

Active Tolerant Testing

Avrim Blum

Toyota Technological Institute at Chicago
Chicago, USA
avrim@ttic.edu

Lunjia Hu

Institute for Interdisciplinary Information Sciences
Tsinghua University
Beijing, China
hulj14@mails.tsinghua.edu.cn

November 2, 2017

Abstract

In this work, we give the first algorithms for tolerant testing of nontrivial classes in the active model: estimating the distance of a target function to a hypothesis class \mathcal{C} with respect to some arbitrary distribution \mathcal{D} , using only a small number of label queries to a polynomial-sized pool of unlabeled examples drawn from \mathcal{D} . Specifically, we show that for the class \mathcal{C} of unions of d intervals on the line, we can estimate the error rate of the best hypothesis in the class to an additive error ϵ from only $O(\frac{1}{\epsilon^6} \log \frac{1}{\epsilon})$ label queries to an unlabeled pool of size $O(\frac{d}{\epsilon^2} \log \frac{1}{\epsilon})$. The key point here is the number of labels needed is independent of the VC-dimension of the class. This extends the work of Balcan et al. [2012] who solved the *non*-tolerant testing problem for this class (distinguishing the zero-error case from the case that the best hypothesis in the class has error greater than ϵ).

We also consider the related problem of estimating the performance of a given learning algorithm \mathcal{A} in this setting. That is, given a large pool of unlabeled examples drawn from distribution \mathcal{D} , can we, from only a few label queries, estimate how well \mathcal{A} would perform if the entire dataset were labeled? We focus on k -Nearest Neighbor style algorithms, and also show how our results can be applied to the problem of hyperparameter tuning (selecting the best value of k for the given learning problem).

1 Introduction

Suppose you are about to embark on a project to label a large quantity of data, such as medical images or street scenes. Your intent is to then feed this data into your favorite learning algorithm for, say, a medical diagnosis or robotic car application. Before embarking on this project, can you, from just a few labeled examples, estimate *how well* your algorithm can be expected to perform when trained on the large sample? We consider this question in two related contexts.

Tolerant testing: The first context we consider is that of tolerant testing, or approximating the distance of a target function f to a hypothesis class \mathcal{C} . Specifically, consider a hypothesis class \mathcal{C} of VC-dimension d , where d should be thought of as large. If we wish to find the ϵ -approximately-best hypothesis in \mathcal{C} , we will need roughly $O(d/\epsilon)$ labeled examples in the (realizable) case that $f \in \mathcal{C}$, or $O(d/\epsilon^2)$ labeled examples in the (agnostic) case that $f \notin \mathcal{C}$. However, if we just want to estimate

the error rate of the best hypothesis in \mathcal{C} (rather than to *find* the hypothesis), can we do this from less data?

The problem of determining *whether* one will be able to learn well using a given hypothesis class \mathcal{C} using substantially less labeled data than needed to actually (attempt to) learn well using that class is the problem of *passive* and *active* property testing, studied by Kearns and Ron [1998] and Balcan et al. [2012]. That work considers the problem of distinguishing the case that (a) the target function f belongs to class \mathcal{C} from (b) the target function f is ϵ -far from any concept in \mathcal{C} with respect to the underlying data distribution \mathcal{D} . For instance, suppose our data consists of points x on the real line, labeled by f as positive or negative, and we are interested in learning using the class \mathcal{C} consisting of unions of d intervals. This class has VC-dimension $2d$ and so would require $\Omega(d)$ labeled examples to learn. However, Balcan et al. [2012] show that in the active testing framework (one can sample *poly*(d) *unlabeled* examples for free and then query for the labels of a small number of those examples), one can solve the testing problem using only a constant number of label queries (when ϵ is constant), independent of d .

One limitation of these results, however, is that they do not guarantee to give a meaningful answer when the target function is “almost” in the class \mathcal{C} . For instance, suppose f can be perfectly described by a union of 10,000 intervals but is $\epsilon/2$ -close to a union of 100 intervals. Then we would like a tester that can say “good enough” at $d = 100$ rather than telling us that we need $d = 10,000$. The tester of Balcan et al. [2012], unfortunately, seems to require f to be $O(\epsilon^3)$ -close to a union of d intervals in order to guarantee an output of YES, which is much less than ϵ .

In this work, we give algorithms for such *tolerant testing* [Parnas et al., 2006] for the case of unions of intervals and a few related classes. We can distinguish the case that the best function in \mathcal{C} has error rate $\geq 2\epsilon$ from the case that the best function in \mathcal{C} has error rate $\leq \epsilon$, and more generally we can estimate the error rate α of the best function in the class up to $\pm\epsilon$. Thus, for the first time, from a small number of label queries, we can solve the property-testing analog of the notion of agnostic learning.

One point we wish to make up front: while the classes of functions we consider are fairly simple, such as unions of intervals on the line, we are operating in a challenging model. We would like algorithms that work for any (unknown) underlying data distribution \mathcal{D} , not just the uniform distribution, *and* we want algorithms that only query for labels from among examples seen in a poly-sized sample of unlabeled data drawn from \mathcal{D} rather than querying arbitrary points in the input space. These are important conditions for being able to use property testing for machine learning problems.

Algorithm estimation: The second context we consider is that we are given a learning algorithm \mathcal{A} and a large unlabeled sample S of N examples drawn from distribution \mathcal{D} . If we were to label all N examples of S and feed them into algorithm \mathcal{A} , then \mathcal{A} would produce some hypothesis (call it h_S) with some error rate α . What we would like to do is, by labeling only very few examples in S , and perhaps a few additional examples drawn from \mathcal{D} , to estimate the value of α (so that we can determine whether our project of labeling all examples in S is worthwhile).

To get a feel for this problem, one algorithm \mathcal{A} for which this task is easy is 1-Nearest Neighbor (1-NN). This algorithm would produce a hypothesis h_S that on any given query point x predicts the label of the example $x' \in S$ that is nearest to x . For this algorithm, we can easily estimate the error rate of h_S from just a few label queries by repeatedly drawing a random x from \mathcal{D} , finding

the point $x' \in S$ that is closest to x , and then requesting the labels of x and x' to see if they agree. We only need to repeat this process $O(1/\epsilon^2)$ times in order to estimate the error rate of h_S to $\pm\epsilon$. This works because h_S is constructed, and makes predictions, in a very local way.¹ In this work, we extend this to different forms of k -Nearest Neighbor algorithms, where the prediction on some point x depends on the k nearest examples in S , developing testers for which the number of queries *does not depend on k* . This then allows us to use this for *hyperparameter tuning*: determining the (approximately) best value of $k \in \{1, \dots, N\}$ for the given application.

We note that there are three natural but somewhat different ways to model the task of estimating the error rate of algorithm \mathcal{A} trained on dataset S . Let $\text{error}(h_S)$ denote the error rate of hypothesis h_S with respect to distribution \mathcal{D} , and let $\hat{\alpha}$ be the output of the tester \mathcal{T} that estimates $\text{error}(h_S)$. In the first model, we require that $\hat{\alpha}$ be a good estimate of $\text{error}(h_S)$ with probability at least $\frac{2}{3}$ for *any* training set S , even sets S not drawn from \mathcal{D} . In the second model, we only require that \mathcal{T} be accurate when S is drawn from \mathcal{D} (that is, the $\frac{2}{3}$ probability is over both the internal randomness in \mathcal{T} and in the draw of S). Finally, in the third model, S is drawn from \mathcal{D} but \mathcal{T} does not have access to it: instead, \mathcal{T} has the ability to draw (a polynomial number of) fresh unlabeled examples and to query points from them. That is,

1. In the first model, we require that $\forall S, \Pr_{\mathcal{T}(S)}[|\hat{\alpha} - \text{error}(h_S)| \leq \epsilon] \geq \frac{2}{3}$.
2. In the second model, we require that $\Pr_{S, \mathcal{T}(S)}[|\hat{\alpha} - \text{error}(h_S)| \leq \epsilon] \geq \frac{2}{3}$.
3. In the third model, we require that $\Pr_{\mathcal{T}}[|\hat{\alpha} - \mathbb{E}_S[\text{error}(h_S)]| \leq \epsilon] \geq \frac{2}{3}$.

Roughly, the first model is the hardest while the third model is the easiest. All our upper bounds and lower bounds in this paper apply to all three models with slight modifications, though for simplicity of presentation we focus on the second model throughout the paper.

2 Our Results

In this paper, we show (Theorem 5) that in the active testing model [Balcan et al., 2012], there is an algorithm that approximates the distance from a function to the class of unions of d intervals on the line up to an additive error ϵ using $O(\frac{1}{\epsilon^8} \log \frac{1}{\epsilon})$ label queries on $O(\frac{d}{\epsilon^2} \log \frac{1}{\epsilon})$ unlabeled examples, even when the data distribution is unknown to the algorithm.

To achieve this result, we propose the notion of *compositions of additive properties* (Section 3.3) and prove the Composition Lemma (Lemma 2) that to approximate the distance to any composition of m additive properties on a semi-uniform distribution up to an additive error ϵ , we only need a distance approximation oracle for compositions of only $O(\frac{1}{\epsilon\mu^2} + \frac{1}{\epsilon^2})$ additive properties, though this may produce a bi-criteria approximation that depends on μ . See Section 3.3 for definitions.

¹In contrast, note that estimating the error rate of this algorithm would require a large labeled sample if we only *passively* receive labeled examples. Specifically, suppose the distribution \mathcal{D} is uniform over c clusters and the 1-NN algorithm aims to use $N = c \log \frac{c}{\delta}$ examples, so that with probability at least $1 - \delta$, every cluster has at least one training example in it. We want to distinguish the following two cases: either every cluster is pure but random so the error rate is roughly 0, or every cluster is 50/50 so the error rate is roughly $\frac{1}{2}$. To distinguish these cases, the tester needs to see at least two labels in the same cluster, implying an $\Omega(\sqrt{c}) = \Omega(\sqrt{N/\log N})$ sample complexity lower bound.

The Composition Lemma implies an (ϵ, μ) -bi-criteria distance approximation algorithm for unions of d intervals on the uniform distribution over $[0, 1]$ using $O((\frac{1}{\epsilon^3 \mu^3} + \frac{1}{\epsilon^4 \mu}) \log \frac{1}{\epsilon})$ queries on $O(\frac{d}{\epsilon^2} \log \frac{1}{\epsilon})$ unlabeled examples; in particular, we estimate the error to an additive $\pm \epsilon$ and the number of intervals to a multiplicative factor $1 + \mu$. We then show (Lemma 8) how to remove the approximation in number of intervals and get a uni-criterion distance approximation algorithm.

To generalize the result to arbitrary unknown distributions, we show a general relationship between query testing and active testing for arbitrary distributions in Theorems 6 and 7, which also improves a previous result in [Balcan et al., 2012] by showing that the unlabeled sample complexity of non-tolerant property testing for unions of d intervals on arbitrary unknown distributions can be reduced to $O(\frac{d}{\epsilon} \log \frac{1}{\epsilon})$, from $O(\frac{d^2}{\epsilon^6})$ in their original paper.

