# Bounds on the Minimax Rate for Estimating a Prior over a VC Class from Independent Learning Tasks

Liu Yang[1]([✉]), Steve Hanneke[2]([✉]), and Jaime Carbonell[3]

[1] IBM T.J. Watson Research Center, Yorktown Heights, NY, USA
{liuy,jgc}@cs.cmu.edu
[2] Princeton, NJ, USA
steve.hanneke@gmail.com
[3] Carnegie Mellon University, Pittsburgh, PA, USA

**Abstract.** We study the optimal rates of convergence for estimating a prior distribution over a VC class from a sequence of independent data sets respectively labeled by independent target functions sampled from the prior. We specifically derive upper and lower bounds on the optimal rates under a smoothness condition on the correct prior, with the number of samples per data set equal the VC dimension. These results have implications for the improvements achievable via transfer learning. We additionally extend this setting to real-valued function, where we establish consistency of an estimator for the prior, and discuss an additional application to a preference elicitation problem in algorithmic economics.

## 1 Introduction

In the *transfer learning* setting, we are presented with a sequence of learning problems, each with some respective target concept we are tasked with learning. The key question in transfer learning is how to leverage our access to past learning problems in order to improve performance on learning problems we will be presented with in the future.

Among the several proposed models for transfer learning, one particularly appealing model supposes the learning problems are independent and identically distributed, with unknown distribution, and the advantage of transfer learning then comes from the ability to estimate this shared distribution based on the data from past learning problems [2,11]. For instance, when customizing a speech recognition system to a particular speaker's voice, we might expect the first few people would need to speak many words or phrases in order for the system to accurately identify the nuances. However, after performing this for many different people, if the software has access to those past training sessions when customizing itself to a new user, it should have identified important properties of the speech patterns, such as the common patterns within each of the major dialects or accents, and other such information about the *distribution* of speech patterns within the user population. It should then be able to leverage this information to reduce the number of words or phrases the next user needs to speak in

order to train the system, for instance by first trying to identify the individual's dialect, then presenting phrases that differentiate common subpatterns within that dialect, and so forth.

In analyzing the benefits of transfer learning in such a setting, one important question to ask is how quickly we can estimate the distribution from which the learning problems are sampled. In recent work, [11] have shown that under mild conditions on the family of possible distributions, if the target concepts reside in a known VC class, then it is possible to estimate this distribution using only a bounded number of training samples per task: specifically, a number of samples equal the VC dimension. However, that work left open the question of quantifying the *rate* of convergence of the estimate, in terms of the number of tasks. This rate of convergence can have a direct impact on how much benefit we gain from transfer learning when we are faced with only a finite sequence of learning problems. As such, it is certainly desirable to derive tight characterizations of this rate of convergence.

The present work continues that of [11], bounding the rate of convergence for estimating this distribution, under a smoothness condition on the distribution. We derive a generic upper bound, which holds regardless of the VC class the target concepts reside in. The proof of this result builds on that earlier work, but requires several interesting innovations to make the rate of convergence explicit, and to dramatically improve the upper bound implicit in the proofs of those earlier results. We further derive a nontrivial lower bound that holds for certain constructed scenarios, which illustrates a lower limit on how good of a general upper bound we might hope for in results expressed only in terms of the number of tasks, the smoothness conditions, and the VC dimension.

We additionally include an extension of the results of [11] to the setting of real-valued functions, establishing consistency (at a uniform rate) for an estimator of a prior over any VC subgraph class. In addition to the application to transfer learning, analogous to the original work of [11], we also discuss an application of this result to a preference elicitation problem in algorithmic economics, in which we are tasked with allocating items to a sequence of customers to approximately maximize the customers' satisfaction, while permitted access to the customer valuation functions only via value queries.

## 2   The Setting

Let $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$ be a measurable space [7] (where $\mathcal{X}$ is called the *instance space*), and let $\mathcal{D}$ be a distribution on $\mathcal{X}$ (called the *data distribution*). Let $\mathbb{C}$ be a VC class of measurable classifiers $h : \mathcal{X} \rightarrow \{-1, +1\}$ (called the *concept space*), and denote by $d$ the VC dimension of $\mathbb{C}$ [9]. We suppose $\mathbb{C}$ is equipped with its Borel $\sigma$-algebra $\mathcal{B}$ induced by the pseudo-metric $\rho(h, g) = \mathcal{D}(\{x \in \mathcal{X} : h(x) \neq g(x)\})$. Though our results can be formulated for general $\mathcal{D}$ (with somewhat more complicated theorem statements), to simplify the statement of results we suppose $\rho$ is actually a *metric*.

For any two probability measures $\mu_1, \mu_2$ on a measurable space $(\Omega, \mathcal{F})$, define the total variation distance

$$\|\mu_1 - \mu_2\| = \sup_{A \in \mathcal{F}} \mu_1(A) - \mu_2(A).$$

For a set function $\mu$ on a *finite* measurable space $(\Omega, \mathcal{F})$, we abbreviate $\mu(\omega) = \mu(\{\omega\})$, $\forall \omega \in \Omega$. Let $\Pi_\Theta = \{\pi_\theta : \theta \in \Theta\}$ be a family of probability measures on $\mathbb{C}$ (called *priors*), where $\Theta$ is an arbitrary index set (called the *parameter space*). We suppose there exists a probability measure $\pi_0$ on $\mathbb{C}$ (the *reference measure*) such that every $\pi_\theta$ is absolutely continuous with respect to $\pi_0$, and therefore has a density function $f_\theta$ given by the Radon-Nikodym derivative $\frac{d\pi_\theta}{d\pi_0}$ [7].

We consider the following type of estimation problem. There is a collection of $\mathbb{C}$-valued random variables $\{h_{t\theta}^* : t \in \mathbb{N}, \theta \in \Theta\}$, where for any fixed $\theta \in \Theta$ the $\{h_{t\theta}^*\}_{t=1}^\infty$ variables are i.i.d. with distribution $\pi_\theta$. For each $\theta \in \Theta$, there is a sequence $\mathcal{Z}^t(\theta) = \{(X_{t1}, Y_{t1}(\theta)), (X_{t2}, Y_{t2}(\theta)), \ldots\}$, where $\{X_{ti}\}_{t,i \in \mathbb{N}}$ are i.i.d. $\mathcal{D}$, and for each $t, i \in \mathbb{N}$, $Y_{ti}(\theta) = h_{t\theta}^*(X_{ti})$. We additionally denote by $\mathcal{Z}_k^t(\theta) = \{(X_{t1}, Y_{t1}(\theta)), \ldots, (X_{tk}, Y_{tk}(\theta))\}$ the first $k$ elements of $\mathcal{Z}^t(\theta)$, for any $k \in \mathbb{N}$, and similarly $\mathbb{X}_{tk} = \{X_{t1}, \ldots, X_{tk}\}$ and $\mathbb{Y}_{tk}(\theta) = \{Y_{t1}(\theta), \ldots, Y_{tk}(\theta)\}$. Following the terminology used in the transfer learning literature, we refer to the collection of variables associated with each $t$ collectively as the $t^{\text{th}}$ *task*. We will be concerned with sequences of estimators $\hat{\theta}_{T\theta} = \hat{\theta}_T(\mathcal{Z}_k^1(\theta), \ldots, \mathcal{Z}_k^T(\theta))$, for $T \in \mathbb{N}$, which are based on only a bounded number $k$ of samples per task, among the first $T$ tasks. Our main results specifically study the case of $d$ samples per task. For any such estimator, we measure the *risk* as $\mathbb{E}\left[\|\pi_{\hat{\theta}_{T\theta_\star}} - \pi_{\theta_\star}\|\right]$, and will be particularly interested in upper-bounding the worst-case risk $\sup_{\theta_\star \in \Theta} \mathbb{E}\left[\|\pi_{\hat{\theta}_{T\theta_\star}} - \pi_{\theta_\star}\|\right]$ as a function of $T$, and lower-bounding the minimum possible value of this worst-case risk over all possible $\hat{\theta}_T$ estimators (called the *minimax risk*).

