathbf{p} - \hat{\mathbf{p}}\|_1^{\text{sorted}} \le A(n) \sqrt{k/n} \;,$$

where $A(n) \leq n^{\gamma}$ for every $\gamma > 0$ and $\|\mathbf{p} - \mathbf{q}\|_1^{\text{sorted}} \stackrel{\text{def}}{=} \min_{\text{permutations }\sigma} \|\mathbf{p} - \mathbf{q}_{\sigma}\|_1$ denotes the sorted ℓ_1 distance between \mathbf{p} and \mathbf{q} . In other words, if an estimator is based on a reasonably good estimate of the true distribution \mathbf{p} (in terms of sorted- ℓ_1 distance), then it cannot be sample optimal for every 1-Lipschitz property. Furthermore, many well-known universal estimators including PML and LLM $\|\mathbf{HJW18}\|$ indeed provide a reasonably good estimate of the true distribution and therefore cannot be sample optimal in the regime $\epsilon \ll n^{-1/3}$. Please refer to $\|\mathbf{Han21}\|$ for further details.

Comparison to approximate PML algorithms: All prior provable approximate PML algorithms [CSS19a, ACSS21] [ACSS20] have two key steps: (Step 1) solve a convex approximation to the PML and (Step 2) round the (fractional) solution to a valid approximate PML distribution.

A convex approximation to PML was first provided in [CSS19a] and a better analysis for it is shown in [ACSS21]. In particular, [CSS19a] and [ACSS21] showed that an integral optimal solution to step 1 approximates the PML up to accuracy $\exp(-n^{2/3}\log n)$ and $\exp(-\min(k,|\mathbf{R}|)\log n)$ respectively, where k and $|\mathbf{R}|$ are the number of distinct frequencies and distinct probability values respectively. In addition to the loss from convex approximation, the previous algorithms also incurred a loss in the rounding step (Step 2). The loss in the rounding step for the previous works is bounded by $\exp(-n^{2/3}\log n)$ [CSS19a], $\exp(-\sqrt{n}\log n)$ [ACSS21] and $\exp(-k\log n)$ [ACSS20].

In our work, we show that there exists a choice of \mathbf{R} (= $\mathbf{R}_{n,1/3}$) that is of small size ($|\mathbf{R}| \le n^{1/3}$) and suffices to get the desired universal estimator. As $|\mathbf{R}| \le n^{1/3}$, our approach only incur a loss of $\exp(-\min(k,|\mathbf{R}|)\log n) \in \exp(-n^{1/3}\log n)$ in the convex approximation step (Step 1). Furthermore for the rounding step (Step 2), we provide a new simpler and a practical rounding algorithm with a better approximation loss of $\exp(-O(\min(k,|\mathbf{R}|)\log n)) \in \exp(-O(n^{-1/3}\log n))$.

Regarding the run times, both [ACSS20] and ours have run times of the form $\mathcal{T}_{\mathrm{solve}} + \mathcal{T}_{\mathrm{sparsify}} + \mathcal{T}_{\mathrm{round}}$, where the terms correspond to the time required to solve the convex program, sparsify and round a solution. In our algorithm, we pay the same cost as [ACSS20] for the first two steps but our run time guarantees are superior to theirs in the rounding step. In particular, the run time of [ACSS20] is shown as a large polynomial and perhaps not practical as their approach requires enumerating all the approximate min cuts. In contrast, our algorithm has a run time that is subquadratic.

Other related work PML was introduced by OSS^+04 . Many heuristic approaches have been proposed to compute approximate PML, such as the EM algorithm in OSS^+04 , an algebraic approaches in ADM^+10 , Bethe approximation in OSS^+04 and OSS^+04 , and a dynamic programming approach in OSS^+04 . For the broad applicability of PML in property testing and to estimate other symmetric properties please refer to OSS^+04 . Please refer to OSS^+04 and to estimate other symmetric properties please refer to OSS^+04 . Please refer to OSS^+04 have designing approaches for designing universal estimators are: OSS^+04 have designing approaches approaches for designing universal estimators are: OSS^+04 have designing approaches approaches in OSS^+04 have designing approaches in OSS^+04 have designed approaches in OSS^+04 have designing approaches in OSS^+04 have designing approaches in OSS^+04 have designed a

Limitations of our work One of the limitations of all the provable approximate PML algorithms [CSS19a] [ACSS21] [ACSS20] (including ours) is that they require the solution of a convex program that approximates the PML objective and all these previous works use the CVX solver which is not practical for large sample instances; note that our results hold for small error regimes which lead to such large sample instances. Therefore, designing a practical algorithm to solve the convex program is an important future research direction. As discussed above, local moment matching (LLM) based approach is another universal approach for property estimation. It is unclear which of the two (PML or LLM) can lead to practical algorithms.

3 Convex relaxation to PML

Here we restate the convex program from [CSS19a] that approximates the PML objective. The current best analysis of this convex program is in [ACSS21]. We first describe the notation and later state several results from [CSS19a] [ACSS21] that capture the guarantees of the convex program.

Notation: For any matrices $\mathbf{X} \in \mathbb{R}^{a \times c}$ and $\mathbf{Y} \in \mathbb{R}^{b \times c}$, we let $\mathrm{concat}(\mathbf{X}, \mathbf{Y})$ denote the matrix $\mathbf{W} \in \mathbb{R}^{(a+b) \times c}$ that satisfies, $\mathbf{W}_{i,j} = \mathbf{X}_{i,j}$ for all $i \in [a]$ and $j \in [c]$ and $\mathbf{W}_{a+i,j} = \mathbf{Y}_{ij}$ for all $i \in [b]$ and $j \in [c]$. Recall we let $\mathbf{R} \stackrel{\mathrm{def}}{=} \{\mathbf{r}_i\}_{i \in [\ell]}$ be a finite discretization of the probability space, where $\mathbf{r}_i \in [0,1]_{\mathbb{R}}$ and $\ell \stackrel{\mathrm{def}}{=} |\mathbf{R}|$. Let $\mathbf{r} \in [0,1]_{\mathbb{R}}^{\ell}$ be a vector whose i'th element is equal to \mathbf{r}_i .

