

1 AN AUTOMATIC INEQUALITY PROVER AND INSTANCE  
2 OPTIMAL IDENTITY TESTING\*

3 GREGORY VALIANT† AND PAUL VALIANT‡

4 **Abstract.** We consider the problem of verifying the identity of a distribution: Given the  
5 description of a distribution over a discrete finite or countably infinite support,  $p = (p_1, p_2, \dots)$ , how  
6 many samples (independent draws) must one obtain from an unknown distribution,  $q$ , to distinguish,  
7 with high probability, the case that  $p = q$  from the case that the total variation distance ( $L_1$  distance)  
8  $\|p - q\|_1 \geq \epsilon$ ? We resolve this question, up to constant factors, on an *instance by instance* basis: there  
9 exist universal constants  $c, c'$  and a function  $f(p, \epsilon)$  on the known distribution  $p$  and error parameter  
10  $\epsilon$ , such that our tester distinguishes  $p = q$  from  $\|p - q\|_1 \geq \epsilon$  using  $f(p, \epsilon)$  samples with success  
11 probability  $> 2/3$ , but no tester can distinguish  $p = q$  from  $\|p - q\|_1 \geq c \cdot \epsilon$  when given  $c' \cdot f(p, \epsilon)$   
12 samples. The function  $f(p, \epsilon)$  is upper-bounded by a multiple of  $\frac{\|p\|_{2/3}}{\epsilon^2}$ , but is more complicated.  
13 This result generalizes and tightens previous results: since distributions of support at most  $n$  have  
14  $L_{2/3}$  norm bounded by  $\sqrt{n}$ , this result immediately shows that for such distributions,  $O(\sqrt{n}/\epsilon^2)$   
15 samples suffice, tightening the previous bound of  $O(\frac{\sqrt{n} \text{polylog } n}{\epsilon^4})$  and matching the (tight) results  
16 for the case that  $p$  is the uniform distribution of support  $n$ .

The analysis of our very simple testing algorithm involves several hairy inequalities. To facilitate  
this analysis, we give a complete characterization of a general class of inequalities—generalizing  
Cauchy-Schwarz, Hölder’s inequality, and the monotonicity of  $L_p$  norms. Specifically, we characterize  
the set of sequences of triples  $(a, b, c)_i = (a_1, b_1, c_1), \dots, (a_r, b_r, c_r)$  for which it holds that for all  
finite sequences of positive numbers  $(x)_j = x_1, \dots$  and  $(y)_j = y_1, \dots$ ,

$$\prod_{i=1}^r \left( \sum_j x_j^{a_i} y_j^{b_i} \right)^{c_i} \geq 1.$$

17 For example, the standard Cauchy-Schwarz inequality corresponds to the triples  $(a, b, c)_i = (1, 0, \frac{1}{2})$ ,  
18  $(0, 1, \frac{1}{2})$ ,  $(\frac{1}{2}, \frac{1}{2}, -1)$ . Our characterization is constructive and algorithmic, leveraging linear program-  
19 ming to prove or refute an inequality, which would otherwise have to be investigated, through trial  
20 and error, by hand. We hope the computational nature of our characterization will be useful to  
21 others, and facilitate analyses like the one here.

22 **Key words.** Hypothesis testing, identity testing, instance optimal, Holder’s inequality

23 **AMS subject classifications.** 68Q32, 26D15, 62G10

24 **1. Introduction.** Suppose you have a detailed record of the distribution of IP  
25 addresses that visit your website. You recently proved an amazing theorem, and are  
26 keen to determine whether this result has changed the distribution of visitors to your  
27 website (or is it simply that the usual crowd is visiting your website more often?). How  
28 many visitors must you observe to decide this, and, algorithmically, how do you decide  
29 this? Formally, given some known distribution  $p$  over a discrete (though possibly  
30 infinite) domain, a parameter  $\epsilon > 0$ , and an unknown distribution  $q$  from which we  
31 may draw independent samples, we would like an algorithm that will distinguish the  
32 case that  $q = p$  from the case that the total variation distance,  $d_{tv}(p, q) > \epsilon$ . We  
33 consider this basic question of verifying the identity of a distribution, also known as

---

\*A preliminary version of this work appeared at the IEEE Symposium on Foundations of Computer Science, 2014.

**Funding:** Gregory’s work was supported by NSF CAREER Award CCF-1351108, and Paul’s work was supported by a Sloan fellowship. Part of this work was done while the authors were visiting the Simons Institute for the Theory of Computing, at UC Berkeley.

†Stanford University, Stanford, CA ([gvaliant@cs.stanford.edu](mailto:gvaliant@cs.stanford.edu), <http://theory.stanford.edu/~valiant/>).

‡Brown University, Providence, RI ([pvaliant@gmail.com](mailto:pvaliant@gmail.com), <http://cs.brown.edu/~pvaliant/>).

34 the problem of “identity testing against a known distribution”. This problem has  
 35 been well studied, and yielded the punchline that it is sometimes possible to perform  
 36 this task using far fewer samples than would be necessary to accurately learn the  
 37 distribution from which the samples were drawn. Nevertheless, previous work on this  
 38 problem either considered only the problem of verifying a uniform distribution (the  
 39 case that  $p = \text{Unif}[n]$ ), or was from the perspective of worst-case analysis—aiming to  
 40 bound the number of samples required to verify a worst-case distribution of a given  
 41 support size.

42 Here, we seek a deeper understanding of this problem. We resolve, up to con-  
 43 stant factors, the sample complexity of this task on an *instance-by-instance* basis—  
 44 determining the optimal number of samples required to verify the identity of a distri-  
 45 bution, *as a function of the distribution in question*.

46 Throughout much of theoretical computer science, the main challenge and goal  
 47 is to characterize problems from a worst-case standpoint, and the efforts to describe  
 48 algorithms that perform well “in practice” are often relegated to the sphere of heuris-  
 49 tics. Still, there is a developing understanding of domains and approaches for which  
 50 one may provide analysis beyond the worst-case (e.g., random instances, smoothed  
 51 analysis, competitive analysis, analysis with respect to various parameterizations of  
 52 the problems, etc.). Against this backdrop, it seems especially exciting when a rich  
 53 setting seems amenable to the development and analysis of *instance optimal* algo-  
 54 rithms, not to mention that instance optimality gives a strong recommendation for  
 55 the practical viability of the proposed algorithms.

56 In the setting of this paper, having the distribution  $p$  explicitly provided to the  
 57 tester enables our approach; nevertheless, it is tantalizing to ask whether this style  
 58 of “instance-by-instance optimal” property testing/estimation or learning is possible  
 59 in more general distributional settings. The authors are optimistic that such strong  
 60 theoretical results are both within our reach, and that pursuing this line may yield  
 61 practical algorithms suited to making the best use of available data. We refer the  
 62 reader to [22] for an example of subsequent work in this direction.

63 To more cleanly present our results, we introduce the following notation.

64 **DEFINITION 1.** *For a probability distribution  $p$  over a discrete support, let  $p^{-\max}$*   
 65 *denote the vector of probabilities obtained from  $p$  by removing the entry corresponding*  
 66 *to the element of largest probability (with ties broken arbitrarily if there are multiple*  
 67 *such elements). For  $\epsilon > 0$ , define  $p_{-\epsilon}$  to be the vector obtained from  $p$  by removing*  
 68 *the domain elements of smallest probability mass under  $p$ , and stopping just before*  
 69 *more than  $\epsilon$  probability mass is removed.*

70 Hence  $p_{-\epsilon}^{-\max}$  is the vector of probabilities corresponding to distribution  $p$ , af-  
 71 ter the largest probability element and the smallest probability elements have been  
 72 removed.

73 Throughout, we use the standard notation for the  $L_p$  norm of a vector: given a  
 74 vector  $x$ , and a real number  $\alpha$  we define the  $\alpha$  norm of  $x$  as

$$75 \quad \|x\|_\alpha = \left( \sum_i x_i^\alpha \right)^{1/\alpha}$$

76 Our main result is the following:

77 **THEOREM 2.** *There exist constants  $c_1, c_2$  such that for any  $\epsilon > 0$  and any known*  
 78 *distribution  $p$ , for any unknown distribution  $q$ , our tester will distinguish  $q = p$  from*

79  $\|p - q\|_1 \geq \epsilon$  with probability  $2/3$  when run on a set of at least  $c_1 \cdot \max \left\{ \frac{1}{\epsilon}, \frac{\|p_{-\epsilon/16}^- \|^{\max}_{2/3}}{\epsilon^2} \right\}$   
 80 samples drawn from  $q$ , and no tester can do this task with probability at least  $2/3$  with  
 81 a set of fewer than  $c_2 \cdot \max \left\{ \frac{1}{\epsilon}, \frac{\|p_{-2\epsilon}^- \|^{\max}_{2/3}}{\epsilon^2} \right\}$  samples.

82 In short, over the entire range of potential distributions  $p$ , our tester is optimal,  
 83 up to constant factors in  $\epsilon$  and the number of samples. The distinction of “con-  
 84 stant factors in  $\epsilon$ ” is needed, as  $\|p_{-\epsilon/16}\|_{2/3}$  might *not* be within a constant factor  
 85 of  $\|p_{-2\epsilon}\|_{2/3}$  if, for example, the vast majority of the  $2/3$ -norm of  $p$  comes from tiny  
 86 domain elements that only comprise an  $\epsilon$  fraction of the 1-norm (and hence would be  
 87 absent from  $p_{-2\epsilon}$ , though not from  $p_{-\epsilon/16}$ ).<sup>1</sup>

88 Because our tester is constant-factor tight, the subscript and superscript on  $p$   
 89 and the max with  $\frac{1}{\epsilon}$  in the sample complexity  $\max \left\{ \frac{1}{\epsilon}, \frac{\|p_{-O(\epsilon)}^- \|^{\max}_{2/3}}{\epsilon^2} \right\}$  all mark real  
 90 phenomena, and are not just artifacts of the analysis. However, except for rather  
 91 pathological distributions, the theorem says that  $\Theta(\frac{\|p\|_{2/3}}{\epsilon^2})$  is the optimal number of  
 92 samples. Additionally, note that the subscript and superscript only reduce the value of  
 93 the norm:  $\|p_{-2\epsilon}^- \|^{\max}_{2/3} < \|p_{-2\epsilon}\|_{2/3} \leq \|p_{-\epsilon/16}\|_{2/3} \leq \|p\|_{2/3}$ , and hence  $O(\|p\|_{2/3}/\epsilon^2)$   
 94 is always an upper bound on the number of samples required. Since  $x^{2/3}$  is concave, for  
 95 distributions  $p$  of support size at most  $n$  the  $L_{2/3}$  norm is maximized on the uniform  
 96 distribution, yielding that  $\|p\|_{2/3} \leq \sqrt{n}$ , with equality if and only if  $p$  is the uniform  
 97 distribution. This immediately yields a worst-case bound of  $O(\sqrt{n}/\epsilon^2)$  on the number  
 98 of samples required to test distributions supported on at most  $n$  elements, tightening  
 99 the previous bound of  $O(\frac{\sqrt{n} \text{ polylog } n}{\epsilon^4})$  from [6], and matching the tight bound on the  
 100 number of samples required for testing the uniform distribution given in [17].

101 The core of our testing algorithm is an extremely simple statistic that is similar to  
 102 Pearson’s chi-squared statistic. Given a set of  $k$  samples, with  $X_i$  denoting the number  
 103 of occurrences of the  $i$ th domain element, and  $p_i$  denoting the probability of drawing  
 104 the  $i$ th domain element from distribution  $p$ , the Pearson chi-squared statistic is given  
 105 as  $\sum_i \frac{(X_i - kp_i)^2 - kp_i}{p_i}$ . Our testing algorithm is, essentially, obtained by modifying this  
 106 statistic in two crucial ways: replacing the second occurrence of  $kp_i$  with  $X_i$  (which  
 107 has expectation  $kp_i$  when drawing samples from  $p$ ), and changing the scaling factor  
 108 from  $1/p_i$  to  $1/p_i^{2/3}$ :

$$109 \quad \sum_i \frac{(X_i - kp_i)^2 - X_i}{p_i^{2/3}}.$$

110 Our simple testing algorithm is stated below:

<sup>1</sup>In the language of the abstract, Theorem 2 defines a function  $f(p, \epsilon)$  characterizing the sample complexity of testing the identity of  $p$ , tight up to a factor of 32 in the error  $\epsilon$  and some constant  $c_1/c_2$  in the number of samples. Interestingly, since the function  $f(p, \epsilon)$  grows at least inversely in  $\epsilon$  as  $\epsilon$  goes to 0, we can merge the two constants into a single multiplicative constant in the error  $\epsilon$  and say that the right number of samples for testing the identity of  $p$  to within  $\epsilon$  must lie between  $f(p, 32 \frac{c_1}{c_2} \cdot \epsilon)$  and  $f(p, \epsilon)$ . This is a cleaner result, in some sense; however, of the two parameters—the accuracy  $\epsilon$  and the sample size  $k$ —it is often perhaps more important to have precise control of the accuracy, so we wanted to emphasize that while our results are constant-factor-tight, the constant, 32, in front of  $\epsilon$  is explicit, and can be made small.

## AN INSTANCE-OPTIMAL TESTER

Given a parameter  $\epsilon > 0$  and a set of  $k$  samples drawn from  $q$ , let  $X_i$  represent the number of times the  $i$ th domain element occurs in the samples. Assume wlog that the domain elements of  $p$  are sorted in non-increasing order of probability. Define  $s = \min\{i : \sum_{j>i} p_j \leq \epsilon/8\}$ , and let  $M = \{2, \dots, s\}$ , and  $S = \{s+1, s+2, \dots\}$ . (Note that  $p_M = p_{-\epsilon/8}^-$ .)

1. If  $\sum_{i \in M} \frac{(X_i - kp_i)^2 - X_i}{p_i^{2/3}} > 4k \|p_M\|_{2/3}^{1/3}$ , or
2. If  $\sum_{i \in S} X_i > \frac{3}{16} \epsilon k$ , then output “ $\|p - q\|_1 \geq \epsilon$ ”, else output “ $p = q$ ”.

While the algorithm we propose is extremely simple, the analysis involves sorting through several messy inequalities. To facilitate this analysis, we give a complete characterization of a general class of inequalities. We characterize the set of sequences of triples  $(a, b, c)_i = (a_1, b_1, c_1), \dots, (a_r, b_r, c_r)$  for which it holds that for all finite sequences of positive numbers  $(x)_j = x_1, \dots$  and  $(y)_j = y_1, \dots$ ,

$$(1) \quad \prod_{i=1}^r \left( \sum_j x_j^{a_i} y_j^{b_i} \right)^{c_i} \geq 1.$$

This is an extremely frequently encountered class of inequalities, and contains the Cauchy-Schwarz inequality and its generalization, the Hölder inequality, in addition to inequalities representing the monotonicity of the  $L_p$  norm, and also clearly contains any finite product of such inequalities. Additionally, we note that the constant 1 on the right hand side cannot be made larger, for all such inequalities are false when the sequences  $x$  and  $y$  consist of a single 1; also, as we show, the class of valid inequalities is unchanged if 1 is replaced by any other constant in the interval  $(0, 1]$ .

**EXAMPLE 1.** *The classic Cauchy-Schwarz inequality can be expressed in the form of Equation 1 as  $\left(\sum_j X_j\right)^{1/2} \left(\sum_j Y_j\right)^{1/2} \left(\sum_j \sqrt{X_j Y_j}\right)^{-1} \geq 1$ , corresponding to the triples  $(a, b, c)_i = (1, 0, \frac{1}{2}), (0, 1, \frac{1}{2}), (\frac{1}{2}, \frac{1}{2}, -1)$ . This inequality is tight when the sequences  $X$  and  $Y$  are proportional to each other. The Hölder inequality generalizes Cauchy-Schwarz by replacing  $\frac{1}{2}$  by  $\lambda \in [0, 1]$ , yielding the inequality defined by the triples  $(a, b, c)_i = (1, 0, \lambda), (0, 1, 1 - \lambda), (\lambda, 1 - \lambda, -1)$ .*

**EXAMPLE 2.** *A fundamentally different inequality that can also be expressed in the form of Equation 1 is the fact that the  $L_p$  norm is a non-increasing function of  $p$ . For  $p \in [0, 1]$  we have the inequality  $\left(\sum_j X_j^p\right) \left(\sum_j X_j\right)^{-p} \geq 1$ , corresponding to the two triples  $(a, b, c)_i = (p, 0, 1), (1, 0, -p)$ . This inequality is tight only when the sequence  $(X)_j$  consists of a single nonzero term.*

We show that the cases where Equation 1 holds are exactly those cases expressible as a product of inequalities of the above two forms, where two arbitrary combinations of  $x$  and  $y$  are substituted for the sequence  $X$  and the sequence  $Y$  in the above examples:

**THEOREM 3.** *For a fixed sequence of triples  $(a, b, c)_i = (a_1, b_1, c_1), \dots, (a_r, b_r, c_r)$ , the inequality  $\prod_{i=1}^r \left(\sum_j x_j^{a_i} y_j^{b_i}\right)^{c_i} \geq 1$  holds for all finite sequences of positive numbers  $(x)_j, (y)_j$  if and only if it can be expressed as a finite product of positive powers of*

Hölder inequalities of the form

$$\left(\sum_j x_j^{a'} y_j^{b'}\right)^\lambda \cdot \left(\sum_j x_j^{a''} y_j^{b''}\right)^{1-\lambda} \geq \sum_j x_j^{\lambda a' + (1-\lambda)a''} y_j^{\lambda b' + (1-\lambda)b''},$$

140 and  $L_p$  monotonicity inequalities of the form  $\left(\sum_j x_j^a y_j^b\right)^\lambda \leq \sum_j x_j^{\lambda a} y_j^{\lambda b}$ , where  $\lambda \in$   
 141  $[0, 1]$ .