In previous work, [11] showed that, if $\Pi_\Theta$ is a totally bounded family, then even with only $d$ number of samples per task, the minimax risk (as a function of the number of tasks $T$) converges to zero. In fact, that work also proved this is not necessarily the case in general for any number of samples less than $d$. However, the actual rates of convergence were not explicitly derived in that work, and indeed the upper bounds on the rates of convergence implicit in that analysis may often have fairly complicated dependences on $\mathbb{C}$, $\Pi_\Theta$, and $\mathcal{D}$, and furthermore often provide only very slow rates of convergence.

To derive explicit bounds on the rates of convergence, in the present work we specifically focus on families of *smooth* densities. The motivation for involving a notion of smoothness in characterizing rates of convergence is clear if we consider the extreme case in which $\Pi_\Theta$ contains two priors $\pi_1$ and $\pi_2$, with $\pi_1(\{h\}) = \pi_2(\{g\}) = 1$, where $\rho(h, g)$ is a very small but nonzero value; in this case, if we have only a small number of samples per task, we would require many tasks (on the order of $1/\rho(h, g)$) to observe any data points carrying any information that would distinguish between these two priors (namely, points $x$ with $h(x) \neq g(x)$); yet $\|\pi_1 - \pi_2\| = 1$, so that we have a slow rate of convergence (at least initially). A total boundedness condition on $\Pi_\Theta$ would limit the number of such pairs

present in $\Pi_\Theta$, so that for instance we cannot have arbitrarily close $h$ and $g$, but less extreme variants of this can lead to slow asymptotic rates of convergence as well. Specifically, in the present work we consider the following notion of smoothness. For $L \in (0, \infty)$ and $\alpha \in (0, 1]$, a function $f : \mathbb{C} \to \mathbb{R}$ is $(L, \alpha)$-Hölder smooth if

$$\forall h, g \in \mathbb{C}, |f(h) - f(g)| \leq L\rho(h, g)^\alpha.$$

## 3    An Upper Bound

We now have the following theorem, holding for an arbitrary VC class $\mathbb{C}$ and data distribution $\mathcal{D}$; it is the main result of this work.

**Theorem 1.** *For $\Pi_\Theta$ any class of priors on $\mathbb{C}$ having $(L, \alpha)$-Hölder smooth densities $\{f_\theta : \theta \in \Theta\}$, for any $T \in \mathbb{N}$, there exists an estimator $\hat{\theta}_{T\theta} = \hat{\theta}_T(\mathcal{Z}_d^1(\theta), \dots, \mathcal{Z}_d^T(\theta))$ such that*

$$\sup_{\theta_\star \in \Theta} \mathbb{E}\|\pi_{\hat{\theta}_T} - \pi_{\theta_\star}\| = \tilde{O}\left(LT^{-\frac{\alpha^2}{2(d+2\alpha)(\alpha+2(d+1))}}\right).$$

*Proof.* By the standard PAC analysis [3,8], for any $\gamma > 0$, with probability greater than $1 - \gamma$, a sample of $k = O((d/\gamma)\log(1/\gamma))$ random points will partition $\mathbb{C}$ into regions of width less than $\gamma$ (under $L_1(\mathcal{D})$). For brevity, we omit the $t$ subscripts and superscripts on quantities such as $\mathcal{Z}_k^t(\theta)$ throughout the following analysis, since the claims hold for any arbitrary value of $t$.

For any $\theta \in \Theta$, let $\pi_\theta'$ denote a (conditional on $X_1, \dots, X_k$) distribution defined as follows. Let $f_\theta'$ denote the (conditional on $X_1, \dots, X_k$) density function of $\pi_\theta'$ with respect to $\pi_0$, and for any $g \in \mathbb{C}$, let $f_\theta'(g) = \frac{\pi_\theta(\{h \in \mathbb{C}: \forall i \leq k, h(X_i) = g(X_i)\})}{\pi_0(\{h \in \mathbb{C}: \forall i \leq k, h(X_i) = g(X_i)\})}$ (or 0 if $\pi_0(\{h \in \mathbb{C} : \forall i \leq k, h(X_i) = g(X_i)\}) = 0$). In other words, $\pi_\theta'$ has the same probability mass as $\pi_\theta$ for each of the equivalence classes induced by $X_1, \dots, X_k$, but conditioned on the equivalence class, simply has a constant-density distribution over that equivalence class. Note that every $h \in \mathbb{C}$ has $f_\theta'(h)$ between the smallest and largest values of $f_\theta(g)$ among $g \in \mathbb{C}$ with $\forall i \leq k, g(X_i) = h(X_i)$; therefore, by the smoothness condition, on the event (of probability greater than $1 - \gamma$) that each of these regions has diameter less than $\gamma$, we have $\forall h \in \mathbb{C}, |f_\theta(h) - f_\theta'(h)| < L\gamma^\alpha$. On this event, for any $\theta, \theta' \in \Theta$,

$$\|\pi_\theta - \pi_{\theta'}\| = (1/2)\int |f_\theta - f_{\theta'}|d\pi_0 < L\gamma^\alpha + (1/2)\int |f_\theta' - f_{\theta'}'|d\pi_0.$$

Furthermore, since the regions that define $f_\theta'$ and $f_{\theta'}'$ are the same (namely, the partition induced by $X_1, \dots, X_k$), we have

$$(1/2)\int |f_\theta' - f_{\theta'}'|d\pi_0 = (1/2)\sum_{y_1, \dots, y_k \in \{-1, +1\}} |\pi_\theta(\{h \in \mathbb{C} : \forall i \leq k, h(X_i) = y_i\})$$
$$- \pi_{\theta'}(\{h \in \mathbb{C} : \forall i \leq k, h(X_i) = y_i\})|$$
$$= \|\mathbb{P}_{\mathbb{Y}_k(\theta)|\mathbb{X}_k} - \mathbb{P}_{\mathbb{Y}_k(\theta')|\mathbb{X}_k}\|.$$

Thus, we have that with probability at least $1 - \gamma$,

$$\|\pi_\theta - \pi_{\theta'}\| < L\gamma^\alpha + \|\mathbb{P}_{\mathbb{Y}_k(\theta)|\mathbb{X}_k} - \mathbb{P}_{\mathbb{Y}_k(\theta')|\mathbb{X}_k}\|.$$

Following analogous to the inductive argument of [11], suppose $I \subseteq \{1, \ldots, k\}$, fix $\bar{x}_I \in \mathcal{X}^{|I|}$ and $\bar{y}_I \in \{-1, +1\}^{|I|}$. Then the $\tilde{y}_I \in \{-1, +1\}^{|I|}$ for which $\|\bar{y}_I - \tilde{y}_I\|_1$ is minimal, subject to the constraint that no $h \in \mathbb{C}$ has $h(\bar{x}_I) = \tilde{y}_I$, has $(1/2)\|\bar{y}_I - \tilde{y}_I\|_1 \le d + 1$; also, for any $i \in I$ with $\bar{y}_i \ne \tilde{y}_i$, letting $\bar{y}'_j = \bar{y}_j$ for $j \in I \setminus \{i\}$ and $\bar{y}'_i = \tilde{y}_i$, we have

$$\mathbb{P}_{\mathbb{Y}_I(\theta)|\mathbb{X}_I}(\bar{y}_I|\bar{x}_I) = \mathbb{P}_{\mathbb{Y}_{I\setminus\{i\}}(\theta)|\mathbb{X}_{I\setminus\{i\}}}(\bar{y}_{I\setminus\{i\}}|\bar{x}_{I\setminus\{i\}}) - \mathbb{P}_{\mathbb{Y}_I(\theta)|\mathbb{X}_I}(\bar{y}'_I|\bar{x}_I),$$

and similarly for $\theta'$, so that

$$\begin{aligned}
&|\mathbb{P}_{\mathbb{Y}_I(\theta)|\mathbb{X}_I}(\bar{y}_I|\bar{x}_I) - \mathbb{P}_{\mathbb{Y}_I(\theta')|\mathbb{X}_I}(\bar{y}_I|\bar{x}_I)| \\
&\le |\mathbb{P}_{\mathbb{Y}_{I\setminus\{i\}}(\theta)|\mathbb{X}_{I\setminus\{i\}}}(\bar{y}_{I\setminus\{i\}}|\bar{x}_{I\setminus\{i\}}) - \mathbb{P}_{\mathbb{Y}_{I\setminus\{i\}}(\theta')|\mathbb{X}_{I\setminus\{i\}}}(\bar{y}_{I\setminus\{i\}}|\bar{x}_{I\setminus\{i\}})| \\
&\quad + |\mathbb{P}_{\mathbb{Y}_I(\theta)|\mathbb{X}_I}(\bar{y}'_I|\bar{x}_I) - \mathbb{P}_{\mathbb{Y}_I(\theta')|\mathbb{X}_I}(\bar{y}'_I|\bar{x}_I)|.
\end{aligned}$$