Lemma 3.1 (Lemma 4.4 in [CSS19a]). Let $\mathbf{R} = \mathbf{R}_{n,\alpha}$ for some $\alpha > 0$. For any profile $\phi \in \Phi^n$ and distribution $\mathbf{p} \in \Delta^{\mathcal{D}}$, there exists a pseudo distribution $\mathbf{q} \in \Delta^{\mathcal{D}}_{\mathbf{R}}$ that satisfies $\mathbb{P}(\mathbf{p}, \phi) \geq \mathbb{P}(\mathbf{q}, \phi) \geq \exp(-\alpha n - 6) \mathbb{P}(\mathbf{p}, \phi)$ and therefore,

$$\max_{\boldsymbol{p} \in \Delta^{\mathcal{D}}} \mathbb{P}(\boldsymbol{p}, \phi) \geq \max_{\boldsymbol{q} \in \Delta^{\mathcal{D}}_{\boldsymbol{p}}} \mathbb{P}(\boldsymbol{q}, \phi) \geq \exp\left(-\alpha n - 6\right) \max_{\boldsymbol{p} \in \Delta^{\mathcal{D}}} \mathbb{P}(\boldsymbol{p}, \phi) \;.$$

For any probability discretization set **R**, profile ϕ and pseudo distribution $\mathbf{q} \in \Delta_{\mathbf{R}}^{\mathcal{D}}$, define:

$$\mathbf{Z}_{\mathbf{R}}^{\phi} \stackrel{\text{def}}{=} \left\{ \mathbf{X} \in \mathbb{R}_{\geq 0}^{\ell \times [0,k]} \mid \mathbf{X} \mathbf{1} \in \mathbb{Z}^{\ell}, [\mathbf{X}^{\top} \mathbf{1}]_{j} = \phi_{j} \text{ for all } j \in [1,k] \text{ and } \mathbf{r}^{\top} \mathbf{X} \mathbf{1} \leq 1 \right\}, \tag{3}$$

$$\mathbf{Z}_{\mathbf{R}}^{\phi, \text{frac}} \stackrel{\text{def}}{=} \left\{ \mathbf{X} \in \mathbb{R}_{\geq 0}^{\ell \times [0, k]} \mid [\mathbf{X}^{\top} \mathbf{1}]_{j} = \phi_{j} \text{ for all } j \in [1, k] \text{ and } \mathbf{r}^{\top} \mathbf{X} \mathbf{1} \leq 1 \right\}. \tag{4}$$

The j'th column corresponds to frequency m_j and we use $\mathbf{m}_0 \stackrel{\text{def}}{=} 0$ to capture the unseen elements. Without loss of generality, we assume $\mathbf{m}_0 < \mathbf{m}_1 < \cdots < \mathbf{m}_k$. Let $\mathbf{C}_{ij} \stackrel{\text{def}}{=} \mathbf{m}_j \log \mathbf{r}_i$ for all $i \in [\ell]$ and $j \in [0, k]$. The objective of the optimization problem is follows: for any $\mathbf{X} \in \mathbb{R}_{>0}^{\ell \times [0, k]}$ define,

$$\mathbf{g}(\mathbf{X}) \stackrel{\text{def}}{=} \exp\left(\sum_{i \in [\ell], j \in [0, k]} \left[\mathbf{C}_{ij} \mathbf{X}_{ij} - \mathbf{X}_{ij} \log \mathbf{X}_{ij}\right] + \sum_{i \in [\ell]} \left[\mathbf{X}\mathbf{1}\right]_i \log\left[\mathbf{X}\mathbf{1}\right]_i\right). \tag{5}$$

For any $\mathbf{q} \in \Delta^{\mathcal{D}}_{\mathbf{R}}$, the function $\mathbf{g}(\mathbf{X})$ approximates the $\mathbb{P}(\mathbf{q}, \phi)$ term and is stated below.

Lemma 3.2 (Theorem 6.7 and Lemma 6.9 in ACSS21). Let \mathbf{R} be a probability discretization set. For any profile $\phi \in \Phi^n$ with k distinct frequencies the following statements hold for $\alpha = \min(k, |\mathbf{R}|) \log n$: $\exp(-O(\alpha)) \cdot C_\phi \cdot \max_{\mathbf{X} \in \mathbf{Z}_R^\phi} \mathbf{g}(\mathbf{X}) \leq \max_{\mathbf{q} \in \Delta_R^\mathcal{D}} \mathbb{P}(\mathbf{q}, \phi) \leq \exp(O(\alpha)) \cdot C_\phi \cdot \max_{\mathbf{X} \in \mathbf{Z}_R^\phi} \mathbf{g}(\mathbf{X})$ and $\max_{\mathbf{q} \in \Delta_R^\mathcal{D}} \mathbb{P}(\mathbf{q}, \phi) \leq \exp(O(\min(k, |\mathbf{R}|) \log n)) \cdot C_\phi \cdot \max_{\mathbf{X} \in \mathbf{Z}_R^\phi, \operatorname{frac}} \mathbf{g}(\mathbf{X})$, where $C_\phi \stackrel{\text{def}}{=} \frac{n!}{\prod_{j \in [1, k]} (m_j!)^{\phi_j}}$ is a term that only depends on the profile.

The proof of concavity for the function g(X) and a running time analysis to solve the convex program are provided in [CSS19a]. For any $X \in \mathbf{Z}_{\mathbf{R}}^{\phi}$, a pseudo-distributions associated with it is defined below.

Definition 3.3. For any $\mathbf{X} \in \mathbf{Z}_{\mathbf{R}}^{\phi}$, the discrete pseudo-distribution $\mathbf{q}_{\mathbf{X}}$ associated with \mathbf{X} and \mathbf{R} is defined as follows: for arbitrary $[\mathbf{X}\mathbf{1}]_i$ number of domain elements assign probability \mathbf{r}_i . Further $\mathbf{p}_{\mathbf{X}} \stackrel{\mathrm{def}}{=} \mathbf{q}_{\mathbf{X}}/\|\mathbf{q}_{\mathbf{X}}\|_1$ is the distribution associated with \mathbf{X} and \mathbf{R} .

Note that $\mathbf{q}_{\mathbf{X}}$ is a valid pseudo-distribution because of the third condition in Equation (3) and these pseudo distributions $\mathbf{p}_{\mathbf{X}}$ and $\mathbf{q}_{\mathbf{X}}$ satisfy the following lemma.

Lemma 3.4 (Theorem 6.7 in [ACSS21]). Let \mathbf{R} and $\phi \in \Phi^n$ be a probability discretization set and a profile with k distinct frequencies. For any $\mathbf{X} \in \mathbf{Z}_{\mathbf{R}}^{\phi}$, the discrete pseudo distribution $\mathbf{q}_{\mathbf{X}}$ and distribution $\mathbf{p}_{\mathbf{X}}$ associated with \mathbf{X} and \mathbf{R} satisfy: $\exp\left(-O(k\log n)\right)C_{\phi} \cdot \mathbf{g}(\mathbf{X}) \leq \mathbb{P}(\mathbf{q}_{\mathbf{X}},\phi) \leq \mathbb{P}(\mathbf{p}_{\mathbf{X}},\phi)$.

²The theorem statement in [ACSS21] is only written with an approximation factor of $\exp(O(k\log n))$. However, their proof provides a stronger approximation factor which is upper bounded by the non-negative rank of the probability matrix, which in turn is upper bounded by the minimum of distinct frequencies and distinct probabilities. Therefore the theorem statement in [ACSS21] holds with a much stronger approximation guarantee of $\exp(O(\min(k, |\mathbf{R}|) \log n))$.

