Near-Optimal Bounds for Testing Histogram Distributions

Anonymous Author(s) Affiliation Address email

Abstract

We investigate the problem of testing whether a discrete probability distribution 1 over an ordered domain is a histogram on a specified number of bins. One of 2 the most common tools for the succinct approximation of data, k-histograms 3 over [n], are probability distributions that are piecewise constant over a set of 4 k intervals. Given samples from an unknown distribution p on [n], we want 5 to distinguish between the cases that p is a k-histogram versus far from any k-6 histogram, in total variation distance. Our main result is a sample near-optimal and 7 computationally efficient algorithm for this testing problem, and a nearly-matching 8 (within logarithmic factors) sample complexity lower bound, showing that the 9 testing problem has sample complexity $\widetilde{\Theta}(\sqrt{nk}/\varepsilon + k/\varepsilon^2 + \sqrt{n}/\varepsilon^2)$. 10

11 **1 Introduction**

12 1.1 Background and Motivation

A classical approach for the efficient exploration of massive datasets involves the construction of 13 succinct data representations, see, e.g., the survey [CGHJ12]. One of the most useful and commonly 14 used compact representations are *histograms*. For a dataset S, whose elements are from the universe 15 $[n] := \{1, \ldots, n\}$, a k-histogram is a function that is piecewise constant over k interval pieces. 16 Histograms constitute the oldest and most popular synopsis structure in databases and have been 17 extensively studied in the database community since their introduction in the 1980s Koo80, see, 18 e.g., [GMP97] JKM⁺98, CMN98, TGIK02, GGI⁺02, GKS06, ILR12, ADH⁺15, Can16], for a 19 partial list of references. In both the statistics and computer science literatures, several methods have 20 been proposed to estimate histogram distributions in a range of natural settings [Sco79, FD81] DL04. 21 LN96, K1e09, CDSS14, ADH+15, ADLS17, DLS18. 22 In this work, we study the algorithmic task of deciding whether a (potentially very large) dataset S23

over the domain [n] is a k-histogram (i.e., it has a succinct histogram representation with k interval 24 pieces) or is "far" from any k-histogram representation (in a well-defined technical sense). Our focus 25 is on sublinear time algorithms [Rub06]. Instead of reading the entire dataset S, which could be 26 highly impractical, one can instead use randomness to sample a small subset of the dataset. Note that 27 sampling a (uniformly) random element from S is equivalent to drawing a sample from the underlying 28 probability distribution **p** of relative empirical frequencies. This observation brings our algorithmic 29 problem of "histogram testing" in the framework of distribution property testing (statistical hypothesis 30 testing) [BFR+00] BFR+13], see, e.g., [Can20] for a survey. 31

Formally, we study the following task: for an integer $1 \le k \le n$, denote by \mathcal{H}_k^n the set of k-histogram distributions over $\{1, 2, ..., n\}$, i.e., the set of all distributions **p** such that there exists a partition of [n] into k consecutive intervals (not necessarily of the same size) with **p** being uniform on each interval. Given access to i.i.d. samples from an unknown distribution **p** on [n] and a desired error

tolerance $0 < \varepsilon < 1$, we want to correctly distinguish (with high probability) between the cases that 36 **p** is a k-histogram versus ε -far from any k-histogram, in total variation distance (or, equivalently, 37 ℓ_1 -norm). It should be noted that the histogram testing problem studied here is not new. Prior work 38 within the algorithms and database theory community has investigated the complexity of the problem 39 in the past ten years (see, e.g., ILR12, ADH⁺15, Can16) and Section 1.4 for a detailed summary of 40 prior work). However, known algorithms for this task are sub-optimal, and in particular there is a 41 polynomial gap between the best known upper and lower bounds on the sample complexity of the 42 problem. At a high level, the difficulty of our histogram testing problem in the sub-linear regime lies 43 in the fact that the location and "length" of the k intervals defining the histogram representation (if 44 one exists) is a priori unknown to the algorithm. 45

We believe that the histogram testing problem is natural and interesting in its own right. Moreover, 46 a sample-efficient algorithm for this testing task can be used as a key primitive in the context of 47 model selection, where the goal is to identify the "most succinct" data representation. Indeed, various 48 algorithms are known for learning k-histograms from samples whose sample complexities (and 49 running times) scale proportionally to the succinctness parameter k (and are completely independent 50 of the domain size n) [CDSS14, ADH+15, ADLS17]. Combined with an efficient tester for the 51 property of being a k-histogram (used to identify the smallest possible value of k such that p is 52 a k-histogram, e.g., via binary search), one can obtain a sketch of the underlying dataset. See 53 Appendix C for a detailed description. 54

55 1.2 Our Results

Our main contribution is a near-characterization of the sample complexity of the histogram testing problem. Specifically, we provide (1) a sample near-optimal and computationally efficient testing algorithm for the problem, and (2) a nearly-matching sample complexity lower bound (within logarithmic factors). In particular, we establish the following theorem:

Theorem 1 (Main Result). There exists a testing algorithm for the class of k-histograms on [n] with sample complexity $m = \tilde{O}(\sqrt{nk}/\varepsilon + k/\varepsilon^2 + \sqrt{n}/\varepsilon^2)$ and running time poly(m). Moreover, for any $k \in [n]$ and $0 < \varepsilon < 1$, any testing algorithm for the class of k-histograms on [n] requires at least $\tilde{\Omega}(\sqrt{nk}/\varepsilon + k/\varepsilon^2 + \sqrt{n}/\varepsilon^2)$ samples.

64 (The $\tilde{O}(\cdot)$ and $\tilde{\Omega}(\cdot)$ notation hides polylogarithmic factors in the argument.) Theorem 1 characterizes 65 the complexity of the histogram testing problem within polylogarithmic factors. Note that there are 66 three terms in the sample complexity; namely, \sqrt{nk}/ε , k/ε^2 , and \sqrt{n}/ε^2 . The sample complexity of 67 the problem is dominated by one of these three different terms, depending on the relative sizes of 68 n, k and $1/\varepsilon$. An illustration is given in Figure 1.

⁶⁹ Prior to our work, the best previous histogram testing algorithm had sample complexity ⁷⁰ $\widetilde{O}(\sqrt{kn}/\varepsilon^3)$ [CDGR18], while the best known lower bound was $\widetilde{\Omega}(\sqrt{n}/\varepsilon^2 + k/\varepsilon)$ [Can16].

We note that previous upper and lower bounds exhibit a polynomial gap, even for constant values 71 of ε or k. For example, in the "large-k" regime where $k = n^c$ for some constant 0 < c < 1, there 72 was a gap between $O(n^{1/2+c/2})$ and $\Omega(n^{1/2})$ in the sample complexity. In this regime, however, 73 Theorem 1 results in the near-optimal bound of $\tilde{\Theta}(n^{1/2+c/2})$. Similarly, in the "high-accuracy" regime where $\varepsilon = 1/n^c$ for some constant c > 0 (and, say, constant k), previous bounds only established that the sample complexity lies between $\tilde{O}(n^{1/2+3c})$ and $\tilde{\Omega}(n^{1/2+2c})$, while our result 74 75 76 shows the (nearly-tight) bound is $\widetilde{\Theta}(n^{1/2+c})$. These are only two specific examples: more generally, 77 the previously known bounds are suboptimal by polynomial factors in $1/\varepsilon$ when $\varepsilon \ge \sqrt{k/n}$; and 78 by polynomial factors in all parameters $k, n, 1/\varepsilon$ when $\varepsilon \leq \sqrt{k/n}$. Theorem 1 settles the sample 79 complexity of the problem, up to logarithmic factors, for *every* parameter setting. 80 At a technical level, our sample complexity lower bound construction conceptually differs from previ-81 ous work in distribution testing, drawing instead from sophisticated techniques from the distribution 82

estimation literature. Our upper bound follows from the "Testing-via-Learning" framework proposed

in [ADK15]. The main technical innovation is a sample- and time- efficient *adaptive* algorithm which

¹As discussed in Section 1.4 while an upper bound of $\tilde{\Omega}(\sqrt{n}/\varepsilon^2 + k/\varepsilon^3)$ is claimed in [Can16], the analysis of their algorithm is flawed; and, indeed, our work shows that the sample complexity bound stated in [Can16] *cannot* hold, as it would contradict our lower bound.



