

# A Lower Bound on the Complexity of Testing Grained Distributions\*

Oded Goldreich<sup>†</sup>      Dana Ron<sup>‡</sup>

July 13, 2022

## Abstract

A distribution is called  $m$ -grained if each element appears with probability that is an integer multiple of  $1/m$ . We prove that, for any constant  $c < 1$ , testing whether a distribution over  $[\Theta(m)]$  is  $m$ -grained requires  $\Omega(m^c)$  samples, where testing a property of distributions means distinguishing between distributions that have the property and distributions that are far (in total variation distance) from any distribution that has the property.

## 1 Introduction

This paper presents an almost linear lower bound on the sample complexity of testing a natural property of distributions. That is, we consider the number of samples that one should get from a distribution over  $[n]$  in order to distinguish between distributions that have the property and distributions that are far from having the property (i.e., far (in total variation distance) from any distribution that has the property). Hence, this paper refers to the model of distribution testing, which is reviewed next.

### 1.1 Distribution testing at a glance

Distribution testing was briefly introduced in [7, Sec. 3.4.3], and was first explicitly studied in [4, 3]. This research direction is a sub-area within property testing [6], which is the study of “approximate decision” procedures of sublinear complexity (i.e., the complexity is sublinear in the size of the object). In the case of distribution testing, the objects are (discrete) distributions, which are “accessed” via samples that are drawn independently (from the tested distribution). Hence, sublinear complexity means a number of samples that is smaller than the support of the distribution (or smaller than the domain on which the distribution is defined).<sup>1</sup> The “approximate” nature of testing is captured by the postulation of a promise problem in which the no-instances are far from the yes-instances (see Definition 1). Here we refer to the total variation distance, where the **total**

---

\*Partially supported by the Israel Science Foundation (grant No. 1041/18).

<sup>†</sup>Department of Computer Science, Weizmann Institute of Science, Rehovot, ISRAEL. E-mail: [oded.goldreich@weizmann.ac.il](mailto:oded.goldreich@weizmann.ac.il). Additional funding received from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No. 819702).

<sup>‡</sup>School of Electrical Engineering, Tel Aviv University, Tel Aviv, ISRAEL. [danaron@tau.ac.il](mailto:danaron@tau.ac.il)

<sup>1</sup>The support of a distribution  $P : \Omega \rightarrow [0, 1]$  is the set  $\{i \in \Omega : P(i) > 0\}$ . Indeed,  $\Omega$  is the domain on which  $P$  is defined (and  $\sum_{i \in \Omega} P(i) = 1$ ).

variation distance between the distributions  $P : \Omega \rightarrow [0, 1]$  and  $Q : \Omega \rightarrow [0, 1]$  equals

$$\frac{1}{2} \cdot \sum_{i \in \Omega} |P(i) - Q(i)|. \quad (1)$$

We say that  $P$  is  $\epsilon$ -far from  $Q$  if the total variation distance between them is greater than  $\epsilon$ , and otherwise we say that  $P$  is  $\epsilon$ -close to  $Q$ .

**Definition 1** (testing properties of distributions): *Let  $\mathcal{D} = \{\mathcal{D}_n\}_{n \in \mathbb{N}}$  be a property of distributions such that  $\mathcal{D}_n$  is a set of distributions over  $[n]$ , and  $s : \mathbb{N} \times (0, 1] \rightarrow \mathbb{N}$ . A tester, denoted  $T$ , of sample complexity  $s$  for the property  $\mathcal{D}$  is a probabilistic machine that, on input parameters  $n$  and  $\epsilon$ , and a sequence of  $s(n, \epsilon)$  samples drawn from an unknown distribution  $P$  over  $[n]$ , satisfies the following two conditions.*

1. The tester accepts distributions that belong to  $\mathcal{D}$ : *If  $P$  is in  $\mathcal{D}_n$ , then*

$$\Pr_{i_1, \dots, i_s \sim P}[T(n, \epsilon; i_1, \dots, i_s) = 1] \geq 2/3,$$

*where  $s = s(n, \epsilon)$  and  $i_1, \dots, i_s$  are drawn independently from the distribution  $P$ .*

2. The tester rejects distributions that are far from  $\mathcal{D}$ : *If  $P$  is  $\epsilon$ -far from any distribution in  $\mathcal{D}_n$  (i.e.,  $P$  is  $\epsilon$ -far from  $\mathcal{D}$ ) with respect to the total variation distance, then*

$$\Pr_{i_1, \dots, i_s \sim P}[T(n, \epsilon; i_1, \dots, i_s) = 0] \geq 2/3,$$

*where  $s = s(n, \epsilon)$  and  $i_1, \dots, i_s$  are as in the previous item.*

We say that testing  $\mathcal{D}$  requires  $s'(n)$  samples, if for some constant  $\epsilon > 0$  the sample complexity  $s$  of any tester of  $\mathcal{D}$  satisfies  $s(n, \epsilon) \geq s'(n)$ .

