
Comparing Apples to Oranges: Learning Similarity Functions for Data Produced by Different Distributions

Leonidas Tsepenekas
JPMorganChase AI Research
New York, USA
leonidas.tsepenekas@jpmchase.com

Ivan Brugere
JPMorganChase AI Research
New York, USA
ivan.brugere@jpmchase.com

Freddy Lecue
JPMorganChase AI Research
New York, USA
freddy.lecue@jpmchase.com

Daniele Magazzeni
JPMorganChase AI Research
London, UK
daniele.magazzeni@jpmorgan.com

Abstract

Similarity functions measure how comparable pairs of elements are, and play a key role in a wide variety of applications, e.g., notions of Individual Fairness abiding by the seminal paradigm of [Dwork et al. \[2012\]](#), as well as Clustering problems. However, access to an accurate similarity function should not always be considered guaranteed, and this point was even raised by [Dwork et al. \[2012\]](#). For instance, it is reasonable to assume that when the elements to be compared are produced by different distributions, or in other words belong to different “demographic” groups, knowledge of their true similarity might be very difficult to obtain. In this work, we present an efficient sampling framework that learns these across-groups similarity functions, using only a limited amount of experts’ feedback. We show analytical results with rigorous theoretical bounds, and empirically validate our algorithms via a large suite of experiments.

1 Introduction

Given a feature space \mathcal{I} , a similarity function $\sigma : \mathcal{I}^2 \mapsto \mathbb{R}_{\geq 0}$ measures how comparable any pair of elements $x, x' \in \mathcal{I}$ are. The function σ can also be interpreted as a distance function, where the smaller $\sigma(x, x')$ is, the more similar x and x' are. Such functions are crucially used in a variety of AI/ML problems, and in each such case σ is assumed to be known.

The most prominent applications where similarity functions have a central role involve considerations of individual fairness. Specifically, all such individual fairness paradigms stem from the seminal work of [Dwork et al. \[2012\]](#), in which fairness is defined as *treating similar individuals similarly*. In more concrete terms, such paradigms interpret the aforementioned abstract definition of fairness as guaranteeing that for every pair of individuals x and y , the difference in the quality of service x and y receive (a.k.a. their received treatment) is upper bounded by their respective similarity value $\sigma(x, y)$; the more similar the individuals are, the less different their quality of service will be. Therefore, any algorithm that needs to abide by such concepts of fairness, should always be able to access the similarity score for every pair of individuals that are of interest.

Another family of applications where similarity functions are vital, involves Clustering problems. In a clustering setting, e.g. the standard k -means task, the similarity function is interpreted as a distance

function, that serves as the metric space in which we need to create the appropriate clusters. Clearly, the aforementioned metric space is always assumed to be part of the input.

Nonetheless, it is not realistic to assume that a reliable and accurate similarity function is always given. This issue was even raised in the work of [Dwork et al. \[2012\]](#), where it was acknowledged that the computation of σ is not trivial, and thus should be deferred to third parties. The starting point of our work here is the observation that there exist scenarios where computing similarity can be assumed as easy (in other words given), while in other cases this task would be significantly more challenging. Specifically, we are interested in scenarios where there are multiple distributions that produce elements of \mathcal{I} . We loosely call each such distribution a “demographic” group, interpreting it as the stochastic way in which members of this group are produced. *In this setting, computing the similarity value of two elements that are produced according to the same distribution, seems intuitively much easier compared to computing similarity values for elements belonging to different groups.* We next present a few motivating examples that clarify this statement.

Individual Fairness: Consider a college admissions committee that needs access to an accurate similarity function for students, so that it provides a similar likelihood of acceptance to similar applicants. Let us focus on the following two demographic groups. The first being students from affluent families living in privileged communities and having access to the best quality schools and private tutoring. The other group would consist of students from low-income families, coming from a far less privileged background. Given this setting, the question at hand is “**Should two students with comparable feature vectors (e.g., SAT scores, strength of school curriculum, number of recommendation letters) be really viewed as similar, when they belong to different groups?**” At first, it appears that directly comparing two students based on their features can be an accurate way to elicit their similarity only if the students belong to the same demographic (belonging to the same group serves as a normalization factor). However, this trivial approach might hurt less privileged students when they are compared to students of the first group. This is because such a simplistic way of measuring similarity does not reflect potential that is undeveloped due to unequal access to resources. Hence, accurate across-groups comparisons that take into account such delicate issues, appear considerably more intricate.

Clustering: Suppose that a marketing company has a collection of user data that wants to cluster, with its end goal being a downstream market segmentation analysis. However, as it is usually the case, the data might come from different sources, e.g., data from private vendors and data from government bureaus. In this scenario, each data source might have its own way of representing user information, e.g., each source might use a unique subset of features. Therefore, eliciting the distance metric required for the clustering task should be straightforward for data coming from the same source, while across-sources distances would certainly require extra care, e.g., how can one extract the distance of two vectors containing different sets of features?

As suggested by earlier work on computing similarity functions for applications of individual fairness [[Ilvento, 2019](#)], when obtaining similarity values is an overwhelming task, one can employ the advice of domain experts. Such experts can be given any pair of elements, and in return produce their true similarity value. However, utilizing this experts’ advice can be thought of as very costly, and hence it should be used sparingly. For example, in the case of comparing students from different economic backgrounds, the admissions committee can reach out to regulatory bodies or civil rights organizations. Nonetheless, resorting to these experts for every student comparison that might arise, is clearly not a sustainable solution. Therefore, our goal in this paper is to learn the across-groups similarity functions, using as few queries to experts as possible.

2 Preliminaries and contribution

For the ease of exposition, we accompany the formal definitions with brief demonstrations on how they could relate to the previously mentioned college applications use-case.

Let \mathcal{I} denote the feature space of elements; for instance each $x \in \mathcal{I}$ could correspond to a valid student profile. We assume that elements come from γ known “demographic” groups, where $\gamma \in \mathbb{N}$, and each group $\ell \in [\gamma]$ is governed by an unknown distribution \mathcal{D}_ℓ over \mathcal{I} . We use $x \sim \mathcal{D}_\ell$ to denote a randomly drawn x from \mathcal{D}_ℓ . Further, we use $x \in \mathcal{D}_\ell$ to denote that x is an element in the support of \mathcal{D}_ℓ , and thus x is a member of group ℓ . In the college admissions scenario where we have two demographic groups, there will be two distributions \mathcal{D}_1 and \mathcal{D}_2 dictating how the profiles of

privileged and non-privileged students are respectively produced. Observe now that for a specific $x \in \mathcal{I}$, we might have $x \in \mathcal{D}^\cdot$ and $x \in \mathcal{D}^{\cdot o}$, for $\ell \neq \ell'$. Hence, in our model group membership is important, and every time we are considering an element $x \in \mathcal{I}$, we know which distribution produced x , e.g., whether the profile x belongs to a privileged or non-privileged student.

For every group $\ell \in [\gamma]$ there is an intra-group similarity function $d^\cdot : \mathcal{I}^2 \mapsto \mathbb{R}_{\geq 0}$, such that for all $x, y \in \mathcal{D}^\cdot$ we have $d^\cdot(x, y)$ representing the true similarity between x, y . In addition, the smaller $d^\cdot(x, y)$ is, the more similar x, y . Note here that the function d^\cdot is only used to compare members of group ℓ (in the college admissions example, the function d_1 would only be used to compare privileged students with each other). Further, a common assumption for functions measuring similarity is that they are metric¹ [Yona and Rothblum, 2018, Kim et al., 2018, Ilvento, 2019, Mukherjee et al., 2020, Wang et al., 2019]. We also adopt the metric assumption for the function d^\cdot . Finally, based on the earlier discussion regarding computing similarity between elements of the same group, we assume that d^\cdot is known, and given as part of the instance.