Figure 1: The x-axis, y-axis are $\log(k)/\log(n)$ and $\log(1/\varepsilon)/\log(n)$ respectively. Each point in the graph corresponds to a setting of n, k, ε , and is colored based on the corresponding dominating term.

can simultaneously learn an unknown histogram distribution *with unknown interval structure* and identify a domain where the learned result is accurate. We elaborate on these aspects next.

87 1.3 Overview of Techniques

Sample Complexity Lower Bound. We follow the typical high-level approach in proving sample complexity lower bound. Namely, we define two ensembles of distributions D_{YES} and D_{NO} such that, with high probability, the following conditions are satisfied: (1) a random distribution from D_{YES} is a *k*-histogram, (2) a random distribution from D_{NO} is ε -far from any *k*-histogram, and (3) given samples of appropriate size, it is information-theoretically impossible to distinguish a random distribution drawn from D_{YES} from a random distribution drawn from D_{NO} .

We start by describing our hard instances for the case that the accuracy parameter ε is a small universal constant. On the one hand we define D_{YES} so that all \mathbf{p}_i 's are the same except for a "small" number of domain elements i.e., $c \cdot k$ for a small constant $c \in (0, 1)$. On the other hand, for a distribution \mathbf{p} drawn from D_{NO} , \mathbf{p}_i will be randomly 0 or roughly 2/n, except for at most a constant fraction of the elements. It is not hard to see that, with high probability, a distribution drawn from D_{YES} (resp. D_{NO}) will be a k-histogram (resp. far from being a k-histogram).

To ensure that the underlying distributions are indistinguishable using a small sample size, we want to 100 guarantee that, for all small values of t, the number of elements with exactly t samples will be roughly 101 the same for $D_{\rm YES}$ and $D_{\rm NO}$, as this rules out any test statistic relying on counting the number of 102 t-way collisions among the samples. Following [Val11, VV13] IVYHW15, WY16] this is essentially 103 equivalent to showing that distributions drawn from D_{YES} and D_{NO} match their low-degree moments. 104 In particular, for a random pair of distributions \mathbf{p}, \mathbf{p}' drawn from D_{YES} and D_{NO} respectively, we want 105 that $\sum_i \mathbf{p}_i^t$ and $\sum_i \mathbf{p}_i'^t$ are roughly the same for all small values t. We note that the non-exceptional elements of a distribution \mathbf{p}' drawn from D_{NO} — which have probability mass either 0 or roughly 106 107 2/n — will have second moment larger than the non-exceptional elements of a distribution p drawn 108 from D_{YES} — which have probability mass roughly 1/n — by approximately 1/n. To counteract this 109 discrepancy, the (fewer than k) exceptional elements in D_{YES} must have average mass at least $1/\sqrt{kn}$. 110 Fortunately, using techniques from $[VV13] WY^+19]$, we are able to construct distributions that match 111 $t = \Theta(\log n)$ moments in which no individual bin has mass more than $O(1/\sqrt{kn})$. Combining this 112 construction with basic information-theoretic arguments gives us an $\Omega(\sqrt{kn})$ sample complexity 113 lower bound. We note that this lower bound is tight in the sense that with more than $\Omega(\sqrt{kn})$ samples 114 one can reliably identify the exceptional elements, as they will each have relatively large numbers 115 of samples with high probability, allowing us to distinguish D_{YES} from D_{NO} simply based on the 116 sub-distributions over these elements. 117

Given the aforementioned construction (for constant ε), it is easy to obtain a sample lower bound of $\tilde{\Omega}(\sqrt{kn}/\varepsilon)$ by mixing our hard instances with the uniform distribution (with mixing weights ε and $1 - \varepsilon$ respectively). In fact, even if the testing algorithm knows in advance which samples come from the uniform part and which samples come from the original hard instance, distinguishing would require $\tilde{\Omega}(\sqrt{kn})$ samples from the original hard instance, and therefore $\tilde{\Omega}(\sqrt{kn}/\varepsilon)$ samples overall.

This sample size lower bound turns out to be tight for ε relatively large, as one can still reliably 123 identify the exceptional bins with only $\Omega(\sqrt{kn/\varepsilon})$ samples. However, when ε becomes sufficiently 124 small, identifying the exceptional bins becomes more challenging. Indeed, if we take m samples, we 125 expect that an exceptional bin has roughly $m\varepsilon/\sqrt{kn}$ more samples than a non-exceptional bin. On 126 the other hand, a non-exceptional bin will have roughly Poi(m/n) samples with standard deviation 127 $\sqrt{m/n}$. When $m/n \gg m\varepsilon/\sqrt{kn}$ (which happens in the regime $\varepsilon \ll \sqrt{k/n}$), in order for the 128 exceptional bins to be distinguishable, we would need that $m\varepsilon/\sqrt{kn} \gg \sqrt{m/n}$ or $m \gg k/\varepsilon^2$ many samples. Using a careful information-theoretic argument, we formalize this intuition to show that 129 130 $\widetilde{\Omega}(k/\varepsilon^2)$ is indeed a sample lower bound in this regime. 131

Sample-Efficient Tester. The starting point of our efficient tester is the Testing-via-Learning ap-132 proach of [ADK15]. Very roughly speaking, we first design a learning procedure which outputs a 133 distribution $\hat{\mathbf{p}}$ that would be close to \mathbf{p} in χ^2 divergence, assuming that \mathbf{p} was in fact a k-histogram. 134 Then we use a χ^2/ℓ_1 tolerant tester, in the spirit of the one introduced in [ADK15], to distinguish 135 between the cases that p is close to \hat{p} in χ^2 divergence versus far from \hat{p} in ℓ_1 -distance. This step is 136 however significantly harder than this simple outline suggests, as it turns out challenging to perform 137 the first step exactly. Instead, we design a specific learning algorithm with an implicit "hybrid" 138 learning guarantee, (see Lemma 5) which in turns requires us to considerably generalize and adapt 139 the "tolerant testing part" to avoid spurious discrepancies (introduced in the imperfect learning stage) 140 which may lead to false negatives. 141

To implement the first step, we follow the general "learn-and-sieve" idea suggested in [Can16], with 142 important modifications to address the flaw in their approach and its analysis. In particular, suppose 143 that **p** is a k-histogram. Then, if we knew the corresponding k intervals (that make up the partition 144 145 for the k-histogram), it suffices to learn the mass of p on each interval, and let \hat{p} be uniform on each interval (with the appropriate total mass). Of course, a key source of difficulty arises from the fact 146 that we do not know the partition a priori. To circumvent this issue, we divide [n] into (roughly) 147 $K = \Theta(k)$ intervals and try to detect if **p** is far in χ^2 divergence from being uniform on any of these 148 intervals. If an interval from our partition incurs large χ^2 error (we call such an interval *bad*), we 149 know that p must not be constant within this interval. Therefore, we proceed to subdivide these bad 150 intervals into roughly equal parts, and recurse on the $\Theta(k)$ intervals in our new partition. Assuming 151 **p** is a k-histogram, we subdivide at most k intervals in each iteration, since there could be at most 152 k intervals from any interval partition of [n] where p is not constant. Hence, in each iteration, we 153 decrease the mass of the bad intervals by at least a constant factor. We repeat the process for at 154 most $O(\log(1/\varepsilon))$ many iterations; after this many iterations, the total mass of the bad intervals will 155 become $O(\varepsilon)$, and thus they may be safely ignored. 156