142 We state this theorem for pairs of sequences  $(x)_j, (y)_j$ , of positive numbers, al-  
 143 though an analogous statement (Theorem 4 stated in Section 2) holds for any number  
 144 of positive sequences and is yielded by a trivial extension of the proof of the above  
 145 theorem. Most commonly encountered instances of inequalities of the above form,  
 146 including those involved in our identity testing result, involve only pairs of sequences.  
 147 Further, the result is nontrivial even for inequalities of the above form that only in-  
 148 involve a single sequence—see Example 3 for a discussion of a single sequence inequality  
 149 with surprising properties.

150 Our proof of Theorem 3 is algorithmic in nature; in fact, we describe an algorithm  
 151 which, when given the sequence of triples  $(a, b, c)_i$  as input, will run in polynomial  
 152 time, and either output a derivation of the desired inequality as a product of a polyno-  
 153 mial number of Hölder and  $L_p$  monotonicity inequalities, or the algorithm will output  
 154 a witness from which a pair of sequences  $(x)_j, (y)_j$  that violate the inequality can be  
 155 constructed. It is worth stressing that the algorithm is efficient despite the fact that  
 156 the shortest counter-example sequences  $(x)_j, (y)_j$  might require a doubly-exponential  
 157 number of terms (doubly-exponential in the number of bits required to represent the  
 158 sequence of triples  $(a, b, c)_i$ —see Example 3).

159 The characterization of Theorem 3 seems to be a useful and general tool, and  
 160 seems absent from the literature, perhaps because linear programming duality is an  
 161 unexpected tool with which to analyze such inequalities. The ability to efficiently  
 162 verify inequalities of the above form greatly simplified the tasks of proving our instance  
 163 optimality results; we believe this tool will prove useful to others and have made a  
 164 Matlab implementation of our inequality prover/refuter publicly available at <http://theory.stanford.edu/~valiant/code.html>.  
 165

166 **1.1. Related work.** The general area of hypothesis testing was launched by  
 167 Pearson in 1900, with the description of Pearson’s chi-squared test. In this cur-  
 168 rent setting of determining whether a set of  $k$  samples was drawn from distribution  
 169  $p = p_1, p_2, \dots$ , that test would correspond to evaluating  $\sum_i \frac{1}{p_i} (X_i - kp_i)^2$ , where  $X_i$   
 170 denotes the number of occurrences of the  $i$ th domain element in the samples, and  
 171 then outputting “yes” if the value of this statistic is sufficiently small. Traditionally,  
 172 such tests are evaluated in the asymptotic regime, for a fixed distribution  $p$  as the  
 173 number of samples tends to infinity. In the current setting of trying to verify the  
 174 identity of a distribution, using this chi-squared statistic might require using many  
 175 more samples than would be necessary even to accurately *learn* the distribution from  
 176 which the samples were drawn (see, e.g., Example 6).

177 Over the past fifteen years, there has been a body of work exploring the general  
 178 question of how to estimate or test properties of distributions *using fewer samples*  
 179 *than would be necessary to learn the distribution in question*. Such properties include  
 180 “symmetric” properties (properties whose value is invariant to relabeling domain ele-  
 181 ments) such as entropy, support size, and distance metrics between distributions (such  
 182 as  $L_1$  distance), with work on both the algorithmic side (e.g., [7, 5, 12, 15, 16, 4, 9]),

183 and on establishing lower bounds [18, 23]. Such problems have been almost exclu-  
 184 sively considered from a worst-case standpoint, with bounds on the sample complexity  
 185 parameterized by an upper bound on the support size of the distribution. The recent  
 186 work [20, 21] resolved the worst-case sample complexities of estimating many of these  
 187 symmetric properties. Also see [19] for a recent survey.

188 The specific question of verifying the identity of a distribution was one of the  
 189 first questions considered in this line of work. Motivated by a connection to testing  
 190 the expansion of graphs, Goldreich and Ron [11] first considered the problem of dis-  
 191 tinguishing whether a set of samples was drawn from the uniform distribution of  
 192 support  $n$  versus from a distribution that is least  $\epsilon$  far from the uniform distribu-  
 193 tion, with the tight bound of  $\Theta(\frac{\sqrt{n}}{\epsilon^2})$  on the number of samples subsequently given by  
 194 Paninski [17]. For the more general problem of verifying the identity of an arbitrary  
 195 distribution, Batu et al. [6], showed that for worst-case distributions of support size  
 196  $n$ ,  $O(\frac{\sqrt{n} \text{polylog } n}{\epsilon^4})$  samples are sufficient. Since the publication of this current paper,  
 197 Diakonikolis et al. [10], considered the problem of identity testing under various as-  
 198 sumptions about the *shape* of the distribution, including, for example, assuming the  
 199 distribution is monotone, unimodal, multimodal, or piecewise constant, etc., relative  
 200 to an ordering of the domain elements; for distributions assumed to be piecewise con-  
 201 stant with  $t$  pieces, they show a tester with  $O(\frac{\sqrt{t}}{\epsilon^2})$  samples, which, letting  $t = n$  yields  
 202 a  $O(\frac{\sqrt{n}}{\epsilon^2})$ -sample tester in our setting, which has worst-case optimal dependence on  $n$   
 203 and  $\epsilon$  (but is not instance-optimal).

204 In a similar spirit to this current paper, motivated by a desire to go beyond worst-  
 205 case analysis, Acharya et al. [1, 2] recently considered the question of identity testing  
 206 with two unknown distributions (i.e., both distributions  $p$  and  $q$  are unknown, and one  
 207 wishes to deduce if  $p = q$  from samples) from the standpoint of *competitive analysis*.  
 208 They asked how many samples are required as a function of the number of samples  
 209 that would be required for the task of distinguishing whether samples were drawn  
 210 from  $p$  versus  $q$  in the case where  $p$  and  $q$  were known to the algorithm. Their main  
 211 results are an algorithm that performs the desired task using  $m^{3/2}$  polylog  $m$  samples,  
 212 and a lower bound of  $\Omega(m^{7/6})$ , where  $m$  represents the number of samples required to  
 213 determine whether a set of samples were drawn from  $p$  versus  $q$  in the setting where  
 214  $p$  and  $q$  are explicitly known. One of the main conceptual messages from Acharya et  
 215 al.’s results is that knowledge of the underlying distributions is extremely helpful—  
 216 without such knowledge one loses a polynomial factor in sample complexity. Our  
 217 results build on this moral, in some sense describing the “right” way that knowledge  
 218 of a distribution can be used to test identity.

219 The form of our tester may be seen as rather similar to those in [1, 2, 8], which  
 220 considered testing whether two distributions were close or not when *both* distributions  
 221 are unknown. The testers in those papers and the tester proposed here consist es-  
 222 sentially of summing up carefully chosen expressions independently evaluated at the  
 223 different domain elements and comparing this sum to a threshold. These testers are  
 224 considerable simpler than many of the proposed testers in other works (including [10]  
 225 and the initial pioneering work [6]), which proceed by subdividing the domain into a  
 226 super-constant number of partitions, and applying tests to each partition separately.  
 227 From a technical perspective, our lower bounds leverage Hellinger distance to intro-  
 228 duce a flexible class of lower bound instances, which yield the tight results of this  
 229 work, and were also employed to give the lower bounds in [8].

230 **1.2. Organization.** We begin with our characterization of the class of inequal-  
 231 ities, as we feel that this tool may be useful to the broader community; this first  
 232 section is entirely self-contained. Section 3.1 contains the definitions and terminology  
 233 relevant to the distribution testing portion of the paper, and Section 3.2 describes  
 234 our very simple instance-optimal distribution identity testing algorithm, and provides  
 235 some context and motivation for the algorithm. Section 4 discusses the lower bounds,  
 236 establishing the optimality of our tester.

237 **2. A class of inequalities generalizing Hölder’s inequality and the mono-**  
 238 **tonicity of  $L_p$  norms.** In this section we characterize under what conditions a large  
 239 class of inequalities holds, showing both how to derive these inequalities when they  
 240 are true and how to refute them when they are false. We encounter such inequalities  
 241 repeatedly in the analysis of our tester in Section 3.

242 The basic question we resolve is: for what sequences of triples  $(a, b, c)_i$  is it true  
 243 that for all sequences of positive numbers  $(x)_j, (y)_j$  we have

$$244 \quad (2) \quad \prod_i \left( \sum_j x_j^{a_i} y_j^{b_i} \right)^{c_i} \geq 1$$

245 We note that the constant 1 on the right hand side cannot be made larger, for all  
 246 such inequalities are false when the sequences  $x$  and  $y$  consist of a single 1; also, as we  
 247 will show later, if this inequality can be violated, it can be violated by an arbitrary  
 248 amount, so if any right hand side constant works, for a given  $(a, b, c)_i$ , then 1 works,  
 249 as stated above.

250 Such inequalities are typically proven by hand, via trial and error. One basic tool  
 251 for this is the Cauchy-Schwarz inequality,  $\left(\sum_j X_j\right)^{1/2} \left(\sum_j Y_j\right)^{1/2} \geq \sum_j \sqrt{X_j Y_j}$ , or  
 252 the slightly more general Hölder inequality, a weighted version of Cauchy-Schwarz,  
 253 where for  $\lambda \in [0, 1]$  we have  $\left(\sum_j X_j\right)^\lambda \left(\sum_j Y_j\right)^{1-\lambda} \geq \sum_j X_j^\lambda Y_j^{1-\lambda}$ . Writing this in  
 254 the form of Equation 2, and substituting arbitrary combinations of  $x$  and  $y$  for  $X$  and  
 255  $Y$  yields families of inequalities of the form:

$$256 \quad \left(\sum_j x_j^{a_1} y_j^{b_1}\right)^\lambda \left(\sum_j x_j^{a_2} y_j^{b_2}\right)^{1-\lambda} \left(\sum_j x_j^{\lambda a_1 + (1-\lambda)a_2} y_j^{\lambda b_1 + (1-\lambda)b_2}\right)^{-1} \geq 1,$$

257 and we can multiply (positive powers of) inequalities of this form together to get  
 258 further cases of the inequality in Equation 2. This inequality is tight when the two  
 259 sequences  $X$  and  $Y$  are proportional to each other.

260 A second and different basic inequality of our general form, for  $\lambda \in [0, 1]$ , is:

261  $\left(\sum_j X_j\right)^\lambda \leq \sum_j X_j^\lambda$ , which is the fact that the  $L_p$  norm is a decreasing function of  $p$ .  
 262 (Intuitively, this is a slight generalization of the trivial fact that  $x^2 + y^2 \leq (x+y)^2$ , and  
 263 follows from the fact that the derivative of  $x^\lambda$  is a decreasing function of  $x$ , for positive  
 264  $x$ ). As above, products of powers of  $x$  and  $y$  may be substituted for  $X$  to yield a more  
 265 general class of inequalities:  $\sum_j x_j^\lambda y_j^{1-\lambda} \left(\sum_j x_j^a y_j^b\right)^{-\lambda} \geq 1$ , for  $\lambda \in [0, 1]$ . Unlike the  
 266 previous case, these inequalities are tight when there is only a single nonzero value of  
 267  $X$ , and the inequality may seem weak for nontrivial cases.

268 The main result of this section is that the cases where Equation 2 holds are  
 269 *exactly* those cases expressible as a product of inequalities of the above two forms,

270 and that such a representation can be efficiently found. While we have been discussing  
 271 inequalities involving two sequences, these results apply to inequalities on  $d$  sequences,  
 272 for any positive integer  $d$ . For completeness, we restate Theorem 3 in this more general  
 273 form. The proof of this more general theorem is similar to that of its two-sequence  
 274 analog, Theorem 3.

275 **THEOREM 4.** *For  $d + 1$  fixed sequences  $(a)_{1,i} = a_{1,1} \dots, a_{1,r}, \dots, (a)_{d,i} =$   
 276  $a_{d,1}, \dots, a_{d,r}$ , and  $(c)_i = c_1, \dots, c_r$ , the inequality  $\prod_{i=1}^r \left( \sum_j \left( \prod_{k=1}^d x_{k,j}^{a_{k,i}} \right)^{c_i} \right) \geq 1$   
 277 holds for all sets of  $d$  finite sequences of positive numbers  $(x)_{k,j}$  if and only if it  
 278 can be expressed as a finite product of positive powers of Hölder inequalities of the  
 279 form  $\left( \sum_j \left( \prod_{k=1}^d x_{k,j}^{a'_k} \right) \right)^\lambda \left( \sum_j \left( \prod_{k=1}^d x_{k,j}^{a''_k} \right) \right)^{1-\lambda} \geq \sum_j \left( \prod_{k=1}^d x_{k,j}^{\lambda a'_k + (1-\lambda) a''_k} \right)$ , and  
 280  $L_p$  monotonicity inequalities of the form  $\left( \sum_j \left( \prod_{k=1}^d x_{k,j}^{a'_k} \right) \right)^\lambda \leq \sum_j \left( \prod_{k=1}^d x_{k,j}^{\lambda a'_k} \right)$ ,  
 281 where  $\lambda \in [0, 1]$ , and where  $a'_k, a''_k$  can be any real numbers.*

282  
 283 *Further, there exists an algorithm which, given  $d + 1$  sequences  $(a)_{1,i} = a_{1,1} \dots, a_{1,r},$   
 284  $\dots, (a)_{d,i} = a_{d,1}, \dots, a_{d,r}$ , and  $(c)_i = c_1, \dots, c_r$  describing the inequality, runs in time  
 285 polynomial in the input description, and either outputs a representation of the desired  
 286 inequality as a product of a polynomial number of positive powers of Hölder and  $L_p$   
 287 monotonicity inequalities, or yields a witness describing  $d$  finite sequences of positive  
 288 numbers  $(x)_{k,j}$  that violate the inequality.*

289 The second portion of the theorem—the existence of an efficient algorithm that  
 290 provides a derivation or refutation of the inequality—is surprising. As the following  
 291 example demonstrates, it is possible that the shortest sequences  $x, y$  that violate the  
 292 inequality have a number of terms that is *doubly exponential* in the description length  
 293 of the sequence of triples  $(a, b, c)_i$  (and exponential in the inverse of the accuracy of the  
 294 sequences). Hence, in the case that the inequality does not hold, our algorithm cannot  
 295 be expected to return a pair of counter-example sequences. Nevertheless, we show that  
 296 it efficiently returns a witness describing such a construction. We observe that the  
 297 existence of this example precludes any efficient algorithm that tries to approach this  
 298 problem by solving some linear or convex program in which the variables correspond  
 299 to the elements of the sequences  $x, y$ .

300 **EXAMPLE 3.** *Consider for some  $\epsilon \geq 0$  the single-sequence inequality*

$$301 \left( \sum_j x_j^{-2} \right)^{-1} \left( \sum_j x_j^{-1} \right)^3 \left( \sum_j x_j^0 \right)^{-2-\epsilon} \left( \sum_j x_j^1 \right)^3 \left( \sum_j x_j^2 \right)^{-1} \geq 1,$$

302 *which can be expressed in the form of Equation 1 via the triples  $(a, b, c)_i = (-2, 0, -1),$   
 303  $(-1, 0, 3), (0, 0, -2 - \epsilon), (1, 0, 3), (2, 0, -1)$ . This inequality is true for  $\epsilon = 0$  but false  
 304 for any positive  $\epsilon$ . However, the shortest counterexample sequences have length that  
 305 grows as  $\exp(\frac{1}{\epsilon})$  as  $\epsilon$  approaches 0. Counterexamples are thus hard to write down,  
 306 though possibly easy to express—for example, letting  $n = 64^{1/\epsilon}$ , the sequence  $x$  of  
 307 length  $2 + n$  consisting of  $n, \frac{1}{n}$ , followed by  $n$  ones violates the inequality.<sup>2</sup>*

308 In the following section we give an overview of the linear programming based  
 309 proof of Theorem 3, and then give the formal proof in Section 2.2. In Section 2.3 we

<sup>2</sup>Showing that counterexample sequences must be essentially this long requires technical machinery from the proof of Theorem 3, however one can glean intuition by evaluating the inequality on the given sequence— $n, \frac{1}{n}$ , followed by  $n$  ones.