4 Approximate PML algorithm

Here we provide a proof sketch of Theorem [2.3] and provide a rounding algorithm that proves it. Our rounding algorithm takes as input a matrix $\mathbf{X} \in \mathbf{Z}_{\mathbf{R}}^{\phi,\mathrm{frac}}$ which may have fractional row sums and round it to integral values. This new rounded matrix $\mathbf{X}_{\mathrm{final}}$ corresponds to our approximate PML distribution (See Definition [3.3]). The description of our algorithm is as follows.

Algorithm 1 Approximate PML(ϕ , $\mathbf{R} = \mathbf{R}_{n,\alpha}$)

- 1: Let \mathbf{X} be any solution that satisfies, $\log \mathbf{g}(\mathbf{X}) \geq \max_{\mathbf{Y} \in \mathbf{Z}^{\phi, \mathrm{frac}}_{\mathbf{R}}} \log \mathbf{g}(\mathbf{Y}) O\left(\min(k, |\mathbf{R}|) \overline{\log n}\right)$.
- 2: $\mathbf{X}' = \text{sparsify}(\mathbf{X})$.
- 3: $(\mathbf{A}, \mathbf{B}) = \text{swapmatrixround}(\mathbf{X}')$.
- 4: $(\mathbf{X}_{\text{final}}, \mathbf{R}_{\text{final}}) = \text{create}(\mathbf{A}, \mathbf{B}, \mathbf{R})$
- 5: Let \mathbf{p}' be the distribution with respect to $\mathbf{X}_{\text{final}}$ and $\mathbf{R}_{\text{final}}$ (See Definition 3.3).
- 6: Return $\mathbf{q} = \text{discretize}(\mathbf{p}', \phi, \mathbf{R})$

We now provide a guarantee for each of these lines of Algorithm $\boxed{1}$ We later use these guarantees to prove our final theorem $\boxed{2.3}$. The guarantees of the approximate maximizer \mathbf{X} computed in the first step of the algorithm are summarized in the following lemma.

Lemma 4.1 ([CSS19a, ACSS21]). Line [I] of the algorithm can be implemented in $\widetilde{O}(|\mathbf{R}|k^2 + |\mathbf{R}|^2k)$ time and the approximate maximizer \mathbf{X} satisfies: $C_{\phi} \cdot \mathbf{g}(\mathbf{X}) \geq \exp\left(-O\left(\min(k,|\mathbf{R}|)\log n\right)\right) \max_{\mathbf{q} \in \Delta_{\mathcal{P}}^{\mathcal{P}}} \mathbb{P}(\mathbf{q},\phi)$.

The guarantees of the second step of our algorithm are summarized in the following lemma. Please refer to [ACSS20] for the description of the procedure sparsify. We use this procedure so that we can assume $|\mathbf{R}| \le k+1$ as we can ignore the zero rows of the matrix \mathbf{X} .

Lemma 4.2 (Lemma 4.3 in [ACSS20]). For any $X \in \mathbb{Z}_R^{\phi, \text{frac}}$, the algorithm sparsify (X) runs in time $\widetilde{O}(|R| \ k^{\omega})$ and outputs $X' \in \mathbb{Z}_R^{\phi, \text{frac}}$ such that: $g(X') \geq g(X)$ and $|\{i \in [\ell] \mid [X'\overrightarrow{1'}]_i > 0\}| \leq k+1$.

To explain our next step, we need to define a new operation called the swap.

Definition 4.3. Given a matrix **A**, indices $i_1 < i_2$, $j_1 < j_2$ and a parameter $\epsilon \ge 0$, the operation swap $(\mathbf{A}, i_1, i_2, j_1, j_2, \epsilon)$ outputs a matrix **A**' that satisfies,

$$\mathbf{A}'_{ij} = \begin{cases} \mathbf{A}_{i,j} + \epsilon \text{ for } i = i_1, \ j = j_1 \\ \mathbf{A}_{i,j} - \epsilon \text{ for } i = i_2, \ j = j_1 \end{cases} \qquad \mathbf{A}_{i,j} - \epsilon \text{ for } i = i_1, \ j = j_2, \\ \mathbf{A}_{i,j} - \epsilon \text{ for } i = i_2, \ j = j_1 \end{cases} \qquad \mathbf{A}_{i,j} + \epsilon \text{ for } i = i_2, \ j = j_2,$$

$$(6)$$

$$\mathbf{A}_{i,j} \text{ otherwise.}$$

Definition 4.4 (Swap distance). \mathbf{A}' is x-swap distance from \mathbf{A} , if \mathbf{A}' can be obtained from \mathbf{A} through a sequence of swap operations and the summation of the value ϵ 's in these operations is at most x, i.e. there is a set of parameters $\{(i_1^{(s)},i_2^{(s)},j_1^{(s)},j_2^{(s)},\epsilon^{(s)})\}_{s\in[t]}$, where $\sum_{s\in[t]}\epsilon^{(s)}\leq x$, such that $\mathbf{A}^{(s)}=\operatorname{swap}(\mathbf{A}^{(s-1)},i_1^{(s)},i_2^{(s)},j_1^{(s)},j_2^{(s)},\epsilon^{(s)})$ for $s\in[t]$, where $\mathbf{A}^{(0)}=\mathbf{A}$ and $\mathbf{A}^{(t)}=\mathbf{A}'$.

The following lemma directly follows from Definition 4.3 and Definition 4.4.

Lemma 4.5. For any matrices $A, A' \in \mathbb{R}^{s \times t}$, if A' is x-swap distance from A for some $x \geq 0$, then $A' \overrightarrow{1} = A \overrightarrow{1}$ and $A' \overrightarrow{1} = A^{\top} \overrightarrow{1}$.

Recall that our objective $\mathbf{g}(\mathbf{X})$ contains two terms: (1) the linear term $\sum_{i \in [\ell], j \in [0,k]} \mathbf{C}_{ij} \mathbf{X}_{ij}$ and (2) the entropy term $\sum_{i \in [\ell]} [\mathbf{X} \overrightarrow{1}]_i \log [\mathbf{X} \overrightarrow{1}]_i - \sum_{i \in [\ell], j \in [0,k]} \mathbf{X}_{ij} \log \mathbf{X}_{ij}$. The swap operation always increases the first term, and in the following lemma we bound the loss due to the second term.

Lemma 4.6. If $A' \in \mathbb{R}^{\ell \times [0,k]}$ is x-swap distance from $A \in \mathbf{Z}_{R}^{\phi, \text{frac}}$, then, $A' \in \mathbf{Z}_{R}^{\phi, \text{frac}}$ and $\mathbf{g}(A') \ge \exp(-O(x \log n))\mathbf{g}(A)$.

One of the main contributions of our work is the following lemma, where we repeatedly apply swap operation to recover a matrix **A** which exhibits several nice properties as stated below.

