

Testing Closeness of Multivariate Distributions via Ramsey Theory

Ilias Diakonikolas*
 UW Madison
 ilias@cs.wisc.edu

Daniel M. Kane†
 UC San Diego
 dakane@ucsd.edu

Sihan Liu
 UC San Diego
 sil046@ucsd.edu

November 23, 2023

Abstract

We investigate the statistical task of closeness (or equivalence) testing for multidimensional distributions. Specifically, given sample access to two unknown distributions \mathbf{p}, \mathbf{q} on \mathbb{R}^d , we want to distinguish between the case that $\mathbf{p} = \mathbf{q}$ versus $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} > \epsilon$, where $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k}$ denotes the generalized \mathcal{A}_k distance between \mathbf{p} and \mathbf{q} — measuring the maximum discrepancy between the distributions over any collection of k disjoint, axis-aligned rectangles. Our main result is the first closeness tester for this problem with *sub-learning* sample complexity in any fixed dimension and a nearly-matching sample complexity lower bound.

In more detail, we provide a computationally efficient closeness tester with sample complexity $O\left(k^{6/7}/\text{poly}_d(\epsilon)\log^d(k)\right)$. On the lower bound side, we establish a qualitatively matching sample complexity lower bound of $\Omega(k^{6/7}/\text{poly}(\epsilon))$, even for $d = 2$. These sample complexity bounds are surprising because the sample complexity of the problem in the univariate setting is $\Theta(k^{4/5}/\text{poly}(\epsilon))$. This has the interesting consequence that the jump from one to two dimensions leads to a substantial increase in sample complexity, while increases beyond that do not.

As a corollary of our general \mathcal{A}_k tester, we obtain d_{TV} -closeness testers for pairs of k -histograms on \mathbb{R}^d over a common unknown partition, and pairs of uniform distributions supported on the union of k unknown disjoint axis-aligned rectangles.

Both our algorithm and our lower bound make essential use of tools from Ramsey theory.

*Supported by NSF Medium Award CCF-2107079, NSF Award CCF-1652862 (CAREER), and a Sloan Research Fellowship.

†Supported by NSF Medium Award CCF-2107547, and NSF Award CCF-1553288 (CAREER), and a grant from CasperLabs.

1 Introduction

Background and Motivation A fundamental statistical task is to ascertain whether a set of samples comes from a given model, where the model may consist of either a single fully specified probability distribution or a family of probability distributions. The study of this broad task was initiated in a field now known as *statistical hypothesis testing* over a century ago [Pea00, NP33]; see, e.g., [LR05] for an introductory textbook on the topic. In the past three decades, hypothesis testing has been extensively studied by the theoretical computer science and information-theory communities — under the name *distribution testing* — in the framework of property testing [RS96, GGR98]. It is instructive to note that the TCS style definition of hypothesis testing is equivalent to the minimax testing definition introduced and studied by Ingster and coauthors [Ing94, Ing97, IS03].

The paradigmatic problem in distribution testing is the following: given sample access to one or more unknown probability distributions, we want to correctly distinguish (with high probability) between the cases that the underlying distributions satisfy some global property \mathcal{P} or are “far” from satisfying the property. The primary objective is to obtain a tester that is statistically efficient, i.e., it has information-theoretically optimal sample complexity. An additional important criterion is computational efficiency; that is, the testing algorithm should run in sample-polynomial time. After the pioneering early works formulating this field [GR00, BFR⁺00] from a TCS perspective, there has been substantial progress on testing a wide range of properties; see, e.g., [BFF⁺01, BDKR02, BKR04, Pan08, Val11, VV11, ADJ⁺11, LRR11, VV14, CDVV14, DKN15a, CDKS17, DDK18, CDKS18, DGK⁺21, CJKL22, CDKL22] for a sample of works, and [Rub12, Can22] for surveys on the topic.

Here we study the problem of *closeness testing* (or equivalence testing) between two unknown probability distributions. Specifically, given independent samples from a pair of distributions \mathbf{p}, \mathbf{q} , we want to determine whether the two distributions are the same versus ϵ -far from each other. Early work on this problem [BFR⁺00] focused on the setting that \mathbf{p}, \mathbf{q} are arbitrary discrete distributions of a given support size n , and the metric used to quantify “closeness” is the ℓ_1 -distance (equivalently, total variation distance). It is now known [CDVV14] that the optimal sample complexity of ℓ_1 -closeness testing for distributions with support of size n is $\Theta(\max\{n^{2/3}/\epsilon^{4/3}, n^{1/2}/\epsilon^2\})$.

In summary, it is known that the complexity measure determining the sample complexity of testing the equivalence (and a range of other related properties) of unstructured (i.e., potentially arbitrary) discrete distributions is the domain size of the underlying distributions. Unfortunately, this implies that if \mathbf{p}, \mathbf{q} are (potentially arbitrary) continuous distributions (even in one dimension!), no closeness tester with finite sample complexity exists. There are two natural approaches to circumvent this bottleneck. The first approach is to assume that \mathbf{p}, \mathbf{q} have some nice structure, in which case the domain size may not be the right complexity measure for the testing problem. The second approach is to make no assumptions on the underlying distributions, but relax the metric under which we measure closeness.

Interestingly, it turns out that these two seemingly orthogonal approaches are intimately related to each other. In particular, for the important special case of *one-dimensional* distributions, a line of works, see, e.g., [DDS⁺13, DKN15b, DKN15a, DKN17], developed a general framework that yields optimal testers (for closeness and other properties) for a range of structured distribution families. The key idea underlying these testers is to design a *single* tester for *arbitrary* one-dimensional distributions *but under a different — carefully selected — metric*; and then appropriately use this metric as a proxy for the total variation distance (for each structured distribution family of interest).

In more detail, for one-dimensional distributions $\mathbf{p}, \mathbf{q} : \mathbb{R} \rightarrow \mathbb{R}_+$, the appropriate metric is known as \mathcal{A}_k -distance [DL01, CDSS14a] and is defined as follows: The \mathcal{A}_k -distance between one-dimensional distributions \mathbf{p} and \mathbf{q} , denoted by $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k}$, is defined as the maximum ℓ_1 -distance

between the reduced distributions¹ obtained from \mathbf{p}, \mathbf{q} over all partitions of the domain in at most k intervals. The motivation for this particular definition of the \mathcal{A}_k -distance [CDSS14a, DKN15b] between one-dimensional distributions comes from the VC-inequality (see, e.g., page 31 of [DL01]).

The positive integer k in the definition of the \mathcal{A}_k -distance is a tunable parameter that is selected appropriately depending on the application. For $k = 2$, the \mathcal{A}_k -distance amounts to the distance between the *cumulative* distribution functions (known as Kolmogorov distance). As k increases, the metric becomes stronger and converges to the total variation distance when $k \rightarrow \infty$ (under mild assumptions on the distributions). Moreover, if the underlying distributions \mathbf{p}, \mathbf{q} belong to some class of shape restricted densities (e.g., univariate histograms or log-concave distributions), a finite value of k suffices so that the \mathcal{A}_k -distance closely approximates the total variation distance.

It is worth noting that, in addition to distribution testing, the one-dimensional \mathcal{A}_k distance has been a crucial ingredient in developing efficient *learning* algorithms for structured univariate distributions [CDSS13, CDSS14a, ADLS17, CLM20].

Testing Closeness of Multivariate Distributions The main motivation behind this work is to generalize the aforementioned framework to the *multivariate* setting with a focus on the task of closeness testing. A first step to achieve this is an appropriate generalization of the notion of \mathcal{A}_k -distance (which applies to one-dimensional distributions) for distributions on \mathbb{R}^d for all $d \geq 1$. Here we study the following natural definition, that has been previously used in the context of learning [DLS18] and uniformity testing [DKP19] for multivariate distributions.

Definition 1.1 (Multidimensional \mathcal{A}_k -distance). For two probability distributions (with densities/mass functions) $\mathbf{p}, \mathbf{q} : \mathbb{R}^d \mapsto \mathbb{R}_+$ and $k \in \mathbb{Z}^+$, we define the *multi-dimensional \mathcal{A}_k -distance* between \mathbf{p} and \mathbf{q} as the maximum value of $\sum_{i=1}^k |\mathbf{p}(R_i) - \mathbf{q}(R_i)|$ for k arbitrarily chosen non-overlapping axis-aligned rectangles $\{R_i\}_{i=1}^k$ in \mathbb{R}^d .

Motivation for Definition 1.1 Recall that the total variation distance between two distributions \mathbf{p}, \mathbf{q} on \mathbb{R}^d is defined as $d_{TV}(\mathbf{p}, \mathbf{q}) = \sup_{A \in \mathcal{S}} |\mathbf{p}(A) - \mathbf{q}(A)|$, where \mathcal{S} is the collection of all measurable subsets on \mathbb{R}^d . Since learning or testing under the total variation distance may be too strong a goal if the underlying distributions lack structure, a reasonable compromise is to consider alternative metrics. The VC-inequality states the following: Let \mathcal{A} be any collection of subsets of \mathbb{R}^d with VC-dimension d . Then for *any* distribution \mathbf{p} on \mathbb{R}^d it holds that $\mathbf{E}[\sup_{A \in \mathcal{A}} |\hat{\mathbf{p}}_n(A) - \mathbf{p}(A)|] = O(\sqrt{d/n})$, where $\hat{\mathbf{p}}_n$ is the empirical distribution obtained after drawing n i.i.d. samples from \mathbf{p} . In other words, for $n \gg d/\epsilon^2$, the empirical distribution is ϵ -close to \mathbf{p} with respect to the \mathcal{A} -metric, defined as $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}} \stackrel{\text{def}}{=} \sup_{A \in \mathcal{A}} |\mathbf{p}(A) - \mathbf{q}(A)|$. For the univariate case, the \mathcal{A}_k -distance defined in the aforementioned works [CDSS14a, DKN15b, DKN15a] is obtained from the \mathcal{A} -metric by considering the family of all unions of at most k intervals (which has VC-dimension $2k$).

Our Definition 1.1 is a natural generalization of the one-dimensional definition, where we consider the family of all unions of at most k rectangles, which has VC-dimension $\tilde{\Theta}(kd)$. Since learning an arbitrary distribution on \mathbb{R}^d under this metric requires $\tilde{\Theta}(kd)/\epsilon^2$ samples, it is natural to ask whether the distribution testing problem has qualitatively lower sample complexity. We also note that the multidimensional \mathcal{A}_k -distance is a strengthening of the Kolmogorov-Sminov (KS) metric and converges to the total variation distance as $k \rightarrow \infty$ (under mild assumptions). It should be noted that a line of work in mathematical statistics — see, e.g., [Bic69, FR79, Hen88, JPZ97] for some classical works — has developed two-sample testers (aka closeness testers) for *non-parametric*

¹The reduced distribution obtained from \mathbf{p} with respect to a partition of the domain into k subsets R_1, \dots, R_k is the discrete distribution with support size k assigning probability mass $\mathbf{p}(R_i)$ to the i -th point.

multivariate distributions under the KS metric. Our work can be viewed as a strengthening and generalization of these results in the minimax setting.

We believe that, in addition to being a potential tool for performing multivariate d_{TV} -closeness testing for structured distributions, the \mathcal{A}_k distance is an interesting metric on its own merits. To see this, we recall that one of the main motivations for considering the total variation distance is the following property: If a decision algorithm is run twice on different inputs that follow two distributions that are close in total variation distance, then the acceptance probabilities will also be approximately the same in the two cases. Hence, for two distributions \mathbf{p}, \mathbf{q} that have passed the d_{TV} -closeness testing, we can be confident that running some downstream decision algorithm on inputs drawn from \mathbf{p} and from \mathbf{q} should give similar results. For the \mathcal{A}_k distance, we have an analogous property if one restricts the algorithm in the above statement to be an axis-aligned decision tree (i.e., a decision tree whose leaf nodes follow the branching rule of $x_i < b$ for some coordinate $i \in [d]$ and some real number $b \in \mathbb{R}$) with at most k leaves. Though being a restricted family of algorithms, axis-aligned decision trees are commonly used in machine learning applications due to their exceptional interpretability; see, e.g., [YA01, BDS10, BPG20]. This suggests that testing in \mathcal{A}_k distance, even though being a weaker test compared to its d_{TV} -counterpart for arbitrary distributions (which is provably impossible without structural information), may be sufficient for certain structured downstream decision-making tasks.

We return to our closeness testing task. One approach to solve the multidimensional² \mathcal{A}_k -closeness testing problem is to learn \mathbf{p} and \mathbf{q} up to \mathcal{A}_k -distance $\epsilon/4$, and then check whether the hypotheses are $\epsilon/4$ -close to each other. Thus, the sample complexity of closeness testing is bounded above by the sample complexity of learning (within constant factors). Since $\tilde{\Theta}(kd/\epsilon^2)$ samples suffice to learn an arbitrary distribution on \mathbb{R}^d up to \mathcal{A}_k -distance ϵ , the naive “testing-by-learning” approach requires $\Omega(k)$ samples (even in one dimension and for constant ϵ).

It is natural to ask whether a better sample size could be achieved for testing, since closeness testing is, in some sense, less demanding than learning. That is, the goal is to develop a closeness tester with sample complexity *strongly sublinear* in k , namely $O(k^c)$ for some constant $c < 1$. The aforementioned line of work on univariate distributions [DKN15b, DKN15a, DKN17] developed identity and closeness testers under the \mathcal{A}_k -distance with strongly sublinear sample complexity. These testers were also applied to give total variation distance testers for classes of “shape constrained” distributions [BBBB72, GJ14], including histograms and logconcave distributions.

Concretely, for the problem of closeness testing of univariate distributions, [DKN15a] developed a sample-optimal \mathcal{A}_k -closeness tester with sample complexity of $\Theta(k^{4/5}/\text{poly}(\epsilon))$ (for not too small ϵ). Interestingly, this bound differs from the sample complexity of closeness testing discrete distributions on k points, which is $\Theta(k^{2/3}/\text{poly}(\epsilon))$ [CDVV14].

This discussion motivates the following natural question:

What is the sample complexity of \mathcal{A}_k -closeness testing for multivariate distributions?

Prior to this work, no closeness tester with sub-learning sample complexity was known even for $d = 2$. The main contribution of this work is a *sample near-optimal and computationally efficient* \mathcal{A}_k -closeness tester in any **fixed**³ dimension. Moreover, we show that the sample complexity of our tester is optimal as a function of k , within logarithmic factors. As an immediate corollary, we

²We will henceforth omit the term “multidimensional” when it is clear from the context, and use the term \mathcal{A}_k -distance for multivariate distributions as well.

³We emphasize here that the focus of our work is in closeness testing of non-parametric families of distributions. In non-parametric estimation/testing, the sample complexity inherently scales exponentially with the dimension; hence, it is standard to consider the dimension as being fixed. For example, estimation/testing for the class of log-concave distributions on \mathbb{R}^d is known to require $2^{\Omega(d)}$ samples (see [KS16]).

obtain the first closeness tester for multivariate histogram distributions (with respect to the same unknown set of axis-aligned rectangles) under the total variation distance.

Specifically, our main result (Theorem 1.2) establishes the following: *For any $k, d \in \mathbb{Z}_+, \epsilon > 0$, and sample access to arbitrary distributions \mathbf{p}, \mathbf{q} on \mathbb{R}^d , there exists a closeness testing algorithm under the \mathcal{A}_k -distance using $O((k^{6/7}/\text{poly}_d(\epsilon)) \log^d(k))$ samples. Moreover, this bound is information-theoretically optimal as a function of k , even for $d = 2$.* We remark that our \mathcal{A}_k -testing algorithm applies to *any* pair of distributions (over both continuous and discrete domains).

As a corollary, we obtain the first closeness tester (with sub-learning sample complexity) between k -histograms with respect to the total variation distance. A probability distribution on \mathbb{R}^d with density \mathbf{p} is called a *k-histogram* if there exists a partition of the support into k axis-aligned rectangles R_1, \dots, R_k such that \mathbf{p} is constant on R_i , for all $i = 1, \dots, k$. This is one of the most basic non-parametric distribution families and have been extensively studied in statistics [Sco79, FD81, Sco92, LN96, DL04, WN07, Kle09] and computer science — including database theory [JKM⁺98, CMN98, TGIK02, GGI⁺02, GKS06, ILR12, ADH⁺15] and theoretical ML [DDS12, CDSS13, CDSS14a, CDSS14b, ADLS17, ADK15, DDS⁺13, DKN15a, DKN15b, DKN17, DKP19, CDKL22]. Prior to this work, no closeness testing algorithm with sub-learning sample complexity was known for k -histograms, even for $d = 2$. As a corollary of our main result, we provide such an algorithm (see Corollary 2.18) for the case that the two histograms are supported on the same unknown partition. In addition, we also obtain d_{TV} -closeness tester for uniform distributions supported on some unknown k disjoint axis-aligned rectangles (see Corollary 2.19). We remark that though histograms and uniform distributions over unions of axis-aligned rectangles are conceptually similar, these two families of distributions are orthogonal to each other.

1.1 Our Results

We study the complexity of closeness testing between two (arbitrary) distributions \mathbf{p}, \mathbf{q} on \mathbb{R}^d with respect to the \mathcal{A}_k distance. Our main result is the following.

Theorem 1.2 (Main Result). *Given $\epsilon > 0$, integer $k \geq 2$, and sample access to distributions with density functions $\mathbf{p}, \mathbf{q} : \mathbb{R}^d \rightarrow \mathbb{R}_+$, there exists a computationally efficient algorithm which draws $C 2^{d/3} k^{6/7} \log^{3d}(k)/\epsilon^{\alpha_d}$ samples from \mathbf{p}, \mathbf{q} , for a sufficiently large universal constant $C > 0$, where $\alpha_d = O(d^2 2^{2^{d+1}})$, and with probability at least $2/3$ correctly distinguishes whether $\mathbf{p} = \mathbf{q}$ versus $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} \geq \epsilon$. Moreover, $\Omega(\min\{k^{6/7}/\epsilon^{8/7}, k\})$ many samples are information-theoretically necessary for this hypothesis testing task, even if \mathbf{p}, \mathbf{q} are two-dimensional discrete distributions on a sufficiently large domain.*

Discussion To interpret Theorem 1.2, some comments are in order. We reiterate that the focus of our work is on the *non-parametric* setting and consequently we view the dimension d as a fixed constant. In this regime, the sample complexity of our algorithm is $\tilde{O}_d(k^{6/7})/\text{poly}_d(\epsilon)$.

The one-dimensional special case of our closeness testing result was solved in [DKN15a], where the authors established a tight sample complexity bound of $\Theta(k^{4/5}/\epsilon^{6/5} + k^{1/2}/\epsilon^2)$. Prior to our work, no $o(k)$ sample upper bound was known for this testing problem even for $d = 2$ and $\epsilon = 0.99$.

For the regime of fixed dimension that we focus on, our upper and lower bounds are essentially optimal in terms of their dependence on k — the main parameter of interest. For simplicity, let us fix ϵ to be a universal constant. Examining the exponent of k in the dominant term of the sample complexity, we observe a surprising pattern: the exponent begins at $4/5$ when $d = 1$ (as follows from the prior work [DKN15a]), jumps to $6/7$ when $d = 2$, and then stays at $6/7$ as d increases (as

follows from Theorem 1.2)! This suggests that the $d = 1$ case is a degenerate case and the essence and complexity of the problem is not entirely revealed until $d = 2$.

Some remarks are in order regarding the dependence of the sample complexity on the parameters ϵ and d . First, we briefly comment on the $\log^d(k)$ term. Perhaps surprisingly, prior work [DKP19] has shown a sample complexity *lower bound* of $(\sqrt{k}/\epsilon^2)\Omega(\log(k)/d)^{d-1}$ for the *easier* problem of \mathcal{A}_k -uniformity testing. This suggests that the $\log^d(k)$ factor is necessary for closeness testing as well, assuming that k is sufficiently large. Finally, we conjecture that the correct dependence on ϵ in the sample complexity of this task should be a fixed degree polynomial, independent of d . We leave this as an interesting technical question for future work (see Question 4.1).

Regarding our sample complexity lower bound, Theorem 1.2 does not specify how large the domain size of the hard distributions needs to be. Due to the application of Ramsey-theoretic arguments in the proof of our lower bound, we need it to be extremely large in terms of k (a tower function of k). In Section 3.4, we show that the domain size can be optimized to be (at most) doubly exponential in k — using a significantly more sophisticated construction (Theorem 3.8).

As immediate corollaries of our main theorem, we obtain d_{TV} -closeness testers (with strongly sub-learning sample complexities) for multivariate structured distributions. In particular, we highlight here the d_{TV} -closeness tester for distributions in \mathbb{R}^d that are k -histograms, i.e., piecewise constant over (the same) k unknown disjoint axis-aligned rectangles. Notably, the sample complexity of this tester is the same as that of our \mathcal{A}_k closeness testing. This implication and additional applications are given in Section 2.4.

1.2 Overview of Techniques

Here we provide a detailed overview of our technical approach to establish our upper and lower bounds.

Closeness Tester By definition of the \mathcal{A}_k distance, there exist k disjoint axis-aligned rectangles $\{R_i\}_{i=1}^k$ on \mathbb{R}^d which witness the \mathcal{A}_k discrepancy between \mathbf{p} and \mathbf{q} ; that is, $\sum_{i=1}^k |\mathbf{p}(R_i) - \mathbf{q}(R_i)| = \|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k}$. If we knew what these rectangles were, the testing task would be easy. Indeed, we could simply consider the reduced measures of \mathbf{p} and \mathbf{q} over $\{R_i\}_{i=1}^k$ (recall that these measures, after normalization, become distributions with support size k that we can simulate access to) and then use an optimal ℓ_1 -closeness tester as a black-box. Given the optimal ℓ_1 -closeness tester of [CDVV14], such an approach would lead to a sample complexity upper bound of $O(k^{2/3})$ (for constant ϵ). Of course, the difficulty is that we are not given these rectangles a priori, which intuitively could make the problem require more samples than ℓ_1 -closeness testing on a domain of size k^4 .

The lack of a priori knowledge of the witnessing rectangles is *the* major obstacle towards developing a closeness tester with sub-learning sample complexity. Overcoming this bottleneck necessitates the bulk of the new technical ideas developed here. To achieve this, at a very high-level, we will proceed to compute *some* small set of rectangles that capture a “non-trivial”⁵ fraction of the discrepancy (i.e., \mathcal{A}_k -distance) between \mathbf{p} and \mathbf{q} .

A simple but important observation in this context is the following: one should not expect that an *obliviously* selected (i.e., without drawing samples from the underlying distributions) set of rectangles suffices for this purpose. Indeed, this holds even for the one-dimensional setting: as was noted in [DKN15a], any obliviously chosen set of intervals may capture *no* discrepancy between a pair of adversarially chosen one-dimensional distributions even though they have large \mathcal{A}_k distance.

⁴In hindsight, given our sample complexity lower bound of $\Omega(k^{6/7})$, the fact that the witnessing rectangles are unknown implies that the \mathcal{A}_k closeness testing problem *provably* requires more samples.

⁵Quantitatively, the term “non-trivial” here means “a function of the form $\text{poly}_d(\epsilon)$ ”.

That is, it appears necessary to select rectangles using samples from the tested distributions. Note that, in any dimension d , one needs at least two points in \mathbb{R}^d to define an axis-aligned rectangle. In particular, given two sample points $x, y \in \mathbb{R}^d$, we consider the following natural rectangle defined by these points, namely

$$R_{x,y} \stackrel{\text{def}}{=} \{z = (z_1, \dots, z_d) \in \mathbb{R}^d \mid \min(x_i, y_i) \leq z_i \leq \max(x_i, y_i) \text{ for all } i \in [d]\} .$$

The main intuition behind this definition is the following. Suppose that we draw two samples x, y from the mixture $(1/2)(\mathbf{p} + \mathbf{q})$ (the uniform mixture of \mathbf{p} and \mathbf{q}), and they both happen to land in some rectangle R such that the discrepancy $|\mathbf{p}(R) - \mathbf{q}(R)|$ is non-trivial. Then, intuitively, the rectangle $R_{x,y}$ will capture (in expectation) a non-trivial fraction of the rectangle R , and therefore also a non-trivial fraction of the discrepancy between \mathbf{p} and \mathbf{q} within R . The latter statement turns out to be true (see Proposition 2.1) and its proof makes essential use of tools from Ramsey theory.

Before we provide an overview of the ideas required to prove Proposition 2.1, we explain how to leverage this statement to develop our closeness tester. Suppose that the \mathcal{A}_k -distance between \mathbf{p}, \mathbf{q} is ϵ . Then at the cost of increasing k and decreasing ϵ by at most a constant factor, we can without loss of generality assume that there exist k rectangles $\{R_i\}_{i=1}^k$, each of which has probability mass approximately $1/k$ and witnesses roughly ϵ/k discrepancy. If we draw m samples from each of \mathbf{p}, \mathbf{q} , approximately m^2/k of these rectangles will contain two samples. Given Proposition 2.1, we know that each pair of samples landing in some R_i can be used to define a rectangle that, with some non-trivial probability, captures a non-trivial fraction of the discrepancy between \mathbf{p} and \mathbf{q} within R_i .