310 provide an intuitive interpretation of the computation being performed by the linear  
311 program.

312 **2.1. Proof overview of Theorem 3.** Our proof is based on constructing and  
313 analyzing a certain linear program, whose variables, which we denote by  $\ell_i$ , represent  
314  $\log \sum_j x_j^{a_i} y_j^{b_i}$  for each  $i$  in the index set of triples  $(a, b, c)_i$ . Letting  $r$  denote the size  
315 of this index set, the linear program will have  $r$  variables, and  $\text{poly}(r)$  constraints.  
316 We will show that if the linear program does *not* have objective value zero then we  
317 can construct a counterexample pair of sequences  $(x)_j, (y)_j$  for which the inequality is  
318 contradicted. Otherwise, if the objective value is zero, then we will consider a solution  
319 to the *dual* of this linear program, and interpret this solution as an explicit (finite)  
320 combination of Hölder and  $L_p$  monotonicity inequalities whose product yields the  
321 desired inequality in question. Combined, these results imply that we can efficiently  
322 either derive or refute the inequality in all cases.

323 Given (finite) sequences  $(x)_j, (y)_j$ , consider the function  $\ell : \mathbb{R}^2 \rightarrow \mathbb{R}$  defined as  
324  $\ell(a, b) = \log \sum_j x_j^a y_j^b$ . We will call this 2-dimensional function  $\ell(a, b)$  the *norm graph*  
325 of the sequences  $(x)_j, (y)_j$ , and will analyze this function for the remainder of this  
326 proof and show how to capture many of its properties via linear programming. The  
327 inequality in question,  $\prod_i \left( \sum_j x_j^{a_i} y_j^{b_i} \right)^{c_i} \geq 1$ , is equivalent (taking logarithms) to  
328 the claim that  $\sum_i c_i \cdot \ell(a_i, b_i) \geq 0$  for every norm graph  $\ell$  that can be realized via  
329 sequences  $(x)_j, (y)_j$ .

330 The Hölder inequalities explicitly represent the fact that norm graphs  $\ell$  must be  
331 convex, namely for each  $\lambda \in (0, 1)$  and each pair  $(a', b'), (a'', b'')$  we have  $\lambda \ell(a', b') +$   
332  $(1 - \lambda) \ell(a'', b'') \geq \ell(\lambda a' + (1 - \lambda) a'', \lambda b' + (1 - \lambda) b'')$ . The  $L_p$  monotonicity inequalities  
333 can correspondingly be expressed in terms of norm graphs  $\ell$ , intuitively as “any secant  
334 of the graph of  $\ell$  (interpreted as a line in 3 dimensions) that intersects the  $z$ -axis must  
335 intersect it at a nonnegative  $z$ -coordinate,” explicitly, for all  $(a', b')$  and all  $\lambda \in (0, 1)$   
336 we have  $\lambda \ell(a', b') \leq \ell(\lambda a', \lambda b')$ .

337 Instead of modeling the class of norm graphs directly, we instead model the class  
338 of functions that are convex and satisfy the secant property, which we call “linearized  
339 norm graphs”: let  $\mathcal{L}$  represent this family of functions from  $\mathbb{R}^2$  to  $\mathbb{R}$ , namely, those  
340 functions that are convex and whose secants through the  $z$ -axis pass through-or-above  
341 the origin. As we will show, this class  $\mathcal{L}$  essentially captures the class of functions  
342  $\ell : \mathbb{R}^2 \rightarrow \mathbb{R}$  that can be realised as  $\ell(a, b) = \log \sum_j x_j^a y_j^b$  for some sequences  $(x)_j, (y)_j$ ,  
343 provided we only care about the values of  $\ell$  at a finite number of points  $(a_i, b_i)$ , and  
344 provided we only care about the  $r$ -tuple  $\ell(a_i, b_i)$  up to scaling by positive numbers.  
345 In other words, the inequality  $\sum_i c_i \cdot \ell(a_i, b_i) \geq 0$  holds for all norm graphs if and only  
346 if it holds for all linearized norm graphs, showing that products of positive powers of  
347 Hölder and  $L_p$  monotonicity inequalities (used to define the class of linearized norm  
348 graphs) exactly capture all norm graph inequalities. In this manner we can reduce  
349 the very complicated combinatorial phenomena surrounding Equation 2 to a linear  
350 program.

351 The proof can be decomposed into four steps:

352 **1)** We construct a homogeneous linear program (“homogeneous” means the con-  
353 straints have no additive constants) which we will analyze in the rest of the proof. The  
354 linear program has  $r$  variables  $(\ell)_i$ , where feasible points will represent valid  $r$ -tuples  
355  $\ell(a_i, b_i)$  for linearized norm graphs  $\ell \in \mathcal{L}$ . As will become important later, we set  
356 the objective function to minimize the expression corresponding to the logarithm of  
357 the desired inequality:  $\min \sum_i c_i \cdot \ell_i$ . Also, as will become important later, we will

358 construct each of the constraints of the linear program so that they are positive linear  
 359 combinations of logarithms of Hölder and  $L_p$  monotonicity inequalities when the  $(\ell)_i$   
 360 are interpreted as the values of a norm graph at the points  $(a_i, b_i)$ .

361 **2)** We show that for each feasible point, an  $r$ -tuple  $(\ell)_i$ , there is a *linearized* norm  
 362 graph  $\ell : \mathbb{R}^2 \rightarrow \mathbb{R}$  that extends  $\ell_i = \ell(a_i, b_i)$  to the whole plane, where, further, the  
 363 function  $\ell$  is the maximum of a finite number of affine functions (functions of the form  
 364  $\alpha a + \beta b + \gamma$ ).

365 **3)** For any desired accuracy  $\epsilon > 0$ , we show that for sufficiently small  $\delta > 0$  there is a  
 366 (regular, not linearized) norm graph  $\ell'$  such that for any  $(a, b) \in \mathbb{R}^2$  the scaled version  
 367  $\delta \cdot \ell'(a, b)$  approximates the linearized norm graph constructed in the previous part,  
 368  $\ell(a, b)$ , to within error  $\epsilon$ .

369 Namely, any feasible point of our linear program corresponds to a (possibly scaled)  
 370 norm graph. Thus, if there exists a feasible point for which the objective function is  
 371 negative,  $\sum_i c_i \cdot \ell_i < 0$ , then we can construct sequences  $(x)_j, (y)_j$  and a corresponding  
 372 norm graph  $\ell'(a, b) = \log \sum_j x_j^a y_j^b$  for which (because  $\ell'$  can be made to approximate  
 373  $\ell$  arbitrarily well at the points  $(a_i, b_i)$ , up to scaling) we have  $\sum_i c_i \cdot \ell'(a_i, b_i) < 0$ ,  
 374 meaning that the sequences  $(x)_j, (y)_j$  violate the desired inequality. Thus we have  
 375 constructed the desired counterexample

376 **4)** In the other case, where the minimum objective function of the linear program  
 377 is nonnegative, we note that because by construction we have a homogeneous linear  
 378 program (each constraint has a right hand side of 0), the optimal objective value must  
 379 be 0. The solution to the *dual* of our linear program gives a proof of optimality, in  
 380 a particularly convenient form: the dual solution describes a nonnegative linear com-  
 381 bination of the constraints that shows the objective function is always nonnegative,  
 382  $\sum_i c_i \cdot \ell_i \geq 0$ . Recall that, by construction, if each  $\ell_i$  is interpreted as the value of a  
 383 norm graph at point  $(a_i, b_i)$  then each of the linear program constraints is a positive  
 384 linear combination of the logarithms of certain Hölder and  $L_p$  monotonicity inequal-  
 385 ities expressed via values of the norm graph. Combining these two facts yields that  
 386 the inequality  $\sum_i c_i \cdot \ell(a_i, b_i) \geq 0$  can be derived as a positive linear combination of  
 387 the logarithms of certain Hölder and  $L_p$  monotonicity inequalities. Exponentiating  
 388 yields that the desired inequality can be derived as the product of positive powers of  
 389 certain Hölder and  $L_p$  monotonicity inequalities, as desired.

390 The following section provides the proof details for the above overview.

391 **2.2. Proof of Theorem 3.** Given  $r$  triples,  $(a_1, b_1, c_1), \dots, (a_r, b_r, c_r)$ , consider  
 392 the linear program with  $r$  variables denoted by  $\ell_1, \dots, \ell_r$  with objective function  
 393  $\min \sum_i c_i \cdot \ell_i$ . For each index  $k \in [r]$  we add linear constraints to enforce that the  
 394 point  $(a_k, b_k, \ell_k)$  in  $\mathbb{R}^3$  lies on the lower convex hull of the points  $(a_i, b_i, \ell_i)$  and the  
 395 extra point  $(2a_k, 2b_k, 2\ell_k)$ . Recall that the parameters  $(a_i, b_i)$  are constants, so we  
 396 may use them arbitrarily to set up the linear program. Explicitly, for each triple,  
 397 pair, or singleton from the set  $\{(a_i, b_i) : i \neq k\} \cup \{(2a_k, 2b_k)\}$  that have a unique  
 398 convex combination that equals  $(a_k, b_k)$ , we add a constraint that the corresponding  
 399 combination of their associated  $z$ -values (i.e. the corresponding  $\ell_i$  or  $2\ell_k$ ) must be  
 400 greater than or equal to  $\ell_k$ . The total number of constraints is thus  $O(r^4)$ . We note  
 401 that these are homogeneous constraints—there are no additive constants. Intuitively,  
 402 we are expressing all our constraints on the linearized norm graph in this convex hull  
 403 form: the Hölder inequalities are naturally convexity constraints, and by adding these  
 404 “fictitious” points  $(2a_k, 2b_k, 2\ell_k)$ , the  $L_p$  monotonicity inequalities can now also be  
 405 treated as convexity constraints.

406 We now begin our proof of one direction of Theorem 3—that if the above linear  
 407 program has objective function value 0, then the desired inequality can be expressed  
 408 as the product of a finite number of Hölder and  $L_p$  monotonicity inequalities. As  
 409 a first step, we establish that each of the above constraints can be expressed as a  
 410 positive linear combination of these two types of inequalities:

411 LEMMA 5. *Each of the above-described constraints can be expressed as a positive*  
 412 *linear combination of the logarithms of Hölder and  $L_p$  monotonicity inequalities.*

413 *Proof.* Consider, first, the case when the convex combination does not involve the  
 414 special point  $(2a_k, 2b_k)$ . Thus there are indices  $i_1, i_2, i_3$  and nonnegative constants  
 415  $\lambda_1, \lambda_2, \lambda_3$  with  $\lambda_1 + \lambda_2 + \lambda_3 = 1$  for which  $\lambda_1(a_{i_1}, b_{i_1}) + \lambda_2(a_{i_2}, b_{i_2}) + \lambda_3(a_{i_3}, b_{i_3}) =$   
 416  $(a_k, b_k)$  and we want to conclude a kind of “three-way Hölder inequality”, that  
 417  $\lambda_1 \ell(a_{i_1}, b_{i_1}) + \lambda_2 \ell(a_{i_2}, b_{i_2}) + \lambda_3 \ell(a_{i_3}, b_{i_3}) \geq \ell(a_k, b_k)$ , for any norm graph  $\ell$ . If two  
 418 of the three  $\lambda$ ’s are 0 (without loss of generality  $\lambda_2 = \lambda_3 = 0$ ) then  $\lambda_1 = 1$  and  
 419  $(a_{i_1}, b_{i_1}) = (a_k, b_k)$  making the inequality trivially  $\ell(a_k, b_k) \geq \ell(a_k, b_k)$ . If only one of  
 420 the  $\lambda$ ’s is 0, without loss of generality  $\lambda_3 = 0$  and  $\lambda_1 + \lambda_2 = 1$ , making the desired  
 421 inequality a standard Hölder inequality,

$$422 \quad (3) \quad \lambda_1 \ell(a_{i_1}, b_{i_1}) + (1 - \lambda_1) \ell(a_{i_2}, b_{i_2}) \geq \ell(\lambda_1 a_{i_1} + (1 - \lambda_1) a_{i_2}, \lambda_1 b_{i_1} + (1 - \lambda_1) b_{i_2}).$$

423 In the case that all three  $\lambda$ ’s are nonzero, we derive the result by replacing  $\lambda_1$  with  
 424  $\bar{\lambda}_1 = \frac{\lambda_1}{\lambda_1 + \lambda_2}$  in Equation 3 and multiplying both sides of the inequality by  $\lambda_1 + \lambda_2$ ,  
 425 and then adding the following Hölder inequality:

$$426 \quad (4) \quad (\lambda_1 + \lambda_2) \ell(\bar{\lambda}_1 a_{i_1} + (1 - \bar{\lambda}_1) a_{i_2}, \bar{\lambda}_1 b_{i_1} + (1 - \bar{\lambda}_1) b_{i_2}) + \lambda_3 \ell(a_{i_3}, b_{i_3}) \geq \ell(a_k, b_k).$$

427 Finally, we consider the case where  $(2a_k, 2b_k, 2\ell(a_k, b_k))$  is used; we only con-  
 428 sider the triple case as the other cases are easily dealt with. Thus we have that  
 429 a convex combination with coefficients  $\lambda_1 + \lambda_2 + \lambda_3 = 1$  of the points  $(a_{i_1}, b_{i_1})$ ,  
 430  $(a_{i_2}, b_{i_2})$ ,  $(2a_k, 2b_k)$  equals  $(a_k, b_k)$ . We thus must derive the somewhat odd inequality  
 431  $\lambda_1 \ell(a_{i_1}, b_{i_1}) + \lambda_2 \ell(a_{i_2}, b_{i_2}) + 2\lambda_3 \ell(a_k, b_k) \geq \lambda(a_k, b_k)$ . As above, substitute  $\bar{\lambda}_1 = \frac{\lambda_1}{\lambda_1 + \lambda_2}$   
 432 for  $\lambda_1$  in Equation 3 and multiply by  $\lambda_1 + \lambda_2$ ; this time, add to it  $\lambda_1 + \lambda_2$  times the  
 433  $L_p$  monotonicity inequality

$$434 \quad (5) \quad \frac{1 - 2\lambda_3}{\lambda_1 + \lambda_2} \ell(a_k, b_k) \leq \ell\left(\frac{1 - 2\lambda_3}{\lambda_1 + \lambda_2} a_k, \frac{1 - 2\lambda_3}{\lambda_1 + \lambda_2} b_k\right).$$

435 Everything is seen to match up since the points at which the  $\ell$  functions on the  
 436 right hand sides of Equations 3 and 5 are evaluated are equal (since  $(1 - 2\lambda_3)a_k =$   
 437  $\lambda_1 a_{i_1} + \lambda_2 a_{i_2}$  from the original interpolation).  $\square$

438 Given the above lemma, the proof of one direction of Theorem 3 now follows  
 439 easily—essentially following from step 4 of the proof overview given in the previous  
 440 section.