Lemma 4.7. For any matrix $\mathbf{A} \in \mathbb{R}^{s \times t}$ $(s \leq t)$ that satisfies $\mathbf{A}^{\top} \overrightarrow{1} \in \mathbb{Z}^{t}_{\geq 0}$. The algorithm swapmatrix round runs in $O(s^{2}t)$ time and returns matrices \mathbf{A}' and \mathbf{B} such that,

- A' is O(s)-swap distance from $A, A'\overrightarrow{1} = A\overrightarrow{1}$ and $A'^{\top}\overrightarrow{1} = A^{\top}\overrightarrow{1}$.
- $0 \le \mathbf{B}_{ij} \le \mathbf{A}'_{ij}$ for all $i \in [s]$ and $j \in [t]$, $\mathbf{B} \overrightarrow{1} \in \mathbb{Z}^s_{\ge 0}$, $\mathbf{B}^\top \overrightarrow{1} \in \mathbb{Z}^t_{\ge 0}$ and $\|\mathbf{A}' \mathbf{B}\|_1 \le O(s)$.

The above lemma helps us modify our matrix X to a new matrix A that we can round using the create procedure. The guarantees of this procedure are summarized below.

Lemma 4.8 (Lemma 6.13 in [ACSS21]). For any $A \in \mathbf{Z}_{R}^{\phi, \text{frac}} \subseteq \mathbb{R}_{\geq 0}^{\ell \times [0, k]}$ and $\mathbf{B} \in \mathbb{R}_{\geq 0}^{\ell \times [0, k]}$ such that $\mathbf{B} \leq \mathbf{A}$, $\mathbf{B} \vdash \mathbf{C} \in \mathbb{R}_{\geq 0}^{\ell \times [0, k]}$ and $\|\mathbf{A} - \mathbf{B}\|_1 \leq t$. The algorithm create $(\mathbf{A}, \mathbf{B}, \mathbf{R})$ runs in time $O(\ell k)$ and returns a solution \mathbf{A}' and a probability discretization set \mathbf{R}' such that $|\mathbf{R}'| \leq |\mathbf{R}| + \min(k+1, t)$, $\mathbf{A}' \in \mathbf{Z}_{\mathbf{R}'}^{\phi}$ and $\mathbf{g}(\mathbf{A}') \geq \exp\left(-O\left(t\log n\right)\right)\mathbf{g}(\mathbf{A})$.

As our final goal is to return a distribution in $\Delta_{\mathbf{R}}^{\mathcal{D}}$, we also use the following discretization lemma.

Lemma 4.9. The function discretize takes as input a distribution $\mathbf{p} \in \Delta^{\mathcal{D}}$ with ℓ' distinct probability values, a profile ϕ , a discretization set of the form $\mathbf{R} = \mathbf{R}_{n,\alpha}$ for some $\alpha > 0$ and outputs a pseudo distribution $\mathbf{q} \in \Delta^{\mathcal{D}}_{\mathbf{R}}$ such that: $\mathbb{P}\left(\frac{\mathbf{q}}{\|\mathbf{q}\|_1}, \phi\right) \geq \exp(-O(\min(k, |\mathbf{R}|) + \min(k, \ell') + \alpha^2 n) \log n) \mathbb{P}\left(\mathbf{p}, \phi\right)$.

In Section 5, we use the guarantees stated above for each line of Algorithm 1 to prove Theorem 2.3. The description of the function discretize is specified in the proof of Lemma 4.9. We describe the procedure swapmatrixround and provide a proof sketch of Lemma 4.7 in Section 4.1.

4.1 Description of swapmatrixround and comparison to ACSS20

Here we describe the procedure swapmatrixround and compare our rounding algorithm to [ACSS20]. Both of [ACSS20] and our approximate PML algorithm have four main lines (I-4); we differ from [ACSS20] in the key Line [3]. This line in [ACSS20] invokes a procedure called matrixround that takes as input a matrix $\mathbf{A} \in \mathbb{R}^{\ell \times [0,k]}$ and outputs a matrix $\mathbf{B} \in \mathbb{R}^{\ell \times [0,k]}$ such that: $\mathbf{B} \leq \mathbf{A}$, $\mathbf{B} \stackrel{?}{\mathbf{I}} \in \mathbb{Z}^{\ell}_{\geq 0}$, $\mathbf{B}^{\top} \stackrel{?}{\mathbf{I}} \in \mathbb{Z}^{[0,k]}_{\geq 0}$ and $\|\mathbf{A} - \mathbf{B}\|_1 \leq O(\ell + k)$. Such a matrix \mathbf{B} is crucial as the procedure create uses \mathbf{B} to round fractional row sums of matrix \mathbf{A} to integral values. The error incurred in these two steps is at most $\exp(O(\|\mathbf{A} - \mathbf{B}\|_1 \log n)) \in \exp(O((\ell + k) \log n))$. As the procedure sparsify allows us to assume $\ell \leq k+1$, we get an $\exp(-k \log n)$ approximate PML using [ACSS20]. However, the setting that we are interested in is when $\ell \ll k$; for instance when $\ell \in O(n^{1/3})$ and $k \in \Theta(\sqrt{n})$. In these settings, we desire an $\exp(-O(\min(\ell,k)\log n)) \in \exp(-O(\ell \log n))$ approximate PML. In order to get such an improved approximation using [ACSS20], we need a matrix \mathbf{B} satisfying the earlier mentioned inequalities along with $|\mathbf{A} - \mathbf{B}|_1 \leq O(\min(k,\ell))$. However, such a matrix \mathbf{B} may not exist for arbitrary matrices \mathbf{A} and the best guarantee any algorithm can achieve is $|\mathbf{A} - \mathbf{B}|_1 \in O(\ell + k)$.

To overcome this, we introduce a new procedure called swapmatrixround that takes as input, a matrix \mathbf{A} and transforms it to a new matrix \mathbf{A}' that satisfies: $\mathbf{g}(\mathbf{A}') \geq \exp(-O(\min(k,\ell)\log n))\mathbf{g}(\mathbf{A})$. Furthermore, this transformed matrix \mathbf{A}' exhibits a matrix \mathbf{B} that satisfies the guarantees: $\mathbf{B} \leq \mathbf{A}'$, $\mathbf{B} \overrightarrow{1} \in \mathbb{Z}_{\geq 0}^{\ell}$, $\mathbf{B}^{\top} \overrightarrow{1} \in \mathbb{Z}_{\geq 0}^{k}$ and $\|\mathbf{A}' - \mathbf{B}\|_{1} \leq O(\ell)$. These matrices \mathbf{A}' and \mathbf{B} are nice in that we can invoke the procedure create, which would output a valid distribution with required guarantees. In the following we provide a description of the algorithm that finds these matrices \mathbf{A}' and \mathbf{B} .