**Label-invariant properties.** A natural class of properties of distributions consist of properties that are insensitive to the names of the elements in the distribution. Such properties are called label-invariant, where a property of distributions over  $[n]$  is called label-invariant if for every bijection  $\pi : [n] \rightarrow [n]$  and every distribution  $P$ , it holds that  $P$  is in the property if and only if  $P \circ \pi$  is in the property, where  $(P \circ \pi)(y) = P(\pi(y))$ . A celebrated result of [10] asserts that any label-invariant property of distributions over  $[n]$  can be tested using  $O(\epsilon^{-2}n/\log n)$  samples (see also [6, Cor. 11.28]). On the other hand,  $\Omega(n/\log n)$  samples are required for testing some label-invariant properties (see [6, Cor. 11.30] as well as [6, Cor. 11.31], which are based on [10]).

## 1.2 Our result

As stated above, a few label-invariant properties are known to require  $\Omega(n/\log n)$  samples for testing. These properties include the set of distributions over  $[n]$  that have support size at most  $n/2$ . Our main result is a similar lower bound on the sample complexity of testing another natural property of distributions. Specifically, we refer to  $m$ -grained distributions, for some  $m = O(n)$ , and to a lower bound of the form  $n^c$ , for any constant  $c < 1$ .

A distribution  $P : \Omega \rightarrow [0, 1]$  is called  $m$ -grained if  $P(x)$  is a multiple of  $1/m$  for every  $x$  in  $\Omega$ ; that is, for each  $x \in \Omega$ , there exists an integer  $m_x$ , such that  $P(x) = m_x/m$  (see [6, Def. 11.7]).

Grained distributions have appeared implicitly in several prior works (most conspicuously in [8]), and were defined and studied explicitly in [5]. In particular, the challenge of determining the sample complexity of testing the set of grained distributions (i.e., the property of being grained) was raised explicitly in [5, Sec. 4].

It is quite easy to prove that testing the set of  $n$ -grained distributions requires  $\Omega(\sqrt{n})$  samples. In particular,  $\Omega(\sqrt{n})$  samples are required in order to distinguish the uniform distribution on  $[n]$  from a generic distribution that assigns probability  $1/2n$  to each of  $n/2$  elements and probability  $3n/2n$  to each of the remaining elements. To the best of our knowledge, this was the best lower bound known till this work.<sup>2</sup> In this work we obtain a lower bound of  $\Omega(n^c)$ , for any constant  $c < 1$ .

**Theorem 2** (main result): *For every constant  $c < 1$ , the sample complexity of testing whether a distribution over  $[n]$  is  $m$ -grained, where  $m = \Theta(n)$ , is  $\Omega(n^c)$ .*

We mention that the sample complexity of testing the foregoing property of distributions is  $O(\epsilon^{-2}n/\log n)$ ; this follows as a special case from the fact that any label-invariant property of distributions can be tested within this complexity [10] (see also [6, Cor. 11.28]). We conjecture that the aforementioned upper bound is tight; that is:

**Conjecture 3** *The sample complexity of testing  $\Theta(n)$ -grained distributions over  $[n]$  is  $\Omega(n/\log n)$ .*

We mention that the techniques used in our proof of Theorem 2 seem inadequate for proving a lower bound of the form  $\Omega(n^{1-o(1)})$ . In particular, our proof holds also when guaranteed that the tested distribution is  $2m$ -grained and assigns probability  $O(1/n)$  to each element in its support. However, under this promise, for any  $m = \Theta(n)$ , one can test if these distributions are  $m$ -grained using  $O(n^{1-\Omega(1)})$  samples. In fact, under this promise, one can even learn the distribution (up to relabeling) using  $O(n^{1-\Omega(1)})$  samples.<sup>3</sup>

## 2 Proof of Theorem 2

Our proof relies on two standard simplifying assumptions:

1. When considering the task of testing a label-invariant property, one may assume, without loss of generality, that the tester is label-invariant [1] (see also [6, Thm. 11.12]); that is, for every bijection  $\pi$  on the potential support, the tester's verdict on the samples  $i_1, \dots, i_s$  is identical to its verdict on the samples  $\pi(i_1), \dots, \pi(i_s)$ .

---

<sup>2</sup>We mention that a lower bound of  $\Omega(n/\log n)$  was known for the tolerant version [6, Thm. 11.31] in which, for some positive constants  $\delta < \epsilon$ , one is required to distinguish distributions that are  $\delta$ -close to being  $n$ -grained from distributions that are  $\epsilon$ -far from being  $n$ -grained. Subsequent to the posting of the first version of the current paper, Clement Canonne suggested to us a simple proof of an  $\Omega(n^{2/3})$  lower bound (for the standard testing problem). His proof is presented in Appendix A.