Moreover, for any two groups ℓ and ℓ' there exists an *unknown* across-groups similarity function $\sigma^{\cdot, \cdot o} : \mathcal{I}^2 \mapsto \mathbb{R}_{\geq 0}$, such that for all $x \in \mathcal{D}^\cdot$ and $y \in \mathcal{D}^{\cdot o}$, $\sigma^{\cdot, \cdot o}(x, y)$ represents the true similarity between x, y . Again, the smaller $\sigma^{\cdot, \cdot o}(x, y)$ is, the more similar the two elements, and for a meaningful use of $\sigma^{\cdot, \cdot o}$ we must make sure that x is a member of group ℓ and y a member of group ℓ' . In the college admissions scenario, $\sigma_{1,2}$ is the way you can accurately compare a privileged and a non-privileged student. Finally, to capture the metric nature of a similarity function, we impose the following mild properties on $\sigma^{\cdot, \cdot o}$, which can be viewed as across-groups triangle inequalities:

1. **Property \mathcal{M}_1 :** $\sigma^{\cdot, \cdot o}(x, y) \leq d^\cdot(x, z) + \sigma^{\cdot, \cdot o}(z, y)$ for every $x, z \in \mathcal{D}^\cdot$ and $y \in \mathcal{D}^{\cdot o}$.
2. **Property \mathcal{M}_2 :** $\sigma^{\cdot, \cdot o}(x, y) \leq \sigma^{\cdot, \cdot o}(x, z) + d^{\cdot o}(z, y)$ for every $x \in \mathcal{D}^\cdot$ and $y, z \in \mathcal{D}^{\cdot o}$.

In terms of the college admissions use-case, \mathcal{M}_1 and \mathcal{M}_2 try to capture reasonable assumptions of the following form. If a non-privileged student x is similar to another non-privileged student z , and z is similar to a privileged student y , then x and y should also be similar to each other.

Observe now that the collection of all similarity values (intra-group and across-groups) in our model does not axiomatically yield a valid metric space. This is due to the following reasons. **1)** If $\sigma^{\cdot, \cdot o}(x, y) = 0$ for $x \in \mathcal{D}^\cdot$ and $y \in \mathcal{D}^{\cdot o}$, then we do not necessarily have $x = y$. **2)** It is not always the case that $d^\cdot(x, y) \leq \sigma^{\cdot, \cdot o}(x, z) + \sigma^{\cdot, \cdot o}(y, z)$ for $x, y \in \mathcal{D}^\cdot, z \in \mathcal{D}^{\cdot o}$. **3)** It is not always the case that $\sigma^{\cdot, \cdot o}(x, y) \leq \sigma^{\cdot, \cdot oo}(x, z) + \sigma^{\cdot oo, \cdot o}(z, y)$ for $x \in \mathcal{D}^\cdot, y \in \mathcal{D}^{\cdot o}, z \in \mathcal{D}^{\cdot oo}$.

However, not having the collection of similarity values necessarily produce a metric space is not a weakness of our model. On the contrary, we view this as one of its strongest aspects. For one thing, imposing a complete metric constraint on the case of intricate across-groups comparisons sounds unrealistic and very restrictive. Further, even though existing literature treats similarity functions as metric ones, the seminal work of Dwork et al. [2012] mentions that this should not always be the case. Hence, our model is more general than the current literature.

Goal of Our Problem: We want for any two groups ℓ, ℓ' to compute a function $f^{\cdot, \cdot o} : \mathcal{I}^2 \mapsto \mathbb{R}_{\geq 0}$, such that $f^{\cdot, \cdot o}(x, y)$ is our estimate of similarity for any $x \in \mathcal{D}^\cdot$ and $y \in \mathcal{D}^{\cdot o}$. Specifically, we seek a PAC (Probably Approximately Correct) guarantee, where for any given accuracy and confidence parameters $\epsilon, \delta \in (0, 1)$ we have:

$$\Pr_{x \sim \mathcal{D}^\cdot, y \sim \mathcal{D}^{\cdot o}} \overset{\text{h}}{f^{\cdot, \cdot o}(x, y) - \sigma^{\cdot, \cdot o}(x, y) > \epsilon} \overset{\text{i}}{\leq \delta}$$

The subscript in the above probability corresponds to two independent random choices, one $x \sim \mathcal{D}^\cdot$ and one $y \sim \mathcal{D}^{\cdot o}$. In other words, we want for any given pair our estimate to be ϵ -close to the real similarity value, with probability at least $1 - \delta$, where ϵ and δ are user-specified parameters.

As for tools to learn $f^{\cdot, \cdot o}$, we only require two things. At first, for each group ℓ we want a set S^\cdot of i.i.d. samples from \mathcal{D}^\cdot . Obviously, the total number of used samples should be polynomial in the input parameters, i.e., polynomial in $\gamma, \frac{1}{\epsilon}$ and $\frac{1}{\delta}$. Secondly, we require access to an expert oracle, which given any $x \in S^\cdot$ and $y \in S^{\cdot o}$ for any ℓ and ℓ' , returns the true similarity value $\sigma^{\cdot, \cdot o}(x, y)$. We refer to a single invocation of the oracle as a query. Since there is a cost to collecting expert feedback, an additional objective in our problem is minimizing the number of oracle queries.

¹A function d is metric if **a)** $d(x, y) = 0$ iff $x = y$, **b)** $d(x, y) = d(y, x)$ and **c)** $d(x, y) \leq d(x, z) + d(z, y)$ for all x, y, z (triangle inequality).

2.1 Outline and discussion of our results

In Section 4 we present our theoretical results. We begin with a simple and very intuitive learning algorithm which achieves the following guarantees.

Theorem 2.1. *For any given parameters $\epsilon, \delta \in (0, 1)$, the simple algorithm produces a similarity approximation function $f_{\cdot, \cdot, \cdot}$ for every ℓ and ℓ' , such that:*

$$\begin{aligned} \Pr[\text{Error}_{(\cdot, \cdot, \cdot)}] &:= \Pr_{\substack{x \sim \mathcal{D}_{\cdot, \cdot} \\ y \sim \mathcal{D}_{\cdot, \cdot}; \mathcal{A}}} [f_{\cdot, \cdot, \cdot}(x, y) - \sigma_{\cdot, \cdot, \cdot}(x, y)] = \omega(\epsilon) \\ &= O(\delta + p_{\cdot}(\epsilon, \delta) + p_{\cdot'}(\epsilon, \delta)) \end{aligned}$$

The randomness here is of three independent sources. The internal randomness \mathcal{A} of the algorithm, a choice $x \sim \mathcal{D}_{\cdot}$, and a choice $y \sim \mathcal{D}_{\cdot'}$. The algorithm requires $\frac{1}{2} \log \frac{1}{2}$ samples from each group, and utilizes $\frac{(-1)}{2} \log^2 \frac{1}{2}$ oracle queries.

In plain English, the above theorem says that with a polynomial number of samples and queries, the algorithm achieves an $O(\epsilon)$ accuracy with high probability, i.e., with probability $\omega(1 - \delta - p_{\cdot}(\epsilon, \delta) - p_{\cdot'}(\epsilon, \delta))$. The definition of the functions p_{\cdot} is presented next.

Definition 2.2. For each group ℓ , we use $p_{\cdot}(\epsilon, \delta)$ to denote the probability of sampling an (ϵ, δ) -rare element of \mathcal{D}_{\cdot} . We define as (ϵ, δ) -rare for \mathcal{D}_{\cdot} , an element $x \in \mathcal{D}_{\cdot}$ for which there is a less than δ chance of sampling $x' \sim \mathcal{D}_{\cdot}$ with $d_{\cdot}(x, x') \leq \epsilon$. Formally, $x \in \mathcal{D}_{\cdot}$ is (ϵ, δ) -rare iff $\Pr_{x' \sim \mathcal{D}_{\cdot}}[d_{\cdot}(x, x') \leq \epsilon] < \delta$, and $p_{\cdot}(\epsilon, \delta) = \Pr_{x \sim \mathcal{D}_{\cdot}}[x \text{ is } (\epsilon, \delta)\text{-rare for } \mathcal{D}_{\cdot}]$. Intuitively, a rare element should be interpreted as an “isolated” member of the group, in the sense that it is at most δ -likely to encounter another element that is ϵ -similar to it. For instance, a privileged student is considered isolated, if only a small fraction of other privileged students have a profile similar to them.