For the k -Nearest Neighbor (k -NN) algorithm with soft predictions and p th-power loss (the prediction on a point x is the average label of the k nearest examples to x in an unlabeled sample of size N , and we use the p th-power loss to penalize mistakes) we show in Theorem 9 that this loss can be estimated up to an additive error ϵ using $O(\frac{p}{\epsilon^2})$ queries on $N + O(\frac{1}{\epsilon^2})$ unlabeled examples, even when the data distribution is unknown to the tester. The same result also holds for Weighted Nearest Neighbor algorithms, where the prediction on a point x is a weighted average of the labels of all the examples depending on their distances to x . For the $O(\frac{p}{\epsilon^2})$ query complexity upper bound, we show a matching lower bound (Theorem 13). In the case where k is a quantity to be optimized, we show an algorithm that finds an approximately-best choice of k up to an additive error ϵ using roughly $O(\frac{p^2 \log N}{\epsilon^3})$ queries on roughly $N + O(\frac{p \log N}{\epsilon^3})$ unlabeled examples. For k -NN with hard predictions (the prediction is a strict majority vote over the k nearest neighbors), there is a simple algorithm for approximating its accuracy up to an additive error ϵ using $O(\frac{k}{\epsilon^2})$ queries on $N + O(\frac{1}{\epsilon^2})$ unlabeled examples. By a reduction (Theorem 14) from *approximating the number of good arms* (AGA, see Section 6.4 for definition) in the stochastic multi-armed bandit setting, we show that the query complexity cannot be improved beyond $O(\frac{k}{\epsilon \log \frac{1}{\epsilon}})$ (Theorem 15), and cannot even be improved beyond the current $O(\frac{k}{\epsilon^2})$ upper bound if we assume the natural algorithm for AGA is optimal in query complexity.

3 Preliminaries and Related Work

3.1 Property Testing, Tolerant Testing and Distance Approximation

Suppose we have a ground set X and a distribution \mathcal{D} over X . For any two binary functions $f, g \in \{0, 1\}^X$, we define their distance to be $\text{dist}_{\mathcal{D}}(f, g) = \Pr_{x \sim \mathcal{D}}[f(x) \neq g(x)]$.

Suppose we also have a concept class $\mathcal{C} \subseteq \{0, 1\}^X$. Given a function $f \in \{0, 1\}^X$ and a margin ϵ as input, the task of property testing $\text{PT}_{\mathcal{D}}(f, \epsilon)$ [Rubinfeld and Sudan, 1996] is to distinguish the case that f belongs to class \mathcal{C} from the case that f is ϵ -far from \mathcal{C} . In other words, $\forall f$,

1. if $f \in \mathcal{C}$, the algorithm outputs “YES” with probability at least $\frac{2}{3}$;
2. if $\forall g \in \mathcal{C}, \text{dist}_{\mathcal{D}}(f, g) > \epsilon$, the algorithm outputs “NO” with probability at least $\frac{2}{3}$.

The function f can be given to the algorithm in many different ways. In the *query testing* framework [Rubinfeld and Sudan, 1996], the algorithm can query the value of $f(x)$ for any $x \in X$.

In this framework, we say the algorithm has *query access* to f , or has access to f_{query} . The query complexity of the algorithm, as a function of $\frac{1}{\epsilon}$, is measured by the maximum number of queries needed by the algorithm.

Balcan et al. [2012] argued that the query testing framework is not realistic for machine learning practice. They proposed the *active testing* framework, in which the algorithm first requests N unlabeled examples $x_1, x_2, \dots, x_N \in X$ sampled independently according to \mathcal{D} and can only choose to query $f(x_i)$ for $1 \leq i \leq N$. In this framework, we say the algorithm has *active access* to f , or has access to f_{active} . The maximum value of N , as a function of $\frac{1}{\epsilon}$, is called the *unlabeled sample complexity*. The query complexity is still defined as a function of $\frac{1}{\epsilon}$ measuring the maximum number of queries needed by the algorithm.

Goldreich et al. [1998] and Kearns and Ron [1998] studied an even more strict way of accessing f , called *passive access*, in which the algorithm is given the label of an example chosen independently at random from \mathcal{D} for each query the algorithm makes.

Tolerant testing $\text{TT}_{\mathcal{D}}(f, \alpha, \epsilon)$ [Parnas et al., 2006] is a similar task to property testing. The only difference is that besides the margin ϵ , we are given another parameter α as input, and we are asked to distinguish the case that f is α -close to \mathcal{C} from the case that f is $(\alpha + \epsilon)$ -far from \mathcal{C} . The query complexity of tolerant testing is still measured as a function of $\frac{1}{\epsilon}$ [Parnas et al., 2006].

A natural generalization of tolerant testing is distance approximation, in which we are only given the function f and the margin ϵ as input and required to output $\hat{\alpha}$ as an approximation of the distance from f to \mathcal{C} up to an additive error ϵ . More specifically, the goal of $\text{DA}_{\mathcal{D}}(f, \epsilon)$ is to output $\hat{\alpha}$ such that $\forall \alpha$,

1. $\forall f$ such that $\exists g \in \mathcal{C}, \text{dist}_{\mathcal{D}}(f, g) \leq \alpha$, it holds with probability at least $\frac{2}{3}$ that $\hat{\alpha} \leq \alpha + \epsilon$;
2. $\forall f$ such that $\forall g \in \mathcal{C}, \text{dist}_{\mathcal{D}}(f, g) > \alpha$, it holds with probability at least $\frac{2}{3}$ that $\hat{\alpha} > \alpha - \epsilon$.

Because for any \mathcal{D} and ϵ , it's clear that a $\text{DA}_{\mathcal{D}}(f, \frac{\epsilon}{2})$ algorithm implies a $\text{TT}_{\mathcal{D}}(f, \alpha, \epsilon)$ algorithm with the same query and unlabeled sample complexity [Parnas et al., 2006], we focus on distance approximation rather than tolerant testing throughout the paper.

Obviously, as pointed out by Balcan et al. [2012], a $\text{PT}_{\mathcal{D}}(f_{\text{active}}, \epsilon)$ algorithm implies a $\text{PT}_{\mathcal{D}}(f_{\text{query}}, \epsilon)$ algorithm with the same query complexity when \mathcal{D} is known to the algorithm, since it can always then create unlabeled data on its own; this also holds for TT and DA. Here, we show a theorem in the reverse direction for bounds that are worst-case over distributions \mathcal{D} . Specifically, we show in Theorem 6 that a $\text{PT}_{\mathcal{D}}(f_{\text{query}}, \frac{\epsilon}{2})$ algorithm can induce a $\text{PT}_{\mathcal{D}}(f_{\text{active}}, \epsilon)$ algorithm with (except for a constant factor) the same query complexity and reasonable unlabeled sample complexity, under the assumption that the $\text{PT}_{\mathcal{D}}(f_{\text{query}}, \frac{\epsilon}{2})$ algorithm never queries examples outside the support of \mathcal{D} , which holds in all normal cases. We extend the theorem to DA in Theorem 7.

3.2 Unions of d Intervals

We use $\mathcal{I}(d) \subseteq \{0, 1\}^{\mathbb{R}}$ to denote the class of functions $f \in \{0, 1\}^{\mathbb{R}}$ satisfying that $f^{-1}(1)$ can be written as a union of at most d intervals. Note that for $d \in \mathbb{N}$, the VC-dimension of $\mathcal{I}(d)$ is $2d$.

We use $\mathcal{I}_{\mathcal{D}}(d, \alpha)$ to denote the class of functions that are α -close to $\mathcal{I}(d)$, i.e. $\mathcal{I}_{\mathcal{D}}(d, \alpha) = \{f \in \{0, 1\}^{\mathbb{R}} : \exists g \in \mathcal{I}(d), \text{dist}_{\mathcal{D}}(f, g) \leq \alpha\}$. Using this notation, property testing for unions of d intervals is to distinguish $f \in \mathcal{I}(d)$ and $f \notin \mathcal{I}_{\mathcal{D}}(d, \epsilon)$.

In previous work, Kearns and Ron [1998] showed that in the query testing framework, when \mathcal{D} is the uniform distribution over $[0, 1]$, there is a bi-criteria testing algorithm that can distinguish $f \in \mathcal{I}(d)$ and $f \notin \mathcal{I}_{\mathcal{D}}(\frac{d}{\epsilon}, \epsilon)$ using $O(\frac{1}{\epsilon})$ queries. Balcan et al. [2012] improved this work by showing that in the active testing framework, there is a uni-criterion testing algorithm that can distinguish $f \in \mathcal{I}(d)$ and $f \notin \mathcal{I}_{\mathcal{D}}(d, \epsilon)$ using $O(\frac{1}{\epsilon^4})$ queries on $O(\frac{\sqrt{d}}{\epsilon^3})$ unlabeled examples. Their algorithm can be generalized to a testing algorithm that distinguishes $f \in \mathcal{I}_{\mathcal{D}}(d, \epsilon_1)$ and $f \notin \mathcal{I}_{\mathcal{D}}(d, \epsilon)$ when $\epsilon_1 = O(\epsilon^3)$ using the same number of queries and unlabeled examples. They generalized the result from uniform distribution on $[0, 1]$ to any unknown distribution by taking the advantage of unlabeled examples to approximate the CDF of the distribution to enough accuracy using $O(\frac{d^2}{\epsilon^6})$ unlabeled examples. This unlabeled sample complexity is improved to $O(\frac{d}{\epsilon} \log \frac{1}{\epsilon})$ in our paper (Section 5.1).

3.3 Composition of Additive Properties

Balcan et al. [2012] showed that disjoint unions of testable properties are testable in the non-tolerant, active model. We extend their result to distance approximation in Appendix A. Here, we propose a more general notion of a certain concept class formed by composing smaller concept classes on disjoint ground sets.

Suppose we have m disjoint ground sets X_1, X_2, \dots, X_m and on each X_i , we have a sequence of concept classes $\mathcal{C}_i^0, \mathcal{C}_i^1, \mathcal{C}_i^2, \dots \subseteq \{0, 1\}^X$. Suppose $\mathcal{C}_i^0 \neq \emptyset$ for all i . We use X to denote the disjoint union $\bigcup_{i=1}^m X_i$. For any $d \geq 0$, we define a concept class $\mathcal{P}(d)$ on X to be the class of functions $f \in \{0, 1\}^X$ satisfying that $\exists k_1, k_2, \dots, k_m \in \mathbb{N}$ s.t.

1. $\sum_{i=1}^m k_i \leq d$;
2. $\forall 1 \leq i \leq m, f|_{X_i} \in \mathcal{C}_i^{k_i}$.

We call \mathcal{P} a *composition of m additive properties*. Note that $\mathcal{P}(0) = \{f \in \{0, 1\}^X : \forall 1 \leq i \leq m, f|_{X_i} \in \mathcal{C}_i^0\}$, matching the definition of a disjoint union of properties in [Balcan et al., 2012]. Also note that $\mathcal{P}(0) \neq \emptyset$ because of the assumption that $\mathcal{C}_i^0 \neq \emptyset$ for all i .

For a given $t \geq 0$, we define a composition \mathcal{P}^t in the same way as \mathcal{P} except that we further require every k_i to be at most t , or, \mathcal{P}^t is a composition of m additive properties *truncated by t* .