<sup>3</sup>Indeed, suppose that a distribution  $P : [n] \rightarrow [0, 1]$  is guaranteed to satisfy  $P(i) \leq t/m$  for every  $i \in [n]$ . Using the promise that  $P$  is  $2m$ -grained, for each  $i \in [n]$  it holds that  $P(i) \in \{j/2m : j \in \{0, 1, \dots, 2t\}\}$ . Note that the histogram  $(h_0, \dots, h_{2t})$  of  $P$ , defined as  $h_j = |\{i \in [n] : P(i) = j/2m\}|$ , is determined by the probabilities of  $k$ -way collisions for  $k \in \{2, \dots, 2t + 2\}$  (cf., the proof of Theorem 2). Observing that the probability of  $k$ -way collisions can be approximated using  $O(n^{(k-1)/k})$  samples of  $P$ , it follows that the histogram of  $P$  can be (approximately) learned using  $O(n^{(2t+1)/(2t+2)})$  samples of  $P$ , where the approximation is up to a constant deviation.

2. To prove a lower bound of  $L$  on the sample complexity of testing, it suffices to describe two distributions  $P$  and  $Q$  that no algorithm of sample complexity  $L - 1$  can distinguish (with gap  $\Omega(1)$ )<sup>4</sup> such that  $P$  has the property and  $Q$  is  $\Omega(1)$ -far from having the property (cf. [6, Thm. 7.2]).

Combining these two observations, we focus on presenting distributions that cannot be distinguished by label-invariant algorithms of low complexity such that one distribution is  $m$ -grained while the other is  $\Omega(1)$ -far from being  $m$ -grained.

Both distributions that we present are specified by their histograms, which specify how many elements are assigned each value of the probability weight. For  $t = O(1/(1 - c))$ , in both distributions, for some  $m \in \{n/2, \dots, tn/2\}$ , each element in  $[n]$  is assigned weight  $\frac{i}{2m}$  such that  $i \in [t]$ . In particular, for  $n_i^P$ 's and  $n_i^Q$ 's that will be fully determined later, we have:

**Distribution  $P$ :** In this distribution,  $n_i^P$  elements in  $[n]$  are assigned the weight  $\frac{i}{2m}$ , and  $n_i^P = 0$  for every odd  $i \in [t]$ .

**Distribution  $Q$ :** In this distribution,  $n_i^Q$  elements in  $[n]$  are assigned the weight  $\frac{i}{2m}$ , and  $n_i^Q = 0$  for every even  $i \in [t]$ .

Note that  $\sum_{i \in [t]} n_i^P \cdot \frac{i}{2m} = 1 = \sum_{i \in [t]} n_i^Q \cdot \frac{i}{2m}$  and  $\sum_{i \in [t]} n_i^P = n = \sum_{i \in [t]} n_i^Q$ . (Indeed  $2m = \sum_{i \in [t]} n_i^P \cdot i$  implies  $2m \in \{n, \dots, tn\}$ .) Furthermore,  $P$  is  $m$ -grained, whereas  $Q$  is  $\frac{1}{3t}$ -far from being  $m$ -grained (because the weight of each element in the support of  $Q$  has to be modified by at least  $\frac{1}{2m}$  units whereas  $\frac{n}{2m} \geq \frac{1}{t}$ ).

Note that the equation  $\sum_{i \in [t]} n_i^P = \sum_{i \in [t]} n_i^Q$  asserts that both distributions have the same support size, whereas  $\sum_{i \in [t]} n_i^P \cdot i = \sum_{i \in [t]} n_i^Q \cdot i$  asserts that they are assigned the same total probability mass (in terms of units of  $\frac{1}{2m}$ ).

Intuitively (as is indeed proved in Lemma 6), a sample complexity lower bound of  $\Omega(n^{\frac{t-2}{t-1}})$  is related to showing that, for every  $k \in \{2, \dots, t-2\}$ , the probability of a  $k$ -way collision is the same in both distributions. Thus, we require that  $\sum_{i \in [t]} n_i^P \cdot \left(\frac{i}{2m}\right)^k = \sum_{i \in [t]} n_i^Q \cdot \left(\frac{i}{2m}\right)^k$  for every  $k \in \{2, \dots, t-2\}$ , which raises the question of whether such a setting of  $n_i^P$ 's and  $n_i^Q$ 's is possible. Before addressing the latter question (as well as the question of why this yields the desired lower bound), we reformulate the foregoing  $t-1$  equations in a uniform manner; that is, for every  $k \in [[t-2]] \stackrel{\text{def}}{=} \{0, 1, \dots, t-2\}$ , we require

$$\sum_{i \in [t]} n_i^P \cdot i^k = \sum_{i \in [t]} n_i^Q \cdot i^k. \quad (2)$$