A potential concern is how one would find the right set of rectangles defined by the sample points (i.e., that capture enough discrepancy). The statement of Proposition 2.1 only ensures the existence of such rectangles, but offers no clues on how one could reliably identify them. Perhaps the most natural approach is to try all possible sets of $\Theta(m^2/k)$ many rectangles defined by the coordinates of the sample points, and then run a standard ℓ_1 -closeness tester (on the corresponding reduced distributions) to compare the probability mass of \mathbf{p} and \mathbf{q} on the selected rectangles. Unfortunately, in addition to its computational intractability, it is not even clear whether this method can lead to *any* sample complexity sublinear in k . In particular, the standard analysis of the above strategy will apply the union bound on the failure probabilities of running the ℓ_1 -closeness tester on each possible reduced distribution (defined by each set of rectangles). Since there are at least $m^{\Omega(m^2/k)}$ many different ways to select the set of rectangles, this increases the sample complexity of the ℓ_1 -closeness tester by a factor of $\Omega(m^2/k)$, making it hopeless to achieve any sublearning sample complexity (even balancing the quantities m^2/k and m directly will give us $m = k$).

To circumvent this obstacle, we leverage an idea from [DKP19], that we term *Grid Covering* (see Definition 2.4). At a high level, we show that we can cover the set of all possible rectangles that can be defined by the sample coordinates — which we refer to as \mathcal{S} — by a carefully chosen subset of these rectangles — which we refer to as \mathcal{F} — such that each rectangle from \mathcal{S} can be expressed as the union of at most *polylogarithmically many* rectangles from \mathcal{F} . Moreover, \mathcal{F} will be constructed to have the subtle property that any point in \mathbb{R}^d is contained in at most polylogarithmically many rectangles within the subset (in sharp contrast, in the worst case, a point may be included in a constant fraction of \mathcal{S}).

To take advantage of this property, we consider the notion of *induced distributions* (see Definition 2.6), $\mathbf{p}^{\mathcal{F}}, \mathbf{q}^{\mathcal{F}}$, on \mathcal{F} : to sample from $\mathbf{p}^{\mathcal{F}}$, we first draw a sample point $x \in \mathbb{R}^d$ from \mathbf{p} and return uniformly at random some rectangle from \mathcal{F} that includes x (and similarly for $\mathbf{q}^{\mathcal{F}}$). As a consequence of the aforementioned properties of \mathcal{F} , the discrepancy (under some appropriate metric)

between $\mathbf{p}^{\mathcal{F}}$ and $\mathbf{q}^{\mathcal{F}}$ will shrink by at most a polylogarithmic factor compared to the discrepancy between \mathbf{p} and \mathbf{q} captured by the best $\Theta(m^2/k)$ rectangles from our original collection of rectangles \mathcal{S} (defined using the sample points); see Lemma 2.7. Importantly, the new pair of distributions are both discrete, and the discrepancies between them will be supported on a small number of domain elements. Therefore, one could hope to apply techniques from “standard” ℓ_1 -closeness testing of discrete distributions from there on. While this turns out to be manageable, we emphasize that the induced distributions still have very large support size. Hence, a direct application of ℓ_1 -closeness testing on an arbitrary discrete domain is not sufficient for our purposes. We will return to this issue when we analyze the sample complexity of our tester in detail.

It remains to show correctness of this scheme. That is, we want to establish that there exists a small set of rectangles defined by the sample points which capture a non-trivial amount of discrepancy between \mathbf{p} and \mathbf{q} with high constant probability. To show this, we return to $\{R_i\}_{i=1}^k$ — a set of k rectangles which witness the \mathcal{A}_k distance between \mathbf{p} and \mathbf{q} . We will prove that for each of these rectangles R_i , if two samples $x, y \in \mathbb{R}^d$ are drawn from R_i , there is a non-negligible probability that $R_{x,y}$ — the rectangle defined by x and y — captures a non-trivial fraction of the discrepancy in R_i . (see Proposition 2.1 and its proof in Section 2.2).

As a starting point to achieve this, we show that if two sample points x, y are drawn from the restriction of $(1/2)(\mathbf{p} + \mathbf{q})$ to a rectangle R , there is a decent probability that $R_{x,y}$ will capture a non-trivial fraction of the mass of $(1/2)(\mathbf{p} + \mathbf{q})$ in R (see Lemma 2.9). This statement turns out to be *essentially equivalent* to a result in Ramsey theory shown by De Bruijn that can be viewed as a generalization of the classical Erdős-Szekeres theorem (see Fact 2.10). In particular, De Bruijn showed that given N points in \mathbb{R}^d (for N at least doubly exponential in d), there exists a triplet (x, y, z) of these points such that one of the points z is inside the rectangle $R_{x,y}$ defined by the other two. This statement provides us with the desired discrepancy result for the special case that one of \mathbf{p}, \mathbf{q} has non-trivial probability mass in the rectangle R while the other has mass zero.

To prove the desired discrepancy result for the general case, we introduce and leverage the notion of *discrepancy density* of a set $S \subset \mathbb{R}^d$, defined to be the discrepancy between \mathbf{p}, \mathbf{q} in S divided by the total mass assigned by \mathbf{p} and \mathbf{q} in S (see Definition 2.12). At a high level, our analysis proceeds as follows. We define an iterative process that selects rectangles with increasing discrepancy density. As the discrepancy density approaches one, the situation qualitatively resembles the case that only one of \mathbf{p}, \mathbf{q} assigns non-zero mass to the rectangle. We now provide some further details of the process. Note that if in expectation the rectangle $R_{x,y}$ — defined by random points x, y from R — captures a non-trivial amount of discrepancy between \mathbf{p}, \mathbf{q} in R , we are done. Otherwise, there exists a rectangle R_{x^*,y^*} such that the probability masses of \mathbf{p} and \mathbf{q} in R_{x^*,y^*} differ by a negligible amount. As a result, since the probability masses of \mathbf{p} and \mathbf{q} within R_{x^*,y^*} are approximately the same, the *complement* of R_{x^*,y^*} , which we denote by $S := R \setminus R_{x^*,y^*}$, must have higher discrepancy density between \mathbf{p} and \mathbf{q} . Since the complement S can be shown to be a union of a small number of axis-aligned rectangles (see Claim 2.13), we can select one of these rectangles to restart the process. By iterating this procedure, we obtain a sequence of rectangles whose discrepancy densities increase monotonically until we reach the case that a random pair of points drawn from one of these rectangles can capture a non-trivial amount of discrepancy between \mathbf{p} and \mathbf{q} in expectation.

Up to this point, we have summarized the key ideas needed for the correctness analysis of our closeness tester. We now proceed to describe the tester in more detail and provide a sketch of its sample complexity. Using Proposition 2.1 and (an adaptation of) the grid-covering approach of [DKP19], we obtain a pair of discrete induced distributions (that we can simulate access to based on the samples drawn) such that they have $\text{poly}_d(\epsilon) m^2/k^2$ discrepancy concentrated over approximately m^2/k domain elements (up to polylogarithmic factors). Leveraging the guarantees

of the pair of induced distributions we have constructed, it is tempting to apply the so called $\ell_{1,k}$ -tester from [DKN17]⁶. In particular, given samples from a pair of discrete distributions, such a tester aims at distinguishing between the cases that the underlying distributions are equal versus far in $\ell_{1,k}$ -distance — i.e., there exist k domain elements such that the ℓ_1 -distance restricted to these elements is large. Due to the “sparsity assumption” on the discrepancies, the sample complexity of $\ell_{1,k}$ -closeness testing is comparable to that of standard ℓ_1 -closeness testing on a domain of size k , even though the actual domain size of the input distributions may be much larger. In particular, using the guarantees of the $\ell_{1,k}$ tester in a black-box manner, we can detect the existing discrepancy between the pair of induced distributions obtained with sample size approximately

$$(m^2/k)^{2/3} / (\text{poly}_d(\epsilon) m^2/k^2)^{4/3} + (m^2/k)^{1/2} / (\text{poly}_d(\epsilon) m^2/k^2)^2 ,$$

where m is the initial number of samples drawn to construct the rectangles. Balancing the number of samples used for defining the rectangles and the number of samples used for detecting the discrepancy, we obtain that $m = \tilde{\Theta}_d(k^{7/8}/\text{poly}_d(\epsilon))$ suffices. This sample upper bound is strongly sub-linear in k , but it turns out (in hindsight) to provide a sub-optimal dependence on k .

Intuitively, the reason that the above guarantee turns out to be sub-optimal is the following. There would be a key property of our underlying discrete distributions left unused if we were to apply the guarantees of $\ell_{1,k}$ -testing in a black-box manner. Specifically, since the k rectangles that witness the \mathcal{A}_k -distance between \mathbf{p} and \mathbf{q} are themselves each of probability mass at most $O(1/k)$, the rectangles defined by our sample points (that capture non-trivial discrepancies) will also each be of mass at most $O(1/k)$. This in turn implies that in the end we only need to detect discrepancies supported on a few number of *light* domain elements (i.e., domain elements with small probability masses in the constructed discrete distributions). By carefully incorporating this additional property (i.e., that the bins witnessing discrepancies are themselves of small probability mass) into the analysis of the $\ell_{1,k}$ -tester, we obtain an improved sample complexity upper bound of

$$(m^2/k)^{2/3} / (\text{poly}_d(\epsilon) m^2/k^2)^{4/3} .$$

See Lemma 2.8 for the new tester and its analysis. We believe that this tester — customized for detecting discrepancies supported on a small number of *light* domain elements — may be applicable in other scenarios, as it allows us to escape from some worst-case scenarios of $\ell_{1,k}$ testing. Finally, balancing the number of samples used for defining the rectangles and that of samples used for detecting the discrepancy gives us a sample bound of approximately $m = \tilde{\Theta}_d(k^{6/7}/\text{poly}_d(\epsilon))$.

Sample Complexity Lower Bound Our sample complexity lower bound applies specifically for 2-dimensional distributions. This suffices for us to conclude that our sample upper bound is nearly optimal as a function of k for any constant dimension $d > 1$.

The starting point of our sample lower bound technique is the lower bound for one-dimensional \mathcal{A}_k closeness testing shown in [DKN15a]. Specifically, we start by showing that it is no loss of generality to establish a lower bound for “order-based” testers, and then prove a lower bound for such testers. In the proceeding discussion, we elaborate on each of these steps.

We start by noting that most reasonable testers seem to only be able to take advantage of the ordering of the x -coordinates and the y -coordinates of the points they observe — and not the precise numerical values of these coordinates (see Definition 3.1). We call such a tester an *order-based tester*. Intuitively, this holds because the \mathcal{A}_k distance is invariant under applying a monotonic

⁶The original $\ell_{1,k}$ -tester is for identity testing; one can adapt these techniques to derive an $\ell_{1,k}$ -tester for closeness testing.

transformation to all of the x -coordinates or all of the y -coordinates, and only the ordering of these coordinates is invariant under all monotonic transformations. In fact, we show that if there exists a non order-based \mathcal{A}_k closeness tester on a domain of size N , we can use it to construct an order-based tester that has almost the same guarantees — albeit on a smaller domain (see Lemma 3.2). Hence, using our reduction, we can translate any sample complexity lower bound against order-based testers into one against general testers at the cost of increasing the domain size. To obtain the reduction, we show that for any 2-dimensional \mathcal{A}_k tester on a sufficiently large domain there exists a large subset of its domain such that if the samples are drawn from the subdomain, the general tester’s output will depend only on the order of the samples. In other words, restricted to this subdomain, the tester becomes exactly an order-based tester.

The argument itself resembles the one in [DKN15a]. The key difference is that, due to the tester being 2-dimensional, the structure of the order information becomes much more complicated. More specifically, there is now order information from both of the dimensions. To deal with this issue, we need to take a two-fold approach. Namely, we need to first select a subset of coordinates in the first dimension to make the tester’s output independent of the samples’ order information in the first dimension, and then adaptively select the subsets of coordinates in the second dimension to hide the remaining order information (see Lemma 3.2).

For order-based testers, we construct families of distributions that are hard to distinguish. Lying in the center of the construction are two small gadgets, each consisting of a pair of distributions. We denote the two gadgets as \mathcal{Y} and \mathcal{N} respectively. In the \mathcal{Y} gadget, the two distributions are both uniform distributions supported on the edges of a square, whose diagonals are parallel to the x and y axis respectively (which we term a “diagonal square”). In the \mathcal{N} gadget, one distribution is distributed uniformly over a randomly chosen pair of parallel edges of the square, and the other one is distributed uniformly over the remaining two edges. The key point is that though the two distributions in the \mathcal{Y} gadget are identical and the distributions in the \mathcal{N} gadget have \mathcal{A}_k distance equal to one (even for $k = 4$), we show that no order-based tester can distinguish between the two gadgets when fewer than three samples are drawn (see Section 3.2).

To construct the full hard instance, we replicate the gadgets many times in a fairly standard way. In particular, we let \mathbf{p}, \mathbf{q} have their supports in several “boxes”. If the tester draws m samples, to introduce “noise”, we produce roughly m heavy boxes on which \mathbf{p}, \mathbf{q} are identical. We also have k light boxes each with mass approximately ϵ/k on which \mathbf{p}, \mathbf{q} either use the construction of the \mathcal{Y} gadget, and are therefore identical (if we want to construct $\mathbf{p} = \mathbf{q}$); or they use the construction of the \mathcal{N} gadget, and are therefore far from each other (if we want to construct $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} = \epsilon$). As we have discussed, observing up to three samples from any of the light boxes gives an order-based tester no information regarding which case one is in. In other words, one will only gain information from light boxes with at least four samples; note that there will only exist approximately m^4/k^3 such boxes if one draws m samples. Additionally, the m heavy boxes will “add noise” on the order of \sqrt{m} , and thus one can only distinguish between the two cases if $m^4/k^3 \gg m^{1/2}$ (or equivalently $m \gg k^{6/7}$). This heuristic argument can be made rigorous with an appropriate use of information theory (see Section 3.3).

A disadvantage of the above proof technique is that the Ramsey theory argument (used in the first step) only applies if the domain is extremely large. Using an enhancement of the technique from [DKN17], we can reduce this to domains of doubly exponential size in k (see Theorem 3.8). To achieve this, we need to modify our square-diagonal construction so that three samples provides little information to the tester even when the numerical values of these samples are also revealed. To do this, we show that by applying carefully chosen random functions to the x - and y - coordinates, we can effectively obscure almost all non-order-based information contained in any set of three samples. For the univariate case, [DKN17] showed that for two samples, applying a random *affine*

transformation can obscure both the difference and the average of a pair of points. However, when there are three points $a < b < c$, applying an affine transformation preserves the value of $(a - c)/(b - c)$. Hence, a non-trivial amount of information may be retrieved from the tester by computing this quantity, even if a random affine transformation is applied. To address this issue in our two-dimensional setting, we will apply an exponential function $x \mapsto \exp(\exp(\lambda)x)$, where λ is a carefully chosen uniform variable. Then, if a, b, c are not too close, $(a - c)/(b - c)$ will be exponentially close to $\exp(\exp(\lambda)(a - b)) = \exp(\exp(\lambda + \log(a - b)))$. When λ is large compared to $\log(a - b)$, the ratio $(a - c)/(b - c)$ will therefore have roughly the same distribution of outputs, independent of a, b, c . As a result, the transformation effectively hides any information encoded by the ratio $(a - b)/(b - c)$. Afterwards, we can mirror the analysis from [DKN17] to apply a suitable random affine transformation to hide all of the remaining information. The details of the construction and its analysis can be found in Section 3.4.

1.3 Basic Notation

For $n \in \mathbb{Z}_+$, we denote $[n] \stackrel{\text{def}}{=} \{1, \dots, n\}$. We will use \mathbb{S}_m for the set of all permutations over m distinct elements. Given $m > 0$, we use $\text{Poi}(m)$ to denote the Poisson distribution with mean m .

An axis-aligned rectangle R is a set in \mathbb{R}^d that can be represented as the product of d intervals I_1, \dots, I_d , i.e., $R = \prod_{i=1}^d I_i$. Given $x, y \in \mathbb{R}^d$, the axis-aligned rectangle defined by x, y is the set $R_{x,y} \stackrel{\text{def}}{=} \{z \in \mathbb{R}^d \mid \min(x_i, y_i) \leq z_i \leq \max(x_i, y_i) \text{ for all } i \in [d]\}$.

We will use \mathbf{p}, \mathbf{q} to denote the probability density functions of our distributions (or probability mass functions for discrete distributions). For discrete distributions \mathbf{p}, \mathbf{q} over $[n]$, their ℓ_1 and ℓ_2 distances are $\|\mathbf{p} - \mathbf{q}\|_1 \stackrel{\text{def}}{=} \sum_{i=1}^n |\mathbf{p}(i) - \mathbf{q}(i)|$ and $\|\mathbf{p} - \mathbf{q}\|_2 \stackrel{\text{def}}{=} \sqrt{\sum_{i=1}^n (\mathbf{p}(i) - \mathbf{q}(i))^2}$. For density functions $\mathbf{p}, \mathbf{q} : \mathbb{R}^d \mapsto \mathbb{R}_+$, we have $\|\mathbf{p} - \mathbf{q}\|_1 \stackrel{\text{def}}{=} \int_{\mathbb{R}^d} |\mathbf{p}(x) - \mathbf{q}(x)| dx$. The total variation distance between distributions \mathbf{p}, \mathbf{q} is defined to be $d_{\text{TV}}(\mathbf{p}, \mathbf{q}) = \frac{1}{2} \|\mathbf{p} - \mathbf{q}\|_1$. Let $R \subset \mathbb{R}^d$ be a subset of the domain of \mathbf{p} . We denote by $\mathbf{p}|_R$ the conditional distribution of \mathbf{p} restricted to R , i.e., $\mathbf{p}|_R(x) = \mathbf{p}(x) / \int_R \mathbf{p}(x) dx$ for $x \in R$. Let $\mathcal{R} = \{R_1, \dots, R_k\}$ be a collection of disjoint sets $R_i \subseteq \mathbb{R}^d$. The *reduced measure* corresponding to \mathbf{p} and \mathcal{R} , which we denote by $\mathbf{p}^{\mathcal{R}}$, is a discrete measure on $[k]$ defined as $\mathbf{p}_i^{\mathcal{R}} = \mathbf{p}(R_i)$ for $i \in [k]$.

1.4 Organization

The structure of this paper is as follows: In Section 2 we develop the analysis tools required to design and analyze our closeness tester. Section 3 contains our sample complexity lower bound. In Section 4, we provide some conclusions and open problems.

2 Closeness Testing Algorithm

In this section, we describe and analyze our multivariate \mathcal{A}_k -closeness tester. The structure of this section is as follows: In Section 2.1, we present our algorithm and its analysis. The proof of our main structural result (Proposition 2.1) which relies on Ramsey theory is given in Section 2.2. In Section 2.3, we describe and analyze our new closeness tester for discrete distributions which detects discrepancies supported on a small number of light domain elements (Lemma 2.8). Finally, Section 2.4 describes some applications of our \mathcal{A}_k closeness tester to test closeness of structured distributions under the total variation distance.

2.1 The Tester and its Analysis

We start with an overview of our algorithmic approach followed by a detailed pseudo-code and analysis of our tester.

Overview of Algorithmic Approach Let $\mathcal{R} = \{R_i\}_{i=1}^k$ be a collection of k disjoint rectangles which witness the \mathcal{A}_k -distance between \mathbf{p}, \mathbf{q} ⁷. The main technical obstacle of \mathcal{A}_k closeness testing is that the algorithm does not know (a priori) such a collection of rectangles. To circumvent this issue, we draw samples from \mathbf{p}, \mathbf{q} and use the obtained information to construct a set of rectangles that capture a non-trivial amount of discrepancy between the underlying distributions. A natural way to construct our rectangles is as follows. Given a collection of sample points from \mathbf{p} and \mathbf{q} , we group these points into disjoint pairs and make our rectangles be those defined by the corresponding pairs.

Note that the number of ways to group the sample points into disjoint pairs scales exponentially with the number of samples drawn. But before we discuss how the grouping is done in our algorithm, we need to prove that this approach can work *in principle*, i.e., that if one draws sufficiently many samples, there exists a *small* set of rectangles (each defined by pairs of sample points) that capture enough discrepancy between \mathbf{p} and \mathbf{q} .

Let x, y be two samples drawn from the mixture $(1/2)(\mathbf{p} + \mathbf{q})$. Conditioned on the event that x, y both land in some rectangle $R \in \mathcal{R}$ of the witnessing partition, we show that $R_{x,y}$ — the rectangle defined by x, y — will in expectation capture a non-trivial amount of the discrepancy in R . The formal statement is specified in Proposition 2.1 and its proof is given in Section 2.2.

By applying Proposition 2.1 to each rectangle $R_i \in \mathcal{R}$, one can show the existence of a collection of $k' = O(k)$ rectangles defined by the sample points which capture enough discrepancy between \mathbf{p}, \mathbf{q} (Lemma 2.3). It then remains to find these rectangles and invokes an appropriate closeness testing procedure to compare the probability mass of \mathbf{p}, \mathbf{q} on them. Trying all possible collections of rectangles defined by the sample points is certainly not computationally feasible. Even worse, the natural analysis of this brute-force strategy would require one to union-bound the failure probabilities of the closeness testing steps executed on each possible collection of rectangles. As the number of possible collections scales exponentially with the size of the collection, i.e., k' , each individual closeness testing routine is only allowed to fail with exponentially small probability, making the sample complexity of this approach at least linear in k .

We instead follow an approach inspired by the idea of a *Good Oblivious Covering* in [DKP19]. In particular, we consider a sub-collection of rectangles defined by the coordinates of the sample points that form a nice “cover” of all possible such rectangles. We then proceed to define the notion of “induced” distributions of \mathbf{p}, \mathbf{q} on the cover such that the two corresponding induced distributions have large ℓ_2 -discrepancy supported on a small number of domain elements if and only if there exists a collection of rectangles defined by the sample points over which the probability mass of \mathbf{p}, \mathbf{q} differ significantly. Then, applying a novel variant of the $\ell_{1,k}$ -tester from [DKN17] (see Lemma 2.8) yields our final tester.

We are now ready to proceed with the details of the proof.

Discrepancy from Random Points Let R be an axis-aligned rectangle such that $\mathbf{p}(R)$ and $\mathbf{q}(R)$ differ substantially and x, y be sample points drawn from $\frac{1}{2}(\mathbf{p} + \mathbf{q})|_R$, the uniform mixture distribution between \mathbf{p}, \mathbf{q} restricted to R . We consider the rectangle defined by x, y , which we denote

⁷Note that such a collection is not necessarily unique.

by $R_{x,y}$. Our main structural result, serving as the direct motivation for our algorithm, shows that $R_{x,y}$ captures non-trivial amount of discrepancy between \mathbf{p} and \mathbf{q} with non-trivial probability.

Proposition 2.1 (Random Point Discrepancy). *Let \mathbf{p}, \mathbf{q} be distributions on \mathbb{R}^d and R be an axis-aligned rectangle $R \subset \mathbb{R}^d$ satisfying $|\mathbf{p}(R) - \mathbf{q}(R)| \geq \epsilon(\mathbf{p}(R) + \mathbf{q}(R))$. Let x, y be random points sampled from $(1/2)(\mathbf{p} + \mathbf{q})|_R$. Then there exists a number $\alpha_d = Cd^22^{2^{d+1}}$, for some sufficiently large universal constant $C > 0$, such that $\mathbf{E}[|\mathbf{p}(R_{x,y}) - \mathbf{q}(R_{x,y})|] \geq \epsilon^{\alpha_d}(\mathbf{p}(R) + \mathbf{q}(R))$.*

The proof of Proposition 2.1 makes essential use of Ramsey theory and is one of the main technical contributions of this work. We defer its proof to Section 2.2.

Here we comment on the quantitative aspects of this result. Specifically, it is not clear whether the ϵ^{α_d} multiplicative factor in the right hand side of the final inequality is best possible. It is a plausible conjecture that the optimal dependence is $\text{poly}(\epsilon)$ — independent of the dimension d (see Question 4.1). Such an improvement would directly improve the sample complexity of our closeness tester, as a function of ϵ .

Existence of Witnessing Grid-aligned Rectangles We begin with an assumption that simplifies our analysis: the cumulative density function of each coordinate of \mathbf{p} or of \mathbf{q} is continuous. We will eventually remove the assumption in the proof of our main theorem. Suppose that $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} \geq \epsilon$. Then there exists a collection of k disjoint axis-aligned rectangles $R_1, R_2, \dots, R_k \subseteq \mathbb{R}^d$ such that $\sum_i |\mathbf{p}(R_i) - \mathbf{q}(R_i)| \geq \epsilon$. By Proposition 2.1, if two sample points x, y happen to land in the same rectangle R_i , the rectangle $R_{x,y}$ they define will capture a non-trivial fraction of discrepancy in R_i . For this reason, we restrict our attention to rectangles lying on the *sample-point grid* defined below.