Clearly, to get a PAC guarantee where the algorithm’s error probability for ℓ and ℓ' is $O(\delta)$, we need $p_{\cdot}(\epsilon, \delta), p_{\cdot'}(\epsilon, \delta) = O(\delta)$. We hypothesize that in realistic distributions each $p_{\cdot}(\epsilon, \delta)$ should indeed be fairly small, and this hypothesis is actually validated by our experiments. The reason we believe this hypothesis to be true, is that very frequently real data demonstrate high concentration around certain archetypal elements. Hence, this sort of distributional density does not leave room for isolated elements in the rest of the space. Nonetheless, we also provide a strong *no free lunch* result for the values $p_{\cdot}(\epsilon, \delta)$, which shows that any practical PAC-algorithm necessarily depends on them. *This result further implies that our algorithm’s error probabilities are indeed almost optimal.*

Theorem 2.3 (No-Free Lunch Theorem). *For any given $\epsilon, \delta \in (0, 1)$, any algorithm using finitely many samples, will yield similarity approximations $f_{\cdot, \cdot, \cdot}$ with $\Pr[|f_{\cdot, \cdot, \cdot}(x, y) - \sigma_{\cdot, \cdot, \cdot}(x, y)| = \omega(\epsilon)] = \Omega(\max\{p_{\cdot}(\epsilon, \delta), p_{\cdot'}(\epsilon, \delta)\} - \epsilon)$; the probability is over the independent choices $x \sim \mathcal{D}_{\cdot}$ and $y \sim \mathcal{D}_{\cdot'}$ as well as any potential internal randomness of the algorithm.*

In plain English, Theorem 2.3 says that any algorithm using a finite amount of samples, can achieve ϵ -accuracy with a probability that is necessarily at most $1 - \max\{p_{\cdot}(\epsilon, \delta), p_{\cdot'}(\epsilon, \delta)\} + \epsilon$, i.e., if $\max\{p_{\cdot}(\epsilon, \delta), p_{\cdot'}(\epsilon, \delta)\}$ is large, learning is impossible.

Moving on, we focus on minimizing the oracle queries. By carefully modifying the earlier simple algorithm, we obtain a new more intricate algorithm with the following guarantees:

Theorem 2.4. *For any given parameters $\epsilon, \delta \in (0, 1)$, the query-minimizing algorithm produces similarity approximation functions $f_{\cdot, \cdot, \cdot}$ for every ℓ and ℓ' , such that:*

$$\Pr[\text{Error}_{(\cdot, \cdot, \cdot)}] = O(\delta + p_{\cdot}(\epsilon, \delta) + p_{\cdot'}(\epsilon, \delta))$$

$\Pr[\text{Error}_{(\cdot, \cdot, \cdot)}]$ is as defined in Theorem 2.1. Let $N = \frac{1}{2} \log \frac{1}{2}$. The algorithm requires N samples from each group, and the number of oracle queries used is at most

$$\begin{array}{ccc} \times & & \times \\ & Q_{\cdot} & \\ & \in [] & \cdot' \in [] \end{array}$$

where $Q_{\cdot} \leq N$ and $\mathbb{E}[Q_{\cdot}] \leq \frac{1}{2} + p_{\cdot}(\epsilon, \delta)N$ for each ℓ .

At first, the confidence, accuracy and sample complexity guarantees of the new algorithm are the same as those of the simpler one described in Theorem 2.1. Furthermore, because $Q_{\cdot} \leq N$, the

queries of the improved algorithm are at most $(1/\epsilon)N$, which is exactly the number of queries in our earlier simple algorithm. However, the smaller the values ϵ are, the fewer queries in expectation. Our experimental results indeed confirm that the improved algorithm always leads to a significant decrease in the used queries.

Our final theoretical result involves a lower bound on the number of queries required for learning.

Theorem 2.5. For all $\epsilon \in (0, 1)$, any learning algorithm producing similarity approximation functions $f: \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$ with $\Pr_{x \sim D_1, y \sim D_2} |f(x, y) - \phi(x, y)| \leq \epsilon$ needs $\Omega(1/\epsilon^2)$ queries.

Combining Theorems 2.4 and 2.5 implies that when ϵ and δ are negligible, i.e. $\epsilon, \delta \rightarrow 0$, the expected queries of the Theorem 2.4 algorithm are asymptotically optimal.

Finally, Section 5 contains our experimental evaluation, where through a large suite of simulations we validate our theoretical findings.

3 Related work

Metric learning is a very well-studied area [Bellet et al., 2013, Kulis, 2013, Moutaris et al., 2017, Suárez-Díaz et al., 2018]. There is also an extensive amount of work on using human feedback for learning metrics in specific tasks, e.g., image similarity and low-dimensional embeddings [Frome et al., 2007, Jamieson and Nowak, 2011, Tamuz et al., 2011, van der Maaten and Weinberger, 2012, Wilber et al., 2014]. However, since these works are either tied to specific applications or specific metrics, they are only distantly related to ours.

Our model is more closely related to the literature on trying to learn the similarity function from the fairness definition of Dwork et al. [2012]. This concept of fairness requires treating similar individuals similarly. Thus, it needs access to a function that returns a non-negative value for any pair of individuals, and this value corresponds to how similar the individuals are. Specifically, the smaller the value the more similar the elements that are compared. Even though the fairness definition of Dwork et al. [2012] is very elegant and intuitive, the main obstacle for adopting it in practice is the inability to easily compute or access the crucial similarity function. To our knowledge, the only papers that attempt to learn this similarity function using expert oracles like us, are Ilvento [2019], Mukherjee et al. [2020] and Wang et al. [2019]. Ilvento [2019] addresses the scenario of learning a general metric function, and gives theoretical PAC guarantees. Mukherjee et al. [2020] give theoretical guarantees for learning similarity functions that are only of a specific Mahalanobis form. Wang et al. [2019] simply provide empirical results. The first difference between our model and these papers is that unlike us, they do not consider elements coming from multiple distributions. However, the most important difference is that these works only learn metric functions. In our case the collection of similarity values (from all D_1 and D_2) does not necessarily yield a complete metric space see the discussion in Section 2. Hence, our problem addresses more general functions.

Regarding the difficulty in computing similarity between members of different groups, we are only aware of a brief result by Dwork et al. [2012]. In particular, given a metric over the whole feature space, they mention that it can only be trusted for comparisons between elements of the same group, and not for across-groups comparisons. In order to achieve the latter for groups, they find a new similarity function that approximates, while minimizing the Earthmover distance between the distributions $D_1; D_2$. This is completely different from our work, since here we assume the existence of across-groups similarity values, which we eventually want to learn. On the other hand, the approach of Dwork et al. [2012] can be seen as an optimization problem, where the across-groups similarity values need to be computed in a way that minimizes some objective. Also, unlike our model, this optimization approach has a serious limitation, and that is requiring D_1, D_2 to be explicitly known (recall that here we only need samples from these distributions).

Finally, since similarity as distance is difficult to compute in practice, there has been a line of research that defines similarity using simpler, yet less expressive structures. Examples include similarity lists Chakrabarti et al. [2022], graphs Lahoti et al. [2019] and ordinal relationships Jung et al. [2019].

4 Theoretical results

All proofs for this section can be found in the supplementary material

4.1 A simple learning algorithm

Given any confidence and accuracy parameters $\epsilon, \delta \in (0, 1)$ respectively, our approach is summarized as follows. At first, for every group \mathcal{S}_j we need a set of samples that are chosen i.i.d. according to \mathcal{D}_j , such that $|\mathcal{S}_j| = \frac{1}{\epsilon} \log \frac{1}{\delta}$. Then, for every distinct \mathcal{S}_j and $\mathcal{S}_{j'}$, and for all $x \in \mathcal{S}_j$ and $y \in \mathcal{S}_{j'}$, we ask the expert oracle for the true similarity value $d_0(x; y)$. The next observation follows trivially.

Observation 4.1. The algorithm uses $\log \frac{1}{\epsilon}$ samples, and $\binom{L}{2} \log^2 \frac{1}{\delta}$ queries to the oracle.

Suppose now that we need to compare $x \in \mathcal{D}_j$ and $y \in \mathcal{D}_{j'}$. Our high level idea is that the properties M_1 and M_2 of d_0 (see Section 2), will actually allow us to use the closest element to x in \mathcal{S}_j and the closest element in $\mathcal{S}_{j'}$ as proxies. Thus, let $\hat{x} = \arg \min_{x' \in \mathcal{S}_j} d_0(x; x')$ and $\hat{y} = \arg \min_{y' \in \mathcal{S}_{j'}} d_0(y; y')$. The algorithm then sets

$$f_{\epsilon, \delta}(x; y) := d_0(\hat{x}; \hat{y})$$

where $d_0(\hat{x}; \hat{y})$ is known from the earlier queries.

Before we proceed with the analysis of the algorithm, we need to recall some notation which was introduced in Section 2.1. Consider any group \mathcal{D}_j with $\Pr_{x \in \mathcal{D}_j} [d_0(x; x^0) > \epsilon] < \delta$ is called an (ϵ, δ) -rare element of \mathcal{D}_j , and also $\rho(\epsilon, \delta) := \Pr_{x \in \mathcal{D}_j} [x \text{ is } (\epsilon, \delta)\text{-rare for } \mathcal{D}_j]$.