441 LEMMA 6. *If the objective value of the linear program is non-negative, then it*  
 442 *must be zero, and the inequality  $\prod_i \left(\sum_j x_j^{a_i} y_j^{b_i}\right)^{c_i}$  can be expressed as a product of at*  
 443 *most  $O(r^4)$  Hölder and  $L_p$  monotonicity inequalities.*

444 *Proof.* Recall that since the linear program is homogeneous (each constraint has  
 445 a right hand side of 0), the optimal objective value cannot be larger than 0, and  
 446 hence if the objective value is not negative, it must be 0. The solution to the *dual*

447 of our linear program gives a proof of optimality, in a particularly convenient form:  
 448 the dual solution describes nonnegative coefficients for each of the primal inequality  
 449 constraints, such that when we add up these constraints scaled by these coefficients,  
 450 we find  $\sum_i c_i \cdot \ell_i \geq 0$ —a lower bound on our primal objective function. Recall that,  
 451 by construction, if each  $\ell_i$  is interpreted as the value of a norm graph at point  $(a_i, b_i)$ ,  
 452 then Lemma 5 shows that each of the linear program constraints is a positive linear  
 453 combination of the logarithms of certain Hölder and  $L_p$  monotonicity inequalities  
 454 expressed via values of the norm graph. Combining these two facts yields that the  
 455 inequality  $\sum_i c_i \cdot \ell(a_i, b_i) \geq 0$  can be derived as a positive linear combination of the  
 456 logarithms of certain Hölder and  $L_p$  monotonicity inequalities. Exponentiating yields  
 457 that the desired inequality can be derived as the product of positive powers of Hölder  
 458 and  $L_p$  monotonicity inequalities, as claimed.  $\square$

459 We now flesh out steps 2 and 3 of the proof overview of the previous section to  
 460 establish the second direction of the theorem—namely that if the solution to the linear  
 461 program is negative, we can construct a pair of sequences  $(x)_j, (y)_j$  that violates the  
 462 inequality. We accomplish this in two steps. The first step is to show that for any  
 463 feasible point,  $(\ell)_i$ , of the linear program, one can construct a function  $\ell(a, b) : \mathbb{R}^2 \rightarrow \mathbb{R}$   
 464 defined on the entire plane with the property that the function is convex and has the  
 465 secants through-or-above the origin property, and satisfies  $\ell(a_i, b_i) = \ell_i$ , where  $\ell_i$  is  
 466 the assignment of the linear program variable corresponding to  $a_i, b_i$ .

467 **LEMMA 7.** *For any feasible point  $(\ell)_i$  of the linear program, we can construct*  
 468 *a linearized norm graph  $\ell(a, b) : \mathbb{R}^2 \rightarrow \mathbb{R}$ , which will be the maximum of  $r$  affine*  
 469 *functions  $z_i(a, b) = \alpha_i a + \beta_i b + \gamma_i$  with  $\gamma_i \geq 0$ , such that the function is convex, and*  
 470 *for any  $i \in [r]$ ,  $\ell(a_i, b_i) = \ell_i$ .*

471 *Proof.* We explicitly construct  $\ell$  as the maximum of  $r$  linear functions. Recall  
 472 that for each index  $k$  we constrained  $(a_k, b_k, \ell_k)$  to lie on the lower convex hull of all  
 473 the points  $(a_i, b_i, \ell_i)$  and the special point  $(2a_k, 2b_k, 2\ell_k)$ . Thus through each point  
 474  $(a_k, b_k, \ell_k)$  construct a plane that passes through or below all these other points; define  
 475  $\ell(a, b)$  to be the maximum of these  $r$  functions. For each  $k \in [r]$  we have  $\ell(a_k, b_k) = \ell_k$   
 476 since the  $k$ th plane passes through this value, and every other plane passes through or  
 477 below this value. The maximum of these planes is clearly a convex function. Finally,  
 478 we note that each plane passes through-or-above the origin since a plane that passes  
 479 through  $(a_k, b_k, \ell_k)$  and through-or-below  $(2a_k, 2b_k, 2\ell_k)$  must pass through or above  
 480 the origin; hence for all  $i \in [r]$ ,  $\gamma_i \geq 0$ .  $\square$

481 The second step of the proof consists of showing that we can use the function  
 482  $\ell(a, b)$  of the above lemma to construct sequences  $(x)_j, (y)_j$  that instantiate solutions  
 483 of the linear program arbitrarily well, up to a scaling factor:

**LEMMA 8.** *For a feasible point of the linear program, expressed as an  $r$ -tuple of*  
*values  $(\ell)_i$ , and any  $\epsilon > 0$ , for sufficiently small  $\delta > 0$  there exist finite sequences*  
 *$(x)_j, (y)_j$  such that for all  $i \in [r]$ ,*

$$|\ell_i - \delta \log \sum_j x_j^{a_i} y_j^{b_i}| < \epsilon.$$

484 *Proof.* Consider the linearized norm graph  $\ell(a, b)$  of Lemma 7 that extends  $\ell(a_i, b_i)$   $\blacksquare$   
 485 to the whole plane, constructed as the maximum of  $r$  planes  $z_i(a, b) = \alpha_i a + \beta_i b + \gamma_i$ ,  
 486 with  $\gamma_i \geq 0$ .

Consider, for parameter  $t_i$  to be defined shortly, the sequences  $(x)_j, (y)_j$  consisting

of  $t_i$  copies respectively of  $e^{\alpha_i/\delta}$  and  $e^{\beta_i/\delta}$ . Hence, for all  $a, b$  we have that

$$\delta \log \sum_j x_j^a y_j^b = \alpha_i a + \beta_i b + \delta \log t_i.$$

487 Since  $\gamma_i \geq 0$ , if we let  $t_i = \text{round}(e^{\gamma_i/\delta})$  we can approximate  $\gamma_i$  arbitrarily well  
 488 for small enough  $\delta$ . Finally, we concatenate this construction for all  $i$ . Namely, let  
 489  $(x)_j, (y)_j$  consist of the concatenation, for all  $i$ , of  $t_i = \text{round}(e^{\gamma_i/\delta})$  copies respectively  
 490 of  $e^{\alpha_i/\delta}$  and  $e^{\beta_i/\delta}$ . The values of  $\sum_j x_j^a y_j^b$  will be the sum of the values of these  $r$   
 491 components, thus at least the maximum of these  $r$  components, and at most  $r$  times  
 492 the maximum. Thus the values of  $\delta \log \sum_j x_j^a y_j^b$  will be within  $\delta \log r$  of  $\delta$  times  
 493 the logarithm of the max of these components. Since each of the  $r$  components  
 494 approximates the corresponding affine function  $z_i$  arbitrarily well, for small enough  $\delta$ ,  
 495 the function  $\delta \log \sum_j x_j^a y_j^b$  is thus an  $\epsilon$ -good approximation to the function  $\ell$ , and in  
 496 particular is an  $\epsilon$ -good approximation to  $\ell(a_i, b_i)$  when evaluated at  $(a_i, b_i)$ , for each  
 497  $i$ . □

498 The following lemma completes the proof of Theorem 3:

499 **LEMMA 9.** *Given a feasible point of the linear program that has a negative objec-*  
 500 *tive function value, there exist finite sequences  $(x)_j, (y)_j$  which falsify the inequality*  
 501  $\prod_i \left( \sum_j x_j^{a_i} y_j^{b_i} \right)^{c_i} \geq 1$ .

*Proof.* Letting  $v > 0$  denote the negative of the objective function value corresponding to feasible point  $(\ell)_i$  of the linear program, define  $\epsilon = \frac{v}{\sum_i |c_i|}$ , and let  $\delta_\epsilon$  and sequences  $(x)_j, (y)_j$  be those guaranteed by Lemma 8 to satisfy  $|\ell_i - \delta_\epsilon \log \sum_j x_j^{a_i} y_j^{b_i}| < \epsilon$ , for all  $i \in r$ . Multiplying this expression by  $c_i$  for each  $i$ , summing, and using the triangle inequality yields

$$\left| \sum_i c_i \ell_i - \delta_\epsilon \left( \sum_i c_i \log \sum_j x_j^{a_i} y_j^{b_i} \right) \right| < v,$$

502 and hence  $\sum_i c_i \log \sum_j x_j^{a_i} y_j^{b_i} < 0$ , and the lemma is obtained by exponentiating both  
 503 sides. □

504 **2.3. A geometric interpretation of inequality derivations.** We provide a  
 505 pleasing and intuitive interpretation of the problem being solved by the linear pro-  
 506 gram in the proof of Theorem 3. This interpretation is most easily illustrated via an  
 507 example, and we use one of the inequalities that we encounter in Section 3 in the the  
 508 analysis of our instance-optimal tester.

509 **EXAMPLE 4.** *The 4th component of Lemma 10 (in Section 3.3) consists of show-*  
 510 *ing the inequality*

$$511 \quad (6) \quad \left( \sum_j x_j^2 y_j^{-2/3} \right)^2 \left( \sum_j x_j^2 y_j^{-1/3} \right)^{-1} \left( \sum_j x_j \right)^{-2} \left( \sum_j y_j^{2/3} \right)^{3/2} \geq 1,$$

512 *where in the notation of the lemma, the sequence  $x$  corresponds to  $\Delta$  and the se-*  
 513 *quence  $y$  corresponds to  $p$ . In the notation of Theorem 3, this inequality corresponds*  
 514 *to the sequence of four triples  $(a_i, b_i, c_i) = (2, -\frac{2}{3}, 2), (2, -\frac{1}{3}, -1), (1, 0, -2), (0, \frac{2}{3}, \frac{3}{2})$ .*  
 515 *How does Theorem 3 help us, even without going through the algorithmic machinery*  
 516 *presented in the proof?*

517 Consider the task of proving this inequality via a combination of Hölder and  $L_p$   
 518 monotonicity inequalities as trying to win the following game. At any moment, the  
 519 game board consists of some numbers written on the plane (with the convention that  
 520 every point without a number is interpreted as having a 0), and you win if you can  
 521 remove all the numbers from the board via a combination of moves of the following  
 522 two types:

- 523 1. Any two positive numbers can be moved to their weighted mean. (Namely,  
 524 we can subtract 1 from one location in the plane, subtract 3 from a second  
 525 location in the plane, and add 4 to a point  $\frac{3}{4}$  of the way from the first location  
 526 to the second location.)
- 527 2. Any negative number can be moved towards the origin by a factor  $\lambda \in (0, 1)$   
 528 and scaled by  $\frac{1}{\lambda}$ . (Namely, we can add 1 to one location in the plane, and  
 529 subtract 2 from a location halfway to the origin.)

530 Thus our desired inequality corresponds to the “game board” having a “2” at location  
 531  $(2, -\frac{2}{3})$ , a “-1” at location  $(2, -\frac{1}{3})$ , a “-2” at location  $(1, 0)$ , and a “ $\frac{3}{2}$ ” at location  
 532  $(0, \frac{2}{3})$ . And the rules of the game allow us to push positive numbers together, and push  
 533 negative numbers towards the origin (scaling them). Our visual intuition is quite good  
 534 at solving these types of puzzles. (Try it!)

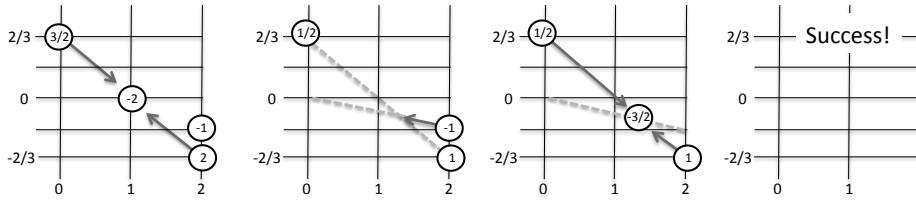


FIG. 1. Depiction of a successful sequence of “moves” in the game corresponding to the inequality  $(\sum_j x_j^2 y_j^{-2/3})^2 (\sum_j x_j^2 y_j^{-1/3})^{-1} (\sum_j x_j)^{-2} (\sum_j y_j^{2/3})^{3/2} \geq 1$ , showing that the inequality is true. The first diagram illustrates the initial configuration of positive and negative weights, together with the “Hölder-type move” that takes one unit of weight from each of the points at  $(0, 2/3)$  and  $(2, -2/3)$  and moves it to the point  $(1, 0)$ , canceling out the weight of  $-2$  that was initially at  $(1, 0)$ . The second diagram illustrates the resulting configuration, together with the “ $L_p$  monotonicity move” that moves the  $-1$  weight at location  $(2, -1/3)$  towards the origin by a factor of  $2/3$  while scaling it by a factor of  $3/2$ , resulting in a point at  $(4/3, -2/9)$  with weight  $-3/2$ , which is now collinear with the remaining two points. The third diagram illustrates the final “Hölder-type move” that moves the two points with positive weight to their weighted average, zeroing out all weights.

535 The answer, as illustrated in Figure 1 is to first realize that 3 of the points lie on  
 536 a line, with the “ $-2$ ” halfway between the “ $\frac{3}{2}$ ” and the “ $2$ ”. Thus we take 1 unit from  
 537 each of the endpoints and cancel out the “ $-2$ ”. No three points are collinear now, so  
 538 we need to move one point onto the line formed by the other two: “ $-1$ ”, being negative,  
 539 can be moved towards the origin, so we move it until it crosses the line formed by the  
 540 two remaining numbers. This moves it  $\frac{1}{3}$  of the way to the origin, thus increasing  
 541 it from “ $-1$ ” to “ $-\frac{3}{2}$ ”; amazingly, this number, at position  $\frac{2}{3}(2, -\frac{1}{3}) = (\frac{4}{3}, -\frac{2}{9})$  is  
 542 now  $\frac{2}{3}$  of the way from the remaining “ $\frac{1}{2}$ ” at  $(0, \frac{2}{3})$  to the number “ $1$ ” at  $(2, -\frac{2}{3})$ ,  
 543 meaning that we can remove the final three numbers from the board in a single move,  
 544 winning the game. We thus made three moves total, two of the Hölder type, one of  
 545 the  $L_p$  monotonicity type. Reexpressing these moves as inequalities yields the desired  
 546 derivation of our inequality (Equation 6) as a product of powers of Hölder and  $L_p$   
 547 monotonicity inequalities, explicitly, as the product of the following three inequalities,

548 which are respectively 1) the square of a Cauchy-Schwarz inequality, 2) the 3/2 power  
 549 of an  $L_p$  monotonicity inequality for  $\lambda = 2/3$ , and 3) the 3/2 power of a Hölder  
 550 inequality for  $\lambda = 2/3$ :

$$\begin{aligned}
 551 \quad & \left( \sum_j x_j^2 y_j^{-2/3} \right) \left( \sum_j x_j^0 y_j^{2/3} \right) \left( \sum_j x_j^1 y_j^0 \right)^{-2} \geq 1 \\
 552 \quad & \left( \sum_j x_j^{4/3} y_j^{-2/9} \right)^{3/2} \left( \sum_j x_j^2 y_j^{-1/3} \right)^{-1} \geq 1 \\
 553 \quad & \left( \sum_j x_j^2 y_j^{-2/3} \right) \left( \sum_j x_j^0 y_j^{2/3} \right)^{1/2} \left( \sum_j x_j^{4/3} y_j^{-2/9} \right)^{-3/2} \geq 1 \\
 554 \quad &
 \end{aligned}$$

555 The above example demonstrates how transformative it is to know that the only possi-  
 556 ble ways of making progress proving a given inequality are by two simple possibilities,  
 557 thus transforming inequality proving into winning a 2d game with two types of moves.  
 558 As we have shown in Theorem 3, this process can be completed automatically in poly-  
 559 nomial time via linear programming; but in practice looking at the “2d game board”  
 560 is often all that is necessary, even for intricate counterintuitive inequalities like the  
 561 one above.

562 **3. An instance-optimal testing algorithm.** In this section we describe our  
 563 instance-by-instance optimal algorithm for verifying the identity of a distribution,  
 564 based on independent draws from the distribution. We begin by providing the defi-  
 565 nitions and terminology that will be used throughout the remainder of the paper. In  
 566 Section 3.2 we describe our very simple tester, and give some intuitions and motiva-  
 567 tions behind its form.

568 **3.1. Definitions.** We use  $[n]$  to denote the set  $\{1, \dots, n\}$ , and denote a distribu-  
 569 tion of support size  $n$  by  $p = p_1, \dots, p_n$ , where  $p_i$  is the probability of the  $i$ th domain  
 570 element. Throughout, we assume that all samples are drawn independently from the  
 571 distribution in question.

572 We denote the Poisson distribution with expectation  $\lambda$  by  $Poi(\lambda)$ , which has  
 573 probability density function  $poi(\lambda, i) = \frac{e^{-\lambda} \lambda^i}{i!}$ . We make heavy use of the standard  
 574 “Poissonization” trick (this goes back to at least Kolmogorov’s 1933 paper [13]; see  
 575 Chapter 5.4 of [14]). That is, rather than drawing  $k$  samples from a fixed distribution  
 576  $p$ , we first select  $k' \leftarrow Poi(k)$ , and then draw  $k'$  samples from  $p$ . Given such a  
 577 process, the number of times each domain element occurs is independent, with the  
 578 distribution of the number of occurrences of the  $i$ th domain element distributed as  
 579  $Poi(k \cdot p_i)$ . The independence yielded from Poissonization significantly simplifies many  
 580 kinds of analysis. Additionally, since  $Poi(k)$  is closely concentrated around  $k$ : from  
 581 both the perspective of upper bounds as well as lower bounds, at the cost of only  
 582 a subconstant factor, one may assume without loss of generality that one is given  
 583  $Poi(k)$  samples rather than exactly  $k$ .

