
A Scalable Tester for Samplers

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper we study the problem of testing of constrained samplers over high-
2 dimensional distributions with $(\varepsilon, \eta, \delta)$ guarantees. Samplers are increasingly used
3 in a wide range of safety-critical ML applications, and hence the testing problem
4 has gained importance. For n -dimensional distributions, the existing state-of-the-art
5 algorithm, Barbarik2, has a worst case query complexity of exponential in n and
6 hence is not ideal for use in practice. Our primary contribution is an exponentially
7 faster algorithm that has a query complexity linear in n and hence can easily scale
8 to larger instances. We demonstrate our claim by implementing our algorithm and
9 then comparing it against Barbarik2. Our experiments on the samplers wUnigen3
10 and wSTS, find that Pacoco requires $10\times$ fewer samples for wUnigen3 and $450\times$
11 fewer samples for wSTS as compared to Barbarik2.

12 1 Introduction

13 The constrained sampling problem is to draw samples from high-dimensional distributions over
14 constrained spaces. A constrained sampler $\mathcal{Q}(\varphi, \mathbf{w})$ takes in a set of constraints $\varphi : \{0, 1\}^n \rightarrow \{0, 1\}$
15 and a weight function $\mathbf{w} : \{0, 1\}^n \rightarrow \mathbb{R}_{>0}$, and returns a sample $\sigma \in \varphi^{-1}(1)$ with probability propor-
16 tional to $\mathbf{w}(\sigma)$. Constrained sampling is a core primitive of many statistical inference methods used
17 in ML, such as Sequential Monte Carlo[29], Markov Chain Monte Carlo(MCMC)[3, 9], Simulated
18 Annealing [4], and Variational Inference [24]. Sampling from real-world distributions is often compu-
19 tationally intractable, and hence, in practice, samplers are heuristical and lack theoretical guarantees.
20 For such samplers, it is an important problem to determine whether the sampled distribution is close
21 to the desired distribution, and this problem is known as *testing of samplers*. The problem was
22 formalised in [14, 27] as follows: Given access to a target distribution \mathcal{P} , a sampler $\mathcal{Q}(\varphi, \mathbf{w})$, and
23 three parameters $(\varepsilon, \eta, \delta)$, with probability at least $1 - \delta$, return (1) Accept if $d_\infty(\mathcal{P}, \mathcal{Q}(\varphi, \mathbf{w})) < \varepsilon$,
24 or (2) Reject if $d_{TV}(\mathcal{P}, \mathcal{Q}(\varphi, \mathbf{w})) > \eta$. Here d_{TV} is the total variation distance, d_∞ the multiplicative
25 distance, and ε, η , and δ are parameters for closeness, fairness and confidence respectively.

26 There is substantial interest in the testing problem due to the increasing use of ML systems in
27 real-world applications where safety is essential, such as medicine [2], transportation [8, 25], and
28 warfare [28]. For the ML systems that incorporate samplers, the typical testing approach has been
29 to show the convergence of the sampler with the target distribution via empirical tests that rely on
30 heuristics and do not provide any guarantees [19, 22, 31, 34]. In a recent work [27], a novel framework,
31 called Barbarik2, was proposed that could test a given sampler while providing $(\varepsilon, \eta, \delta)$ guarantees,
32 using $\tilde{O}\left(\frac{\text{tilt}(\mathcal{P})^2}{\eta(\eta-3\varepsilon)^3}\right)$ queries, where $\text{tilt}(\mathcal{P}) := \max_{\sigma_1, \sigma_2 \in \{0, 1\}^n} \frac{\mathcal{P}(\sigma_1)}{\mathcal{P}(\sigma_2)}$ for $\mathcal{P}(\sigma_2) > 0$. Since the $\text{tilt}(\mathcal{P})$
33 can be take arbitrary values, we observe that the query complexity can be prohibitively large¹. On
34 the other hand, the best known lower bound for the problem, derived from [30], is $\tilde{\Omega}\left(\frac{\sqrt{n/\log(n)}}{\eta^2}\right)$.

¹A simple modification reveals that in terms of n, η, ε , the bound is $\tilde{O}\left(\frac{4^n}{\eta(\eta-3\varepsilon)^3}\right)$

In this work, we take a step towards bridging this gap with our algorithm, Pacoco, that has a query complexity of $\tilde{O}\left(\frac{\sqrt{n} \log n}{(\eta - 11.6\varepsilon)\eta^3} + \frac{n}{\eta^2}\right)$, representing an exponential improvement over the state of the art.

To be of any real value, testing tools must be able to scale to larger instances. In the case of constrained samplers, the only existing testing tool, Barbarik2, is not scalable owing to its query complexity. The lack of scalability is illustrated by the following fact: product distributions are the simplest possible constrained distributions, and given a union of two n -dimensional product distributions, Barbarik2 requires more than 10^8 queries for $n > 30$. On the other hand, the query complexity of Pacoco scales linearly with n , the number of dimensions, thus making it more appropriate for practical use.

We implement Pacoco and compare it against Barbarik2 to determine their relative performance. In our experiments, we consider two sets of problems, (1) constrained sampling benchmarks, (2) scalable benchmarks and two constrained samplers wSTS and wUnigen3. We found that to complete the test Pacoco required at least $450\times$ fewer samples from wSTS and $10\times$ fewer samples from wUnigen3 as compared to Barbarik2. Moreover, Pacoco terminates with a result on at least $3\times$ more benchmarks than Barbarik2 in each experiment.

Our contributions can be summarized as follows:

1. For the problem of testing of samplers, we provide an exponential improvement in query complexity over the current state of the art test Barbarik2. Our test, Pacoco, makes a total of $\tilde{O}\left(\frac{\sqrt{n} \log n}{(\eta - 11.6\varepsilon)\eta^3} + \frac{n}{\eta^2}\right)$ queries, where \tilde{O} hides polylog factors of ε , η and δ .
2. We present an extensive empirical evaluation of Pacoco and contrast it with Barbarik2. The results indicate that Pacoco requires far fewer samples and terminates on more benchmarks when compared to Barbarik2.

We define the notation and discuss related work in Section 2. We then present the main contribution of the paper, the test Pacoco, and its proof of correctness in Section 3. We present our experimental findings in Section 4 and then we conclude the paper and discuss some open problems in Section 5. Due to space constraints, we defer some proofs and the full experimental results to the supplementary Section A and B respectively.

2 Notation and preliminaries

Probability distributions In this paper we deal with samplers that sample from discrete probability distributions over high-dimensional spaces. We consider the sample space to be the n -dimensional Boolean hypercube $\{0, 1\}^n$. A constrained sampler \mathcal{Q} takes in a set of constraints $\varphi : \{0, 1\}^n \rightarrow \{0, 1\}$ and a weight function $\mathbf{w} : \{0, 1\}^n \rightarrow \mathbb{R}_{>0}$, and samples from the distribution $\mathcal{Q}(\varphi, \mathbf{w})$ defined as

$$\mathcal{Q}(\varphi, \mathbf{w})(\sigma) = \begin{cases} \mathbf{w}(\sigma)/\mathbf{w}(\varphi) & \sigma \in \varphi^{-1}(1) \\ 0 & \sigma \in \varphi^{-1}(0) \end{cases}, \text{ where } \mathbf{w}(\varphi) = \sum_{\sigma \in \varphi^{-1}(1)} \mathbf{w}(\sigma). \text{ To improve readability,}$$

we use \mathcal{Q} to refer to the distribution $\mathcal{Q}(\varphi, \mathbf{w})$. For an element i , $\mathcal{D}(i)$ denotes its probability in distribution \mathcal{D} and $i \sim \mathcal{D}$ represents that i is sampled from \mathcal{D} . For any non-empty set $S \subseteq \{0, 1\}^n$, \mathcal{D}_S is the distribution \mathcal{D} conditioned on set S , and $\mathcal{D}(S)$ is the probability of S in \mathcal{D} i.e., $\mathcal{D}(S) = \sum_{i \in S} \mathcal{D}(i)$.

The total variation (TV) distance of two probability distributions \mathcal{D}_1 and \mathcal{D}_2 is defined as: $d_{TV}(\mathcal{D}_1, \mathcal{D}_2) = \frac{1}{2} \sum_{i \in \{0, 1\}^n} |\mathcal{D}_1(i) - \mathcal{D}_2(i)|$. For $S \subseteq \{0, 1\}^n$, we define $d_{TV(S)}(\mathcal{D}_1, \mathcal{D}_2) = \frac{1}{2} \sum_{i \in S} |\mathcal{D}_1(i) - \mathcal{D}_2(i)|$. The multiplicative distance of \mathcal{D}_2 from \mathcal{D}_1 is defined as: $d_\infty(\mathcal{D}_1, \mathcal{D}_2) = \max_{i \in \{0, 1\}^n} |\mathcal{D}_2(i)/\mathcal{D}_1(i) - 1|$. The two notions of distance obey the identity: $2d_{TV}(\mathcal{D}_1, \mathcal{D}_2) \leq d_\infty(\mathcal{D}_1, \mathcal{D}_2)$.

In the rest of the paper, $\mathbb{E}[v]$ represents the expectation of random variable v and $[k]$ represents the set $\{1, 2, \dots, k\}$.

Tools used in the analysis

Proposition 1 (Hoeffding). *For independent 0-1 random variables X_i , $X = \sum_{i=1}^k X_i$, and $t \geq 0$, $\Pr(X - \mathbb{E}X > t) \leq \exp(-2t^2/k)$ and $\Pr(\mathbb{E}X - X > t) \leq \exp(-2t^2/k)$*

82 **Proposition 2** (Chebyshev). *Given bounded r.v. X , we have $\Pr(|X - \mathbb{E}[X]| < \mathbb{E}[X]) > \frac{\mathbb{E}[X]^2}{\mathbb{E}[X^2]}$*

83 **Proposition 3.** *Given distributions \mathcal{D}_1 and \mathcal{D}_2 supported on $\{0, 1\}^n$, and a set $S \subseteq \{0, 1\}^n$,*

$$\sum_{i \in S} \mathcal{D}_1(i) \mathcal{D}_2(i) > \frac{(\mathcal{D}_1(S) + \mathcal{D}_2(S) - 2d_{TV(S)}(\mathcal{D}_1, \mathcal{D}_2))^2}{4|S|}$$

84 The proof can be found in the Appendix A.1 □

85 If we are given samples $\{s_1, s_2, \dots, s_n\}$ from a distribution \mathcal{D} over $[k]$, then the empirical distribution

86 $\widehat{\mathcal{D}}$ is defined to be $\widehat{\mathcal{D}}(i) = \frac{1}{n} \sum_{j=1}^k \mathbb{1}_{\{s_j=i\}}$.