Theorem 4.2. For any given parameters $\epsilon, \delta \in (0, 1)$, the simple algorithm produces similarity approximation functions $f_{\epsilon, \delta}$ for every \mathcal{S}_j and $\mathcal{S}_{j'}$, such that

$$\Pr[\text{Error}_{(\epsilon, \delta)}] = O(\epsilon + \rho(\epsilon, \delta) + \rho_0(\epsilon, \delta))$$

where $\Pr[\text{Error}_{(\epsilon, \delta)}]$ is as in Theorem 2.1.

Observation 4.1 and Theorem 4.2 directly yield Theorem 2.1.

A potential criticism of the algorithm presented here, is that its error probabilities depend on ρ . However, Theorem 2.3 shows that such a dependence is unavoidable (see appendix for proof).

4.2 Optimizing the number of expert queries

Here we modify the earlier algorithm in a way that improves the number of queries used. The idea behind this improvement is the following. Given the sets of samples \mathcal{S}_j instead of asking the oracle for all possible similarity values $d_0(x; y)$ for every \mathcal{S}_j and every $x \in \mathcal{S}_j$ and $y \in \mathcal{S}_{j'}$, we would rather choose a set \mathcal{R}_j of representative elements for each group \mathcal{S}_j . Then, we would ask the oracle for the values $d_0(x; y)$ for every \mathcal{S}_j , but this time only for every $x \in \mathcal{R}_j$ and $y \in \mathcal{R}_{j'}$. The choice of the representatives is inspired by the center algorithm of Hochbaum and Shmoys [1985]. Intuitively, the representatives of group \mathcal{S}_j will serve as similarity proxies for the elements of \mathcal{S}_j , such that each $x \in \mathcal{S}_j$ is assigned to a nearby $r \in \mathcal{R}_j$ via a mapping function $\pi_j : \mathcal{S}_j \rightarrow \mathcal{R}_j$. Hence, if $d_0(x; \pi_j(x))$ is small enough x and $r \in \mathcal{R}_j$ are highly similar, and thus r acts as a good approximation of x . The full details for the construction of \mathcal{R}_j, π_j are presented in Algorithm 1.

Suppose now that we need to compare some $x \in \mathcal{D}_j$ and $y \in \mathcal{D}_{j'}$. Our approach will be almost identical to that of Section 4.1. Once again, let $\hat{x} = \arg \min_{x' \in \mathcal{S}_j} d_0(x; x')$ and $\hat{y} = \arg \min_{y' \in \mathcal{S}_{j'}} d_0(y; y')$. However, unlike the simple algorithm of Section 4.1 that directly uses \hat{x} and \hat{y} , the more intricate algorithm here will rather use their proxies $\pi_j(\hat{x})$ and $\pi_{j'}(\hat{y})$. Our prediction will then be

$$f_{\epsilon, \delta}(x; y) := d_0(\pi_j(\hat{x}); \pi_{j'}(\hat{y}))$$

where $d_0(\pi_j(\hat{x}); \pi_{j'}(\hat{y}))$ is known from the earlier queries.

Theorem 4.3. For any given parameters $\epsilon, \delta \in (0, 1)$, the new query optimization algorithm produces similarity approximation functions $f_{\epsilon, \delta}$ for every \mathcal{S}_j and $\mathcal{S}_{j'}$, such that

$$\Pr[\text{Error}_{(\epsilon, \delta)}] = O(\epsilon + \rho(\epsilon, \delta) + \rho_0(\epsilon, \delta))$$

where $\Pr[\text{Error}_{(\epsilon, \delta)}]$ is as in Theorem 2.1.

Since the number of samples used by the algorithm is easily seen to be $\frac{1}{\epsilon}$, the only thing left in order to prove Theorem 2.4 is analyzing the number of oracle queries. To that end, for every group \mathcal{S}_j with its sampled set \mathcal{S}_j , we define the following Set Cover problem.

Algorithm 1 Training Phase

Accuracy and confidence parameters. For every group $\mathcal{G} \in \mathcal{G}$, a set S of i.i.d. samples chosen according to D , such that $|S| = \frac{1}{\epsilon} \log \frac{1}{\delta}$.

```

1: for each  $\mathcal{G} \in \mathcal{G}$  do
2:    $H_x := \{x \in S : d(x; x^0) \leq \epsilon\}$  for each  $x \in S$ .
3:    $U = S, R = \emptyset$ ; and  $r(x) = x$  for each  $x \in S$ .
4:   while  $U \neq \emptyset$ ; do
5:     Choose an arbitrary  $x \in U$ .
6:     Set  $R = R \cup \{x\}$ .
7:      $W_x := \{x^0 \in U : H_x \setminus H_{x^0} \neq \emptyset\}$ .
8:      $r(x^0) = x$  for every  $x^0 \in W_x$ .
9:      $U = U \setminus W_x$ .
10:  end while
11: end for
12: For every distinct  $\mathcal{G}^0$  and  $\mathcal{G}$ , and for every  $x \in R$  and  $y \in R^0$ , ask the oracle for  $\phi(x; y)$  and store all these values.
13: For every  $\mathcal{G}$ , return the set  $R$  and the function  $r$ .

```

Definition 4.4. Let $H_x := \{x^0 \in S : d(x; x^0) \leq \epsilon\}$ for all $x \in S$. Find $C \subseteq S$ minimizing $|C|$, with $\bigcup_{x \in C} H_x = S$. We use OPT to denote the optimal value of this problem. Using standard terminology, we say S is covered by C if $x \in \bigcup_{x^0 \in C} H_{x^0}$, and C is feasible if it covers all $x \in S$.

Lemma 4.5. For every $\mathcal{G} \in \mathcal{G}$ we have $|R| \leq OPT$.

Lemma 4.6. Let $N = \frac{1}{\epsilon} \log \frac{1}{\delta}$. For each group $\mathcal{G} \in \mathcal{G}$ we have $OPT \leq N$ with probability 1, and $E[|OPT|] \leq \frac{1}{\epsilon} \log \frac{1}{\delta} + \epsilon N$. The randomness here is over the samples.

The proof of Theorem 2.4 follows since all pairwise queries for the elements in the sets

$$\sum_{\mathcal{G} \in \mathcal{G}} \sum_{x \in R} \sum_{x^0 \in R^0} \sum_{x^0 \in W_x} OPT_{\mathcal{G}} \leq OPT_{\mathcal{G}} \quad (1)$$

where the inequality follows from Lemma 4.5. Combining (1) and Lemma 4.6 proves Theorem 2.4.

Remark 4.7. The factor 8 in the definition of H_x at line 2 of Algorithm 1 is arbitrary. Actually, any factor $\epsilon = O(1)$ would yield the same asymptotic guarantees, with any changes in accuracy and queries being only of $\Theta(1)$ order of magnitude. Specifically, the smaller ϵ , the better the achieved accuracy and the more queries we are using.

Finally, we are interested in lower bounds on the queries required for learning. Thus, we give Theorem 2.5, showing that any algorithm with accuracy ϵ and confidence $1 - \delta$ needs $\Omega(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$ queries.

This immediately leads to the following corollary characterizing the optimality of our algorithm.

Corollary 4.8. When for every $\mathcal{G} \in \mathcal{G}$ the value $\epsilon(\mathcal{G})$ is arbitrarily close to 0, the algorithm presented in this section achieves an expected number of oracle queries that is asymptotically optimal.

5 Experimental evaluation

We implemented all algorithms in Python 3.10.6 and ran our experiments on a personal laptop with Intel(R) Core(TM) i7-7500U CPU @ 2.70GHz 2.90 GHz and 16.0 GB memory.

Algorithms: We implemented the simple algorithm from Section 4.1, and the more intricate algorithm of Section 4.2. We refer to the former as `SIMPLE`, and to the latter as `CLUSTER`. For the training phase of `CLUSTER`, we set the dilation factor at line 2 of Algorithm 1 to 1 instead of 8. The reason for this, is that a minimal experimental investigation revealed that this choice leads to a good balance between accuracy guarantees and oracle queries. As explained in Section 3, neither the existing similarity learning algorithms [Ilvento, 2019, Mukherjee et al., 2020, Wang et al., 2019] nor the Earthmover minimization approach of Dwork et al. [2012] address the problem of finding similarity values for heterogeneous data. Furthermore, if the across-groups functions are not metric, then no approach from the metric learning literature can be used. Hence, as baselines we used three general regression

models; an MLP, a Random Forest regressor (RF) and an XGBoost regressor (XGB). The MLP uses 4 hidden layers of 32 relu activation nodes, and both RF and XGB use 200 estimators.