Recalling the  $t$  initial equalities (i.e.,  $n_i^P = 0$  for odd  $i \in [t]$  and  $n_i^Q = 0$  for even  $i \in [t]$ ), we write the foregoing linear system in a matrix form as  $Ax = 0$ , where  $x = (n_1^P, \dots, n_t^P, n_1^Q, \dots, n_t^Q)^\top$ . For  $i \in [t]$ , the  $i^{\text{th}}$  row of the  $(2t-1)$ -by- $2t$  matrix  $A$  is  $(0^{i-1}10^{2t-i})$  if  $i$  is odd, and  $(0^{t+i-1}10^{t-i})$  if  $i$  is even, whereas (for  $k \in \{0, 1, \dots, t-2\}$ ) the  $(t+k+1)^{\text{th}}$  row of  $A$  is  $(1^k, 2^k, \dots, t^k, -1^k, -2^k, \dots, -t^k)$ . Figure 1 depicts  $A$  in the case of  $t = 5$ .

<sup>4</sup>We say that  $A$  distinguishes  $s$  samples of  $P$  from  $s$  samples of  $P$  with gap  $\gamma$  if

$$|\Pr_{i_1, \dots, i_s \sim P}[A(i_1, \dots, i_s) = 1] - \Pr_{i_1, \dots, i_s \sim P}[A(i_1, \dots, i_s) = 1]| \geq \gamma.$$

<b>1</b>	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	<b>1</b>	0	0	0
0	0	<b>1</b>	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	<b>1</b>	0
0	0	0	0	<b>1</b>	0	0	0	0	0
<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>-1</b>	<b>-1</b>	<b>-1</b>	<b>-1</b>	<b>-1</b>
<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>-1</b>	<b>-2</b>	<b>-3</b>	<b>-4</b>	<b>-5</b>
<b>1<sup>2</sup></b>	<b>2<sup>2</sup></b>	<b>3<sup>2</sup></b>	<b>4<sup>2</sup></b>	<b>5<sup>2</sup></b>	<b>-1<sup>2</sup></b>	<b>-2<sup>2</sup></b>	<b>-3<sup>2</sup></b>	<b>-4<sup>2</sup></b>	<b>-5<sup>2</sup></b>
<b>1<sup>3</sup></b>	<b>2<sup>3</sup></b>	<b>3<sup>3</sup></b>	<b>4<sup>3</sup></b>	<b>5<sup>3</sup></b>	<b>-1<sup>3</sup></b>	<b>-2<sup>3</sup></b>	<b>-3<sup>3</sup></b>	<b>-4<sup>3</sup></b>	<b>-5<sup>3</sup></b>

Figure 1: The matrix  $A$  and the submatrix considered in the analysis.

We seek a solution  $x$  that is *positive*, which means that each of the entries of  $x$  is non-negative, and at least one of the entries is positive. It turns out that such a solution exists if and only if for every  $v \in \mathbb{R}^{2t}$  it holds that  $vA$  is *not* strongly positive [9, Thm. 15.1(2)], where  $u$  is strongly positive if all its entries are positive.

**Claim 4** ( $vA$  is not strongly positive): *For every  $v \in \mathbb{R}^{2t}$ , not all entries of  $vA$  are positive.*

**Proof:** It suffice to show that it not possible that the entries of  $vA$  that correspond to even  $i$ 's in  $[t]$  and to  $t+i$ 's for odd  $i$ 's (in  $[t]$ ) are all positive. To verify this, observe that the first  $t$  rows in the corresponding columns (i.e., the  $i^{\text{th}}$  column for even  $i \in [t]$  and the  $t+i^{\text{th}}$  column for odd  $i \in [t]$ ) are all-zero. Hence, for even  $i \in [t]$  the value of the  $i^{\text{th}}$  entry (in  $vA$ ) is  $\sum_{k \in [[t-2]]} v_{t+k+1} i^k$ , whereas for odd  $i \in [t]$  the value of the  $(t+i)^{\text{th}}$  entry is  $-\sum_{k \in [[t-2]]} v_{t+k+1} i^k$ . It follows that  $\sum_{k \in [[t-2]]} v_{t+k+1} i^k$  should be positive if  $i \in [t]$  is even, and negative otherwise. But this is impossible since the degree of this polynomial (in  $i$ ) is  $t-2$  (and so its sign cannot alternate  $t-1$  times). ■