Definition 2.2 (Sample-Point Grid). Let $S = \{x^{(1)}, \dots, x^{(m)}\} \subset \mathbb{R}^d$ be a set of sample points such that no two points overlap in any of their coordinates, i.e., $x_\ell^{(i)} \neq x_\ell^{(j)}$ for all $i \neq j \in [m]$ and $\ell \in [d]$. The *sample-point grid* G_S (with respect to S) is the set of all points $z \in \mathbb{R}^d$ such that the i -th coordinate z_i is chosen from the set $\{x_i^{(1)}, \dots, x_i^{(m)}\}$. Given an axis-aligned rectangle R , we say that R is a *grid-aligned rectangle* with respect to G_S if all its vertices are grid-points from G_S .

Let G_S be a sample-point grid with respect to a collection of sufficiently many i.i.d. samples from $(1/2)(\mathbf{p} + \mathbf{q})$. We first show that, with high constant probability, there exist $O(k)$ many rectangles aligned with G_S that capture enough discrepancy between \mathbf{p}, \mathbf{q} in ℓ_2 distance.

Lemma 2.3 (Existence of a Small Set of Witnessing Grid-aligned Rectangles). *Let $\alpha_d > 0$ be as defined in Proposition 2.1. Let \mathbf{p}, \mathbf{q} be distributions over \mathbb{R}^d satisfying $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} \geq \epsilon$. Let S be a set of $\text{Poi}(m)$ i.i.d. samples from $(1/2)(\mathbf{p} + \mathbf{q})$, where $k > m \geq C\sqrt{k}/(\epsilon/4)^{2\alpha_d}$, for some sufficiently large universal constant $C > 0$, and G_S be the sample-point grid defined by these points. With probability at least $9/10$, there exist $k' \leq 3k$ disjoint grid-aligned rectangles $\tilde{R}_1, \dots, \tilde{R}_{k'}$ with respect to G_S satisfying the following:*

- (i) $\mathbf{p}(\tilde{R}_i) + \mathbf{q}(\tilde{R}_i) \leq O(1/k)$ for all $i \in [k']$, and
- (ii) $\sum_{i=1}^{k'} \left(\mathbf{p}(\tilde{R}_i) - \mathbf{q}(\tilde{R}_i) \right)^2 \geq \Omega((\epsilon/4)^{2\alpha_d} m^2/k^3)$.

Proof. Let R_1, \dots, R_k be a collection of k axis-aligned rectangles which realize the \mathcal{A}_k -distance between \mathbf{p}, \mathbf{q} . Namely, it holds $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} = \sum_{i=1}^k |\mathbf{p}(R_i) - \mathbf{q}(R_i)|$. For convenience, for each rectangle R_i , we will denote $v_i \stackrel{\text{def}}{=} (\mathbf{p}(R_i) + \mathbf{q}(R_i)), \epsilon_i \stackrel{\text{def}}{=} |\mathbf{p}(R_i) - \mathbf{q}(R_i)|/v_i$. We first perform

some preliminary simplifications to make sure that $v_i = O(1/k)$ and $\epsilon_i \geq \epsilon/4$. Given $v_i > 1/k$, we can subdivide R_i into $\lfloor v_i k \rfloor$ sub-rectangles evenly along the first coordinate according to the cumulative density function of the first coordinate of $(1/2)(\mathbf{p} + \mathbf{q})$. We next discard any rectangles R_i such that $\epsilon_i < \epsilon/4$, which leads to us losing at most $\sum_i (\mathbf{p}(R_i) + \mathbf{q}(R_i))\epsilon/4 \leq \epsilon/2$ discrepancy. In summary, after these operations, we will have a collection of $\tilde{k} \leq 3k$ rectangles $R_1, \dots, R_{\tilde{k}}$ such that for each rectangle R_i in the collection we have that $v_i \leq 1/k$, $\epsilon_i \geq \epsilon/4$, and $\sum_{i=1}^{\tilde{k}} v_i \epsilon_i \geq \epsilon/2$.

Let S be the set of $\text{Poi}(m)$ many i.i.d. samples drawn and G_S be the corresponding sample-point grid. We define the random variable Y_i as follows: if exactly two samples $x, y \in S$ fall in the same rectangle R_i for $i \in [\tilde{k}]$, then $Y_i = (\mathbf{p}(R_{x,y}) - \mathbf{q}(R_{x,y}))^2$; otherwise, $Y_i = 0$. By the definition of Y_i , we know that if $Y_i > 0$, then there exists some rectangle $\tilde{R} \subset R_i$ aligned with G_S such that $Y_i = (\mathbf{p}(\tilde{R}) - \mathbf{q}(\tilde{R}))^2$. Hence, $\sum_{i=1}^{\tilde{k}} Y_i$ is always a lower bound on the discrepancy collected by the best collection of at most $\tilde{k} \leq 3k$ rectangles aligned with the grid G_S for any instance of the set S . Consequently, to prove the lemma, it suffices to show that $\sum_{i=1}^{\tilde{k}} Y_i \geq \Omega((\epsilon/4)^{2\alpha_d} m^2/k^3)$ with probability at least $9/10$.

Consider the event E_i that exactly two sample points land inside R_i . Then it is easy to see that

$$\Pr[E_i] = \Pr[\text{Poi}(mv_i/2) = 2] = \Theta(1) (mv_i)^2.$$

Conditioned on the event E_i , Y_i is equal to $(\mathbf{p}(R_{x,y}) - \mathbf{q}(R_{x,y}))^2$, where x, y are two random points from $\frac{1}{2}(\mathbf{p} + \mathbf{q})|_{R_i}$. By our preliminary simplification, we have that $|\mathbf{p}(R_i) - \mathbf{q}(R_i)| \geq (\epsilon/4)(\mathbf{p}(R_i) + \mathbf{q}(R_i))$. Hence, applying Proposition 2.1, we obtain

$$\mathbf{E}_{x,y \sim \frac{1}{2}(\mathbf{p} + \mathbf{q})|_{R_i}} [|\mathbf{p}(R_{x,y}) - \mathbf{q}(R_{x,y})|] \geq \epsilon_i^{\alpha_d} v_i.$$

Combining this with Jensen's inequality then gives that

$$\mathbf{E}_{x,y \sim \frac{1}{2}(\mathbf{p} + \mathbf{q})|_{R_i}} [(\mathbf{p}(R_{x,y}) - \mathbf{q}(R_{x,y}))^2] \geq \left(\mathbf{E}_{x,y \sim \frac{1}{2}(\mathbf{p} + \mathbf{q})|_{R_i}} [|\mathbf{p}(R_{x,y}) - \mathbf{q}(R_{x,y})|] \right)^2 \geq \epsilon_i^{2\alpha_d} v_i^2.$$

Since Y_i conditioned on the event E_i is distributed as $(\mathbf{p}(R_{x,y}) - \mathbf{q}(R_{x,y}))^2$ and Y_i is always non-negative, we thus have

$$\mathbf{E}[Y_i] \geq \mathbf{E}[Y_i|E_i] \Pr[E_i] \geq \Omega(1) (m v_i)^2 \epsilon_i^{2\alpha_d} v_i^2 \geq \Omega(1) \epsilon_i^{2\alpha_d} m^2 v_i^4. \quad (1)$$

Summing over all Y_i 's, we obtain

$$\sum_{i=1}^{\tilde{k}} \mathbf{E}[Y_i] \geq \Omega(1) \sum_{i=1}^{\tilde{k}} \epsilon_i^{2\alpha_d-4} m^2 (v_i \epsilon_i)^4 \geq \Omega(m^2) (\epsilon/4)^{2\alpha_d-4} \sum_{i=1}^{\tilde{k}} (v_i \epsilon_i)^4 \geq \Omega(m^2) (\epsilon/4)^{2\alpha_d} / k^3,$$

where the first inequality uses (Equation (1)), in the second inequality we bound from below ϵ_i by $\epsilon/4$, and in the third inequality we use the fact that $\sum_{i=1}^{\tilde{k}} a_i^4$ subject to $\sum_i a_i = A$, $a_i \geq 0$, is minimized at $a_i = A/b$.

On the other hand, since Y_i is defined to be non-zero only when there exist two points landing in R_i , and takes values at most v_i^2 , we have that $\mathbf{Var}[Y_i] \leq \Pr[E_i] v_i^4 = O(1)m^2 v_i^6$. Furthermore, since the Y_i 's are independently distributed, it follows that

$$\mathbf{Var} \left[\sum_{i=1}^{\tilde{k}} Y_i \right] = \sum_{i=1}^{\tilde{k}} \mathbf{Var}[Y_i] \leq O(1) m^2 \sum_{i=1}^{\tilde{k}} v_i^6 \leq O(m^2/k^5),$$

where in the last inequality we use that $v_i \leq O(1/k)$. We then have that $\mathbf{Var}[\sum_{i=1}^{\tilde{k}} Y_i] \leq (1/20) \left(\mathbf{E} \left[\sum_{i=1}^{\tilde{k}} Y_i \right] \right)^2$ as long as $m \geq C\sqrt{k}/(\epsilon/4)^{2\alpha_d}$ for some sufficiently large universal constant $C > 0$. Then, by Chebyshev's inequality, it follows that

$$\Pr \left[\sum_{i=1}^{\tilde{k}} Y_i \geq \Omega \left((\epsilon/4)^{2\alpha_d} m^2/k^3 \right) \right] \geq 9/10 .$$

This concludes the proof of Lemma 2.3. \square

Existence of Good Grid Covering By Lemma 2.3, there exist $O(k)$ grid-aligned rectangles that capture $\Omega_{d,\epsilon}(m^2/k^3)$ discrepancy between \mathbf{p}, \mathbf{q} . A naive tester may proceed as follows: Choose a set of $k' = O(k)$ disjoint rectangles aligned with the sample-point grid, and then perform closeness testing between the reduced distributions of \mathbf{p}, \mathbf{q} on the chosen rectangles. Then, with non-trivial probability, the chosen rectangles will capture enough discrepancy between \mathbf{p}, \mathbf{q} , and a standard closeness tester would suffice. Unfortunately, the number of ways to choose k' disjoint grid-aligned rectangles from a grid containing m^d grid points is at least $m^{\Omega(d k')}$. If we were to try all possible collections of k' disjoint grid-aligned rectangles, the resulting tester would likely be inefficient, as discussed in our techniques overview (Section 1.2), in terms of both sample complexity and computational complexity. To circumvent this issue, we will instead consider a carefully chosen subset of all grid-aligned rectangles with respect to the sample-point grid such that any grid-aligned rectangle can be decomposed into the union of a small number of rectangles from the family. Moreover, the subset is carefully constructed to have the subtle property that any point $x \in \mathbb{R}^d$ is contained in a small number of rectangles from the subset. This leads us to the concept of *Grid Covering*, which is based on the idea of *Good Oblivious Covering* (Definition 2 from [DKP19]).

Definition 2.4 (Grid Covering). Let m be a power of 2 and S be a set of $(m + 1)$ points in \mathbb{R}^d and G_S be the corresponding sample-point grid. A *grid covering* is a family of rectangles aligned with the sample-point grid, which we denote by $\mathcal{F}(G_S)$, satisfying the following:

- Any rectangle aligned with the grid can be represented as the union of at most $2^d \log^d m$ disjoint rectangles from $\mathcal{F}(G_S)$.
- Any point in \mathbb{R}^d is contained in exactly $\log^d m$ rectangles.

With a construction similar to that in [DKP19], we show that a Grid Covering always exists.

Lemma 2.5 (Existence of Grid Covering). *Let m be a power of 2, S be a set of $(m + 1)$ points from \mathbb{R}^d and G_S be the corresponding sample-point grid. Then there exists a grid covering $\mathcal{F}(G_S)$.*

Proof. For each coordinate $j \in [d]$, let $x_j^{(1)}, \dots, x_j^{(m+1)}$ be the j -th coordinates of the samples collected sorted in increasing order. We will refer to these numbers as the “grid values”. For each $i \in [\log m]$, we will define $\mathcal{I}_{j,i}$ as the partition of the interval $[x_j^{(1)}, x_j^{(m)}]$ into 2^i many sub-intervals such that each sub-interval in the partition contains an equal number of grid values. Then the rectangles in $\mathcal{F}(G_S)$ are those of the following form: for $j \in [d]$, an interval $I_j \in \bigcup_i \mathcal{I}_{j,i}$ is chosen and the rectangle is simply the product of the d selected intervals I_j .

Then it is easy to see that for any value $z \in [x_j^{(1)}, x_j^{(m+1)}]$, z is within $\log m$ intervals from $\bigcup_i \mathcal{I}_{j,i}$ (one interval from each partition). As a result, any point in \mathbb{R}^d is within $\log^d m$ rectangles from $\mathcal{F}(G_S)$.

Let R be a grid-aligned rectangle that is the product of the intervals I_1, \dots, I_d . Notice that the interval I_j can be decomposed into at most $2 \log m$ intervals from $\bigcup_i \mathcal{I}_{j,i}$ (at most 2 intervals from each partition $\mathcal{I}_{j,i}$). Thus, R can be decomposed into at most $2^d \log^d m$ rectangles from $\mathcal{F}(G_S)$. This completes the proof. \square

We next define the notion of the *induced distribution* of \mathbf{p}, \mathbf{q} on $\mathcal{F}(G_S)$.

Definition 2.6 (Induced Distribution). Given a distribution \mathbf{p} on \mathbb{R}^d and a family of sets \mathcal{F} whose elements are non-empty sets in \mathbb{R}^d that are not necessarily disjoint, the *induced distribution* $\mathbf{p}^{\mathcal{F}}$ is defined as follows. To draw a random sample from $\mathbf{p}^{\mathcal{F}}$, one first draws a random sample x from \mathbf{p} . If x does not belong to any set in \mathcal{F} , we return the special element \emptyset . Otherwise, we return a uniformly random set $S \in \mathcal{F}$ such that $x \in S$.

Notice that for a rectangle $R \in \mathcal{F}(G_S)$, we have that $\mathbf{p}^{\mathcal{F}(G_S)}(R) = \mathbf{p}(R) / \log^d m$, since each point appears in exactly $\log^d m$ rectangles from $\mathcal{F}(G_S)$. This then allows us to show that the ℓ_2 -discrepancy between the induced distributions $\mathbf{p}^{\mathcal{F}(G_S)}, \mathbf{q}^{\mathcal{F}(G_S)}$ must be non-trivial if the grid G satisfies the conclusion in Lemma 2.3. Specifically, we show:

Lemma 2.7. *Let m be a power of 2, S be a set of $(m+1)$ points in \mathbb{R}^d , and G_S be the corresponding sample-point grid. Moreover, suppose that the conclusion of Lemma 2.3 holds for G_S . Then there exists a subset of rectangles $H \subset \mathcal{F}(G_S)$ such that the following conditions hold:*

- (i) $|H| \leq 3k 2^d \log^d m$,
- (ii) $\mathbf{p}^{\mathcal{F}(G_S)}(R) + \mathbf{q}^{\mathcal{F}(G_S)}(R) \leq O(\log^{-d}(m)/k)$ for all $R \in H$, and
- (iii) $\sum_{R \in H} (\mathbf{p}^{\mathcal{F}(G_S)}(R) - \mathbf{q}^{\mathcal{F}(G_S)}(R))^2 \geq \Omega(1) 2^{-d} \log^{-3d}(k) (\epsilon/4)^{2\alpha_d} m^2/k^3$.

Proof. Since we assume that the conclusion in Lemma 2.3 is satisfied, there exist $k' \leq 3k$ many grid-aligned rectangles $R_1, \dots, R_{k'}$ (with respect to G_S) satisfying $\mathbf{p}(R_i) + \mathbf{q}(R_i) \leq O(1/k)$ for all $i \in [k']$, and

$$\sum_{i=1}^{k'} |\mathbf{p}(R_i) - \mathbf{q}(R_i)|^2 \geq \Omega((\epsilon/4)^{2\alpha_d} m^2/k^3). \quad (2)$$

By the definition of the grid covering, each R_i can be decomposed into at most $2^d \log^d(m)$ rectangles from $\mathcal{F}(G_S)$. Let H_i be the set of rectangles in $\mathcal{F}(G_S)$ into which R_i is decomposed. We will consider $H = \bigcup_i H_i$. It is clear that $|H| \leq 3k 2^d \log^d(m)$, which shows (i). Moreover, by the definition of the induced distribution, for any rectangle $R \in \mathcal{F}(G_S)$, we have $\mathbf{p}^{\mathcal{F}(G_S)}(R) = \log^{-d}(m) \mathbf{p}(R)$. Therefore, for each $R \in H$, it holds $\mathbf{p}^{\mathcal{F}(G_S)}(R) + \mathbf{q}^{\mathcal{F}(G_S)}(R) \leq \log^{-d}(m) O(1/k)$, which shows (ii).

It remains to show (iii). For each H_i , we have

$$\sum_{R \in H_i} (\mathbf{p}(R) - \mathbf{q}(R))^2 \geq 2^{-d} \log^{-d}(m) (\mathbf{p}(R_i) - \mathbf{q}(R_i))^2,$$

since $\sum_{R \in H_i} \mathbf{p}(R) = \mathbf{p}(R_i)$ (and the same for \mathbf{q}) and $|H_i| \leq 2^d \log^d(m)$. Combining this with the fact that $\mathbf{p}^{\mathcal{F}(G_S)}(R) = \log^{-d}(m) \mathbf{p}(R)$ and Equation (2) gives (iii). This completes the proof. \square

Unlike the naive testing approach (running an ℓ_1 -closeness tester on many different pairs of reduced distributions), we can now run a closeness tester just on the induced distributions $\mathbf{p}^{\mathcal{F}(G_S)}, \mathbf{q}^{\mathcal{F}(G_S)}$. A technical issue is that the domain size of $\mathcal{F}(G_S)$ is still very large. This makes

the black-box application of any ℓ_1 -closeness tester sample inefficient. Instead, we need to leverage the fact that a non-trivial fraction of the discrepancy between the two distribution is supported on a small number of elements. Interestingly, a tester with similar guarantees was developed in [DKN17] (see Lemma 2.5). However, as is, that tester is not sufficient for our purposes. More specifically, we essentially need to develop an ℓ_2 -version of it. The reason is that we need to distinguish between the cases $\mathbf{p} = \mathbf{q}$ versus the case that a non-trivial amount of ℓ_2 -discrepancy is supported on a few elements *that are themselves not too heavy*. In particular, using tools developed in [DK16] and [CDVV14], we show the following:

Lemma 2.8. *Let \mathbf{p}, \mathbf{q} be discrete distributions on $[n]$ and $s \in [n]$. Given $\epsilon > 0$ and $\text{Poi}(m)$ many i.i.d. samples from \mathbf{p}, \mathbf{q} , for $m = \Theta(\max(\epsilon^{-4/3}, \epsilon^{-2}/\sqrt{s}))$, there exists a tester FLATTEN-CLOSENESS that distinguishes between the following cases with probability at least 9/10: (a) $\mathbf{p} = \mathbf{q}$ versus (b) there exists a set of elements H of size s such that (i) $\sum_{i \in H} (p_i - q_i)^2 \geq \epsilon^2$ and (ii) $\max_{i \in S} (p_i + q_i)/2 \leq 1/s$.*

The proof of Lemma 2.8 builds on the approach of the $\ell_{1,k}$ tester. An important difference is that we now need to carefully incorporate the upper bound on the mass of the elements witnessing the discrepancy into the analysis. We defer the proof to Section 2.3.

We are now ready to present the pseudo-code of our testing algorithm and provide its proof of correctness.

Algorithm 1 Multidimensional \mathcal{A}_k Closeness Tester

Require: sample access to \mathbf{p}, \mathbf{q} on \mathbb{R}^d ; accuracy ϵ .

- 1: Set $m \leftarrow C' k^{6/7} \epsilon^{-2\alpha_d/3} \log^d(k) 2^{d/3}$, where C' is a sufficiently large constant and α_d is defined in Proposition 2.1.
 - 2: Draw $\text{Poi}(m)$ samples from $(1/2)(\mathbf{p} + \mathbf{q})$ and denote the set of samples by S .
 - 3: Add arbitrarily some distinct points to S such that $|S|$ is a power of 2.
 - 4: Construct the grid G_S (Definition 2.2) and the grid covering $\mathcal{F}(G_S)$ (Definition 2.4).
 - 5: Run the ℓ_2 -closeness tester of Lemma 2.8 on the induced distributions $\mathbf{p}^{\mathcal{F}(G_S)}, \mathbf{q}^{\mathcal{F}(G_S)}$ with accuracy parameter $\kappa = c 2^{-d} \log^{-3d} k (\epsilon/4)^{2\alpha_d} m^2/k^3$ for some sufficient small constant $c > 0$.
 - 6: Accept if that closeness tester accepts; otherwise Reject.
-

Proof of Upper Bound in Theorem 1.2. We first present the analysis assuming that \mathbf{p}, \mathbf{q} are continuous distributions and in the end give a preprocessing step to make sure the algorithm works for general distributions. Let $F(G_S)$ be defined as in Algorithm 1. If $\mathbf{p} = \mathbf{q}$, we have $\mathbf{p}^{\mathcal{F}(G_S)} = \mathbf{q}^{\mathcal{F}(G)}$. Therefore, the tester will accept with probability at least 2/3 by Lemma 2.8.

Next, we consider the case $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} > \epsilon$. We claim that with probability at least 9/10 there exist $k' \leq k$ grid-aligned rectangles (with respect to G_S) such that the conclusion of Lemma 2.3 is satisfied. Without the operation of adding extra points into S in Line 3 of Algorithm 1, the claim just follows from Lemma 2.3. Now it is easy to see that $R_1, \dots, R_{k'}$ are still grid-aligned rectangles with the extra points. Hence, the claim follows.

Condition on the event that the conclusion of Lemma 2.3 holds. We can then apply Lemma 2.7, which gives us that there exists a set of elements $H \subset \mathcal{F}(G_S)$ such that (i) $|H| \leq 2^d \log^d m k$ (ii) $\mathbf{p}(R) + \mathbf{q}(R) \leq O(1/k)$ for all $R \in H$, and (iii)

$$\sum_{R \in H} \left(\mathbf{p}^{\mathcal{F}(G_S)}(R) - \mathbf{q}^{\mathcal{F}(G_S)}(R) \right)^2 \geq \Omega(1) 2^{-d} \log^{-3d}(k) (\epsilon/4)^{2\alpha_d} m^2/k^3.$$

Then, applying Lemma 2.8, gives that the tester rejects with probability at least 9/10 given that

$$m \geq C \max \left(\kappa^{-2/3}, \kappa^{-1}/\sqrt{k} \right),$$

where $\kappa = \Theta(1) 2^{-d} \log^{-3d}(k) (\epsilon/4)^{2\alpha_d} m^2/k^3$ and $C > 0$ is a sufficiently large constant. One can verify that $m = C' k^{6/7} \epsilon^{-2\alpha_d/3} \log^d(k) 2^{d/3}$ suffices, where C' is a sufficiently large constant.

Now let us relax the assumption that the marginal distributions of \mathbf{p}, \mathbf{q} in each coordinate have continuous cumulative density functions. We begin with the observation that the algorithm's output essentially depends only on the order information of the sample points. That is, given two different sets of samples $S = \{x^{(1)}, \dots, x^{(m)}\}, \tilde{S} = \{\tilde{x}^{(1)}, \dots, \tilde{x}^{(m)}\}$ such that the relative orders of $x_j^{(1)}, \dots, x_j^{(m)}$ and $\tilde{x}_j^{(1)}, \dots, \tilde{x}_j^{(m)}$ are the same for each coordinate $j \in [d]$, the output of the algorithm will always be the same.

Based on this observation, we know that the algorithm will satisfy the same guarantee if we give it only the “rank” information of the samples. For $j \in [d]$, we sort $x_j^{(1)}, \dots, x_j^{(m)}$ in increasing order. We will denote by $\pi(j)_i$ the rank of $x_j^{(i)}$ in the sorted sequence. Then, for each $i \in [m]$, we replace the original sample with the new sample $\hat{x}^{(i)}$ defined as $\hat{x}_j^{(i)} = \pi(j)_i$.

If the marginal distributions of \mathbf{p} or \mathbf{q} are not continuous, we may observe multiple samples sharing the same value at some coordinates. Then, when computing the rank information of the samples, we will break ties uniformly at random. Now consider the distributions \mathbf{p}', \mathbf{q}' obtained by stretching any point-mass of their marginal distributions at any coordinate into an interval. If the algorithm takes samples from \mathbf{p}', \mathbf{q}' instead, the guarantees are satisfied, since \mathbf{p}', \mathbf{q}' are both continuous distributions and $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} = \|\mathbf{p}' - \mathbf{q}'\|_{\mathcal{A}_k}$. On the other hand, the order of samples taken from \mathbf{p}', \mathbf{q}' has the same distribution as the order of samples taken from \mathbf{p}, \mathbf{q} after we break ties uniformly at random. This then concludes the proof. \square

2.2 Proof of Proposition 2.1

Let R be an axis-aligned rectangle such that $|\mathbf{p}(R) - \mathbf{q}(R)| \geq \epsilon(\mathbf{p}(R) + \mathbf{q}(R))$. Let x, y be samples from $(1/2)(\mathbf{p} + \mathbf{q})|_R$ — the uniform mixture of \mathbf{p}, \mathbf{q} restricted to R . We want to show that in expectation over x, y the discrepancy $|\mathbf{p}(R_{x,y}) - \mathbf{q}(R_{x,y})|$ is large.