87 **Proposition 4** (See [11] for a simple proof). *Suppose \mathcal{D} is a distribution over $[k]$, and $\widehat{\mathcal{D}}$ is constructed*
 88 *using $\max\left(\frac{k}{\eta^2}, \frac{2 \ln(2/\delta)}{\eta^2}\right)$ samples from \mathcal{D} . Then $d_{TV}(\mathcal{D}, \widehat{\mathcal{D}}) \leq \eta$ with probability at least $1 - \delta$.*

89 Testing with the help of oracles

90 In distribution testing, we are given samples from an unknown distribution \mathcal{P} over a large support
 91 $\{0, 1\}^n$, and the task is to test whether \mathcal{P} satisfies some property of interest. One of the important
 92 properties we care about is whether \mathcal{P} is close to another distribution \mathcal{Q} , and this subfield of testing is
 93 known as *closeness testing*. It was shown by Valiant and Batu et al. that the ability to draw samples
 94 from \mathcal{P} and \mathcal{Q} is not powerful enough, as at least $\Omega(2^{2n/3})$ samples are required to provide any sort
 95 of probabilistic guarantee for closeness testing. Since n is usually large, it was desirable to find tests
 96 that could solve the closeness testing problem using polynomially many samples in n .

97 Motivated by the above requirement, Canonne et al. and Chakraborty et al. introduced the *conditional*
 98 *sampling oracle* (COND), that is a more powerful way to access distributions. A COND oracle for
 99 distribution \mathcal{D} over $\{0, 1\}^n$ takes as input a set $S \subseteq \{0, 1\}^n$ with $\mathcal{D}(S) > 0$, and returns a sample
 100 $i \in S$ with probability $\mathcal{D}(i)/\mathcal{D}(S)$. It has been shown that the use of the COND oracle, and its
 101 variants, drastically reduces the sample complexity of many tasks in distribution testing [1, 21, 12, 15,
 102 6, 26, 7, 17, 13, 30] (see [10] for an extensive survey). In this paper, we consider the pair-conditioning
 103 (PCOND) oracle, which is a special case of the COND oracle with the restriction that $|S| = 2$ i.e.,
 104 the size of the conditioning set has to be two. To engineer practical PCOND oracle access into
 105 constrained samplers, we use the chain formula construction introduced in [14].

106 With the same goal of designing tests with polynomial sample complexity, a different kind of oracle,
 107 known as the DUAL oracle, was proposed by Canonne et al.. The DUAL oracle allows one to sample
 108 from a given distribution and also query the distribution for the probability of arbitrary elements of the
 109 support. Tractable DUAL oracle access is supported by a number of distribution representations, such
 110 as the fragments of probabilistic circuits (PC) that support the EVI query [18]. In our experimental
 111 evaluation, we use distributions from one such fragment: weighted d-DNNFs. Weighted d-DNNFs
 112 are a class of arithmetic circuits with properties that enable DUAL oracle access in time linear in the
 113 size of the circuit [16, 23].

114 3 Pacoco: an algorithm for testing samplers

115 We start by providing a brief overview of our testing algorithm before providing the full analysis.

116 3.1 Algorithm outline

117 The pseudocode of Pacoco is given in Algorithm 1. We adapt the definition of bucketing of distribu-
 118 tions from [30] for use in our analysis.

119 **Definition 1.** *For a given $k \in \mathbb{N}_{>0}$, the bucketing of $\{0, 1\}^n$ with respect to \mathcal{P} is defined as follows:*
 120 *For $1 \leq i \leq k$, let $S_i = \{b : 2^{-i} < \mathcal{P}(b) \leq 2^{-i+1}\}$ and let $S_0 = \{0, 1\}^n \setminus \bigcup_{i \in [k]} S_i$. Given*
 121 *any distribution \mathcal{D} over $\{0, 1\}^n$, we define a distribution $B_{\mathcal{D}}$ over $[k] \cup \{0\}$ as: for $0 \leq i \leq k$,*
 122 *$B_{\mathcal{D}}(i) = \mathcal{D}(S_i)$. We call $B_{\mathcal{D}}$ the bucket distribution of \mathcal{D} and S_i the i^{th} bucket.*

Algorithm 1 Pacoco($\mathcal{P}, \mathcal{Q}, \eta, \varepsilon, \delta$)

```
1:  $k \leftarrow n + \lceil \log_2(100/\eta) \rceil$ 
2: for  $i = 1$  to  $k$  do
3:    $S_i = \{b : 2^{-i} < \mathcal{P}(b) \leq 2^{-i+1}\}$ 
4:  $S_0 = \{0, 1\}^n \setminus \bigcup_{i \in [k]} S_i$ 
5:  $B_{\mathcal{P}}$  is the distribution over  $[k] \cup \{0\}$  where we sample  $i \sim B_{\mathcal{P}}$  if we sample  $j \sim \mathcal{P}$  and  $j \in S_i$ 
6:  $B_{\mathcal{Q}}$  is the distribution over  $[k] \cup \{0\}$  where we sample  $i \sim B_{\mathcal{Q}}$  if we sample  $j \sim \mathcal{Q}$  and  $j \in S_i$ 
7:  $\theta \leftarrow \eta/20$ 
8:  $\hat{d} \leftarrow \text{OutBucket}(B_{\mathcal{P}}, B_{\mathcal{Q}}, k, \theta, \delta/2)$ 
9: if  $\hat{d} > \varepsilon/2 + \theta$  then
10:   Return Reject
11:  $\varepsilon_2 \leftarrow \hat{d} + \theta$ 
12: Return InBucket( $\mathcal{P}, \mathcal{Q}, k, \varepsilon, \varepsilon_2, \eta, \delta/2$ )
```

123 Pacoco takes as input two distributions \mathcal{P} and \mathcal{Q} defined over the support $\{0, 1\}^n$, along with the
124 parameters for closeness(ε), fairness(η), and confidence(δ). On Line 1, Pacoco computes the value of
125 k using η and the number of dimensions n . Then, using DUAL access to \mathcal{P} , and SAMP access to \mathcal{Q} ,
126 Pacoco creates bucket distributions $B_{\mathcal{P}}$ and $B_{\mathcal{Q}}$ as in Defn. 1, in the following way: To sample from
127 $B_{\mathcal{P}}$, Pacoco first draws a sample $j \sim \mathcal{P}$, then using the DUAL oracle, determines the value of $\mathcal{P}(j)$.
128 Then, if j lies in the i^{th} bucket i.e., $2^{-i} < \mathcal{P}(j) \leq 2^{-i+1}$, the algorithm takes sample i as the sample
129 from $B_{\mathcal{P}}$. Similarly, to draw a sample from $B_{\mathcal{Q}}$, Pacoco draws a sample $j \sim \mathcal{Q}$ and then, using the
130 DUAL oracle to find $\mathcal{P}(j)$, finds i such that j lies in the i^{th} bucket, and then uses i as the sample.

131 Pacoco then calls two subroutines, OutBucket (Section 3.3) and InBucket (Section 3.2). The
132 OutBucket subroutine returns an θ -multiplicative estimate of the TV distance between $B_{\mathcal{P}}$ and
133 $B_{\mathcal{Q}}$, the two bucket distributions of \mathcal{P} and \mathcal{Q} , with an error of at most $\delta/2$. If it is found on Line 9 that
134 the estimate \hat{d} is greater than $\varepsilon/2 + \theta$, we know that $d_{TV}(\mathcal{P}, \mathcal{Q}) > \varepsilon/2$ and also that $d_{\infty}(\mathcal{P}, \mathcal{Q}) > \varepsilon$,
135 and hence the algorithm returns Reject. Otherwise, the algorithm calls the InBucket subroutine.

136 Now suppose that $d_{TV}(\mathcal{P}, \mathcal{Q}) \geq \eta$. Then, for ε_2 (Line 11), it is either the case that $d_{TV}(B_{\mathcal{P}}, B_{\mathcal{Q}}) >$
137 ε_2 or else $d_{TV}(B_{\mathcal{P}}, B_{\mathcal{Q}}) \leq \varepsilon_2$. In the former case, the algorithm returns Reject on Line 10, and in
138 the latter case the InBucket subroutine returns Reject. In both cases, the failure probability is at most
139 $\delta/2$. Thus Pacoco returns Reject on given η -far input distributions with probability at least $1 - \delta$.

140 We will now prove the following theorem:

141 **Theorem 1.** Pacoco($\mathcal{P}, \mathcal{Q}, \eta, \varepsilon, \delta$) takes in distributions \mathcal{P} and \mathcal{Q} defined over $\{0, 1\}^n$, and pa-
142 rameters $\eta \in (0, 1]$, $\varepsilon \in [0, \eta/11.6]$ and $\delta \in (0, 1/2]$. With probability at least $1 - \delta$, Pacoco
143 returns

- 144 • Accept if $d_{\infty}(\mathcal{P}, \mathcal{Q}) \leq \varepsilon$
- 145 • Reject if $d_{TV}(\mathcal{P}, \mathcal{Q}) > \eta$

146 Pacoco has query complexity $\tilde{O}\left(\frac{\sqrt{n} \log(n)}{\eta^3(\eta-11.6\varepsilon)} + \frac{n}{\eta^2}\right)$, where \tilde{O} hides polylog factors of ε, η and δ .

147 3.2 The InBucket subroutine

148 In this section, we present the InBucket subroutine, whose behavior is stated in the following lemma.

149 **Lemma 1.** InBucket($\mathcal{P}, \mathcal{Q}, k, \varepsilon, \varepsilon_2, \eta, \delta$) takes as input two distributions \mathcal{P}, \mathcal{Q} , an integer k and
150 parameters $\varepsilon, \varepsilon_2, \eta, \delta$. If $d_{\infty}(\mathcal{P}, \mathcal{Q}) \leq \varepsilon$, InBucket returns Accept. If $d_{TV}(\mathcal{P}, \mathcal{Q}) \geq \eta$ and
151 $d_{TV}(B_{\mathcal{P}}, B_{\mathcal{Q}}) < \varepsilon_2$, then InBucket returns Reject. InBucket errs with probability at most δ .

152 InBucket makes extensive use of the PCOND oracle access to \mathcal{Q} via the Bias subroutine, which we
153 describe in the following subsection.