For any distribution \mathcal{D} over X , we use $\mathcal{P}_{\mathcal{D}}(d, \alpha)$ to denote functions that are α -close to $\mathcal{P}(d)$ with respect to \mathcal{D} , i.e. $\mathcal{P}_{\mathcal{D}}(d, \alpha) = \{f \in \{0, 1\}^X : \exists g \in \mathcal{P}(d), \text{dist}_{\mathcal{D}}(f, g) \leq \alpha\}$. Similarly, we define $\mathcal{P}_{\mathcal{D}}^t(d, \alpha) = \{f \in \{0, 1\}^X : \exists g \in \mathcal{P}^t(d), \text{dist}_{\mathcal{D}}(f, g) \leq \alpha\}$. We say \mathcal{D} is *semi-uniform* if $\forall 1 \leq i \leq m, \Pr_{x \sim \mathcal{D}}[x \in X_i] = \frac{1}{m}$.

Lemma 1. *Suppose the distribution \mathcal{D} is semi-uniform. We have $\mathcal{P}_{\mathcal{D}}^t(d, \alpha) \subseteq \mathcal{P}_{\mathcal{D}}(d, \alpha) \subseteq \mathcal{P}_{\mathcal{D}}^t(d, \alpha + \frac{d}{tm})$.*

Proof. $\mathcal{P}_{\mathcal{D}}^t(d, \alpha) \subseteq \mathcal{P}_{\mathcal{D}}(d, \alpha)$ is obvious. To see $\mathcal{P}_{\mathcal{D}}(d, \alpha) \subseteq \mathcal{P}_{\mathcal{D}}^t(d, \alpha + \frac{d}{tm})$, we note that for any $g \in \mathcal{P}(d)$, for each i such that $k_i > t$, substituting a function in \mathcal{C}_i^0 for $g|_{X_i}$ causes at most a $\frac{1}{m}$ increase in the distance from $f \in \mathcal{P}_{\mathcal{D}}(d, \alpha)$ to g . An easy observation that $|\{i : k_i > t\}| \leq \frac{d}{t}$ given $\sum_{i=1}^m k_i \leq d$ completes the proof. \square

An (ϵ, μ) -bi-criteria distance approximation algorithm $\text{Comp}_{\mathcal{D}}(f, (\epsilon, \mu), d)$ for composition \mathcal{P} of additive properties, is an algorithm that takes f, ϵ, μ and d as input and outputs $\hat{\alpha}$ such that $\forall \alpha$

1. $\forall f \in \mathcal{P}_{\mathcal{D}}(d, \alpha)$, it holds with probability at least $\frac{2}{3}$ that $\hat{\alpha} \leq \alpha + \epsilon$;
2. $\forall f \notin \mathcal{P}_{\mathcal{D}}((1 + \mu)d, \alpha)$, it holds with probability at least $\frac{2}{3}$ that $\hat{\alpha} > \alpha - \epsilon$.

Suppose we have a sequence of indices $1 \leq i_1 < i_2 < \dots < i_l \leq m$ denoted by \mathbf{i} for short. Let $\mathcal{D}_{\mathbf{i}}$ denote the conditional distribution of \mathcal{D} on $\bigcup_{j=1}^l X_{i_j}$. A (d, l, t, ϵ) distance approximation oracle is an algorithm taking a length- l sequence \mathbf{i} of indices and $f \in \{0, 1\}^X$ as input, and performing $\text{Comp}_{\mathcal{D}_{\mathbf{i}}}(f_{\text{active}}, (\epsilon, 0), d)$ on composition \mathcal{P}^t . In other words, this algorithm performs distance approximation on any given l -sub-union (l is typically small) of the m ground sets. For convenience of use, we require the success probability of the oracle to be at least $\frac{11}{12}$.

4 The Composition Lemma

Lemma 2 (Composition Lemma). *Suppose \mathcal{P} is the composition of m additive properties defined in Section 3.3. Let \mathcal{D} be a semi-uniform distribution. For parameters $\lambda > 0, \alpha \in [0, 1]$ and $\mu, \epsilon \in (0, 1)$ taken as input, there exists $l = O(\frac{1}{\epsilon\mu^2} + \frac{1}{\epsilon^2})$ such that we have an algorithm that performs $\text{Comp}_{\mathcal{D}}(f_{\text{active}}, (\epsilon, \mu), \lambda m)$ by calling once a $((1 + \frac{\mu}{2})\lambda l, l, \frac{4\lambda}{\epsilon}, \frac{\epsilon}{2})$ distance approximation oracle. Suppose the query complexity and the unlabeled sample complexity of the oracle are q and N , respectively. Then the query complexity and the unlabeled sample complexity of the algorithm are q and $O(\frac{Nm}{l})$, respectively.*

Proof. The algorithm first picks indices $1 \leq i_1 < i_2 < \dots < i_l \leq m$ uniformly at random for $l = O(\frac{1}{\epsilon\mu^2} + \frac{1}{\epsilon^2})$. Then the algorithm asks for $O(\frac{Nm}{l})$ unlabeled examples to make sure with probability at least $\frac{11}{12}$, there are at least N examples lying in $\bigcup_{j=1}^l X_{i_j}$. These examples can be treated as drawn independently at random according to $\mathcal{D}_{\mathbf{i}}$, where $\mathbf{i} = (i_1, i_2, \dots, i_l)$. Finally, the algorithm calls the oracle to approximate the distance from f to $\mathcal{P}^t((1 + \frac{\mu}{2})\lambda l)$ truncated by $t = \frac{4\lambda}{\epsilon}$ on distribution $\mathcal{D}_{\mathbf{i}}$ up to an additive error $\frac{\epsilon}{2}$ using these unlabeled examples and outputs what the oracle outputs.

The correctness of the algorithm follows from the following two lemmas and the Union Bound. \square

Lemma 3. *Suppose $t = \frac{4\lambda}{\epsilon}$. If $f \in \mathcal{P}_{\mathcal{D}}(\lambda m, \alpha)$, then choosing $l = O(\frac{1}{\epsilon\mu^2} + \frac{1}{\epsilon^2})$ is enough to make sure that with probability at least $\frac{5}{6}$, $f \in \mathcal{P}_{\mathcal{D}_{\mathbf{i}}}^t((1 + \frac{\mu}{2})\lambda l, \alpha + \frac{\epsilon}{2})$.*

Lemma 4. *Suppose $t = \frac{4\lambda}{\epsilon}$. If $f \notin \mathcal{P}_{\mathcal{D}}((1 + \mu)\lambda m, \alpha)$, then choosing $l = O(\frac{1}{\epsilon\mu^2} + \frac{1}{\epsilon^2})$ is enough to make sure that with probability at least $\frac{5}{6}$, $f \notin \mathcal{P}_{\mathcal{D}_{\mathbf{i}}}^t((1 + \frac{\mu}{2})\lambda l, \alpha - \frac{\epsilon}{2})$.*

Proof of Lemma 3. By the choice of truncation $t = \frac{4\lambda}{\epsilon}$, according to Lemma 1, we know $f \in \mathcal{P}_{\mathcal{D}}^t(\lambda m, \alpha + \frac{\epsilon}{4})$. Suppose $\text{dist}_{\mathcal{D}}(f, g) \leq \alpha + \frac{\epsilon}{4}$ for some $g \in \mathcal{P}^t(\lambda m)$. According to the Multiplicative Chernoff Bound for sampling without replacement, choosing $l = O(\frac{1}{\epsilon\mu^2})$ is enough to make sure

that with probability at least $\frac{11}{12}$, $\exists g'$ s.t. $g' \in \mathcal{P}^t((1 + \frac{\mu}{2})\lambda l)$ and $\text{dist}_{\mathcal{D}_i}(g, g') = 0$.² According to the Chernoff Bound for sampling without replacement, choosing $l = O(\frac{1}{\epsilon^2})$ is enough to make sure that with probability at least $\frac{11}{12}$, $\text{dist}_{\mathcal{D}_i}(f, g) \leq \alpha + \frac{\epsilon}{2}$. By the Union Bound, these two events happen at the same time with probability at least $\frac{5}{6}$, and in this case, $f \in \mathcal{P}_{\mathcal{D}_i}^t((1 + \frac{\mu}{2})\lambda l, \alpha + \frac{\epsilon}{2})$. \square

Proof of Lemma 4. According to Lemma 1, we know $f \notin \mathcal{P}_{\mathcal{D}}^t((1 + \mu)\lambda m, \alpha)$. Therefore, by definition, there exists $g \in \mathcal{P}^t((1 + \mu)\lambda m)$ with the following two properties:³

1. $\text{dist}_{\mathcal{D}}(f, g) > \alpha$;
2. $\forall g' \in \mathcal{P}^t((1 + \mu)\lambda m), \text{dist}_{\mathcal{D}}(f, g') > \text{dist}_{\mathcal{D}}(f, g) - \frac{\epsilon}{4} \cdot \frac{l}{m}$.

Suppose $g|_{X_i} \in \mathcal{C}_i^{k_i}$ for $k_i \leq t = \frac{4\lambda}{\epsilon}$ satisfying $k := \sum_{i=1}^m k_i \leq (1 + \mu)\lambda m$. We enlarge k_i to $k'_i \in [k_i, t]$ to make sure that $k' := \sum_{i=1}^m k'_i = (1 + \mu)\lambda m$.⁴ According to the Multiplicative Chernoff Bound for sampling without replacement, choosing $l = O(\frac{1}{\epsilon\mu^2})$ is enough to make sure that with probability at least $\frac{11}{12}$, $\sum_{j=1}^l k'_{i_j} \geq (1 + \frac{\mu}{2})\lambda l$.

Now suppose it's the case that $\sum_{j=1}^l k'_{i_j} \geq (1 + \frac{\mu}{2})\lambda l$. Then, according to the second property of g , we know

$$\forall g' \in \mathcal{P}^t((1 + \frac{\mu}{2})\lambda l), \text{dist}_{\mathcal{D}_i}(f, g') > \text{dist}_{\mathcal{D}_i}(f, g) - \frac{\epsilon}{4}.$$

Otherwise, we can swap g' for g on $\bigcup_{j=1}^l X_{i_j}$ causing a violation of the second property of g .

Finally, according to the Chernoff Bound for sampling without replacement, choosing $l = O(\frac{1}{\epsilon^2})$ is enough to make sure that with probability at least $\frac{11}{12}$, $\text{dist}_{\mathcal{D}_i}(f, g) > \alpha - \frac{\epsilon}{4}$. Therefore, by the Union Bound, with probability at least $\frac{5}{6}$,

$$\forall g' \in \mathcal{P}^t((1 + \frac{\mu}{2})\lambda l), \text{dist}_{\mathcal{D}_i}(f, g') > \text{dist}_{\mathcal{D}_i}(f, g) - \frac{\epsilon}{4} > \alpha - \frac{\epsilon}{2},$$

a completion of the proof. \square

5 Distance Approximation for Unions of d Intervals

In this section, we consider the concept class $\mathcal{I}(d)$ of binary functions f defined on \mathbb{R} such that $f^{-1}(1)$ is a union of at most d intervals. We use $\text{Int}_{\mathcal{D}}(f, \epsilon, d)$ to denote the distance approximation task $\text{DA}_{\mathcal{D}}(f, \epsilon)$ when the underlying hypothesis class is the class $\mathcal{I}(d)$ of unions of d intervals. We show (Theorem 5) that there is an $\text{Int}_{\mathcal{D}}(f, \epsilon, d)$ algorithm in the active model with query complexity independent of d even when the data distribution is unknown to the algorithm.