Now consider that these two terms inductively define a binary tree. Every time the tree branches left once, it arrives at a difference of probabilities for a set $I$ of one less element than that of its parent. Every time the tree branches right once, it arrives at a difference of probabilities for a $\bar{y}_I$ one closer to an unrealized $\tilde{y}_I$ than that of its parent. Say we stop branching the tree upon reaching a set $I$ and a $\bar{y}_I$ such that either $\bar{y}_I$ is an unrealized labeling, or $|I| = d$. Thus, we can bound the original (root node) difference of probabilities by the sum of the differences of probabilities for the leaf nodes with $|I| = d$. Any path in the tree can branch left at most $k - d$ times (total) before reaching a set $I$ with only $d$ elements, and can branch right at most $d + 1$ times in a row before reaching a $\bar{y}_I$ such that both probabilities are zero, so that the difference is zero. So the depth of any leaf node with $|I| = d$ is at most $(k - d)d$. Furthermore, at any level of the tree, from left to right the nodes have strictly decreasing $|I|$ values, so that the maximum width of the tree is at most $k - d$. So the total number of leaf nodes with $|I| = d$ is at most $(k - d)^2 d$. Thus, for any $\bar{y} \in \{-1, +1\}^k$ and $\bar{x} \in \mathcal{X}^k$,

$$\begin{aligned}
&|\mathbb{P}_{\mathbb{Y}_k(\theta)|\mathbb{X}_k}(\bar{y}|\bar{x}) - \mathbb{P}_{\mathbb{Y}_k(\theta')|\mathbb{X}_k}(\bar{y}|\bar{x})| \\
&\le (k - d)^2 d \cdot \max_{\bar{y}^d \in \{-1, +1\}^d} \max_{D \in \{1, \ldots, k\}^d} |\mathbb{P}_{\mathbb{Y}_d(\theta)|\mathbb{X}_d}(\bar{y}^d|\bar{x}_D) - \mathbb{P}_{\mathbb{Y}_d(\theta')|\mathbb{X}_d}(\bar{y}^d|\bar{x}_D)|.
\end{aligned}$$

Since

$$\|\mathbb{P}_{\mathbb{Y}_k(\theta)|\mathbb{X}_k} - \mathbb{P}_{\mathbb{Y}_k(\theta')|\mathbb{X}_k}\| = (1/2) \sum_{\bar{y}^k \in \{-1, +1\}^k} |\mathbb{P}_{\mathbb{Y}_k(\theta)|\mathbb{X}_k}(\bar{y}^k) - \mathbb{P}_{\mathbb{Y}_k(\theta')|\mathbb{X}_k}(\bar{y}^k)|,$$

and by Sauer's Lemma this is at most

$$(ek)^d \max_{\bar{y}^k \in \{-1, +1\}^k} |\mathbb{P}_{\mathbb{Y}_k(\theta)|\mathbb{X}_k}(\bar{y}^k) - \mathbb{P}_{\mathbb{Y}_k(\theta')|\mathbb{X}_k}(\bar{y}^k)|,$$

we have that

$$
\begin{aligned}
&\|\mathbb{P}_{\mathbb{Y}_k(\theta)|\mathbb{X}_k} - \mathbb{P}_{\mathbb{Y}_k(\theta')|\mathbb{X}_k}\| \\
&\le (ek)^d k^2 d \max_{\bar{y}^d \in \{-1,+1\}^d} \max_{D \in \{1,\dots,k\}^d} |\mathbb{P}_{\mathbb{Y}_d(\theta)|\mathbb{X}_D}(\bar{y}^d) - \mathbb{P}_{\mathbb{Y}_d(\theta')|\mathbb{X}_D}(\bar{y}^d)|.
\end{aligned}
$$

Thus, we have that

$$
\|\pi_\theta - \pi_{\theta'}\| = \mathbb{E}\|\pi_\theta - \pi_{\theta'}\|
$$

$$
< \gamma + L\gamma^\alpha + (ek)^d k^2 d \mathbb{E}\left[ \max_{\bar{y}^d \in \{-1,+1\}^d} \max_{D \in \{1,\dots,k\}^d} \mathbb{P}_{\mathbb{Y}_d(\theta)|\mathbb{X}_D}(\bar{y}^d) - \mathbb{P}_{\mathbb{Y}_d(\theta')|\mathbb{X}_D}(\bar{y}^d)| \right].
$$

Note that

$$
\begin{aligned}
&\mathbb{E}\left[ \max_{\bar{y}^d \in \{-1,+1\}^d} \max_{D \in \{1,\dots,k\}^d} |\mathbb{P}_{\mathbb{Y}_d(\theta)|\mathbb{X}_D}(\bar{y}^d) - \mathbb{P}_{\mathbb{Y}_d(\theta')|\mathbb{X}_D}(\bar{y}^d)| \right] \\
&\le \sum_{\bar{y}^d \in \{-1,+1\}^d} \sum_{D \in \{1,\dots,k\}^d} \mathbb{E}\left[ |\mathbb{P}_{\mathbb{Y}_d(\theta)|\mathbb{X}_D}(\bar{y}^d) - \mathbb{P}_{\mathbb{Y}_d(\theta')|\mathbb{X}_D}(\bar{y}^d)| \right] \\
&\le (2k)^d \max_{\bar{y}^d \in \{-1,+1\}^d} \max_{D \in \{1,\dots,k\}^d} \mathbb{E}\left[ |\mathbb{P}_{\mathbb{Y}_d(\theta)|\mathbb{X}_D}(\bar{y}^d) - \mathbb{P}_{\mathbb{Y}_d(\theta')|\mathbb{X}_D}(\bar{y}^d)| \right],
\end{aligned}
$$

and by exchangeability, this last line equals

$$
(2k)^d \max_{\bar{y}^d \in \{-1,+1\}^d} \mathbb{E}\left[ |\mathbb{P}_{\mathbb{Y}_d(\theta)|\mathbb{X}_d}(\bar{y}^d) - \mathbb{P}_{\mathbb{Y}_d(\theta')|\mathbb{X}_d}(\bar{y}^d)| \right].
$$

[11] showed that $\mathbb{E}\left[ |\mathbb{P}_{\mathbb{Y}_d(\theta)|\mathbb{X}_d}(\bar{y}^d) - \mathbb{P}_{\mathbb{Y}_d(\theta')|\mathbb{X}_d}(\bar{y}^d)| \right] \le 4\sqrt{\|\mathbb{P}_{\mathcal{Z}_d(\theta)} - \mathbb{P}_{\mathcal{Z}_d(\theta')}\|}$, so that in total we have $\|\pi_\theta - \pi_{\theta'}\| < (L+1)\gamma^\alpha + 4(2ek)^{2d+2}\sqrt{\|\mathbb{P}_{\mathcal{Z}_d(\theta)} - \mathbb{P}_{\mathcal{Z}_d(\theta')}\|}$. Plugging in the value of $k = c(d/\gamma)\log(1/\gamma)$, this is

$$
(L+1)\gamma^\alpha + 4\left( 2ec\frac{d}{\gamma}\log\left(\frac{1}{\gamma}\right) \right)^{2d+2} \sqrt{\|\mathbb{P}_{\mathcal{Z}_d(\theta)} - \mathbb{P}_{\mathcal{Z}_d(\theta')}\|}.
$$

Thus, it suffices to bound the rate of convergence (in total variation distance) of some estimator of $\mathbb{P}_{\mathcal{Z}_d(\theta_\star)}$. If $N(\varepsilon)$ is the $\varepsilon$-covering number of $\{\mathbb{P}_{\mathcal{Z}_d(\theta)} : \theta \in \Theta\}$, then taking $\hat{\theta}_{T\theta_\star}$ as the minimum distance skeleton estimate of [5,13] achieves expected total variation distance $\varepsilon$ from $\mathbb{P}_{\mathcal{Z}_d(\theta_\star)}$, for some $T = O((1/\varepsilon^2)\log N(\varepsilon/4))$. We can partition $\mathbb{C}$ into $O((L/\varepsilon)^{d/\alpha})$ cells of diameter $O((\varepsilon/L)^{1/\alpha})$, and set a constant density value within each cell, on an $O(\varepsilon)$-grid of density values, and every prior with $(L,\alpha)$-Hölder smooth density will have density within $\varepsilon$ of some density so-constructed; there are then at most $(1/\varepsilon)^{O((L/\varepsilon)^{d/\alpha})}$ such densities, so this bounds the covering numbers of $\Pi_\Theta$. Furthermore, the covering number of $\Pi_\Theta$ upper bounds $N(\varepsilon)$ [11], so that $N(\varepsilon) \le (1/\varepsilon)^{O((L/\varepsilon)^{d/\alpha})}$.