154 **The Bias subroutine** The Bias subroutine takes in distribution \mathcal{Q} , two elements p, q and a positive
155 integer r . Then, using the PCOND oracle, Bias draws r samples from the conditional distribution
156 $\mathcal{Q}_{\{p, q\}}$ and returns the number of times it sees p in the r samples. It can be seen that the returned

Algorithm 2 $\text{InBucket}(\mathcal{P}, \mathcal{Q}, k, \varepsilon, \varepsilon_2, \eta, \delta)$

```
1:  $\varepsilon_1 \leftarrow (0.99\eta - 3.25\varepsilon_2 - 2\varepsilon/(1 - \varepsilon))/1.05 + 2\varepsilon/(1 - \varepsilon)$ 
2:  $m \leftarrow \lceil \sqrt{k}/(0.99\eta - 3.25\varepsilon_2 - \varepsilon_1) \rceil$ 
3:  $\alpha \leftarrow (\varepsilon_1 + 2\varepsilon/(1 - \varepsilon))/2$ 
4:  $t \leftarrow \left\lceil \frac{\ln(4/\delta)}{\ln(10/(10 - \varepsilon_1 + \alpha))} \right\rceil$ 
5: for  $t$  iterations do
6:    $\Gamma_{\mathcal{P}} \leftarrow m$  samples from  $\mathcal{P}$ 
7:    $\forall_{i \in [k]} \Gamma_{\mathcal{P}}^i \leftarrow \Gamma_{\mathcal{P}} \cap S_i$   $\triangleright S_i$  is defined in Defn. 1
8:    $\Gamma_{\mathcal{Q}} \leftarrow m$  samples from  $\mathcal{Q}$ 
9:    $\forall_{i \in [k]} \Gamma_{\mathcal{Q}}^i \leftarrow \Gamma_{\mathcal{Q}} \cap S_i$ 
10:  for all  $j \in [k]$  s.t.  $|\Gamma_{\mathcal{P}}^j|, |\Gamma_{\mathcal{Q}}^j| > 0$  do
11:     $p \leftarrow \Gamma_{\mathcal{P}}^j$   $\triangleright p$  is an arbitrary sample from the set  $\Gamma_{\mathcal{P}}^j$ 
12:     $q \leftarrow \Gamma_{\mathcal{Q}}^j$   $\triangleright q$  is an arbitrary sample from the set  $\Gamma_{\mathcal{Q}}^j$ 
13:     $h \leftarrow \frac{\mathcal{P}(p)}{\mathcal{P}(p) + \mathcal{P}(q)(1 + \frac{2\varepsilon}{1 - \varepsilon})}$ 
14:     $\ell \leftarrow \frac{\mathcal{P}(p)}{\mathcal{P}(p) + \mathcal{P}(q)(1 + \alpha)}$ 
15:     $r \leftarrow \left\lceil \frac{2 \ln(4mt/\delta)}{(h - \ell)^2} \right\rceil$ 
16:     $\widehat{c} \leftarrow \text{Bias}(\mathcal{Q}, p, q, r)$ 
17:    if  $\widehat{c} \leq (h + \ell)/2$  then
18:      Return Reject
19: Return Accept
```

Algorithm 3 $\text{Bias}(\mathcal{Q}, p, q, r)$

```
1: if  $p$  and  $q$  are identical then
2:   Return 0.5
3:  $\Gamma_{\mathcal{Q}_{\{p, q\}}} \leftarrow r$  samples from  $\mathcal{Q}_{\{p, q\}}$ 
4: Return # of times  $p$  appears in  $\Gamma_{\mathcal{Q}_{\{p, q\}}}$ 
```

157 value is an empirical estimate of $\frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)}$. Let the estimate be \widehat{c}_{pq} . We use the Hoeffding bound in
158 Prop. 1, and the value of r from Line 15 of Alg. (2) to show that:

$$\Pr \left[\frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} - \widehat{c}_{pq} \geq \frac{h - \ell}{2} \right] \leq \frac{\delta}{4mt} \quad \Pr \left[\widehat{c}_{pq} - \frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} \geq \frac{h - \ell}{2} \right] \leq \frac{\delta}{4mt}$$

159 Here t represents the number of iterations of the outer loop (Line 4), and m is the number of samples
160 drawn from $B_{\mathcal{P}}$ and $B_{\mathcal{Q}}$. Together, there are at most mt pairs of samples that are passed to the
161 Bias oracle. Since in each invocation of Bias, the probability of error is $\delta/4mt$, using the union
162 bound we find that the probability that all mt Bias calls return correctly is at least $1 - \delta/4$ and thus
163 with probability at least $1 - \delta/4$, the empirical estimate \widehat{c}_{pq} is closer than $(h - \ell)/4$ to $\frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)}$.
164 Henceforth we assume:

$$\left| \widehat{c}_{pq} - \frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} \right| \leq \frac{h - \ell}{2} \tag{1}$$

165 3.2.1 The Accept case

166 In this section we will provide an analysis of the case when $d_{\infty}(\mathcal{P}, \mathcal{Q}) < \varepsilon$. We will now state a
167 proposition required for the remaining proofs, the proof of which we relegate to Appendix A.4.

168 **Proposition 5.** *Let \mathcal{P}, \mathcal{Q} be distributions and let $p \sim \mathcal{P}$ and $q \sim \mathcal{Q}$. Then,*

169 1. *If $d_{\infty}(\mathcal{P}, \mathcal{Q}) < \varepsilon$ then*

$$\frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} \geq \frac{\mathcal{P}(p)}{\mathcal{P}(p) + (1 + \frac{2\varepsilon}{1 - \varepsilon})\mathcal{P}(q)}$$

170 2. If $d_{TV}(\mathcal{P}, \mathcal{Q}) > \varepsilon_1$, then for $0 \leq \alpha < \varepsilon_1$, with probability at least $(d_{TV}(\mathcal{P}, \mathcal{Q}) - \alpha)/2$,

$$\frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} < \frac{\mathcal{P}(p)}{\mathcal{P}(p) + (1 + \alpha)\mathcal{P}(q)}$$

171 From our assumption (1), we know that for all invocations of Bias, with probability at least $1 - \delta/4$,
 172 $\left| \widehat{c}_{pq} - \frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} \right| \leq (h - \ell)/2$. Using Prop. 5, and using the value of h given on Line 13, we can
 173 see that $\frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} > h$. From this we can observe that for all invocations of Bias, $\widehat{c}_{pq} > (h + \ell)/2$
 174 and the test does not return Reject in any iteration, hence eventually returning Accept. Thus, in the
 175 case that $d_\infty(\mathcal{P}, \mathcal{Q}) < \varepsilon$, the InBucket subroutine returns Accept with probability at least $1 - \delta/4$.

176 3.2.2 The Reject case

177 In this section we analyse the case when $d_{TV}(\mathcal{P}, \mathcal{Q}) \geq \eta$ and $d_{TV}(B_{\mathcal{P}}, B_{\mathcal{Q}}) \leq \varepsilon_2$ and we will show
 178 that the algorithm returns Reject with probability at least $1 - \delta$. For the purpose of the proof we will
 179 define a set of bad buckets $Bad \subseteq [k]$. Note that bucket $\{0\}$ is not in Bad .

180 **Definition 2.** $Bad = \{i \in [k] : d_{TV}(\mathcal{P}_{S_i}, \mathcal{Q}_{S_i}) > \varepsilon_1 \wedge B_{\mathcal{P}}(i)/B_{\mathcal{Q}}(i) \in [5^{-1}, 2]\}$

181 Suppose we have an indicator variable $X_{r,s}$ constructed as follows: draw m samples from \mathcal{P} and \mathcal{Q} ,
 182 and if the r^{th} sample from \mathcal{P} and the s^{th} sample from \mathcal{Q} both belong to some bucket $b \in Bad$, then
 183 $X_{r,s} = 1$ else $X_{r,s} = 0$. Then,

$$\mathbb{E}[X_{r,s}] = \sum_{b \in Bad} B_{\mathcal{P}}(b)B_{\mathcal{Q}}(b) > \frac{(B_{\mathcal{P}}(Bad) + B_{\mathcal{Q}}(Bad) - 2d_{TV}(Bad)(B_{\mathcal{P}}, B_{\mathcal{Q}}))^2}{4K}$$

184 The inequality is by the application of Prop. 3.

185 We analyse the expression for the expectation in the following lemma, the proof of which we relegate
 186 to Appendix A.2

Lemma 2.

$$B_{\mathcal{Q}}(Bad) + B_{\mathcal{P}}(Bad) - 2d_{TV}(Bad)(B_{\mathcal{Q}}, B_{\mathcal{P}}) > 2 \left(0.99\eta - \frac{13}{4}\varepsilon_2 - \varepsilon_1 \right)$$

187 Using Lemma 2 we immediately derive the fact that $\mathbb{E}[X_{r,s}] > (0.99\eta - \frac{13}{4}\varepsilon_2 - \varepsilon_1)^2 / K$. Let
 188 $X = \sum_{r,s \in [m]} X_{r,s}$. Given m samples from \mathcal{P} and \mathcal{Q} , $\Pr(X \geq 1)$ is the probability that there is at
 189 least one bucket in Bad that is sampled at least once each in both sets of samples.

190 **Lemma 3.** $\Pr(X \geq 1) > 1/5$

191 The proof can be found in Appendix A.3. □

192 Henceforth we will condition on the the event that $X \geq 1$. In such a case, we know that for some
 193 $k \in Bad$, there is a sample $p \sim \mathcal{P}_{S_k}$ and a sample $q \sim \mathcal{Q}_{S_k}$. Then for such a pair of samples (p, q) ,
 194 and some α , Prop. 5 tells us that with probability at least $(d_{TV}(\mathcal{P}, \mathcal{Q}) - \alpha)/2$ we have

$$\frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} < \frac{\mathcal{P}(p)}{\mathcal{P}(p) + (1 + \alpha)\mathcal{P}(q)}$$

195 Using the assumption made in (1), we immediately have that $\widehat{c}_{pq} \leq \frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} + \frac{h - \ell}{2}$. But from
 196 Prop. 5 we have that $\frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} < \ell$ and hence $\widehat{c}_{pq} < (h + \ell)/2$. Since $d_{TV}(\mathcal{P}, \mathcal{Q}) \geq \varepsilon_1$, we see that
 197 if $X \geq 1$, then with probability at least $(\varepsilon_1 - \alpha)/2$, the iteration returns Reject.

198 Then, using Lemma 3 we see that in every iteration, with probability at least $(\varepsilon_1 - \alpha)/10$, InBucket
 199 returns Reject. There are t iterations, where t (line 4) is chosen such that the overall probability of
 200 the test returning Reject is at least $1 - \delta/2$.