**The actual  $n_i^{\text{P}}$ 's and  $n_i^{\text{Q}}$ 's.** The foregoing discussion and claim establishes the existence of non-negative and non-trivial  $n_i^{\text{P}}$ 's and  $n_i^{\text{Q}}$ 's that satisfy Eq. (2) for every  $k \in [[t-2]]$  as well as  $n_i^{\text{P}} = 0$  for odd  $i \in [t]$  and  $n_i^{\text{Q}} = 0$  for even  $i \in [t]$ , where non-trivial means that not all  $\xi$ 's and  $n_i^{\text{Q}}$ 's are zero. These  $n_i^{\text{P}}$ 's and  $n_i^{\text{Q}}$ 's may be assumed to be rational, but they do not necessarily sum-up to  $n$  nor are integers. In fact, these  $n_i^{\text{P}}$ 's and  $n_i^{\text{Q}}$ 's are independent of  $n$ , and so by multiplying them with an adequate number (e.g., the least common multiplier of their denominators) we obtain integers. Hence, we can fit any  $n$  that is an integer multiple of the sum of the resulting  $n_i^{\text{P}}$ 's (and, we can handle other  $n$ 's by “padding”).

We have thus established that distributions  $P$  and  $Q$  as postulated above do exist; that is,  $P$  and  $Q$  are  $2m$ -grained, and it holds that  $n_i^{\text{P}} = |\{j \in [n] : P(j) = \frac{i}{2m}\}|$  and  $n_i^{\text{Q}} = |\{j \in [n] : Q(j) = \frac{i}{2m}\}|$  satisfy Eq. (2) for every  $k \in [[t-2]]$  as well as  $n_i^{\text{P}} = 0$  for odd  $i \in [t]$  and  $n_i^{\text{Q}} = 0$  for even  $i \in [t]$ . In order to proceed, we restate the features of the  $n_i^{\text{P}}$ 's and  $n_i^{\text{Q}}$ 's in terms of the (probability) histograms of  $P$  and  $Q$  (or rather their “normalized” forms). Specifically, consider the following random variable:  $X = i$  with probability  $\frac{n_i^{\text{P}}}{n}$  (resp.,  $Y = i$  with probability  $\frac{n_i^{\text{Q}}}{n}$ ), representing the fact

that there are  $n_i^P$  (resp.,  $n_i^Q$ ) elements in the support of  $P$  (resp.,  $Q$ ) that are assigned probability  $\frac{i}{2m}$ . Observe that  $\mathbb{E}[X^k] = \sum_{i \in [t]} \frac{n_i^P}{n} \cdot i^k$  (resp.,  $\mathbb{E}[Y^k] = \sum_{i \in [t]} \frac{n_i^Q}{n} \cdot i^k$ ). Hence, we have established the following:

**Lemma 5** (main lemma): *For every constant  $t \in \mathbb{N}$  and  $m, n \in \mathbb{N}$  such that  $m \in \{0.5n, \dots, 0.5tn\}$ , there exist  $2m$ -grained distributions  $P$  and  $Q$  over  $[n]$  such that the following conditions hold.*

1.  $P$  is  $m$ -grained, whereas  $Q$  is  $\frac{1}{3t}$ -far from being  $m$ -grained.
2. For every  $k \in [t - 2]$ , it holds that  $\mathbb{E}[X^k] = \mathbb{E}[Y^k]$ , where  $X$  and  $Y$  represent the histograms of  $P$  and  $Q$  (i.e.,  $\Pr[X = i] = n_i^P/n$  and  $\Pr[Y = i] = n_i^Q/n$ ).

At this point we can apply a result of [8], which we modify and rephrase as follows.<sup>5</sup>

**Lemma 6** (a variant of [8, Thm. 5.6]): *Let  $P$  and  $Q$  be  $2m$ -grained distributions over  $[n]$  such that their support equals  $[n]$ , and  $a_1, \dots, a_t \in \mathbb{N}$  such that for every  $j \in [n]$  it holds that  $P(j) \in \{\frac{a_i}{2m} : i \in [t]\}$  and  $Q(j) \in \{\frac{a_i}{2m} : i \in [t]\}$ . Define a random variable  $X$  (resp.,  $Y$ ) over  $[t]$  such that  $X = i$  (resp.,  $Y = i$ ) with probability that represents the fraction of elements in  $[n]$  that are assigned probability  $\frac{a_i}{2m}$  by  $P$  (resp.,  $Q$ ). If, for every  $k \in [t - 2]$ , it holds that  $\mathbb{E}[X^k] = \mathbb{E}[Y^k]$ , then the distinguishing gap of any label-invariant algorithm between  $s$  samples of  $P$  and  $s$  samples of  $Q$  is upper-bounded by*

$$O\left(\frac{t^2 \cdot s}{m/a} + \frac{s^{t-1}}{(m/a)^{t-2}}\right) + \exp(-\Omega(s)), \quad (3)$$

where  $a = \max_{i \in [t]} \{a_i\}$ .