Number of demographic groups: All our experiments are performed for two groups, i.e., 2. The following reasons justify this decision. At first, this case captures the essence of our algorithmic results; the > 2 case can be viewed as running the algorithm for 2 multiple times, one for each pair of groups. Secondly, as the theoretical guarantees suggest, the achieved confidence and accuracy of our algorithms are completely independent of

Similarity functions: In all our experiments the feature space \mathcal{X} is where $2 \leq N$ is case-specific. In line with our motivation which assumes that the intra-group similarity functions are simple, we define d_1 and d_2 to be the Euclidean distance. Specifically, for $\mathcal{X} \subseteq \mathbb{R}^d$, the similarity between any $x, y \in \mathcal{X}$ is given by

$$d(x; y) = \sqrt{\sum_{i=1}^d (x_i - y_i)^2}$$

For the across-groups similarities, we aim for a difficult-to-learn metric function; we purposefully chose a non-trivial function in order to challenge both our algorithms and the baselines. Namely, for any $x \in \mathcal{X}_1$ and $y \in \mathcal{X}_2$, we assume

$$d(x; y) = \sqrt{\sum_{i=1}^d \left(\frac{x_i}{\alpha_i} - \frac{y_i}{\beta_i} \right)^2 + \sum_{i=1}^d \alpha_i^3 \beta_i^3} \quad (2)$$

where the vectors $\alpha; \beta \in \mathbb{R}^d$ are basically the hidden parameters to be learned (of course none of the learners we use has any insight on the specific structure of

The function implicitly adopts a paradigm of feature importance [Niño-Adan et al., 2021]. Specifically, when x and y are to be compared, their features are scaled accordingly by the expert using the parameters $\alpha; \beta$, while some offsetting via might also be necessary. For example, in the college admissions use-case, certain features may have to be properly adjusted (increase a feature for a non-privileged student and decrease it for the privileged one). In the end, the similarity is calculated in an ℓ_3 -like manner. To see why is not a metric and why it satisfies the necessary properties and M_2 from Section 2, refer to the last theorem in the Appendix.

In each experiment we choose all $\alpha_i; \beta_i$ independently. The values $\alpha_i; \beta_i$ are chosen uniformly at random from $[0; 1]$, while the $\alpha_i^3; \beta_i^3$ are chosen uniformly at random from $[0; 0.01; 0.01]$. Obviously, the algorithms do not have access to the vectors $\alpha; \beta$, which are only used to simulate the oracle and compare our predictions with the corresponding true values.

Datasets: We used 2 datasets from the UCI ML Repository, namely Adult (48,842 points - 14 features) [Kohavi, 1996] and Credit Card Default (30,000 points - 23 features) [Yeh and Lien, 2009], and the publicly available Give Me Some Credit dataset (150,000 points - 11 features) [Credit Fusion, 2011]. We chose these datasets because this type of data is frequently used in applications of issuing credit scores, and in such cases fairness considerations are of utmost importance. For Adult, where categorical features are not encoded as integers, we assigned each category to an integer in $\{1; \dots; \# \text{ categories}\}$, and this integer is used in place of the category in the feature vector [Ding et al., 2021]. Finally, in every dataset we standardized all features through a MinMax re-scaler. Due to space constraints, the figures for Adult and Give me Some Credit are moved to the supplementary material. However, the observed traits there the same as the ones shown here for Credit Card Default.

Choosing the two groups: For Credit Card Default and Adult, we defined groups based on marital status. Specifically, the first group corresponds to points that are married individuals, and the second to points that are not married (singles, divorced and widowed are merged together). In Give Me Some Credit, we partition individuals into two groups based on whether or not they have dependents.

Choosing the accuracy parameter: To use a meaningful value for ϵ we need to know the order of magnitude of $d(x; y)$. Thus, we calculated the value of $d(x; y)$ over 10,000 trials, where the randomness was of multiple factors, i.e., the random choices for the $\alpha; \beta$, and the sampling of x and y . Figure 1 shows histograms for the empirical frequency of $d(x; y)$ over the 10,000 runs. In addition, in those trials the minimum value of $d(x; y)$ observed was 0.1149 for Credit Card Default, 0.1155 for Adult, and 0.0132 for Give Me Some Credit. Thus, aiming for an accuracy parameter that is at least an order of magnitude smaller than the value to be learned, we choose $\epsilon = 0.01$ for Credit Card Default and Adult, and $\epsilon = 0.001$ for Give Me Some Credit.

(a) Credit Card Default

(b) Adult

(c) Give Me Some Credit

Figure 1: Frequency counts for $(x; y)$

Table 1: Error statistics for Credit Card Default

Algorithm	Average Relative Error %	SD of Relative Error %	Average (Absolute Error)/	SD of (Absolute Error)/
NAIVE	1.558	3.410	0.636	1.633
CLUSTER	1.593	3.391	0.646	1.626
MLP	3.671	4.275	1.465	1.599
RF	3.237	4.813	1.306	2.318
XGB	3.121	4.606	1.273	1.869

Confidence and number of samples: All our experiments are performed with $\epsilon = 0.001$. In NAIVE and CLUSTER, we follow our theoretical results and sample $n = \frac{1}{\epsilon} \log \frac{1}{\delta}$ points from each group. Choosing the training samples for the baselines is a bit more tricky. For these regression tasks a training point is a tuple $(x; y; (x; y))$. In other words, these tasks do not distinguish between queries and samples. To be as fair as possible when comparing against our algorithms, we provide the baselines with $N + Q$ training points of the form $(x; y; (x; y))$, where Q is the maximum number of queries used in any of our algorithms.

Testing: We test our algorithms over $r = 1000$ trials, where each trial consists of independently sampling two elements $x; y$, one for each group, and then inputting those to the predictors. We are interested in two metrics. The first is the relative error percentage $\frac{|p - t|}{t}$ where p is the prediction for element x and t is their true similarity value, the relative error percentage $\frac{|p - t|}{p}$ where p is the prediction for two elements and t is their true similarity value, this metric is $\frac{|p - t|}{p}$. We are interested in the latter metric because our theoretical guarantees are of the form $\frac{|p - t|}{p} = O(\epsilon)$.

Table 1 shows the average relative and absolute value error, together with the standard deviation for these metrics across $r = 1000$ runs. It is clear that our algorithms dominate the baselines since they exhibit smaller errors with smaller standard deviations. In addition, NAIVE appears to have a tiny edge over CLUSTER, and this is expected (the queries of CLUSTER are a subset of NAIVE's). However, as shown in Table 2, CLUSTER leads to a significant decrease in oracle queries compared to NAIVE, thus justifying its superiority. To further demonstrate the statistical behavior of the errors, we present Figure 2. This depicts the empirical CDF of each error metric through a bar plot. Specifically, the height of each bar corresponds to the fraction of test instances whose error is at most the value in the x-axis directly underneath the bar. Once again, we see that our algorithms outperform the baselines, since their corresponding bars are always higher than those of the baselines.

Disclaimer

This paper was prepared for informational purposes in part by the Artificial Intelligence Research group of JPMorganChase and its affiliates ("JP Morgan"), and is not a product of the Research Department of JPMorganChase. JPMorganChase makes no representation and warranty whatsoever and

Table 2: Percent of decrease in queries when using CLUSTER instead of NAIVE

Credit Card Default	Adult	Give Me Some Credit
81.01%	80.40%	84.87%

(a) Relative Error Percentage

(b) (Absolute Error) /

Figure 2: Empirical CDF for Credit Card Default

disclaims all liability, for the completeness, accuracy or reliability of the information contained herein. This document is not intended as investment research or investment advice, or a recommendation, offer or solicitation for the purchase or sale of any security, financial instrument, financial product or service, or to be used in any way for evaluating the merits of participating in any transaction, and shall not constitute a solicitation under any jurisdiction or to any person, if such solicitation under such jurisdiction or to such person would be unlawful.