Algorithm 2 swapmatrixround(**A**)

1: Let $\mathbf{A}^{(0)} = \mathbf{A}$ and $\mathbf{D}^{(0)} = 0$. 2: **for** $r = 1 \dots \ell$ **do** 3: $(\mathbf{Y}, j) = \text{partialRound}(\mathbf{A}^{(r-1)}, r)$ 4: $\mathbf{A}^{(r)} = \text{roundiRow}(\mathbf{Y}, j, r)$. 5: $\mathbf{D}^{(r)} = \mathbf{D}^{(r-1)} + \mathbf{Y} - \mathbf{A}^{(r)}$. 6: **end for** 7: Return $\mathbf{A}' = \mathbf{D}^{(\ell)} + \mathbf{A}^{(\ell)}$ and $\mathbf{B} = \mathbf{A}^{(\ell)}$. Our algorithm includes two main subroutines: partialRound and roundiRow. At each iteration i, the procedure partialRound considers row i and modifies it by repeatedly applying the swap operation. This modified row is nice as the procedure roundiRow can round this row to have an integral row sum while not affecting the rows in [i-1]. By iterating through all rows, we get the required matrices \mathbf{A}' and \mathbf{B} that satisfy the required guarantees. In the remainder, we formally state the guarantees achieved by the procedures partialRound and roundiRow.

Lemma 4.10. The algorithm partialRound takes as inputs $\mathbf{X} \in \mathbb{R}_{\geq 0}^{\ell \times [0,k]}$ and $i \in [\ell-1]$ that satisfies the following, $[\mathbf{X}\overrightarrow{1}]_{i'} \in \mathbb{Z}_{\geq 0}$ for all $i' \in [1,i-1]$ and $[\mathbf{X}^{\top}\overrightarrow{1}]_{j} \in \mathbb{Z}_{\geq 0}$ for all $j \in [0,k]$, and outputs a matrix $\mathbf{Y} \in \mathbb{R}_{\geq 0}^{\ell \times [0,k]}$ and an index j' such that:

- Y is within 3-swap distance from X.
- $Y_{ij'} \geq o$ and $\sum_{i'=1}^{i-1} Y_{i'j'} + Y_{ij'} o \in \mathbb{Z}_{\geq 0}$, where $o = [X\overrightarrow{1'}]_i \lfloor [X\overrightarrow{1'}]_i \rfloor$.

Furthermore, the running time of the algorithm is $O(\ell k)$.

Note that by Lemma 4.5 if **Y** is within 3-swap distance from **X**, then $\overrightarrow{\mathbf{Y}} = \overrightarrow{\mathbf{X}}$ and $\overrightarrow{\mathbf{Y}}$ $\overrightarrow{\mathbf{1}} = \overrightarrow{\mathbf{X}}$ and $\overrightarrow{\mathbf{Y}}$

Lemma 4.11. The algorithm roundiRow takes as inputs $\mathbf{Y} \in \mathbb{R}_{\geq 0}^{\ell \times [0,k]}$, an column index $j \in [0,k]$ and a row index $i \in [\ell-1]$ such that: $\mathbf{Y}^{\top}\overrightarrow{1} \in \mathbb{Z}_{\geq 0}^{[0,k]}$, $\mathbf{Y}_{ij} \geq o$ and $\sum_{i'=1}^{i-1} \mathbf{Y}_{i'j} + \mathbf{Y}_{ij} - o \in \mathbb{Z}_{\geq 0}$, where $o = [\overrightarrow{\mathbf{Y}}\overrightarrow{1}]_i - \lfloor [\overrightarrow{\mathbf{Y}}\overrightarrow{1}]_i \rfloor$. Outputs a matrix $\mathbf{X} \in \mathbb{R}_{\geq 0}^{\ell \times [0,k]}$ such that,

- $X \le Y$ and $||X Y||_1 \le 1$.
- $[X\overrightarrow{1}]_{i'} = [Y\overrightarrow{1}]_{i'}$ for all $i' \in [i-1]$, $[X\overrightarrow{1}]_i \in \mathbb{Z}_{\geq 0}$, and $X^{\top}\overrightarrow{1} \in \mathbb{Z}_{\geq 0}^{[0,k]}$.

We defer the description of all the missing procedures and proofs to appendix.

5 Proof of Main Result (Theorem 2.3)

Here we put together the results from the previous sections to prove, Theorem 2.3

Proof of Theorem 2.3 Algorithm 1 achieves the guarantees of Theorem 2.3 In the remainder of the proof, we combine the guarantees of each step of the algorithm to prove the theorem. Toward this end, we first show the following two inequalities: $\mathbf{X}_{\text{final}} \in \mathbf{Z}_{\mathbf{R}_{\text{final}}}^{\phi}$ and $\mathbf{g}(\mathbf{X}_{\text{final}}) \geq \exp(-O(\min(k,|\mathbf{R}|)\log n))\mathbf{g}(\mathbf{X})$. By Lemma 4.1 the Line 1 of Algorithm 1 returns a solution $\mathbf{X} \in \mathbf{Z}_{\mathbf{R}}^{\phi,\text{frac}}$ that satisfies,

$$C_{\phi} \cdot \mathbf{g}(\mathbf{X}) \ge \exp\left(-O\left(\min(k, |\mathbf{R}|) \log n\right)\right) \max_{\mathbf{q} \in \Delta_{\mathbf{p}}^{D}} \mathbb{P}(\mathbf{q}, \phi)$$
 (7)

By Lemma 4.2, the Line 2 of Algorithm 1 takes input \mathbf{X} and outputs \mathbf{X}' such that

$$\mathbf{X}' \in \mathbf{Z}_{\mathbf{R}}^{\phi, \text{frac}} \text{ and } \mathbf{g}(\mathbf{X}') \ge \mathbf{g}(\mathbf{X}),$$
 (8)

and $|\{i \in [\ell] \mid [\mathbf{X}'\overrightarrow{1}]_i > 0\}| \le k+1$. As the matrix \mathbf{X}' has at most k+1 non-zero rows, without loss of generality we can assume $|\mathbf{R}| \le k+1$ (by discarding zero rows).

As matrix $\mathbf{X}' \in \mathbf{Z}_{\mathbf{R}}^{\phi,\mathrm{frac}}$, we have that \mathbf{X}' has integral column sums and by invoking Lemma 4.7 with parameters $s = |\mathbf{R}|$ and t = k+1, we get matrices \mathbf{A} and \mathbf{B} that satisfy guarantees of Lemma 4.7 As $[\mathbf{A} \ \overrightarrow{1}]_i = [\mathbf{X}' \ \overrightarrow{1}]_i$ for all $i \in [\ell]$, $[\mathbf{A}^\top \ \overrightarrow{1}]_j = [\mathbf{X}'^\top \ \overrightarrow{1}]_j$ for all $j \in [0,k]$ and $\mathbf{X}' \in \mathbf{Z}_{\mathbf{R}}^{\phi,\mathrm{frac}}$, we immediately get that $\mathbf{A} \in \mathbf{Z}_{\mathbf{R}}^{\phi,\mathrm{frac}}$. Further note that \mathbf{A} is within $O(|\mathbf{R}|) = O(\min(|\mathbf{R}|,k))$ -swap distance from \mathbf{X}' and by Lemma 4.6 we get that $\mathbf{g}(\mathbf{A}) \geq \exp(-O(\min(|\mathbf{R}|,k)\log n))\mathbf{g}(\mathbf{X}')$. To summarize, we showed the following inequalities,

$$\mathbf{A} \in \mathbf{Z}_{\mathbf{R}}^{\phi, \text{frac}} \text{ and } \mathbf{g}(\mathbf{A}) \ge \exp(-O(\min(|\mathbf{R}|, k) \log n))\mathbf{g}(\mathbf{X}')$$
. (9)