3.3 The OutBucket subroutine

The OutBucket subroutine takes as input two distributions $\mathcal{D}_1, \mathcal{D}_2$ over $k + 1$ elements and two parameters θ and δ . Then with probability at least $1 - \delta$, InBucket returns a θ -multiplicative estimate for $d_{TV}(\mathcal{D}_1, \mathcal{D}_2)$.

The OutBucket starts by drawing $\max\left(\frac{4(k+1)}{\theta^2}, \frac{8 \ln(4/\delta)}{\theta^2}\right)$ samples from the two distributions \mathcal{D}_1 and \mathcal{D}_2 , and constructs the empirical distributions $\widehat{\mathcal{D}}_1$ and $\widehat{\mathcal{D}}_2$. Then from Prop. 4, we know that with probability at least $1 - \delta$, both $d_{TV}(\mathcal{D}_1, \widehat{\mathcal{D}}_1) \leq \theta/2$ and $d_{TV}(\mathcal{D}_2, \widehat{\mathcal{D}}_2) \leq \theta/2$.

From the triangle inequality we have that,

$$d_{TV}(\widehat{\mathcal{D}}_1, \widehat{\mathcal{D}}_2) \leq d_{TV}(\mathcal{D}_1, \widehat{\mathcal{D}}_1) + d_{TV}(\mathcal{D}_2, \widehat{\mathcal{D}}_2) + d_{TV}(\mathcal{D}_1, \mathcal{D}_2) < \theta + d_{TV}(\mathcal{D}_1, \mathcal{D}_2)$$

and also that,

$$d_{TV}(\mathcal{D}_1, \mathcal{D}_2) \leq d_{TV}(\mathcal{D}_1, \widehat{\mathcal{D}}_1) + d_{TV}(\mathcal{D}_2, \widehat{\mathcal{D}}_2) + d_{TV}(\widehat{\mathcal{D}}_1, \widehat{\mathcal{D}}_2) < \theta + d_{TV}(\widehat{\mathcal{D}}_1, \widehat{\mathcal{D}}_2)$$

Thus with probability at least $1 - \delta$, the returned estimate $d_{TV}(\widehat{\mathcal{D}}_1, \widehat{\mathcal{D}}_2)$ satisfies $|d_{TV}(\widehat{\mathcal{D}}_1, \widehat{\mathcal{D}}_2) - d_{TV}(\mathcal{D}_1, \mathcal{D}_2)| < \theta$.

Query and runtime complexity The number of queries made by OutBucket to \mathcal{P} and \mathcal{Q} is given by $\tilde{O}\left(\frac{n}{\eta^2}\right)$, where \tilde{O} hides polylog factors of ε, η and δ . The number of queries required by InBucket is given by mtr . Bounding the terms individually, we see that $m = \tilde{O}\left(\frac{\sqrt{n}}{\eta - 11.6\varepsilon}\right)$, $t = \tilde{O}\left(\frac{1}{\eta}\right)$ and $r = \tilde{O}\left(\frac{\log n}{\eta^2}\right)$. Thus $mtr = \tilde{O}\left(\frac{\sqrt{n} \log n}{(\eta - 11.6\varepsilon)\eta^3}\right)$ and hence the total query complexity is $\tilde{O}\left(\frac{\sqrt{n} \log n}{(\eta - 11.6\varepsilon)\eta^3} + \frac{n}{\eta^2}\right)$.

4 Evaluation

To evaluate the performance of Pacoco and test the quality of publicly available samplers, we implemented Pacoco in Python. Our evaluation took inspiration from the experiments presented in previous work [14, 27], and we used the same framework to evaluate our proposed algorithm. The role of target distribution \mathcal{P} was played by the exact constrained sampler WAPS² [23]. For the role of sampler $\mathcal{Q}(\varphi, \mathbf{w})$, we used the state-of-the-art samplers wSTS and wUnigen3. wUnigen3 [32] is a hashing-based sampler that provides (ε, δ) guarantees on the quality of the samples. wSTS [20] is a sampler designed for sampling over challenging domains such as energy barriers and highly asymmetric spaces. wSTS generates samples much faster than wUnigen3, albeit without any guarantees on the quality of the samples.

For the closeness(ε), fairness(η), and confidence(δ) parameters, we choose the values 0.05, 0.9 and 0.2. This setting implies that for a given distribution \mathcal{P} , and for a given sampler $\mathcal{Q}(\varphi, \mathbf{w})$, Pacoco returns (1) Accept if $d_\infty(\mathcal{P}, \mathcal{Q}(\varphi, \mathbf{w})) < 0.05$, and (2) Reject if $d_{TV}(\mathcal{P}, \mathcal{Q}(\varphi, \mathbf{w})) > 0.9$, with probability at least 0.8. Our empirical evaluation sought to answer the question: How does the performance of Pacoco compare with the state-of-the-art tester Barbarik2?

Our experiments were conducted on a high-performance compute cluster with Intel Xeon(R) E5-2690v3@2.60GHz CPU cores. We use a single core with 4GB memory with a timeout of 16 hours for each benchmark. We set a sample limit of 10^8 samples for our experiments due to our limited computational resources.

4.1 Setting A - scalable benchmarks

Dataset Our dataset consists of the union of two n -dimensional product distributions, for $n \in \{4, 7, 10, \dots, 118\}$. We have 39 problems in the dataset. We represent the union of two product distributions as the constraint: $\varphi(\sigma) = \bigwedge_{i=1}^{2k} (\sigma_{3k+1} \vee \sigma_i) \wedge \bigwedge_{i=2k+1}^{3k} (\neg \sigma_{3k+1} \vee \sigma_i)$, and the weight function: $\mathbf{w}(\sigma) = \prod_{i=2k+1}^{3k} 3^{\sigma_i}$, where σ_i is the value of σ at position i .

²<https://github.com/meelgroup/WAPS>

Results We observe that in the case of wSTS, Barbarik2 can handle only 12 instances within the sample limit of 10^8 . On the other hand, Pacoco can handle all 39 instances using at the most 10^6 samples. In the case of wUnigen3, Barbarik2 solves 5 instances, and Pacoco can handle 17 instances.

Figure 1 shows a cactus plot comparing the sample requirement of Pacoco and Barbarik2. The x -axis represents the number of benchmarks and y -axis represents the number of samples, a point (x, y) implies that the relevant tester took less than y number of samples to distinguish between $d_{TV}(\mathcal{P}, \mathcal{Q}(\varphi, \mathbf{w})) > \eta$ and $d_{\infty}(\mathcal{P}, \mathcal{Q}(\varphi, \mathbf{w})) < \varepsilon$, for x many benchmarks. We display the set of benchmarks for which at least one of the two tools terminated within the sample limit of 10^8 . We want to highlight that the y -axis is in log-scale, thus showing the sample efficiency of Pacoco compared to Barbarik2. For every benchmark, we compute the ratio of the number of samples required by Barbarik2 to test a sampler and the number of samples required by Pacoco. The geometric mean of these ratios indicates the mean speedup. We find that the Pacoco’s speedup on wSTS is $451\times$ and on wUnigen3 is $10\times$.

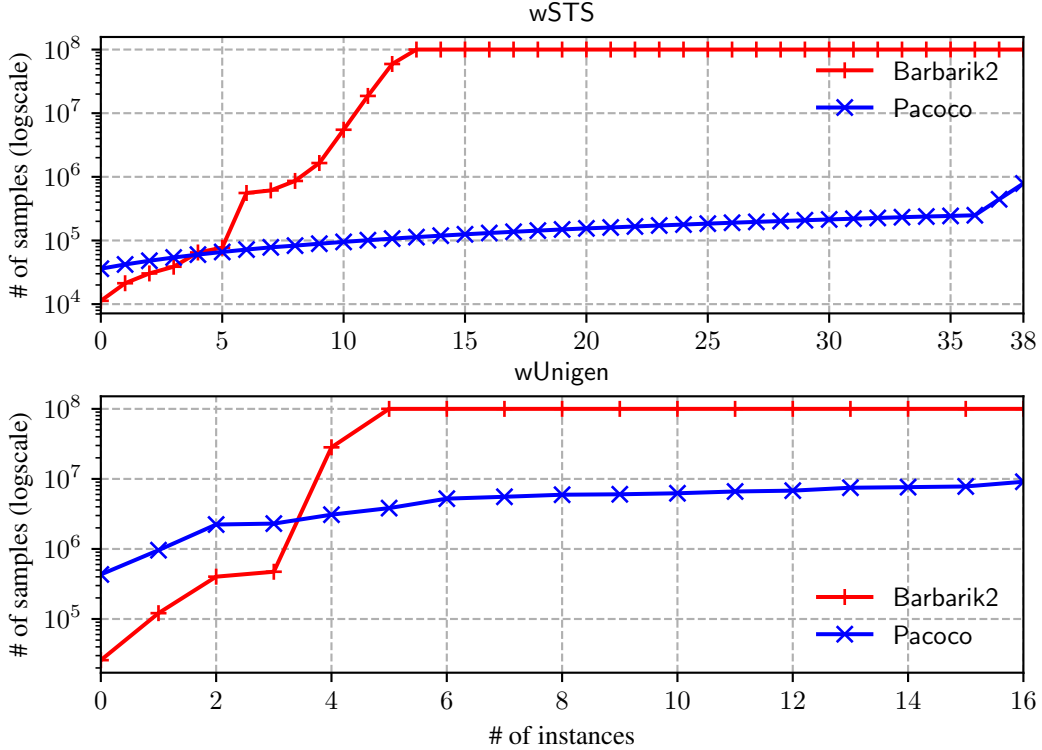


Figure 1: Cactus plot: Pacoco vs. Barbarik2. We set the sample limit to be 10^8 , and our dataset consists of 39 benchmarks. The plot shows all the instances where at least one of the two tools terminated within the time limit of 16 hours and sample limit of 10^8 .

4.2 Setting B - real-life benchmarks

Dataset We experiment on 87 constraints drawn from a collection of publicly available benchmarks arising from sampling and counting tasks³. We use distributions from the log-linear family. In a log-linear distribution, the probability of an element $\sigma \in \varphi^{-1}(1)$ is given as: $\Pr[\sigma] \propto \exp(\sum_{i=1}^n \sigma_i \theta_i)$, where $\theta_i \in \mathbb{R}_{\geq 0}^n$. We found that wUnigen3 was not able to sample from most of the benchmarks in the dataset within the given time limit, and hence we present the results only for wSTS.