Warm-up: Special case $\mathbf{p}(R) > 0$ and $\mathbf{q}(R) = 0$. Towards establishing the desired statement, we first analyze the special case that $\mathbf{p}(R) > 0$ and $\mathbf{q}(R) = 0$. The proof for this case also serves as intuition regarding why selecting the interval $R_{x,y}$ is a good choice.

In this case, the discrepancy between \mathbf{p} and \mathbf{q} is simply $\mathbf{p}(R_{x,y})$ — the probability mass of $R_{x,y}$ with respect to \mathbf{p} . Therefore, whether $R_{x,y}$ captures enough discrepancy boils down to the following question: *Let x, y be random points drawn from an arbitrary distribution D over \mathbb{R}^d . What is the minimum amount of mass captured by the rectangle $R_{x,y}$ in expectation?* We show the quantity is indeed non-trivial.

Interestingly, the proof of this statement relies on a certain generalized version of the famous Erdős-Szekeres theorem. In particular, the generalized Erdős-Szekeres theorem bounds from above the minimum length of a sequence consisting of points in \mathbb{R}^d such that there exists a subsequence of points that is monotonic in each coordinate.

Lemma 2.9. *Let x, y be random samples independently drawn from a distribution D on \mathbb{R}^d . Then it holds $\mathbf{E}_{x,y \sim D}[D(R_{x,y})] \geq \beta_d$, where $\beta_d = \left(2^{2^{d-1}} + 1\right)^{-3}$. Moreover, there exists a distribution D such that $\mathbf{E}_{x,y \sim D}[D(R_{x,y})] \leq \frac{2}{2^{2^{d-1}}}$.*

We note that the above statement is qualitatively nearly tight as a function of d .

Proof of Lemma 2.9. To prove the lemma, we make essential use of the following generalized version of the Erdős-Szekeres theorem proved by De Bruijn.

Fact 2.10 (De Bruijn's Generalized Erdős-Szekeres Theorem, see [Kru53]). *Let $\psi(n, d)$ denote the least integer N such that every sequence of points $x^{(1)}, \dots, x^{(N)}$ in \mathbb{R}^d contains a monotonic subsequence $x^{(i_1)}, \dots, x^{(i_n)}$ of length n satisfying the following: for each coordinate $j \in [d]$, we have either that $x_j^{(i_1)} \leq \dots \leq x_j^{(i_n)}$ or that $x_j^{(i_1)} \geq \dots \geq x_j^{(i_n)}$. Then it holds $\psi(n, d) = (n - 1)^{2^d} + 1$.*

As an immediate corollary, we obtain the following:

Corollary 2.11. *Let $S \subset \mathbb{R}^d$ be a set of points with size $|S| \geq 2^{2^{d-1}} + 1$. Then there exists a triple $x, y, z \in S$ such that $z \in R_{x,y}$. Furthermore, there exists a set of points S of size $|S| = 2^{2^{d-1}}$ such that there is no triple $x, y, z \in S$ satisfying $z \in R_{x,y}$.*

Proof. Let S be an arbitrary set of m points in \mathbb{R}^d . Let $x^{(1)}, \dots, x^{(m)}$ be a sequence of points in \mathbb{R}^{d-1} obtained by (1) sorting the points in S based on their first coordinates, and (2) throwing away their first coordinates. Applying Fact 2.10 with $n = 3$ gives us that we will have a subsequence $x^{(i_1)}, x^{(i_2)}, x^{(i_3)}$ such that the points are either monotonically increasing or monotonically decreasing in each of the $(d - 1)$ coordinates *if and only if* $m \geq 2^{2^{d-1}} + 1$. Furthermore, by our construction of the sequence of $x^{(i)}$'s, the first coordinates of the corresponding points in S are always monotonically increasing. This concludes the proof. \square

It is worth noting that for a set of points $S \subset \mathbb{R}^d$ to contain a triple x, y, z such that $z \in R_{x,y}$, the size of S needs to be doubly exponential in d ; this bound is tight since the corollary is essentially equivalent to Fact 2.10, which is itself quantitatively tight. Now let S be the set of $2^{2^{d-1}}$ points such that there is no triple $x, y, z \in S$ satisfying $z \in R_{x,y}$ (Corollary 2.11 ensures the existence of such a set of points). Let D be the uniform distribution over S . One can see that $D(R_{x,y}) \leq 2/2^{2^{d-1}}$ for any $x, y \in S$. It hence follows that $\mathbf{E}_{x,y \sim D} [D(R_{x,y})] \leq 2/2^{2^{d-1}}$, showing that our lower bound is qualitatively tight.

To relate Lemma 2.9 to the Generalized Erdős-Szekeres theorem (Fact 2.10), we make the following observations: (i) the probability mass of $R_{x,y}$ under D is equal to the probability that a third random point z drawn from D happens to land in $R_{x,y}$, and (ii) drawing three random samples from D is equivalent to first drawing N random samples from D and then choosing 3 *distinct* points from these N points uniformly at random.

Let D_N be the empirical distribution obtained after drawing N i.i.d. samples from D . The observations above allow us to conclude that

$$\mathbf{E}_{x,y \sim D} [D(R_{x,y})] = \mathbf{Pr}_{x,y,z \sim D} [z \in R_{x,y}] \geq \mathbf{E}_{D_N} \left[\mathbf{Pr}_{\substack{x,y,z \sim D_N \\ \text{without replacement}}} [z \in R_{x,y}] \right].$$

If we have $N \geq 2^{2^{d-1}} + 1$, Corollary 2.11 guarantees the existence of a triple $x, y, z \in D_N$ such that $z \in R_{x,y}$. Hence, the probability in the last equation above is at least $1/N^3$. This completes the proof of Lemma 2.9. \square

General Case. We are now ready to handle the general case and complete the proof of Proposition 2.1. To do so, we leverage the concept of the *discrepancy density* defined below.

Definition 2.12 (Discrepancy Density). Let \mathbf{p}, \mathbf{q} be distributions over \mathbb{R}^d . For a set $S \subseteq \mathbb{R}^d$, we define the discrepancy density of S with respect to \mathbf{p}, \mathbf{q} as follows:

$$\rho(S; \mathbf{p}, \mathbf{q}) \stackrel{\text{def}}{=} 2 |\mathbf{p}(S) - \mathbf{q}(S)| / (\mathbf{p}(S) + \mathbf{q}(S)) .$$

The high-level intuition is the following. Let $R_{x,y}$ be the axis-aligned rectangle defined by x, y , where x, y are independent random samples drawn from $(1/2)(\mathbf{p} + \mathbf{q})|_R$. By Lemma 2.9, the probability mass of $R_{x,y}$ (with respect to the mixture distribution $(1/2)(\mathbf{p} + \mathbf{q})$) is a non-trivial fraction of the mass of R in expectation. If the discrepancy between \mathbf{p}, \mathbf{q} within $R_{x,y}$ is a non-trivial fraction of the mass of $R_{x,y}$, we are done. Otherwise, if we were to “remove” the region $R_{x,y}$ from R , we would discard about approximately equal amounts of \mathbf{p} mass and \mathbf{q} mass. Therefore, the discrepancy density of the remaining space, $\rho(R \setminus R_{x,y}; \mathbf{p}, \mathbf{q})$, must have increased. We can then carve the remaining space into at most $2d$ many axis-aligned sub-rectangles and pick a sub-rectangle with significantly higher discrepancy density to restart the process. When the discrepancy density approaches one, the situation qualitatively resembles the special case where $\mathbf{q}(R) = 0$ (or $\mathbf{p}(R) = 0$); and if the mass of $R_{x,y}$ is non-trivial, the discrepancy captured will also be non-trivial. The formal proof follows.

Proof of Proposition 2.1. For notational convenience, we will denote $D := (1/2)(\mathbf{p} + \mathbf{q})$. Let x, y be two sample points drawn from $D|_R$, the restriction of D to R . Then, by Lemma 2.9, it holds

$$\mathbf{E}_{x,y \sim D|_R} [D(R_{x,y})] \geq \beta_d D(R) ,$$

for some β_d depending only on d . We will use E to denote the event $\{D(R_{x,y}) \geq \beta_d D(R)/2\}$. Then we must have that $\mathbf{Pr}_{x,y \sim D|_R} [E] \geq \beta_d/2$, since otherwise $\mathbf{E}_{x,y \sim D|_R} [D(R_{x,y})]$ will be no more than $\beta_d D(R)$.

We consider two complementary cases. First, if

$$\mathbf{E}_{x,y \sim D|_R} [|\mathbf{p}(R_{x,y}) - \mathbf{q}(R_{x,y})||E] > \frac{\epsilon}{2} \mathbf{E}_{x,y \sim D|_R} [D(R_{x,y})|E] , \quad (3)$$

we will have

$$\begin{aligned} \mathbf{E}_{x,y \sim D|_R} [|\mathbf{p}|_R(R_{x,y}) - \mathbf{q}|_R(R_{x,y})|] &\geq \mathbf{E}_{x,y \sim D|_R} [|\mathbf{p}|_R(R_{x,y}) - \mathbf{q}|_R(R_{x,y})||E] \mathbf{Pr}_{x,y \sim D|_R} [E] \\ &\geq \epsilon/2 \mathbf{E}_{x,y \sim D|_R} [D|_R(R_{x,y})|E] \mathbf{Pr}_{x,y \sim D|_R} [E] \\ &\geq \epsilon \beta_d^2 D(R)/8 \end{aligned}$$

and we are done.

Otherwise, it holds

$$\mathbf{E}_{x,y \sim D|_R} [\epsilon/2 D(R_{x,y}) - |\mathbf{p}(R_{x,y}) - \mathbf{q}(R_{x,y})||E] \geq 0 . \quad (4)$$

Since we also condition on the event E , we know that there exists a rectangle $\tilde{R} \subset R$ such that

$$(\epsilon/2) D(\tilde{R}) - |\mathbf{p}(\tilde{R}) - \mathbf{q}(\tilde{R})| \geq 0 , D(\tilde{R}) \geq \beta_d D(R)/2 .$$

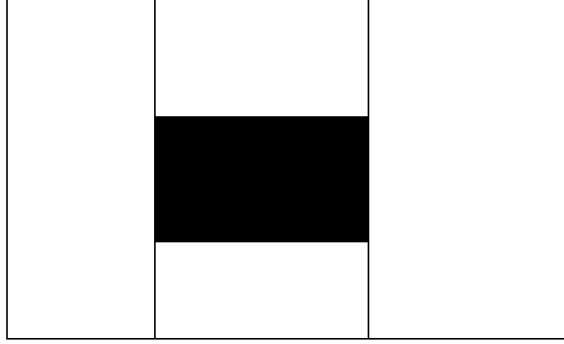


Figure 1: The black rectangle represents $\tilde{R} \subset R$ in \mathbb{R}^2 . One can see that there is a natural way to carve the remaining space $R \setminus \tilde{R}$ into four axis-aligned rectangles.

We consider the remaining space $R \setminus \tilde{R}$. We have that its discrepancy density satisfies

$$|\mathbf{p}(R \setminus \tilde{R}) - \mathbf{q}(R \setminus \tilde{R})| / D(R \setminus \tilde{R}) \geq \frac{\epsilon D(R) - \epsilon/2 D(\tilde{R})}{D(R) - D(\tilde{R})} = \epsilon(1 + \gamma),$$

where we denote $\gamma = \frac{D(R) - D(\tilde{R})/2}{D(R) - D(\tilde{R})} - 1$ for convenience. Notice that

$$\gamma = \frac{D(R) - D(\tilde{R})/2}{D(R) - D(\tilde{R})} - 1 = \frac{D(\tilde{R})/2}{D(R) - D(\tilde{R})} \geq \frac{(\beta_d/4) D(R)}{D(R)} = \beta_d/4,$$

where in the inequality above we bound below $D(\tilde{R})$ by $\beta_d D(R)/2$ by our choice of \tilde{R} . This then gives that $\gamma \geq \beta_d/4$.

It turns out that remaining space $R \setminus \tilde{R}$, can be carved into $2d$ many axis-aligned rectangles. An illustration of the $d = 2$ case is given in Figure 1.

Specifically, we show the following:

Claim 2.13. *Let $\tilde{R} \subseteq R$ be an axis-aligned rectangle. The set $R \setminus \tilde{R}$ can be decomposed into $2d$ axis-aligned rectangles R_1, \dots, R_{2d} .*

Proof. The proof proceeds via induction. The base case ($d = 2$) is clear, as shown in Figure 1. Assume that the statement holds for $d = k$. We proceed to show that it still holds for $d = k + 1$. Suppose that R is defined by points $x, y \in \mathbb{R}^{k+1}$ and \tilde{R} is defined by points $\tilde{x}, \tilde{y} \in \mathbb{R}^{k+1}$. We let R_{2k+1} be the rectangle that occupies the interval $[x_1, \tilde{x}_1]$ in the first dimension and occupies the same intervals as R in the other dimensions; similarly, let R_{2k+2} be the rectangle that occupies the interval $[\tilde{y}_1, y_1]$ in the first dimension and occupies the same intervals as R in the other dimensions. Then the remaining space $R \setminus (R_1 \cup R_2)$ lies entirely in the interval $[\tilde{x}_1, \tilde{y}_1]$ in the first dimension. We can then discard the first dimension. We denote the projection of $R \setminus (R_1 \cup R_2)$ into the remaining subspace \mathbb{R}^k as R' , and the projection of \tilde{R} as \tilde{R}' . We can then apply our inductive hypothesis on R' and \tilde{R}' to obtain $2k$ rectangles R'_1, \dots, R'_{2k} that live in the subspace of \mathbb{R}^k . Then let R_i to be the product of R'_i and the interval $[\tilde{x}_1, \tilde{y}_1]$ for $i = [2k]$. It is easy to verify that R_1, \dots, R_{2k+2} partitions the space $R \setminus \tilde{R}$. \square

We denote the rectangles obtained by applying the above claim to $R \setminus \tilde{R}$ as R_1, \dots, R_{2d} . Furthermore, we will denote $\gamma_i = \frac{|\mathbf{p}(R_i) - \mathbf{q}(R_i)|}{D(R_i) \epsilon} - 1$. In other words, $|\mathbf{p}(R_i) - \mathbf{q}(R_i)| / D(R_i) = (1 + \gamma_i) \epsilon_i$.

We claim that there exists a rectangle R_{i^*} in the remaining space such that

$$D(R_{i^*}) \geq \frac{\gamma \left(D(R) - D(\tilde{R}) \right)}{2d \gamma_{i^*}}, \gamma_{i^*} \geq \frac{\gamma}{2d}. \quad (5)$$

By definition of $D(R_i)$ and γ_i , we have

$$\sum_{i=1}^{2d} D(R_i) = D(R) - D(\tilde{R}), \quad (6)$$

$$\sum_{i=1}^{2d} D(R_i) \epsilon (1 + \gamma_i) \geq \left(D(R) - D(\tilde{R}) \right) \epsilon (1 + \gamma), \quad (7)$$

where the first equality follows from the fact that the rectangles form a partition of the remaining space, and the second equality follows from the fact that the sum of discrepancies in each rectangle must be at least the total discrepancy in the remaining space.

Substituting Equation (6) into Equation (7) and simplifying the result gives $\sum_{i=1}^{2d} D(R_i) \gamma_i \geq \left(D(R) - D(\tilde{R}) \right) \gamma$. Therefore, there exists i^* such that

$$D(R_{i^*}) \gamma_{i^*} \geq \left(D(R) - D(\tilde{R}) \right) \gamma / (2d).$$

Since $R_{i^*} \subseteq R \setminus \tilde{R}$, we must have that $D(R_{i^*}) \leq D(R) - D(\tilde{R})$, which implies that $\gamma_{i^*} \geq \gamma/2d$. On the other hand, we also have that

$$D(R_{i^*}) \geq \frac{\gamma}{2d} \frac{D(R) - D(\tilde{R})}{\gamma_{i^*}}.$$

This then establishes the existence of an i^* such that Equation (5) is satisfied.

We can inductively restart the process with $R' = R_{i^*}$ and $\epsilon' = \epsilon(1 + \lambda_{i^*})$. In each iteration, the discrepancy density must increase by at least a multiplicative factor of $(1 + \gamma/2d) \geq (1 + C \beta_d/2d)$, for some universal constant $C > 0$. Since the discrepancy density is at most one, the process must terminate in $O(d \beta_d^{-1} \log(1/\epsilon))$ many iterations, and we will eventually find some rectangle R^* such that Equation (3) is satisfied.

It remains to show that the mass $D(R^*)$ is bounded below. Suppose that in the t -th iteration, we start from the rectangle $R^{(t)}$ with discrepancy density $\epsilon^{(t)} := \frac{|\mathbf{p}(R^{(t)}) - \mathbf{q}(R^{(t)})|}{D(R^{(t)})}$ and end with the rectangle $R^{(t+1)} := R_{i^*}^{(t)}$ with discrepancy density $\epsilon^{(t+1)} := \epsilon^{(t)} \left(1 + \gamma_{i^*}^{(t)} \right)$. Denote by $\tilde{R}^{(t)}$ the rectangle discarded and $\epsilon^{(t+1/2)} := \epsilon^{(t)} \left(1 + \gamma^{(t)} \right)$ the discrepancy density of the remaining space $R^{(t)} \setminus \tilde{R}^{(t)}$. We analyze how much the mass of the rectangle can shrink in each iteration, as follows:

$$\begin{aligned} \frac{D(R_{i^*}^{(t)})}{D(R^{(t)})} &\geq \frac{\gamma^{(t)}}{2d} \left(1 - D(\tilde{R}^{(t)})/D(R^{(t)}) \right) \frac{1}{\gamma_{i^*}^{(t)}} \\ &\geq \frac{\gamma^{(t)}}{2d} \left(\frac{1 - D(\tilde{R}^{(t)})/D(R^{(t)})}{1 - D(\tilde{R}^{(t)})/(2D(R^{(t)}))} \right)^2 \frac{1}{\gamma_{i^*}^{(t)}} \\ &= \frac{\gamma^{(t)}}{2d} \left(\frac{1}{1 + \gamma^{(t)}} \right)^2 \frac{1}{\gamma_{i^*}^{(t)}} \\ &\geq \Omega(1) \frac{\beta_d}{d^3} \left(\frac{1}{1 + \gamma_{i^*}^{(t)}} \right)^3, \end{aligned}$$

where the first line uses our choice of $R_{i^*}^{(t)}$ such that Equation (5) is satisfied, the second line uses the elementary inequality $(1-x) \geq (1-x)^2/(1-x/2)^2$ for any $x \leq 1$, the third line uses the definition of $\gamma^{(t)}$, and the last line uses the facts $\gamma_{i^*}^{(t)} \geq \gamma^{(t)}/2d$ by our choice of $R_{i^*}^{(t)}$ such that Equation (5) is satisfied and $\gamma^{(t)} \geq \Omega(\beta_d)$ by our choice of $\tilde{R}^{(t)}$.

Notice that the discrepancy density increases by a multiplicative factor of $1 + \gamma_{i^*}^{(t)}$ in the t -th iteration. Thus, we have

$$\epsilon \prod_t \left(1 + \gamma_{i^*}^{(t)}\right) \leq 1. \quad (8)$$

Therefore, $D(R^*)$ is at least

$$\prod_t C \frac{\beta_d}{d^3} \left(\frac{1}{1 + \gamma_{i^*}^{(t)}}\right)^3 \geq \left(C \frac{\beta_d}{d^3}\right)^{O(d\beta_d^{-1} \log(1/\epsilon))} \epsilon^3 \geq \epsilon^{\tilde{O}(d\beta_d^{-1})},$$

where we used the fact that the process terminates in at most $O(d\beta_d^{-1} \log(1/\epsilon))$ iterations and Equation (8). This then shows that there exists a rectangle $R^* \subseteq R$ such that $D(R^*) \geq \epsilon^{\alpha_d}$ for some $\alpha_d = \tilde{O}(d\beta_d^{-1})$ and Equation (3) is satisfied. Then it holds

$$\mathbf{E}_{x,y \sim D|_{R^*}} \left[|\mathbf{p}(\tilde{R}) - \mathbf{q}(\tilde{R})| \right] \geq \epsilon/2 \frac{1}{2^{2d}} D(R^*)^2 \geq \epsilon/2 \frac{1}{2^{2d}} \epsilon^{2\alpha_d} D(R) \geq \epsilon^{\alpha'_d} D(R),$$

where $\alpha'_d = \tilde{O}(d\beta_d^{-1}) \leq C d^2 2^{2d+1}$, for some sufficiently large universal constant C . This concludes the proof of Proposition 2.1. □

2.3 Proof of Lemma 2.8

Given a discrete distribution \mathbf{p} , flattening [DK16] is the technique of using a small set of samples from \mathbf{p} to appropriately subdivide its bins (domain elements) aiming to reduce the ℓ_2 -norm of the distribution. Formally, the flattening technique yields what was described in [DK16] as a *split distribution*.

Definition 2.14 (Definition 2.4 from [DK16]). Given a distribution \mathbf{p} on $[n]$ and a multiset S of elements of $[n]$, define the *split distribution* p_S on $[n + |S|]$ as follows: For $1 \leq i \leq n$, let a_i denote 1 plus the number of elements of S that are equal to i . Thus, $\sum_{i=1}^n a_i = n + |S|$. We can therefore associate the elements of $[n + |S|]$ to elements of the set $B = \{(i, j) : i \in [n], 1 \leq j \leq a_i\}$. We now define a distribution \mathbf{p}_S with support B , by letting a random sample from \mathbf{p}_S be given by (i, j) , where i is drawn randomly from p and j is drawn randomly from $[a_i]$.

We will use the following basic facts about split distributions.

Fact 2.15 (Fact 2.5 and Lemma 2.6 from [DK16]). *Let \mathbf{p} and \mathbf{q} be probability distributions on $[n]$, and S a given multiset of $[n]$. Then: (i) We can simulate a sample from the split distributions \mathbf{p}_S or \mathbf{q}_S by taking a single sample from \mathbf{p} or \mathbf{q} , respectively. (ii) It holds $\|\mathbf{p}_S - \mathbf{q}_S\|_1 = \|\mathbf{p} - \mathbf{q}\|_1$. (iii) For any multisets $S \subseteq S' \subseteq [n]$, $\|\mathbf{p}_{S'}\|_2 \leq \|\mathbf{p}_S\|_2$. (iv) If S is obtained by drawing $\text{Poi}(m)$ samples from \mathbf{p} , then $\mathbf{E} \left[\|\mathbf{p}_S\|_2^2 \right] \leq 1/m$.*

We will also leverage the following ℓ_2 -distance estimator to develop our final tester.

Lemma 2.16 (Proposition 6 from [CDVV14]). *Let \mathbf{p} and \mathbf{q} be unknown distributions on $[n]$. There exists an algorithm that on input $n, \epsilon > 0$, and $b \geq \max(\|\mathbf{p}\|_2^2, \|\mathbf{q}\|_2^2)$, it draws $\text{Poi}(m)$ samples from \mathbf{p}, \mathbf{q} , where $m = \Theta(\sqrt{b}/\epsilon^2 + \sqrt{b} \|\mathbf{p} - \mathbf{q}\|_4^2/\epsilon^4)$, and with probability $3/4$ estimates $\|\mathbf{p} - \mathbf{q}\|_2$ up to accuracy $\pm\epsilon$.*

We can easily convert the above ℓ_2 -distance estimator to an ℓ_2 -closeness tester, which is more applicable to our setting.

Corollary 2.17. *Let \mathbf{p} and \mathbf{q} be unknown distributions on $[n]$. There exists an algorithm that on input $n, \epsilon > 0$, and $b \geq \max(\|\mathbf{p}\|_2^2, \|\mathbf{q}\|_2^2)$, the algorithm draws $\text{Poi}(m)$ samples from \mathbf{p}, \mathbf{q} , where $m = \Theta(\sqrt{b}/\epsilon^2)$, and with probability $3/4$ distinguishes between the cases $\mathbf{p} = \mathbf{q}$ versus $\|\mathbf{p} - \mathbf{q}\|_2 > \epsilon$.*

Let S be a multiset of $\text{Poi}(m)$ i.i.d. samples from $(1/2)(\mathbf{p} + \mathbf{q})$ for some $m \leq s/100$. First, we argue that the ℓ_2 -distance between \mathbf{p}, \mathbf{q} will not decrease by too much after flattening, by taking advantage of the fact that the ℓ_2 -discrepancy between \mathbf{p}, \mathbf{q} is supported on a few light elements.