A significant difference between our method and the approach from [Can16] lies in the method 157 of sieving. In [Can16], it was only said that the algorithm would filter out a subset of breakpoint 158 intervals based on the χ^2 statistics ([ADK15]) with the goal of reducing discrepancy; this is where 159 the main gap in their analysis lies, and the particular (flawed) approach they suggested does not seem 160 161 to be fixable [Can22]. On the contrary, we characterize the exact set of intervals that need to (and can) be removed with the formal definition of *bad* intervals with respect to a given partition \mathcal{I} of [n]162 (See Definition 2). Based on that, our approach is to search for any sub-intervals J (not necessarily 163 an interval in \mathcal{I} on which the χ^2 divergence between **p** and $\hat{\mathbf{p}}$, an approximation of **p** assuming **p** is 164 uniform over intervals within the given partition, is more than $\Omega(\varepsilon^2/k)$. For an interval I from the 165 partition \mathcal{I} , we show the inclusion of such "bad sub-interval" $J \subseteq I$ then certifies the "badness" of I 166 itself. To find such a J, we need a technique for accurately approximating p(J) simultaneously for 167 all intervals $J \subseteq [n]$, in both absolute and relative error; a notion of approximation much stronger 168 that what classical tools from statistical learning theory such as the VC inequality or the Dvoretzky-169 Kiefer–Wolfowitz (DKW) inequality provide. Notice that, for a *fixed* interval $J \subseteq [n]$, taking the 170 empirical distribution over b samples gives an estimate **q** of **p** such that $|\mathbf{q}(J) - \mathbf{p}(J)| < \sqrt{\mathbf{p}(J)/b}$ 171 with constant probability. By taking $\Theta(\log(n))$ batches of samples (each containing b i.i.d. samples 172 from p), and computing the median value of all of the q(J)'s, with high probability for each J, we 173 then obtain an estimate $\hat{\varphi}(J)^2$ for which the above condition holds. Using the sub-routine, as long 174 as b is at least $\Omega(k/\varepsilon^2)$, we can ensure that $|\hat{\varphi}(J) - \mathbf{p}(J)|^2/\mathbf{p}(J) \ll \varepsilon^2/k$, and we can then safely 175 use our estimate $\hat{\varphi}(J)$ as a proxy for $\mathbf{p}(J)$ for the detection of those "bad sub-intervals" for which 176

²Notice that $\hat{\varphi}$ is neither a distribution nor a measure, but just a map from intervals to positive real values.

177 $|\mathbf{p}(J) - \hat{\mathbf{p}}(J)|^2 / \hat{\mathbf{p}}(J)$ is large, which in turns certify the bad intervals from a given partition. This 178 suffices *unless* $\mathbf{p}(J)$ is substantially larger than our estimate $\hat{\mathbf{p}}(J)$.

Unfortunately, the ratio between $|\hat{\varphi}(J) - \mathbf{p}(J)|^2 / \hat{\mathbf{p}}(J)$ and $|\hat{\varphi}(J) - \mathbf{p}(J)|^2 / \mathbf{p}(J)$ (in particular 179 $\mathbf{p}(J)/\hat{\mathbf{p}}(J)$ can be unbounded when $\mathbf{p}(J)$ is smaller than 1/b. In such a case, in a collection of 180 b samples from p, we are likely to see no samples in J, and thus our empirical estimate $\hat{\mathbf{p}}(J)$ will 181 be 0. We can fix this issue (i.e., the case where $\hat{\mathbf{p}}(J)$ is actually 0) by mixing both \mathbf{p} and $\hat{\mathbf{p}}$ with 182 the uniform distribution, thus allowing us to assume that $\hat{\mathbf{p}}(J) > |J|/2n > 1/(2n)$. Yet, this still 183 leaves a potential gap of roughly n/b between the ratio of $\mathbf{p}(J)$ and $\hat{\mathbf{p}}(J)$. Fortunately, if we select 184 $b \gg \sqrt{nk}/\varepsilon$, we will have that $|\hat{\varphi}(J) - \mathbf{p}(J)|^2/\mathbf{p}(J) \ll \varepsilon/\sqrt{nk}$, and even accounting for losing a factor of b/n, we will still have that $|\hat{\varphi}(J) - \mathbf{p}(J)|^2/\mathbf{p}(J) \ll \varepsilon^2/k$. This implies that we will 185 186 successfully detect any bad intervals and achieve our learning guarantees. 187

188 1.4 Prior Work

Motivated by the question of building provably good succinct representations of a dataset from 189 only a small sub-sample of the data, [ILR12] first introduced histogram testing as a preliminary, 190 ultra-efficient decision subroutine to find the best parameter k for the number of bins. They provided 191 an algorithm for this task which required $\hat{O}(\sqrt{kn}/\varepsilon^5)$ samples from the dataset, a sample complexity 192 which beats the naïve approach (reading and processing the whole dataset) for small values of k193 and relatively large values of the accuracy parameter ε . Subsequent work [CDGR18] reduced the 194 dependence on ε from quintic to cubic, giving an algorithm with sample complexity $O(\sqrt{kn/\varepsilon^3})$. 195 This bound was, however, still quite far from the "trivial" lower bound of $\Omega(\sqrt{n}/\varepsilon^2)$, which follows 196 from a reduction to uniformity testing (i.e., the case k = 1) Pan08. 197

Prior to the current work, an $\widetilde{O}(\sqrt{n}/\varepsilon^2 + k/\varepsilon^3)$ upper bound and an $\widetilde{\Omega}(\sqrt{n}/\varepsilon^2 + k/\varepsilon)$ lower bound 198 were obtained in Can16. While the lower bound is theoretically sound (albeit, as it turns out, 199 suboptimal), as pointed out in [Can22], the upper bound does not hold due to a technical flaw in the 200 analysis, leaving the optimal sample complexity of the problem open for even constant ε . Moreover, 201 the lower bound of Can16, based on a reduction of histogram testing to the well-studied problem of 202 support size estimation, provably cannot be improved to provide either (i) a quadratic dependence on 203 ε , i.e., $\hat{\Omega}(k/\varepsilon^2)$ or (ii) coupling between the two domain parameters k, n, i.e., $\hat{\Omega}(\sqrt{nk/\varepsilon})$. Our work 204 remedies all those issues, fully resolving the question of histogram testing, for the whole parameter 205 range, within logarithmic factors. 206

Finally, we note that a number of works have obtained algorithms and lower bounds for related, yet significantly different, testing problems. Specifically, [DK16] gave a sample-optimal testing algorithm for the special case of our problem where the *k* intervals are known *a priori*. Moreover, a number of works [DKN15b], DKN15a, DKN17 have obtained identity and equivalence testers *under the assumption* that the input distributions are *k*-histograms.

Preliminaries.We denote by $\operatorname{TV}(\mathbf{p}, \mathbf{q})$ the total variation (TV) distance between probability distributions \mathbf{p}, \mathbf{q} over $[n] \coloneqq \{1, 2, \ldots, n\}$, defined as $\operatorname{TV}(\mathbf{p}, \mathbf{q}) \coloneqq \sup_{S \subseteq [n]}(\mathbf{p}(S) - \mathbf{q}(S)) = \frac{1}{2} \sum_{i=1}^{n} |\mathbf{p}(i) - \mathbf{q}(i)|$, where $\mathbf{p}(S) \coloneqq \sum_{i \in S} \mathbf{p}(i)$. We will make essential use of the χ^2 -divergence of \mathbf{p} with respect to \mathbf{q} , defined as $d_{\chi^2}(\mathbf{p} \| \mathbf{q}) \coloneqq \sum_{i=1}^{n} (\mathbf{p}_i - \mathbf{q}_i)^2 / \mathbf{q}_i$. We will also require generalizations of these definitions on restrictions of the domain. In particular, given $S \subseteq [n]$, we use the notation $\operatorname{TV}^S(\mathbf{p}, \mathbf{q}) \coloneqq (1/2) \sum_{i \in S} |\mathbf{p}(i) - \mathbf{q}(i)|$ and $d_{\chi^2}^S(\mathbf{p} \| \mathbf{q}) \coloneqq \sum_{i \in S} (\mathbf{p}_i - \mathbf{q}_i)^2 / \mathbf{q}_i$. We note that for any $S \subseteq [n]$, it holds that $\operatorname{TV}^S(\mathbf{p}, \mathbf{q})^2 \leq \frac{1}{4} d_{\chi^2}^S(\mathbf{p} \| \mathbf{q})$.