Solving $T = O(\varepsilon^{-2}(L/\varepsilon)^{d/\alpha}\log(1/\varepsilon))$ for $\varepsilon$, we have $\varepsilon = O\left( L\left(\frac{\log(TL)}{T}\right)^{\frac{\alpha}{d+2\alpha}} \right)$.

So this bounds the rate of convergence for $\mathbb{E}\|\mathbb{P}_{\mathcal{Z}_d(\hat{\theta}_T)} - \mathbb{P}_{\mathcal{Z}_d(\theta_\star)}\|$, for $\hat{\theta}_T$ the

minimum distance skeleton estimate. Plugging this rate into the bound on the priors, combined with Jensen's inequality, we have

$$\mathbb{E}\|\pi_{\hat{\theta}_T} - \pi_{\theta_\star}\| < (L+1)\gamma^\alpha + 4\left(2ec\frac{d}{\gamma}\log\left(\frac{1}{\gamma}\right)\right)^{2d+2} \times O\left(L\left(\frac{\log(TL)}{T}\right)^{\frac{\alpha}{2d+4\alpha}}\right).$$

This holds for any $\gamma > 0$, so minimizing this expression over $\gamma > 0$ yields a bound on the rate. For instance, with $\gamma = \tilde{O}\left(T^{-\frac{\alpha}{2(d+2\alpha)(\alpha+2(d+1))}}\right)$, we have

$$\mathbb{E}\|\pi_{\hat{\theta}_T} - \pi_{\theta_\star}\| = \tilde{O}\left(LT^{-\frac{\alpha^2}{2(d+2\alpha)(\alpha+2(d+1))}}\right).$$

$\square$

## 4 A Minimax Lower Bound

One natural quesiton is whether Theorem 1 can generally be improved. While we expect this to be true for some fixed VC classes (e.g., those of finite size), and in any case we expect that some of the constant factors in the exponent may be improvable, it is not at this time clear whether the general form of $T^{-\Theta(\alpha^2/(d+\alpha)^2)}$ is sometimes optimal. One way to investigate this question is to construct specific spaces $\mathbb{C}$ and distributions $\mathcal{D}$ for which a lower bound can be obtained. In particular, we are generally interested in exhibiting lower bounds that are worse than those that apply to the usual problem of density estimation based on direct access to the $h_{t\theta_\star}^*$ values (see Theorem 3 below).

Here we present a lower bound that is interesting for this reason. However, although larger than the optimal rate for methods with direct access to the target concepts, it is still far from matching the upper bound above, so that the question of tightness remains open. Specifically, we have the following result.

**Theorem 2.** *For any integer $d \geq 1$, any $L > 0, \alpha \in (0, 1]$, there is a value $C(d, L, \alpha) \in (0, \infty)$ such that, for any $T \in \mathbb{N}$, there exists an instance space $\mathcal{X}$, a concept space $\mathbb{C}$ of VC dimension $d$, a distribution $\mathcal{D}$ over $\mathcal{X}$, and a distribution $\pi_0$ over $\mathbb{C}$ such that, for $\Pi_\Theta$ a set of distributions over $\mathbb{C}$ with $(L, \alpha)$-Hölder smooth density functions with respect to $\pi_0$, any estimator $\hat{\theta}_T = \hat{\theta}_T(\mathcal{Z}_d^1(\theta_\star), \ldots, \mathcal{Z}_d^T(\theta_\star))$ has*

$$\sup_{\theta_\star \in \Theta} \mathbb{E}\left[\|\pi_{\hat{\theta}_T} - \pi_{\theta_\star}\|\right] \geq C(d, L, \alpha)T^{-\frac{\alpha}{2(d+\alpha)}}.$$

*Proof.* (Sketch) We proceed by a reduction from the task of determining the bias of a coin from among two given possibilities. Specifically, fix any $\gamma \in (0, 1/2)$, $n \in \mathbb{N}$, and let $B_1(p), \ldots, B_n(p)$ be i.i.d Bernoulli($p$) random variables, for each $p \in [0, 1]$; then it is known that, for any (possibly nondeterministic) decision rule $\hat{p}_n : \{0, 1\}^n \to \{(1+\gamma)/2, (1-\gamma)/2\}$,

$$\frac{1}{2} \sum_{p \in \{(1+\gamma)/2, (1-\gamma)/2\}} \mathbb{P}(\hat{p}_n(B_1(p), \ldots, B_n(p)) \neq p)$$

$$\geq (1/32) \cdot \exp\left\{-128\gamma^2 n/3\right\}. \quad (1)$$

This easily follows from the results of [1], combined with a result of [6] bounding the KL divergence (see also [10])

To use this result, we construct a learning problem as follows. Fix some $m \in \mathbb{N}$ with $m \geq d$, let $\mathcal{X} = \{1, \ldots, m\}$, and let $\mathbb{C}$ be the space of all classifiers $h : \mathcal{X} \to \{-1, +1\}$ such that $|\{x \in \mathcal{X} : h(x) = +1\}| \leq d$. Clearly the VC dimension of $\mathbb{C}$ is $d$. Define the distribution $\mathcal{D}$ as uniform over $\mathcal{X}$. Finally, we specify a family of $(L, \alpha)$-Hölder smooth priors, parameterized by $\Theta = \{-1, +1\}^{\binom{m}{d}}$, as follows. Let $\gamma_m = (L/2)(1/m)^\alpha$. First, enumerate the $\binom{m}{d}$ distinct $d$-sized subsets of $\{1, \ldots, m\}$ as $\mathcal{X}_1, \mathcal{X}_2, \ldots, \mathcal{X}_{\binom{m}{d}}$. Define the reference distribution $\pi_0$ by the property that, for any $h \in \mathbb{C}$, letting $q = |\{x : h(x) = +1\}|$, $\pi_0(\{h\}) = (\frac{1}{2})^d \binom{m-q}{d-q} / \binom{m}{d}$. For any $\mathbf{b} = (b_1, \ldots, b_{\binom{m}{d}}) \in \{-1, 1\}^{\binom{m}{d}}$, define the prior $\pi_{\mathbf{b}}$ as the distribution of a random variable $h_{\mathbf{b}}$ specified by the following generative model. Let $i^* \sim \mathrm{Uniform}(\{1, \ldots, \binom{m}{d}\})$, let $C_{\mathbf{b}}(i^*) \sim \mathrm{Bernoulli}((1 + \gamma_m b_{i^*})/2)$; finally, $h_{\mathbf{b}} \sim \mathrm{Uniform}(\{h \in \mathbb{C} : \{x : h(x) = +1\} \subseteq \mathcal{X}_{i^*}, \mathrm{Parity}(|\{x : h(x) = +1\}|) = C_{\mathbf{b}}(i^*)\})$, where $\mathrm{Parity}(n)$ is 1 if $n$ is odd, or 0 if $n$ is even. We will refer to the variables in this generative model below. For any $h \in \mathbb{C}$, letting $H = \{x : h(x) = +1\}$ and $q = |H|$, we can equivalently express $\pi_{\mathbf{b}}(\{h\}) = (\frac{1}{2})^d \binom{m}{d}^{-1} \sum_{i=1}^{\binom{m}{d}} \mathbb{1}[H \subseteq \mathcal{X}_i](1 + \gamma_m b_i)^{\mathrm{Parity}(q)} (1 - \gamma_m b_i)^{1 - \mathrm{Parity}(q)}$. From this explicit representation, it is clear that, letting $f_{\mathbf{b}} = \frac{d\pi_{\mathbf{b}}}{d\pi_0}$, we have $f_{\mathbf{b}}(h) \in [1 - \gamma_m, 1 + \gamma_m]$ for all $h \in \mathbb{C}$. The fact that $f_{\mathbf{b}}$ is Hölder smooth follows from this, since every distinct $h, g \in \mathbb{C}$ have $\mathcal{D}(\{x : h(x) \neq g(x)\}) \geq 1/m = (2\gamma_m/L)^{1/\alpha}$.