584 Much of the analysis in this paper centers on  $L_p$  norms, where for a vector  $q$ , we  
 585 use the standard notation  $\|q\|_c$  to denote  $(\sum_i q_i^c)^{1/c}$ . The notation  $\|q\|_c^b$  is just the  
 586  $b$ th power of  $\|q\|_c$ . For example,  $\|q\|_{2/3}^{2/3} = \sum_i q_i^{2/3}$ .

587 **3.2. An optimal tester.** Our testing algorithm is extremely simple, and takes  
 588 the form of a simple statistic that is similar to Pearson’s chi-squared statistic, though  
 589 differs in two crucial ways. Given a set of  $k$  samples, with  $X_i$  denoting the number  
 590 of occurrences of the  $i$ th domain element, and  $p_i$  denoting the probability of drawing  
 591 the  $i$ th domain element from distribution  $p$ , the Pearson chi-squared statistic is given  
 592 as  $\sum_i \frac{1}{p_i} (X_i - kp_i)^2$ . Adding a constant does not change the behavior of the statistic,  
 593 and it will prove easier to compare with our statistic if we subtract  $k$  from each term,  
 594 yielding the following:

$$595 \quad (7) \quad \sum_i \frac{(X_i - kp_i)^2 - kp_i}{p_i}.$$

596 In the Poissonized setting (where the number of samples is drawn from a Poisson  
 597 distribution of expectation  $k$ ), if the samples are drawn from distribution  $p$ , then the  
 598 expectation of this chi-squared statistic is 0 because in that case  $X_i$  is distributed  
 599 according to a Poisson distribution of expectation  $kp_i$ , and hence has variance  $kp_i$ .  
 600 Our testing algorithm is, essentially, obtained by modifying this statistic in two ways:  
 601 replacing the second occurrence of  $kp_i$  with  $X_i$  (which has expectation  $kp_i$  when  
 602 drawing samples from  $p$  and thus does not change the statistic in expectation), and  
 603 changing the scaling factor from  $1/p_i$  to  $1/p_i^{2/3}$ :

$$604 \quad (8) \quad \sum_i \frac{(X_i - kp_i)^2 - X_i}{p_i^{2/3}}.$$

605 Note that this statistic still has the property that its expectation is 0 if the samples are  
 606 drawn from distribution  $p$ . The following examples motivate these two modifications.

607 **EXAMPLE 5.** *Let  $p$  be the distribution with  $p_1 = p_2 = 1/4$ , and where the re-*  
 608 *maining half of its probability mass composed of  $n/2$  domain elements, each oc-*  
 609 *curring with probability  $1/n$ . If we draw  $k = n^{2/3}$  samples from  $p$ , the contribu-*  
 610 *tion of the  $n/2$  small elements to the variance of Pearson’s statistic (Equation 7)*  
 611 *is  $\approx \frac{n}{2}(n^{-1/3}n^2) = \Omega(n^{8/3})$ , and the standard deviation would be  $\Omega(n^{4/3})$ . If the  $k$*   
 612 *samples were not drawn from  $p$ , and instead were drawn from distribution  $q$  that is*  
 613 *identical to  $p$ , except with  $p_1 = 1/8$  and  $p_2 = 3/8$ , then the expectation of Pearson’s*  
 614 *statistic would be  $O(n^{4/3})$ , though this signal might be buried by the  $\Omega(n^{4/3})$  standard*  
 615 *deviation due to the small domain elements.*

616 The above example illustrates that the scaling factor  $1/p_i$  in Pearson’s chi-squared  
 617 statistic places too much weight on the small elements, burying a drastic change in  
 618 the distribution (that could be detected with  $O(1)$  samples). Thus we are motivated  
 619 to consider a smoother scaling factor. There does not seem to be a simple intuition for  
 620 the  $2/3$  exponent in our statistic—it comes out of optimizing the interplay between  
 621 various inequalities in the analysis, and is cleanly revealed by our inequality prover  
 622 of Section 2. Intuitive reasoning from the perspective of the tester seems to lead  
 623 to a scaling factor of  $p_i^{1/2}$ , whereas intuitive reasoning from the perspective of the  
 624 lower bounds seems to lead to a scaling factor of  $p_i^{3/4}$ . Both intuitions turn out to be  
 625 misleading, and the correct scaling of  $p_i^{2/3}$ —resulting from balancing the upper and  
 626 lower bound desiderata—was unexpected.

627 The following example illustrates a second benefit of our statistic of Equation 8  
 628 over the chi-squared statistic, resulting from changing  $kp_i$  to  $X_i$ :



629 **EXAMPLE 6.** Let  $p$  be the distribution with  $p_1 = 1 - 1/n$ , and where the remain-  
 630 ing  $1/n$  probability mass is evenly split among  $n$  domain elements each with prob-  
 631 ability  $1/n^2$ . If we draw  $100 \cdot n$  samples, we are likely to see roughly  $100 \pm 10$  of  
 632 the “rare” domain elements, each exactly once. Such domain elements will have a  
 633 huge contribution to the variance of Pearson’s chi-squared statistic—a contribution  
 634 of  $\Omega(n^2)$ . On the other hand, these domain elements contribute almost nothing to  
 635 the variance of our statistic, because the contribution of such domain elements is  
 636  $((X_i - kp_i)^2 - X_i)p_i^{-2/3} \approx (X_i^2 - X_i)p_i^{-2/3}$ , which is 0 if  $X_i$  is 0 or 1 and with  
 637 overwhelming probability, none of these “rare” domain elements will occur more than  
 638 once. Hence our statistic is extremely robust to seeing rare things either 0 or 1 times,  
 639 and this significantly reduces the variance of our statistic.

640 We now formally define our tester and prove Theorem 2. The tester essentially  
 641 just computes the statistic of Equation 8, though one also needs to shave off a small  
 642  $O(\epsilon)$  portion of the distribution  $p$  before computing it, and also verify that not too  
 643 much probability mass lies on this supposedly small portion that was removed.

AN INSTANCE-OPTIMAL TESTER

644 Given a parameter  $\epsilon > 0$  and a set of  $k$  samples drawn from  $q$ , let  $X_i$  represent the  
 number of times the  $i$ th domain element occurs in the samples. Assume wlog that  
 the domain elements of  $p$  are sorted in non-increasing order of probability. Define  
 $s = \min\{i : \sum_{j>i} p_j \leq \epsilon/8\}$ , and let  $M = \{2, \dots, s\}$ , and  $S = \{s + 1, s + 2, \dots\}$ .

(Note that  $p_M = p_{-\epsilon/8}^{\max}$ .)

1. If  $\sum_{i \in M} \frac{(X_i - kp_i)^2 - X_i}{p_i^{2/3}} > 4k \|p_M\|_{2/3}^{1/3}$ , or
2. If  $\sum_{i \in S} X_i > \frac{3}{16} \epsilon k$ , then output “ $\|p - q\|_1 \geq \epsilon$ ”, else output “ $p = q$ ”.

645 For convenience, we restate Theorem 2, characterizing the performance of the  
 646 above tester.

647 **Theorem 2.** There exist constants  $c_1, c_2$  such that for any  $\epsilon > 0$  and any known  
 648 distribution  $p$ , for any unknown distribution  $q$ , our tester will distinguish  $q = p$  from  
 649  $\|p - q\|_1 \geq \epsilon$  with probability  $2/3$  when run on a set of at least  $c_1 \cdot \max\left\{\frac{1}{\epsilon}, \frac{\|p_{-\epsilon/16}^{\max}\|_{2/3}}{\epsilon^2}\right\}$   
 650 samples drawn from  $q$ , and no tester can do this task with probability at least  $2/3$  with  
 651 a set of fewer than  $c_2 \cdot \max\left\{\frac{1}{\epsilon}, \frac{\|p_{-2\epsilon}^{\max}\|_{2/3}}{\epsilon^2}\right\}$  samples.

652 Before proving the theorem, we provide some intuition behind the form of the  
 653 sample complexity,  $\max\left\{\frac{1}{\epsilon}, \frac{\|p_{-c\epsilon}^{\max}\|_{2/3}}{\epsilon^2}\right\}$ . The maximum with  $\frac{1}{\epsilon}$  only very rarely  
 654 comes into play: the  $\frac{2}{3}$  norm of a vector is always at least its 1 norm, so the max with  
 655  $\frac{1}{\epsilon}$  only takes over from  $\|p_{-c\epsilon}^{\max}\|_{2/3}/\epsilon^2$  if  $p$  is of the very special form where removing  
 656 its max element and its smallest  $c\epsilon$  mass leaves less than  $\epsilon$  probability mass remaining;  
 657 the max expression thus prevents the sample size in the theorem from going to 0 in  
 658 extreme versions of this case.

659 The subscript and superscript in  $\|p_{-c\epsilon}^{\max}\|_{2/3}$  each reduce the final value, and  
 660 mark two ways in which the problem might be “unexpectedly easy”. To see the  
 661 intuition behind these two modifications in the vector of probabilities, note that if the  
 662 distribution  $p$  contains a single domain element  $p_m$  that comprises the majority of the  
 663 probability mass, then in some sense it is hard to hide changes in  $p$ : at least half of  
 664 the discrepancy between  $p$  and  $q$  must lie in other domain elements, and if these other

665 domain elements comprise just a tiny fraction of the total probability mass, then the  
 666 fact that half the discrepancy is concentrated on a tiny fraction of the distribution  
 667 makes recognizing such discrepancy easier.

668 On the other hand, having many small domain elements makes the identity testing  
 669 problem harder, as indicated by the  $L_{2/3}$  norm, however only “harder up to a point”.  
 670 If most of the  $L_{2/3}$  norm of  $p$  comes from a portion of the distribution with tiny  $L_1$   
 671 norm, then it is also hard to “hide” much discrepancy in this region: if a portion  
 672 of the domain consisting of  $\epsilon/3$  total mass in  $p$  has discrepancy  $\epsilon$  between  $p$  and  $q$ ,  
 673 then the probability mass of these elements in  $q$  must total at least  $\frac{2}{3}\epsilon$  by the triangle  
 674 inequality, namely at least twice what we would expect if  $q = p$ ; this discrepancy is  
 675 thus easy to detect in  $O(\frac{1}{\epsilon})$  samples. Thus discrepancy cannot hide in the very small  
 676 portion of the distribution, and we may effectively ignore the small portion of the  
 677 distribution when figuring out how hard it is to test discrepancy.

678 In these two ways—represented by the subscript and superscript of  $p_{-c\epsilon}^{-\max}$  in our  
 679 results—the identity testing problem may be “easier” than the simplified  $O(\frac{\|p\|_{2/3}}{\epsilon^2})$   
 680 bound. But our corresponding lower bound shows that these are the only ways.

681 **Remark on “tolerant testing”.** We note that the “yes” case of the theorem, where  
 682  $q = p$ , can always be relaxed to a “tolerant testing” condition  $\|p - q\|_1 \leq O(\frac{1}{k})$  where  
 683  $k = c_1 \cdot \max\left\{\frac{1}{\epsilon}, \frac{\|p_{-\epsilon/16}^{-\max}\|_{2/3}}{\epsilon^2}\right\}$  is the number of samples used. This kind of tolerant  
 684 testing result is true for *any* tester, because statistical distance is subadditive on  
 685 product distributions, so a change of  $\frac{\epsilon}{k}$  in the distribution  $p$  can induce a change of at  
 686 most  $c$  on the distribution of the output of any testing algorithm that uses  $k$  samples.  
 687 A more refined analysis of our tester (or a tester tailored to the tolerant regime) yields  
 688 better bounds in some cases. However, the problem of distinguishing  $\|p - q\|_1 \leq \epsilon_1$   
 689 from  $\|p - q\|_1 \geq \epsilon_2$  enters a very different regime when  $\epsilon_1$  is not much smaller than  $\epsilon_2$ ,  
 690 and many more samples are required. (These problems are very related to the task  
 691 of *estimating* the distance from  $q$  to the known distribution  $p$ .) For any constants  
 692  $\epsilon_1 < \epsilon_2$ , it requires  $\Theta(\frac{n}{\log n})$  samples to distinguish  $\|p - q\|_1 \leq \epsilon_1$  from  $\|p - q\|_1 \geq \epsilon_2$   
 693 when  $p$  is the uniform distribution on  $n$  elements, many more than the  $\sqrt{n}$  needed  
 694 here [20, 21].

695 **3.3. Analysis of the tester.** The core of the proof of the algorithmic direction  
 696 of Theorem 2 is an application of Chebyshev’s inequality: first arguing that if the  
 697 samples were drawn from a distribution  $q$  with  $\|p - q\|_1 \geq \epsilon$ , then the expectation of  
 698 the statistic in question is large in comparison to its standard deviation, whereas if the  
 699 samples were drawn from  $q = p$ , then the expectation is 0 and the standard deviation  
 700 is sufficiently small so that the distribution of the statistic will not overlap significantly  
 701 with the previous case (where  $\|p - q\|_1 \geq \epsilon$ ). In order to prove the desired inequalities  
 702 relating the expectation and the variance, we reexpress these inequalities in terms  
 703 of the two sequences of positive numbers  $p = p_1, p_2, \dots$ , and  $\Delta = \Delta_1, \Delta_2, \dots$ , with  
 704  $\Delta_i := |p_i - q_i|$ , leading to an expression that is the sum of five inequalities essentially of  
 705 the canonical form  $\prod_i \left(\sum_j p_j^{a_i} \Delta_j^{b_i}\right)^{c_i} \geq 1$ . The machinery of Section 2 thus yields an  
 706 easily verifiable derivation of the desired inequalities as a product of positive powers of  
 707 Hölder type inequalities, and  $L_p$  monotonicity inequalities. For the sake of presenting  
 708 a self-contained complete proof of Theorem 2, we write out these derivations explicitly  
 709 below.

710 We now begin the analysis of the performance of the above tester, establishing

711 the upper bounds of Theorem 2. When  $\|p - q\|_1 \geq \epsilon$ , we note that at most half of  
 712 the discrepancy is accounted for by the most frequently occurring domain element of  
 713  $p$ , since the total probability masses of  $p$  and  $q$  must be equal (to 1), and thus  $\geq \epsilon/2$   
 714 discrepancy must occur on the remaining elements. We split the analysis into two  
 715 cases: when a significant portion of the remaining  $\epsilon/2$  discrepancy falls above  $s$  then we  
 716 show that case 1 of the algorithm will recognize it; otherwise, if  $\|p_{<s} - q_{<s}\|_1 \geq (3/8)\epsilon$ ,  
 717 then case 2 of the algorithm will recognize it.

We first analyze the mean and variance of the left hand side of the first condition  
 of the tester, under the assumption (as discussed in Section 3.1) that a Poisson-  
 distributed number of samples,  $Poi(k)$  is used. This makes the number of times each  
 domain element is seen,  $X_i$ , be distributed as  $Poi(kq_i)$ , and makes all  $X_i$  independent  
 of each other. It is thus easy to calculate the mean and variance of each term.  
 Explicitly, defining  $\Delta_i = p_i - q_i$  we have

$$E_{X_i \leftarrow Poi(kq_i)} \left[ [(X_i - kp_i)^2 - X_i] p_i^{-2/3} \right] = k^2 \Delta_i^2 p_i^{-2/3}$$

and

$$Var_{X_i \leftarrow Poi(kq_i)} \left[ [(X_i - kp_i)^2 - X_i] p_i^{-2/3} \right] = [2k^2(p_i - \Delta_i)^2 + 4k^3(p_i - \Delta_i)\Delta_i^2] p_i^{-4/3}$$

718 Note that when  $p = q$ , the expectation is 0, since  $\Delta_i \equiv 0$ . However, in the case  
 719 that a significant portion of the  $\epsilon$  deviation between  $p$  and  $q$  occurs in the region above  
 720  $s$ , we show that for suitable  $k$ , the variance is somewhat less than the square of the  
 721 expectation, leading to a reliable test for distinguishing this case from the  $p = q$  case.

722 The motivation for the convoluted steps in the derivations in the following lemma  
 723 comes entirely from the general inequality result of Theorem 3, though as guaranteed  
 724 by that theorem, the resulting inequalities can all be derived by elementary means  
 725 without reference to the theorem.

726 As defined in the tester, considering the elements of  $p$  to be sorted in decreasing  
 727 order by probability, we let  $s$  be the smallest integer so that  $\sum_{i>s} \leq \epsilon/8$ . For  
 728 notational convenience, we define the set  $M = \{2, \dots, s\}$ , so that  $p_M$  consists of those  
 729 elements of  $p$  that have “medium” probabilities—not the largest element, and not  
 730 the smallest elements that comprise  $\leq \epsilon/8$  probability. We define  $M$  so that we may  
 731 explicitly analyze the corresponding discrepancies  $\Delta_M$ . (Note that the probabilities  
 732 in the distribution  $q$  will typically not be sorted, and may not be similar to the  
 733 corresponding probabilities in  $p$ ).