The asymptotic notation \tilde{O} (resp. $\tilde{\Omega}$) suppresses logarithmic factors in its argument, i.e., $\tilde{O}(f(n)) = O(f(n) \log^c f(n))$ and $\tilde{\Omega}(f(n)) = \Omega(f(n)/\log^c f(n))$, where c > 0 is a universal constant. The notations \ll and \gg intuitively mean "much less than" and "much greater than" respectively. Formally, we write $f(n) \ll g(n)$ to denote that $f(n) < c \cdot g(n)$, for some universal constant 0 < c.

223 2 Near-Optimal Tester

A preliminary simplification. Without loss of generality, we will assume that $\mathbf{p}(i) \ge \frac{1}{2n}$ for every *i* $\in [n]$. Indeed, this can be ensured by mixing the unknown distribution with the uniform distribution u_n on [n] beforehand, i.e., $\mathbf{p}' \coloneqq \frac{1}{2}(\mathbf{p} + \mathbf{u}_n)$ (see Fact 3 in Appendix for how to sample from \mathbf{p}' efficiently). It is easy to see that \mathbf{p}' remains a histogram after mixing: $\mathbf{p}' \in \mathcal{H}_k^n$ if $\mathbf{p} \in \mathcal{H}_k^n$, and \mathbf{p}' is at least $(\varepsilon/2)$ -far away from every histogram if \mathbf{p} is ε -far from every histogram.

Testing via Learning. The main approach is to follow the Testing-via-Learning framework pro-229 posed in [ADK15]. In particular, suppose we have a learning algorithm capable of constructing $\hat{\mathbf{p}}$ 230 that is close to \mathbf{p} in χ^2 divergence when $\mathbf{p} \in \mathcal{H}_n^k$. Then, (1) if $\mathbf{p} \in \mathcal{H}_k^n$, we will have that \mathbf{p} and $\hat{\mathbf{p}}$ are close *and* (as a consequence of this) that $\hat{\mathbf{p}}$ is close to being a k-histogram. Yet, (2) if \mathbf{p} is far from 231 232 being a k-histogram, then by the triangle inequality we must have either that $\hat{\mathbf{p}}$ is far from being a 233 k-histogram, or that p and $\hat{\mathbf{p}}$ are far from each other in ℓ_1 distance. We can use dynamic programming 234 to check the explicit description is indeed close to a k-histogram in ℓ_1 distance efficiently (see Lemma 235 4.11 of [CDGR18]). To verify p and \hat{p} are close, we will use a result of [ADK15] on tolerant identity 236 testing. In particular, given an explicit description $\hat{\mathbf{p}}$, the tester takes sample from the unknown 237 distribution **p** and decides whether **p** and $\hat{\mathbf{p}}$ are closed in χ^2 divergence or far in ℓ_1 distance. We 238 remark that $\hat{\mathbf{p}}$ can be relaxed to be a positive measure. 239

Lemma 1 (Adapted from Lemmas 2 and 3 [ADK15]). Let \mathbf{p} and $\hat{\mathbf{p}}$ be a distribution and a positive measure defined on [n] respectively. Fix $\varepsilon \in (0, 1)$ and let $\mathcal{A} = \{i \in [n] : \hat{\mathbf{p}}(i) \ge \varepsilon/(50n)\}$. There exists a tester Tolerance-Identity-Test, which takes $\operatorname{Poi}(m)$ i.i.d. samples from \mathbf{p} and outputs Accept if $d_{\chi^2}^{\mathcal{A}}(\mathbf{p} \| \hat{\mathbf{p}}) \le \varepsilon^2/500$ and Reject if $TV^{\mathcal{A}}(\mathbf{p}, \hat{\mathbf{p}}) \ge \varepsilon$ with constant probability.

Outline for Learning. If $\mathbf{p} \in \mathcal{H}_k^n$ and we know the partition of \mathbf{p} in advance, one can learn \mathbf{p} up 244 to ε^2 in χ^2 divergence with $\Theta(k/\varepsilon^2)$ samples (following the analysis of Laplace estimator from 245 KOPS15). Without the partition information, we can nonetheless achieve a weaker guarantee. That 246 is, we can output a fully specified measure $\hat{\mathbf{p}}$ on [n], together with a sub-domain $\mathcal{G} \subseteq [n]$, such that 247 $d_{\gamma^2}^{\mathcal{G}}(\mathbf{p}\|\hat{\mathbf{p}})$ is small. In particular, we can achieve the guarantee in three steps. (i) Equally divide the 248 domain [n] into $K \gg k$ many intervals (Lemma 2). (ii) Output a succinct measure $\hat{\mathbf{p}}$ that is constant on each interval specified by Step (i) (Section 2.1). (iii) Identify the intervals \mathcal{B} where $d_{\chi^2}^{\mathcal{B}}(\mathbf{p} \| \hat{\mathbf{p}})$ is large (Section 2.2) and take $\mathcal{G} = [n] \setminus \mathcal{B}$. The fact that we only have \mathbf{p} and $\hat{\mathbf{p}}$ close in χ^2 divergence 249 250 251 on a sub-domain \mathcal{G} is a reasonable compromise as long as $\mathbf{p}(\mathcal{B}), \hat{\mathbf{p}}(\mathcal{B}) \ll \varepsilon$: if \mathbf{p} is ε -far away from 252 $\hat{\mathbf{p}}$ in ℓ_1 distance on [n], \mathbf{p} is at least $(\varepsilon - \mathbf{p}(\mathcal{B}) - \hat{\mathbf{p}}(\mathcal{B}))$ -far away from $\hat{\mathbf{p}}$ on $[n] \setminus \mathcal{B}$. Otherwise, we 253 may take more samples from p restricted to \mathcal{B} and sub-divide the problematic intervals identified in 254 Step (iii). Repeating the above steps iteratively then brings us to the case $p(\mathcal{B}) \ll \varepsilon$. 255

Equitable Partition. The first step is to divide the domain into $\Theta(k)$ many intervals over which the masses of p are approximately equal. As shown in [ADK15], this can be done with $\widetilde{\Theta}(k)$ many samples through a routine we denote as **Approx-Divide**. We also need a routine for sub-dividing a set of disjoint intervals into even lighter sub-intervals. Nonetheless, one can reduce the sub-dividing task to domain partitioning by running **Approx-Divide** on the sub-distribution restricted to the set of disjoint intervals. Proofs are provided in Appendix [A.1].

Lemma 2. There exists an algorithm Approx-Sub-Divide that, given parameters $\delta \in (0, 1]$ and integer B > 1, as well as a set of disjoint intervals $\mathcal{I} = \{I_1, I_2, \dots, I_q\}$, given sample access to \mathbf{p} on [n], outputs a list of partitions S_1, \dots, S_q , where S_i is the partition of the interval $I_i \in \mathcal{I}$, such that the following holds with probability at least $1 - \delta$. (i) The algorithm uses $O\left(B/\mathbf{p}(\mathcal{I}) \cdot \log(B/\delta)\right)$ samples. (ii) The output contains at most (8B+q) intervals in total. (iii) Every non-singleton interval $S \in \bigcup_{i=1}^q S_j$ satisfies $\mathbf{p}(S) \leq \mathbf{p}(\mathcal{I}) \cdot 16/B$.

268 2.1 Simultaneously Estimating Mass of Intervals

In this section, we first introduce **Interval-Mass-Estimate**, a sub-routine that can accurately approximate the mass of $\mathbf{p}(J)$ for all intervals $J \subseteq [n]$ simultaneously and then show how we can use it to learn \mathbf{p} (assuming $\mathbf{p} \in \mathcal{H}_n^k$).