Next we set up the reduction as follows. For any estimator $\hat{\pi}_T = \hat{\pi}_T(\mathcal{Z}_d^1(\theta_\star), \ldots, \mathcal{Z}_d^T(\theta_\star))$, and each $i \in \{1, \ldots, \binom{m}{d}\}$, let $h_i$ be the classifier with $\{x : h_i(x) = +1\} = \mathcal{X}_i$; also, if $\hat{\pi}_T(\{h_i\}) > (\frac{1}{2})^d / \binom{m}{d}$, let $\hat{b}_i = 2\mathrm{Parity}(d) - 1$, and otherwise $\hat{b}_i = 1 - 2\mathrm{Parity}(d)$. We use these $\hat{b}_i$ values to estimate the original $b_i$ values. Specifically, let $\hat{p}_i = (1 + \gamma_m \hat{b}_i)/2$ and $p_i = (1 + \gamma_m b_i)/2$, where $\mathbf{b} = \theta_\star$. Then

$$\|\hat{\pi}_T - \pi_{\theta_\star}\| \geq (1/2) \sum_{i=1}^{\binom{m}{d}} |\hat{\pi}_T(\{h_i\}) - \pi_{\theta_\star}(\{h_i\})|$$

$$\geq (1/2) \sum_{i=1}^{\binom{m}{d}} \frac{\gamma_m}{2^d \binom{m}{d}} |\hat{b}_i - b_i|/2 = (1/2) \sum_{i=1}^{\binom{m}{d}} \frac{1}{2^d \binom{m}{d}} |\hat{p}_i - p_i|.$$

Thus, we have reduced from the problem of deciding the biases of these $\binom{m}{d}$ independent Bernoulli random variables. To complete the proof, it suffices to lower bound the expectation of the right side for an *arbitrary* estimator.

Toward this end, we in fact study an even easier problem. Specifically, consider an estimator $\hat{q}_i = \hat{q}_i(\mathcal{Z}_d^1(\theta_\star), \ldots, \mathcal{Z}_d^T(\theta_\star), i_1^*, \ldots, i_T^*)$, where $i_t^*$ is the $i^*$ random variable in the generative model that defines $h_{t\theta_\star}^*$; that is, $i_t^* \sim \mathrm{Uniform}(\{1, \ldots, \binom{m}{d}\})$, $C_t \sim \mathrm{Bernoulli}((1 + \gamma_m b_{i_t^*})/2)$, and $h_{t\theta_\star}^* \sim \mathrm{Uniform}(\{h \in \mathbb{C} : \{x : h(x) = +1\} \subseteq \mathcal{X}_{i_t^*}, \mathrm{Parity}(|\{x : h(x) = +1\}|) = C_t\})$, where the $i_t^*$ are

independent across $t$, as are the $C_t$ and $h^*_{t\theta_\star}$. Clearly the $\hat{p}_i$ from above can be viewed as an estimator of this type, which simply ignores the knowledge of $i^*_t$. The knowledge of these $i^*_t$ variables simplifies the analysis, since given $\{i^*_t : t \leq T\}$, the data can be partitioned into $\binom{m}{d}$ disjoint sets, $\{\{\mathcal{Z}^t_d(\theta_\star) : i^*_t = i\} : i = 1, \ldots, \binom{m}{d}\}$, and we can use only the set $\{\mathcal{Z}^t_d(\theta_\star) : i^*_t = i\}$ to estimate $p_i$. Furthermore, we can use only the subset of these for which $\mathbb{X}_{td} = \mathcal{X}_i$, since otherwise we have zero information about the value of Parity$(|\{x : h^*_{t\theta_\star}(x) = +1\}|)$. That is, given $i^*_t = i$, any $\mathcal{Z}^t_d(\theta_\star)$ is conditionally independent from every $b_j$ for $j \neq i$, and is even conditionally independent from $b_i$ when $\mathbb{X}_{td}$ is not completely contained in $\mathcal{X}_i$; specifically, in this case, regardless of $b_i$, the conditional distribution of $\mathbb{Y}_{td}(\theta_\star)$ given $i^*_t = i$ and given $\mathbb{X}_{td}$ is a product distribution, which deterministically assigns label $-1$ to those $Y_{tk}(\theta_\star)$ with $X_{tk} \notin \mathcal{X}_i$, and gives uniform random values to the subset of $\mathbb{Y}_{td}(\theta_\star)$ with their respective $X_{tk} \in \mathcal{X}_i$. Finally, letting $r_t = $ Parity$(|\{k \leq d : Y_{tk}(\theta_\star) = +1\}|)$, we note that given $i^*_t = i$, $\mathbb{X}_{td} = \mathcal{X}_i$, and the value $r_t$, $b_i$ is conditionally independent from $\mathcal{Z}^t_d(\theta_\star)$. Thus, the set of values $C_{iT}(\theta_\star) = \{r_t : i^*_t = i, \mathbb{X}_{td} = \mathcal{X}_i\}$ is a sufficient statistic for $b_i$ (hence for $p_i$). Recall that, when $i^*_t = i$ and $\mathbb{X}_{td} = \mathcal{X}_i$, the value of $r_t$ is equal to $C_t$, a Bernoulli$(p_i)$ random variable. Thus, we neither lose nor gain anything (in terms of risk) by restricting ourselves to estimators $\hat{q}_i$ of the type $\hat{q}_i = \hat{q}_i(\mathcal{Z}^1_d(\theta_\star), \ldots, \mathcal{Z}^T_d(\theta_\star), i^*_1, \ldots, i^*_T) = \hat{q}'_i(C_{iT}(\theta_\star))$, for some $\hat{q}'_i$ [7]: that is, estimators that are a function of the $N_{iT}(\theta_\star) = |C_{iT}(\theta_\star)|$ Bernoulli$(p_i)$ random variables, which we should note are conditionally i.i.d. given $N_{iT}(\theta_\star)$.

Thus, by (1), for any $n \leq T$,

$$\frac{1}{2} \sum_{b_i \in \{-1,+1\}} \mathbb{E}\left[|\hat{q}_i - p_i| \,\Big|\, N_{iT}(\theta_\star) = n\right] = \frac{1}{2} \sum_{b_i \in \{-1,+1\}} \gamma_m \mathbb{P}\left(\hat{q}_i \neq p_i \,\Big|\, N_{iT}(\theta_\star) = n\right)$$
$$\geq (\gamma_m/32) \cdot \exp\left\{-128\gamma_m^2 N_i/3\right\}.$$

Also, $\forall i$, $\mathbb{E}[N_i] = \frac{d!(1/m)^d}{\binom{m}{d}} T \leq (d/m)^{2d} T = d^{2d}(2\gamma_m/L)^{2d/\alpha} T$. Thus, Jensen's inequality, linearity of expectation, and the law of total expectation imply

$$\frac{1}{2} \sum_{b_i \in \{-1,+1\}} \mathbb{E}\left[|\hat{q}_i - p_i|\right] \geq (\gamma_m/32) \cdot \exp\left\{-43(2/L)^{2d/\alpha} d^{2d} \gamma_m^{2+2d/\alpha} T\right\}.$$

Thus, by linearity of the expectation,

$$\left(\frac{1}{2}\right)^{\binom{m}{d}} \sum_{\mathbf{b} \in \{-1,+1\}^{\binom{m}{d}}} \mathbb{E}\left[\sum_{i=1}^{\binom{m}{d}} \frac{1}{2^d \binom{m}{d}} |\hat{q}_i - p_i|\right] = \sum_{i=1}^{\binom{m}{d}} \frac{1}{2^d \binom{m}{d}} \frac{1}{2} \sum_{b_i \in \{-1,+1\}} \mathbb{E}\left[|\hat{q}_i - p_i|\right]$$
$$\geq (\gamma_m/(32 \cdot 2^d)) \cdot \exp\left\{-43(2/L)^{2d/\alpha} d^{2d} \gamma_m^{2+2d/\alpha} T\right\}.$$