Note that, Lemma 4.7 also outputs a matrix **B** that satisfies: $\mathbf{B} \leq \mathbf{A}$, $\mathbf{B} \overrightarrow{1} \in \mathbb{Z}^{\ell}$, $\mathbf{B}^{\top} \overrightarrow{1} \in \mathbb{Z}^{[0,k]}$ and $\|\mathbf{A} - \mathbf{B}\|_{1} \leq O(\min(|\mathbf{R}|,k))$. These matrices **A** and **B** satisfy the conditions of Lemma 4.8 with parameter value $t = O(\min(|\mathbf{R}|, k))$. Therefore, the procedure create takes in input matrices \mathbf{A}, \mathbf{B} and returns a solution $(\mathbf{X}_{\text{final}}, \mathbf{R}_{\text{final}})$ such that $|\mathbf{R}_{\text{final}}| \leq |\mathbf{R}| + \min(\mathbf{R}, k) \leq 2|\mathbf{R}|$ and,

$$\mathbf{X}_{\text{final}} \in \mathbf{Z}_{\mathbf{R}_{\text{final}}}^{\phi} \text{ and } \mathbf{g}(\mathbf{X}_{\text{final}}) \ge \exp(-O(\min(|\mathbf{R}|, k) \log n))\mathbf{g}(\mathbf{A})$$
. (10)

As $\mathbf{X}_{\text{final}} \in \mathbf{Z}_{\mathbf{R}_{\text{final}}}^{\phi}$, by definition Definition 3.3 and Lemma 3.2, the distribution \mathbf{p}' satisfies,

$$\mathbb{P}(\mathbf{p}', \phi) \ge \exp(-O(\min(k, |\mathbf{R}_{\text{final}}|) \log n)) C_{\phi} \mathbf{g}(\mathbf{X}_{\text{final}}) \ge \exp(-O(\min(k, |\mathbf{R}|)) \log n)) C_{\phi} \mathbf{g}(\mathbf{A})$$

$$\ge \exp(-O(\min(k, |\mathbf{R}|) \log n)) C_{\phi} \mathbf{g}(\mathbf{X}') \ge \exp(-O(\min(k, |\mathbf{R}|) \log n)) C_{\phi} \mathbf{g}(\mathbf{X})$$

$$\ge \exp(-O(\min(k, |\mathbf{R}|) \log n)) \max_{\mathbf{q} \in \Delta_{\mathbf{R}}^{\mathcal{D}}} \mathbb{P}(\mathbf{q}, \phi) .$$

In the second inequality we used Equation $\overline{(10)}$ and $|{f R}_{
m final}| \le 2|{f R}|$. In the third, fourth and fifth inequalities, we used Equation (9), Equation (8) and Equation (7) respectively.

Recall we need a distribution that approximately maximizes $\max_{\mathbf{q}\in\Delta_\mathbf{R}^\mathcal{D}}\mathbb{P}(\frac{\mathbf{q}}{\|\mathbf{q}\|_1},\phi)$ instead of just $\max_{\mathbf{q}\in\Delta_\mathbf{R}^\mathcal{D}}\mathbb{P}(\mathbf{q},\phi)$. In the remainder of the proof we provide a procedure to output such a distribution.

For any constant c>0, let $c\cdot \mathbf{R}\stackrel{\mathrm{def}}{=}\{c\cdot \mathbf{r}_i\mid \mathbf{r}_i\in \mathbf{R}\}$. For any $\mathbf{q}\in\Delta^{\mathcal{D}}_{\mathbf{R}}$, as $\|\mathbf{q}\|_1$ satisfies: $\mathbf{r}_{\min} \leq \|\mathbf{q}\|_1 \leq 1$, we get that,

$$\max_{\mathbf{q} \in \Delta_{\mathbf{R}}^{\mathcal{D}}} \mathbb{P}(\frac{\mathbf{q}}{\|\mathbf{q}\|_{1}}, \phi) = \max_{c \in [1, 1/\mathbf{r}_{\min}]_{\mathbb{R}}} \max_{\mathbf{q} \in \Delta_{c, \mathbf{R}}^{\mathcal{D}}} \mathbb{P}(\mathbf{q}, \phi) . \tag{11}$$

The above expression holds as the maximizer \mathbf{q}^* of the left hand side satisfies: $\mathbf{q}^* \in \Delta^{\mathcal{D}}_{(1/\|\mathbf{q}*\|_1) \cdot \mathbf{R}}$ Define $C \stackrel{\text{def}}{=} \{(1+\beta)^i\}_{i\in[a]}$ for some $\beta \in o(1)$, where $a \in O(\frac{1}{\beta}\log(1/\mathbf{r}_{\min}))$ is such that $\mathbf{r}_{\min}(1+\beta)^a=1$. For any constant $c\in[1,1/\mathbf{r}_{\min}]_{\mathbb{R}}$, note that there exists a constant $c'\in C$ such that $c(1-\beta)\leq c'\leq c$. Furthermore, for any distribution $\mathbf{q}\in\Delta^{\mathcal{D}}_{\mathbf{R}}$ with $\|\mathbf{q}\|_1=1/c$, note that the distribution $\mathbf{q}'=c'\mathbf{q}\in\Delta^{\mathcal{D}}_{c'\cdot\mathbf{R}}$ and satisfies: $\mathbb{P}(\frac{\mathbf{q}}{\|\mathbf{q}\|_1},\phi)=\mathbb{P}(c\cdot\mathbf{q},\phi)=\mathbb{P}(\frac{c}{c'}\mathbf{q}',\phi)=\frac{(\frac{c}{c'})^n\mathbb{P}(\mathbf{q}',\phi)}{(\frac{c}{c'})^n\mathbb{P}(\mathbf{q}',\phi)}$. Therefore we get that, $\mathbb{P}(\mathbf{q}',\phi)=\frac{(\frac{c}{c'})^n\mathbb{P}(\frac{\mathbf{q}}{\|\mathbf{q}\|_1},\phi)\geq (1-\beta)^n\mathbb{P}(\frac{\mathbf{q}}{\|\mathbf{q}\|_1},\phi)\geq (1-\beta)^n\mathbb{P}(\frac{\mathbf{q}}{\|\mathbf{q}\|_1},\phi)$ $\exp(-2\beta n)\mathbb{P}(\frac{\mathbf{q}}{\|\mathbf{q}\|_1},\phi)$. Combining this analysis with Equation (11) we get that,

$$\max_{c \in C} \max_{\mathbf{q} \in \Delta_{c \cdot \mathbf{R}}^{\mathcal{D}}} \mathbb{P}(\mathbf{q}, \phi) \ge \exp(-2\beta n) \max_{\mathbf{q} \in \Delta_{\mathcal{D}}^{\mathcal{D}}} \overline{\mathbb{P}}(\frac{\mathbf{q}}{\|\mathbf{q}\|_{1}}, \phi).$$
(12)

For each c>0 as $|\mathbf{R}|=|c\cdot\mathbf{R}|$, our algorithm (Algorithm 1) returns a distribution \mathbf{p}_c that satisfies, $\mathbb{P}\left(\mathbf{p}_c,\phi\right)\geq \exp(-O(\min(k,|\mathbf{R}|)\log n))\max_{c\in A\mathcal{D}}\mathbb{P}(\mathbf{q},\phi)$.