Interval-Mass-Estimate first divides the number of samples drawn into $\Theta(\log(n/\delta))$ batches. For an interval *I*, we compute the estimate (number of samples falling in *I* divided by the batch size) for each batch separately and compute the median over the statistics. This is often referred as the "Median Trick" and is crucial in achieving the learning guarantees with high probability. Pseudo-code and analysis are provided in Appendix [A.2]

Lemma 3. Let be **p** be supported on [n] such that $\mathbf{p}(i) \ge 1/(2n)$. Fix $b \in \mathbb{Z}^+$ and $\delta \in (0, 1]$. The algorithm **Interval-Mass-Estimate** takes $3b \log(n/\delta)$ i.i.d. samples from **p** and outputs $\hat{\varphi}$, a map

from sub-intervals of [n] to real values, such that, with probability at least $1 - \delta$, for every sub-interval $I \subseteq [n]$ it holds that $\mathbf{p}(I)/\hat{\varphi}(I) \leq \max(2, 8n/b), \hat{\varphi}(I)/\mathbf{p}(I) \leq 3$ and $|\hat{\varphi}(I) - \mathbf{p}(I)| \leq \sqrt{\mathbf{p}(I)/b}$.

Let \mathcal{I} be a partition of [n]. We try to learn \mathbf{p} pretending that \mathbf{p} is constant over each interval within \mathcal{I} with the routine **Empirical-Learning**. In particular, the algorithm uses **Interval-Mass-Estimate** to obtain estimations of the mass of $I \in \mathcal{I}$ and then flatten the mass uniformly among elements $i \in I$. Notice that, due to the application of the median trick, the output is not necessarily a distribution but rather a positive measure³ $\hat{\mathbf{p}}$ on [n] which is constant over each interval within \mathcal{I} .

²⁸⁶ If **p** is indeed a k-histogram, errors are only incurred on a special type of intervals (of which there ²⁸⁷ are at most k) which we refer to as the *breakpoint intervals*.

Definition 1 (Breakpoint Intervals). *Given a k-histogram* \mathbf{p} *on* [n], we say that $i \in [n]$ is a breakpoint with respect to \mathbf{p} if $\mathbf{p}(i) \neq \mathbf{p}(i+1)$; and that an interval $I \subseteq [n]$ is a breakpoint interval (with respect to \mathbf{p}) if I contains at least one breakpoint.

With Definition i in mind, we now specify the formal learning guarantees. Pseudo-code and proofs are provided in Appendix A.3

Lemma 4. Suppose $\mathbf{p} \in \mathcal{H}_k^n$. Let \mathcal{I} be a partition of [n] into K intervals. Let $b \in \mathbb{Z}^+$, $\delta \in (0, 1]$ and $T := 3 \log(K/\delta)$. There exists an algorithm **Empirical-Learning**, given (Tb) i.i.d. samples from \mathbf{p} , outputs a positive measure $\hat{\mathbf{p}}$, which satisfies the following with probability at least $1 - \delta$. (i) $\hat{\mathbf{p}}$ is constant within each interval in \mathcal{I} . (ii) For every sub-intervals $J \subseteq I$ where $I \in \mathcal{I}$ is a non-breakpoint interval with respect to \mathbf{p} , we have $\mathbf{p}(J)/\hat{\mathbf{p}}(J) \leq \max(2, 8 \cdot n/b)$ and $|\hat{\mathbf{p}}(J) - \mathbf{p}(J)| \leq \sqrt{\mathbf{p}(J)/b}$.

By combining the two guarantees in (ii) in Lemma 4, one can see the χ^2 divergence between p 298 and $\hat{\mathbf{p}}$, restricted to the non-breakpoint intervals, will be at most ε^2 with high probability if taking 299 $\Theta(KT/\varepsilon^2)$ many samples. However, following a result from [KOPS15, Can16], one only need $\Theta(K/\varepsilon^2)$ samples to learn a K-histogram up to ε^2 error in this restricted notion of χ^2 divergence. 300 301 One may wonder whether this is enough for us, and if the stronger (but less natural) guarantees 302 provided by Lemma 4, which end up increasing the number of samples required, are necessary. As 303 we will see in the next section, we indeed need not only that the χ^2 divergence is small, but also that 304 the ratio $\mathbf{p}(I)/\hat{\mathbf{p}}(I)$ is bounded for all non-breakpoint intervals. In particular, this latter property 305 enables us to compute relatively accurate estimates of the χ^2 divergence restricted to sub-intervals 306 and (consequently) to tell whether p is constant or from far from being constant on an interval. 307

308 2.2 Bad Interval Detection

While large contributions to the χ^2 divergence (assuming the learning phase was successful) will only come from breakpoint intervals, not all of them will necessarily contribute significantly to the χ^2 divergence. In particular, a breakpoint interval is only considered "bad" and needs to be filtered out if the error incurred is proportional to the number of breakpoints within.

Definition 2 (ε -Bad-Interval). *Fix a partition* \mathcal{I} *of* [n] *containing* K *intervals. Let* $I \in \mathcal{I}$ *be a breakpoint interval of* \mathbf{p} . *Furthermore, suppose* I *contains* j - 1 *breakpoints i.e.* \mathbf{p} *is* j-*piece-wise uniform in* I. We say $I \in \mathcal{I}$ *is an* ε -bad interval with respect to $\hat{\mathbf{p}}$ and \mathcal{I} if $d_{\chi^2}^I(\mathbf{p} \| \hat{\mathbf{p}}) \geq j \cdot \varepsilon^2/K$.

The definition suits our purpose for two reasons. (i) The total χ^2 error between \mathbf{p} and $\hat{\mathbf{p}}$ on the set of "good" intervals (complement of the set of "bad" intervals) is small. Indeed, let $\mathcal{G} \in \mathcal{I}$ be a set containing no ε -bad intervals. Since there are at most K intervals contained in \mathcal{G} and k breakpoints contained in the intervals in \mathcal{G} , it is easy to see that $d_{\chi^2}^{\mathcal{G}}(\mathbf{p} \| \hat{\mathbf{p}}) \leq O(\varepsilon^2)$. (ii) One can reliably separate bad intervals from non-breakpoint intervals assuming the learning phase was successful. To see why, note that in that case every non-breakpoint interval I satisfies $d_{\chi^2}^{J}(\mathbf{p} \| \hat{\mathbf{p}}) \ll \varepsilon^2/K$ for all $J \subseteq I$ with high probability. On the contrary, for any bad interval I, we claim there must be a sub-interval $Q \subseteq I$ where $d_{\chi^2}^{Q}(\mathbf{p} \| \hat{\mathbf{p}}) \ge \varepsilon^2/K$ and both \mathbf{p} and $\hat{\mathbf{p}}$ are constant within. In particular, if I is an ε -bad interval that contains (j-1) breakpoints, we then have a partition $\{Q_1, \dots, Q_j\}$ of Iover which \mathbf{p} is piece-wise constant and at least one of them will have χ^2 error at least ε^2/K .

Our next step is to show how we can leverage the separating condition to design an efficient bad interval detection mechanism. This is where our method *significantly differs* from [Can16]. At a high level, we take another set of independent samples to get an estimate $\hat{\varphi}(Q)$ of $\mathbf{p}(Q)$ for all $Q \subseteq [n]$

³That is, $\hat{\mathbf{p}}$ might not sum to one, and thus is not itself a probability distribution.