² g' is chosen such that $g'|_{X_i} \in \mathcal{C}_i^0$ for all $i \notin \{i_1, i_2, \dots, i_l\}$ and $g'|_{X_i} = g|_{X_i}$ for all $i \in \{i_1, i_2, \dots, i_l\}$. The fact that the k_i 's of g are bounded between 0 and $t = \frac{4\lambda}{\epsilon}$ allows us to use the Multiplicative Chernoff Bound.

³E.g., choose g to be the closest or approximately-closest function in the class to f . Note that $\mathcal{P}^t((1 + \mu)\lambda m)$ can't be empty, because $\mathcal{P}^t((1 + \mu)\lambda m) \supseteq \mathcal{P}^t(0) = \mathcal{P}(0) \neq \emptyset$.

⁴ k'_i doesn't have to be an integer. Also note that $mt = \frac{4\lambda}{\epsilon} \cdot m > 4\lambda m > (1 + \mu)\lambda m$.

Theorem 5 (main theorem). *Suppose $d > 0$ and $\epsilon \in (0, \frac{1}{2})$ are given as input. Let \mathcal{D} be an unknown distribution on \mathbb{R} . Suppose we have active access to an input function $f : \mathbb{R} \rightarrow \{0, 1\}$ with respect to \mathcal{D} . There is an $\text{Int}_{\mathcal{D}}(f_{\text{active}}, \epsilon, d)$ algorithm using $O(\frac{1}{\epsilon^6} \log \frac{1}{\epsilon})$ queries on $O(\frac{d}{\epsilon^2} \log \frac{1}{\epsilon})$ unlabeled examples.*

Before proving Theorem 5, we first introduce two helper results, Theorems 6 and 7.

5.1 Relationship between Query Testing and Active Testing

As Balcan et al. [2012] have pointed out, in the task of testing unions of d intervals in the query testing framework, any known distribution can be reduced to uniform distribution on $[0, 1]$ by its CDF. Our following theorem shows that once we can deal with arbitrary distributions for query testing, we can automatically deal with unknown distributions for active testing, improving a previous upper bound on unlabeled sample complexity in [Balcan et al., 2012].

Theorem 6. *Let \mathcal{C} be a concept class on ground set X with VC-dimension d . Suppose $\epsilon \in (0, \frac{1}{2})$. Suppose there is a $\text{PT}_{\mathcal{D}}(f_{\text{query}}, \frac{\epsilon}{2})$ algorithm \mathcal{A} using at most q queries on arbitrarily given distribution \mathcal{D} with finite support. Suppose all the queries algorithm \mathcal{A} makes lie in the support of \mathcal{D} .⁵ Then, there is a $\text{PT}_{\mathcal{D}}(f_{\text{active}}, \epsilon)$ algorithm \mathcal{B} using at most $O(q)$ queries on $O(\frac{d}{\epsilon} \log \frac{1}{\epsilon})$ unlabeled examples, even when distribution \mathcal{D} is unknown to algorithm \mathcal{B} .*

Proof of Theorem 6. Algorithm \mathcal{B} first draws $N = O(\frac{d}{\epsilon} \log \frac{1}{\epsilon})$ unlabeled examples: x_1, x_2, \dots, x_N . We use \mathcal{S} to denote the uniform distribution over these unlabeled examples. By VC Theory, we know if $\forall g \in \mathcal{C}, \text{dist}_{\mathcal{D}}(f, g) > \epsilon$, then with probability at least $\frac{5}{6}$, $\forall g \in \mathcal{C}, \text{dist}_{\mathcal{S}}(f, g) > \frac{\epsilon}{2}$. So algorithm \mathcal{B} only needs to call algorithm \mathcal{A} to distinguish $f \in \mathcal{C}$ and $\forall g \in \mathcal{C}, \text{dist}_{\mathcal{S}}(f, g) > \frac{\epsilon}{2}$ with probability at least $\frac{5}{6}$. By the Union Bound, algorithm \mathcal{B} succeeds with probability at least $\frac{2}{3}$. \square

Therefore, since Balcan et al. [2012] have an algorithm in the query testing framework that can distinguish $f \in \mathcal{I}(d)$ and $f \notin \mathcal{I}(d, \epsilon)$ using $O(\frac{1}{\epsilon^4})$ queries, there is an algorithm in the active testing framework that can distinguish $f \in \mathcal{I}(d)$ and $f \notin \mathcal{I}_{\mathcal{D}}(d, \epsilon)$ using $O(\frac{1}{\epsilon^4})$ queries on $O(\frac{d}{\epsilon} \log \frac{1}{\epsilon})$ unlabeled examples, even when the distribution \mathcal{D} is unknown, according to Theorem 6. Here, the unlabeled sample complexity is $O(\frac{d}{\epsilon} \log \frac{1}{\epsilon})$, an improvement from $O(\frac{d^2}{\epsilon^6})$ in their original paper. The theorem can be easily generalized to distance approximation.

Theorem 7. *Let \mathcal{C} be a concept class on ground set X with VC-dimension d . Suppose $\epsilon \in (0, \frac{1}{2})$. Suppose there is a $\text{DA}_{\mathcal{D}}(f_{\text{query}}, \frac{\epsilon}{2}, \mathcal{D})$ algorithm \mathcal{A} using at most q queries on arbitrarily given distribution \mathcal{D} with finite support. Then, there is a $\text{DA}_{\mathcal{D}}(f_{\text{active}}, \epsilon)$ algorithm \mathcal{B} using at most $O(q)$ queries on $O(\frac{d}{\epsilon^2} \log \frac{1}{\epsilon})$ unlabeled examples, even when distribution \mathcal{D} is unknown to algorithm \mathcal{B} .*

5.2 Proof of Theorem 5

By Theorem 7, we only need to show a distance approximation algorithm with an additive error bounded below ϵ using $O(\frac{1}{\epsilon^6} \log \frac{1}{\epsilon})$ queries on arbitrary *known* distribution \mathcal{D} . As pointed out by

⁵For TT and DA, we can assume without loss of generality that the algorithm never queries examples outside the support of the distribution, but this is not without loss of generality for PT, because $f \in \mathcal{C}$ is stronger than $\exists g \in \mathcal{C}, \text{dist}_{\mathcal{D}}(f, g) = 0$.

Balcan et al. [2012], any known distribution can be reduced to uniform distribution on $[0, 1]$ by its CDF. Therefore, we only consider \mathcal{D} as the uniform distribution on $[0, 1]$ in this section and we omit it for simplicity.

We first reveal a basic property of unions of d intervals.

Lemma 8. $\forall \epsilon \in (0, \frac{1}{2}), \forall \alpha \in [0, 1], \forall d > \frac{2}{\epsilon}, \mathcal{I}((1 + \frac{\epsilon}{2})d, \alpha - \epsilon) \subseteq \mathcal{I}(d, \alpha)$.

Proof. $\forall f \in \mathcal{I}((1 + \frac{\epsilon}{2})d, \alpha - \epsilon), \exists g \in \mathcal{I}((1 + \frac{\epsilon}{2})d)$ s.t. $\text{dist}(f, g) \leq \alpha - \epsilon$. Assume g uses $k \leq (1 + \frac{\epsilon}{2})d$ intervals. Without loss of generality, we can assume that $k \geq d$. We remove $\lceil \frac{\epsilon k}{2} \rceil$ shortest intervals from the k intervals. The number of remaining intervals is at most $(1 - \frac{\epsilon}{2})k \leq (1 - \frac{\epsilon}{2})(1 + \frac{\epsilon}{2})d \leq d$. The distance increase is upper bounded by $(\frac{\epsilon k}{2} + 1) \cdot \frac{1}{k} \leq \frac{\epsilon}{2} + \frac{1}{d} \leq \epsilon$. \square

Now we come back and prove Theorem 5. If $d \leq \frac{8}{\epsilon}$, we can simply do agnostic learning using $O(\frac{d}{\epsilon^2} \log \frac{1}{\epsilon}) = O(\frac{1}{\epsilon^3} \log \frac{1}{\epsilon})$ queries and unlabeled examples. So in the rest of the proof, we assume $d > \frac{8}{\epsilon}$. We pick the largest positive integer m satisfying $m \leq \frac{cd}{8}$ and we define $\lambda = \frac{d}{m} = O(\frac{1}{\epsilon})$.

Since the data distribution is assumed uniform on $[0, 1]$, we can assume without loss of generality that our ground set X is $[0, 1]$ and $f \in \{0, 1\}^X$. We evenly cut X into m pieces: $X_1 = [0, \frac{1}{m}], X_2 = (\frac{1}{m}, \frac{2}{m}], X_3 = (\frac{2}{m}, \frac{3}{m}], \dots, X_m = (\frac{m-1}{m}, 1]$. $\forall 1 \leq i \leq m, \forall k \in \mathbb{N}$, we define \mathcal{C}_i^k to be the class of binary functions f on X_i such that $f^{-1}(1)$ is a union of at most k intervals. Note that $\mathcal{C}_i^0 \neq \emptyset$. Therefore, we can define \mathcal{P} , the composition of m additive properties as in Section 3.3.

Note that for any $d' > 0$ and any truncation $t > 0$, the concept class $\mathcal{P}^t(d')$ has VC-dimension at most $2d'$. Therefore, simply by VC Theory for agnostic learning, for any $\mu, \epsilon' \in (0, \frac{1}{2}), l = O(\frac{1}{\epsilon' \mu^2} + \frac{1}{\epsilon'^2})$, we have a $((1 + \frac{\mu}{2})(1 + \frac{\epsilon}{8})\lambda l, l, \frac{4(1 + \frac{\epsilon}{8})\lambda}{\epsilon'}, \frac{\epsilon'}{2})$ distance approximation oracle using $O(\frac{(1 + \frac{\mu}{2})(1 + \frac{\epsilon}{8})\lambda l}{(\frac{\epsilon'}{2})^2} \log \frac{1}{\frac{\epsilon'}{2}}) = O(\frac{l}{\epsilon'^2 \epsilon} \log \frac{1}{\epsilon'}) = O((\frac{1}{\epsilon'^3 \epsilon \mu^2} + \frac{1}{\epsilon'^4 \epsilon}) \log \frac{1}{\epsilon'})$ queries and unlabeled examples. By the Composition Lemma (Lemma 2), we have an algorithm that outputs $\hat{\alpha}$ such that $\forall \alpha$,

1. $\forall f \in \mathcal{P}((1 + \frac{\epsilon}{8})\lambda m, \alpha)$, it holds with probability at least $\frac{2}{3}$ that $\hat{\alpha} \leq \alpha + \epsilon'$;
2. $\forall f \notin \mathcal{P}((1 + \mu)(1 + \frac{\epsilon}{8})\lambda m, \alpha)$, it holds with probability at least $\frac{2}{3}$ that $\hat{\alpha} > \alpha - \epsilon'$.