Let p* be the distribution that achieves the maximum objective value to our convex program among the distributions $\{\mathbf{p}_c\}_{c\in C}$. Then note that \mathbf{p}^* satisfies: $\mathbb{P}(\mathbf{p}^*,\phi) \geq \exp(-O(\min(k,|\mathbf{R}|)\log n) - 2\beta n)\max_{\mathbf{q}\in\Delta_{\mathbf{R}}^{\mathcal{D}}}\mathbb{P}(\frac{\mathbf{q}}{\|\mathbf{q}\|_1},\phi)$. Substituting $\beta=\frac{\min(k,|\mathbf{R}|)}{n}$ in the previous expression, we get, $\mathbb{P}(\mathbf{p}^*,\phi) \geq \exp(-O(\min(k,|\mathbf{R}|)\log n))\max_{\mathbf{q}\in\Delta_{\mathbf{R}}^{\mathcal{D}}}\mathbb{P}(\frac{\mathbf{q}}{\|\mathbf{q}\|_1},\phi) \ .$

$$\mathbb{P}\left(\mathbf{p}^*, \phi\right) \geq \exp(-O(\min(k, |\mathbf{R}|) \log n)) \max_{\mathbf{q} \in \Delta_{\mathbf{p}}^{\mathbf{p}}} \mathbb{P}(\frac{\mathbf{q}}{\|\mathbf{q}\|_{1}}, \phi) .$$

As each of our distributions \mathbf{p}_c (including \mathbf{p}^*) have the number of distinct probability values upper bounded by $2|\mathbf{R}|$, by applying Lemma 4.9 we get a pseudo distribution $\mathbf{q} \in \Delta^{\mathcal{D}}_{\mathbf{R}}$ with the desired guarantees. The final run time of our algorithm is $O(|C|\mathcal{T}_1) \in O(\frac{n}{\min(k,|\mathbf{R}|)} \cdot \mathcal{T}_1)$, where \mathcal{T}_1 is the time to implement Algorithm [1]. Further note that by Lemma [3.1] without loss of generality we can assume $|\mathbf{R}| \leq n/k$. As all the lines of Algorithm 1 are polynomial in n, our final running time follows from the run times of each line and we conclude the proof.

Acknowledgments and Disclosure of Funding

We would like to thank the reviewers for their valuable feedback. Researchers on this project were supported by an Amazon Research Award, a Dantzig-Lieberman Operations Research Fellowship, a Google Faculty Research Award, a Microsoft Research Faculty Fellowship, NSF CAREER Award CCF-1844855, NSF Grant CCF-1955039, a PayPal research gift, a Simons-Berkeley Research Fellowship, a Simons Investigator Award, a Sloan Research Fellowship and a Stanford Data Science Scholarship.

References

- [ACSS20] Nima Anari, Moses Charikar, Kirankumar Shiragur, and Aaron Sidford. Instance based approximations to profile maximum likelihood. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.*
- [ACSS21] Nima Anari, Moses Charikar, Kirankumar Shiragur, and Aaron Sidford. The bethe and sinkhorn permanents of low rank matrices and implications for profile maximum likelihood. In Mikhail Belkin and Samory Kpotufe, editors, *Conference on Learning Theory, COLT 2021, 15-19 August 2021, Boulder, Colorado, USA*, volume 134 of *Proceedings of Machine Learning Research*, pages 93–158. PMLR, 2021.
- [ADM+10] J. Acharya, H. Das, H. Mohimani, A. Orlitsky, and S. Pan. Exact calculation of pattern probabilities. In 2010 IEEE International Symposium on Information Theory, pages 1498–1502, June 2010.
- [ADOS17] Jayadev Acharya, Hirakendu Das, Alon Orlitsky, and Ananda Theertha Suresh. A unified maximum likelihood approach for estimating symmetric properties of discrete distributions. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 11–21. PMLR, 06–11 Aug 2017.
- [AOST14] Jayadev Acharya, Alon Orlitsky, Ananda Theertha Suresh, and Himanshu Tyagi. The complexity of estimating rényi entropy. In *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1855–1869, 2014.
- [AOST17] Jayadev Acharya, Alon Orlitsky, Ananda Theertha Suresh, and Himanshu Tyagi. Estimating renyi entropy of discrete distributions. *IEEE Trans. Inf. Theor.*, 63(1):38–56, January 2017.
 - [AW21] Josh Alman and Virginia Vassilevska Williams. A refined laser method and faster matrix multiplication. In *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms, SODA 2021, Virtual Conference, January 10 13, 2021*, pages 522–539. SIAM, 2021.
 - [BF93] John Bunge and Michael Fitzpatrick. Estimating the number of species: a review. *Journal of the American Statistical Association*, 88(421):364–373, 1993.
- [BZLV16] Y. Bu, S. Zou, Y. Liang, and V. V. Veeravalli. Estimation of kl divergence between large-alphabet distributions. In 2016 IEEE International Symposium on Information Theory (ISIT), pages 1118–1122, July 2016.
- [CCG+12] Robert K Colwell, Anne Chao, Nicholas J Gotelli, Shang-Yi Lin, Chang Xuan Mao, Robin L Chazdon, and John T Longino. Models and estimators linking individual-based and sample-based rarefaction, extrapolation and comparison of assemblages. *Journal of plant ecology*, 5(1):3–21, 2012.
 - [Cha84] A Chao. Nonparametric estimation of the number of classes in a population. scandinavianjournal of statistics11, 265-270. *Chao26511Scandinavian Journal of Statistics1984*, 1984.
- [CSS19a] Moses Charikar, Kirankumar Shiragur, and Aaron Sidford. Efficient profile maximum likelihood for universal symmetric property estimation. In *Proceedings of the 51st* Annual ACM SIGACT Symposium on Theory of Computing, STOC 2019, pages 780– 791, New York, NY, USA, 2019. ACM.
- [CSS19b] Moses Charikar, Kirankumar Shiragur, and Aaron Sidford. A general framework for symmetric property estimation. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 12447–12457. Curran Associates, Inc., 2019.