Results We find that Pacoco terminated with a result on all 87 instances from the set of real-life benchmarks, while Barbarik2 could only terminate on 16. We present the results of our experiments

³<https://zenodo.org/record/3793090>

in Table 1. The first column indicates the benchmark’s name, and the second column has the number of dimensions of the space the distribution is defined on. The third and fifth columns indicate the number of samples required by Barbarik2 and Pacoco. The fourth and sixth columns report the output of Barbarik2 and Pacoco.

Table 1: Runtime performance of Pacoco. We experiment with 87 benchmarks, and out of the 87 benchmarks we display 15 in the table and we display the full data in Appendix B. In the table ‘A’ represents Accept, ‘R’ represents Reject and ‘TO’ represents that the tester either asked for more than 10^8 samples or did not terminate in the given time limit of 16 hours.

Benchmark	Dimensions	Barbarik2		Pacoco	
		Result	# of samples	Result	# of samples
SetTest.sk_9_21	21	R	2817	R	58000
Pollard.sk_1_10	10	R	7606	R	36000
s444_3_2	24	R	848148	R	64000
s526a_3_2	24	R	848148	R	64000
s510_15_7	25	R	12708989	R	66000
s27_new_7_4	7	A	23997012	R	30000
s298_15_7	17	R	38126967	R	50000
s420_3_2	34	TO	-	R	83000
s382_3_2	24	TO	-	R	64000
s641_3_2	54	TO	-	R	123000
111.sk_2_36	36	TO	-	R	87000
7.sk_4_50	50	TO	-	R	115000
56.sk_6_38	38	TO	-	R	91000
s820a_15_7	23	TO	-	R	62000
ProjectService3.sk_12_55	55	TO	-	R	125000

5 Conclusion

In this paper, we studied the problem of testing constrained samplers over high-dimensional distributions with $(\varepsilon, \eta, \delta)$ guarantees. For n -dimensional distributions, the existing state-of-the-art testing algorithm, Barbarik2, has a worst-case query complexity that is exponential in n and hence is not ideal for use in practice. We provided an exponentially faster algorithm, Pacoco, that has a query complexity linear in n and hence can easily scale to larger instances. We implemented Pacoco and tested the samplers wSTS and wUnigen3 to determine their sample complexity in practice. The results demonstrate that Pacoco is significantly more sample efficient than Barbarik2, requiring $450\times$ fewer samples when it tested wSTS and $10\times$ fewer samples when it tested wUnigen3. Since there is a \sqrt{n} gap between the upper bound provided by our work and the lower bound shown in [30], the problem of designing a more sample efficient algorithm or finding a stronger lower bound, remains open.

Limitations For a given fairness parameter η , Pacoco requires the value of the closeness parameter ε to lie in the interval $[0, \eta/11.6)$. In the case of Barbarik2, the previous state-of-the-art test, the permissible values of ε for a given η lie in the interval $[0, \eta/3)$. Thus, Pacoco supports testing with only a subset of parameter values that Barbarik2 supports.

References

- [1] Jayadev Acharya, Clément L. Canonne, and Gautam Kamath. A chasm between identity and equivalence testing with conditional queries. *Electron. Colloquium Comput. Complex.*, 2014.
- [2] Filippo Amato, Alberto López, Eladia María Peña-Méndez, Petr Vaňhara, Aleš Hampl, and Josef Havel. Artificial neural networks in medical diagnosis, 2013.
- [3] Christophe Andrieu, Nando De Freitas, Arnaud Doucet, and Michael I Jordan. An introduction to mcmc for machine learning. *Machine learning*, 2003.
- [4] Christophe Andrieu, Nando De Freitas, and Arnaud Doucet. Reversible jump mcmc simulated annealing for neural networks. *arXiv preprint arXiv:1301.3833*, 2013.
- [5] Tuğkan Batu, Lance Fortnow, Ronitt Rubinfeld, Warren D Smith, and Patrick White. Testing closeness of discrete distributions. *Journal of the ACM (JACM)*, 2013.
- [6] Rishiraj Bhattacharyya and Sourav Chakraborty. Property testing of joint distributions using conditional samples. *ACM Transactions on Computation Theory (TOCT)*, 2018.
- [7] Eric Blais, Clément L Canonne, and Tom Gur. Distribution testing lower bounds via reductions from communication complexity. *ACM Transactions on Computation Theory (TOCT)*, 2019.
- [8] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016.
- [9] Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng. *Handbook of markov chain monte carlo*. CRC press, 2011.
- [10] Clément L Canonne. A survey on distribution testing: Your data is big. but is it blue? *Theory of Computing*, 2020.
- [11] Clément L Canonne. A short note on learning discrete distributions. *arXiv preprint arXiv:2002.11457*, 2020.
- [12] Clément L Canonne, Dana Ron, and Rocco A Servedio. Testing probability distributions using conditional samples. *SIAM Journal on Computing*, 2015.
- [13] Clément L Canonne, Xi Chen, Gautam Kamath, Amit Levi, and Erik Waingarten. Random restrictions of high dimensional distributions and uniformity testing with subcube conditioning. In *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms (SODA)*. SIAM, 2021.
- [14] Sourav Chakraborty and Kuldeep S Meel. On testing of uniform samplers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019.
- [15] Sourav Chakraborty, Eldar Fischer, Yonatan Goldhirsh, and Arie Matsliah. On the power of conditional samples in distribution testing. *SIAM Journal on Computing*, 2016.
- [16] Mark Chavira and Adnan Darwiche. On probabilistic inference by weighted model counting. *Artificial Intelligence*, 2008.
- [17] Xi Chen, Rajesh Jayaram, Amit Levi, and Erik Waingarten. Learning and testing junta distributions with sub cube conditioning. In *Conference on Learning Theory*. PMLR, 2021.
- [18] YooJung Choi, Antonio Vergari, and Guy Van den Broeck. Probabilistic circuits: A unifying framework for tractable probabilistic models. 2020.
- [19] Mary Kathryn Cowles and Bradley P. Carlin. Markov chain monte carlo convergence diagnostics: A comparative review. *Journal of the American Statistical Association*, 1996.
- [20] Stefano Ermon, Carla P. Gomes, and Bart Selman. Uniform solution sampling using a constraint solver as an oracle. In *UAI*, 2012.

- [21] Moein Falahatgar, Ashkan Jafarpour, Alon Orlitsky, Venkatadheeraj Pichapati, and Ananda Theertha Suresh. Faster algorithms for testing under conditional sampling. In *Conference on Learning Theory*. PMLR, 2015.
- [22] Andrew Gelman and Donald B. Rubin. Inference from iterative simulation using multiple sequences. *Statistical Science*, 1992.
- [23] Rahul Gupta, Shubham Sharma, Subhajit Roy, and Kuldeep S Meel. Waps: Weighted and projected sampling. In *TACAS*, 2019.
- [24] Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational inference. *Journal of Machine Learning Research*, 2013.
- [25] Kyle D Julian, Mykel J Kochenderfer, and Michael P Owen. Deep neural network compression for aircraft collision avoidance systems. *Journal of Guidance, Control, and Dynamics*, 2019.
- [26] Gautam Kamath and Christos Tzamos. Anaconda: A non-adaptive conditional sampling algorithm for distribution testing. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM, 2019.
- [27] Kuldeep S. Meel^(r), Yash Pote^(r), and Sourav Chakraborty. On testing of samplers. In *Proceedings of Advances in Neural Information Processing Systems(NeurIPS)*, 12 2020.
- [28] Forrest E. Morgan, Benjamin Boudreaux, Andrew J. Lohn, Mark Ashby, Christian Curriden, Kelly Klima, and Derek Grossman. *Military Applications of Artificial Intelligence: Ethical Concerns in an Uncertain World*. RAND Corporation, Santa Monica, CA, 2020.
- [29] Christian A Naesseth, Fredrik Lindsten, and Thomas B Schön. Elements of sequential monte carlo. *arXiv preprint arXiv:1903.04797*, 2019.
- [30] Shyam Narayanan. On tolerant distribution testing in the conditional sampling model. In *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms (SODA)*. SIAM, 2021.
- [31] Christian P Robert, George Casella, and George Casella. *Monte Carlo statistical methods*, volume 2. Springer, 1999.
- [32] Mate Soos, Stephan Gocht, and Kuldeep S Meel. Tinted, detached, and lazy cnf-xor solving and its applications to counting and sampling. In *International Conference on Computer Aided Verification*, 2020.
- [33] Paul Valiant. Testing symmetric properties of distributions. *SIAM Journal on Computing*, 40, 2011.
- [34] Aki Vehtari, Andrew Gelman, Daniel Simpson, Bob Carpenter, and Paul-Christian Bürkner. Rank-normalization, folding, and localization: An improved \hat{R} for assessing convergence of MCMC (with Discussion). *Bayesian analysis*, 2021.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
 - (b) Did you describe the limitations of your work? [\[Yes\]](#)
 - (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#)
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#)
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#)
3. If you ran experiments...

- 370 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 371 mental results (either in the supplemental material or as a URL)? [Yes]
 372 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 373 were chosen)? [Yes]
 374 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 375 ments multiple times)? [N/A]
 376 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 377 of GPUs, internal cluster, or cloud provider)? [Yes]
 378 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 379 (a) If your work uses existing assets, did you cite the creators? [Yes]
 380 (b) Did you mention the license of the assets? [Yes]
 381 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
 382 (d) Did you discuss whether and how consent was obtained from people whose data you're
 383 using/curating? [N/A]
 384 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 385 information or offensive content? [N/A]
 386 5. If you used crowdsourcing or conducted research with human subjects...
 387 (a) Did you include the full text of instructions given to participants and screenshots, if
 388 applicable? [N/A]
 389 (b) Did you describe any potential participant risks, with links to Institutional Review
 390 Board (IRB) approvals, if applicable? [N/A]
 391 (c) Did you include the estimated hourly wage paid to participants and the total amount
 392 spent on participant compensation? [N/A]