References

- A. Bellet, A. Habrard, and M. Sebban. A Survey on Metric Learning for Feature Vectors and Structured Data. Research report, Laboratoire Hubert Curien UMR 5516, 2013. <https://hal.inria.fr/hal-01666935>.
- D. Chakrabarti, J. P. Dickerson, S. A. Esmaeili, A. Srinivasan, and L. Tsepenekas. A new notion of individually fair clustering: ϵ -equitable k -center. In G. Camps-Valls, F. J. R. Ruiz, and I. Valera, editors, Proceedings of The 25th International Conference on Artificial Intelligence and Statistics volume 151 of Proceedings of Machine Learning Research, pages 6387–6408. PMLR, 28–30 Mar 2022.
- W. C. Credit Fusion. Give me some credit, 2011. <https://kaggle.com/competitions/GiveMeSomeCredit>.
- F. Ding, M. Hardt, J. Miller, and L. Schmidt. Retiring adult: New datasets for fair machine learning. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, Advances in Neural Information Processing Systems volume 34, pages 6478–6490. Curran Associates, Inc., 2021.
- C. Dwork, M. Hardt, T. Pitassi, O. Reingold, and R. Zemel. Fairness through awareness. In Proceedings of the 3rd Innovations in Theoretical Computer Science Conference, 2012.
- A. Frome, Y. Singer, F. Sha, and J. Malik. Learning globally-consistent local distance functions for shape-based image retrieval and classification. 2007 IEEE 11th International Conference on Computer Vision, pages 1–8, 2007. doi: 10.1109/ICCV.2007.4408839.
- D. S. Hochbaum and D. B. Shmoys. A best possible heuristic for the k -center problem. *Mathematics of Operations Research* 10(2):180–184, 1985. ISSN 0364765X, 15265471. <https://www.jstor.org/stable/3689371>.
- C. Ilvento. Metric learning for individual fairness, 2019. <https://arxiv.org/abs/1906.00250>.
- K. G. Jamieson and R. D. Nowak. Low-dimensional embedding using adaptively selected ordinal data. In 2011 49th Annual Allerton Conference on Communication, Control, and Computing (Allerton), pages 1077–1084, 2011. doi: 10.1109/Allerton.2011.6120287.

- C. Jung, M. Kearns, S. Neel, A. Roth, L. Stapleton, and Z. S. Wu. An algorithmic framework for fairness elicitation, 2019. URL <https://arxiv.org/abs/1905.10660> .
- M. P. Kim, O. Reingold, and G. N. Rothblum. Fairness through computationally-bounded awareness. In Proceedings of the 32nd International Conference on Neural Information Processing Systems NIPS'18, page 4847–4857, Red Hook, NY, USA, 2018. Curran Associates Inc.
- R. Kohavi. Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid. Proceedings of the Second International Conference on Knowledge Discovery and Data Mining, 1996, page 202–207. AAAI Press, 1996.
- B. Kulis. Metric learning: A survey. 2013.
- P. Lahoti, K. P. Gummadi, and G. Weikum. Operationalizing individual fairness with pairwise fair representations. Proc. VLDB Endow. 13(4):506–518, dec 2019. ISSN 2150-8097. doi: 10.14778/3372716.3372723. URL <https://doi.org/10.14778/3372716.3372723> .
- P. Moutaris, M. Leng, and I. A. Kakadiaris. An overview and empirical comparison of distance metric learning methods. IEEE Transactions on Cybernetics, 47(3):612–625, 2017. doi: 10.1109/TCYB.2016.2521767.
- D. Mukherjee, M. Yurochkin, M. Banerjee, and Y. Sun. Two simple ways to learn individual fairness metrics from data. In Proceedings of the 37th International Conference on Machine Learning ICML'20. JMLR.org, 2020.
- I. Niño-Adan, D. Manjarres, I. Landa-Torres, and E. Portillo. Feature weighting methods: A review. Expert Systems with Applications, 184:115424, 2021. ISSN 0957-4174.
- J. L. Suárez-Díaz, S. García, and F. Herrera. A tutorial on distance metric learning: Mathematical foundations, algorithms, experimental analysis, prospects and challenges (with appendices on mathematical background and detailed algorithms explanation), 2018. URL <https://arxiv.org/abs/1812.05944> .
- O. Tamuz, C. Liu, S. Belongie, O. Shamir, and A. T. Kalai. Adaptively learning the crowd kernel. In Proceedings of the 28th International Conference on International Conference on Machine Learning ICML'11, page 673–680. Omnipress, 2011. ISBN 9781450306195.
- L. van der Maaten and K. Weinberger. Stochastic triplet embedding. 2012 IEEE International Workshop on Machine Learning for Signal Processing, pages 1–6, 2012. doi: 10.1109/MLSP.2012.6349720.
- H. Wang, N. Grgic-Hlaca, P. Lahoti, K. P. Gummadi, and A. Weller. An empirical study on learning fairness metrics for compas data with human supervision, 2019. URL <https://arxiv.org/abs/1910.10255> .
- M. Wilber, I. Kwak, and S. Belongie. Cost-effective hits for relative similarity comparison. Proceedings of the AAAI Conference on Human Computation and Crowdsourcing, 2014, page 227–233, Sep. 2014. URL <https://ojs.aaai.org/index.php/HCOMP/article/view/13152> .
- I. Yeh and C.-H. Lien. The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. Expert Systems with Applications, 36:2473–2480, 03 2009. doi: 10.1016/j.eswa.2007.12.020.
- G. Yona and G. Rothblum. Probably approximately metric-fair learning. International Conference on Machine Learning, pages 5680–5688. PMLR, 2018.

A Missing proofs

Proof of Theorem 4.2. For two distinct groups and ρ , consider what will happen when we are asked to compare some $x \in D$ and $y \in D$. Properties M_1 and M_2 of ρ imply

$$\begin{aligned} Q &= \rho(x; y) \leq P; \text{ where} \\ P &:= d(x; (x)) + \rho(x; (y)) + d(y; (y)) \\ Q &:= \rho(x; (y)) + d(x; (x)) + d(y; (y)) \end{aligned}$$

Note that when $d(x; (x)) \geq 3$ and $d(y; (y)) \geq 3$, the above inequalities and the definition of $\rho(x; y)$ yield $\rho(x; y) \leq \rho(x; y) - 6$. Thus, we just need upper bounds for $A := \Pr_{x \in D} [8x \in S : d(x; x^0) > 3]$ and $B := \Pr_{y \in D} [8y \in S : d(y; y^0) > 3]$, since the previous analysis and a union bound give $\text{Error}(\rho) \leq A + B$. In what follows we present an upper bound for A . The same analysis gives an identical bound for B .

Before we proceed to the rest of the proof, we have to provide an existential construction. For the sake of simplicity we will be using the term dense for elements that are not $(;)$ -rare. For every $x \in D$ that is dense, we define $B_x := \{x^0 \in D : d(x; x^0) \leq \epsilon\}$. Observe that the definition of dense elements implies $\Pr_{x^0 \in D} [x^0 \in B_x] \geq \epsilon$ for every dense x . Next, consider the following process. We start with an empty set $S = \emptyset$, and we assume that all dense elements are unmarked. Then, we choose an arbitrary unmarked dense element and we place it in the set R . Further, for every dense $x \in D$ that is unmarked and has $B_x \cap B_{x^0} \neq \emptyset$, we mark x^0 and set $(x^0) = x$. Here the function ρ maps dense elements to elements of D . We continue this picking process until all dense elements have been marked. Since $B_{z^0} = \emptyset$ for any two $z; z^0 \in R$ and $\Pr_{x^0 \in D} [x^0 \in B_z] \geq \epsilon$ for $z \in R$, we have $|R| \leq \frac{1}{\epsilon}$. Also, for every dense x we have $d(x; (x)) \leq 2$ due to $B_x \cap B_{(x)} \neq \emptyset$.

Now we are ready to upper bound A .

$$\begin{aligned} C &:= \Pr_{S \subseteq X} [8x \in S : d(x; x^0) > 3 \wedge x \text{ is } (;)\text{-rare}] \\ &= \Pr_{x \in D} [x \text{ is } (;)\text{-rare}] = p((;)) \\ D &:= \Pr_{S \subseteq X} [8x \in S : d(x; x^0) > 3 \wedge x \text{ is dense}] \\ &= \Pr_{r \in R} [9r \in S : B_r \cap S = \emptyset] \\ &= \sum_{r \in R} \Pr_{S \subseteq X} [B_r \cap S = \emptyset] = \sum_{r \in R} (1 - \epsilon)^{|B_r|} \leq \sum_{r \in R} e^{-\epsilon |B_r|} \end{aligned}$$

The upper bound for C is trivial. We next explain the computations for D . For the transition between the first and the second line we use a proof by contradiction. Hence, suppose that $\epsilon > 0$; for every $r \in R$, and let i_r denote an arbitrary element of B_r . Then, for any dense element $x \in D$ we have $d(x; (x)) = d(x; i_{(x)}) + d(i_{(x)}; (x)) \leq 2 + \epsilon = 3$. Back to the computations for D , to get the third line we simply used a union bound. To get from the third to the fourth line, we used the definition of R as a dense element, which implies that the probability of sampling any element of R in one try is at least ϵ . The final bound is a result of numerical calculations using $|R| \leq \frac{1}{\epsilon}$ and $\sum_{j \in S} j = \frac{1}{2}$.