Choose $1 + \mu = \frac{1 + \frac{\epsilon}{4}}{1 + \frac{\epsilon}{8}}$ and note that $\lambda m = d, \mathcal{I}(d, \alpha) \subseteq \mathcal{P}(d + m, \alpha) \subseteq \mathcal{P}((1 + \frac{\epsilon}{8})d, \alpha)$ and $\mathcal{P}((1 + \frac{\epsilon}{4})d, \alpha) \subseteq \mathcal{I}((1 + \frac{\epsilon}{4})d, \alpha)$, we have $\forall \alpha$,

1. $\forall f \in \mathcal{I}(d, \alpha)$, it holds with probability at least $\frac{2}{3}$ that $\hat{\alpha} \leq \alpha + \epsilon'$;
2. $\forall f \notin \mathcal{I}((1 + \frac{\epsilon}{4})d, \alpha)$, it holds with probability at least $\frac{2}{3}$ that $\hat{\alpha} > \alpha - \epsilon'$.

This is an $(\epsilon', \frac{\epsilon}{4})$ -bi-criteria tester for unions of d intervals. According to the Composition Lemma (Lemma 2), the query complexity and the unlabeled sample complexity of the algorithm are $O((\frac{1}{\epsilon'^3 \epsilon \mu^2} + \frac{1}{\epsilon'^4 \epsilon}) \log \frac{1}{\epsilon'}) = O((\frac{1}{\epsilon'^3 \epsilon^3} + \frac{1}{\epsilon'^4 \epsilon}) \log \frac{1}{\epsilon'})$ and $O((\frac{l}{\epsilon'^2 \epsilon} \log \frac{1}{\epsilon'}) \cdot \frac{m}{l}) = O(\frac{d}{\epsilon'^2} \log \frac{1}{\epsilon'})$.

Now we define $\epsilon' = \frac{\epsilon}{2}$ and by rewriting the second statement in an equivalent way, we get $\forall \alpha$,

1. $\forall f \in \mathcal{I}(d, \alpha)$, it holds with probability at least $\frac{2}{3}$ that $\hat{\alpha} \leq \alpha + \frac{\epsilon}{2} < \alpha + \epsilon$;
2. $\forall f \notin \mathcal{I}((1 + \frac{\epsilon}{4})d, \alpha - \frac{\epsilon}{2})$, it holds with probability at least $\frac{2}{3}$ that $\hat{\alpha} > (\alpha - \frac{\epsilon}{2}) - \frac{\epsilon}{2} = \alpha - \epsilon$.

Finally, by Lemma 8, we have $\mathcal{I}((1 + \frac{\epsilon}{4})d, \alpha - \frac{\epsilon}{2}) \subseteq \mathcal{I}(d, \alpha)$, which completes the proof.

6 Estimating the Performance of k -Nearest Neighbor Algorithms

In this section, we develop testers for estimating the performance of k -Nearest Neighbor (k -NN) algorithms [Fix and Hodges Jr, 1951, Fix and Hodges, 1989, Cover and Hart, 1967].

Let \mathcal{D} be a distribution on a ground set X . Suppose that every point $x \in X$ has a (true) label $f(x) \in \{0, 1\}$. In addition, we have a distance metric $d : X \times X \rightarrow \mathbb{R}_{\geq 0}$ that is symmetric, nonnegative and satisfies the triangle inequality. The k -Nearest Neighbor algorithm with soft predictions (k -NN^{soft}) is given a pool S of unlabeled examples, sampled iid from \mathcal{D} , and for any input $x \in X$, finds its k nearest examples $x_1, x_2, \dots, x_k \in S$ with respect to the distance metric d and outputs $\hat{f}(x) = \frac{1}{k} \sum_{i=1}^k f(x_i)$ as an approximation of $f(x)$. In this paper, we assume the k nearest examples are calculated by an oracle M , i.e., when given x and S , M calculates the k nearest examples to x in S . There may be ties when distances to x are compared and we assume M breaks ties according to some (probably random) mechanism.

The k -Nearest Neighbor algorithm with hard predictions (k -NN^{hard}) does the same thing as k -NN^{soft}, except that $\hat{f}(x)$ is chosen as the majority vote $I[\frac{1}{k} \sum_{i=1}^k f(x_i) > 0.5]$.⁶

For both algorithms, we use $\text{err}_1(x) = |\hat{f}(x) - f(x)|$ to denote the L^1 error on point $x \in X$. For soft prediction, we will penalize the algorithm by taking the p th power of the L^1 error for positive integer p .

6.1 Estimating the Performance of k -NN^{soft}

Given a loss function $\text{loss}(\cdot)$, we can measure the performance of k -NN^{soft} by its expected loss $\mathbb{E}_x[\text{loss}(\text{err}_1(x))]$. The expectation is over the random draw of x with respect to distribution \mathcal{D} and the randomness of the oracle M when ties occur. In this paper, we focus on the p th-power loss $\mathbb{E}_x[(\text{err}_1(x))^p]$ for positive integer p . Let $\mathcal{T}_{\mathcal{D}}^{\text{soft}}(f, \epsilon, S, k)$ denote the testing task of approximating the expected loss of a k -NN^{soft} algorithm up to an additive error ϵ with success probability at least $\frac{2}{3}$. We consider the testing task in the active model, in which the tester is only allowed to query labels of examples in an unlabeled pool sampled iid from \mathcal{D} . In addition to the given unlabeled pool S from which k -NN^{soft} would learn, we allow the $\mathcal{T}_{\mathcal{D}}^{\text{soft}}(f_{\text{active}}, \epsilon, S, k)$ tester to sample fresh unlabeled examples and query their labels. We assume the tester has access to the oracle M .

Theorem 9. *Suppose we consider the p th-power loss for $p \in \mathbb{N}^*$. There is a tester $\mathcal{T}_{\mathcal{D}}^{\text{soft}}(f_{\text{active}}, \epsilon, S, k)$ using $O(\frac{p}{\epsilon^2})$ queries on $N + O(\frac{1}{\epsilon^2})$ unlabeled examples when the unlabeled pool S has size N . The underlying distribution \mathcal{D} is assumed unknown to the tester. Moreover, the tester has success probability at least $\frac{2}{3}$ for any unlabeled pool S .*

Before proving the theorem, we first show a simple tester that works for any loss function $\text{loss}(\cdot)$ bounded in $[0, 1]$ with L -Lipschitz property⁷ using $O(\frac{L^2}{\epsilon^4} \cdot \log \frac{1}{\epsilon})$ queries on $N + O(\frac{1}{\epsilon^2})$ unlabeled examples. The tester runs for $O(\frac{1}{\epsilon^2})$ iterations and in each i th iteration, the tester samples a fresh unlabeled example x and then queries the labels of $w = O(\frac{L^2}{\epsilon^2} \log \frac{1}{\epsilon})$ examples x_1, x_2, \dots, x_w

⁶ $I[\cdot]$ is the indicator function of a statement, which takes value 1 if the statement is true and value 0 if the statement is false.

⁷We say $\text{loss}(\cdot)$ has L -Lipschitz property if $\forall x_1, x_2 \in [0, 1], |\text{loss}(x_1) - \text{loss}(x_2)| \leq L|x_1 - x_2|$.

sampled independently at random uniformly from the k nearest neighbors of x in S . The estimator for this iteration is $E_i = \text{loss}(|\frac{1}{w} \sum_{j=1}^w f(x_j) - f(x)|)$. The final output of the tester is the average of all E_i 's for all iterations i .

We prove Theorem 9 by slightly modifying the above tester's each iteration for p th-power loss. Instead of looking at the labels of w examples, we only need to look at p labels of x_1, x_2, \dots, x_p , still sampled independently at random uniformly from the k nearest neighbors of x in S . In this case, E_i is defined to be $\prod_{j=1}^p |f(x_j) - f(x)|$. The final output of the tester is still the average of E_i 's.

Proof of Theorem 9. We use e_j to denote $|f(x_j) - f(x)|$. To show the above tester works, we first look at the value we want to estimate: $\mathbb{E}_x[(\text{err}_1(x))^p] = \mathbb{E}_x[(\mathbb{E}_{x_1}[e_1])^p]$, where x_1 is sampled uniformly from the k nearest neighbors of x in T . Note that x_1, x_2, \dots, x_p are iid, so we know $\mathbb{E}_x[(\mathbb{E}_{x_1}[e_1])^p] = \mathbb{E}_x[\mathbb{E}_{x_1, x_2, \dots, x_p}[e_1 e_2 \dots e_p]] = \mathbb{E}_{x, x_1, x_2, \dots, x_p}[\prod_{j=1}^p |f(x_j) - f(x)|]$. The Chernoff Bound thus completes the proof. \square

Theorem 9 also holds naturally for Weighted Nearest Neighbor algorithms [Royall, 1966] with soft predictions, in which $\hat{f}(x)$ is a weighted average of $f(x')$ for all $x' \in S$ where the weights depend on the distances $d(x', x)$, simply by sampling x_1, x_2, \dots, x_p iid from S according to the weights.

In Theorem 13 (Section 6.5), we will show a matching lower bound for the $O(\frac{p}{\epsilon^2})$ query complexity.

6.2 Finding an Approximately-Best Choice of k

Based on the result in Section 6.1, we are able to construct an algorithm that approximately optimizes the choice of k in the k -NN^{soft} algorithm.

Suppose we have active access to the true label f with respect to distribution \mathcal{D} over ground set X with distance metric d . Suppose the size of the unlabeled pool S is fixed to be N . We use loss_k to denote the expected loss of the k -NN^{soft} algorithm and consider how the k -NN^{soft} algorithm performs with different values of k . We assume the oracle M uses the same tie-breaking mechanism for different values of k . Specifically, given x and S , M arranges the examples in S as x_1, x_2, \dots, x_N so that $\forall i, d(x_i, x) \leq d(x_{i+1}, x)$. x_1, x_2, \dots, x_k are taken by k -NN^{soft} as the k nearest neighbors of x for any $k \in \{1, 2, \dots, N\}$.

Lemma 10. *Suppose $k_1 \leq k_2$ and the loss function $\text{loss}(\cdot)$ is L -Lipschitz. Then, $|\text{loss}_{k_1} - \text{loss}_{k_2}| \leq L \cdot (1 - \frac{k_1}{k_2})$.*

Proof. When the test point x is chosen, we use x_1, x_2, \dots, x_{k_2} to denote the closest k_2 points to x in S , arranged in non-decreasing order of their distances to x . Each x_i might be random because ties might be broken randomly. We use e_i to denote $|f(x_i) - f(x)|$. Note that we have

$\text{loss}_{k_1} = \mathbb{E}_{x, x_1, x_2, \dots, x_{k_1}} [\text{loss}(\frac{1}{k_1} \sum_{i=1}^{k_1} e_i)]$ and $\text{loss}_{k_2} = \mathbb{E}_{x, x_1, x_2, \dots, x_{k_2}} [\text{loss}(\frac{1}{k_2} \sum_{i=1}^{k_2} e_i)]$. Therefore,

$$\begin{aligned}
& |\text{loss}_{k_1} - \text{loss}_{k_2}| \\
& \leq \mathbb{E}_{x, x_1, x_2, \dots, x_{k_2}} [|\text{loss}(\frac{1}{k_1} \sum_{i=1}^{k_1} e_i) - \text{loss}(\frac{1}{k_2} \sum_{i=1}^{k_2} e_i)|] \\
& \leq L \cdot \mathbb{E}_{x, x_1, x_2, \dots, x_{k_2}} [|\frac{1}{k_1} \sum_{i=1}^{k_1} e_i - \frac{1}{k_2} \sum_{i=1}^{k_2} e_i|] \\
& = L \cdot \mathbb{E}_{x, x_1, x_2, \dots, x_{k_2}} [|\frac{1}{k_1} \sum_{i=1}^{k_1} e_i - \frac{1}{k_2} \sum_{i=k_1+1}^{k_2} e_i|] \tag{1} \\
& \leq L \cdot \mathbb{E}_{x, x_1, x_2, \dots, x_{k_2}} [\max\{(\frac{1}{k_1} - \frac{1}{k_2}) \sum_{i=1}^{k_1} e_i, \frac{1}{k_2} \sum_{i=k_1+1}^{k_2} e_i\}] \\
& \leq L \cdot \max\{(\frac{1}{k_1} - \frac{1}{k_2}) \cdot k_1, \frac{1}{k_2} \cdot (k_2 - k_1)\} \\
& = L \cdot (1 - \frac{k_1}{k_2})
\end{aligned}$$

□

We say k is ϵ -approximately-best, if $\forall k' \in \{1, 2, \dots, N\}, \text{loss}_{k'} \geq \text{loss}_k - \epsilon$. The following theorem states that we can find an ϵ -approximately-best k using a small number of queries.