In particular, taking $m = \left\lceil (L/2)^{1/\alpha} \left( 43(2/L)^{2d/\alpha} d^{2d} T \right)^{\frac{1}{2(d+\alpha)}} \right\rceil$, we have $\gamma_m = \Theta\left( \left( 43(2/L)^{2d/\alpha} d^{2d} T \right)^{-\frac{\alpha}{2(d+\alpha)}} \right)$, so that

$$\left( \frac{1}{2} \right)^{\binom{m}{d}} \sum_{\mathbf{b} \in \{-1,+1\}^{\binom{m}{d}}} \mathbb{E}\left[ \sum_{i=1}^{\binom{m}{d}} \frac{1}{2^d \binom{m}{d}} |\hat{q}_i - p_i| \right] = \Omega\left( 2^{-d} \left( 43(2/L)^{2d/\alpha} d^{2d} T \right)^{-\frac{\alpha}{2(d+\alpha)}} \right).$$

In particular, this implies there exists some $\mathbf{b}$ for which

$$\mathbb{E}\left[ \sum_{i=1}^{\binom{m}{d}} \frac{1}{2^d \binom{m}{d}} |\hat{q}_i - p_i| \right] = \Omega\left( 2^{-d} \left( 43(2/L)^{2d/\alpha} d^{2d} T \right)^{-\frac{\alpha}{2(d+\alpha)}} \right).$$

Applying this lower bound to the estimator $\hat{p}_i$ above yields the result.  □

It is natural to wonder how these rates might potentially improve if we allow $\hat{\theta}_T$ to depend on more than $d$ samples per data set. To establish limits on such improvements, we note that in the extreme case of allowing the estimator to depend on the full $\mathcal{Z}^t(\theta_\star)$ data sets, we may recover the known results lower bounding the risk of density estimation from i.i.d. samples from a smooth density, as indicated by the following result.

**Theorem 3.** *For any integer $d \geq 1$, there exists an instance space $\mathcal{X}$, a concept space $\mathbb{C}$ of VC dimension $d$, a distribution $\mathcal{D}$ over $\mathcal{X}$, and a distribution $\pi_0$ over $\mathbb{C}$ such that, for $\Pi_\Theta$ the set of distributions over $\mathbb{C}$ with $(L, \alpha)$-Hölder smooth density functions with respect to $\pi_0$, any sequence of estimators, $\hat{\theta}_T = \hat{\theta}_T(\mathcal{Z}^1(\theta_\star), \ldots, \mathcal{Z}^T(\theta_\star))$ $(T = 1, 2, \ldots)$, has*

$$\sup_{\theta_\star \in \Theta} \mathbb{E}\left[ \|\pi_{\hat{\theta}_T} - \pi_{\theta_\star}\| \right] = \Omega\left( T^{-\frac{\alpha}{d+2\alpha}} \right).$$

The proof is a simple reduction from the problem of estimating $\pi_{\theta_\star}$ based on direct access to $h^*_{1\theta_\star}, \ldots, h^*_{T\theta_\star}$, which is essentially equivalent to the standard model of density estimation, and indeed the lower bound in Theorem 3 is a well-known result for density estimation from $T$ i.i.d. samples from a Hölder smooth density in a $d$-dimensional space [5].

## 5    Real-Valued Functions and an Application in Algorithmic Economics

In this section, we present results generalizing the analysis of [11] to classes of real-valued functions. We also present an application of this generalization to a preference elicitation problem.

### 5.1    Consistent Estimation of Priors over Real-Valued Functions at a Bounded Rate

In this section, we let $\mathcal{B}$ denote a $\sigma$-algebra on $\mathcal{X} \times \mathbb{R}$, and again let $\mathcal{B}_{\mathcal{X}}$ denote the corresponding $\sigma$-algebra on $\mathcal{X}$. Also, for measurable functions $h, g : \mathcal{X} \to \mathbb{R}$, let $\rho(h, g) = \int |h - g| dP_X$, where $P_X$ is a distribution over $\mathcal{X}$. Let $\mathcal{F}$ be a class of functions $\mathcal{X} \to \mathbb{R}$ with Borel $\sigma$-algebra $\mathcal{B}_{\mathcal{F}}$ induced by $\rho$. Let $\Theta$ be a set, and for each $\theta \in \Theta$, let $\pi_\theta$ denote a probability measure on $(\mathcal{F}, \mathcal{B}_{\mathcal{F}})$. We suppose $\{\pi_\theta : \theta \in \Theta\}$ is totally bounded in total variation distance, and that $\mathcal{F}$ is a uniformly bounded VC subgraph class with pseudodimension $d$. We also suppose $\rho$ is a *metric* when restricted to $\mathcal{F}$.

As above, let $\{X_{ti}\}_{t,i\in\mathbb{N}}$ be i.i.d. $P_X$ random variables. For each $\theta \in \Theta$, let $\{h_{t\theta}^*\}_{t\in\mathbb{N}}$ be i.i.d. $\pi_\theta$ random variables, independent from $\{X_{ti}\}_{t,i\in\mathbb{N}}$. For each $t \in \mathbb{N}$ and $\theta \in \Theta$, let $Y_{ti}(\theta) = h_{t\theta}^*(X_{ti})$ for $i \in \mathbb{N}$, and let $\mathcal{Z}^t(\theta) = \{(X_{t1}, Y_{t1}(\theta)), (X_{t2}, Y_{t2}(\theta)), \ldots\}$; for each $k \in \mathbb{N}$, define $\mathcal{Z}_k^t(\theta) = \{(X_{t1}, Y_{t1}(\theta)), \ldots, (X_{tk}, Y_{tk}(\theta))\}$, $\mathbb{X}_{tk} = \{X_{t1}, \ldots, X_{tk}\}$, and $\mathbb{Y}_{tk}(\theta) = \{Y_{t1}(\theta), \ldots, Y_{tk}(\theta)\}$.

We have the following result. The proof parallels that of [11] (who studied the special case of binary functions), with a few important twists (in particular, a significantly different approach in the analogue of their Lemma 3). Due to space restrictions, the formal details are omitted; we refer the interested reader to the full version of this article online [12].

**Theorem 4.** *There exists an estimator* $\hat{\theta}_{T\theta_\star} = \hat{\theta}_T(\mathcal{Z}_d^1(\theta_\star), \ldots, \mathcal{Z}_d^T(\theta_\star))$, *and functions* $R : \mathbb{N}_0 \times (0, 1] \to [0, \infty)$ *and* $\delta : \mathbb{N}_0 \times (0, 1] \to [0, 1]$ *such that, for any* $\alpha > 0$, $\lim_{T\to\infty} R(T, \alpha) = \lim_{T\to\infty} \delta(T, \alpha) = 0$ *and for any* $T \in \mathbb{N}_0$ *and* $\theta_\star \in \Theta$,

$$\mathbb{P}\left(\|\pi_{\hat{\theta}_{T\theta_\star}} - \pi_{\theta_\star}\| > R(T, \alpha)\right) \le \delta(T, \alpha) \le \alpha.$$

### 5.2    Maximizing Customer Satisfaction in Combinatorial Auctions

Theorem 4 has a clear application in the context of transfer learning, following analogous arguments to those given in the special case of binary classification by [11]. In addition to that application, we can also use Theorem 4 in the context of the following problem in algorithmic economics, where the objective is to serve a sequence of customers so as to maximize their satisfaction.

Consider an online travel agency, where customers go to the site with some idea of what type of travel they are interested in; the site then poses a series of questions to each customer, and identifies a travel package that best suits their desires, budget, and dates. There are many options of travel packages, with options on location, site-seeing tours, hotel and room quality, etc. Because of this, serving the needs of an *arbitrary* customer might be a lengthy process, requiring many detailed questions. Fortunately, the stream of customers is typically not a worst-case sequence, and in particular obeys many statistical regularities: in particular, it is not too far from reality to think of the customers as being independent and identically distributed samples. With this assumption in mind, it becomes desirable to identify some of these statistical regularities so that we

can pose the questions that are typically most relevant, and thereby more quickly identify the travel package that best suits the needs of the typical customer. One straightforward way to do this is to directly *estimate* the distribution of customer value functions, and optimize the questioning system to minimize the expected number of questions needed to find a suitable travel package.