- 329
- 330
- simultaneously. Then, we compare $\hat{\varphi}(Q)$ with $\hat{\mathbf{p}}(Q)$ to see whether we have $d_{\chi^2}^Q(\mathbf{p} \| \hat{\mathbf{p}}) \geq \varepsilon^2 / K$, which would in turns imply the interval $I \supseteq Q$ from the given partition is ε -bad. We next provide the pseudo-code for **Learn-And-Sieve**, which finds a positive measure $\hat{\mathbf{p}}$ on [n] and a domain \mathcal{B} such 331 that $d_{v^2}^{[n]\setminus\mathcal{B}}\left(\mathbf{p}\|\hat{\mathbf{p}}\right) \leq O(\varepsilon^2)$ provided $\mathbf{p} \in \mathcal{H}_k^n$. Its detailed analysis can be found in Appendix A.4

Algorithm 1 Learn-And-Sieve

- **Require:** Sample access to p; a partition \mathcal{I} of [n] containing K intervals; accuracy ε ; failure probability δ .
- 1: Let $m = C \cdot (K/\varepsilon^2 + \sqrt{Kn}/\varepsilon) \cdot \log(n/\delta)$ for a sufficiently large constant C.
- 2: Draw 2m i.i.d. samples from p and split the samples evenly into S_1, S_2 .
- 3: $\hat{\mathbf{p}} \leftarrow \text{Empirical-Learning}(\mathcal{S}_1, \mathcal{I}, \delta/4), \hat{\varphi} \leftarrow \text{Interval-Mass-Estimate}(\mathcal{S}_2, \delta/4), \mathcal{B} \leftarrow \{\}.$
- 4: for all intervals $Q \subseteq I$ for some $I \in \mathcal{I}$ do
- if $\hat{\varphi}(Q)/\hat{\mathbf{p}}(Q) > 6 \cdot \max(1, \varepsilon \cdot \sqrt{n/K})$ or $|\hat{\varphi}(Q) \hat{\mathbf{p}}(Q)| > 0.5\sqrt{\hat{\mathbf{p}}(Q)\varepsilon^2/K}$ then 5:
- Add I to \mathcal{B} . 6:
- 7: Output Reject if \mathcal{B} contains more than k intervals. Otherwise, return \mathcal{B} , $\hat{\mathbf{p}}$.

332 **Lemma 5** (Sieving Lemma). Given a partition \mathcal{I} containing K intervals, sample access to \mathbf{p} on [n]333 and $\delta \in (0, 1)$. Then, the output of **Learn-And-Sieve** (Algorithm 1) satisfies the following. (i) Suppose 334 $\mathbf{p} \in \mathcal{H}_k^n$. Then the algorithm returns a positive measure $\hat{\mathbf{p}}$ and \mathcal{B} such that $d_{\chi^2}^{[n]\setminus\mathcal{B}}(\mathbf{p}\|\hat{\mathbf{p}}) \leq \varepsilon^2$ with probability at least $1 - \delta$. (ii) The output \mathcal{B} contains at most k intervals (if the algorithm does not 335 336 reject). (iii) At most $O((K/\varepsilon^2 + \sqrt{Kn}/\varepsilon) \cdot \log(n/\delta))$ samples are used. 337

Proof Sketch. We claim that, if $\mathbf{p} \in \mathcal{H}_k^n$, \mathcal{B} contains all the ε -bad intervals and no non-338 breakpoint intervals with high probability. Let I be a non-breakpoint interval. For b =339 $\Theta(m/\log(n/\delta)) = \Theta(K/\varepsilon^2 + \sqrt{Kn}/\varepsilon)$, we have, with high probability, $|\hat{\varphi}(Q) - \mathbf{p}(Q)| \leq |\hat{\varphi}(Q)| \leq |\hat{\varphi}(Q)| + |\hat{\varphi}(Q)|$ 340 $\sqrt{\mathbf{p}(Q)/b}, |\hat{\mathbf{p}}(Q) - \mathbf{p}(Q)| \le \sqrt{\mathbf{p}(Q)/b}$ and $\mathbf{p}(Q)/\hat{\mathbf{p}}(Q) \le \max(2, 8 \cdot n/b)$ which follow from 341 Lemmas $\overline{3}$ and $\overline{4}$. Combining this with triangle inequality and our choice of b implies the sec-342 ond condition of Line 5 will be false. The first condition can be shown to be false by rewriting 343 $\hat{\varphi}(Q)/\hat{\mathbf{p}}(Q)$ as $\hat{\varphi}(Q)/\mathbf{p}(Q) \cdot \mathbf{p}(Q)/\hat{\mathbf{p}}(Q)$, which are themselves bounded, with high probability, by 344 3 and $\Theta(1) \cdot \max(1, \varepsilon \sqrt{n/K})$ again by Lemmas 3 and 4 and our choice of b. 345

Let I be a breakpoint interval. We then have $|\mathbf{p}(Q) - \hat{\mathbf{p}}(Q)| \ge \sqrt{\hat{\mathbf{p}}(Q) \cdot \varepsilon^2/K}$ for some sub-interval 346 $Q \subset I$. If $\mathbf{p}(Q)$ is light $(\mathbf{p}(Q) \leq 2\varepsilon/\sqrt{Kn})$, we can show $\mathbf{p}(Q)/b \leq 1/4 \cdot \hat{\mathbf{p}}(Q) \cdot \varepsilon^2/K$, making 347 $\hat{\varphi}(Q)$, our estimation for $\mathbf{p}(Q)$, sufficiently accurate such that the second condition of Line 5 will 348 be true. Otherwise, as $b \gg \sqrt{Kn}/\varepsilon$, the estimation $\hat{\varphi}(Q)$ will be within multiplicative factors of 349 $\mathbf{p}(Q)$. If $\hat{\mathbf{p}}(Q)$ is not much lighter than $\mathbf{p}(Q)$, we can again show $\mathbf{p}(Q)/b \leq 1/4 \cdot \hat{\mathbf{p}}(Q) \cdot \varepsilon^2/K$. 350 Otherwise, the first condition of Line 5 will be true. Conditioned on that \mathcal{B} includes all ε -bad intervals 351 and no non-breakpoint intervals, it is easy to see that \mathcal{B} will contain no more than k intervals and 352 $d_{\chi^2}^{\mathcal{I}\setminus\mathcal{B}}(\mathbf{p}\|\hat{\mathbf{p}}) \leq O(\varepsilon^2)$. We note that points (i) and (iii) follow from the definition of the algorithm. \Box 353 **Learn-and-Sieve** (Algorithm 1) outputs a fully specified description $\hat{\mathbf{p}}$ and a sub-domain $\mathcal{G} := [n] \setminus \mathcal{B}$ such that $d_{\chi^2}^{\mathcal{G}}(\mathbf{p} \| \hat{\mathbf{p}})$ is small given $\mathbf{p} \in \mathcal{H}_n^k$. For testing purposes, this is a reasonable divergence 354 355 from the ideal guarantee that $d_{\chi^2}(\mathbf{p} \| \hat{\mathbf{p}})$ is small as long as $\mathbf{p}(\mathcal{B})$ is also small. If so, we can set 356 $\hat{\mathbf{p}}(i) = 0$ for $i \in \mathcal{B}$ and invoke **Tolerant-Identity-Test** with \mathbf{p} and $\hat{\mathbf{p}}$. If the test passes, we then know 357 that $\mathrm{TV}^{\mathcal{G}}(\mathbf{p}, \hat{\mathbf{p}}) \leq \varepsilon/2$: this together with $\mathbf{p}(\mathcal{B}) \leq \varepsilon/2$ then gives $\mathrm{TV}(\mathbf{p}, \hat{\mathbf{p}}) \leq \varepsilon$. 358

Unfortunately, running **Learn-and-Sieve** only once we may have $\mathbf{p}(\mathcal{B}) = \Omega(1)$. To handle this, we 359 will need more fine-grained sieving procedure, which uses Approx-Sub-Divide to further partition 360 the bad intervals detected and invokes Learn-and-Sieve *iteratively*. In each iteration, the total mass of 361 the bad intervals shrinks by a constant factor, allowing us to reach $\mathbf{p}(\mathcal{B}) \ll \varepsilon$ in at most $O(\log(1/\varepsilon))$ 362 iterations. The pseudo-code (Algorithm 4) and detailed analysis are provided in Appendix A.5. 363

Sample Complexity Lower Bound 3 364

In this section, we describe the hard instance of histogram testing, which leads to an $\Omega(\sqrt{kn}/\varepsilon + k/\varepsilon^2)$ 365 lower bound. We will apply the so-called Poissonization trick: we will relax P, the unknown object 366

- being tested, to be a positive measure with total mass $\Theta(1)$. We denote such a measure as an 367 approximate probability vector and give the corresponding notion of histogram. 368
- Definition 3 (Approximate Probability Vector). We define the set of ν -approximate probability 369 vectors (APV) on the domain [n] by $\tilde{\mathcal{P}}^n(\nu) \coloneqq \{P : P_i \in [0,\infty) \, \forall i \in [n], |||P||_1 - 1| \leq \nu\}.$ 370
- Accordingly, the set of histogram APV is given by $\tilde{\mathcal{H}}_{k}^{n}(\nu) := \{P \in \tilde{\mathcal{P}}^{n}(\nu) : P/||P||_{1} \in \mathcal{H}_{k}^{n}\}.$ 371