Let X_i be the random variable denoting the number of samples from S landing in the i -th element, i.e., $X_i \sim \text{Poi}(m(\mathbf{p}_i + \mathbf{q}_i)/2)$. Then the expected discrepancy restricted to the elements from the set of elements witnessing the discrepancy, i.e., H , after flattening is at least

$$\sum_{i \in H} \mathbf{E} \left[\frac{(\mathbf{p}_i - \mathbf{q}_i)^2}{X_i + 1} \right] = \sum_{i \in H} (\mathbf{p}_i - \mathbf{q}_i)^2 (1 - e^{-\lambda_i})/\lambda_i,$$

where $\lambda_i \stackrel{\text{def}}{=} m(\mathbf{p}_i + \mathbf{q}_i)/2$. Notice that $(1 - e^{-x})/x$ is a decreasing function with respect to x . Since $m \leq s/100$ and $(\mathbf{p}_i + \mathbf{q}_i)/2 \leq s$, it follows that $(1 - e^{-\lambda_i})/\lambda_i$ is bounded below by $1 - e^{-0.01}/0.01 \geq 0.99$. This then gives us

$$\mathbf{E} \left[\sum_{i \in H} \frac{(\mathbf{p}_i - \mathbf{q}_i)^2}{X_i + 1} \right] \geq 0.99 \sum_{i \in H} (\mathbf{p}_i - \mathbf{q}_i)^2. \quad (9)$$

On the other hand, the variance of the discrepancy restricted to the elements in H , after flattening, is bounded above by

$$\begin{aligned} \mathbf{Var} \left[\sum_{i \in H} \frac{(\mathbf{p}_i - \mathbf{q}_i)^2}{X_i + 1} \right] &= \mathbf{E} \left[\left(\sum_{i \in H} \frac{(\mathbf{p}_i - \mathbf{q}_i)^2}{X_i + 1} \right)^2 \right] - \mathbf{E}^2 \left[\left(\sum_{i \in H} \frac{(\mathbf{p}_i - \mathbf{q}_i)^2}{X_i + 1} \right)^2 \right] \\ &\leq \left(\sum_{i \in H} (\mathbf{p}_i - \mathbf{q}_i)^2 \right)^2 - (0.99)^2 \left(\sum_{i \in H} (\mathbf{p}_i - \mathbf{q}_i)^2 \right)^2, \end{aligned}$$

where in the last inequality we use the fact that $\frac{1}{1+X_i}$ is at most 1 and (9). This then gives us

$$\sqrt{\mathbf{Var} \left[\sum_{i \in H} \frac{(\mathbf{p}_i - \mathbf{q}_i)^2}{X_i + 1} \right]} \leq 0.15 \sum_{i \in H} (\mathbf{p}_i - \mathbf{q}_i)^2. \quad (10)$$

Combining (9) and (10), we obtain

$$\mathbf{Pr} \left[\|\mathbf{p}_S - \mathbf{q}_S\|_2^2 < \frac{1}{3} \sum_{i \in H} (\mathbf{p}_i - \mathbf{q}_i)^2 \right] \leq 1/4. \quad (11)$$

On the other hand, by Fact 2.15 and Markov's inequality, it holds

$$\Pr \left[\max \left(\|\mathbf{p}\|_2^2, \|\mathbf{q}\|_2^2 \right) > 40/m \right] \leq 1/10. \quad (12)$$

By the union bound, (11), (12) and Corollary 2.17, it follows that the ℓ_2 -closeness tester of Corollary 2.17 succeeds with probability at least $2/3$, if we take $m' = C\sqrt{\frac{1}{m}} \epsilon^{-2}$ many samples, for a sufficiently large constant C . Balancing m and m' (with the restriction that $m \leq s/100$ in mind) then gives us that the overall tester succeeds with probability at least $2/3$ if we draw $\text{Poi}(m)$ many i.i.d. samples with

$$m = \Theta \left(\max \left(\epsilon^{-4/3}, \epsilon^{-2}/\sqrt{s} \right) \right).$$

This concludes the proof of Lemma 2.8. \square

2.4 Applications: Closeness Testing of Multivariate Structured Distributions under Total Variation Distance

The most direct application of our multivariate \mathcal{A}_k -closeness tester is for the problem of testing closeness of multivariate histogram distributions — distributions that are piecewise constant over (the same) *unknown* collection of axis-aligned rectangles — with respect to the total variation distance.

This follows directly from our main theorem, since for any pair of k -histogram distributions \mathbf{p}, \mathbf{q} with respect to the same set of rectangles, we have $d_{\text{TV}}(\mathbf{p}, \mathbf{q}) = \frac{1}{2} \sum_{i=1}^k |\mathbf{p}(R_i) - \mathbf{q}(R_i)| = \frac{1}{2} \|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k}$. Formally, we have the following:

Corollary 2.18. *Let $\{R_i\}_{i=1}^k$ be a set of axis-aligned rectangles in \mathbb{R}^d . Suppose \mathbf{p}, \mathbf{q} are distributions over \mathbb{R}^d that are piecewise constant over each of $\{R_i\}_{i=1}^k$, i.e., $\mathbf{p}(x) = \mathbf{p}(y)$ for any $x, y \in R_i$ and the same for \mathbf{q} . Then there exists a tester which distinguishes between $\mathbf{p} = \mathbf{q}$ and $d_{\text{TV}}(\mathbf{p}, \mathbf{q}) > \epsilon$ with sample complexity $C k^{6/7} \epsilon^{-2\alpha_d/3} \log^d(k) 2^{d/3}$, where C is a sufficiently large universal constant and $\alpha_d = O(d^2 2^{d+1})$.*

We now proceed with our second application. We consider the binary hypothesis class H consisting of all possible k -unions of axis-aligned rectangles within the unit cube $[0, 1]^d$. Given two hypotheses $h_1, h_2 \in H$, we can test whether h_1 is equivalent to h_2 or they are far from each other under the uniform distribution over the unit cube $[0, 1]^d$.

Corollary 2.19. *Let H be the class of all possible k -unions of axis-aligned rectangles within the unit cube $[0, 1]^d$, i.e.,*

$$H = \left\{ h \mid h = \bigcup_{i=1}^k R_i \text{ where } \{R_i\}_{i=1}^k \subset [0, 1]^d \text{ are disjoint axis-aligned rectangles over } [0, 1]^d \right\}.$$

Let h_1, h_2 be two unknown hypotheses from H . Given $\epsilon > 0$ and sample access to $(x, h_i(x))$, where x follows the uniform distribution over $[0, 1]^d$, there exists an efficient algorithm which distinguishes with probability at least $2/3$ between (i) $h_1(x) = h_2(x)$ for all x , and (ii) $\mathbf{E}_{x \sim U} [\mathbb{1}\{h_1(x) \neq h_2(x)\}] > \epsilon$, where U is uniform distribution over $[0, 1]^d$. Moreover, the algorithm has sample complexity $C k^{6/7} \epsilon^{-2\alpha_d/3} \log^d(k) 2^{d/3}$, where C is a sufficiently large constant and $\alpha_d = O(d^2 2^{d+1})$.

Proof. Consider the distributions \mathbf{p}, \mathbf{q} defined as follows. To draw a sample from \mathbf{p} , we take a sample $(x, h_1(x))$ where $x \sim U$. If $h_1(x) = 1$, we return x . Otherwise, we return some arbitrarily

chosen point $s \notin [0, 1]^d$. We define \mathbf{q} similarly based on h_2 . If h_1 and h_2 are identical, it is easy to see that $\mathbf{p} = \mathbf{q}$. If $\mathbf{E}_{x \sim U} [\mathbb{1}\{h_1(x) \neq h_2(x)\}]$, we claim that $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} \geq \epsilon/2$. Suppose that h_1 is the union of the rectangles $\{R_i\}_{i=1}^k$ and h_2 is the union of the rectangles $\{R'_i\}_{i=1}^k$. Then we have that

$$\begin{aligned} & \mathbf{E}_{x \sim U} [\mathbb{1}\{h_1(x) \neq h_2(x)\}] \\ &= \int_{x \in [0,1]^d} U(x) \mathbb{1} \left\{ x \in \bigcup_{i=1}^k R_i \setminus \bigcup_{i=1}^k R'_i \right\} + \int_{x \in [0,1]^d} U(x) \mathbb{1} \left\{ x \in \bigcup_{i=1}^k R'_i \setminus \bigcup_{i=1}^k R_i \right\} \\ &\leq 2 \max \left(\sum_{i=1}^k \int_{x \in [0,1]^d} U(x) \mathbb{1} \left\{ x \in R_i \setminus \bigcup_{i=1}^k R'_i \right\}, \sum_{i=1}^k \int_{x \in [0,1]^d} U(x) \mathbb{1} \left\{ x \in R'_i \setminus \bigcup_{i=1}^k R_i \right\} \right). \end{aligned}$$

Without loss of generality, we assume that the first term is larger. Then we have that

$$\sum_{i=1}^k \int_{x \in [0,1]^d} U(x) \mathbb{1} \left\{ x \in R_i \setminus \bigcup_{i=1}^k R'_i \right\} \geq \frac{1}{2} \epsilon,$$

if $\mathbf{E}_{x \sim U} [\mathbb{1}\{h_1(x) \neq h_2(x)\}] > \epsilon$. On the other hand, we also have

$$\sum_{i=1}^k \int_{x \in [0,1]^d} U(x) \mathbb{1} \left\{ x \in R_i \setminus \bigcup_{i=1}^k R'_i \right\} = \sum_{i=1}^k \mathbf{p}(R_i) - \mathbf{q}(R_i).$$

Thus, this gives $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} \geq \epsilon/2$. Therefore, we can distinguish between the two cases by performing \mathcal{A}_k -closeness testing between \mathbf{p}, \mathbf{q} with accuracy parameter $\epsilon/2$. \square

3 Sample Complexity Lower Bound

In this section, we prove our sample complexity lower bound. Specifically, we show that the task of \mathcal{A}_k -closeness testing gets information-theoretically harder as we go from one dimension to two dimensions. For the one-dimensional case, it was shown in [DKN15a] that the sample complexity of \mathcal{A}_k -closeness testing is $\Theta(\max(k^{4/5}\epsilon^{-6/5}, k^{1/2}\epsilon^{-2}))$. Perhaps surprisingly, for two-dimensional distributions, we prove a sample complexity lower bound of $\Omega(k^{6/7}/\epsilon^{8/7})$ in the sublinear regime, where $\epsilon > k^{-1/8}$. This lower bound clearly dominates the sample complexity of one-dimensional \mathcal{A}_k testing in the same regime.

At a very high level, we build on the lower bound framework of [DKN15a]. In particular, our lower bound proof consists of two steps. First, we argue that, if the domain size is a sufficiently large function of d, k , we can assume without loss of generality that the output of the tester only depends on the *relative order* of samples ranked in each coordinate. This is shown in Section 3.1.

Then, for such “order-based” testers, we present two explicit families of pairs of two-dimensional distributions such that a random pair of distributions from the first family are identical, and a random pair of distributions from the second family are far from each other in \mathcal{A}_k -distance. Moreover, a random pair of distributions from the first family is hard (i.e., requires many samples) to distinguish from a random pair from the second. This step requires a carefully designed gadget consisting of distributions over \mathbb{R}^2 supported on the edges of a square. We present the construction and analyze its key properties in Section 3.2.

Next we appropriately replicate the gadget many times to create the full hard-instance of 2-dimensional \mathcal{A}_k -closeness testing. The description of the hard instance and its detailed analysis can be found in Section 3.3.

Finally, we provide an alternative way to prove a sample complexity lower bound against general \mathcal{A}_k testers, while requiring the domain size to be at most doubly exponential in k . This involves a careful application of randomly chosen monotonic transformations to the x and y coordinates of all points in order to hide extra “non-order based” information that a tester can retrieve from the numerical values of the sample coordinates. This more refined construction and its analysis are presented in Section 3.4.

3.1 Order-Based Testers

Here we define the class of order-based testers and show that we can translate lower bounds against order-based testers to general testers at the cost of increasing the domain size. More formally, we consider algorithms which are restricted to obtain information from what we call the *Order Sampling* process, as opposed to the usual direct sampling. This can be thought of as follows. We first draw i.i.d. samples from the unknown distributions. Then, instead of feeding them directly to the algorithm, we perform an appropriate pre-processing to extract only the information related to the order of the coordinates of the samples, and reveal only the order information to the algorithm.

Definition 3.1 (Order Sampling). Let \mathbf{p}, \mathbf{q} be a pair of distributions in \mathbb{R}^2 . Let $\{(x_i, y_i), \ell_i\}_{i=1}^m$ be m i.i.d. samples, where (x_i, y_i) are sampled from $(1/2)(\mathbf{p} + \mathbf{q})$ and ℓ_i records whether the sample comes from \mathbf{p} or \mathbf{q} . Let $\sigma(x), \sigma(y) \in \mathbb{S}_m$ be the permutation representing the rank of the x -coordinates and y -coordinates accordingly. The *Order Tuple* associated with the m samples is given by $\text{Order}(\{(x_i, y_i, \ell_i)\}_{i=1}^m) = (\sigma(x), \sigma(y), \ell)$. Furthermore, we will use $\mathcal{D}(\mathbf{p}, \mathbf{q}, m)$ to denote the distribution over the tuple $(\sigma(x), \sigma(y), \ell)$ obtained through this process.

As our first structural lemma, we show that if an algorithm is able to perform \mathcal{A}_k -closeness testing with direct sample access on a domain of size $N \times N$, then we can always use it to build another algorithm which performs the test with only the order tuple of the same number of samples — albeit on a smaller domain of size $n \times n$. The proof uses a Ramsey-theoretic argument and generalizes Theorem 13 in [DKN15a].

Lemma 3.2. *For all $n, m, k \in \mathbb{Z}^+$ where $m < n$ and $\epsilon > 0$, there exist $N_1, N_2 \in \mathbb{Z}^+$ such that the following holds: If there exists an algorithm A that for every pair of distributions \mathbf{p}, \mathbf{q} over $[N_1] \times [N_2]$ distinguishes the case $\mathbf{p} = \mathbf{q}$ from the case $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} > \epsilon$ with probability at least $4/5$ while taking m samples from \mathbf{p} and \mathbf{q} , then there exists an algorithm A' that for every pair of distributions \mathbf{p}', \mathbf{q}' over $[n] \times [n]$ distinguishes the case $\mathbf{p}' = \mathbf{q}'$ versus $\|\mathbf{p}' - \mathbf{q}'\|_{\mathcal{A}_k} > \epsilon$ with probability at least $2/3$ given a tuple T from the order sampling process $\mathcal{D}(\mathbf{p}', \mathbf{q}', m)$.*

Proof. Suppose we are given the algorithm A which can perform \mathcal{A}_k -closeness testing over the domain $[N_1] \times [N_2]$ given direct i.i.d. sample access to \mathbf{p}, \mathbf{q} . We show that we can use A to construct another algorithm A' which performs the test with only tuples obtained from the order sampling process over the domain $[n] \times [n]$.

Let $\{(x_i, y_i), \ell_i\}_{i=1}^m$ be the samples drawn by A . We will write $A(\{(x_i, y_i), \ell_i\}_{i=1}^m)$ to denote the probability that A outputs “YES” given these samples. Before we specify our construction, we remark that we can without loss of generality assume that the image of $A(\{(x_i, y_i), \ell_i\}_{i=1}^m)$ has size at most 11. This is because we can always round the probability to the nearest multiples of $1/10$ and lose only $1/10$ in the overall success probability.

Let \mathbf{p}, \mathbf{q} be the unknown distributions supported on $[n] \times [n]$. The key step is to argue the existence of two monotonic transformations $f_x : [n] \mapsto [N_1]$, $f_y : [n] \mapsto [N_2]$, where N_2 is chosen to be a sufficiently large function of n , and N_1 is chosen to be a sufficiently large function of n and N_2 ,

such that if one feeds the samples $\{(f_x(x_i), f_y(y_i)), \ell_i\}$ to A , the output of A becomes a function only of $\text{Order}(\{(x_i, y_i), \ell_i\})$. In other words, we want to find two mappings f_x, f_y such that

$$A(\{(f_x(x_i), f_y(y_i)), \ell_i\}) = A(\{(f_x(x'_i), f_y(y'_i)), \ell'_i\}),$$

as long as $\text{Order}(\{(x_i, y_i), \ell_i\}) = \text{Order}(\{(x'_i, y'_i), \ell'_i\})$. Given such mappings, we can then define $A'(\{(x_i, y_i), \ell_i\}) := A(\{(f_x(x_i), f_y(y_i)), \ell_i\})$. Then, it is easy to see that A' is an order-based tester. Furthermore, since f_x, f_y are both monotonic, the domain transformation will preserve the A_k distance between \mathbf{p} and \mathbf{q} . Hence, A' enjoys the same guarantee and gives the correct answer with probability at least $2/3$.

We next show the existence of such a pair of transformations f_x, f_y . We do so in two steps. First, we show the existence of the transformation f_x which will make the output of algorithm A independent of the actual values of the x -coordinate. This then allows us to construct an algorithm A_x that depends only on the rank information of the x -coordinates, the y -coordinates and the labels. Then we show the existence of f_y , which is defined with respect to A_x , that makes the output of A_x independent of the actual values of the y -coordinates. This then allows us to conclude the existence of the algorithm A' .

For convenience, we will rewrite the tuples $\{(x_i, y_i, \ell_i)\}_{i=1}^m$ as $(\mathcal{X}, \sigma(x), \{y_i\}_{i=1}^m, \{\ell_i\}_{i=1}^m)$, where \mathcal{X} is the set of x -coordinates and $\sigma(x)$ is the permutation which maps $i \in [m]$ to the rank of x_i among $\{x_i\}_{i=1}^m$. For each $\mathcal{X} \in [N_1]^m$, we can define a mapping $g_{\mathcal{X}} : \mathbb{S}_m \times [N_2]^m \times \{0, 1\}^m \mapsto [0, 1]$ induced by the algorithm A as $g_{\mathcal{X}}(\sigma(x), \{y_i\}_{i=1}^m, \{\ell_i\}_{i=1}^m) := A(\mathcal{X}, \sigma(x), \{y_i\}_{i=1}^m, \{\ell_i\}_{i=1}^m)$. Notice that the set of values that $g_{\mathcal{X}}$ has size at most 11, since we assume the acceptance probability of A conditioned on any input can take at most 11 different values. We note that there can be at most $11^{m!N^m2^m}$ many different types of mapping $g_{\mathcal{X}}$.

If we view \mathcal{X} as a hyper-edge of the hypergraph $\binom{[N_1]}{m}$ and the associated mapping $g_{\mathcal{X}}$ as the coloring of the hyper-edge, by Ramsey's theorem, there exists a subset of vertices V of size n such that the coloring of the hyper-edges in the sub-graph $\binom{V}{m}$ are all the same as long as N_1 is sufficiently large compared to N_2 and m . In other words, there exists a subdomain $V \subset [N_1]$ such that if the x coordinates of the samples are all from this subdomain, the acceptance probability of algorithm A becomes a function of only $y_i, \ell_i, \sigma(x)$ and independent of the actual x -coordinates $\mathcal{X} \subseteq V$. We will then choose f_x as the order-preserving mapping from $[n]$ to $[N_1]$, where the image is exactly V .

We next consider the algorithm A_x which first applies the transformation f_x and then runs the testing algorithm A on the resulting samples. From the argument above, we know that algorithm A_x depends only on $\sigma(x), \{y_i\}_{i=1}^m, \{\ell_i\}_{i=1}^m$. Similarly, we can rewrite the tuple as $\sigma(x), \sigma(y), \mathcal{Y}, \{\ell_i\}_{i=1}^m$, where \mathcal{Y} is the set of y -coordinates and $\sigma(y)$ is the permutation which maps i to the rank of y_i . With a similar argument, as long as N_2 is sufficiently large compared to n, m , we can show the existence of an order-preserving mapping f_y such that if we apply the mapping f_y first and then run A_x , the output of A_x becomes only a function of $\sigma_x, \sigma_y, \{\ell_i\}_{i=1}^m$ and independent of the actual set of y coordinates \mathcal{Y} . Notice that $\sigma_x, \sigma_y, \{\ell_i\}_{i=1}^m$ is exactly the order tuple $\text{Order}(\{(x_i, y_i), \ell_i\}_{i=1}^m)$. Hence, such a pair of transformations f_x, f_y are exactly what we need to construct algorithm A' . Setting $A'(\{(x_i, y_i), \ell_i\}) := A(\{(f_x(x_i), f_y(y_i)), \ell_i\})$ then concludes the proof. \square

3.2 Square-Edge Distributions

We now present the building block of our lower bound construction, which consists of distributions supported on the edges of a square. Notice that though the domain is \mathbb{R}^2 , the supports of such distributions are lower-dimensional. We will use \mathbf{t}, \mathbf{r} to represent such distributions and one can refer to Figure 2 for a visual illustration.

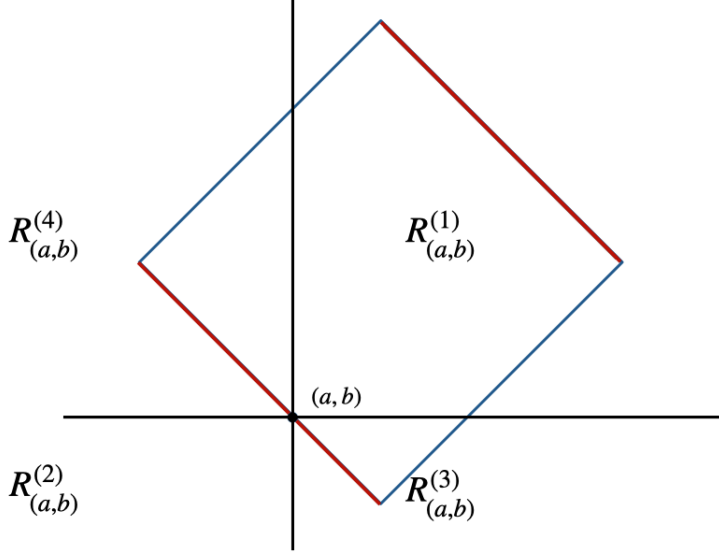


Figure 2: Square Edge Distributions

The red lines represent the distribution \mathbf{t} and the blue lines represent the distribution \mathbf{r} . For any point (a, b) on the edges of the square, it is easy to verify that the four regions $R_{(a,b)}^{(i)}$ in Fact 3.4 have the same probability mass under \mathbf{t} as under \mathbf{r} .

Definition 3.3 (Square-Edge Distributions). Consider a square in \mathbb{R}^2 whose diagonals are parallel to the x -axis and y -axis. We define \mathbf{t} as the uniform distribution supported on the upper-left and lower-right edges and \mathbf{r} as the uniform distribution supported on the remaining two edges.

Let (a, b) be a point lying on the edges of the square. The space can be divided into four regions by drawing one horizontal and one vertical lines across (a, b) . The most important property that we will rely on in our analysis is the following: For any such point (a, b) , any of the resulting four regions have the same mass under \mathbf{t} as under \mathbf{r} .

Fact 3.4. Let \mathbf{t}, \mathbf{r} be the square-edge distributions defined as in Definition 3.3. Consider a point $(a, b) \in \text{supp}(\mathbf{t}) \cup \text{supp}(\mathbf{r})$. Denote the four regions as $R_{a,b}^{(1)} = \{x > a, y > b | (x, y) \in \mathbb{R}^2\}$, $R_{a,b}^{(2)} = \{x < a, y < b | (x, y) \in \mathbb{R}^2\}$, $R_{a,b}^{(3)} = \{x > a, y < b | (x, y) \in \mathbb{R}^2\}$, $R_{a,b}^{(4)} = \{x < a, y > b | (x, y) \in \mathbb{R}^2\}$. Then, it holds $\mathbf{t}(R_{a,b}^{(i)}) = \mathbf{r}(R_{a,b}^{(i)})$ for all i .

Intuitively, the above fact says that if one partitions the space based on one sample (a, b) , the tester cannot distinguish between \mathbf{t} and \mathbf{r} simply based on their mass on any of the regions $R_{a,b}^{(i)}$. As a consequence, to distinguish \mathbf{t} and \mathbf{r} , one needs to take more samples to partition the space into finer pieces (for example, taking two samples and considering the rectangle formed by the two samples).

To formalize this intuition, we will consider the distribution obtained by performing order sampling under a pair of distributions composed of the square-edge distributions. In particular, imagine the following scenario, which can be thought of as a toy example of \mathcal{A}_k closeness testing for $k = 4$. In the YES case, we have $\mathbf{p}_{\text{Yes}} = \mathbf{q}_{\text{Yes}} = (\mathbf{t} + \mathbf{r})/2$. Then we obtain order sampling with m samples drawn from $\mathbf{p}_{\text{Yes}}, \mathbf{q}_{\text{Yes}}$, according to Definition 3.1. The resulting order tuple will then have the distribution $\mathcal{D}((\mathbf{t} + \mathbf{r})/2, (\mathbf{t} + \mathbf{r})/2, m)$ over $\mathbb{S}_m \times \mathbb{S}_m \times \{0, 1\}^m$. In the NO case, with

probability $1/2$, we have $\mathbf{p}_{\text{No}} = \mathbf{t}, \mathbf{q}_{\text{No}} = \mathbf{r}$. Otherwise, we have $\mathbf{p}_{\text{No}} = \mathbf{r}$ and $\mathbf{q}_{\text{No}} = \mathbf{t}$. Then, if we perform order sampling with m samples from $\mathbf{p}_{\text{No}}, \mathbf{q}_{\text{No}}$, we obtain an order tuple following the uniform mixture of $\frac{1}{2}(\mathcal{D}(\mathbf{t}, \mathbf{r}, m) + \mathcal{D}(\mathbf{r}, \mathbf{t}, m))$. Notice that in the YES case, we have $\mathbf{p}_{\text{Yes}} = \mathbf{q}_{\text{Yes}}$; in the NO case, we have $\|\mathbf{p}_{\text{Yes}} - \mathbf{q}_{\text{Yes}}\|_{\mathcal{A}_k} = 1$ deterministically, even for $k = 4$. Yet, we show in the next lemma that the distributions over order-tuples in the two cases are the same when m is no more than 3. This immediately gives us that no order-based algorithm can distinguish between the two cases with fewer than 4 samples.