Note that the upper bound of Eq. (3) holds vacuously for  $s > m/a$ . On the other hand, for non-constant  $s = o(m/(t^2 a))$ , Eq. (3) yields  $\delta \stackrel{\text{def}}{=} O((a \cdot s)^{t-1}/m^{t-2}) + o(1)$ ; that is, for any label-invariant algorithm, the distinguishing gap between  $s$  samples of  $P$  and  $s$  samples of  $Q$  is at most  $\delta$ . In particular, for any non-constant  $s = o(m/a)^{(t-2)/(t-1)}$ , it holds that  $\delta = o(1)$ . Hence, combining Lemmas 5 and 6, while setting  $t = O(1)$  and  $a_i = i$  for every  $i \in [t]$ , we infer that  $s = \Omega(m/t)^{(t-2)/(t-1)}$  samples are required for testing  $m$ -grained distributions. Theorem 2 follows by setting  $t = \lceil 1/(1 - c) \rceil + 1$ .

**Digest.** The proof of Theorem 2 is pivoted at a pair of distributions (over  $[\Theta(m)]$ ) that have identical  $k$ -way collision probabilities, for every  $k \in \{2, \dots, t - 2\}$ . On the one hand, these distributions cannot be distinguished by a label-invariant algorithm that takes  $o(m/t)^{(t-2)/(t-1)}$  samples. On the other hand, the first distribution is  $m$ -grained and the second distribution is  $1/3t$ -far from being  $m$ -grained.

While Theorem 2 follows by setting  $t = \lceil 1/(1 - c) \rceil + 1$ , where  $c < 1$  is a constant, the same argument holds for arbitrary  $t = \omega(1)$ . Specifically, using  $t = \log_2 m + 1$ , we get a lower bound of  $\Omega(m/\log m)$ , but this lower bound refers to testing  $m$ -grained distributions with proximity parameter set to  $1/3t = \Theta(1/\log m)$ . In contrast, the corresponding upper bound in this case (i.e., proximity parameter set to  $1/3t = \Theta(1/\log m)$ ) is  $O(m \log m)$ .

<sup>5</sup>Putting aside the many notational modifications, the actual difference is that Lemma 6 refers to the first  $t - 2$  powers of  $X$  and  $Y$ , whereas [8, Thm. 5.6] refers to the first  $t - 1$  powers. In fact, we present a generalization of [8, Thm. 5.6] in which the number of powers is a free parameter. In Appendix B we outline how this generalization (and in particular Lemma 6) follows from the proof of [8, Thm. 5.6].

## Acknowledgements

We are grateful to Clement Canonne for communicating to us the lower bound presented in Appendix A, and allowing us to include it in this paper.

## References

- [1] Tugkan Batu. *Testing properties of distributions*. PhD thesis, Computer Science department, Cornell University, 2001.
- [2] Tugkan Batu and Clement L. Canonne. Generalized Uniformity Testing. In *58th FOCS*, pages 880–889, 2017.
- [3] Tugkan Batu, Lance Fortnow, Eldar Fischer, Ravi Kumar, Ronitt Rubinfeld, and Patrick White. Testing random variables for independence and identity. In *42nd IEEE Symposium on Foundations of Computer Science*, pages 442–451, 2001.
- [4] Tugkan Batu, Lance Fortnow, Ronitt Rubinfeld, Warren D. Smith, and Patrick White. Testing that distributions are close. In *41st IEEE Symposium on Foundations of Computer Science*, pages 259–269, 2000.
- [5] Oded Goldreich. The Uniform Distribution is Complete with respect to Testing Identity to a Fixed Distribution. *ECCC*, TR16-015, February 2016.
- [6] Oded Goldreich. *Introduction to Property Testing*. Cambridge University Press, 2017.
- [7] O. Goldreich, S. Goldwasser, and D. Ron. Property testing and its connection to learning and approximation. *Journal of the ACM*, pages 653–750, July 1998.
- [8] Sofya Raskhodnikova, Dana Ron, Amir Shpilka, and Adam Smith. Strong Lower Bounds for Approximating Distribution Support Size and the Distinct Elements Problem. *SIAM Journal on Computing*, Vol. 39 (3), pages 813–842, 2009. Extended abstract in *48th FOCS*, 2007.
- [9] Steven Roman. *Advanced Linear Algebra*. Graduate Texts in Mathematics, Vol. 135, Springer, 2005.
- [10] 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 *43rd ACM Symposium on the Theory of Computing*, pages 685–694, 2011.

## Appendix A: An $\Omega(n^{2/3})$ lower bound (by Canonne)

The following proof was communicated to us by Clement Canonne, after the initial posting of this paper on *ECCC* (TR21-129). It is pivoted in the following result that is implicit in [2, Sec. IV], where [2, Thm. I.2] is proved.