393 A Missing proofs and algorithm

394 A.1 Proof of Lemma 3

395 *Proof.* The Hellinger distance of distributions \mathcal{P}, \mathcal{Q} restricted to a set $S \subseteq \{0, 1\}^n$, is defined as

$$\begin{aligned}
 396 \quad d_{H(S)}(\mathcal{P}, \mathcal{Q}) &= \frac{1}{\sqrt{2}} \sqrt{\sum_{i \in S} (\sqrt{\mathcal{Q}(i)} - \sqrt{\mathcal{P}(i)})^2}, \\
 d_{H(S)}(\mathcal{P}, \mathcal{Q}) &= \frac{1}{\sqrt{2}} \sqrt{\sum_{i \in S} (\sqrt{\mathcal{Q}(i)} - \sqrt{\mathcal{P}(i)})^2} \\
 d_{H(S)}^2(\mathcal{P}, \mathcal{Q}) &= \frac{1}{2} \sum_{i \in S} (\sqrt{\mathcal{Q}(i)} - \sqrt{\mathcal{P}(i)})^2 \\
 &= \frac{1}{2} \sum_{i \in S} (\mathcal{Q}(i) + \mathcal{P}(i) - 2\sqrt{\mathcal{P}(i)\mathcal{Q}(i)}) \\
 &= \frac{\mathcal{P}(S) + \mathcal{Q}(S)}{2} - \sum_{i \in S} \sqrt{\mathcal{P}(i)\mathcal{Q}(i)}
 \end{aligned}$$

397 Then using the fact that $d_{H(S)}^2(\mathcal{P}, \mathcal{Q}) \leq d_{TV(S)}(\mathcal{P}, \mathcal{Q})$ we see that, $\sum_{i \in S} \sqrt{\mathcal{P}(i)\mathcal{Q}(i)} \geq$
 398 $\frac{\mathcal{P}(S) + \mathcal{Q}(S)}{2} - d_{TV(S)}(\mathcal{P}, \mathcal{Q})$. Then we use the Cauchy-Schwarz inequality:

$$\sum_{i \in S} \mathcal{P}(i)\mathcal{Q}(i) \geq \frac{(\mathcal{P}(S) + \mathcal{Q}(S) - 2d_{TV(S)}(\mathcal{P}, \mathcal{Q}))^2}{4|S|}$$

399

□

400 A.2 Proof of Lemma 2

Lemma 2.

$$B_{\mathcal{Q}}(Bad) + B_{\mathcal{P}}(Bad) - 2d_{TV(Bad)}(B_{\mathcal{Q}}, B_{\mathcal{P}}) > 2 \left(0.99\eta - \frac{13}{4}\varepsilon_2 - \varepsilon_1 \right)$$

401 *Proof.* Let \mathcal{PQ} be a distribution constructed from \mathcal{P} and \mathcal{Q} , where we first sample $j \sim B_{\mathcal{Q}}$ and then
 402 sample $i \sim \mathcal{P}_{S_j}$, thus $\mathcal{PQ}(i) = \sum_{j \in [k] \cup \{0\}} B_{\mathcal{Q}}(j) \mathcal{P}_{S_j}(i)$. We know that if $i \in S_j$, then $i \notin S_{j'}$ for
 403 $j' \neq j$. This allows us to simplify and write $\mathcal{PQ}(i) = B_{\mathcal{Q}}(j) \mathcal{P}_{S_j}(i)$. Then,

$$\begin{aligned}
 d_{TV}(B_{\mathcal{P}}, B_{\mathcal{Q}}) &= \frac{1}{2} \sum_{j \in [k] \cup \{0\}} |B_{\mathcal{P}}(j) - B_{\mathcal{Q}}(j)| \\
 &= \frac{1}{2} \sum_{j \in [k] \cup \{0\}} \sum_{i \in S_j} \mathcal{P}_{S_j}(i) |B_{\mathcal{P}}(j) - B_{\mathcal{Q}}(j)| \\
 &= \frac{1}{2} \sum_{j \in [k] \cup \{0\}} \sum_{i \in S_j} |\mathcal{P}(i) - \mathcal{PQ}(i)| \\
 &= \frac{1}{2} \sum_{i \in \{0,1\}^n} |\mathcal{P}(i) - \mathcal{PQ}(i)| = d_{TV}(\mathcal{P}, \mathcal{PQ})
 \end{aligned}$$

404 Since $d_{TV}(B_{\mathcal{P}}, B_{\mathcal{Q}}) < \varepsilon_2$, we have $d_{TV}(\mathcal{P}, \mathcal{PQ}) < \varepsilon_2$.

405 From the definition of TV, we have

$$\begin{aligned}
d_{TV}(\mathcal{Q}, \mathcal{P}\mathcal{Q}) &= \frac{1}{2} \sum_{i \in \{0,1\}^n} |\mathcal{Q}(i) - \mathcal{P}\mathcal{Q}(i)| \\
&= \frac{1}{2} \sum_{j \in [k] \cup \{0\}} \sum_{i \in S_j} |\mathcal{Q}(i) - \mathcal{P}\mathcal{Q}(i)| \\
&= \frac{1}{2} \sum_{j \in [k] \cup \{0\}} \sum_{i \in S_j} |B_{\mathcal{Q}}(j) \mathcal{Q}_{S_j}(i) - B_{\mathcal{Q}}(j) \mathcal{P}_{S_j}(i)| \\
&= \frac{1}{2} \sum_{j \in [k] \cup \{0\}} B_{\mathcal{Q}}(j) \sum_{i \in S_j} |\mathcal{Q}_{S_j}(i) - \mathcal{P}_{S_j}(i)| \\
&= \sum_{j \in ([k] \cup \{0\})} B_{\mathcal{Q}}(j) d_{TV}(\mathcal{P}_{S_j}, \mathcal{Q}_{S_j}) \\
&= \sum_{j \in ([k] \cup \{0\}) \setminus \text{Bad}} B_{\mathcal{Q}}(j) d_{TV}(\mathcal{P}_{S_j}, \mathcal{Q}_{S_j}) + \sum_{j \in \text{Bad}} B_{\mathcal{Q}}(j) d_{TV}(\mathcal{P}_{S_j}, \mathcal{Q}_{S_j})
\end{aligned}$$

406 We will need the following sets: $R_1 = \{j : B_{\mathcal{P}}(j) > 2B_{\mathcal{Q}}(j)\}$, $R_2 = \{j : B_{\mathcal{Q}}(j) > 5B_{\mathcal{P}}(j)\}$.

407 From the triangle inequality we have $d_{TV}(\mathcal{P}, \mathcal{Q}) \leq d_{TV}(\mathcal{P}, \mathcal{P}\mathcal{Q}) + d_{TV}(\mathcal{P}\mathcal{Q}, \mathcal{Q})$. We also know
408 that $d_{TV}(\mathcal{P}, \mathcal{P}\mathcal{Q}) < \varepsilon_2$ and $d_{TV}(\mathcal{P}, \mathcal{Q}) > \eta$. Thus we have:

$$\begin{aligned}
\eta - \varepsilon_2 &< d_{TV}(\mathcal{Q}, \mathcal{P}\mathcal{Q}) \\
\eta - \varepsilon_2 &< \sum_{j \in ([k] \cup \{0\}) \setminus \text{Bad}} B_{\mathcal{Q}}(j) d_{TV}(\mathcal{P}_{S_j}, \mathcal{Q}_{S_j}) + \sum_{j \in \text{Bad}} B_{\mathcal{Q}}(j) d_{TV}(\mathcal{P}_{S_j}, \mathcal{Q}_{S_j}) \\
\eta - \varepsilon_2 &< \sum_{j \in \{0\} \cup R_1 \cup R_2} B_{\mathcal{Q}}(j) d_{TV}(\mathcal{P}_{S_j}, \mathcal{Q}_{S_j}) + \sum_{j \in [k] \setminus \{R_1 \cup R_2 \cup \text{Bad}\}} B_{\mathcal{Q}}(j) d_{TV}(\mathcal{P}_{S_j}, \mathcal{Q}_{S_j}) \\
&\quad + \sum_{j \in \text{Bad}} B_{\mathcal{Q}}(j) d_{TV}(\mathcal{P}_{S_j}, \mathcal{Q}_{S_j})
\end{aligned}$$

409 By definition, if $j \in [k] \setminus \{ \text{Bad} \cup R_1 \cup R_2 \}$, then j has the property that $d_{TV}(\mathcal{P}_{S_j}, \mathcal{Q}_{S_j}) \leq \varepsilon_1$. Then,

$$\begin{aligned}
\eta - \varepsilon_2 &< \sum_{j \in \{0\} \cup R_1 \cup R_2} B_{\mathcal{Q}}(j) + \sum_{j \in [k] \setminus \{R_1 \cup R_2 \cup \text{Bad}\}} B_{\mathcal{Q}}(j) \varepsilon_1 + \sum_{j \in \text{Bad}} B_{\mathcal{Q}}(j) \\
\eta - \varepsilon_2 - \varepsilon_1 &< B_{\mathcal{Q}}(\{0\} \cup R_1 \cup R_2) + B_{\mathcal{Q}}(\text{Bad}) \\
\eta - \varepsilon_2 - \varepsilon_1 - B_{\mathcal{Q}}(\{0\} \cup R_1 \cup R_2) &< B_{\mathcal{Q}}(\text{Bad}) \tag{2}
\end{aligned}$$

410 If $i \in R_1$, then $B_{\mathcal{P}}(i) > 2B_{\mathcal{Q}}(i)$, and thus $B_{\mathcal{P}}(i) - B_{\mathcal{Q}}(i) > B_{\mathcal{Q}}(i)$. And thus,

$$B_{\mathcal{Q}}(R_1) < \sum_{i \in R_1} (B_{\mathcal{P}}(i) - B_{\mathcal{Q}}(i)) \tag{3}$$

411 And if $i \in R_2$, then $B_{\mathcal{Q}}(i) > 5B_{\mathcal{P}}(i)$, and thus $B_{\mathcal{Q}}(i) - B_{\mathcal{P}}(i) > 4B_{\mathcal{P}}(i)$, giving

$$\begin{aligned}
B_{\mathcal{P}}(R_2) &< \frac{1}{4} \sum_{i \in R_2} (B_{\mathcal{Q}}(i) - B_{\mathcal{P}}(i)) \\
B_{\mathcal{P}}(R_2) + \sum_{i \in R_2} (B_{\mathcal{Q}}(i) - B_{\mathcal{P}}(i)) &< \frac{5}{4} \sum_{i \in R_2} (B_{\mathcal{Q}}(i) - B_{\mathcal{P}}(i)) \\
B_{\mathcal{Q}}(R_2) &< \frac{5}{4} \sum_{i \in R_2} (B_{\mathcal{Q}}(i) - B_{\mathcal{P}}(i)) \tag{4}
\end{aligned}$$