Lemma 3.5. *We have that $\mathcal{D}((\mathbf{t} + \mathbf{r})/2, (\mathbf{t} + \mathbf{r})/2, m) = (\mathcal{D}(\mathbf{t}, \mathbf{r}, m) + \mathcal{D}(\mathbf{r}, \mathbf{t}, m)) / 2$ for $m = 1, 2, 3$.*

Proof. Let $(\sigma(x), \sigma(y), \ell)$ be an order tuple. We remark that the tuple can be decomposed into two parts: (i) the permutation patterns $\sigma(x), \sigma(y) \in \mathbb{S}_3$, which encodes the “geometric pattern” of the three points sampled and (ii) a bit string $\ell \in \{0, 1\}^3$, which indicates whether the samples come from \mathbf{p} or \mathbf{q} . Now let $(\sigma(x)_{\text{Yes}}, \sigma(y)_{\text{Yes}}, \ell_{\text{Yes}}) \sim \mathcal{D}((\mathbf{t} + \mathbf{r})/2, (\mathbf{t} + \mathbf{r})/2, m)$ and $(\sigma(x)_{\text{No}}, \sigma(y)_{\text{No}}, \ell_{\text{No}}) \sim (\mathcal{D}(\mathbf{t}, \mathbf{r}, m) + \mathcal{D}(\mathbf{r}, \mathbf{t}, m)) / 2$. We begin with the following observations.

1. The marginal distribution over the “geometric pattern” is identical for the two cases, i.e. $\Pr[\sigma(x)_{\text{Yes}} = \pi, \sigma(y)_{\text{Yes}} = \pi'] = \Pr[\sigma(x)_{\text{No}} = \pi, \sigma(y)_{\text{No}} = \pi']$ for all $\pi, \pi' \in \mathbb{S}_m$. This is because the samples, ignoring the labels, in both cases come from the distribution supported uniformly on the four edges of the square.
2. The distribution of ℓ_{Yes} conditioned on any “geometric pattern” will be uniform over all possible bit strings, i.e. $\Pr[\ell_{\text{Yes}} = \beta | \sigma(x)_{\text{Yes}} = \pi, \sigma(y)_{\text{Yes}} = \pi']$ is the same for all $\beta \in \{0, 1\}^m$ and $\pi, \pi' \in \mathbb{S}_m$. This is because $(\sigma(x)_{\text{Yes}}, \sigma(y)_{\text{Yes}}, \ell_{\text{Yes}})$ is obtained by performing order sampling from two identical distributions (both are $(\mathbf{t} + \mathbf{r})/2$).

Hence, it suffices to show that the distribution over the label vector ℓ_{No} conditioned on any geometric patterns $\sigma(x)_{\text{No}}, \sigma(y)_{\text{No}}$ is uniform.

With this observation in mind, the $m = 1$ case is trivial since there is only 1 geometric pattern and it is clear that the label ℓ_{No} is uniform. For $m = 2$, let the coordinate of the first sample be (a, b) , which divides the space into four quadrants. Then, by Fact 3.4, it holds that no matter which of the four quadrants the second sample fall into, the probability that the point comes from \mathbf{t} is the same as it comes from \mathbf{r} . Hence, the uniformity of ℓ_{No} follows.

For $m = 3$, we make some preliminary simplifications. Let $\{(x_i, y_i), \ell_i\}_{i=1}^m$ be three i.i.d. samples drawn. Since they are all identically distributed and independent, the sampling order does not matter. Hence, we can without loss of generality just examine the case $x_1 < x_2 < x_3$ (and accordingly $\sigma(x) = (1, 2, 3)$). Secondly, observe that our construction is invariant under reflections over x - or y -axis, and rotations of angle $\pi/4, \pi/2, 3\pi/4$. After reflection over the x -axis, any three points that have the pattern $\sigma(y) = (1, 2, 3)$ ($x_1 < x_2 < x_3, y_1 < y_2 < y_3$) then becomes $\sigma(y) = (3, 2, 1)$ ($x_1 < x_2 < x_3, y_1 > y_2 > y_3$). After rotations, the pattern $(1, 3, 2)$ yields $(2, 3, 1)$, $(2, 1, 3)$ and $(3, 1, 2)$. Hence, by symmetry, we can simply focus on the argument for $\sigma(x) = (1, 2, 3), \sigma(y) = (1, 2, 3)$ and $\sigma(x) = (1, 2, 3), \sigma(y) = (1, 3, 2)$.

We will begin with $\sigma(x) = (1, 2, 3), \sigma(y) = (1, 2, 3)$ and show that ℓ_{No} is uniform conditioned on that. We claim that this is true even if we further condition on the coordinates of the “middle point”: we will condition on that $x_2 = x, y_2 = y$ for some arbitrarily chosen point (x, y) from the support. It is easy to see that the marginal distribution of ℓ_2 is uniform since it only depends on whether we are sampling from $\mathcal{D}(\mathbf{t}, \mathbf{r}, 3)$ or $\mathcal{D}(\mathbf{r}, \mathbf{t}, 3)$. For the same reason, further conditioning on the value of ℓ_2 then completely determines whether we are sampling from $\mathcal{D}(\mathbf{t}, \mathbf{r})$ or $\mathcal{D}(\mathbf{r}, \mathbf{t}, 3)$. Consequently, $(x_1, y_1, \ell_1), (x_3, y_3, \ell_3)$ are now independent samples from the lower left quadrant $R_{x,y}^{(2)}$

and upper right quadrant $R_{x,y}^{(4)}$ of the point (x, y) respectively. By Fact 3.4, the amount of mass from $\mathbf{p}_{\text{No}}^{(i)}$ and from $\mathbf{q}_{\text{No}}^{(i)}$ in $R_{x,y}^{(2)}$ is the same. Hence, the conditional distribution for ℓ_1 is uniform (and similarly for ℓ_3).

Next, we will show that ℓ_{No} conditioned on $\sigma(x)_{\text{No}} = (1, 2, 3), \sigma(y)_{\text{No}} = (1, 3, 2)$ is also uniform. Notice that it actually suffices for us to show the uniformity of ℓ_{No} conditioned on the more general event $\sigma(x)_{\text{No}} = (1, 2, 3)$ and $\sigma(y)_{\text{No}}$ is either $(1, 3, 2)$ (the case we are analyzing now) or $(1, 2, 3)$ (the case analyzed in the previous paragraph). If this is true, we can then combine it with the fact that ℓ_{No} is uniform conditioned on $\sigma(x)_{\text{No}} = (1, 2, 3), \sigma(y)_{\text{No}} = (1, 2, 3)$ to conclude that ℓ_{No} must be uniform conditioned on $\sigma(x)_{\text{No}} = (1, 2, 3), \sigma(y)_{\text{No}} = (1, 3, 2)$. Notice that this more general event happens if and only if $x_1 < \min(x_2, x_3)$ and $y_1 < \min(y_2, y_3)$. We can then use techniques similar to the analysis of the last case. In particular, we claim that this is true even if we further condition on the coordinates of the first point: $x_1 = x, y_1 = y$ for some arbitrary point (x, y) from the support. The analysis is then almost the same: After we have conditioned on the value of $(x_1, y_1, \ell_1), (x_2, y_2, \ell_2), (x_3, y_3, \ell_3)$ now both become independent samples from the upper right quadrant $R_{x,y}^{(4)}$. Applying Fact 3.4 then allows us to conclude the uniformity of ℓ_2, ℓ_3 after the conditioning. This finishes the argument that ℓ_{No} conditioned on any geometric patterns $\sigma(x)_{\text{No}}, \sigma(y)_{\text{No}}$ and concludes the proof. \square

3.3 \mathcal{A}_k Closeness Lower Bound Construction

We will now use \mathbf{t}, \mathbf{r} as building blocks to construct the full hard instance of 2-dimensional \mathcal{A}_k closeness testing and establish the desired sample complexity lower bound $\Omega(\min(k^{6/7}\epsilon^{-8/7}, k))$.

We will readily apply the ‘‘Poissonization trick’’, which is a standard technique in proving lower bounds for distribution testing problems. In particular, instead of drawing a fixed number of m samples, we make the testers draw $\text{Poi}(m)$ many samples. It is easy to translate any lower bound in the Poisson sampling model to the standard sampling model where the testers draw a fixed number of samples, since with probability at least 99% the testers will receive at least $\Omega(m)$ many samples.

Furthermore, we will relax \mathbf{p}, \mathbf{q} to be non-negative measures whose total mass is $\Theta(1)$ rather than equal to 1. Clearly, taking samples from a non-negative measure μ is no longer a sensible concept. Instead, we can take $\text{Poi}(m \|\mu\|_1)$ samples from the normalized distribution $\mu / \|\mu\|_1$. We will slightly abuse the definition of sampling to describe the above the process as ‘‘taking $\text{Poi}(m)$ samples from μ ’’.

Lastly, since we are only proving a sample complexity lower bound that is sublinear with respect to k , we can safely assume $m < k/2$ throughout the section.

Now we are ready to describe the hard instance. We will first partition the domain into r^2 squares with equal size, for some $r = \Theta(k)$ that will be specified later. Most of the squares will be left blank: \mathbf{p}, \mathbf{q} will have all their probability mass supported within the squares along one diagonal of the square grids. For each square on the diagonal, we will make it a ‘‘heavy’’ square with probability m/k (this is a well-defined probability since $m < k$) and a ‘‘light’’ square otherwise, whose purpose will become clear later.

Now consider the following random process for generating a pair of measures \mathbf{p}, \mathbf{q} . Let X be a random variable that takes 0 or 1 each with probability 1/2. If $X = 0$, we will randomly generate a pair of measures $\mathbf{p} = \mathbf{q}$, which belongs to the YES instance. If $X = 1$, we randomly generate a pair of measures satisfying $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} > \Omega(\epsilon)$, which belong to the NO instance.

When $X = 1$, \mathbf{p}, \mathbf{q} restricted to one square (after normalization) will be both $(\mathbf{t} + \mathbf{r})/2$, which is the uniform distribution supported on the edges of a diagonal square. Moreover, the mass of \mathbf{p} will be $1/m$ if the square is ‘‘heavy’’ and ϵ/k if the square is ‘‘light’’ (and the same for \mathbf{q} as well).

When $X = 0$, the mass of \mathbf{p}, \mathbf{q} restricted to a square will be the same as the case $X = 1$. Yet, the conditional distributions within a square for \mathbf{p}, \mathbf{q} will be different.

- For a “heavy” square, the conditional distributions of \mathbf{p}, \mathbf{q} restricted to the square are still both $(\mathbf{t} + \mathbf{r})/2$. Intuitively, samples produced by the “heavy” squares behave the same in the NO instance as in the YES instance, serving as noise to “confuse” the algorithm.
- For a “light” square, the conditional distributions of \mathbf{p}, \mathbf{q} restricted to the square are respectively \mathbf{t}, \mathbf{r} with probability $1/2$ and \mathbf{r}, \mathbf{t} otherwise. These squares contribute to the \mathcal{A}_k discrepancy between \mathbf{p}, \mathbf{q} but remain hard to distinguish from the YES case.

We first argue that the measures \mathbf{p}, \mathbf{q} constructed from the random process described above qualify for basic properties of \mathcal{A}_k closeness testing.

Lemma 3.6. *Suppose $m < k/2$. It holds that \mathbf{p}, \mathbf{q} are positive measures with mass $\Theta(1)$ with probability 99%. Moreover, if $X = 1$, we have $\mathbf{p} = \mathbf{q}$. If $X = 0$, we have $\|\mathbf{p}/\|\mathbf{p}\|_1 - \mathbf{q}/\|\mathbf{q}\|_1\|_{\mathcal{A}_k} > \Omega(\epsilon)$ with probability 99%.*

Proof. We first verify that \mathbf{p}, \mathbf{q} are both measures with mass $\Theta(1)$ with probability 99%. By Chebyshev’s inequality, we have that the number of heavy squares is $r \frac{m}{k} \pm \Theta(1) \sqrt{r \frac{m}{k}} = \Theta(1) \frac{rm}{k}$ with probability 99%. Conditioned on that, the contribution of the heavy squares to mass is $\Theta(1) \frac{rm}{k} \frac{1}{m} = \Theta(1)$ given that $r = \Theta(k)$. The contribution of the light squares is at most $r \frac{\epsilon}{k} = O(\epsilon)$. Hence, we have the total mass will be $\Theta(1)$.

If $X = 1$, it is easy to see that $\mathbf{p} = \mathbf{q}$. If $X = 0$, for each light square R , recall that \mathbf{p}, \mathbf{q} restricted to R are exactly the square edge distributions after normalization. By the definition of the square edge distribution, there exists 4 sub-squares R_1, R_2, R_3, R_4 such that for each R_i , exactly one of $\mathbf{p}(R_i), \mathbf{q}(R_i)$ is 0 and the other one is $\epsilon/(2k)$. We have seen that $\|\mathbf{p}\|_1, \|\mathbf{q}\|_1$ are both $\Theta(1)$. Hence, we have $\sum_{i=1}^4 |\mathbf{p}(R_i)/\|\mathbf{p}\|_1 - \mathbf{q}(R_i)/\|\mathbf{q}\|_1| \geq \Omega(\epsilon/k)$. Moreover, with probability 99%, the number of light squares is $r (1 - \frac{m}{k}) \pm \Theta(1) \sqrt{r (1 - \frac{m}{k})} = \Theta(1) r$ since $m/k < 1/2$. Conditioned on this, if we choose $r = c k$ for a sufficiently small constant c , we ensure that there are $r' = \Theta(1) c k$ light squares. Notice that if c is chosen appropriately, we can ensure $\Omega(k) < r' < k/4$. Therefore, there exists $k' = 4 r' < k$ rectangles such that $\sum_{i=1}^{k'} |\mathbf{p}(R_i)/\|\mathbf{p}\|_1 - \mathbf{q}(R_i)/\|\mathbf{q}\|_1| = \epsilon/k r' = \Omega(\epsilon)$. \square

Let T be the tuple obtained from the order sampling process $\mathcal{D}(\mathbf{p}, \mathbf{q}, m')$, where \mathbf{p}, \mathbf{q} are the pair of random measures described above and $m' \sim \text{Poi}(m)$. We will bound above the mutual information $I(X : T)$, implying that T reveals little information of the random variable X . The implication argument is standard, see, e.g., the proof of Theorem 16 from [DKN15a]. In particular, we try to bound the information about X obtained from samples falling in each of the squares. In [DKN15a], we have that squares with fewer than two samples are uninformative. By Lemma 3.5, we can further ignore the squares with three samples, therefore allowing us to obtain a stronger lower bound.

Our key technical lemma is the following:

Lemma 3.7. *We have that $I(X : T) = O(m^7 \epsilon^8 / k^6)$.*

Proof. Let $Y = \{(x_i, y_i), \ell_i\}_{i=1}^{m'}$, where $m' \sim \text{Poi}(m)$, be the sample points drawn. Namely, $T = \text{Order}(Y)$. Denote by Y_i the set of points in the i -th square along the diagonal and define the tuple $T_i = \text{Order}(Y_i)$. One can easily reconstruct T from $\{T_1, \dots, T_r\}$: Given $i < j$, all points from the i -th square will be ranked after points from the j -th square in both x and y coordinates in T . This hence gives us that $I(X : T) \leq \sum_{i=1}^r I(X : T_i)$. Next, we will bound $I(X : T_i)$ by $O(m^7 \epsilon^8 / k^7)$. Our

lemma easily follows from that since we also have $r = \Theta(k)$. We first bound the mutual information as a summation over all possible order tuples grouped by the size of the order tuple (recall that for an order tuple t , the size of the order tuple, denoted as $|t|$, is simply the number of samples from which the order tuple is derived). We have that

$$I(X : T_i) \leq O(1) \sum_{\lambda=0}^{\infty} \sum_{\text{order tuple } t:|t|=\lambda} \frac{(\Pr [T_i = t|X = 0] - \Pr [T_i = t|X = 1])^2}{\Pr [T_i = t]}.$$

We will use the indicator variable H_i to denote whether the i -th square is chosen to be a ‘‘heavy’’ square. Notice that $\Pr[H_i = 0] = 1 - \frac{m}{k} = O(1)$ and H_i is independent of X . Furthermore, if the i -th square is chosen to be a heavy square, the distribution of T_i conditioned on $X = 0$ and $X = 1$ is exactly the same. This gives us that

$$I(X : T_i) \leq O(1) \cdot \sum_{\lambda=0}^{\infty} \sum_{t:|t|=\lambda} \frac{(\Pr [T_i = t|H_i = 0, X = 0] - \Pr [T_i = t|H_i = 0, X = 1])^2}{\Pr [T_i = t]}.$$

Next, we note that $\Pr [T_i = t|X = 0, H_i = 0]$ for $|t| = \lambda$ is given by the distribution $\frac{1}{2}\mathcal{D}(\mathbf{t}, \mathbf{r}, \lambda) + \frac{1}{2}\mathcal{D}(\mathbf{r}, \mathbf{t}, \lambda)$. On the other hand, $\Pr [T_i = t|X = 1, H_i = 0]$ for $|t| = \lambda$ is given by the distribution $\mathcal{D}((\mathbf{t} + \mathbf{r})/2, (\mathbf{t} + \mathbf{r})/2, \lambda)$. Hence, by Lemma 3.5, it holds

$$\Pr [T = t|H_i = 0, X = 0] = \Pr [T = t|H_i = 0, X = 1]$$

for any t satisfying $|t| \leq 3$. This allows us to discard the summation over any t with $|t| \leq 3$. Hence, the expression can be further upper bounded by

$$\begin{aligned} & O(1) \sum_{\lambda=4}^{\infty} \sum_{t:|t|=\lambda} \frac{(\Pr [T_i = t|H_i = 0, X = 0] - \Pr [T_i = t|H_i = 0, X = 1])^2}{\Pr [T_i = t]} \\ & \leq O(1) \sum_{\lambda=4}^{\infty} \sum_{t:|t|=\lambda} \frac{(\Pr [T_i = t|H_i = 0])^2}{\Pr [T_i = t, H_i = 1]} \\ & \leq O(1) \sum_{\lambda=4}^{\infty} \max_{t:|t|=\lambda} \frac{\Pr [T_i = t|H_i = 0]}{\Pr [T_i = t, H_i = 1]} \sum_{t:|t|=\lambda} \Pr [T_i = t|H_i = 0] \\ & = O(1) \sum_{\lambda=4}^{\infty} \max_{t:|t|=\lambda} \frac{\Pr [T_i = t|H_i = 0]}{\Pr [T_i = t, H_i = 1]} \Pr [|T_i| = \lambda|H_i = 0], \end{aligned} \quad (13)$$

where in the second line above we upper bound the difference in the numerator by their sum and upper bound the denominator by $\Pr [T_i = t, H_i = 1]$, in the third line above we use that $\sum_i a_i b_i \leq (\max_i a_i) (\sum_i b_i)$ when $a_i, b_i \geq 0$ and in the final equality we note that the summation over the probability of $T_i = t$ for each $|t| = \lambda$ is exactly that of $|T_i| = \lambda$. Next we claim that

$$\max_{t:|t|=\lambda} \frac{\Pr [T_i = t|H_i = 0]}{\Pr [T_i = t, H_i = 1]} \leq O(1) 2^\lambda \frac{\Pr [|T_i| = \lambda|H_i = 0]}{\Pr [|T_i| = \lambda, H_i = 1]}. \quad (14)$$

To show this, we first remark that

$$\max_{t:|t|=\lambda} \frac{\Pr [T_i = t|H_i = 0]}{\Pr [T_i = t|H_i = 1]} = \max_{t:|t|=\lambda} \frac{\Pr [T_i = t|H_i = 0, |T_i| = \lambda]}{\Pr [T_i = t|H_i = 1, |T_i| = \lambda]} \frac{\Pr [|T_i| = \lambda|H_i = 0]}{\Pr [|T_i| = \lambda|H_i = 1]}. \quad (15)$$

Then recall that T_i can be decomposed into a binary vector ℓ_i representing the labels and a permutation tuple $\sigma_i \in S_\lambda \times S_\lambda$ representing the rank information of x and y coordinates. We note that $\sigma_i|H_i = 0, |T_i| = \lambda$ has the same distribution as $\sigma_i|H_i = 1, |T_i| = \lambda$. Then, conditioned on σ_i , the distribution of ℓ_i is uniform when $H_i = 1$. This then gives

$$\max_{t:|t|=\lambda} \frac{\Pr[T_i = t|H_i = 0, |T_i| = \lambda]}{\Pr[T_i = t|H_i = 1, |T_i| = \lambda]} \leq O(2^\lambda).$$

Combining this with Equation (15) and multiplying both sides by $\frac{1}{\Pr[H_i=1]}$ then gives (14). Substituting Equation (14) into Equation (13) then gives us

$$I(X : T_i) \leq O(1) \cdot \sum_{\lambda=4}^{\infty} 2^\lambda \frac{(\Pr[|T_i| = \lambda|H_i = 0])^2}{\Pr[|T_i| = \lambda, H_i = 1]}.$$

Finally, notice that $\Pr[|T_i| = \lambda|H_i = 1] = \text{Poi}(1, \lambda) = \Theta(1)/\lambda!$, $\Pr[H_i = 1] = m/k$, and $\Pr[|T| = \lambda|H_i = 0] = \text{Poi}(\epsilon m/k, \lambda) \leq (\epsilon m/k)^\lambda/\lambda!$. This further gives

$$\begin{aligned} I(X : T_i) &\leq O(1) \sum_{\lambda=4}^{\infty} \frac{2^\lambda}{\lambda!} \frac{k}{m} \left(\frac{\epsilon m}{k}\right)^{2\lambda} \\ &= O(1) \frac{k}{m} \sum_{\lambda=4}^{\infty} \left(\sqrt{2} \frac{\epsilon m}{k}\right)^{2\lambda} \leq O(1) \left(\frac{m}{k}\right)^7 \epsilon^8. \end{aligned}$$

This concludes the proof of Lemma 3.7. \square

We are now ready to conclude the proof of our main lower bound result.

Proof of Lower Bound in Theorem 1.2. By Lemma 3.6, given that $m < k/2$, it holds that both \mathbf{p}, \mathbf{q} are measures of mass $\Theta(1)$ with probability at least 99% and if $X = 0$, it holds $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} > \Omega(\epsilon)$ with probability at least 99%. By Lemma 3.7, we have that the mutual information between the random bit X and the ordering tuple $T \sim \mathcal{D}(\mathbf{p}, \mathbf{q}, m')$, for $m' \sim \text{Poi}(m)$, is at most $O(m^7 \epsilon^8 / k^6)$. This means that no algorithm, given T as input, can reliably predict the value of X with probability more than $2/3$ unless $m > \Omega(1) \min(k^{6/7}/\epsilon^{8/7}, k)$. By Lemma 3.6, it holds that $\mathbf{p}/\|\mathbf{p}\|_1, \mathbf{q}/\|\mathbf{q}\|_1$ are a pair of identical distributions if $X = 0$ and a pair of distributions that are $\Omega(\epsilon)$ far in \mathcal{A}_k distance with probability at least 99% if $X = 1$. Furthermore, with probability 99%, T is an order-tuple of at most $O(m)$ many samples. Therefore, we conclude that the sample complexity of \mathcal{A}_k testing is at least $\Omega(1) \min(k^{6/7}/\epsilon^{8/7}, k)$.

Even though the distributions \mathbf{p}, \mathbf{q} used in the construction are continuous, we next show that they can be easily “rounded” to discrete distributions that remain hard for the testing algorithm. In particular, we can construct a grid \mathcal{G} which splits the domain into $\Theta(m^6)$ squares such that the mass of any square $R \in \mathcal{G}$ under $(1/2)(\mathbf{p} + \mathbf{q})$ is bounded by m^3 . Then, we consider the discrete distributions \mathbf{p}', \mathbf{q}' which round the points falling in the square $R \in \mathcal{G}$ to its top-left vertex. It is easy to see that if $\mathbf{p} = \mathbf{q}$, then $\mathbf{p}' = \mathbf{q}'$. Moreover, for an arbitrary rectangle $R \subset \mathbb{R}^2$, we have $\mathbf{p}(R) - \mathbf{q}(R) = \mathbf{p}'(R) - \mathbf{q}'(R) \pm \Theta(\frac{1}{m^3})$. Hence, the effect of rounding to the \mathcal{A}_k distance between \mathbf{p}, \mathbf{q} is at most $\Theta(k/m^3)$, which can be safely ignored when $m < k$. On the other hand, $\mathcal{D}(\mathbf{p}, \mathbf{q}, m')$ is nearly the same as $\mathcal{D}(\mathbf{p}', \mathbf{q}', m')$, since the distributions over the order tuples are the same as long as no two points fall in the same square (which happens with probability at most $O(1/m)$). Hence, the cases $\mathbf{p}' = \mathbf{q}'$ and $\|\mathbf{p}' - \mathbf{q}'\|_{\mathcal{A}_k} > \epsilon$ are also hard to distinguish given tuples from the order sampling process unless $m > \Omega(1) \min(k^{6/7}/\epsilon^{8/7}, k)$. Finally, by Lemma 3.2, we can translate any lower bound under order sampling back to the usual direct sampling. \square

3.4 Domain Size Optimization

The lower bound of Theorem 1.2 holds only when the domain size N is substantially larger than the other parameters. In particular, the statement does not quantitatively characterize the sample complexity as a function of the domain size. The bottleneck of the analysis lies in Lemma 3.2, which offers an inefficient (in terms of the size of the domain after the transformation) way of transforming the domain to “hide” the extra information that an algorithm can extract from the samples in addition to their relative order. In this section, we provide a more efficient and constructive way to disguise the information in the values of each samples’ coordinates and build on it to provide a tighter lower bound in terms of the domain size.