- [DS13] Timothy Daley and Andrew D Smith. Predicting the molecular complexity of sequencing libraries. *Nature methods*, 10(4):325, 2013.
- [ET76] Bradley Efron and Ronald Thisted. Estimating the number of unseen species: How many words did shakespeare know? *Biometrika*, 63(3):435–447, 1976.
- [Für05] Johannes Fürnkranz. Web mining. In *Data mining and knowledge discovery handbook*, pages 899–920. Springer, 2005.
- [Gal14] François Le Gall. Powers of tensors and fast matrix multiplication. In *International Symposium on Symbolic and Algebraic Computation, ISSAC '14, Kobe, Japan, July 23-25, 2014*, pages 296–303. ACM, 2014.
- [GTPB07] Zhan Gao, Chi-hong Tseng, Zhiheng Pei, and Martin J Blaser. Molecular analysis of human forearm superficial skin bacterial biota. *Proceedings of the National Academy of Sciences*, 104(8):2927–2932, 2007.
 - [Han21] Yanjun Han. On the high accuracy limitation of adaptive property estimation. In Arindam Banerjee and Kenji Fukumizu, editors, *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 910–918. PMLR, 13–15 Apr 2021.
- [HHRB01] Jennifer B Hughes, Jessica J Hellmann, Taylor H Ricketts, and Brendan JM Bohannan. Counting the uncountable: statistical approaches to estimating microbial diversity. *Appl. Environ. Microbiol.*, 67(10):4399–4406, 2001.
 - [HJW16] Yanjun Han, Jiantao Jiao, and Tsachy Weissman. Minimax estimation of KL divergence between discrete distributions. *ArXiv e-prints*, abs/1605.09124, 2016.
 - [HJW18] Yanjun Han, Jiantao Jiao, and Tsachy Weissman. Local moment matching: A unified methodology for symmetric functional estimation and distribution estimation under wasserstein distance. *ArXiv e-prints, arXiv:1802.08405*, 2018.
 - [HO19] Yi Hao and Alon Orlitsky. The Broad Optimality of Profile Maximum Likelihood. *ArXiv e-prints*, page arXiv:1906.03794, Jun 2019.
 - [HO20] Yi Hao and Alon Orlitsky. Profile entropy: A fundamental measure for the learnability and compressibility of distributions. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS'20, Red Hook, NY, USA, 2020. Curran Associates Inc.
 - [HS21] Yanjun Han and Kirankumar Shiragur. On the competitive analysis and high accuracy optimality of profile maximum likelihood. In Dániel Marx, editor, *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms, SODA 2021, Virtual Conference, January 10 13, 2021*, pages 1317–1336. SIAM, 2021.
 - [JHW16] J. Jiao, Y. Han, and T. Weissman. Minimax estimation of the 11 distance. In 2016 IEEE International Symposium on Information Theory (ISIT), pages 750–754, July 2016.
- [JVHW15] J. Jiao, K. Venkat, Y. Han, and T. Weissman. Minimax estimation of functionals of discrete distributions. *IEEE Transactions on Information Theory*, 61(5):2835–2885, May 2015.
 - [KLR99] Ian Kroes, Paul W Lepp, and David A Relman. Bacterial diversity within the human subgingival crevice. *Proceedings of the National Academy of Sciences*, 96(25):14547– 14552, 1999.
- [OSS⁺04] A. Orlitsky, S. Sajama, N. P. Santhanam, K. Viswanathan, and Junan Zhang. Algorithms for modeling distributions over large alphabets. In *International Symposium on Information Theory*, 2004. ISIT 2004. Proceedings., pages 304–304, 2004.
- [OSW16] Alon Orlitsky, Ananda Theertha Suresh, and Yihong Wu. Optimal prediction of the number of unseen species. *Proceedings of the National Academy of Sciences*, 113(47):13283–13288, 2016.

- [PBG⁺01] Bruce J Paster, Susan K Boches, Jamie L Galvin, Rebecca E Ericson, Carol N Lau, Valerie A Levanos, Ashish Sahasrabudhe, and Floyd E Dewhirst. Bacterial diversity in human subgingival plaque. *Journal of bacteriology*, 183(12):3770–3783, 2001.
 - [PJW17] D. S. Pavlichin, J. Jiao, and T. Weissman. Approximate Profile Maximum Likelihood. *ArXiv e-prints, arXiv:1712.07177*, December 2017.
- [RCS⁺09] Harlan S Robins, Paulo V Campregher, Santosh K Srivastava, Abigail Wacher, Cameron J Turtle, Orsalem Kahsai, Stanley R Riddell, Edus H Warren, and Christopher S Carlson. Comprehensive assessment of t-cell receptor β -chain diversity in $\alpha\beta$ t cells. Blood, 114(19):4099–4107, 2009.
 - [TE87] Ronald Thisted and Bradley Efron. Did shakespeare write a newly-discovered poem? *Biometrika*, 74(3):445–455, 1987.
 - [Von12] Pascal O. Vontobel. The bethe approximation of the pattern maximum likelihood distribution. In 2012 IEEE International Symposium on Information Theory Proceedings, pages 2012–2016, 2012.
 - [Von14] P. O. Vontobel. The bethe and sinkhorn approximations of the pattern maximum likelihood estimate and their connections to the valiant-valiant estimate. In 2014 Information Theory and Applications Workshop (ITA), pages 1–10, Feb 2014.
 - [VV11a] G. Valiant and P. Valiant. The power of linear estimators. In 2011 IEEE 52nd Annual Symposium on Foundations of Computer Science, pages 403–412, Oct 2011.
 - [VV11b] Gregory Valiant and Paul Valiant. Estimating the unseen: An n/log(n)-sample estimator for entropy and support size, shown optimal via new clts. In *Proceedings of the Forty-third Annual ACM Symposium on Theory of Computing*, STOC '11, pages 685–694, New York, NY, USA, 2011. ACM.
 - [Wil12] Virginia Vassilevska Williams. Multiplying matrices faster than coppersmith-winograd. In *Proceedings of the Forty-Fourth Annual ACM Symposium on Theory of Computing*, STOC '12, page 887–898, New York, NY, USA, 2012. Association for Computing Machinery.
 - [WY15] Y. Wu and P. Yang. Chebyshev polynomials, moment matching, and optimal estimation of the unseen. *ArXiv e-prints, arXiv:1504.01227*, April 2015.
 - [WY16] Y. Wu and P. Yang. Minimax rates of entropy estimation on large alphabets via best polynomial approximation. *IEEE Transactions on Information Theory*, 62(6):3702–3720, June 2016.
- [ZVV⁺16] James Zou, Gregory Valiant, Paul Valiant, Konrad Karczewski, Siu On Chan, Kaitlin Samocha, Monkol Lek, Shamil Sunyaev, Mark Daly, and Daniel G. MacArthur. Quantifying unobserved protein-coding variants in human populations provides a roadmap for large-scale sequencing projects. *Nature Communications*, 7:13293 EP–, Oct 2016.

Checklist

- 1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] See Section [2]
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A]
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
 - (b) Did you include complete proofs of all theoretical results? [Yes]
- 3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [N/A]
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [N/A]
- 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [N/A]
 - (b) Did you mention the license of the assets? [N/A]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
- 5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]