To conclude the proof, observe that $A = C + D$, and using a similar reasoning as the one in upper-bounding A we also get $B = p((;)) + \dots$ \square

Proof of Theorem 2.3 Given any ϵ , consider the following instance of the problem. We have two groups represented by the distributions D_1 and D_2 . For the first group we have only one element belonging to it, and let that element be x . In other words, every time we draw an element from D_1 that element turns out to be x , i.e., $\Pr_{x^0 \in D_1} [x^0 = x] = 1$. For the second group we have that every $y \in D_2$ appears with probability $\frac{1}{D_2}$, and D_2 is a huge constant $\gg 0$, with $\frac{1}{D_2} \geq \epsilon$.

Now we define all similarity values. At first, the similarity function for D_1 will trivially be $d_1(x; x) = 0$. For the second group, for every distinct $y^0 \in D_2$ we define $d_2(y; y^0) = 1$. Obviously, for every $y \in D_2$ we set $d_2(y; y) = 0$. Observe that d_1 and d_2 are metric functions for their respective groups. As for the across-groups similarities, each $d(x; y)$ for $y \in D_2$ is chosen independently, and it is

drawn uniformly at random from $[0, 1]$. Note that this choice of satisfies the necessary metric-like properties M_1 and M_2 that were introduced in Section 2.

Further, since $\epsilon \in (0, 1)$, any D_2 will be $(\epsilon; \epsilon)$ -rare:

$$\Pr_{y^0 \in D_2} [d_2(y; y^0) \leq \epsilon] = \Pr_{y^0 \in D_2} [y^0 = y] = \frac{1}{|D_2|} < \epsilon$$

The first equality is because the only element within distance ϵ from y is y itself. The last inequality is because $\frac{1}{|D_2|} < \epsilon$. Therefore, since all elements are $(\epsilon; \epsilon)$ -rare, we have $\rho_2(\epsilon; \epsilon) = 1$.

Consider now any learning algorithm that produces an estimate function \hat{f} for any D_2 , let us try to analyze the probability of having $\hat{f}(x; y) \neq f(x; y) \forall (x; y) \in S_2$. At first, note that when $y \in S_2$, we can always get the exact value $f(x; y)$, since x will always be in S_1 . The probability of having $y \notin S_2$ is $1 - \frac{1}{|D_2|}$, where N is the number of used samples. Since we have control over $|D_2|$ when constructing this instance, we can always set it to a large enough value that will give $1 - \frac{1}{|D_2|} \leq \epsilon$; note that this is possible because $\frac{1}{|D_2|}$ is decreasing in $|D_2|$ and $\lim_{|D_2| \rightarrow \infty} \frac{1}{|D_2|} = 0$. Hence,

$$\Pr[\hat{f}(x; y) \neq f(x; y) \forall (x; y) \in S_2] = (1 - \frac{1}{|D_2|})^N \Pr[\hat{f}(x; y) \neq f(x; y) \forall (x; y) \in S_2] + \epsilon \quad (3)$$

When y will not be among the samples, the algorithm needs to learn $f(x; y)$ via some other value $f(x; y^0)$, for y^0 being a sampled element of the second group. However, due to the construction of the values $(x; y)$ and $(x; y^0)$ are independent. This means that knowledge of $f(x; y^0)$ (with $y \neq y^0$) provides no information at all on $f(x; y)$. Thus, the best any algorithm can do is guess $f(x; y)$ uniformly at random from $[0, 1]$. This yields $\Pr[\hat{f}(x; y) \neq f(x; y) \forall (x; y) \in S_2] = 1 - \Pr[\hat{f}(x; y) = f(x; y) \forall (x; y) \in S_2] = 1 - \frac{1}{|D_2|} = \epsilon$. Combining this with (3) gives the desired result. \square

Proof of Theorem 4.3. For two distinct groups \mathcal{X}^0 and \mathcal{X}^1 , consider comparing some D^0 and D^1 . To begin with, let us assume that $d(\mathcal{X}^0; \mathcal{X}^1) \geq \epsilon$ and $d_0(\mathcal{Y}^0; \mathcal{Y}^1) \geq \epsilon$. Furthermore, the execution of the algorithm implies $\mathcal{H}_{r^0}(\mathcal{X}^0) \cap \mathcal{H}_{r^1}(\mathcal{X}^1) = \emptyset$, and thus the triangle inequality and the definitions of the sets $\mathcal{H}_{r^0}(\mathcal{X}^0); \mathcal{H}_{r^1}(\mathcal{X}^1)$ give $d(\mathcal{X}^0; \mathcal{X}^1) \geq \epsilon$. Similarly $d_0(\mathcal{Y}^0; \mathcal{Y}^1) \geq \epsilon$. Eventually:

$$\begin{aligned} d(\mathcal{X}^0; \mathcal{X}^1) &\geq d(\mathcal{X}^0; \mathcal{H}_{r^0}(\mathcal{X}^0)) + d(\mathcal{H}_{r^0}(\mathcal{X}^0); \mathcal{X}^1) \geq \epsilon \\ d_0(\mathcal{Y}^0; \mathcal{Y}^1) &\geq d_0(\mathcal{Y}^0; \mathcal{H}_{r^0}(\mathcal{Y}^0)) + d_0(\mathcal{H}_{r^0}(\mathcal{Y}^0); \mathcal{Y}^1) \geq \epsilon \end{aligned} \quad (19)$$

For notational convenience, let $A := d(\mathcal{X}^0; \mathcal{X}^1)$ and $B := d_0(\mathcal{Y}^0; \mathcal{Y}^1)$. Then, the metric properties M_1 and M_2 of ρ_0 and the definition of $\rho_0(x; y)$ yield

$$\rho_0(x; y) \leq f_0(x; y) \leq A + B \quad (38)$$

Overall, we proved that when $d(\mathcal{X}^0; \mathcal{X}^1) \geq \epsilon$ and $d_0(\mathcal{Y}^0; \mathcal{Y}^1) \geq \epsilon$, we have $f_0(x; y) \leq A + B$. Finally, as shown in the proof of Theorem 4.2, the probability of not having $d(\mathcal{X}^0; \mathcal{X}^1) \geq \epsilon$ and $d_0(\mathcal{Y}^0; \mathcal{Y}^1) \geq \epsilon$, is at most $2\epsilon + \epsilon^2$. \square

Proof of Lemma 4.5. Consider a group \mathcal{R} , and let C be its optimal solution for the problem of Definition 4.4. We first claim that each H_c with $c \in C$ contains at most one element of \mathcal{R} . This is due to the following. For any $c \in C$, we have $d(z; z^0) = d(z; c) + d(c; z^0) \leq \epsilon$ for all $z; z^0 \in H_c$. In addition, the construction of \mathcal{R} trivially implies $d(x; x^0) > \epsilon$ for all $x; x^0 \in \mathcal{R}$. Thus, no two elements of \mathcal{R} can be in the same H_c with $c \in C$. Finally, since Definition 4.4 requires all \mathcal{R} to be covered, we have $\sum_{c \in C} |H_c| = |\mathcal{R}| = \text{OPT}$. \square

Proof of Lemma 4.6 Consider a group \mathcal{R} . Initially, through Definition 4.4 it is clear that $\text{OPT} \leq |\mathcal{R}| = N$. For the second statement of the lemma we need to analyze OPT in a more clever way.

Recall the classification of elements $x \in \mathcal{R}$ that was first introduced in the proof of Theorem 4.2. According to this, an element can either be $(\epsilon; \epsilon)$ -rare, or dense. Now we will construct a solution to the problem of Definition 4.4 as follows.