412 Since $|S_0| \leq 2^n$ and all elements $i \in S_0$ satisfy $\mathcal{P}(i) \leq 2^{-k}$, we have $B_{\mathcal{P}}(0) \leq 2^{n-k}$, where we
413 substitute $k = n + \log_2(100/\eta)$ to get

$$B_{\mathcal{P}}(0) \leq \frac{\eta}{100} \tag{5}$$

414 Then,

$$\begin{aligned}
& B_{\mathcal{Q}}(\{0\} \cup R_1 \cup R_2) = B_{\mathcal{Q}}(\{0\}) + B_{\mathcal{Q}}(R_1) + B_{\mathcal{Q}}(R_2) \\
\text{Using (3),(4) and (5)} & \leq \frac{\eta}{100} + \sum_{i \in \{0\}} (B_{\mathcal{Q}}(i) - B_{\mathcal{P}}(i)) + \sum_{i \in R_1} (B_{\mathcal{P}}(i) - B_{\mathcal{Q}}(i)) + \frac{5}{4} \sum_{i \in R_2} (B_{\mathcal{Q}}(i) - B_{\mathcal{P}}(i))
\end{aligned} \tag{6}$$

415 Here we partition the set $Bad \cup \{0\}$ into two sets Bad^+ and Bad^- , where $Bad^+ = \{i \in Bad \cup$
416 $\{0\} | B_{\mathcal{P}}(i) \geq B_{\mathcal{Q}}(i)\}$ and similarly $Bad^- = \{i \in Bad \cup \{0\} | B_{\mathcal{P}}(i) < B_{\mathcal{Q}}(i)\}$.

$$\begin{aligned}
& B_{\mathcal{Q}}(Bad) + B_{\mathcal{P}}(Bad) - 2d_{TV(Bad)}(B_{\mathcal{Q}}, B_{\mathcal{P}}) \\
& \geq 2(B_{\mathcal{Q}}(Bad) - 2d_{TV(Bad)}(B_{\mathcal{P}}, B_{\mathcal{Q}})) \\
& \text{(From 2)} > 2 \left(\eta - \varepsilon_2 - \varepsilon_1 - B_{\mathcal{Q}}(\{0\} \cup R_1 \cup R_2) - \sum_{i \in Bad} |B_{\mathcal{P}}(i) - B_{\mathcal{Q}}(i)| \right) \\
& \text{(From 6)} > 2 \left(.99\eta - \varepsilon_2 - \varepsilon_1 - \sum_{i \in R_1 \cup Bad^+} (B_{\mathcal{P}}(i) - B_{\mathcal{Q}}(i)) - \frac{5}{4} \sum_{i \in R_2 \cup Bad^-} (B_{\mathcal{Q}}(i) - B_{\mathcal{P}}(i)) \right)
\end{aligned}$$

417 And using the fact that $d_{TV}(B_{\mathcal{P}}, B_{\mathcal{Q}}) = \sum_{i: B_{\mathcal{P}}(i) \leq B_{\mathcal{Q}}(i)} (B_{\mathcal{Q}}(i) - B_{\mathcal{P}}(i)) =$
418 $\sum_{i: B_{\mathcal{P}}(i) > B_{\mathcal{Q}}(i)} (B_{\mathcal{P}}(i) - B_{\mathcal{Q}}(i)) = \varepsilon_2$, and the fact that $\forall_{i \in R_1} B_{\mathcal{P}}(i) > B_{\mathcal{Q}}(i)$ and
419 $\forall_{i \in R_2} B_{\mathcal{Q}}(i) > B_{\mathcal{P}}(i)$ we have,

$$B_{\mathcal{Q}}(Bad) + B_{\mathcal{P}}(Bad) - 2d_{TV(Bad)}(B_{\mathcal{Q}}, B_{\mathcal{P}}) > 2 \left(0.99\eta - \frac{13}{4}\varepsilon_2 - \varepsilon_1 \right)$$

420 □

421 A.3 Proof of Lemma 3

422 Recall that for all $r, s \in [m]$, $\mathbb{E}[X_{r,s}] = \sum_{b \in Bad} B_{\mathcal{P}}(b)B_{\mathcal{Q}}(b)$. Then since $X = \sum_{r,s \in [m]} X_{r,s}$,
423 $\mathbb{E}[X] = \sum_{r,s \in [m]} \mathbb{E}[X_{r,s}] = m^2 \mathbb{E}[X_{r,s}]$. Then for $i, j, k, l \in [m]$,

- 424 • if $i = k, j = l$ then $\mathbb{E}[X_{i,j}X_{k,l}] = \sum_{b \in Bad} B_{\mathcal{P}}(b)B_{\mathcal{Q}}(b) = \mathbb{E}[X_{r,s}]$
- 425 • if $i = k, j \neq l$ then $\mathbb{E}[X_{i,j}X_{k,l}] = \sum_{b \in Bad} B_{\mathcal{P}}(b)B_{\mathcal{Q}}^2(b)$
- 426 • if $i \neq k, j = l$ then $\mathbb{E}[X_{i,j}X_{k,l}] = \sum_{b \in Bad} B_{\mathcal{P}}^2(b)B_{\mathcal{Q}}(b)$
- 427 • if $i \neq k, j \neq l$ then $\mathbb{E}[X_{i,j}X_{k,l}] = \left(\sum_{b \in Bad} B_{\mathcal{P}}(b)B_{\mathcal{Q}}(b) \right)^2 = \mathbb{E}[X_{r,s}]^2$

$$\begin{aligned}
\mathbb{E}[X^2] &= \mathbb{E} \left[\sum_{i,j,k,l \in [m]} X_{i,j}X_{k,l} \right] \\
&= \mathbb{E} \left[\sum_{\substack{a \neq c, b \neq d \\ i,j,k,l \in [m]}} X_{i,j}X_{k,l} \right] + \mathbb{E} \left[\sum_{\substack{a=c, b \neq d \\ i,j,k,l \in [m]}} X_{i,j}X_{k,l} \right] + \mathbb{E} \left[\sum_{\substack{a \neq c, b=d \\ i,j,k,l \in [m]}} X_{i,j}X_{k,l} \right] + \mathbb{E} \left[\sum_{\substack{a=c, b=d \\ i,j,k,l \in [m]}} X_{i,j}X_{k,l} \right] \\
&= m^2(m-1)^2 \mathbb{E}[X_{r,s}]^2 + m^2(m-1) \left(\sum_{b \in Bad} (B_{\mathcal{P}}(b) + B_{\mathcal{Q}}(b)) B_{\mathcal{P}}(b)B_{\mathcal{Q}}(b) \right) + m^2 \mathbb{E}[X_{r,s}] \\
&\leq m^4 \mathbb{E}[X_{r,s}]^2 + m^3 \left(\sum_{b \in Bad} (B_{\mathcal{P}}(b) + B_{\mathcal{Q}}(b)) B_{\mathcal{P}}(b)B_{\mathcal{Q}}(b) \right) + m^2 \mathbb{E}[X_{r,s}]
\end{aligned}$$

428 Then,

$$\begin{aligned} \frac{\mathbb{E}[X]^2}{\mathbb{E}[X^2]} &> \frac{m^4 \mathbb{E}[X_{r,s}]^2}{m^4 \mathbb{E}[X_{r,s}]^2 + m^3 \left(\sum_{b \in \text{Bad}} (B_{\mathcal{P}}(b) + B_{\mathcal{Q}}(b)) B_{\mathcal{P}}(b) B_{\mathcal{Q}}(b) \right) + m^2 \mathbb{E}[X_{r,s}]} \\ &= \frac{1}{1 + m^{-1} \left(\frac{\sum_{b \in \text{Bad}} (B_{\mathcal{P}}(b) + B_{\mathcal{Q}}(b)) B_{\mathcal{P}}(b) B_{\mathcal{Q}}(b)}{\left(\sum_{b \in \text{Bad}} B_{\mathcal{P}}(b) B_{\mathcal{Q}}(b) \right)^2} \right) + m^{-2} \mathbb{E}[X_{r,s}]^{-1}} \end{aligned}$$

429 We will now focus on finding the maximum for ratio of summations:

$$\begin{aligned} \frac{\sum_{b \in \text{Bad}} (B_{\mathcal{P}}(b) + B_{\mathcal{Q}}(b)) B_{\mathcal{P}}(b) B_{\mathcal{Q}}(b)}{\left(\sum_{b \in \text{Bad}} B_{\mathcal{P}}(b) B_{\mathcal{Q}}(b) \right)^2} &= \frac{\sum_{b \in \text{Bad}} \left(\sqrt{\frac{B_{\mathcal{P}}(b)}{B_{\mathcal{Q}}(b)}} + \sqrt{\frac{B_{\mathcal{Q}}(b)}{B_{\mathcal{P}}(b)}} \right) (B_{\mathcal{P}}(b) B_{\mathcal{Q}}(b))^{3/2}}{\left(\sum_{b \in \text{Bad}} B_{\mathcal{P}}(b) B_{\mathcal{Q}}(b) \right)^2} \\ &\leq \frac{(\sqrt{1/5} + \sqrt{5}) \sum_{b \in \text{Bad}} (B_{\mathcal{P}}(b) B_{\mathcal{Q}}(b))^{3/2}}{\left(\sum_{b \in \text{Bad}} B_{\mathcal{P}}(b) B_{\mathcal{Q}}(b) \right)^2} \\ &\quad \text{(If } b \in \text{Bad then } B_{\mathcal{P}}(b)/B_{\mathcal{Q}}(b) \in [5^{-1}, 2]) \\ &\leq \frac{3}{\left(\sum_{b \in \text{Bad}} B_{\mathcal{P}}(b) B_{\mathcal{Q}}(b) \right)^{1/2}} = 3 \mathbb{E}[X_{r,s}]^{-1/2} \\ &\quad \text{(Using the monotonicity of } \ell_p \text{ norms)} \end{aligned}$$

430 Thus,

$$\begin{aligned} \frac{\mathbb{E}[X]^2}{\mathbb{E}[X^2]} &> \frac{1}{1 + 3m^{-1} \mathbb{E}[X_{r,s}]^{-1/2} + m^{-2} \mathbb{E}[X_{r,s}]^{-1}} \\ &> \frac{1}{5} \quad \text{(Since } m^2 \mathbb{E}[X_{r,s}] \geq 1) \end{aligned}$$

431 The Chebyshev bound from 2 tells us that $\Pr[|X - \mathbb{E}[X]| < \mathbb{E}[X]] > \mathbb{E}[X]^2 / \mathbb{E}[X^2] > \frac{1}{5}$. Thus,

$$\begin{aligned} \Pr[X > 0] &> \frac{1}{5} \\ \Pr[X \geq 1] &> \frac{1}{5} \quad \text{(Since } X \text{ takes only integer values)} \end{aligned}$$

432 A.4 Proof of Proposition 5

433 **Proposition 5.** Let \mathcal{P}, \mathcal{Q} be distributions and let $p \sim \mathcal{P}$ and $q \sim \mathcal{Q}$. Then,

434 1. If $d_{\infty}(\mathcal{P}, \mathcal{Q}) < \varepsilon$ then

$$\frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} \geq \frac{\mathcal{P}(p)}{\mathcal{P}(p) + (1 + \frac{2\varepsilon}{1-\varepsilon})\mathcal{P}(q)}$$

435 2. If $d_{TV}(\mathcal{P}, \mathcal{Q}) > \varepsilon_1$, then for $0 \leq \alpha < \varepsilon_1$, with probability at least $(d_{TV}(\mathcal{P}, \mathcal{Q}) - \alpha)/2$,

$$\frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} < \frac{\mathcal{P}(p)}{\mathcal{P}(p) + (1 + \alpha)\mathcal{P}(q)}$$

436 *Proof.* If $d_{\infty}(\mathcal{P}, \mathcal{Q}) < \varepsilon$ then

$$\begin{aligned} \frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} &\geq \frac{\mathcal{P}(p)(1 - \varepsilon)}{\mathcal{P}(p)(1 - \varepsilon) + (1 + \varepsilon)\mathcal{P}(q)} \\ &= \frac{\mathcal{P}(p)}{\mathcal{P}(p) + (1 + \frac{2\varepsilon}{1-\varepsilon})\mathcal{P}(q)} \end{aligned}$$

437 and hence we show the first part of the claim.