Theorem 11. *Suppose k -NN^{soft} algorithms with an unlabeled pool S of size N are measured by p th-power loss for $p \in \mathbb{N}^*$. Suppose $\epsilon \in (0, \frac{1}{2})$. There is an algorithm that finds an ϵ -approximately-best k w.p. at least $\frac{2}{3}$ using $O(\frac{p^2 \log N}{\epsilon^3} (\log \log N + \log p + \log \frac{1}{\epsilon}))$ queries on $N + O(\frac{p \log N}{\epsilon^3} (\log \log N + \log p + \log \frac{1}{\epsilon}))$ unlabeled examples.*

Proof. If we apply Lemma 10 to p th-power loss, which is p -Lipschitz, we know for any $1 \leq \frac{k_2}{k_1} \leq \frac{p}{p-\epsilon}$, it holds that $|\text{loss}_{k_1} - \text{loss}_{k_2}| \leq \epsilon$. If we define $t = \lfloor \log_{\frac{p}{p-\epsilon}} N \rfloor$, $k_{2i} = \lfloor (\frac{p}{p-\epsilon})^i \rfloor$, $k_{2i+1} = \lceil (\frac{p}{p-\epsilon})^i \rceil$ for $i = 0, 1, 2, \dots, t$, then we know $\exists 0 \leq i \leq 2t + 1$ such that k_i is $\frac{\epsilon}{3}$ -approximately-best. By Theorem 9, we can approximate loss_{k_i} for every $0 \leq i \leq 2t + 1$ up to an additive error $\frac{\epsilon}{3}$ using $O(\frac{pt \log t}{\epsilon^2})$ queries on $N + O(\frac{t \log t}{\epsilon^2})$ unlabeled examples.⁸ The k_i yielding the smallest approximation of loss_{k_i} is ϵ -approximately-best. Note that $t = O(\frac{p \log N}{\epsilon})$, so the query complexity is $O(\frac{p^2 \log N}{\epsilon^3} (\log \log N + \log p + \log \frac{1}{\epsilon}))$ and the unlabeled sample complexity is $N + O(\frac{p \log N}{\epsilon^3} (\log \log N + \log p + \log \frac{1}{\epsilon}))$. □

6.3 Estimating the Performance of k -NN^{hard}

The performance of k -NN^{hard} is naturally measured by its error rate $\mathbb{E}_x[\text{err}_1(x)]$ and we use $\mathcal{T}_{\mathcal{D}}^{\text{hard}}(f, \epsilon, S, k)$ to denote the corresponding testing task of estimating the error rate of k -NN^{hard} up to an additive error ϵ with success probability at least $\frac{2}{3}$.

⁸Repeat the tester $O(\log t)$ times and take the median to boost its success probability to $1 - O(\frac{1}{t})$.

A trivial tester achieving this goal using $O(\frac{k}{\epsilon^2})$ queries on $N + O(\frac{1}{\epsilon^2})$ unlabeled examples is to use the empirical mean of $\text{err}_1(x)$ as an estimator of $\mathbb{E}_x[\text{err}_1(x)]$. This tester is not satisfactory because its query complexity grows with respect to k . In Section 6.5, we will show (Theorem 15) that this linear growth with respect to k can't be eliminated. Also, we will show (Theorem 14) that the $O(\frac{k}{\epsilon^2})$ query complexity is optimal if we assume a natural algorithm for *approximating the fraction of good arms* (AGA) in the stochastic multi-armed bandit setting has the optimal query complexity. Before we show our lower bound results, we first define the problems of *counting and approximating the number of good arms*.

6.4 Counting and Approximating the Number of Good Arms

To show query complexity lower bound results for estimating the performance of k -Nearest Neighbor algorithms, we show reductions from two related problems in the stochastic multi-armed bandit setting: counting the number of good arms (CGA) and approximating the number of good arms (AGA).

The setting of stochastic multi-armed bandit problems [Robbins, 1985] is as follows. The algorithm is given n arms, denoted by $\mathbf{A} = (A_1, A_2, \dots, A_n)$. Each arm is a distribution over \mathbb{R} unknown to the algorithm. The algorithm adaptively accesses these arms to receive values independently sampled according to the distributions.

In this paper, we only consider arms with Bernoulli distributions. When given $\gamma \in (0, \frac{1}{2}]$, we define good arms to be arms with mean at least $\frac{1}{2} + \gamma$ and bad arms to be arms with mean at most $\frac{1}{2} - \gamma$.

The problem of $\text{CGA}(\mathbf{A}, \gamma)$ is, when given \mathbf{A} in which every A_i is either good or bad, to output the number of good arms among the given n arms. The algorithm should output the correct answer with probability at least $\frac{2}{3}$.

The problem of $\text{AGA}(\mathbf{A}, \gamma, \epsilon)$ is a similar task to $\text{CGA}(\mathbf{A}, \gamma)$, except that we only need to approximate the correct answer up to an additive error ϵn .

The following lemma is developed by Kaufmann et al. [2016] as a useful tool for proving lower bounds in the stochastic multi-armed bandit setting.

Lemma 12 (Change of measure). *Suppose $\mathbf{A} = (A_1, A_2, \dots, A_n)$ and $\mathbf{A}' = (A'_1, A'_2, \dots, A'_n)$ are two sequences of arms. Suppose an algorithm \mathcal{A} taking n arms as input almost-surely terminates within finite time. Suppose \mathcal{E} is an event defined in the probability space induced by the randomness of the arms and the internal randomness of algorithm \mathcal{A} . Suppose τ_i is the number of queries on A_i made by the algorithm. Then,*

$$\sum_{i=1}^n \mathbb{E}_{\mathcal{A}, \mathbf{A}}[\tau_i] \text{KL}(A_i, A'_i) \geq D(\Pr_{\mathcal{A}, \mathbf{A}}[\mathcal{E}], \Pr_{\mathcal{A}, \mathbf{A}'}[\mathcal{E}]).^9$$

A simple special case ($n = 1$) of the lemma is that to distinguish a coin with mean μ_1 from a coin with mean μ_2 with success probability at least $1 - \delta$, an algorithm needs at least $\frac{D(1-\delta, \delta)}{D(\mu_1, \mu_2)} = \Omega(\frac{1}{D(\mu_1, \mu_2)} \log \frac{1}{\delta})$ queries in expectation for $\mu_1 \neq \mu_2$ and $0 < \delta \leq \frac{2}{5}$.

⁹ $\text{KL}(X, Y)$ denotes the Kullback-Leibler divergence from distribution Y to distribution X . If the two distributions X and Y are Bernoulli with means x and y , their Kullback-Leibler divergence is the relative entropy $D(x, y) = x \log \frac{x}{y} + (1-x) \log \frac{1-x}{1-y}$.

6.5 Lower Bound Results

Our lower bound results in this section are stronger in the sense that the tester has query access to f , knows the distribution to be the uniform distribution \mathcal{U} over a finite ground set X and is only supposed to work on some fixed tie-breaking mechanism. Moreover, we don't require the tester to have success probability at least $\frac{2}{3}$ for *any* S ; instead, the success probability is calculated over the random draw of S and the internal randomness of the tester.

Theorem 13. *Let \mathcal{U} be the uniform distribution over a finite ground set X . There exists a positive constant c such that for any fixed $p \geq 1$, $\epsilon \in (0, \frac{1}{6\sqrt{e}})$ and oracle M using any fixed tie-breaking mechanism, $\mathcal{T}_{\mathcal{U}}^{\text{soft}}(f_{\text{query}}, \epsilon, S, k)$ for p th-power loss requires at least $c \cdot \frac{p}{\epsilon^2}$ queries in the worst case over all finite metric spaces (X, d) .*

Proof. We define $\epsilon' = 6\sqrt{e}\epsilon$. Note that $D(\frac{1-\epsilon'}{2p}, \frac{1}{2p}) = O(\frac{\epsilon'^2}{p})$ for $p \geq 1$ and $\epsilon' \in (0, 1)$. Therefore, we only need to show that a $\mathcal{T}_{\mathcal{U}}^{\text{soft}}(f_{\text{query}}, \epsilon, S, k)$ tester implies an algorithm that distinguishes a coin of mean $\frac{1-\epsilon'}{2p}$ from a coin of mean $\frac{1}{2p}$ with success probability at least $\frac{3}{5}$ using at most the same number of queries. We construct the algorithm in the following way.

The algorithm first constructs a k - NN^{soft} instance with ground set X and distance metric d . We first choose $k = \lceil \frac{c'p^2}{\epsilon^2} \rceil$, $b = \lceil \frac{6}{\epsilon} \rceil$, $N = \lceil c'' \cdot (1+b)k \rceil$ and $m = \lceil \frac{c'''N^2}{1+b} \rceil \geq \frac{6N}{(1+b)\epsilon}$. Here, c' , c'' and c''' are sufficiently large constants. X consists of a star with m centers and bm leaves. Each center C has a distance $d_C \in (1, 2)$ to every leaf in the star and different centers have different values of d_C to avoid ties. The distance between each pair of leaves is 2 and the distance between each pair of centers is 1.

The algorithm then simulates the tester $\mathcal{T}_{\mathcal{U}}^{\text{soft}}(f_{\text{query}}, \epsilon, S, k)$ on this k - NN^{soft} instance without knowing f beforehand. Every time the tester queries the label of a new example, it simulates the result as follows. If the example being queried is a leaf, the result is 1. If the example being queried is a center, the result is obtained to be the same result of an independent toss of the coin we want to distinguish. Finally, if the output of $\mathcal{T}_{\mathcal{U}}^{\text{soft}}(f_{\text{query}}, \epsilon, S, k)$ is above $\frac{1}{2}[(1 - \frac{1}{2p})^p + (1 - \frac{1-\epsilon'}{2p})^p]$, the algorithm then guesses the coin to have mean $\frac{1-\epsilon'}{2p}$. Otherwise, the algorithm guesses the coin to have mean $\frac{1}{2p}$.