**Proposition 7** (implicit in [2, Sec. IV]): *For any distribution  $P : [n] \rightarrow [0, 1]$ , distinguishing (say with gap  $1/4$ ) between a distribution that equals  $P$  up to relabeling and a distribution that is uniform on a set of size  $1/\|P\|_2^2$  requires  $\Omega(1/\|P\|_3)$  samples.*

For  $m = n/2$ , we pick  $P$  such that  $P(i) = 1/2m$  for every  $i \in [m]$ , and  $P(m+i) = 3/2m$  for every  $i \in [m/3]$ . Note that  $P$  is 0.6-far from being  $m$ -grained, whereas  $\|P\|_2 = (m \cdot (1/2m)^2 + (m/3) \cdot (3/2m)^2)^{1/2} = m^{-1/2}$  and  $\|P\|_3 = (m \cdot (1/2m)^3 + (m/3) \cdot (3/2m)^3)^{1/3} = (5/4m^2)^{-1/3}$ . Hence, distinguishing between the set of relabeling of  $P$  (which are all far from  $m$ -grained) and the set of distributions that are uniform on  $m$  elements (which hence are all  $m$ -grained) requires  $\Omega(m^{2/3})$  samples.

## Appendix B: Deriving Lemma 6 from the proof of [8, Thm. 5.6]

There are several differences between Lemma 6 and [8, Thm. 5.6].

1. Lemma 6 refers to algorithms that obtain samples drawn from ( $2m$ -grained) distributions whereas [8, Thm. 5.6] refers to algorithms that see the colors of balls drawn uniformly and independently (with replacement) among  $N$  balls.

Note that samples drawn from a  $2m$ -grained distribution over  $[n]$  correspond to the colors of uniformly selected balls, where the number of balls equals  $2m$  and the number of colors is  $n$ . That is, a  $2m$ -grained distribution  $D$  corresponds to a collection of  $2m$  balls such that (for every  $\chi \in [n]$ ) exactly  $2m \cdot D(\chi)$  balls are assigned the color  $\chi$ .

2. Lemma 6 refers to algorithms that obtain  $s$  samples, whereas [8, Thm. 5.6] refers to algorithms that obtain  $\text{Poi}(s)$  balls, where  $\text{Poi}(s)$  denotes the Poisson distribution with parameter  $s$ .

Recall that  $\Pr[\text{Poi}(s) < s/2] = \exp(-\Omega(s))$ , which means that an algorithm that gets  $\text{Poi}(s)$  samples can emulate an algorithm that expects  $s/2$  samples, with error probability  $\exp(-\Omega(s))$ . The latter error term is accounted for by the last term in Eq. (3).

3. In Lemma 6 the distribution  $P$  and  $Q$  play the main role while their histograms  $X$  and  $Y$  appear as secondary players, whereas in [8, Thm. 5.6] the histograms appear as main players and the corresponding distributions of colors appear in the second role.
4. Most importantly, Lemma 6 presupposes equality between the first  $t - 2$  powers of  $X$  and  $Y$ , whereas in [8, Thm. 5.6] the hypothesis refers to the first  $t - 1$  powers (but merely presupposes that they are at a fixed proportion). Recall that  $t$  denotes the number of values assigned to  $X$  and  $Y$  (i.e.,  $X, Y \in \{a_1, \dots, a_t\}$ ).

However, we observe (see details below) that the actual proof of [8, Thm. 5.6] supports a generalization in which the number of powers is  $d - 1$ , where  $d$  and  $t$  are free parameters. Hence, we may use  $d = t - 1$  (for our application) rather than  $d = t$  (as in [8, Thm. 5.6]).

We now turn to reviewing the presentation of [8], but we do so while using slightly different notation.<sup>6</sup> This presentation refers to  $N$  balls, where each ball has a *color*, and there are  $n$  colors. The presentation starts from a histogram that describes the frequencies of colors that appear in a specific number of balls; that is, for natural numbers  $a_1 < a_2 < \dots < a_t$  and non-negative  $p_1, \dots, p_t$  that sum-up to 1, a  $p_i$  fraction of the colors each occur in  $a_i$  balls (i.e.,  $|C_i| = p_i \cdot n$ , and for each  $\chi \in C_i$  there are  $a_i$  balls that have color  $\chi$ ).

<sup>6</sup>For example, we replace  $n$  by  $N$  (as denoting the number of balls), replace  $k$  by  $t$ , and  $(a_1, \dots, a_t)$  by  $(a_0, \dots, a_{k-1})$ . The number of colors is implicit in [8], but is explicit here.