438 For the second part of the proof we introduce the some sets. Let $H_0 = \{h | 1 \leq \frac{\mathcal{Q}(h)}{\mathcal{P}(h)} < 1 + \alpha\}$
 439 and $H_1 = \{h | 1 + \alpha \leq \frac{\mathcal{Q}(h)}{\mathcal{P}(h)}\}$ and $H = H_0 \cup H_1$. Similarly define, $L_0 = \{\ell | 1 - \alpha < \frac{\mathcal{Q}(\ell)}{\mathcal{P}(\ell)} < 1\}$,
 440 $L_1 = \{\ell | \frac{\mathcal{Q}(\ell)}{\mathcal{P}(\ell)} \leq 1 - \alpha\}$ and $L = L_0 \cup L_1$.

441 Now consider that we have a pair of samples, $p \sim \mathcal{P}$ and $q \sim \mathcal{Q}$. We know that either $\mathcal{P}(L) \geq 1/2$
 442 or $\mathcal{P}(H) > 1/2$.

443 $\mathcal{P}(L) \geq 1/2$: We see that $\Pr[p \in L] \geq 1/2$. Then from the definition of H_0 , $\mathcal{Q}(h_0) - \mathcal{P}(h_0) < \alpha$
 444 and recall that $\mathcal{Q}(H) - \mathcal{P}(H) = d_{TV}(\mathcal{P}, \mathcal{Q})$. Thus we have that $\mathcal{Q}(H_1) - \mathcal{P}(H_1) > d_{TV}(\mathcal{P}, \mathcal{Q}) - \alpha$
 445 and hence $\Pr[q \in H_1] > d_{TV}(\mathcal{P}, \mathcal{Q}) - \alpha$. We can now confirm that $q \in H_1 \wedge p \in L$ with probability
 446 at least $(d_{TV}(\mathcal{P}, \mathcal{Q}) - \alpha)/2$. Then,

$$\begin{aligned} \frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} &< \frac{\mathcal{P}(p)}{\mathcal{P}(p) + \mathcal{Q}(q)} \quad (\text{From } \mathcal{P}(p) > \mathcal{Q}(p)) \\ &< \frac{\mathcal{P}(p)}{\mathcal{P}(p) + (1 + \alpha)\mathcal{P}(q)} \quad (\text{Since } q \in H_1) \end{aligned}$$

447 $\mathcal{P}(H) > 1/2$: We see that $\Pr[q \in H] \geq 1/2$. Then we have that $\mathcal{P}(L_0) - \mathcal{Q}(L_0) < \alpha$ and also that
 448 $\mathcal{P}(L) - \mathcal{Q}(L) = d_{TV}(\mathcal{P}, \mathcal{Q})$, we have that $\mathcal{P}(L_1) - \mathcal{Q}(L_1) < d_{TV}(\mathcal{P}, \mathcal{Q}) - \alpha$. Then, we deduce
 449 that probability at least $(d_{TV}(\mathcal{P}, \mathcal{Q}) - \alpha)/2$, $q \in H \wedge p \in L_1$. Then,

$$\begin{aligned} \frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{Q}(q)} &< \frac{\mathcal{Q}(p)}{\mathcal{Q}(p) + \mathcal{P}(q)} \quad (\text{From } \mathcal{P}(q) < \mathcal{Q}(q)) \\ &< \frac{\mathcal{P}(p)(1 - \alpha)}{\mathcal{P}(p)(1 - \alpha) + \mathcal{P}(q)} \quad (\text{Since } p \in L_1) \\ &< \frac{\mathcal{P}(p)}{\mathcal{P}(p) + (1 + \alpha)\mathcal{P}(q)} \end{aligned}$$

450

□

451 A.5 The outbucket subroutine

Algorithm 4 OutBucket($B_{\mathcal{P}}, B_{\mathcal{Q}}, k, \theta, \delta$)

- 1: Sample $\max\left(\frac{4(k+1)}{\theta^2}, \frac{8 \ln(4/\delta)}{\theta^2}\right)$ times from $B_{\mathcal{P}}$ and $B_{\mathcal{Q}}$ and construct empirical distributions $\widehat{B}_{\mathcal{P}}$ and $\widehat{B}_{\mathcal{Q}}$.
 - 2: **Return** $d_{TV}(\widehat{B}_{\mathcal{P}}, \widehat{B}_{\mathcal{Q}})$
-

452 **B Data missing from the main paper**

Benchmark	Dimensions	Barbarik2		Pacoco	
		Result	# of samples	Result	# of samples
SetTest.sk_9_21	21	R	2817	R	58000
s27_7_4	7	R	4789	R	30000
polynomial.sk_7_25	25	R	4789	R	66000
s27_15_7	7	R	4789	R	30000
Pollard.sk_1_10	10	R	7606	R	36000
s298_3_2	17	R	57431	R	50000
s27_3_2	7	R	62220	R	30000
s27_new_15_7	7	R	128264	R	30000
s444_3_2	24	R	848148	R	64000
s526a_3_2	24	R	848148	R	64000
s27_new_3_2	7	R	905579	R	30000
s298_7_4	17	R	12708989	R	50000
s510_15_7	25	R	12708989	R	66000
s1488_15_7	14	R	12708989	R	44000
s27_new_7_4	7	A	23997012	R	30000
s298_15_7	17	R	38126967	R	50000
s526_3_2	24	TO	-	R	64000
s420_3_2	34	TO	-	R	83000
s420_new1_3_2	34	TO	-	R	83000
s382_3_2	24	TO	-	R	64000
s641_3_2	54	TO	-	R	123000
111.sk_2_36	36	TO	-	R	87000
s526_7_4	24	TO	-	R	64000
s510_3_2	25	TO	-	R	66000
7.sk_4_50	50	TO	-	R	115000
56.sk_6_38	38	TO	-	R	91000
s820a_15_7	23	TO	-	R	62000
ProjectService3.sk_12_55	55	TO	-	R	125000
s420_7_4	34	TO	-	R	83000
s832a_7_4	23	TO	-	R	62000
s420_new1_7_4	34	TO	-	R	83000
s420_15_7	34	TO	-	R	83000
s420_new_7_4	34	TO	-	R	83000
s713_3_2	54	TO	-	R	123000
s526a_15_7	24	TO	-	R	64000
s1196a_7_4	32	TO	-	R	80000
81.sk_5_51	51	TO	-	R	117000
s420_new_3_2	34	TO	-	R	83000
s349_15_7	24	TO	-	R	64000
s344_15_7	24	TO	-	R	64000
s713_7_4	54	TO	-	R	123000
77.sk_3_44	44	TO	-	R	103000
s420_new_15_7	34	TO	-	R	83000
s832a_3_2	23	TO	-	R	62000
UserServiceImpl.sk_8_32	32	TO	-	R	80000
19.sk_3_48	48	TO	-	R	111000
s953a_7_4	45	TO	-	R	105000
s349_7_4	24	TO	-	R	64000
s444_15_7	24	TO	-	R	64000
LoginService2.sk_23_36	36	TO	-	R	87000
29.sk_3_45	45	TO	-	R	105000
s1238a_3_2	32	TO	-	R	80000
s1488_3_2	14	TO	-	R	44000

s344_7_4	24	TO	-	R	64000
s1196a_3_2	32	TO	-	R	80000
s444_7_4	24	TO	-	R	64000
51.sk_4_38	38	TO	-	R	91000
57.sk_4_64	64	TO	-	R	143000
53.sk_4_32	32	TO	-	R	80000
s832a_15_7	23	TO	-	R	62000
s953a_3_2	45	TO	-	R	105000
63.sk_3_64	64	TO	-	R	143000
s526_15_7	24	TO	-	R	64000
110.sk_3_88	88	TO	-	R	190000
s349_3_2	24	TO	-	R	64000
s820a_3_2	23	TO	-	R	62000
s1196a_15_7	32	TO	-	R	80000
10.sk_1_46	46	TO	-	R	107000
s1238a_7_4	32	TO	-	R	80000
s420_new1_15_7	34	TO	-	R	83000
s344_3_2	24	TO	-	R	64000
s953a_15_7	45	TO	-	R	105000
s526a_7_4	24	TO	-	R	64000
80.sk_2_48	48	TO	-	R	111000
32.sk_4_38	38	TO	-	R	91000
s820a_7_4	23	TO	-	R	62000
s382_15_7	24	TO	-	R	64000
17.sk_3_45	45	TO	-	R	105000
s382_7_4	24	TO	-	R	64000
s641_7_4	54	TO	-	R	123000
s1238a_15_7	32	TO	-	R	80000
s838_7_4	66	TO	-	R	147000
27.sk_3_32	32	TO	-	R	80000
55.sk_3_46	46	TO	-	R	107000
109.sk_4_36	36	TO	-	R	87000
70.sk_3_40	40	TO	-	R	95000
s838_3_2	66	TO	-	R	147000

Table 2: Performance of Pacoco. We experiment with 87 benchmarks, and out of the 87 benchmarks. In the table ‘TO’ represents that either the tester timed out or asked for more than 10^8 samples. The value of the parameter for closeness is $\varepsilon = 0.05$, for fairness is $\eta = 0.9$ and for confidence is $\delta = 0.2$.