734 The following lemma shows that the variance of case 1 of our estimator can be  
 735 made arbitrarily smaller than the square of its expectation, which we will use for a  
 736 Chebyshev bound proof in Proposition 11 below.

LEMMA 10. For any  $c \geq 1$ , if  $k = c \cdot \max\left\{\frac{\|p_M\|_{2/3}^{1/3}}{p_s^{1/3} \cdot (\epsilon/8)}, \frac{\|p_M\|_{2/3}}{(\epsilon/8)^2}\right\}$  and if at least  $\epsilon/8$   
 of the discrepancy falls in the medium region, namely  $\sum_{i \in M} |\Delta_i| \geq \epsilon/8$ , then

$$\sum_{i \in M} [2k^2(p_i - \Delta_i)^2 + 4k^3(p_i - \Delta_i)\Delta_i^2] p_i^{-4/3} < \frac{16}{c} \left[ \sum_{i \in M} k^2 \Delta_i^2 p_i^{-2/3} \right]^2$$

737 *Proof.* Dividing both sides by  $k^4$ , the left hand side has terms proportional to  
 738  $(p_i - \Delta_i)/k$  and its square. We bound such terms via the triangle inequality and the

739 definition of  $k$  as  $(p_i - \Delta_i)/k \leq \left( p_i \frac{(\epsilon/8)^2}{\|p_M\|_{2/3}} + |\Delta_i| \frac{p_i^{1/3}(\epsilon/8)}{\|p_M\|_{2/3}^{1/3}} \right) / c$ . Expanding, yields the  
 740 left hand side divided by  $k^4$  bounded as the sum of 5 terms:

$$741 \quad \sum_{i \in M} \frac{2}{c^2} \left( p_i^{2/3} \frac{(\epsilon/8)^4}{\|p_M\|_{2/3}^2} + 2|\Delta_i| p_i^{-1/3} \frac{p_s^{1/3}(\epsilon/8)^3}{\|p_M\|_{2/3}^{4/3}} + \Delta_i^2 p_i^{-4/3} \frac{p_s^{2/3}(\epsilon/8)^2}{\|p_M\|_{2/3}^{2/3}} \right)$$

$$742 \quad + \frac{4}{c} \left( \Delta_i^2 p_i^{-1/3} \frac{(\epsilon/8)^2}{\|p_M\|_{2/3}} + |\Delta_i^3| p_i^{-4/3} \frac{p_s^{1/3}(\epsilon/8)}{\|p_M\|_{2/3}^{1/3}} \right).$$

743 We bound each of the five terms separately by  $\left[ \sum_{i \in M} \Delta_i^2 p_i^{-2/3} \right]^2$ , using the fact  
 744 that  $\frac{1}{c^2} \leq \frac{1}{c}$ , and sum the constants  $2(1 + 2 + 1) + 4(1 + 1)$  to yield 16 on the right  
 745 hand side.

746 1. Cauchy-Schwarz yields  $\sum_{i \in M} \Delta_i^2 p_i^{-2/3} \geq (\sum_{i \in M} |\Delta_i|)^2 / (\sum_{i \in M} p_i^{2/3}) \geq$   
 747  $(\frac{\epsilon}{8})^2 / \|p_M\|_{2/3}^{2/3}$ . Squaring this inequality and noting that, by definition,  $\sum_{i \in M} p_i^{2/3} =$   
 748  $\|p_M\|_{2/3}^{2/3}$  bounds the first term as desired.

749 2. We bound  $\frac{\epsilon}{p_s^{1/3}} = \frac{\epsilon}{\|\Delta_M\|_1} \sum_{i \in M} |\Delta_i| p_s^{-1/3} \geq \frac{\epsilon}{\|\Delta_M\|_1} \sum_{i \in M} |\Delta_i| p_i^{-1/3}$  since  $p_i \geq$   
 750  $p_s$  for  $i \in M$ . Multiplying this inequality by the square of the Cauchy-Schwarz  
 751 inequality of the previous case:  $\left( \sum_{i \in M} \Delta_i^2 p_i^{-2/3} \right)^2 \geq \|\Delta_M\|_1^4 / \|p_M\|_{2/3}^{4/3}$  and the bound  
 752  $\|\Delta_M\|_1^3 \geq (\frac{\epsilon}{8})^3$  yields the desired bound on the second term.

753 3. Simplifying the third term via  $p_i^{-4/3} p_s^{2/3} \leq p_i^{-2/3}$  lets us bound this term as  
 754 the product of the Cauchy-Schwarz inequality of the first case:  $\sum_{i \in M} \Delta_i^2 p_i^{-2/3} \geq$   
 755  $\|\Delta_M\|_1^2 / \|p_M\|_{2/3}^{2/3}$  and the bound  $\|\Delta_M\|_1^2 \geq (\frac{\epsilon}{8})^2$ .

756 4. Here and in the next case we use the basic fact that for  $\beta > \alpha > 0$  and  
 757 a (nonnegative) vector  $z$  we have  $\|z\|_\beta \leq \|z\|_\alpha$  (with equality only when  $z$  has at  
 758 most one nonzero entry). Thus  $\sum_{i \in M} \Delta_i^2 p_i^{-1/3} \leq \left( \sum_{i \in M} \Delta_i^{4/3} p_i^{-2/9} \right)^{3/2}$ , and this  
 759 last expression is bounded via (the  $3/2$  power of) Hölder's inequality for  $\lambda = 2/3$   
 760 by  $\left( \sum_{i \in M} \Delta_i^2 p_i^{-2/3} \right) \left( \sum_{i \in M} p_i^{2/3} \right)^{1/2}$ . Multiplying this inequality by the Cauchy-  
 761 Schwarz inequality of the first case:  $\|\Delta_M\|_1^2 / \|p_M\|_{2/3}^{2/3} \leq \sum_{i \in M} \Delta_i^2 p_i^{-2/3}$  and the bound  
 762  $(\frac{\epsilon}{8})^2 \leq \|\Delta_M\|_1^2$  yields the desired bound on the fourth term.

5. The norm inequality from the previous case also yields

$$\sum_{i \in M} \Delta_i^3 p_i^{-4/3} \leq \left( \sum_{i \in M} \Delta_i^2 p_i^{-8/9} \right)^{3/2} \leq p_s^{-1/3} \left( \sum_{i \in M} \Delta_i^2 p_i^{-2/3} \right)^{3/2}.$$

Multiplying by the square root of the Cauchy-Schwarz bound of the first case,

$$\|\Delta_M\|_1 / \|p_M\|_{2/3}^{1/3} \leq \left( \sum_{i \in M} \Delta_i^2 p_i^{-2/3} \right)^{1/2}$$

763 and the bound  $\frac{\epsilon}{8} \leq \|\Delta_M\|_1$  yields the desired bound on the fifth term.

764 We now prove the upper bound portion of Theorem 2.

765 PROPOSITION 11. *There exists a constant  $c_1$  such that for any  $\epsilon > 0$  and any*  
 766 *known distribution  $p$ , for any unknown distribution  $q$  on the same domain, our tester*  
 767 *will distinguish  $q = p$  from  $\|p - q\|_1 \geq \epsilon$  with probability  $2/3$  using a set of  $k =$*   
 768  $c_1 \cdot \max \left\{ \frac{1}{\epsilon}, \frac{\|p_{-\epsilon/16}^{-\max}\|_{2/3}}{\epsilon^2} \right\}$  *samples.*

769 *Proof.* We first show that if  $p = q$  then the tester will recognize this fact with  
 770 high probability.

Consider the first test of the algorithm, whether

$$\sum_{i \in M} [(X_i - kp_i)^2 - X_i] p_i^{-2/3} > 4k \|p_M\|_{2/3}^{1/3}.$$

771 As calculated above, the expectation of the left hand side is 0 in this case, and the  
 772 variance is  $2k^2 \|p_M\|_{2/3}^{2/3}$ . Thus Chebyshev's inequality yields that this random variable  
 773 will be greater than  $2\sqrt{2}$  standard deviations from its mean with probability at most  
 774  $1/8$ , and thus the first test will be accurate with probability at least  $7/8$  in this case.

775 For the second test, whether  $\sum_{i \in S} X_i > \frac{3}{16} \epsilon k$ , recall that  $S$  was defined to contain  
 776 those elements of  $p$  with probabilities smaller than the ‘‘medium’’ elements  $M$ , and,  
 777 explicitly, have total probability mass  $\|p_S\| \leq \epsilon/8$ . Denote this total mass by  $m$ . Thus  
 778  $\sum_{i \in S} X_i$  is distributed as  $Poi(mk)$ , which has mean and variance both  $mk \leq \frac{\epsilon k}{8}$ .

779 Thus Chebyshev's inequality yields that the probability that this quantity exceeds  
 780  $\frac{3}{16} \epsilon k$  is at most  $\left( \frac{\sqrt{mk}}{(3/16)\epsilon k - mk} \right)^2 \leq \left( \frac{\sqrt{\epsilon k}}{\sqrt{8}(1/16)\epsilon k} \right)^2 = \frac{2^5}{\epsilon k}$ . Hence provided  $k \geq \frac{2^8}{\epsilon}$ , this  
 781 probability will be at most  $1/8$ . For the sake of what follows, we actually make  $k$  at  
 782 least twice as large as this, setting  $c_1 \geq 2^9$  so that, from the definition of  $k$ , we have

$$783 k = c_1 \cdot \max \left\{ \frac{1}{\epsilon}, \frac{\|p_{-\epsilon/16}^{-\max}\|_{2/3}}{\epsilon^2} \right\} \geq \frac{2^9}{\epsilon}.$$

784 We now consider the case when  $\|p - q\|_1 \geq \epsilon$ , and show that the tester is also  
 785 correct in this setting. Consider the element with largest probability under distri-  
 786 bution  $p$ , and note that at most half of the discrepancy  $\|p - q\|_1$  can be due to the  
 787 difference in probabilities assigned to this one element, since the total probability  
 788 masses of  $p$  and  $q$  are equal (to 1). Thus at least half the discrepancy between  $p$  and  
 789  $q$  occurs on the remaining elements, which consist of the elements in  $S \cup M$ . Hence  
 790  $\|(p - q)_{S \cup M}\|_1 \geq \epsilon/2$ . We consider two cases. If  $\|(p - q)_S\|_1 \geq \frac{3}{8} \epsilon$ , namely if most of  
 791 the at least  $\epsilon/2$  discrepancy occurs on the small elements, then since  $\|p_S\|_1 \leq \frac{1}{8} \epsilon$  by  
 792 assumption, the triangle inequality yields that  $\|q_S\|_1 \geq \frac{1}{4} \epsilon$ . Consider the second test  
 793 in this case. Analogously to the argument above, Chebyshev's inequality shows that  
 794 this test will pass except with probability at most  $\frac{64}{\epsilon k}$ . Hence since  $k \geq \frac{2^9}{\epsilon}$  from the  
 795 previous paragraph, we have that the algorithm will be successful in this case with  
 796 probability at least  $7/8$ .

797 In the remaining case,  $\|(p - q)_M\|_1 \geq \frac{1}{8} \epsilon$ , we apply Lemma 10. We first show  
 798 that the number of samples  $k = c_1 \frac{\|p_{-\epsilon/16}^{-\max}\|_{2/3}}{\epsilon^2}$  is at least as many as needed for the  
 799 lemma,  $c \cdot \max \left\{ \frac{\|p_M\|_{2/3}^{1/3}}{p_s^{1/3}(\epsilon/8)}, \frac{\|p_M\|_{2/3}}{(\epsilon/8)^2} \right\}$ , provided  $c_1 \geq 128c$ . The second component  
 800 of this maximum is trivially less than or equal to  $k$ , since by definition  $\|p_M\|_{2/3} =$   
 801  $\|p_{-\epsilon/8}^{-\max}\|_{2/3} \leq \|p_{-\epsilon/16}^{-\max}\|_{2/3}$ . To bound the first component, we let  $r$  (analogously to  
 802  $s$ ) be defined as the smallest integer such that  $\sum_{i > r} p_i \leq \epsilon/16$ , recalling that the  
 803 probabilities  $p_i$  are sorted in decreasing order. Since  $\sum_{i \geq s} p_i = \sum_{i \in S \cup \{s\}} p_i \geq \epsilon/8$ ,  
 804 the difference of these expressions yields  $\sum_{i=s}^r p_i \geq \epsilon/16$ . Since each  $p_i$  in this last

805 sum is at most  $p_s$ , we have that  $p_i^{-1/3} \geq p_s^{-1/3}$  for such  $i$ , which yields  $\sum_{i=s}^r p_i^{2/3} \geq$   
 806  $\frac{\epsilon}{16p_s^{1/3}}$ . Thus  $\|p_{-\epsilon/16}^{-\max}\|_{2/3}^{2/3} = \sum_{i=2}^r p_i^{2/3} \geq \sum_{i=s}^r p_i^{2/3} \geq \frac{\epsilon}{16p_s^{1/3}}$ , where the second-to-last  
 807 inequality assumes  $s \neq 1$ . Multiplying by the inequality  $\|p_{-\epsilon/16}^{-\max}\|_{2/3}^{1/3} \geq \|p_{-\epsilon/8}^{-\max}\|_{2/3}^{1/3}$   
 808 yields the bound. (In the unusual case that  $s = 1$ , the set  $M = \{2, \dots, s\}$  is empty,  
 809 and thus Lemma 10 is trivially true, requiring 0 samples, which we trivially have.)

810 We thus invoke Lemma 10, which shows that, for any  $c \geq 1$ , the expectation of  
 811 the left hand side of the first test,  $\sum_{i \in M} [(X_i - kp_i)^2 - X_i] p_i^{-2/3}$ , is at least  $\sqrt{c/16}$   
 812 times its standard deviation; further, we note that the triangle-inequality expression  
 813 by which we bounded the standard deviation is minimized when  $p = q$ , in which case,  
 814 as noted above, the standard deviation is  $\sqrt{2}k\|p_M\|_{2/3}^{1/3}$ . Thus the expression on the  
 815 right hand side of the first test,  $4k\|p_M\|_{2/3}^{1/3}$ , is always at least  $\sqrt{c/16} - 2\sqrt{2}$  standard  
 816 deviations away from the mean of the left hand side. Thus for  $c \geq 512$ , Chebyshev's  
 817 inequality yields that the first test will correctly report that  $p$  and  $q$  are different with  
 818 probability at least  $7/8$ .

819 Thus by the union bound, in either case  $p = q$  or  $\|p - q\|_1 \geq \epsilon$ , the tester will  
 820 correctly report it with probability at least  $\frac{3}{4}$ .  $\square$

821 **4. Lower bounds.** In this section we show how to construct distributions that  
 822 are very hard to distinguish from a given distribution  $p$  despite being far from  $p$ ,  
 823 establishing the lower bound portion of Theorem 2. Explicitly, we will construct  
 824 a distribution over distributions, that we will call  $Q_\epsilon$ , such that most distributions  
 825 in  $Q_\epsilon$  are far from  $p$ , yet  $k$  samples from a randomly chosen member of  $Q_\epsilon$  will be  
 826 distributed very close to the distribution of  $k$  samples from  $p$ . Analyzing the statistics  
 827 of such sampling processes can be enormously involved (see for example the lower  
 828 bounds of [20], which involve deriving new and general central limit theorems in high  
 829 dimensions).

830 In this paper, however, we show that the statistics of  $k$  samples from a ran-  
 831 domly chosen distribution from  $Q_\epsilon$  can be captured much more directly, by a product  
 832 distribution over univariate distributions that are a ‘‘coin flip between Poisson dis-  
 833 tributions.’’ Thus we can analyze this process dimension-by-dimension and sum the  
 834 distances. That is, if  $d_i$  is the distance between what happens for the  $i$ th domain  
 835 element given  $k$  samples from  $p$  versus  $k$  samples from the product distribution ‘‘cap-  
 836 turing’’  $Q_\epsilon$ , we can sum these up to bound the probability of distinguishing  $p$  from  
 837  $Q_\epsilon$  by  $\sum_i d_i$ . However, this is not good enough for us since the actual probability of  
 838 distinguishing these two cases for an ideal tester is more like the  $L_2$  norm of these  $d_i$   
 839 distances instead of the  $L_1$  norm—to achieve a tight result we need something like  
 840  $\sqrt{\sum_i d_i^2}$  instead of  $\sum_i d_i$ .