The actual presentation of [8] starts with a random variable  $\Phi$  that ranges over  $\{a_1, \dots, a_t\} \subset \mathbb{N}$ , and lets  $p_i = \Pr[\Phi = a_i]$ . Given  $\Phi$  and an integer  $N$ , it defines the following instance of the *colored balls* problem, denoted  $B_{\Phi, N}$ : For each  $i \in [t]$ , there are  $\lfloor Np_i/E[\Phi] \rfloor$  colors of type  $i$  such each color of type  $i$  occurs in  $a_i$  balls. In our case, the  $p_i$ 's are multiples of  $1/n$  and  $N = \sum_{i \in [t]} p_i \cdot n \cdot a_i$  is an integer, which implies that

$$\frac{Np_i}{E[\Phi]} = p_i \cdot \frac{\sum_{j \in [t]} p_j \cdot n \cdot a_j}{\sum_{j \in [t]} p_j \cdot a_j} = p_i \cdot n$$

is an integer (and there is no need for additional tweaks as in [8]). That is, there are  $n_i = p_i n$  colors of type  $i$ , and the total number of balls is  $\sum_{i \in [t]} n_i \cdot a_i$ , which equals  $2m$  in our case. We next state a generalization of [8, Thm. 5.6], in which the hypothesis refers to the first  $d-1$  powers of  $\Phi_1$  and  $\Phi_2$ , while noting that in [8, Thm. 5.6]  $d = t$  (whereas in our application  $d = t-1$ ).

**Lemma 8** (a generalization of [8, Thm. 5.6], slightly rephrased):<sup>7</sup> *Let  $\Phi_1$  and  $\Phi_2$  be random variables over positive integers  $a_1 < a_2 < \dots < a_t$  such that*

$$\frac{E[\Phi_1]}{E[\Phi_2]} = \frac{E[\Phi_1^2]}{E[\Phi_2^2]} = \dots = \frac{E[\Phi_1^{d-1}]}{E[\Phi_2^{d-1}]} \quad (4)$$

*Then, for natural numbers  $N$  and  $s \leq \frac{N}{2a_t}$ , the distinguishing gap between  $B_{\Phi_1, N}$  and  $B_{\Phi_2, N}$  as judged by any label-invariant algorithm that takes  $\text{Poi}(2s)$  samples is upper-bounded by*

$$O\left(\frac{t \cdot d \cdot 2s}{N/a_t} + \frac{d}{\lfloor d/2 \rfloor! \cdot \lceil d/2 \rceil!} \cdot \frac{(2s)^d}{(N/a_t)^{d-1}}\right) \quad (5)$$

Lemma 6 follows from Lemma 8 by using  $\Phi_1 = X$  and  $\Phi_2 = Y$ , observing that  $N = 2m$  and  $B_{\Phi_1, N} \equiv P$  (resp.,  $B_{\Phi_2, N} \equiv Q$ ), setting  $d = t-1$ , simplifying Eq. (5), and using  $\Pr[\text{Poi}(2s) < s] = \exp(-\Omega(s))$ .

Recall that Lemma 8 generalizes [8, Thm. 5.6] by allowing  $d$  and  $t$  to be arbitrary natural numbers rather than mandating that  $d = t$ . However, the proof of [8, Thm. 5.6] does not use  $d = t$  in an essential manner, and so going over that proof one merely needs to keep track of when  $k$  (i.e., the parameter used in [8]) stands for  $t$  and when  $k$  stands for  $d$  (and observing that in all places  $a_{k-1}$  merely stands for the maximal  $a_i$ ).<sup>8</sup> In particular, denoting  $a = \max_{i \in [t]} \{a_i\}$ , the upper bound in [8, Lem. 5.9] is  $\delta_1 \stackrel{\text{def}}{=} O\left(\frac{a^{d-1}}{d!} \cdot \frac{(2s)^d}{N^{d-1}}\right)$ , the upper bound in [8, Lem. 5.10] is  $\delta_2 \stackrel{\text{def}}{=} \frac{2t \cdot a \cdot 2s}{N}$ , the upper bound  $\delta_3$  in [8, Lem. 5.12] is a  $\Theta(1/d)$  factor of the bound in Eq. (5), and the final upper bound is  $2 \cdot \delta_1 + 2 \cdot \delta_2 + (d-1) \cdot \delta_3$ , which matches Eq. (5).

---

<sup>7</sup>In the case of  $d = t$ , our rephrasing is merely notational (e.g.,  $(a_1, \dots, a_t)$  replaces  $(a_0, \dots, a_{k-1})$ , and  $N$  replaces  $n$ ). In addition, we incorporate Eq. (4) in our formulation of the lemma rather than referring to a notion (i.e. “proportional moments”) defined before, and avoid a notation for the gap of an algorithm (i.e., a notation as in Footnote 4 is avoided in Eq. (5)).

<sup>8</sup>Recall that the parameter  $s$  in [8] is replaced here by  $2s$ , and  $n$  is replaced by  $N$ .