Now we show that the above algorithm correctly distinguishes the coins with success probability at least $\frac{3}{5}$. The process of the algorithm, by interchanging the randomness of the labels (coin tosses) and the internal randomness of the $\mathcal{T}_{\mathcal{U}}^{\text{soft}}(f_{\text{query}}, \epsilon, S, k)$ tester, can be viewed in the way that the true labels f are determined before we run the $\mathcal{T}_{\mathcal{U}}^{\text{soft}}(f_{\text{query}}, \epsilon, S, k)$ tester. The leaves all have label 1 and each center is independently labeled 0 or 1 according to the result of a toss of the coin. After the labels f are decided, the $\mathcal{T}_{\mathcal{U}}^{\text{soft}}(f_{\text{query}}, \epsilon, S, k)$ tester is then simulated to approximate the p th-power loss of the k - NN^{soft} instance up to additive error ϵ with success probability at least $\frac{2}{3}$.

Suppose the coin to be distinguished has mean μ . Note that the total number of points in the ground set is $(1+b)m = \Omega(N^2)$, therefore we can make sure with probability at least $1 - \frac{1}{40}$ that no two unlabeled examples lie on the same point. Because each random example has probability $\frac{1}{1+b}$ to lie in the centers and $N \geq c'' \cdot (1+b)k$, therefore by choosing a sufficiently large c'' , we can make sure with probability at least $1 - \frac{1}{40}$ that in the unlabeled sample pool, there are at least k examples lying at the centers. These two events happen at the same time with probability at least $1 - \frac{1}{20}$ by the Union Bound. Conditioned on these two events happening, by a sufficiently

large choice of c' , among those unlabeled examples lying at the centers, we can make sure that with probability at least $1 - \frac{1}{20}$, the average of the labels of the k examples with smallest d_C is contained in $(\mu - \frac{\epsilon}{6p}, \mu + \frac{\epsilon}{6p})$. All these events happen at the same time with probability at least $(1 - \frac{1}{20})^2 \geq 1 - \frac{1}{10}$, and in this case, every leaf outside the unlabeled pool S has L^1 error in $(1 - \mu - \frac{\epsilon}{6p}, 1 - \mu + \frac{\epsilon}{6p})$ and thus has p th-power loss in $((1 - \mu)^p - \frac{\epsilon}{6}, (1 - \mu)^p + \frac{\epsilon}{6})$. The total number of leaves in the unlabeled pool S and centers is upper bounded by the size N of the pool plus m , which contributes only a $\frac{N+m}{(b+1)m} \leq \frac{\epsilon}{3}$ fraction of the total number of points. Therefore, with probability at least $1 - \frac{1}{10}$, the average p th-power loss of all points is contained in $((1 - \mu)^p - \frac{\epsilon}{2}, (1 - \mu)^p + \frac{\epsilon}{2})$.

Note that $(1 - \frac{1-c'}{2p})^p - (1 - \frac{1}{2p})^p > 3\epsilon$, therefore the algorithm correctly guesses the mean of the coin with probability at least $(1 - \frac{1}{10}) \cdot \frac{2}{3} = \frac{3}{5}$. \square

Theorem 14. *There exists a positive constant c such that for any fixed $k \in \mathbb{N}^*$, $\epsilon \in (0, \frac{1}{4})$ and oracle M using any fixed tie-breaking mechanism, if there is a $\mathcal{T}_U^{\text{hard}}(f_{\text{query}}, \epsilon, S, k)$ tester using at most q queries in the worst case, then there is an $\text{AGA}(\mathbf{A}, \gamma, 2\epsilon)$ algorithm using at most $O(q)$ queries in the worst case where $\gamma = \min \left\{ \frac{1}{2}, c \cdot \sqrt{\frac{\log \frac{1}{\epsilon}}{k}} \right\}$.*

Proof. Since the success probability can be boosted by repetition, we only show an $\text{AGA}(\mathbf{A}, \gamma, 2\epsilon)$ algorithm with success probability at least $\frac{3}{5}$. Given any instance of $\text{AGA}(\mathbf{A}, \gamma, 2\epsilon)$ with total number of arms equal to n , the algorithm constructs a ground set X and the distance metric d on it to form a k - NN^{hard} instance in the following way. We first choose $b = \lceil \frac{3}{\epsilon} \rceil$, $N = \lceil c' \cdot (1+b)n(k + \log \frac{1}{\epsilon}) \rceil$ and $m = \lceil \frac{c''N^2}{(1+b)n} \rceil \geq \frac{3N}{(1+b)n\epsilon}$. Here, c' and c'' are sufficiently large constants. X consists of n identical stars, each corresponds to an arm, with the distances between stars to be very large. Each star consists of m centers and bm leaves. Each center C has a distance $d_C \in (1, 2)$ to every leaf in the same star and different centers have different values of d_C to avoid ties. The distance between each pair of leaves in the same star is 2 and the distance between each pair of centers in the same star is 1.

The algorithm then simulates the tester $\mathcal{T}_U^{\text{hard}}(f_{\text{query}}, \epsilon, S, k)$ on this k - NN^{hard} instance without knowing f beforehand. Every time the tester queries the label of a new example, it simulates the result as follows. If the example being queried is a leaf, the result is 0. If the example being queried is a center, the result is obtained to be the same result of an independent query to the corresponding arm. Finally, the algorithm outputs $\hat{\alpha}n$ when the $\mathcal{T}_U^{\text{hard}}(f_{\text{query}}, \epsilon, S, k)$ tester outputs α .

Now we show that the above is an $\text{AGA}(\mathbf{A}, \gamma, 2\epsilon)$ algorithm with success probability at least $\frac{3}{5}$. The process of the algorithm, by interchanging the randomness of the labels (arms) and the internal randomness of the $\mathcal{T}_U^{\text{hard}}(f_{\text{query}}, \epsilon, S, k)$ tester, can be viewed in the way that the true labels f are determined before we run the $\mathcal{T}_U^{\text{hard}}(f_{\text{query}}, \epsilon, S, k)$ tester. The leaves all have labels 0 and each center is independently labeled 0 or 1 according to the result of a query to the corresponding arm. After the labels f are decided, the $\mathcal{T}_U^{\text{hard}}(f_{\text{query}}, \epsilon, S, k)$ tester is then simulated to approximate the error rate of the k - NN^{hard} instance up to additive error ϵ with success probability at least $\frac{2}{3}$.

Let's say a star is good (bad) if it corresponds to a good (bad) arm. Suppose there are ξn good arms, and thus ξn good stars. Note that there are $(1+b)mn = \Omega(N^2)$ points in the ground set, we can make sure with probability at least $1 - \frac{1}{20}$ that no two unlabeled examples lie on the same point, on which the following discussion is conditioned. Let's first fix a star R whose corresponding arm has mean μ . Because each random example has probability $\frac{1}{(1+b)n}$ to lie in the centers of R

and $N \geq c' \cdot (1+b)n(k + \log \frac{1}{\epsilon})$, therefore by choosing a sufficiently large c' , we can make sure with probability at most $\frac{\frac{\epsilon}{120}}{1-\frac{1}{20}}$ that in the unlabeled sample pool, there are less than k examples lying at the centers of R . Therefore, by a sufficiently large choice of c , among those unlabeled examples lying at the centers of R , we can make sure that with probability at least $(1 - \frac{\frac{\epsilon}{120}}{1-\frac{1}{20}})(1 - \frac{\epsilon}{200}) \geq 1 - \frac{\frac{\epsilon}{60}}{1-\frac{1}{20}}$, the average of the labels of the k examples with smallest d_C is contained in $(\mu - \gamma, \mu + \gamma)$, or R is *satisfied*. By Markov's Inequality, with probability at least $1 - \frac{\frac{1}{20}}{1-\frac{1}{20}}$, or $1 - \frac{1}{10}$ if we unwrap the conditional probability of $1 - \frac{1}{20}$, at least a $(1 - \frac{\epsilon}{3})$ fraction of all the n stars are satisfied. In a satisfied star, any leaf that is not in the unlabeled pool has L^1 error 1 if the star is good and L^1 error 0 if the star is bad. Note that there are at most N leaves in the unlabeled pool, contributing at most an $\frac{N}{(1+b)mn} \leq \frac{\epsilon}{3}$ fraction of the total number of points. Also there are only mn centers in total, contributing at most an $\frac{mn}{(1+b)mn} \leq \frac{\epsilon}{3}$ fraction of the total number of points. Therefore, with probability at least $1 - \frac{1}{10}$, the average error of all points is contained in $[\xi - \epsilon, \xi + \epsilon]$, which implies that with probability at least $(1 - \frac{1}{10}) \cdot \frac{2}{3} = \frac{2}{5}$, $\hat{\alpha} \in [\xi - 2\epsilon, \xi + 2\epsilon]$. \square

The above theorem shows that a query complexity lower bound for $\text{AGA}(\mathbf{A}, \gamma, \epsilon)$ can imply a query complexity lower bound for $\mathcal{T}_{\mathcal{U}}^{\text{hard}}(f_{\text{query}}, \epsilon, S, k)$. $\text{AGA}(\mathbf{A}, \gamma, \epsilon)$ has a simple algorithm requiring $O(\frac{1}{\gamma^2 \epsilon^2} \log \frac{1}{\epsilon})$ queries as follows. The algorithm randomly picks $O(\frac{1}{\epsilon^2})$ arms. For each of the picked arms, the algorithm queries it $O(\frac{1}{\gamma^2} \log \frac{1}{\epsilon})$ times and thinks of it as “good” if more than half of the results are positive and “bad” otherwise. The algorithm outputs the fraction of “good” arms among the picked arms.

If we assume the simple $O(\frac{1}{\gamma^2 \epsilon^2} \log \frac{1}{\epsilon})$ query complexity for AGA is not improvable, then Theorem 14 implies that the $O(\frac{k}{\epsilon^2})$ query complexity for $\mathcal{T}^{\text{hard}}$ is also not improvable. In other words, if for every sequences $\epsilon_n \rightarrow 0$ and $\gamma_n \rightarrow 0$, there exists a positive constant c such that $\text{AGA}(\mathbf{A}, \epsilon_i, \gamma_i)$ needs at least $c \cdot \frac{1}{\gamma_i^2 \epsilon_i^2} \log \frac{1}{\epsilon_i}$ queries in the worst case, then according to Theorem 14, we know for any sequences $\{k_n\}, \{\epsilon_n\}$ such that $\epsilon_n \rightarrow 0, \frac{k_n}{\log \frac{1}{\epsilon_n}} \rightarrow \infty$, there exists a positive constant c' such that the tester $\mathcal{T}_{\mathcal{U}}^{\text{hard}}(f_{\text{query}}, \epsilon, S, k)$ for k - NN^{hard} algorithms needs at least $c' \cdot \frac{k_i}{\epsilon_i^2}$ queries in the worst case.

The following theorem states an unconditional lower bound $\Omega(\frac{k}{\epsilon \log \frac{1}{\epsilon}})$ for the query complexity of $\mathcal{T}^{\text{hard}}$, implying that the linear growth with respect to k in the query complexity of $\mathcal{T}^{\text{hard}}$ can't be improved.

Theorem 15. *There exists a positive constant c such that for any fixed $k \in \mathbb{N}^*, \epsilon \in (0, \frac{1}{4})$ and oracle M using any fixed tie-breaking mechanism, $\mathcal{T}_{\mathcal{U}}^{\text{hard}}(f_{\text{query}}, \epsilon, S, k)$ requires at least $c \cdot \frac{k}{\epsilon \log \frac{1}{\epsilon}}$ queries in the worst case.*

Before proving Theorem 15, we first show a query complexity lower bound for CGA.