To accomplish this, we analyze all distances below via the *Hellinger distance*,

$$H(p, q) = \frac{1}{\sqrt{2}} \sqrt{\sum_i (\sqrt{p_i} - \sqrt{q_i})^2}.$$

841 Hellinger distance has two properties perfectly suited for our task: its *square* is sub-  
 842 additive on product distributions (meaning it combines via the  $L_2$  norm instead of  
 843 the  $L_1$  norm), and the Hellinger distance (times  $\sqrt{2}$ ) bounds the statistical distance.  
 844 See [3] for a more in-depth discussion of Hellinger distance and its applications to  
 845 hypothesis testing lower bounds.

846 We first prove a technical but ultimately straightforward lemma characterizing the  
 847 Hellinger distance between the ‘‘coin flip between Poisson distributions’’ mentioned

848 above and a regular Poisson distribution. We then show how a product distribution  
 849 of these coin flip distributions forms a powerful class of testing lowerbounds, Theo-  
 850 rem 13, which has already found use in [8]. We then assemble the pieces using some  
 851 inequalities, to show the lowerbound portion of Theorem 2.

852 Let  $Poi(\lambda \pm \epsilon)$  denote the probability distribution with pdf over nonnegative  
 853 integers  $i$ :  $\frac{1}{2}poi(\lambda + \epsilon) + \frac{1}{2}poi(\lambda - \epsilon)$ , which is only defined for  $\epsilon \leq \lambda$ .

854 LEMMA 12.  $H(Poi(\lambda), Poi(\lambda \pm \epsilon)) \leq c \cdot \frac{\epsilon^2}{\lambda}$  for constant  $c$ .

855 *Proof.* Assume throughout this proof that  $\epsilon \leq \frac{1}{2}\sqrt{\lambda}$ , for otherwise the lemma is  
 856 trivially true.

We bound

$$H(Poi(\lambda), Poi(\lambda \pm \epsilon))^2 = \frac{1}{2} \sum_{i \geq 0} \left( \sqrt{\frac{e^{-\lambda} \lambda^i}{i!}} - \sqrt{\frac{1}{2} \left[ \frac{e^{-\lambda - \epsilon} (\lambda + \epsilon)^i}{i!} + \frac{e^{-\lambda + \epsilon} (\lambda - \epsilon)^i}{i!} \right]} \right)^2$$

term-by-term via the inequality  $|\sqrt{a} - \sqrt{b}| \leq \frac{|a-b|}{\sqrt{b}}$ . We let  $a = \frac{e^{-\lambda} \lambda^i}{i!}$  and  $b =$   
 $\frac{1}{2} \left[ \frac{e^{-\lambda - \epsilon} (\lambda + \epsilon)^i}{i!} + \frac{e^{-\lambda + \epsilon} (\lambda - \epsilon)^i}{i!} \right]$  for some specific  $i$ , and sum over  $i$  later. We bound the  
 numerator of  $\frac{|a-b|}{\sqrt{b}}$  by noting that

$$|a - b| = \left| \frac{e^{-\lambda} \lambda^i}{i!} - \frac{1}{2} \frac{e^{-\lambda - \epsilon} (\lambda + \epsilon)^i}{i!} - \frac{1}{2} \frac{e^{-\lambda + \epsilon} (\lambda - \epsilon)^i}{i!} \right|$$

857 is bounded by  $\frac{1}{2}\epsilon^2$  times the maximum magnitude of the second derivative with respect  
 858 to  $x$  of  $poi(x, i)$  for  $x \in [\lambda - \epsilon, \lambda + \epsilon]$ . Explicitly,  $\frac{d^2}{dx^2} \frac{e^{-x} x^i}{i!} = poi(x, i) \frac{(i-x)^2 - i}{x^2}$ .

859 For the denominator of  $\frac{|a-b|}{\sqrt{b}}$  we will first bound it in the case when  $\lambda \geq 1$ , in which  
 860 case since  $\epsilon \leq \frac{1}{2}\sqrt{\lambda}$ , there is an absolute constant  $c$  such that for any  $x \in [\lambda - \epsilon, \lambda + \epsilon]$   
 861 we have  $poi(x, i) \leq c \cdot b = \frac{1}{2}c[Poi(\lambda - \epsilon) + Poi(\lambda + \epsilon)]$ . Let  $x^*$  be the value of  $x$  in  
 862 the interval  $[\lambda - \epsilon, \lambda + \epsilon]$  where  $poi(x, i)$  is maximized. Thus the denominator  $\sqrt{b}$  is at  
 863 least  $\sqrt{\frac{1}{c}poi(x^*, i)}$ .

864 We combine the bounds of the previous two paragraphs to conclude the case  $\lambda \geq 1$ .  
 865 Thus we have  $\frac{|a-b|}{\sqrt{b}} \leq \frac{\sqrt{c}}{2}\epsilon^2 \sqrt{poi(x^*, i)} \max_{x \in [\lambda - \epsilon, \lambda + \epsilon]} \left| \frac{(i-x)^2 - i}{x^2} \right|$ . Since  $\lambda - \epsilon \geq \frac{1}{2}$  in  
 866 our case, this last expression is thus bounded as  $c_2 \epsilon^2 \sqrt{poi(x^*, i)} \frac{(i-\lambda)^2 + i}{\lambda^2}$  for some  
 867 constant  $c_2$ . We thus sum the square of this expression, over all  $i \geq 0$ , to obtain our  
 868 bound on the (square of the) Hellinger distance. Since  $poi(x^*, i)$  dies off exponentially  
 869 outside an interval of width  $O(\sqrt{\lambda})$ , we may bound the sum over all  $i$  as just a constant  
 870 times the sum over an interval of width  $\sqrt{\lambda}$  centered at  $x^*$ . We note that  $poi(x^*, i)$  is  
 871 bounded by a constant multiple of  $\frac{1}{\sqrt{\lambda}}$ ; since we are considering  $i$  within  $\frac{1}{2}\sqrt{\lambda}$  of  $x^*$ ,  
 872 which is within  $\frac{1}{2}\sqrt{\lambda}$  of  $\lambda$  by definition, we have that  $i$  is bounded by a constant times  
 873  $\lambda$ , as is  $(i - \lambda)^2$ . Thus, in total for the square of the Hellinger distance, we have  $\sqrt{\lambda}$   
 874 terms that are each bounded as  $\left( c_2 \epsilon^2 \sqrt{poi(x^*, i)} \frac{(i-\lambda)^2 + i}{\lambda^2} \right)^2 \leq c_3 \epsilon^4 \frac{1}{\sqrt{\lambda}} \frac{\lambda^2}{\lambda^4} = c_3 \frac{\epsilon^4}{\lambda^2 \sqrt{\lambda}}$   
 875 for some constant  $c_3$ . Multiplying by the number of terms,  $\sqrt{\lambda}$ , yields the desired  
 876 bound.

877 For the case  $\lambda < 1$ , we note that the second derivative of  $poi(x, i)$  is globally  
 878 bounded by a constant, bounding the numerator of  $\frac{|a-b|}{\sqrt{b}}$  by  $O(\epsilon^2)$ . To bound the  
 879 denominator, we note that, for  $\lambda < 1$ , the value  $b = \frac{1}{2} \left[ \frac{e^{-\lambda - \epsilon} (\lambda + \epsilon)^i}{i!} + \frac{e^{-\lambda + \epsilon} (\lambda - \epsilon)^i}{i!} \right]$  is

880  $\Omega(1)$  for  $i = 0$ , it is  $\Omega(\lambda)$  for  $i = 1$ , and it is  $\Omega(\lambda^2)$  for  $i = 2$ , thus yielding a bound of  
 881  $O(\frac{\epsilon^4}{\lambda^2})$  on each of the first three terms in the expression for  $H^2$ . For  $i \geq 3$  we have,  
 882 for  $x \in (0, 2\lambda]$  that  $\frac{d^2}{dx^2} \text{poi}(x, i) = \text{poi}(x, i) \frac{(i-x)^2 - i}{x^2} = O(\frac{\lambda^{i-2} i^2}{i!})$ . Thus the numerator  
 883 of  $\frac{|a-b|}{\sqrt{b}}$  is bounded by  $\epsilon^2$  times this. To bound the denominator, we have that  $b \geq$   
 884  $\frac{1}{2} \text{poi}(\lambda + \epsilon, i) = \Omega(\frac{\lambda^i}{i!})$ , leading to a combined bound of  $\frac{|a-b|}{\sqrt{b}} = O(\epsilon^2 \lambda^{i/2-2} \frac{i^2}{\sqrt{i!}})$ , which  
 885 is bounded as  $O(\frac{\epsilon^2}{\lambda} \frac{i^2}{\sqrt{i!}})$  since  $i \geq 3$  and  $\lambda < 1$ . Summing up the square of this over  
 886 all  $i \geq 3$  clearly yields  $O(\frac{\epsilon^4}{\lambda^2})$ , the desired bound.

887 Thus in all cases the square of the Hellinger distance is  $O(\frac{\epsilon^4}{\lambda^2})$ , yielding the lemma.

888 This lemma is a crucial ingredient in the proof of the following general lower  
 889 bound.

890 **THEOREM 13.** *Given a distribution  $p$ , and associated values  $\epsilon_i$  such that  $\epsilon_i \in$   
 891  $[0, p_i]$  for each domain element  $i$ , define the distribution over distributions  $Q_\epsilon$  by the  
 892 process: for each domain element  $i$ , randomly choose  $q_i = p_i \pm \epsilon_i$ , and then normalize  
 893  $q$  to be a distribution. Then there exists a constant  $c$  such that it takes at least  
 894  $c \left( \sum_i \frac{\epsilon_i^4}{p_i^2} \right)^{-1/2}$  samples to distinguish  $p$  from  $Q_\epsilon$  with success probability  $2/3$ . Further,  
 895 with probability at least  $1/2$ , the  $L_1$  distance between a random distribution from  $Q_\epsilon$   
 896 and  $p$  is at least  $\min\{(\sum_{i \neq \arg \max \epsilon_i} \epsilon_i), \frac{1}{2} \sum_i \epsilon_i\}$ .*

897 The lower bound portion of Theorem 2 follows from the above theorem by appro-  
 898 priately choosing the sequence  $\epsilon_i$ .

899 *Proof of Theorem 13.* For the first part of the theorem, we first analyze the trivial  
 900 case where  $\sum_i \epsilon_i^2 \geq \frac{1}{64}$ . The inequality  $\sum_i p_i^2 \leq 1$  ( $L_p$  monotonicity) and Cauchy-  
 901 Schwarz yield that  $\sum_i \frac{\epsilon_i^4}{p_i^2} \geq \sum_i p_i^2 \sum_i \frac{\epsilon_i^4}{p_i^2} \geq (\sum_i \epsilon_i^2)^2 \geq \frac{1}{64^2}$ , which means the number  
 902 of samples requested by the theorem can be made 1 by setting  $c \leq \frac{1}{64}$ ; and clearly at  
 903 least 1 sample is needed to distinguish different distributions, yielding the theorem in  
 904 this case.

905 Otherwise, we assume  $\sum_i \epsilon_i^2 < \frac{1}{64}$ . Consider the following distributions, which  
 906 emulate the number of times each domain element is seen in  $Q_\epsilon$  and  $p$  if we take  
 907  $\text{Poi}(2k)$  samples: first randomly generate  $\bar{q}_i = p_i \pm \epsilon_i$  without normalizing, and then  
 908 for each  $i$  draw a sample from  $\text{Poi}(\bar{q}_i \cdot 2k)$ ; compare this to, for each  $i$ , drawing a sample  
 909 from  $\text{Poi}(p_i \cdot 2k)$ . Since  $\sum_i \bar{q}_i$  has mean 1 and variance  $\sum_i \epsilon_i^2 < \frac{1}{64}$ , by Chebyshev's  
 910 inequality, we have  $\sum_i \bar{q}_i \geq \frac{1}{2}$  with probability at least  $\frac{15}{16}$ . Provided  $\sum_i \bar{q}_i \geq \frac{1}{2}$ , then  
 911 the expected number of samples drawn (when, as described above, for each  $i$  we draw  
 912 a sample from from  $\text{Poi}(\bar{q}_i \cdot 2k)$ ) is at least  $k$ , and thus with probability at least  $\frac{1}{2}$ ,  
 913 at least  $k$  samples will be drawn. Thus via this Poisson process, with probability  $\frac{1}{2}$ ,  
 914 we have emulated drawing a sample of size  $k$  from a distribution that corresponds to  
 915  $Q_\epsilon$  at least  $\frac{15}{16}$  of the time.

916 Correspondingly, we emulate  $p$  by the simple Poisson process of drawing  $\text{Poi}(2k)$   
 917 samples from  $p$ , and throwing out all but  $k$  samples; there will be at least  $k$  samples  
 918 with probability greater than  $\frac{1}{2}$ .

919 Assume for the sake of contradiction that there is a hypothetical tester that could  
 920 distinguish  $p$  from  $Q_\epsilon$  in  $k$  samples with probability  $2/3$ , then this tester could be  
 921 used to distinguish the following two processes with probability  $\frac{1/2+2/3}{2} = \frac{7}{12}$ :

- 922 1. Draw  $\bar{q}_i = p_i \pm \epsilon_i$
- 923 (a) If  $\sum_i \bar{q}_i < \frac{1}{2}$  then with probability  $\frac{1}{2}$  output "FAIL" and with probability  
 924  $\frac{1}{2}$  output "Q"



- 925 (b) Otherwise, for each  $i$  generate a sample from  $Poi(\bar{q}_i \cdot 2k)$ ; if fewer than  
 926  $k$  total samples are generated, output “FAIL”, otherwise flip a biased  
 927 coin and either output a randomly chosen  $k$  of the generated samples, or  
 928 “FAIL” so that the total probability of outputting “FAIL” in this case  
 929 equals  $\frac{1}{2}$ .
- 930 2. Or, draw a sample of size  $Poi(2k)$  from  $p$ , and if fewer than  $k$  total samples  
 931 are generated, output “FAIL”, otherwise flip a biased coin and either output  
 932 a randomly chosen  $k$  of the generated samples, or “FAIL” so that the total  
 933 probability of outputting “FAIL” in this case equals  $\frac{1}{2}$ .

934 The tester is simulated on the samples if the chosen process above outputs sam-  
 935 ples, yielding an opinion “P” or “Q”; if the chosen process above outputs “FAIL”,  
 936 then a random one of “P” or “Q” is chosen; and if the (first) process outputs “Q”,  
 937 then this is output overall. This tester succeeds with probability at least the average  
 938 of  $\frac{1}{2}$  and  $\frac{2}{3}$ , since the above processes outputs “FAIL” with probability  $\frac{1}{2}$  yielding a  
 939 random guess about “P” or “Q”, and otherwise either generate a faithful sample from  
 940 the corresponding distribution, or in Case 1a outputs the answer directly, and is thus  
 941 at least as accurate as the  $\frac{2}{3}$ -accurate tester.

942 The same tester will perform within  $\frac{1}{32}$  of the success rate above if we remove  
 943 Case 1a and replace it with Case 1b, since this change affects the outcome only if  
 944  $\sum_i \bar{q}_i < \frac{1}{2}$  and simultaneously “FAIL” is not chosen, which happens with probability  
 945  $\frac{1}{16} \cdot \frac{1}{2} = \frac{1}{32}$ , yielding an accuracy at least  $\frac{7}{12} - \frac{1}{32} > \frac{1}{2}$ .

946 We thus derive a contradiction by showing that we cannot distinguish the fol-  
 947 lowing two processes with constant probability bounded above  $1/2$ : 1) for each  $i$ ,  
 948 draw a sample from  $Poi((p_i \pm \epsilon_i) \cdot 2k)$ ; versus 2) for each  $i$ , draw a sample from  
 949  $Poi(p_i \cdot 2k)$ . These two Poisson processes are both product distributions, and we can  
 950 thus compare them from the fact that the squared Hellinger distance is subadditive  
 951 on product distributions. For each component  $i$ , the squared Hellinger distance is  
 952  $H(Poi(kp_i), Poi(k[p_i \pm \epsilon_i]))^2$  which by Lemma 12 is at most  $c_1 k^2 \frac{\epsilon_i^4}{p_i^2}$ . Summing over  $i$

953 and taking the square root yields a bound on the Hellinger distance of  $k \left( c_1 \sum_i \frac{\epsilon_i^4}{p_i^2} \right)^{1/2}$ ,  
 954 which thus bounds the  $L_1$  distance. Thus when  $k$  satisfies the bound of the theorem,  
 955 the statistical distance between a set of  $k$  samples drawn from  $p$  versus drawn from a  
 956 random distribution of  $Q_\epsilon$  is bounded as  $O(c)$ , and thus for small enough constant  $c$   
 957 the two cannot be distinguished.