One can model this problem in the style of Bayesian combinatorial auctions, in which each customer has a value function for each possible bundle of items. However, it is slightly different, in that we do not assume the distribution of customers is known, but rather are interested in estimating this distribution; the obtained estimate can then be used in combination with methods based on Bayesian decision theory. In contrast to the literature on Bayesian auctions (and subjectivist Bayesian decision theory in general), this technique is able to maintain general guarantees on performance that hold under an objective interpretation of the problem, rather than merely guarantees holding under an arbitrary assumed prior belief. This general idea is sometimes referred to as *Empirical Bayesian* decision theory in the machine learning and statistics literatures. The ideal result for an Empirical Bayesian algorithm is to be competitive with the corresponding Bayesian methods based on the *actual* distribution of the data (assuming the data are random, with an unknown distribution); that is, although the Empirical Bayesian methods only operate with a data-based estimate of the distribution, the aim is to perform nearly as well as methods based on the true (unobservable) distribution. In this work, we present results of this type, in the context of an abstraction of the aforementioned online travel agency problem, where the measure of performance is the expected number of questions to find a suitable package.

The specific application we are interested in here may be expressed abstractly as a kind of combinatorial auction with preference elicitation. Specifically, we suppose there is a collection of items on a menu, and each possible bundle of items has an associated fixed price. There is a stream of customers, each with a valuation function that provides a value for each possible bundle of items. The objective is to serve each customer a bundle of items that nearly-maximizes his or her surplus value (value minus price). However, we are not permitted direct observation of the customer valuation functions; rather, we may query for the value of any given bundle of items; this is referred to as a *value query* in the literature on preference elicitation in combinatorial auctions (see Chapter 14 of [4], [14]). The objective is to achieve this near-maximal surplus guarantee, while making only a small number of queries per customer. We suppose the customer valuation function are sampled i.i.d. according to an unknown distribution over a known (but arbitrary) class of real-valued functions having finite pseudo-dimension. Reasoning that knowledge of this distribution should allow one to make a smaller number of value queries per customer, we are interested in estimating this unknown distribution, so that as we serve more and more customers, the number of queries per customer required to identify a near-optimal bundle should decrease. In this context, we in fact prove that in the limit, the

expected number of queries per customer converges to the number required of a method having direct knowledge of the true distribution of valuation functions.

Formally, suppose there is a menu of $n$ items $[n] = \{1, \ldots, n\}$, and each bundle $B \subseteq [n]$ has an associated price $p(B) \geq 0$. Suppose also there is a sequence of customers, each with a valuation function $v_t : 2^{[n]} \to \mathbb{R}$. We suppose these $v_t$ functions are i.i.d. samples. We can then calculate the satisfaction function for each customer as $s_t(x)$, where $x \in \{0,1\}^n$, and $s_t(x) = v_t(B_x) - p(B_x)$, where $B_x \subseteq [n]$ contains element $i \in [n]$ iff $x_i = 1$.

Now suppose we are able to ask each customer a number of questions before serving up a bundle $B_{\hat{x}_t}$ to that customer. More specifically, we are able to ask for the value $s_t(x)$ for any $x \in \{0,1\}^n$. This is referred to as a *value query* in the literature on preference elicitation in combinatorial auctions (see Chapter 14 of [4], [14]). We are interested in asking as few questions as possible, while satisfying the guarantee that $\mathbb{E}[s_t(\hat{x}_t) - \max_x s_t(x)] \leq \varepsilon$.

Now suppose, for every $\pi$ and $\varepsilon$, we have a method $A(\pi, \varepsilon)$ such that, given that $\pi$ is the actual distribution of the $s_t$ functions, $A(\pi, \varepsilon)$ guarantees that the $\hat{x}_t$ value it selects has $\mathbb{E}[\max_x s_t(x) - s_t(\hat{x}_t)] \leq \varepsilon$; also let $\hat{N}_t(\pi, \varepsilon)$ denote the actual (random) number of queries the method $A(\pi, \varepsilon)$ would ask for the $s_t$ function, and let $Q(\pi, \varepsilon) = \mathbb{E}[\hat{N}_t(\pi, \varepsilon)]$. We suppose the method never queries any $s_t(x)$ value twice for a given $t$, so that its number of queries for any given $t$ is bounded.

Also suppose $\mathcal{F}$ is a VC subgraph class of functions mapping $\mathcal{X} = \{0,1\}^n$ into $[-1, 1]$ with pseudodimension $d$, and that $\{\pi_\theta : \theta \in \Theta\}$ is a known totally bounded family of distributions over $\mathcal{F}$ such that the $s_t$ functions have distribution $\pi_{\theta_\star}$ for some unknown $\theta_\star \in \Theta$. For any $\theta \in \Theta$ and $\gamma > 0$, let $\mathrm{B}(\theta, \gamma) = \{\theta' \in \Theta : \|\pi_\theta - \pi_{\theta'}\| \leq \gamma\}$.

Suppose, in addition to $A$, we have another method $A'(\varepsilon)$ that is not $\pi$-dependent, but still provides the $\varepsilon$-correctness guarantee, and makes a bounded number of queries (e.g., in the worst case, we could consider querying all $2^n$ points, but in most cases there are more clever $\pi$-independent methods that use far fewer queries, such as $O(1/\varepsilon^2)$). Consider the method described in Algorithm 1; the quantities $\hat{\theta}_{T\theta_\star}$, $R(T, \alpha)$, and $\delta(T, \alpha)$ from Theorem 4 are here considered with respect $P_X$ taken as the uniform distribution on $\{0,1\}^n$.

The following theorem indicates that Algorithm 1 is correct, and furthermore that the long-run average number of queries is not much worse than that of a method that has direct knowledge of $\pi_{\theta_\star}$. The proof of this result parallels that of [11] for the transfer learning setting, but is included here for completeness.

**Theorem 5.** *In Algorithm 1, $\forall t \leq T, \mathbb{E}[\max_x s_t(x) - s_t(\hat{x}_t)] \leq \varepsilon$. Furthermore, if $S_T(\varepsilon)$ is the total number of queries made by the method, then*

$$\limsup_{T \to \infty} \frac{\mathbb{E}[S_T(\varepsilon)]}{T} \leq Q(\pi_{\theta_\star}, \varepsilon/4) + d.$$

*Proof.* By Theorem 4, for any $t \leq T$, if $R(t-1, \varepsilon/2) \leq \varepsilon/8$, then with probability at least $1 - \varepsilon/2$, $\|\pi_{\theta_\star} - \pi_{\hat{\theta}_{(t-1)\theta_\star}}\| \leq R(t-1, \varepsilon/2)$, so that a triangle inequality

**Algorithm 1.** An algorithm for sequentially maximizing expected customer satisfaction.

> **for** $t = 1, 2, \ldots, T$ **do**
>> Pick points $X_{t1}, X_{t2}, \ldots, X_{td}$ uniformly at random from $\{0,1\}^n$
>> **if** $R(t-1, \varepsilon/2) > \varepsilon/8$ **then**
>>> Run $A'(\varepsilon)$
>>> Take $\hat{x}_t$ as the returned value
>> **else**
>>> Let $\check{\theta}_{t\theta_\star} \in \mathrm{B}\left(\hat{\theta}_{(t-1)\theta_\star}, R(t-1, \varepsilon/2)\right)$ be such that
>>> $$Q(\pi_{\check{\theta}_{t\theta_\star}}, \varepsilon/4) \le \min_{\theta \in \mathrm{B}\left(\hat{\theta}_{(t-1)\theta_\star}, R(t-1, \varepsilon/2)\right)} Q(\pi_\theta, \varepsilon/4) + \frac{1}{t}$$
>>> Run $A(\pi_{\check{\theta}_{t\theta_\star}}, \varepsilon/4)$ and let $\hat{x}_t$ be its return value
>> **end if**
> **end for**