Lemma 16. *There exists a universal constant c such that for any fixed $\gamma \in (0, \frac{1}{2}]$ and $n \in \mathbb{N}^*$, $\text{CGA}(\mathbf{A}, \gamma)$ requires at least $c \cdot \frac{n}{\gamma^2}$ queries in the worst case, where n is the number of arms in \mathbf{A} .*

Proof of Lemma 16. Obviously, n is a query complexity lower bound since we need to query each arm at least once. So in the rest of the proof, we assume $\gamma < \frac{1}{4}$. We use G to denote the good

arm with mean $\frac{1}{2} + \gamma$ and B to denote the bad arm with mean $\frac{1}{2} - \gamma$. Then, $\text{KL}(G, B) = O(\gamma^2)$. We claim a stronger fact that for any $0 \leq q \leq n$ and any instance consisting of q G 's and $n - q$ B 's, $\text{CGA}(\mathbf{A}, \gamma)$ needs at least $c \cdot \frac{1}{\gamma^2}$ queries on *every* of the n arms. By symmetry between “good” and “bad”, we only show that every G arm needs to be queried at least $c \cdot \frac{1}{\gamma^2}$ times. The reason is as follows. Suppose $\mathbf{A} = (A_1, A_2, \dots, A_n)$ in which $A_i = G$ for $1 \leq i \leq q$ and $A_i = B$ otherwise. We define $\mathbf{A}' = (A'_1, A'_2, \dots, A'_n)$ in which $A'_i = G$ for $1 \leq i \leq p - 1$ and $A'_i = B$ otherwise. The only difference between \mathbf{A} and \mathbf{A}' is that $A_p = G$ while $A'_p = B$. We use \mathcal{E} to denote the event that $\text{CGA}(\mathbf{A}, \gamma)$ outputs p . By Lemma 12, we know $\mathbb{E}[\tau_p] \cdot O(\gamma^2) \geq D(\frac{2}{3}, \frac{1}{3}) = \Omega(1)$ and thus $\mathbb{E}[\tau_p] = \Omega(\frac{1}{\gamma^2})$. For similar reasons, we can show for all $1 \leq i \leq p$ that $\mathbb{E}[\tau_i] = \Omega(\frac{1}{\gamma^2})$, which completes the proof. \square

Proof of Theorem 15. Lemma 16 immediately implies the existence of a positive constant c' such that for any fixed $\epsilon \in (0, \frac{1}{2})$ and $\gamma \in (0, \frac{1}{2}]$, $\text{AGA}(\mathbf{A}, \gamma, \epsilon)$ requires at least $c' \cdot \frac{1}{\gamma^2 \epsilon}$ queries in the worst case by choosing $n = \lceil \frac{1}{2\epsilon} \rceil - 1$. Then, by Theorem 14, we get an $\Omega(\frac{1}{\left(\min\left\{\frac{1}{2}, \sqrt{\frac{\log \frac{1}{\epsilon}}{k}}\right\}\right)^2 \cdot \frac{1}{2\epsilon}}) = \Omega(\frac{k}{\epsilon \log \frac{1}{\epsilon}})$

lower bound for $\mathcal{T}_U^{\text{hard}}(f_{\text{query}}, \epsilon, S, k)$ for $k \in \mathbb{N}^*$ and $\epsilon \in (0, \frac{1}{4})$. \square

Acknowledgements

This work was supported in part by the National Science Foundation under grant CCF-1525971. This work was conducted in part while Avrim Blum was at Carnegie Mellon University and while Lunjia Hu was visiting at Carnegie Mellon University and TTI-Chicago.

References

- Maria-Florina Balcan, Eric Blais, Avrim Blum, and Liu Yang. Active property testing. In *53rd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 21–30. IEEE, 2012.
- Thomas Cover and Peter Hart. Nearest neighbor pattern classification. *IEEE transactions on information theory*, 13(1):21–27, 1967.
- Evelyn Fix and Joseph Lawson Hodges. Discriminatory analysis. nonparametric discrimination: consistency properties. *International Statistical Review/Revue Internationale de Statistique*, 57(3):238–247, 1989.
- Evelyn Fix and Joseph L Hodges Jr. Discriminatory analysis-nonparametric discrimination: consistency properties. Technical report, California Univ Berkeley, 1951.
- Oded Goldreich, Shari Goldwasser, and Dana Ron. Property testing and its connection to learning and approximation. *Journal of the ACM (JACM)*, 45(4):653–750, 1998.
- Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best arm identification in multi-armed bandit models. *Journal of Machine Learning Research*, 17(1):1–42, 2016.

Michael Kearns and Dana Ron. Testing problems with sub-learning sample complexity. In *Proceedings of the eleventh annual conference on Computational learning theory*, pages 268–279. ACM, 1998.

Michal Parnas, Dana Ron, and Ronitt Rubinfeld. Tolerant property testing and distance approximation. *Journal of Computer and System Sciences*, 72(6):1012–1042, 2006.

Herbert Robbins. Some aspects of the sequential design of experiments. In *Herbert Robbins Selected Papers*, pages 169–177. Springer, 1985.

Richard Miles Royall. *A class of non-parametric estimates of a smooth regression function*. PhD thesis, Department of Statistics, Stanford University, 1966.

Ronitt Rubinfeld and Madhu Sudan. Robust characterizations of polynomials with applications to program testing. *SIAM Journal on Computing*, 25(2):252–271, 1996.

A Distance Approximation for Disjoint Unions of Properties

In this section, we extend the theorem of Balcan et al. [2012] that disjoint unions of testable properties are testable from property testing to distance approximation.

We first introduce the definition of disjoint unions of properties in [Balcan et al., 2012]. Suppose the ground set X is partitioned as a disjoint union $\bigcup_{i=1}^m X_i$. For every X_i , there is a property (concept class) $\mathcal{C}_i \neq \emptyset$. The disjoint union of these properties is defined to be $\mathcal{C} = \{f \in \{0, 1\}^X : \forall 1 \leq i \leq m, f|_{X_i} \in \mathcal{C}_i\}$.

Let \mathcal{D} be a distribution over X . Suppose the conditional distribution of \mathcal{D} on X_i is denoted by \mathcal{D}_i and the probability $\Pr_{x \sim \mathcal{D}}[x \in X_i]$ is denoted by p_i .

Theorem 17. *Suppose $\epsilon \in (0, \frac{1}{2})$. Suppose there is a $\text{DA}_{\mathcal{D}_i}(f_{\text{active}}, \frac{\epsilon}{2})$ algorithm for every $1 \leq i \leq m$ using at most q queries on N unlabeled examples. Then, there is a $\text{DA}_{\mathcal{D}}(f_{\text{active}}, \epsilon)$ algorithm using at most $O(\frac{q}{\epsilon^2} \log \frac{1}{\epsilon})$ queries on $O(\frac{mN}{\epsilon} \log \frac{1}{\epsilon})$ unlabeled examples. If the $\text{DA}_{\mathcal{D}_i}(f_{\text{active}}, \frac{\epsilon}{2})$ algorithm can perform on unknown distributions, then the $\text{DA}_{\mathcal{D}}(f_{\text{active}}, \epsilon)$ algorithm can also perform on unknown distributions, though we need extra $O(\frac{1}{\epsilon^2})$ unlabeled examples.*

Proof. The $\text{DA}_{\mathcal{D}}(f_{\text{active}}, \epsilon)$ algorithm is constructed as follows. The algorithm chooses $s = O(\frac{1}{\epsilon^2})$, receives an unlabeled pool of size $O(\frac{mN}{\epsilon} \log s)$ and independently chooses s indices i_1, i_2, \dots, i_s from $\{1, 2, \dots, m\}$ according to distribution $\{p_i\}_{1 \leq i \leq m}$. This can be achieved by looking at on which X_i 's the extra s unlabeled examples are, when the distribution \mathcal{D} is unknown. Then for each $1 \leq j \leq s$, if there are enough ($O(N \log s)$) unlabeled examples lying in X_{i_j} , the algorithm repeats $\text{DA}_{\mathcal{D}_{i_j}}(f_{\text{active}}, \frac{\epsilon}{2})$ for $O(\log s)$ times to calculate an estimator $\widehat{\text{dist}}_{i_j}$ of the distance from f to \mathcal{C} on \mathcal{D}_{i_j} up to an additive error $\frac{\epsilon}{2}$ with success probability at least $1 - \frac{1}{9s}$; otherwise, define $\widehat{\text{dist}}_{i_j} = 0$. The final output of the algorithm is $\frac{1}{s} \cdot \sum_{j=1}^s \widehat{\text{dist}}_{i_j}$.

To prove the correctness of the above algorithm, we first define $\text{dist}_i := \inf_{g \in \mathcal{C}} \text{dist}_{\mathcal{D}_i}(f, g)$ and $\text{dist} := \inf_{g \in \mathcal{C}} \text{dist}_{\mathcal{D}}(f, g)$. Note that $\text{dist} = \sum_{i=1}^m p_i \text{dist}_i$.

For every $1 \leq i \leq m$, we further define $\text{dist}'_i = \begin{cases} \text{dist}_i, & \text{if } p_i \geq \frac{\epsilon}{4m} \\ 0, & \text{if } p_i < \frac{\epsilon}{4m} \end{cases}$ and $\text{dist}''_i = \begin{cases} \text{dist}_i, & \text{if } p_i \geq \frac{\epsilon}{4m} \\ 1, & \text{if } p_i < \frac{\epsilon}{4m} \end{cases}$.

Then $\text{dist} - \frac{\epsilon}{4} \leq \sum_{i=1}^m p_i \text{dist}'_i \leq \sum_{i=1}^m p_i \text{dist}''_i \leq \text{dist} + \frac{\epsilon}{4}$. By the Chernoff Bound, $s = O(\frac{1}{\epsilon^2})$ is enough to make sure with probability at least $1 - \frac{1}{9}$ that $\text{dist} - \frac{\epsilon}{2} < \frac{1}{s} \sum_{j=1}^s \text{dist}'_{i_j} \leq \frac{1}{s} \sum_{j=1}^s \text{dist}''_{i_j} < \text{dist} + \frac{\epsilon}{2}$.

Note that the unlabeled pool has size $O(\frac{mN}{\epsilon} \log s)$, which is enough to make sure that with probability at least $1 - \frac{1}{9}$, for every i_j with $p_{i_j} \geq \frac{\epsilon}{4m}$, there are enough ($O(N \log s)$) unlabeled examples lying in X_{i_j} . Therefore, with probability at least $(1 - \frac{1}{9})(1 - s \cdot \frac{1}{9s}) \geq 1 - \frac{2}{9}$, for all i_j such that $p_{i_j} \geq \frac{\epsilon}{4m}$, it holds that $|\widehat{\text{dist}}_{i_j} - \text{dist}_{i_j}| \leq \frac{\epsilon}{2}$.

Finally, by the Union Bound, we know with probability at least $1 - \frac{1}{3}$, it holds that $\text{dist} - \epsilon < \frac{1}{s} \sum_{j=1}^s \text{dist}'_{i_j} - \frac{\epsilon}{2} \leq \frac{1}{s} \sum_{j=1}^s \widehat{\text{dist}}_{i_j} \leq \frac{1}{s} \sum_{j=1}^s \text{dist}''_{i_j} + \frac{\epsilon}{2} < \text{dist} + \epsilon$. \square