At first, let S_{ϵ} be the set of ϵ -rare elements of \mathcal{D} . We will include all of S_{ϵ} to C , so that all ϵ -rare elements of \mathcal{D} are covered by C . Further:

$$E |S_{\epsilon}| = \epsilon N \quad (4)$$

Moving on, recall the construction shown in the proof of Theorem 4.2. According to that, there exists a set R of at most $1/\epsilon$ dense elements from \mathcal{D} , and a function ϕ that maps every dense element $x \in \mathcal{D}$ to an element $\phi(x) \in R$, such that $d(x; \phi(x)) \leq \epsilon$. Let us now define for each $x \in R$ a set $G_x := \{x^0 \in \mathcal{D} : x^0 \text{ is dense and } d(x^0) = \epsilon\}$, and note that $d(z; z^0) \leq d(z; x) + d(x; \phi(x)) + d(\phi(x); x^0) \leq 4\epsilon$ for all $z; z^0 \in G_x$. Thus, for each $x \in R$ with $G_x \cap S_{\epsilon} \neq \emptyset$, we place in C an arbitrary $x^0 \in G_x \cap S_{\epsilon}$, and that gets all of $G_x \cap S_{\epsilon}$ covered. Finally, since the sets G_x induce a partition of the dense elements of \mathcal{D} , C covers all dense elements of \mathcal{D} .

Equation (4) and $\epsilon N \leq |C| \leq \epsilon N + \epsilon N$. Also, since C is shown to be a feasible solution for problem of Definition 4.4, we get

$$\text{OPT} \leq |C| \leq \epsilon N + \epsilon N = 2\epsilon N \quad \square$$

Proof of Theorem 2.5 We are given accuracy and confidence parameters $\epsilon \in (0, 1)$ respectively. For the sake of simplifying the exposition in the proof, let us assume that n is an integer; all later arguments can be generalized in order to handle the case of n .

We construct the following problem instance. We have two groups represented by the distributions D_1 and D_2 . In addition, for both of these groups we assume that the support of the corresponding distribution contains n elements, and $\Pr_{x^0 \in D_1}[x = x^0] = 1/n$ for every $x \in D_1$ as well as $\Pr_{y^0 \in D_2}[y = y^0] = 1/n$ for every $y \in D_2$.

For every $x; x^0 \in D_1$ let $d_1(x; x^0) = 1/n$, and $d_1(x; x) = 0$ for every $x \in D_1$. Similarly, for every $y; y^0 \in D_2$ we set $d_2(y; y^0) = 1/n$, and for every $y \in D_2$ we set $d_2(y; y) = 0$. The functions d_1 and d_2 are clearly metrics. As for the across-groups similarity values, each $x \in D_1$ and $y \in D_2$ is chosen independently, and it is drawn uniformly at random from $[0, 1]$. Note that this choice of d satisfies the necessary properties $M_1; M_2$ introduced in Section 2.

In this proof we are also focusing on a more special learning model. In particular, we assume that the distributions D_1 and D_2 are known. Hence, there is no need for sampling. The only randomness here is over the random arrivals $x \in D_1$ and $y \in D_2$, where x and y are the elements that need to be compared. Obviously, the similarity function would still remain unknown to any learner. Finally, the queries required for learning in this model cannot be more than the queries required in the original model, and this is because this model is a special case of the original.

Consider now an algorithm with the error guarantees mentioned in the Theorem statement, and focus on a fixed pair $(x; y)$ with $x \in D_1$ and $y \in D_2$. If the algorithm has queried the oracle $f(x; y)$, it knows $d(x; y)$ with absolute certainty. Let us study what happens when the algorithm has not queried the oracle $f(x; y)$. In this case, because the values $(x^0; y^0)$ with $x^0 \in D_1$ and $y^0 \in D_2$ are independent, no query the algorithm has performed can provide any information about $(x; y)$. Thus, the best the algorithm can do is uniformly at random guess a value $\hat{d}(x; y)$ and return that as the estimate for $d(x; y)$. If $Q_{x; y}$ denotes the event where no query is performed for $(x; y)$, then

$$\begin{aligned} P &:= \Pr [|f(x; y) - \hat{d}(x; y)| \geq \epsilon] \\ &= \Pr [|f(x; y) - \hat{d}(x; y)| \geq \epsilon \mid Q_{x; y}] \\ &= \Pr [\text{random guess} \geq \epsilon] \\ &= 1/2 \end{aligned}$$

where the randomness comes only from the algorithm.

For the sake of contradiction, suppose the algorithm uses $(1/\epsilon^2)$ queries. Since each $(x; y)$ is equally likely to appear for a comparison, the overall error probability is

$$P = \frac{1/\epsilon^2}{1/\epsilon^2} \cdot \frac{1}{2} = \frac{1}{2} > \epsilon$$

Contradiction; the error probability was assumed to be ϵ . □

Proof of Corollary 4.8. When every $p \cdot (\epsilon, \delta)$ is very close to 0, Lemma 4.6 gives $\mathbb{E}[OPT] \leq \frac{1}{2}$. Thus, by inequality (1) the expected queries are $\frac{(-1)}{2}$. Theorem 2.5 concludes the proof. \square

Theorem A.1. The function σ defined in Equation (2) is not metric, and satisfies properties $\mathcal{M}_1, \mathcal{M}_2$ for $d_1(x, y) = d_2(x, y) = \frac{1}{3} \sum_{i \in [d]} (x_i - y_i)^2$, when $\alpha_i, \beta_i \in [0, 1]$ for all $i \in [d]$.

Proof. The easiest way to see that $\sigma(x, y)$ is not a metric, is by realizing that it does not satisfy the symmetry property and also $x = y$ does not necessarily imply $\sigma(x, y) = 0$. Both of these issues stem from the offsets θ_i . To verify that symmetry is violated let x, y be 1-dimensional, and let $x = 1, y = 2, \alpha = 1, \beta = 1, \theta = 2$. Then $\sigma(1, 2) = 1$, while $\sigma(2, 1) = 3$. Next, let $x = 1, y = 1, \alpha = 1, \beta = 1, \theta = 1$. In this example, although $x = y$ we have $\sigma(x, y) = 1 \neq 0$.

In the following we are going to show that σ satisfies \mathcal{M}_1 . The proof for \mathcal{M}_2 is identical. At first, take $x, y, z \in \mathbb{R}^d$, such that $x, z \in \mathcal{D}_1$ and $y \in \mathcal{D}_2$. Then:

$$\begin{aligned}
 \sigma(x, y) &= \frac{1}{3} \sum_{i \in [d]} |\alpha_i \cdot x_i - \beta_i \cdot y_i + \theta_i|^3 = \frac{1}{3} \sum_{i \in [d]} |(\alpha_i \cdot x_i - \alpha_i \cdot z_i) + (\alpha_i \cdot z_i - \beta_i \cdot y_i + \theta_i)|^3 \\
 &\leq \frac{1}{3} \sum_{i \in [d]} |(\alpha_i \cdot x_i - \alpha_i \cdot z_i)|^3 + \frac{1}{3} \sum_{i \in [d]} |(\alpha_i \cdot z_i - \beta_i \cdot y_i + \theta_i)|^3 \\
 &= \frac{1}{3} \sum_{i \in [d]} \alpha_i^3 |x_i - z_i|^3 + \sigma(z, y) \\
 &\leq \frac{1}{3} \sum_{i \in [d]} |x_i - z_i|^3 + \sigma(z, y) \\
 &\leq d_1(x, z) + \sigma(z, y)
 \end{aligned}$$

To get the first inequality we used the triangle inequality for ℓ_3 . The second to last inequality is because $\alpha_i^3 \leq 1$ for all i , and the last inequality is because the ℓ_3 norm of a vector is always smaller than its ℓ_2 norm. \square

B Additional experimental results

Table 3: Error statistics for Adult

Algorithm	Average Relative Error %	SD of Relative Error %	Average (Absolute Error)/ ϵ	SD of (Absolute Error)/ ϵ
NAIVE	1.319	3.381	0.847	2.382
CLUSTER	1.321	3.383	0.849	2.383
MLP	4.119	6.092	2.306	1.977
RF	3.519	6.657	2.020	2.987
XGB	2.248	5.351	1.215	1.655

Table 4: Error statistics for Give Me Some Credit

Algorithm	Average Relative Error %	SD of Relative Error %	Average (Absolute Error)/ ϵ	SD of (Absolute Error)/ ϵ
NAIVE	1.630	4.106	2.644	7.710
CLUSTER	1.645	4.118	2.648	7.665
MLP	5.819	9.073	7.588	8.750
RF	5.692	12.404	7.258	34.508
XGB	5.614	13.355	6.516	7.754