The main result of this section is the following:

Theorem 3.8 (Stronger Lower Bound for Discrete Distributions). *Fix an integer $V > 0$. Let \mathbf{p} and \mathbf{q} be distributions on $[V] \times [V]$ and let $\epsilon > 0$ be less than a sufficiently small constant. Any tester that distinguishes between $\mathbf{p} = \mathbf{q}$ and $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} \geq \epsilon$ for some $k \leq V$ with probability at least $2/3$ must use at least m many samples for some m with*

$$m \geq \Omega(1) \cdot \min \left(k^{2/3} \epsilon^{-4/3} \cdot \left(\frac{\log \log V}{\log \log \log V} \right)^{1/3}, k^{6/7} \epsilon^{-8/7}, k \right).$$

Before presenting the transformation formally, we provide some high level intuition. Recall that in the lower bound construction from Section 3.3 the domain is partitioned into $r \times r$ many squares where $r = \Theta(k)$ and the distributions are supported on squares lying on the diagonal. The argument then proceeds to bound the order information of samples coming from each of the squares. Now suppose that the algorithm is allowed to look at the absolute coordinates of the samples. If only 1 or 2 points fall in some square, the only extra information we need to hide is its absolute position and the distance between the points. To do so, we can generalize the techniques developed in [DKN17] to randomly scale and shift the square in both the x and y -axis.

For 2-dimensional \mathcal{A}_k closeness testing, if the algorithm takes $\Theta(k^{6/7} \epsilon^{-8/7})$ many samples, since there are $\Theta(k)$ many squares in total, 3 or more samples could fall in the same square. Then the algorithm also gets to see the ratio of distances between different pairs of points, which remains invariant even if the coordinates of the points are scaled uniformly within the square. To handle this, we will instead apply an uneven scaling on different parts of the square. In particular, we map points with x -coordinate a to $\exp(a \lambda)$ with some randomly chosen λ (and the same for the y -coordinate), which then makes the ratio of distances also noisy.

To formalize this idea, we first define a distribution over monotonic mappings, which we will then use to transform the points.

Definition 3.9 (Distribution over monotonic mappings). Let $W > 0$. We define $\mathcal{M}(W)$ as a distribution over monotonic mappings of the form $f : [0, 1] \mapsto \mathbb{R}_+$. To sample a mapping from \mathcal{M} , we first sample three parameters $\lambda_1, \lambda_2, \lambda_3$ which are uniform variables over the intervals $[\log \log W, 2 \log \log W]$, $[0, \log^3 W]$, $[0, \exp(2 \log^3 W)]$ respectively. Then, the mapping $f \sim \mathcal{M}(W)$ is given by $f(x) = \exp(x \exp(\lambda_1)) \exp(\lambda_2) + \lambda_3$.

Let $a < b < c$ be three points lying on $[0, 1]$. Here we show that, as long as a, b, c are sufficiently separated, transforming the points by some random mapping f from $\mathcal{M}(W)$ helps obfuscate the information a tester can retrieve from them. In particular, we argue the distribution of $(f(a), f(b), f(c))$ (where the randomness is over f) is close to some fixed distribution D for any choice of well-separated points a, b, c .

Lemma 3.10. *Let $f \sim \mathcal{M}(W)$ such that $f(x) = \exp(\exp(\lambda_1) x) \exp(\lambda_2) + \lambda_3$. Then, there exists some fixed distributions D over the domain \mathbb{R}_+^3 such that for any three points $a < b < c$ from $[0, 1]$ satisfying*

$$\min(c - b, b - a) > 1/\log \log W, \quad (16)$$

we always have

$$d_{\text{TV}}((f(a), f(b), f(c)), D) \leq O\left(\frac{\log \log \log W}{\log \log W}\right).$$

Proof. Define $A \stackrel{\text{def}}{=} \frac{f(c)-f(a)}{f(b)-f(a)}$, $B \stackrel{\text{def}}{=} f(b) - f(a)$, $C \stackrel{\text{def}}{=} f(a)$. First, we note that it suffices to show $(\log \log A, \log B, C)$ is close in total variation distance to some distribution D' for an arbitrary choice of a, b, c satisfying the condition in Equation (16) since is a bijection between $(f(a), f(b), f(c))$ and $(\log \log A, \log B, C)$.

In particular, let U_1, U_2, U_3 are uniform distributions over the intervals $[\log \log W, 2 \log \log W]$, $[0, \log^3 W]$ and $[0, \exp(2 \log^3 W)]$ respectively. We argue $(\log \log A, \log B, C)$ is close to the distribution $U_1 \times U_2 \times U_3$. The proof strategy is the following. We first bound the total variation distance between $\log \log A$ and U_1 . Then, conditioned on $\log \log A$ and U_1 , we show $\log B$ is close to U_2 . Finally, conditioned on everything other variables, we show C is close to U_3 .

Suppose $f(x) = \exp(\exp(\lambda_1) x + \lambda_2) + \lambda_3$. We have that $\log \log A = g_{a,b,c}(\lambda_1)$, where

$$g_{a,b,c}(x) = \log \log \left(\frac{\exp(c \exp(x)) - \exp(a \exp(x))}{\exp(b \exp(x)) - \exp(a \exp(x))} \right).$$

It is easy to verify that $g_{a,b,c}$ is monotonically increasing as a function of x for any $a < b < c$. Since λ_1 is uniform over $[\log \log W, 2 \log \log W]$, the support of $\log \log A$ will be $[g_{a,b,c}(\log \log W), g_{a,b,c}(2 \log \log W)]$. By the change of variable rule of probability density functions, we have

$$\Pr[\log \log A = x] = \begin{cases} \Pr[\lambda_1 = g_{a,b,c}^{-1}(x)] \frac{1}{g'_{a,b,c}(g_{a,b,c}^{-1}(x))}, & \text{if } x \in [g_{a,b,c}(\log \log W), g_{a,b,c}(2 \log \log W)], \\ 0 & \text{otherwise.} \end{cases}$$

Before we bound the total variation distance between $\log \log A$ and U_1 , we discuss some useful properties of $g_{a,b,c}$.

Claim 3.11. *Given $a < b < c \in [0, 1]$ are well separated (satisfying Equation (16)), it holds that (i) $g_{a,b,c}(x) \leq x$ (ii) $|g_{a,b,c}(x) - x| \leq \log \log \log W + O(1)$ (iii) $|g'_{a,b,c}(x) - 1| \leq O\left(\frac{1}{\log \log W}\right)$ for $x \in [\log \log W, 2 \log \log W]$.*

Proof. For the proof of this claim, we will temporarily drop the subscript of $g_{a,b,c}$ and write only g . For property (i), we have

$$g(x) \leq \log \log \left(\frac{\exp(c \exp(x))}{\exp(b \exp(x))} \right) = \log(c - b) + x \leq x,$$

where the last inequality is true since $c - b \in [0, 1]$, which follows from $b < c$ and $c, b \in [0, 1]$.

For properties (ii) and (iii), our strategy is to show that $g(x)$ is approximately just $x + \log(c - b)$ for sufficiently large W . To do so, we consider the function $h(x) := \exp(g(x)) =$

$\log\left(\frac{\exp(c\exp(x))-\exp(a\exp(x))}{\exp(b\exp(x))-\exp(a\exp(x))}\right)$. Our goal now is to show $h(x)$ is approximately $(c-b)\exp(x)$. Denote $L_\theta(x) = \log(1 - \exp(-\theta \exp(x)))$. We then have

$$\begin{aligned} h(x) &= \log\left(\exp(c\exp(x)) - \exp(a\exp(x))\right) - \log\left(\exp(b\exp(x)) - \exp(a\exp(x))\right) \\ &= \log\left(\exp(c\exp(x))(1 - \exp((a-c)\exp(x)))\right) - \log\left(\exp(b\exp(x))(1 - \exp((a-b)\exp(x)))\right) \\ &= c\exp(x) + \log\left(1 - \exp((a-c)\exp(x))\right) - b\exp(x) - \log\left(1 - \exp((a-b)\exp(x))\right) \\ &= (c-b)\exp(x) + L_{c-a}(x) - L_{b-a}(x). \end{aligned} \tag{17}$$

For $\theta \in [1/\log\log W, 1]$ and $x \in [\log\log W, 2\log\log W]$, we claim $L_\theta(x)$ becomes almost the 0 function in terms of its function values and its derivative when W grows. Using the inequality $-\log(1 - \exp(-z)) \leq 1/z$ for $z > 0$, we have

$$|L_\theta(x)| \leq \frac{1}{\theta \exp(x)} < \frac{1}{\log\log W}. \tag{18}$$

Furthermore, the derivative of L_θ can be bounded by

$$L'_\theta(x) = \frac{\exp(x)\theta}{\exp(\exp(x)\theta) + 1} \leq \frac{\log^2 W}{W^{1/\log\log W}} < \frac{1}{\log\log W}. \tag{19}$$

Combining Equations (17) and (18), we then have

$$h(x) = (c-b)\exp(x) \pm O\left(\frac{1}{\log\log W}\right).$$

Notice that $h(x)$ is at most $\exp(x) + O(1/\log\log W)$. Then, we have

$$g(x) \leq \log\left(\exp(x) + O\left(\frac{1}{\log\log W}\right)\right) = x + \log\left(1 + O(1)\frac{1}{\exp(x)\log\log W}\right) \leq x + O(1).$$

On the other hand, since $(c-b)$ is at least $\frac{1}{\log\log W}$, $h(x)$ is at least $\frac{\exp(x)}{\log\log W} - O(1/\log\log W)$. Then, we have

$$\begin{aligned} g(x) &\geq \log\left(\frac{\exp(x)}{\log\log W} - O\left(\frac{1}{\log\log W}\right)\right) \\ &= \log(\exp(x) - O(1)) - \log\log W \\ &= x + \log(1 - O(\exp(-x))) - \log\log\log W \\ &\geq x + \log((1 - O(\exp(-\log\log W)))) - \log\log\log W \\ &\geq x - O(1) - \log\log\log W, \end{aligned}$$

where the last inequality holds since for sufficiently large W , we have $O(\exp(-\log\log W)) \leq 1/2$. This then gives us property (ii).

Using Equations (19) and (17), we then have

$$h'(x) = (c-b)\exp(x) + L'_{c-a}(x) + L'_{b-a}(x) = (c-b)\exp(x) \pm O\left(\frac{1}{\log\log W}\right).$$

Hence, we can bound the derivative of $g(x)$ as

$$\begin{aligned}
g'(x) &= \frac{1}{h(x)} h'(x) = \frac{(c-b) \exp(x) \pm O\left(\frac{1}{\log \log W}\right)}{(c-b) \exp(x) \pm O\left(\frac{1}{\log \log W}\right)} \\
&= \frac{1 \pm O\left(\frac{1}{\log \log W} \frac{1}{(c-b) \exp(x)}\right)}{1 \pm O\left(\frac{1}{\log \log W} \frac{1}{(c-b) \exp(x)}\right)} = \frac{1 \pm O\left(\frac{1}{\log \log W}\right)}{1 \pm O\left(\frac{1}{\log \log W}\right)} \\
&= 1 \pm O\left(\frac{1}{\log \log W}\right),
\end{aligned}$$

where in the second last equality, we use the fact that $(c-b) \exp(x)$ is at least $\log W / \log \log W$ and hence lower bounded by a constant for sufficiently large W . This concludes the proof of Claim 3.11. \square

To bound the total variation distance between $\log \log A$ and U_1 , we will introduce $U_{a,b,c}$, which denotes the uniform distribution over $[g_{a,b,c}(\log \log W), g_{a,b,c}(2 \log \log W)]$. Then, by the triangle inequality, we have that $d_{TV}(\log \log A, U_1) \leq d_{TV}(\log \log A, U_{a,b,c}) + d_{TV}(U_{a,b,c}, U_1)$. Notice that the second term is just the total variation distance between two uniform variables - one over the interval $[\log \log W, 2 \log \log W]$ and the other over $[g_{a,b,c}(\log \log W), g_{a,b,c}(2 \log \log W)]$. By Claim 3.11, it holds $|g_{a,b,c}(x) - x| = O(\log \log \log W)$ and $g_{a,b,c}(x) \leq x$. We thus have

$$g_{a,b,c}(\log \log W) \leq \log \log W \leq g_{a,b,c}(2 \log \log W) \leq 2 \log \log W.$$

Hence, the total variation distance between U_1 and $U_{a,b,c}$ is exactly

$$\begin{aligned}
&\frac{1}{2} \left(\frac{\log \log W - g_{a,b,c}(\log \log W)}{g_{a,b,c}(2 \log \log W) - g_{a,b,c}(\log \log W)} + \frac{2 \log \log W - g_{a,b,c}(2 \log \log W)}{\log \log W} \right. \\
&\quad \left. + (g_{a,b,c}(2 \log \log W) - \log \log W) \left| \frac{1}{\log \log W} - \frac{1}{g_{a,b,c}(2 \log \log W) - g_{a,b,c}(\log \log W)} \right| \right),
\end{aligned}$$

where the first two terms capture the difference between U_1 and $U_{a,b,c}$ on the domain such that exactly one of U_1 and $U_{a,b,c}$ is supported on, and the last term captures the difference on the domain they are commonly supported on. For the first two terms, the numerators are of size $O(\log \log \log W)$ and the denominators are at least $\log \log W - O(\log \log \log W)$ since $|g_{a,b,c}(x) - x| = O(\log \log \log W)$. Therefore, both of them are of order $O(\log \log \log W / \log \log W)$. For the last term, we have $g_{a,b,c}(2 \log \log W) - \log \log W \leq \log \log W + O(\log \log \log W)$ and

$$\left| \frac{1}{\log \log W} - \frac{1}{g_{a,b,c}(2 \log \log W) - g_{a,b,c}(\log \log W)} \right| \leq O(\log \log \log W).$$

Hence, in total, we have $d_{TV}(U_1, U_{a,b,c}) \leq O(\log \log \log W / \log \log W)$.

For the term $d_{TV}(\log \log A, U_{a,b,c})$, one can see that the two variables have the same support. We will first show the PDF of $\log \log A$ and $U_{a,b,c}$ are point-wise close. In particular, for $x \in$

$[g(\log \log W), g(2 \log \log W)]$, we have

$$\begin{aligned}
& |\Pr[\log \log A = x] - \Pr[U_{a,b,c} = x]| \\
&= \left| \Pr[\lambda_1 = g_{a,b,c}^{-1}(x)] \frac{1}{g'_{a,b,c}(g_{a,b,c}^{-1}(x))} - \frac{1}{g_{a,b,c}(2 \log \log W) - g_{a,b,c}(\log \log W)} \right| \\
&= \left| \frac{1}{\log \log W} \frac{1}{1 \pm O\left(\frac{1}{\log \log W}\right)} - \frac{1}{\log \log W \pm O(\log \log \log W)} \right| \\
&= O\left(\frac{\log \log \log W}{(\log \log W)^2}\right),
\end{aligned}$$

where in the second equality we use the fact λ_1 is a uniform variable over an interval of length $\log \log W$ and that $g'(a, b, c) = 1 \pm O(1/\log \log W)$ by Claim 3.11. Then, since the interval where $\log \log A, U_{a,b,c}$ are supported on is of length at most $O(\log \log W)$. We then have $d_{TV}(\log \log A, U_{a,b,c}) \leq O(\log \log \log W / \log \log W)$. Hence, overall, we then have

$$d_{TV}(\log \log A, U_1) \leq O(\log \log \log W / \log \log W).$$

Next, we show $\log B$ conditioned on $\log \log A$ is close to U_2 . We can simplify the expression of $\log B$ and arrive at

$$\log B = \lambda_2 + \log(\exp(a \exp(\lambda_1)) - \exp(b \exp(\lambda_1))).$$

Notice that that since $\log \log A$ depends only on a, b, c, λ_1 (since λ_2, λ_3 are cancelled in the expression of A), conditioning on $\log \log A$ only makes λ_1 fixed while λ_2 is still the uniform distribution over $[0, \log^3 W]$, which is the same as U_2 . Hence, to show that $\log B$ is close to U_2 , it suffices to show $\log B - \lambda_2$ is small after fixing any valid choice of λ_1, a, b . We can write

$$\begin{aligned}
& |\log(\exp(a \exp(\lambda_1)) - \exp(b \exp(\lambda_1)))| = a \exp(\lambda_1) + |\log(1 - \exp((a - b) \exp(\lambda_1)))| \\
&\leq a \exp(\lambda_1) + \frac{\exp(-\lambda_1)}{b - a} \leq O(\log^2 W) + O(\log \log W / \log W) \leq O(\log^2 W),
\end{aligned}$$

where in the first inequality we again use that $|\log(1 - \exp(-z))| \leq 1/z$ for $z > 0$, and in the second inequality we use $a \leq 1$, $\log \log W \leq \lambda_1 \leq 2 \log \log W$, $b - a \geq 1/\log \log W$. Then, recall that $\log B$ and U_2 are both uniform variables supported on intervals with the same lengths but different offsets (differ by $O(\log^2 W)$). Thus, conditioned on any value of λ_1 , we have

$$d_{TV}(U_2, \log B) \leq O(1/\log W).$$

Lastly, consider the random variables $C \stackrel{\text{def}}{=} \exp(a \exp(\lambda_1)) \exp(\lambda_2) + \lambda_3$. Again, we remark that conditioning on B and A only fixes λ_1, λ_2 . So λ_3 is still a uniform random variable over $[0, \exp(2 \log^3 W)]$, just like U_3 . Hence, the total variation distance between C and U_3 can be bounded by

$$\frac{1}{\exp(2 \log^3 W)} \exp(a \exp(\lambda_1)) \exp(\lambda_2) \leq \exp(\log^3 W + \log^2 W - 2 \log^3 W) \leq O(1/W),$$

where we use the fact $a \leq 1, \lambda_1 \leq 2 \log \log W, \lambda_2 \leq \log^3 W$. This concludes the proof of Lemma 3.10. \square

Let \mathbf{t}, \mathbf{r} be the square edge distributions defined in Definition 3.3. In the lower bound construction from Section 3.3, within each square, (\mathbf{p}, \mathbf{q}) is either $((\mathbf{t} + \mathbf{r})/2, (\mathbf{t} + \mathbf{r})/2)$, (\mathbf{t}, \mathbf{r}) or (\mathbf{r}, \mathbf{t}) . Lemma 3.5 states that, if the tester is only given the order information of three samples, it cannot tell whether the samples are taken from $((\mathbf{t} + \mathbf{r})/2, (\mathbf{t} + \mathbf{r})/2)$ or a random pair from (\mathbf{t}, \mathbf{r}) and (\mathbf{r}, \mathbf{t}) . In order to hide the extra information, one need to apply the transformation specified in Lemma 3.2, which increases the domain size substantially. Here, we argue that applying transformations sampled from $\mathcal{M}(W)$ also eliminates most of the extra information in addition to the order information.

Lemma 3.12. *Let \mathbf{t}, \mathbf{r} be the square edge distributions defined in Definition 3.3. Let $\{u_i, v_i, b_i\}_{i=1}^m$ be samples drawn from the pair of distributions $((\mathbf{t} + \mathbf{r})/2, (\mathbf{t} + \mathbf{r})/2)$. With probability $1/2$, we draw $\{x_i, y_i, \ell_i\}_{i=1}^m$ from (\mathbf{t}, \mathbf{r}) . Otherwise, we draw $\{x_i, y_i, \ell_i\}_{i=1}^m$ from (\mathbf{r}, \mathbf{t}) . Let f_1, f_2, f_3, f_4 be four random mappings drawn independently from $\mathcal{M}(W)$. Then, the quantity*

$$d_{\text{TV}}(\{f_1(u_i), f_2(v_i), b_i\}_{i=1}^m, \{f_3(x_i), f_4(y_i), \ell_i\}_{i=1}^m)$$

is 0 for $m = 1$ and $O(\log \log \log W / \log \log W)$ for $m = 2, 3$.

Proof. We first analyze the case for $m = 1$. We claim that the tuple (u_1, v_1, b_1) has the same distribution as (x_1, y_1, ℓ_1) . since conditioned on any values of (u_1, v_1) , the distribution of b_1 is uniform (and similarly for x_1, y_1, ℓ_1). Then, since f_1, f_2, f_3, f_4 are all identically distributed, it follows the distributions in the two cases are the same.

We then proceed to prove the cases $m = 2, 3$. We remark that the total variation distance for $m = 2$ is at most that for $m = 3$ since one can always explicitly drop the extra sample and this operation will only decrease the total variation distance. Thus, we only need to consider the case $m = 3$. By Lemma 3.5, we have $\text{Order}(\{u_i, v_i, b_i\}_{i=1}^m)$ has the same distribution as $\text{Order}(\{x_i, y_i, \ell_i\}_{i=1}^m)$. Hence, there exists a coupling J between $\{u_i, v_i, b_i\}_{i=1}^m$ and $\{x_i, y_i, \ell_i\}_{i=1}^m$ such that if we sample from J we always have $\text{Order}(\{u_i, v_i, b_i\}_{i=1}^m) = \text{Order}(\{x_i, y_i, \ell_i\}_{i=1}^m)$. Hence, we can bound the overall total variation distance by

$$\mathbb{E}[d_{\text{TV}}(\{f_1(u_i), f_2(v_i), b_i\}_{i=1}^m, \{f_3(x_i), f_4(y_i), \ell_i\}_{i=1}^m)].$$

We remark that the total variation distance inside the expectation is now for fixed values of $\{u_i, v_i, b_i\}_{i=1}^m$ and $\{x_i, y_i, \ell_i\}_{i=1}^m$ that share the same order information and over the random choice of the transformations f_1, f_2, f_3, f_4 . Since the transformations along the two dimensions are picked independently and $b_i = \ell_i$ under the coupling J , we thus have

$$\begin{aligned} & d_{\text{TV}}(\{f_1(u_i), f_2(v_i), b_i\}_{i=1}^m, \{f_3(x_i), f_4(y_i), \ell_i\}_{i=1}^m) \\ &= d_{\text{TV}}(\{f_1(u_i)\}_{i=1}^m, \{f_3(x_i)\}_{i=1}^m) + d_{\text{TV}}(\{f_2(v_i)\}_{i=1}^m, \{f_4(y_i)\}_{i=1}^m). \end{aligned}$$

The arguments for bounding the total variation distance over the two different dimensions are identical. We will therefore just focus on the first dimension. Then consider the event E such that $\min(u_i - u_{i-1}, x_i - x_{i-1}) \geq 1/\log \log W$ for any i . By the union bound, it is easy to see that E does not hold under J with probability at most $O(1/\log \log W)$. By the triangle inequality, we have

$$d_{\text{TV}}(\{f_1(u_i)\}_{i=1}^m, \{f_3(x_i)\}_{i=1}^m) \leq d_{\text{TV}}(\{f_1(u_i)\}_{i=1}^m, D) + d_{\text{TV}}(\{f_3(x_i)\}_{i=1}^m, D)$$

where D is the distribution defined in Lemma 3.10. Conditioned on the event E , we then have that the expression is bounded by $O(\log \log \log W / \log \log W)$. Since the total variation distance is bounded by 1 and E does not hold with probability at most $O(1/\log \log W)$. The overall total variation distance is at most $O(\log \log \log W / \log \log W)$. This completes the proof of Lemma 3.12 \square

Now, let X be an unbiased binary variable. Let \mathbf{p}, \mathbf{q} be a pair of measures generated by the random process described in Section 3.3. Recall that in the construction from the last section, we divide the domain into $\Theta(k^2)$ squares and \mathbf{p}, \mathbf{q} are only supported on the $\Theta(k)$ squares along the diagonal. We will then apply the following domain transformation. For each square along the diagonal, we will independently generate two monotonic mappings $f_1, f_2 \sim \mathcal{M}$. Then, we stretch the square along the x -axis by f_1 and stretch it along the y -axis by f_2 . We will denote the transformed measures as \mathbf{p}', \mathbf{q}' . Let P be the set of samples obtained by taking $\text{Poi}(m)$ samples from \mathbf{p}', \mathbf{q}' . We claim that P reveals little information about X .