958 We now analyze the second part of the theorem, bounding the distance between  
 959 a distribution  $q \leftarrow Q_\epsilon$  and  $p$ . We note that the total excess probability mass in the  
 960 process of generating  $q$  that must subsequently be removed (or added, if it is negative)  
 961 by the normalization step is distributed as  $\sum_i \pm \epsilon_i$ , and thus by the triangle inequality,  
 962 the  $L_1$  distance between  $q$  and  $p$  is at least as large as a sample from  $\sum_i \epsilon_i - |\sum_i \pm \epsilon_i|$ .  
 963 We thus show that with probability at least  $1/2$ , a random value from  $|\sum_i \pm \epsilon_i|$  is at  
 964 most either  $\max_i \epsilon_i$  or  $\frac{1}{2} \sum_i \epsilon_i$ .

965 Consider the sequence  $\epsilon_i$  as sorted in descending order. We have two cases.  
 966 Suppose  $\epsilon_1 \geq \frac{1}{2} \sum_i \epsilon_i$ . Consider the random number  $|\sum_i \pm \epsilon_i|$ , where without loss of  
 967 generality the plus sign is chosen for  $\epsilon_1$ . With probability at least  $1/2$ , the sum of  
 968 the remaining elements will be  $\leq 0$ ; further, by the assumption of this case, this sum  
 969 cannot be smaller than  $-2\epsilon_1$ . Thus the sum of all the elements has magnitude at  
 970 most  $\epsilon_1$  with probability at least  $1/2$ .

971 In the other case,  $\epsilon_1 < \frac{1}{2} \sum_i \epsilon_i$ . Consider randomly choosing signs  $s_i \in \{-1, +1\}$   
 972 for the elements iteratively, stopping *before* choosing the sign for the first element

973  $j$  for which it would be possible for  $\left|(\sum_{i<j} s_i \epsilon_i) \pm \epsilon_j\right|$  to exceed  $\frac{1}{2} \sum_i \epsilon_i$ . Since  
 974 by assumption  $\epsilon_1 < \frac{1}{2} \sum_i \epsilon_i$ , we have  $j \geq 2$ . Without loss of generality, assume  
 975  $\sum_{i<j} s_i \epsilon_i \geq 0$ . We have  $\sum_{i<j} s_i \epsilon_i < \frac{1}{2} \sum_i \epsilon_i$ , and (by symmetry) with probabil-  
 976 ity at most  $1/2$  the sum of the remaining elements with randomly chosen signs will  
 977 be positive. Further, since  $s_1 \epsilon_1 + s_2 \epsilon_2 + \dots + s_{j-1} \epsilon_{j-1} + \epsilon_j \geq \frac{1}{2} \sum_i \epsilon_i$ , we have  
 978  $s_1 \epsilon_1 + s_2 \epsilon_2 + \dots + s_{j-1} \epsilon_{j-1} - \sum_{i \geq j} \epsilon_i \geq -\frac{1}{2} \sum_i \epsilon_i$ , for otherwise if this last inequality  
 979 was “ $<$ ” we could subtract these last two equations to conclude  $\epsilon_j + \sum_{i \geq j} \epsilon_i > \sum_i \epsilon_i$ ,  
 980 which contradicts the facts that  $s_1 \geq s_j$  and  $j \geq 2$ . Thus a random choice of the re-  
 981 maining signs starting with  $s_j$  will yield a total sum at most  $\frac{1}{2} \sum_i \epsilon_i$ , with probability  
 982 at least  $1/2$ , as desired.  $\square$

983 We apply this result as follows.

984 **COROLLARY 14.** *There is a constant  $c'$  such that for all probability distributions*  
 985  *$p$  and each  $\alpha > 0$ , there is no tester that, via a set of  $c' \cdot \left(\sum_{i \neq m} \frac{\min\{p_i, \alpha p_i^{2/3}\}^4}{p_i^2}\right)^{-1/2}$*   
 986 *samples can distinguish  $p$  from distributions with  $L_1$  distance  $\frac{1}{2} \sum_{i \neq m} \min\{p_i, \alpha p_i^{2/3}\}$*   
 987 *from  $p$  with probability 0.6, where  $m$  is the index of the element of  $p$  with maximum*  
 988 *probability.*

989 Note that for sufficiently small  $\alpha$ , the min is superfluous and the bound on  
 990 the number of samples becomes  $\frac{c'}{\alpha^2 \|p^{-\max}\|_{2/3}^{1/3}}$  and the  $L_1$  distance bound becomes  
 991  $\frac{1}{2} \alpha \|p^{-\max}\|_{2/3}^{2/3}$ , which more intuitively rephrases the result in terms of basic norms,  
 992 for this range of parameters.

993 *Proof.* Consider defining the vector of  $\epsilon_i$ 's by letting  $\epsilon_i = \min\{p_i, \alpha p_i^{2/3}\}$  for  
 994  $i \neq m$ , and  $\epsilon_m = \max_{i \neq m} \epsilon_i$ ; hence if the domain is sorted with  $p_1 \geq p_2 \geq \dots$ ,  
 995 then for  $i \geq 2$  we set  $\epsilon_i = \min\{p_i, \alpha p_i^{2/3}\}$ , and then set  $\epsilon_1 \epsilon_2$ . Theorem 13 yields  
 996 that  $p$  and  $Q_\epsilon$  cannot be distinguished given a set of  $\sqrt{2}c \cdot \left(\sum_{i \neq m} \frac{\min\{p_i, \alpha p_i^{2/3}\}^4}{p_i^2}\right)^{-1/2}$   
 997 samples where  $c$  is the constant from Theorem 13. Also from Theorem 13, with  
 998 probability at least  $1/2$ , the distance between  $p$  and an element of  $Q_\epsilon$  is at least the  
 999 min of  $\sum_{i \neq m} \min\{p_i, \alpha p_i^{2/3}\}$  and  $\frac{1}{2} \sum_i \min\{p_i, \alpha p_i^{2/3}\}$ , which we trivially bound by  
 1000  $\frac{1}{2} \sum_{i \neq m} \min\{p_i, \alpha p_i^{2/3}\}$ . We derive a contradiction as follows. If a tester with the  
 1001 parameters of this corollary existed, then repeating it a constant number of times  
 1002 and taking the majority output would amplify its success probability to at least 0.9;  
 1003 such a tester could be used to violate Theorem 13 via the procedure: given a set of  
 1004 samples drawn from either  $p$  or  $Q_\epsilon$ , run the tester, and if it outputs “ $Q_\epsilon$ ” then output  
 1005 “ $Q_\epsilon$ ”, and if it outputs “ $p$ ” then flip a coin and with probability 0.7 output “ $p$ ” and  
 1006 otherwise output “ $Q_\epsilon$ ”. If the distribution is  $p$  then our tester will correctly output  
 1007 this with  $0.9 \cdot 0.7 > 0.6$  probability. If the distribution was drawn from  $Q_\epsilon$  then with  
 1008 probability at least  $1/2$  the distribution will be far enough from  $p$  for the tester to  
 1009 apply (as noted above, by Theorem 13) and report this with probability 0.9; otherwise  
 1010 the tester will report “ $Q_\epsilon$ ” with probability at least  $1 - 0.7 = 0.3$ . Thus the tester will  
 1011 correctly report “ $Q_\epsilon$ ” with probability at least  $\frac{0.9+0.3}{2} = 0.6$  in all cases, the desired  
 1012 contradiction.  $\square$

1013 We now prove the lower bound portion of Theorem 2.

1014 **PROPOSITION 15.** *There exists a constant  $c_2$  such that for any  $\epsilon \in (0, 1)$  and any*  
 1015 *known distribution  $p$ , no tester can distinguish for an unknown distribution  $q$  whether*

1016  $q = p$  or  $\|p - q\|_1 \geq \epsilon$  with probability  $\geq 2/3$  when given a set of samples of size  
 1017  $c_2 \cdot \max \left\{ \frac{1}{\epsilon}, \frac{\|p - 2\epsilon\|_{2/3}^{\max}}{\epsilon^2} \right\}$ .

1018 *Proof.* We note, trivially, that the distributions of the vectors of  $k$  samples from  
 1019 two distributions that are  $\epsilon$  far apart are themselves at most  $k\epsilon$  far apart; thus for  
 1020 an appropriate constant  $c_2$ , at least  $c_2 \cdot \frac{1}{\epsilon}$  samples are needed to distinguish such  
 1021 distributions, showing the first part of our max bound.

1022 To show that the second term in the maximum is also a lower bound on the  
 1023 necessary sample size, we apply Corollary 14. Consider the probabilities  $p_i$  to be  
 1024 sorted in decreasing order, so that  $p_1$  is the maximum probability element. Define  $\alpha$   
 1025 to be the value which satisfies  $\frac{1}{2} \sum_{i \geq 2} \min\{p_i, \alpha p_i^{2/3}\} = \epsilon$ , and let  $s$  be the smallest  
 1026 integer such that  $\sum_{i > s} p_i \leq 2\epsilon$ . We note that for  $i \in \{2, \dots, s\}$  the min is never  
 1027  $p_i$ , or else (since  $p_i$  are sorted in descending order and the inequality  $p_i \leq \alpha p_i^{2/3}$  gets  
 1028 stronger for smaller  $p_i$ ), the sum would be at least  $\sum_{i > s} p_i$  which is greater than  $2\epsilon$  by  
 1029 definition of  $s$ . Thus  $\alpha \sum_{i=2}^s p_i^{2/3} = \sum_{i=2}^s \min\{p_i, \alpha p_i^{2/3}\} \leq \sum_{i \geq 2} \min\{p_i, \alpha p_i^{2/3}\} =$   
 1030  $2\epsilon$ , which yields  $\alpha \leq 2\|p_{\{2, \dots, s\}}\|_{2/3}^{-2/3} \epsilon$ . The lower bound on  $k$  from Corollary 14 is  
 1031 thus bounded (since the min of two quantities can only increase if we replace one  
 1032 by a weighted geometric mean of both of them) as  $c' \cdot \left( \sum_{i \geq 2} \frac{\min\{p_i, \alpha p_i^{2/3}\}^4}{p_i^2} \right)^{-1/2} =$   
 1033  $c' \cdot \left( \sum_{i \geq 2} \min\{p_i^2, \alpha^4 p_i^{2/3}\} \right)^{-1/2} \geq c' \cdot \left( \alpha^3 \sum_{i \geq 2} \min\{p_i, \alpha p_i^{2/3}\} \right)^{-1/2}$ . We bound this  
 1034 last expression by bounding  $\alpha^3$  by the cube of our bound  $\alpha \leq 2\|p_{\{2, \dots, s\}}\|_{2/3}^{-2/3} \epsilon$  and  
 1035 then plugging in the definition  $\frac{1}{2} \sum_{i \geq 2} \min\{p_i, \alpha p_i^{2/3}\} = \epsilon$  to yield a lower bound on  
 1036  $k$  of  $c' \cdot \left( 16\|p_{\{2, \dots, s\}}\|_{2/3}^{-2} \epsilon^4 \right)^{-1/2} = \frac{c'}{4} \cdot \frac{\|p_{\{2, \dots, s\}}\|_{2/3}}{\epsilon^2}$ . A constant number of repetitions  
 1037 lets us amplify the accuracy of the tester from the 0.6 of Corollary 14 to the 2/3 of  
 1038 this theorem.  $\square$

1039

## REFERENCES

- 1040 [1] J. ACHARYA, H. DAS, A. JAFARPOUR, A. ORLITSKY, AND S. PAN, *Competitive closeness testing*,  
 1041 in Conference on Learning Theory (COLT), 2011.  
 1042 [2] J. ACHARYA, H. DAS, A. JAFARPOUR, A. ORLITSKY, AND S. PAN, *Competitive classification and*  
 1043 *closeness testing*, Proc. 25th Conference on Learning Theory (COLT), 23 (2012), pp. 22.1–  
 1044 22.18.  
 1045 [3] Z. BAR-YOSSEF, *The Complexity of Massive Data Set Computations*, PhD thesis, Berkeley,  
 1046 CA, USA, 2002. AAI3183783.  
 1047 [4] Z. BAR-YOSSEF, R. KUMAR, AND D. SIVAKUMAR, *Sampling algorithms: lower bounds and*  
 1048 *applications*, in Symposium on Theory of Computing (STOC), 2001.  
 1049 [5] T. BATU, S. DASGUPTA, R. KUMAR, AND R. RUBINFELD, *The complexity of approximating the*  
 1050 *entropy*, SIAM Journal on Computing, (2005).  
 1051 [6] T. BATU, E. FISCHER, L. FORTNOW, R. KUMAR, R. RUBINFELD, AND P. WHITE, *Testing random*  
 1052 *variables for independence and identity*, in IEEE Symposium on Foundations of Computer  
 1053 Science (FOCS), 2001.  
 1054 [7] T. BATU, L. FORTNOW, R. RUBINFELD, W. D. SMITH, AND P. WHITE, *Testing closeness of*  
 1055 *discrete distributions*, J. ACM, 60 (2013), p. 4.  
 1056 [8] S. CHAN, I. DIAKONIKOLAS, G. VALIANT, AND P. VALIANT, *Optimal algorithms for testing close-*  
 1057 *ness of discrete distributions*, in Proceedings of the ACM-SIAM Symposium on Discrete  
 1058 Algorithms (SODA), 2014, pp. 1193–1203.  
 1059 [9] M. CHARIKAR, S. CHAUDHURI, R. MOTWANI, AND V. NARASAYYA, *Towards estimation error*  
 1060 *guarantees for distinct values*, in Symposium on Principles of Database Systems (PODS),  
 1061 2000.

- 1062 [10] I. DIAKONIKOLAS, D. M. KANE, AND V. NIKISHKIN, *Testing identity of structured distribu-*  
1063 *tions*, in Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete  
1064 Algorithms, SODA '15, Philadelphia, PA, USA, 2015, Society for Industrial and Applied  
1065 Mathematics, pp. 1841–1854, <http://dl.acm.org/citation.cfm?id=2722129.2722252>.
- 1066 [11] O. GOLDBREICH AND D. RON, *On testing expansion in bounded-degree graphs*, in Technical  
1067 Report TR00-020, Electronic Colloquium on Computational Complexity, 2000.
- 1068 [12] S. GUHA, A. MCGREGOR, AND S. VENKATASUBRAMANIAN, *Streaming and sublinear approxima-*  
1069 *tion of entropy and information distances*, in Proceedings of the ACM-SIAM Symposium  
1070 on Discrete Algorithms (SODA), 2006.
- 1071 [13] A. N. KOLMOGOROV, *On the empirical determination of a distribution law*, *Giornale*  
1072 *dell’Istituto Italiano degli Attuari*, 4 (1933), pp. 83–91.
- 1073 [14] M. MITZENMACHER AND E. UPFAL, *Probability and computing: Randomized algorithms and*  
1074 *probabilistic analysis*, Cambridge University Press, 2005.
- 1075 [15] L. PANINSKI, *Estimation of entropy and mutual information*, *Neural Computation*, 15 (2003),  
1076 pp. 1191–1253.
- 1077 [16] L. PANINSKI, *Estimating entropy on  $m$  bins given fewer than  $m$  samples*, *IEEE Trans. on*  
1078 *Information Theory*, 50 (2004), pp. 2200–2203.
- 1079 [17] L. PANINSKI, *A coincidence-based test for uniformity given very sparsely-sampled discrete data*,  
1080 *IEEE Transactions on Information Theory*, 54 (2008), pp. 4750–4755.
- 1081 [18] S. RASKHODNIKOVA, D. RON, A. SHPILKA, AND A. SMITH, *Strong lower bounds for approx-*  
1082 *imating distribution support size and the distinct elements problem*, *SIAM Journal on*  
1083 *Computing*, 39 (2009), pp. 813–842.
- 1084 [19] R. RUBINFELD, *Taming big probability distributions*, *XRDS*, 19 (2012), pp. 24–28.
- 1085 [20] G. VALIANT AND P. VALIANT, *Estimating the unseen: an  $n/\log(n)$ -sample estimator for en-*  
1086 *tropy and support size, shown optimal via new CLTs*, in Proceedings of the ACM Sympo-  
1087 sium on Theory of Computing (STOC), 2011.
- 1088 [21] G. VALIANT AND P. VALIANT, *The power of linear estimators*, in IEEE Symposium on Foun-  
1089 dations of Computer Science (FOCS), 2011.
- 1090 [22] G. VALIANT AND P. VALIANT, *Instance optimal learning of discrete distributions*, in Proceedings  
1091 of the ACM Symposium on Theory of Computing (STOC), 2016.
- 1092 [23] P. VALIANT, *Testing symmetric properties of distributions*, in Symposium on Theory of Com-  
1093 puting (STOC), 2008.