Under the Poisson sampling model, given an unknown $P \in \tilde{\mathcal{P}}^n(\nu)$, the goal it to decide whether 372 $P \in \tilde{\mathcal{H}}_k^n(\nu)$ or P is at least $\varepsilon(1+\nu)$ -fat⁴ from any $P' \in \tilde{\mathcal{H}}_k^n(\nu)$ in ℓ_1 distance when given the vector $\{M_1, M_2, \cdots, M_n\}$ where $M_i \sim \operatorname{Poi}(m \cdot P_i)$. We denote the sample complexity of the problem by 373 374 $m_{\text{hist}}^{\text{poi}}(n,k,\varepsilon,\nu)$ and provide its formal definition in Appendix B. 375

- To lower bound $m_{\text{hist}}^{\text{poi}}(n, k, \varepsilon, \nu)$, we follow the idea of *moment matching* illustrated in [Val11] VV13] 376 **WY16**]. In particular, one first constructs two discrete non-negative random variables U, U' whose 377 first few moments are identical. Moreover, U and U' will be designed to have different properties 378 such that one can use i.i.d. copies of U (and U') to generate random measures that are histograms 379 (and far-away from histograms respectively). 380

Our construction of such a pair of random variables is based on *Chebyshev's polynomials*, a standard 381 tool in approximation theory and the parameter estimation literature. The two variables will be 382 supported on the roots of the polynomial $p(x) = x \left(x - \frac{1}{n}\right) \left(x - \frac{2}{n}\right) T_d \left(1 - \frac{\sqrt{kn}}{C \cdot \log^2 n} \cdot x\right)$, where 383 $T_d(\cdot)$ is the Chebyshev's polynomial (of the first kind) and C is a sufficiently large constant. More 384 precisely, U will be supported on roots r where the derivatives p'(r) < 0, U' will be on roots where 385 p'(r) > 0, and the probabilities will be proportional to p'(r) accordingly. Consequently, U will most 386 likely be 1/n (hence useful for histogram construction) and U' will most likely be 0 or 2/n, each 387 with non-trivial probabilities (hence appropriate for non-histogram construction). Besides, they will 388 have their maximums bounded by $O(1/\sqrt{kn})$, which is crucial to achieve the nearly optimal lower 389 bounds. The detailed construction and analysis are provided in Appendix B.1. 390

Lemma 6. Given positive integers k, n where k < n, there exists a pair of non-negative random 391 variable U, U' supported on [0, 1) and absolute constants c, c' > 0 satisfying (i) $\Pr\left[U \neq \frac{1}{n}\right] \ll \frac{k}{n}$ 392 (*ii*) $\Pr[U'=0] > 1/3$ and $\Pr\left[U'=\frac{2}{n}\right] > 1/3$. (*iii*) $U, U' \le c' \log^2 n/\sqrt{kn}$. (*iv*) $\mathbf{E}[U] = \mathbf{E}[U'] = \mathbf{E}[U']$ 393 $\frac{1}{n}(1+O(\sqrt{k/n})). (v) \mathbf{E}[U^t] = \mathbf{E}[U^{\prime t}] \text{ for } 1 \le t \le c \cdot \log n.$ 394

We the proceed to construct two families of Approximate Probability Vectors, one of which belongs 395 to $\tilde{\mathcal{H}}_k^n$ and the other far from it using the random variables stated in Lemma 6. To do so, we 396 define $H = (1/n + \varepsilon U^{(1)}, \dots, 1/n + \varepsilon U^{(n)}), H' = (1/n + \varepsilon U'^{(1)}, \dots, 1/n + \varepsilon U'^{(n)})$ where $U^{(i)}, U'^{(i)}$ are n i.i.d. copies of U, U' in Lemma 6. 397 398

We address the two regimes $\sqrt{k/n} \le \varepsilon \log^2 n$ and $\sqrt{k/n} \ge \varepsilon \log^2 n$ separately. In the former case, 399 the heaviest element among H and H' are roughly $\widetilde{\Theta}(\varepsilon/\sqrt{kn})$. Hence, when the algorithm takes 400 $\widetilde{o}(\sqrt{kn}/\varepsilon)$ samples, it rarely sees any element appearing a large number of times. By the moment-401 matching property of U and U', the probabilities of seeing some elements appearing for t times for 402 $t \leq \log n$ are almost identical under H and H', therefore making H and H' indistinguishable. In the 403 latter case, we have $\varepsilon U \ll \frac{1}{n}$, implying that no elements in the measures are significantly heavier than 404 the rest. As a result, H and H' are both almost uniform except with a different number of "bumps" 405 (elements that are slightly heavier). Subsequently, the algorithm needs more samples (about $\widetilde{\Omega}(k/\varepsilon^2)$) 406 to tell whether a certain element is heavier than the rest, leading to a phase transition in the sample 407 complexity of the problem. We remark that whether $\widetilde{\Omega}(k/\varepsilon^2)$ or $\widetilde{\Omega}(\sqrt{nk}/\varepsilon)$ dominates depends 408 exactly on the relationship between $\sqrt{k/n}$ and ε (omitting poly-logarithmic factors). Combining the 409 two regimes then gives us the following lower bound, whose proof is provided in Appendix B.2. 410

Proposition 2. There exists a constant $\nu \in (0,1)$ such that for any sufficiently large n and $\varepsilon \in$ 411 $(0, 1/10), \text{ it holds } m_{\text{hist}}^{\text{poi}}(n, k, \varepsilon, \nu) \geq \Omega(\max(\sqrt{kn}/(\varepsilon \log n), k/(\varepsilon^2 \log^3 n))).$ 412

Finally, we can easily translate our lower bound result in the Poissonized sampling model to the 413 Multinomial (standard fixed-size) sampling model by a standard reduction. Combining it with the 414

known $\Omega(\sqrt{n}/\varepsilon^2)$ bound (see [Can16, Proposition 4.1]) then concludes our lower bound argument. 415

Formal proofs are given in Appendix B.3. 416

⁴The extra $(1 + \nu)$ factor accommodates the fact that P may not be a distribution, i.e., $1 \le ||P||_1 < (1 + \nu)$.