Lemma 3.13. *Suppose $m < k$. The mutual information between X and P is at most*

$$I(X : P) \leq O\left(\frac{m^3 \epsilon^2}{k^2} \frac{\log \log \log W}{\log \log W} + \frac{m^7 \epsilon^8}{k^6}\right).$$

Proof. Let P_i be the samples taken from the i -th square. Notice that P_i, P_j for $i \neq j$ are conditionally independent on X . Hence, we have $I(X : P) \leq O(k) I(X : P_1)$. We will use a multiset s made up of elements from \mathbb{R}_+^2 to represent the possible values P_1 can take. Besides, we write $|s|$ to represent the size of the multiset. Then, it holds that

$$I(X : P_1) = O(1) \sum_{\gamma=1}^{\infty} \int_{|s|=\gamma} \frac{(\Pr[P_1 = s|X=0] - \Pr[P_1 = s|X=1])^2}{\Pr[P_1 = s]}.$$

Let H_1 be the indicator variable of whether the first square is selected as a heavy square. We can use techniques similar to the proof of Lemma 3.7 to show that $I(X : P_1)$ can be bounded by

$$\sum_{\gamma=1}^{\infty} O(2^\gamma) \frac{\Pr[|P_1| = \gamma | H_1 = 0]}{\Pr[|P_1| = \gamma, H_1 = 1]} \int_{|s|=\gamma} |\Pr[P_1 = s | X=0, H_1=0] - \Pr[P_1 = s | X=1, H_1=0]|.$$

We will take a closer look at the integral in the expression. Given the observation that $|P_1|$ and X is conditionally independent on H_1 , it is not hard to see that

$$\begin{aligned} & \int_{|s|=\gamma} |\Pr[P_1 = s | X=0, H_1=0] - \Pr[P_1 = s | X=1, H_1=0]| \\ &= 2 d_{\text{TV}}(P_1|(X=0, H_1=0, |P_1|=\gamma), P_1|(X=1, H_1=0, |P_1|=\gamma)) \Pr[|P_1|=\gamma | H_1=0]. \end{aligned}$$

Notice that $P_1|(X=0, H_1=0, |P_1|=\gamma)$ and $P_1|(X=1, H_1=0, |P_1|=\gamma)$ correspond exactly to the distributions of $\{f_1(x_i), f_2(y_i), \ell_i\}_{i=1}^\gamma$ and $\{f_3(u_i), f_4(v_i), b_i\}_{i=1}^\gamma$ specified in Lemma 3.12. Therefore, for $1 \leq \gamma \leq 3$, we can apply Lemma 3.12 and bound the total variation distance by 0 for $\gamma=1$ and $O(\log \log \log W / \log \log W)$ for $\gamma=2, 3$. For $\gamma \geq 4$, we will simply bound the total variation distance by 1. This then allows us to bound $I(X : P_1)$ by

$$O(1) \sum_{\gamma=2}^3 \frac{(\Pr[|P_1|=\gamma | H_i=0])^2}{\Pr[|P_1|=\gamma, H_i=1]} \frac{\log \log \log W}{\log \log W} + O(1) \sum_{\gamma=4}^{\infty} 2^\gamma \frac{(\Pr[|P_1|=\gamma | H_i=0])^2}{\Pr[|P_1|=\gamma, H_i=1]}.$$

Now, recall that $\Pr[H_1=1] = m/k$. When $H_1=1$, the mass of the square will be $1/m$ and hence $|P_1||H_1=1$ will be distributed as $\text{Poi}(1)$. Therefore, we have $\Pr[|P_1|=\gamma | H_i=1] = \text{Poi}(1, \gamma) = \Theta(1)/\gamma!$. On the other hand, when $H_1=0$, the mass of the square is ϵ/k . Hence, $|P_1||H_1=0$ is

distributed as $\text{Poi}(\epsilon m/k)$. Therefore, we have $\Pr[|P_1| = \gamma | H_i = 0] = \text{Poi}(\epsilon m/k, \gamma) \leq (\epsilon m/k)^\gamma / \gamma!$. Together with our assumption $m < k$, we can simplify the bound as

$$I(X : P_1) \leq O(1) \frac{m^3 \epsilon^4}{k^3} \frac{\log \log \log W}{\log \log W} + O(1) \frac{m^7 \epsilon^8}{k^7}.$$

This concludes the proof of Lemma 3.13. \square

We are now ready to conclude the proof of Theorem 3.8.

Proof of Theorem 3.8. Throughout the proof, we assume that $m < k/2$ as this is the regime where we can use the random process described in Section 3.3 to generate measures.

Let X be an unbiased binary variable and \mathbf{p}, \mathbf{q} be a pair of measures generated according to the random process described in Section 3.3 and \mathbf{p}', \mathbf{q}' be the measures obtained after applying the random transformation defined by mappings sampled from $\mathcal{M}(W)$. Since the transformation is monotonic in both x and y axis, we thus have $\|\mathbf{p} - \mathbf{q}\|_{\mathcal{A}_k} = \|\mathbf{p}' - \mathbf{q}'\|_{\mathcal{A}_k}$. Therefore, when $X = 0$, we have $\mathbf{p}' = \mathbf{q}'$; when $X = 1$, we have $\|\mathbf{p}' - \mathbf{q}'\|_{\mathcal{A}_k} > \epsilon$. By Lemma 3.13, we have the mutual information between the random bit X and the output of any algorithm that uses $\text{Poi}(m)$ samples is at most $O\left(\frac{m^3 \epsilon^4}{k^2} \frac{\log \log \log W}{\log \log W} + \frac{m^7 \epsilon^8}{k^6}\right)$. Hence, no tester can reliably distinguish between the case that $\mathbf{p}' = \mathbf{q}'$ and $\|\mathbf{p}' - \mathbf{q}'\|_{\mathcal{A}_k} > \epsilon$ with probability more than $2/3$ unless

$$m \geq \Omega(1) \min \left(k^{2/3} \epsilon^{-4/3} \left(\frac{\log \log \log W}{\log \log W} \right)^{1/3}, k^{6/7} \epsilon^{-8/7} \right). \quad (20)$$

Note that the measures \mathbf{p}', \mathbf{q}' are continuous. The remaining step is to turn them into discrete measures $\tilde{\mathbf{p}}', \tilde{\mathbf{q}}'$ such that distinguishing between $\tilde{\mathbf{p}}' = \tilde{\mathbf{q}}'$ versus $\|\tilde{\mathbf{p}}' - \tilde{\mathbf{q}}'\|_{\mathcal{A}_k} \geq \epsilon$ is about as hard as $\mathbf{p}' = \mathbf{q}'$ versus $\|\mathbf{p}' - \mathbf{q}'\|_{\mathcal{A}_k} \geq \epsilon$.

First, we argue that, for any horizontal or vertical strip of width at most $\frac{\epsilon}{8}$, the mass of \mathbf{p}', \mathbf{q}' is at most $\frac{\epsilon}{8k}$. It is easy to see the claim is true for \mathbf{p}, \mathbf{q} since their marginal distributions in any dimension is uniform over intervals whose lengths add up to at least k . For the transformed distribution \mathbf{p}', \mathbf{q}' , the bound still holds since the transformation only stretches the distribution along x, y axis.

Then, we can construct a grid \mathcal{G} which splits the domain into small unit squares, each of size $\frac{\epsilon}{8} \times \frac{\epsilon}{8}$. Then, consider $\tilde{\mathbf{p}}', \tilde{\mathbf{q}}'$ which round the points falling in each square in \mathcal{G} to its top-left vertex. Then, for an arbitrary rectangle R , $|\mathbf{p}'(R) - \mathbf{q}'(R)| - |\tilde{\mathbf{p}}'(R) - \tilde{\mathbf{q}}'(R)|$ is at most the mass of \mathbf{p}' or \mathbf{q}' in the two vertical strips and the two horizontal strips, each of width at most $\frac{\epsilon}{8}$. Thus, for any R , it holds

$$|\tilde{\mathbf{p}}'(R) - \tilde{\mathbf{q}}'(R)| \geq |\mathbf{p}'(R) - \mathbf{q}'(R)| - \epsilon/(2k).$$

Consequently, it holds $\|\tilde{\mathbf{p}}' - \tilde{\mathbf{q}}'\|_{\mathcal{A}_k} \geq \epsilon/2$ if $\|\mathbf{p}' - \mathbf{q}'\|_{\mathcal{A}_k} \geq \epsilon$. On the other hand, if $\mathbf{p}' = \mathbf{q}'$, it is easy to see that we still have $\tilde{\mathbf{p}}' = \tilde{\mathbf{q}}'$ after the rounding. Thus, if there is an algorithm which can distinguish between the cases $\tilde{\mathbf{p}}' = \tilde{\mathbf{q}}'$ and $\|\tilde{\mathbf{p}}' - \tilde{\mathbf{q}}'\|_{\mathcal{A}_k} > \epsilon/2$, we can use it to distinguish between the cases $\mathbf{p}' = \mathbf{q}'$ and $\|\mathbf{p}' - \mathbf{q}'\|_{\mathcal{A}_k} > \epsilon$ as well by simulating the rounding process. Hence, the sample complexity lower bound in Equation (20) applies to $\tilde{\mathbf{p}}', \tilde{\mathbf{q}}'$ as well.

Finally, we note the supports of the transformed measures \mathbf{p}', \mathbf{q}' are always contained in some $W' \times W'$ square where $W' = k \exp(4 \log^3 W)$. Hence, $\tilde{\mathbf{p}}', \tilde{\mathbf{q}}'$ are over a $V \times V$ discrete grid where $V = \Theta(k \exp(4 \log^3 W)/\epsilon)$. If the first term in the sample complexity bound (Equation (20)) is dominating, we must have $\log \log W \geq \log(k/\epsilon)$. This then implies that V is at most $\exp(5 \log^3 W)$,

which further implies that $\log \log W \geq \Omega(1) \log \log V$. On the other hand, it is easy to see that $V > W$ and so $1/\log \log \log W > 1/\log \log \log V$. We can then rewrite Equation (20) as

$$m \geq \Omega(1) \min \left(k^{2/3} \epsilon^{-4/3} \left(\frac{\log \log V}{\log \log \log V} \right)^{1/3}, k^{6/7} \epsilon^{-8/7} \right),$$

which is indeed the desired lower bound. □

4 Conclusions and Open Problems

In this work, we studied the problem of closeness testing between two multidimensional distributions under the \mathcal{A}_k distance. Our main contribution is the first tester for this task with sublinear sample complexity. The sample complexity of our tester is provably near-optimal as a function of the parameter k (within logarithmic factors) for any fixed dimension $d \geq 2$.

Conceptually, our sample complexity lower bound implies that the testing problem is provably harder in the multidimensional setting. In particular, there is a “phase transition” between the one-dimensional and the two-dimensional cases. On the positive side, we show that as the dimension d further increases the dependency of the sample complexity on k — the main parameter of our interest — stays approximately the same.

As immediate corollaries of our \mathcal{A}_k closeness tester, we also obtain the first closeness tester for families of structured multidimensional distributions — including k -histograms and uniform distributions over unions of axis-aligned rectangles — under the total variation distance.

While Theorem 1.2 implies that our upper and lower bounds are nearly optimal in terms of their dependence on k , their dependence on ϵ do not match. In particular, the upper bound scales polynomially with $1/\epsilon$, where the degree of the polynomial depends on the dimension d . On the other hand, the lower bound applies to 2-dimensional distributions, and hence has a constant exponent in its (polynomial) ϵ -dependence. This leads to the following question.

Question 4.1. What is the optimal sample complexity as a function of ϵ for multidimensional \mathcal{A}_k closeness testing?

In the current and prior works, the multidimensional \mathcal{A}_k -distance is defined as the maximum discrepancy between two distributions over k disjoint axis-aligned *rectangles*. On the other hand, the \mathcal{A}_k -distance for univariate distributions is defined with respect to *intervals*. This definition inherently uses axis-aligned rectangles in \mathbb{R}^d , as the natural generalization of intervals in \mathbb{R} . Yet, rectangles are not necessarily the only valid choice. More specifically, one can replace axis-aligned rectangles in the definition of multidimensional \mathcal{A}_k distance with other geometric shapes whose 1-dimensional projection corresponds to intervals. For example, we can use shapes like unit-balls, simplices, or any other convex set. Such natural variants of multidimensional \mathcal{A}_k distance can be used to build d_{TV} -closeness testers of other families of structured distributions, such as log-concave distributions. This leads to the following question.

Question 4.2. Are there alternative definitions of multidimensional \mathcal{A}_k distance for multivariate distributions that can lead to optimal d_{TV} -closeness/identity testers for other multivariate shape-restricted distributions?

Exploring other notions of multidimensional \mathcal{A}_k distance is of significant interest and may lead to a unified theory of testing multivariate structured distributions.

References

- [ADH⁺15] J Acharya, I. Diakonikolas, C. Hegde, J. Li, and L. Schmidt. Fast and near-optimal algorithms for approximating distributions by histograms. In Tova Milo and Diego Calvanese, editors, *Proceedings of the 34th ACM Symposium on Principles of Database Systems, PODS 2015*, pages 249–263. ACM, 2015.
- [ADJ⁺11] J. Acharya, H. Das, A. Jafarpour, A. Orlitsky, and S. Pan. Competitive closeness testing. *Journal of Machine Learning Research - Proceedings Track*, 19:47–68, 2011.
- [ADK15] J. Acharya, C. Daskalakis, and G. Kamath. Optimal testing for properties of distributions. In *NeurIPS*, pages 3591–3599, 2015.
- [ADLS17] J. Acharya, I. Diakonikolas, J. Li, and L. Schmidt. Sample-optimal density estimation in nearly-linear time. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2017*, pages 1278–1289, 2017. Full version available at <https://arxiv.org/abs/1506.00671>.
- [BBBB72] R.E. Barlow, D.J. Bartholomew, J.M. Bremner, and H.D. Brunk. *Statistical Inference under Order Restrictions*. Wiley, New York, 1972.
- [BDKR02] T. Batu, S. Dasgupta, R. Kumar, and R. Rubinfeld. The complexity of approximating entropy. In *ACM Symposium on Theory of Computing*, pages 678–687, 2002.
- [BDS10] Y. Bengio, O. Delalleau, and C. Simard. Decision trees do not generalize to new variations. *Computational Intelligence*, 26(4):449–467, 2010.
- [BFF⁺01] T. Batu, E. Fischer, L. Fortnow, R. Kumar, R. Rubinfeld, and P. White. Testing random variables for independence and identity. In *Proc. 42nd IEEE Symposium on Foundations of Computer Science*, pages 442–451, 2001.
- [BFR⁺00] T. Batu, L. Fortnow, R. Rubinfeld, W. D. Smith, and P. White. Testing that distributions are close. In *IEEE Symposium on Foundations of Computer Science*, pages 259–269, 2000.
- [Bic69] P. J. Bickel. A Distribution Free Version of the Smirnov Two Sample Test in the p -Variate Case. *The Annals of Mathematical Statistics*, 40(1):1 – 23, 1969.
- [BKR04] T. Batu, R. Kumar, and R. Rubinfeld. Sublinear algorithms for testing monotone and unimodal distributions. In *ACM Symposium on Theory of Computing*, pages 381–390, 2004.
- [BPGB20] S. Bruch, J. Pfeifer, and M. Guillame-Bert. Learning representations for axis-aligned decision forests through input perturbation. *arXiv preprint arXiv:2007.14761*, 2020.
- [Can22] C. L. Canonne. Topics and techniques in distribution testing: A biased but representative sample. *Found. Trends Commun. Inf. Theory*, 19(6):1032–1198, 2022.
- [CDKL22] C. L. Canonne, I. Diakonikolas, D. M. Kane, and S. Liu. Near-optimal bounds for testing histogram distributions. *CoRR*, abs/2207.06596, 2022. Conference version in *NeurIPS’22*.

- [CDKS17] C. L. Canonne, I. Diakonikolas, D. M. Kane, and A. Stewart. Testing bayesian networks. In *Proceedings of the 30th Conference on Learning Theory, COLT 2017*, pages 370–448, 2017.
- [CDKS18] C. L. Canonne, I. Diakonikolas, D. M. Kane, and A. Stewart. Testing conditional independence of discrete distributions. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018, Los Angeles, CA, USA, June 25-29, 2018*, pages 735–748. ACM, 2018.
- [CDSS13] S. Chan, I. Diakonikolas, R. Servedio, and X. Sun. Learning mixtures of structured distributions over discrete domains. In *SODA*, pages 1380–1394, 2013.
- [CDSS14a] S. Chan, I. Diakonikolas, R. Servedio, and X. Sun. Efficient density estimation via piecewise polynomial approximation. In *STOC*, pages 604–613, 2014.
- [CDSS14b] S. Chan, I. Diakonikolas, R. Servedio, and X. Sun. Near-optimal density estimation in near-linear time using variable-width histograms. In *NIPS*, pages 1844–1852, 2014.
- [CDVV14] S. O. Chan, I. Diakonikolas, P. Valiant, and G. Valiant. Optimal algorithms for testing closeness of discrete distributions. In *Proceedings of the twenty-fifth annual ACM-SIAM symposium on Discrete algorithms*, pages 1193–1203. SIAM, 2014.
- [CJKL22] C. L. Canonne, A. Jain, G. Kamath, and J. Li. The price of tolerance in distribution testing. In *Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pages 573–624. PMLR, 2022.
- [CLM20] S. Chen, J. Li, and A. Moitra. Efficiently learning structured distributions from untrusted batches. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing, STOC 2020, 2020*, pages 960–973. ACM, 2020.
- [CMN98] S. Chaudhuri, R. Motwani, and V. R. Narasayya. Random sampling for histogram construction: How much is enough? In *SIGMOD Conference*, pages 436–447, 1998.
- [DDK18] C. Daskalakis, N. Dikkala, and G. Kamath. Testing ising models. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2018, 2018*, pages 1989–2007. SIAM, 2018.
- [DDS12] C. Daskalakis, I. Diakonikolas, and R.A. Servedio. Learning k -modal distributions via testing. In *SODA*, pages 1371–1385, 2012.
- [DDS⁺13] C. Daskalakis, I. Diakonikolas, R. Servedio, G. Valiant, and P. Valiant. Testing k -modal distributions: Optimal algorithms via reductions. In *SODA*, pages 1833–1852, 2013.
- [DGK⁺21] I. Diakonikolas, T. Gouleakis, D. M. Kane, J. Peebles, and E. Price. Optimal testing of discrete distributions with high probability. In *STOC '21: 53rd Annual ACM SIGACT Symposium on Theory of Computing, 2021*, pages 542–555. ACM, 2021.
- [DK16] I. Diakonikolas and D. M. Kane. A new approach for testing properties of discrete distributions. In *2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 685–694. IEEE, 2016.
- [DKN15a] I. Diakonikolas, D. M. Kane, and V. Nikishkin. Optimal algorithms and lower bounds for testing closeness of structured distributions. In *2015 IEEE 56th Annual Symposium on Foundations of Computer Science*, pages 1183–1202. IEEE, 2015.

- [DKN15b] I. Diakonikolas, D. M. Kane, and V. Nikishkin. Testing identity of structured distributions. In *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2015*, pages 1841–1854, 2015.
- [DKN17] I. Diakonikolas, D. M. Kane, and V. Nikishkin. Near-optimal closeness testing of discrete histogram distributions. In *44th International Colloquium on Automata, Languages, and Programming, ICALP 2017*, pages 8:1–8:15, 2017.
- [DKP19] I. Diakonikolas, D. M. Kane, and J. Peebles. Testing identity of multidimensional histograms. In *Conference on Learning Theory*, pages 1107–1131. PMLR, 2019.
- [DL01] L. Devroye and G. Lugosi. *Combinatorial methods in density estimation*. Springer Series in Statistics, Springer, 2001.
- [DL04] L. Devroye and G. Lugosi. Bin width selection in multivariate histograms by the combinatorial method. *Test*, 13(1):129–145, 2004.
- [DLS18] I. Diakonikolas, J. Li, and L. Schmidt. Fast and sample near-optimal algorithms for learning multidimensional histograms. In *Conference On Learning Theory, COLT 2018*, volume 75 of *Proceedings of Machine Learning Research*, pages 819–842. PMLR, 2018.
- [FD81] D. Freedman and P. Diaconis. On the histogram as a density estimator: theory. *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 57(4):453–476, 1981.
- [FR79] J. H. Friedman and L. C. Rafsky. Multivariate Generalizations of the Wald-Wolfowitz and Smirnov Two-Sample Tests. *The Annals of Statistics*, 7(4):697 – 717, 1979.
- [GGI⁺02] A. C. Gilbert, S. Guha, P. Indyk, Y. Kotidis, S. Muthukrishnan, and M. Strauss. Fast, small-space algorithms for approximate histogram maintenance. In *STOC*, pages 389–398, 2002.
- [GGR98] O. Goldreich, S. Goldwasser, and D. Ron. Property testing and its connection to learning and approximation. *Journal of the ACM*, 45:653–750, 1998.
- [GJ14] P. Groeneboom and G. Jongbloed. *Nonparametric Estimation under Shape Constraints: Estimators, Algorithms and Asymptotics*. Cambridge University Press, 2014.
- [GKS06] S. Guha, N. Koudas, and K. Shim. Approximation and streaming algorithms for histogram construction problems. *ACM Trans. Database Syst.*, 31(1):396–438, 2006.
- [GR00] O. Goldreich and D. Ron. On testing expansion in bounded-degree graphs. Technical Report TR00-020, Electronic Colloquium on Computational Complexity, 2000.
- [Hen88] N. Henze. A Multivariate Two-Sample Test Based on the Number of Nearest Neighbor Type Coincidences. *The Annals of Statistics*, 16(2):772 – 783, 1988.
- [ILR12] P. Indyk, R. Levi, and R. Rubinfeld. Approximating and Testing k -Histogram Distributions in Sub-linear Time. In *PODS*, pages 15–22, 2012.
- [Ing94] Y. I. Ingster. Minimax detection of a signal in ℓ_p -metrics. *Journal of Mathematical Sciences*, 68(4):503–515, 1994.
- [Ing97] Y. I. Ingster. Adaptive chi-square tests. *Zapiski Nauchnykh Seminarov POMI*, 244:150–166, 1997.

- [IS03] Y. I. Ingster and I. A. Suslina. *Nonparametric Goodness-of-fit Testing Under Gaussian Models*, volume 169. Springer Series in Statistics, Springer, 2003.
- [JKM⁺98] H. V. Jagadish, N. Koudas, S. Muthukrishnan, V. Poosala, K. C. Sevcik, and T. Suel. Optimal histograms with quality guarantees. In *VLDB*, pages 275–286, 1998.
- [JPZ97] A. Justel, D. Pena, and R. Zamar. A multivariate kolmogorov-smirnov test of goodness of fit. *Statistics & Probability Letters*, 35(3):251–259, 1997.
- [Kle09] J. Klemela. Multivariate histograms with data-dependent partitions. *Statistica Sinica*, 19(1):159–176, 2009.
- [Kru53] J. B. Kruskal. Monotonic subsequences. *Proceedings of the American Mathematical Society*, 4(2):264–274, 1953.
- [KS16] A. K. H. Kim and R. J. Samworth. Global rates of convergence in log-concave density estimation. *The Annals of Statistics*, 44(6):2756–2779, 2016.
- [LN96] G. Lugosi and A. Nobel. Consistency of data-driven histogram methods for density estimation and classification. *Ann. Statist.*, 24(2):687–706, 04 1996.
- [LR05] E. L. Lehmann and J. P. Romano. *Testing statistical hypotheses*. Springer Texts in Statistics. Springer, 2005.
- [LRR11] R. Levi, D. Ron, and R. Rubinfeld. Testing properties of collections of distributions. In *ICS*, pages 179–194, 2011.
- [NP33] J. Neyman and E. S. Pearson. On the problem of the most efficient tests of statistical hypotheses. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, 231(694-706):289–337, 1933.
- [Pan08] L. Paninski. A coincidence-based test for uniformity given very sparsely-sampled discrete data. *IEEE Transactions on Information Theory*, 54:4750–4755, 2008.
- [Pea00] K. Pearson. On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. *Philosophical Magazine Series 5*, 50(302):157–175, 1900.
- [RS96] R. Rubinfeld and M. Sudan. Robust characterizations of polynomials with applications to program testing. *SIAM Journal on Computing*, 25:252–271, 1996.
- [Rub12] R. Rubinfeld. Taming big probability distributions. *XRDS*, 19(1):24–28, 2012.
- [Sco79] D. W. Scott. On optimal and data-based histograms. *Biometrika*, 66(3):605–610, 1979.
- [Sco92] D.W. Scott. *Multivariate Density Estimation: Theory, Practice and Visualization*. Wiley, New York, 1992.
- [TGIK02] N. Thaper, S. Guha, P. Indyk, and N. Koudas. Dynamic multidimensional histograms. In *SIGMOD Conference*, pages 428–439, 2002.

- [Val11] P. Valiant. Testing symmetric properties of distributions. *SIAM J. Comput.*, 40(6):1927–1968, 2011.
- [VV11] G. Valiant and P. Valiant. Estimating the unseen: an $n/\log(n)$ -sample estimator for entropy and support size, shown optimal via new CLTs. In *STOC*, pages 685–694, 2011.
- [VV14] G. Valiant and P. Valiant. An automatic inequality prover and instance optimal identity testing. In *FOCS*, 2014.
- [WN07] R. Willett and R. D. Nowak. Multiscale poisson intensity and density estimation. *IEEE Transactions on Information Theory*, 53(9):3171–3187, 2007.
- [YA01] C. T. Yildiz and E. Alpaydin. Omnivariate decision trees. *IEEE Transactions on Neural Networks*, 12(6):1539–1546, 2001.