implies $\|\pi_{\theta_\star} - \pi_{\check{\theta}_{t\theta_\star}}\| \le 2R(t-1, \varepsilon/2) \le \varepsilon/4$. Thus, $\mathbb{E}\left[\max_x s_t(x) - s_t(\hat{x}_t)\right] \le \varepsilon/2 + \mathbb{E}\left[\mathbb{E}\left[\max_x s_t(x) - s_t(\hat{x}_t)\Big|\check{\theta}_{t\theta_\star}\right]\mathbb{1}\left[\|\pi_{\check{\theta}_{t\theta_\star}} - \pi_{\theta_\star}\| \le \varepsilon/2\right]\right]$. For $\theta \in \Theta$, let $\hat{x}_{t\theta}$ denote the point $x$ that would be returned by $A(\pi_{\check{\theta}_{t\theta_\star}}, \varepsilon/4)$ when queries are answered by some $s_{t\theta} \sim \pi_\theta$ instead of $s_t$ (and supposing $s_t = s_{t\theta_\star}$). If $\|\pi_{\check{\theta}_{t\theta_\star}} - \pi_{\theta_\star}\| \le \varepsilon/4$, then

$$\mathbb{E}\left[\max_x s_t(x) - s_t(\hat{x}_t)\Big|\check{\theta}_{t\theta_\star}\right] = \mathbb{E}\left[\max_x s_{t\theta_\star}(x) - s_{t\theta_\star}(\hat{x}_t)\Big|\check{\theta}_{t\theta_\star}\right]$$
$$\le \mathbb{E}\left[\max_x s_{t\check{\theta}_{t\theta_\star}}(x) - s_{t\check{\theta}_{t\theta_\star}}(\hat{x}_{t\check{\theta}_{t\theta_\star}})\Big|\check{\theta}_{t\theta_\star}\right] + \|\pi_{\check{\theta}_{t\theta_\star}} - \pi_{\theta_\star}\| \le \varepsilon/4 + \varepsilon/4 = \varepsilon/2.$$

Plugging into the above bound, we have $\mathbb{E}\left[\max_x s_t(x) - s_t(\hat{x}_t)\right] \le \varepsilon$.

For the result on $S_T(\varepsilon)$, first note that $R(t-1, \varepsilon/2) > \varepsilon/8$ only finitely many times (due to $R(t, \alpha) = o(1)$), so that we can ignore those values of $t$ in the asymptotic calculation (as the number of queries is always bounded), and rely on the correctness guarantee of $A'$. For the remaining values $t$, let $N_t$ denote the number of queries made by $A(\pi_{\check{\theta}_{t\theta_\star}}, \varepsilon/4)$. then $\limsup_{T\to\infty} \frac{\mathbb{E}[S_T(\varepsilon)]}{T} \le d + \limsup_{T\to\infty} \sum_{t=1}^T \frac{\mathbb{E}[N_t]}{T}$. Since $\lim_{T\to\infty} \frac{1}{T}\sum_{t=1}^T \mathbb{E}\left[N_t\mathbb{1}[\|\pi_{\hat{\theta}_{(t-1)\theta_\star}} - \pi_{\theta_\star}\| > R(t-1, \varepsilon/2)]\right] \le \lim_{T\to\infty} \frac{1}{T}\sum_{t=1}^T 2^n\mathbb{P}\left(\|\pi_{\hat{\theta}_{(t-1)\theta_\star}} - \pi_{\theta_\star}\| > R(t-1, \varepsilon/2)\right) \le 2^n \lim_{T\to\infty} \frac{1}{T}\sum_{t=1}^T \delta(t-1, \varepsilon/2) = 0$, we have $\limsup_{T\to\infty} \sum_{t=1}^T \frac{\mathbb{E}[N_t]}{T} = \limsup_{T\to\infty} \frac{1}{T}\sum_{t=1}^T \mathbb{E}\left[N_t\mathbb{1}[\|\pi_{\hat{\theta}_{(t-1)\theta_\star}} - \pi_{\theta_\star}\| \le R(t-1, \varepsilon/2)]\right]$. For $t \le T$, let $N_t(\check{\theta}_{t\theta_\star})$ denote the number of queries $A(\pi_{\check{\theta}_{t\theta_\star}}, \varepsilon/4)$ would make if queries were answered with $s_{t\check{\theta}_{t\theta_\star}}$ instead of $s_t$. On the event $\|\pi_{\hat{\theta}_{(t-1)\theta_\star}} - \pi_{\theta_\star}\| \le R(t-1, \varepsilon/2)$, we have $\mathbb{E}\left[N_t\Big|\check{\theta}_{t\theta_\star}\right] \le \mathbb{E}\left[N_t(\check{\theta}_{t\theta_\star})\Big|\check{\theta}_{t\theta_\star}\right] + 2R(t-1, \varepsilon/2) = Q(\pi_{\check{\theta}_{t\theta_\star}}, \varepsilon/4) + 2R(t-1, \varepsilon/2) \le Q(\pi_{\theta_\star}, \varepsilon/4) + 2R(t-1, \varepsilon/2) + 1/t$. Therefore, $\limsup_{T\to\infty} \frac{1}{T}\sum_{t=1}^T \mathbb{E}\left[N_t\mathbb{1}[\|\pi_{\hat{\theta}_{(t-1)\theta_\star}} - \pi_{\theta_\star}\| \le R(t-1, \varepsilon/2)]\right] \le Q(\pi_{\theta_\star}, \varepsilon/4) + \limsup_{T\to\infty} \frac{1}{T}\sum_{t=1}^T 2R(t-1, \varepsilon/2) + 1/t = Q(\pi_{\theta_\star}, \varepsilon/4)$. $\qquad \square$

In many cases, this result will even continue to hold with an infinite number of goods ($n = \infty$), since Theorem 4 has no dependence on the cardinality of $\mathcal{X}$.

## 6    Open Problems

There are several interesting questions that remain open at this time. Can either the lower bound or upper bound be improved in general? If, instead of $d$ samples per task, we instead use $m \geq d$ samples, how does the minimax risk vary with $m$? Related to this, what is the optimal value of $m$ to optimize the rate of convergence as a function of $mT$, the total number of samples? More generally, if an estimator is permitted to use $N$ total samples, taken from however many tasks it wishes, what is the optimal rate of convergence as a function of $N$?

## References

1. Bar-Yossef, Z.: Sampling lower bounds via information theory. In: Proceedings of the 35th Annual ACM Symposium on the Theory of Computing, pp. 335–344 (2003)
2. Baxter, J.: A Bayesian/information theoretic model of learning to learn via multiple task sampling. Machine Learning **28**, 7–39 (1997)
3. Blumer, A., Ehrenfeucht, A., Haussler, D., Warmuth, M.: Learnability and the Vapnik-Chervonenkis dimension. Journal of the Association for Computing Machinery **36**(4), 929–965 (1989)
4. Cramton, P., Shoham, Y., Steinberg, R.: Combinatorial Auctions. The MIT Press (2006)
5. Devroye, L., Lugosi, G.: Combinatorial Methods in Density Estimation. Springer, New York (2001)
6. Poland, J., Hutter, M.: MDL convergence speed for Bernoulli sequences. Statistics and Computing **16**, 161–175 (2006)
7. Schervish, M.J.: Theory of Statistics. Springer, New York (1995)
8. Vapnik, V.: Estimation of Dependencies Based on Empirical Data. Springer-Verlag, New York (1982)
9. Vapnik, V., Chervonenkis, A.: On the uniform convergence of relative frequencies of events to their probabilities. Theory of Probability and its Applications **16**, 264–280 (1971)
10. Wald, A.: Sequential tests of statistical hypotheses. The Annals of Mathematical Statistics **16**(2), 117–186 (1945)
11. Yang, L., Hanneke, S., Carbonell, J.: A theory of transfer learning with applications to active learning. Machine Learning **90**(2), 161–189 (2013)
12. Yang, L., Hanneke, S., Carbonell, J.: Bounds on the minimax rate for estimating a prior over a vc class from independent learning tasks. arXiv:1505.05231 (2015)
13. Yatracos, Y.G.: Rates of convergence of minimum distance estimators and Kolmogorov's entropy. The Annals of Statistics **13**, 768–774 (1985)
14. Zinkevich, M., Blum, A., Sandholm, T.: On polynomial-time preference elicitation with value queries. In: Proceedings of the 4th ACM Conference on Electronic Commerce, pp. 175–185 (2003)