417 **References**

418 419 420 421	[ADH+15]	J Acharya, I. Diakonikolas, C. Hegde, J. Li, and L. Schmidt. Fast and near-optimal algorithms for approximating distributions by histograms. In Tova Milo and Diego Calvanese, editors, <i>Proceedings of the 34th ACM Symposium on Principles of Database Systems, PODS 2015</i> , pages 249–263. ACM, 2015.
422 423	[ADK15]	J. Acharya, C. Daskalakis, and G. Kamath. Optimal testing for properties of distributions. In <i>NeurIPS</i> , pages 3591–3599, 2015.
424 425 426 427	[ADLS17]	J. Acharya, I. Diakonikolas, J. Li, and L. Schmidt. Sample-optimal density estimation in nearly-linear time. In <i>Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2017</i> , pages 1278–1289, 2017. Full version available at https://arxiv.org/abs/1506.00671.
428 429 430	[BFR+00]	T. Batu, L. Fortnow, R. Rubinfeld, W. D. Smith, and P. White. Testing that distributions are close. In <i>IEEE Symposium on Foundations of Computer Science</i> , pages 259–269, 2000.
431 432	[BFR ⁺ 13]	T. Batu, L. Fortnow, R. Rubinfeld, W. D. Smith, and P. White. Testing closeness of discrete distributions. <i>J. ACM</i> , 60(1):4, 2013.
433	[Can]	C. L. Canonne. A short note on poisson tail bounds.
434 435 436 437	[Can16]	C. L. Canonne. Are few bins enough: Testing histogram distributions. In <i>Proceedings</i> of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS '16, page 455–463, New York, NY, USA, 2016. Association for Computing Machinery.
438 439	[Can20]	C. L. Canonne. A Survey on Distribution Testing: Your Data is Big. But is it Blue? Number 9 in Graduate Surveys. Theory of Computing Library, 2020.
440 441	[Can22]	C. L. Canonne, 2022. Personal communication. Corrigendum for [Can16] sent to the conference.
442 443 444	[CDGR18]	C. L. Canonne, I. Diakonikolas, T. Gouleakis, and R. Rubinfeld. Testing shape restrictions of discrete distributions. <i>Theory Comput. Syst.</i> , 62(1):4–62, 2018. Invited issue for STACS'16.
445 446	[CDSS14]	S. Chan, I. Diakonikolas, R. Servedio, and X. Sun. Near-optimal density estimation in near-linear time using variable-width histograms. In <i>NIPS</i> , pages 1844–1852, 2014.
447 448	[CGHJ12]	G. Cormode, M. Garofalakis, P. J. Haas, and C. Jermaine. Synopses for massive data: Samples, histograms, wavelets, sketches. <i>Found. Trends databases</i> , 4:1–294, 2012.
449 450	[CMN98]	S. Chaudhuri, R. Motwani, and V. R. Narasayya. Random sampling for histogram construction: How much is enough? In <i>SIGMOD Conference</i> , pages 436–447, 1998.
451 452 453	[DK16]	I. Diakonikolas and D. M. Kane. A new approach for testing properties of discrete distributions. In 2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS), pages 685–694. IEEE, 2016.
454 455 456	[DKN15a]	I. Diakonikolas, D. M. Kane, and V. Nikishkin. Optimal algorithms and lower bounds for testing closeness of structured distributions. In <i>IEEE 56th Annual Symposium on Foundations of Computer Science, FOCS 2015</i> , pages 1183–1202, 2015.
457 458 459	[DKN15b]	I. Diakonikolas, D. M. Kane, and V. Nikishkin. Testing identity of structured distribu- tions. In <i>Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete</i> <i>Algorithms, SODA 2015</i> , pages 1841–1854, 2015.
460 461 462	[DKN17]	I. Diakonikolas, D. M. Kane, and V. Nikishkin. Near-optimal closeness testing of discrete histogram distributions. In 44th International Colloquium on Automata, Languages, and Programming, ICALP 2017, pages 8:1–8:15, 2017.

463 464	[DL04]	L. Devroye and G. Lugosi. Bin width selection in multivariate histograms by the combinatorial method. <i>Test</i> , 13(1):129–145, 2004.
465 466 467 468	[DLS18]	I. Diakonikolas, J. Li, and L. Schmidt. Fast and sample near-optimal algorithms for learning multidimensional histograms. In <i>Conference On Learning Theory, COLT 2018</i> , volume 75 of <i>Proceedings of Machine Learning Research</i> , pages 819–842. PMLR, 2018.
469 470 471	[FD81]	D. Freedman and P. Diaconis. On the histogram as a density estimator:12 theory. <i>Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete</i> , 57(4):453–476, 1981.
472 473 474	[GGI+02]	A. C. Gilbert, S. Guha, P. Indyk, Y. Kotidis, S. Muthukrishnan, and M. Strauss. Fast, small-space algorithms for approximate histogram maintenance. In <i>STOC</i> , pages 389–398, 2002.
475 476	[GKS06]	S. Guha, N. Koudas, and K. Shim. Approximation and streaming algorithms for histogram construction problems. <i>ACM Trans. Database Syst.</i> , 31(1):396–438, 2006.
477 478	[GMP97]	P. B. Gibbons, Y. Matias, and V. Poosala. Fast incremental maintenance of approximate histograms. In <i>VLDB</i> , pages 466–475, 1997.
479	[Han19]	Y. Han. Mixture vs. mixture and moment matching, 2019.
480 481	[ILR12]	P. Indyk, R. Levi, and R. Rubinfeld. Approximating and Testing <i>k</i> -Histogram Distributions in Sub-linear Time. In <i>PODS</i> , pages 15–22, 2012.
482 483	[JKM ⁺ 98]	H. V. Jagadish, N. Koudas, S. Muthukrishnan, V. Poosala, K. C. Sevcik, and T. Suel. Optimal histograms with quality guarantees. In <i>VLDB</i> , pages 275–286, 1998.
484 485	[JVYHW15]	J. Jiao, K. Venkat, Yanjun Y. Han, and T. Weissman. Minimax estimation of functionals of discrete distributions. <i>IEEE Trans. Inform. Theory</i> , 61(5):2835–2885, 2015.
486 487	[Kle09]	J. Klemela. Multivariate histograms with data-dependent partitions. <i>Statistica Sinica</i> , 19(1):159–176, 2009.
488	[Koo80]	R. P. Kooi. The Optimization of Queries in Relational Databases. PhD thesis, 1980.
489 490 491 492	[KOPS15]	S. Kamath, A. Orlitsky, D. Pichapati, and A. T. Suresh. On learning distributions from their samples. In Peter Grünwald, Elad Hazan, and Satyen Kale, editors, <i>Proceedings of The 28th Conference on Learning Theory</i> , volume 40 of <i>Proceedings of Machine Learning Research</i> , pages 1066–1100, Paris, France, 03–06 Jul 2015. PMLR.
493 494	[LN96]	G. Lugosi and A. Nobel. Consistency of data-driven histogram methods for density estimation and classification. <i>Ann. Statist.</i> , 24(2):687–706, 04 1996.
495 496	[Pan08]	L. Paninski. A coincidence-based test for uniformity given very sparsely sampled discrete data. <i>IEEE Transactions on Information Theory</i> , 54(10):4750–4755, 2008.
497 498	[Rub06]	R. Rubinfeld. Sublinear time algorithms. Proceedings of the International Congress of Mathematicians (ICM), 2006.
499	[Sco79]	D. W. Scott. On optimal and data-based histograms. <i>Biometrika</i> , 66(3):605–610, 1979.
500 501	[TGIK02]	N. Thaper, S. Guha, P. Indyk, and N. Koudas. Dynamic multidimensional histograms. In <i>SIGMOD Conference</i> , pages 428–439, 2002.
502 503	[Val11]	P. Valiant. Testing symmetric properties of distributions. <i>SIAM J. Comput.</i> , 40(6):1927–1968, 2011.
504 505	[VV13]	G. Valiant and P. Valiant. Estimating the unseen: improved estimators for entropy and other properties. <i>Advances in Neural Information Processing Systems</i> , 26, 2013.
506 507	[WY16]	Y. Wu and P. Yang. Minimax rates of entropy estimation on large alphabets via best polynomial approximation. <i>IEEE Trans. Inform. Theory</i> , 62(6):3702–3720, 2016.
508 509	[WY ⁺ 19]	Y. Wu, P. Yang, et al. Chebyshev polynomials, moment matching, and optimal estimation of the unseen. <i>Annals of Statistics</i> , 47(2):857–883, 2019.

510 Checklist

511	1. For all authors
512 513	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
514	(b) Did you describe the limitations of your work? [Yes]
515	(c) Did you discuss any potential negative societal impacts of your work? [N/A]
516 517	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
518	2. If you are including theoretical results
519	(a) Did you state the full set of assumptions of all theoretical results? [Yes] See Theorem 1.
520 521	(b) Did you include complete proofs of all theoretical results? [Yes] See Appendix (Supplementary Materials).
522	3. If you ran experiments
523 524	(a) Did you include the code, data, and instructions needed to reproduce the main experi- mental results (either in the supplemental material or as a URL)? [N/A]
525 526	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]
527 528	(c) Did you report error bars (e.g., with respect to the random seed after running experi- ments multiple times)? [N/A]
529 530	(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]
531	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
532	(a) If your work uses existing assets, did you cite the creators? [N/A]
533	(b) Did you mention the license of the assets? [N/A]
534 535	(c) Did you include any new assets either in the supplemental material or as a URL? $\left[\mathrm{N/A}\right]$
536 537	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
538 539	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
540	5. If you used crowdsourcing or conducted research with human subjects
541 542	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
543 544	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
545 546	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]