

---

# Nearly Optimal Algorithms for Level Set Estimation

---

**Blake Mason**  
Rice University

**Romain Camilleri**  
University of Washington

**Subhojyoti Mukherjee**  
University of Wisconsin–Madison

**Kevin Jamieson**  
University of Washington

**Robert Nowak**  
University of Wisconsin–Madison

**Lalit Jain**  
University of Washington

## Abstract

The level set estimation problem seeks to find all points in a domain  $\mathcal{X}$  where the value of an unknown function  $f : \mathcal{X} \rightarrow \mathbb{R}$  exceeds a threshold  $\alpha$ . The estimation is based on noisy function evaluations that may be acquired at sequentially and adaptively chosen locations in  $\mathcal{X}$ . The threshold value  $\alpha$  can either be *explicit* and provided a priori, or *implicit* and defined relative to the optimal function value, i.e.  $\alpha = (1 - \epsilon)f(\mathbf{x}_*)$  for a given  $\epsilon > 0$  where  $f(\mathbf{x}_*)$  is the maximal function value and is unknown. In this work we provide a new approach to the level set estimation problem by relating it to recent adaptive experimental design methods for linear bandits in the Reproducing Kernel Hilbert Space (RKHS) setting. We assume that  $f$  can be approximated by a function in the RKHS up to an unknown misspecification and provide novel algorithms for both the implicit and explicit cases in this setting with strong theoretical guarantees. Moreover, in the linear (kernel) setting, we show that our bounds are nearly optimal, namely, our upper bounds match existing lower bounds for threshold linear bandits. To our knowledge this work provides the first instance-dependent, non-asymptotic upper bounds on sample complexity of level-set estimation that match information theoretic lower bounds.

## 1 INTRODUCTION

The level-set of a function is a subset of its domain where it exceeds a specific value. Level set estimation is the problem of identifying a subset that approximates the true level-set based on a finite set of potentially noisy function evaluations. As an example, consider the goal of detecting a region in a body of water, such as a channel, that is at least  $20m$  deep for ships to safely pass. Given that we can obtain noisy estimates of depth using a sonar device at the locations of our choosing, where should we measure in order to acquire the most accurate level-set estimation while using as few total measurements as possible? Level-set estimation can also be interpreted as a kind of classification rule. For example, using as few total experiments as possible, we may want to identify all compounds among a given finite set under consideration that have some property (e.g., binding affinity) that exceeds some target threshold.

While level-set estimation is somewhat of a well-studied problem, to date there is a lack of theoretical understanding of the limits and tradeoffs of estimation accuracy and number of measurements. Most algorithms proceed by sequentially and greedily optimizing an *acquisition function* that is constructed using all the measurements observed up to the current time. These heuristics are known to work very well in practice, but their guarantees are ad hoc and, at best, worst-case (minimax). In this work we are interested in understanding the instance-dependent sample complexity of level-set estimation. That is, we would like for an algorithm to output a satisfactory estimate of the level-set as fast as *any* algorithm could for *this particular* instance, not some worst-case instance.

In contrast to prior works that propose a sampling heuristic—usually based on identifying an informative point—and bound its sample complexity, we work backwards. Namely, we first consider an information theoretic lower bound for the level-set estimation problem

that suggests an “optimal” sampling strategy. Because this ideal sampling strategy is a function of the true (unknown) function, it is a priori impossible to realize. Instead, we propose a series of sampling strategies, based on *experimental designs*, that mimic this optimal sampling strategy given the information available at the current time. By the end, these strategies provably achieve the optimal sample complexity with minimal overhead. Furthermore, we show that our sampling strategy leads to an upper bound on the sample complexity that is tighter than those in the existing literature. In what follows, we first formally state the problem and our desired objectives. We then review the related work in context before proceeding to our lower bounds and algorithms. We finish with experiments contrasting with existing work.

### 1.1 Problem Statement

We assume there exists an unknown function  $f : \mathbb{R}^d \rightarrow [-B, B]$  and a subset of allowable sampling locations  $\mathcal{X} \subset \mathbb{R}^d$  which span  $\mathbb{R}^d$ . Though the function  $f$  is unknown, we may query its value for any  $\mathbf{x} \in \mathcal{X}$  and receive a noisy estimate  $f(\mathbf{x}) + \eta$  where  $\eta$  is iid,  $\mathbb{E}[\eta] = 0$ , and  $\mathbb{E}[\eta^2] \leq \sigma^2$ . We define two objectives.

**Explicit Level Set Estimation:** Given a specified threshold  $\alpha \in \mathbb{R}$ , the goal is to identify  $G_\alpha := \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) > \alpha\}$ .

**Implicit Level Set Estimation:** Let  $\mathbf{x}_* \in \arg \max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ . Given  $\epsilon > 0$ , the goal is to identify  $G_\epsilon := \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) > (1 - \epsilon)f(\mathbf{x}_*)\}$ <sup>1</sup>.

Consider an algorithm that at each time  $t$  selects an arm  $\mathbf{x}_t \in \mathcal{X}$  that is measurable with respect to a  $\sigma$ -algebra  $\mathcal{F}_{t-1} = \sigma(\mathbf{x}_1, y_1, \dots, \mathbf{x}_{t-1}, y_{t-1})$  and receives a value  $y_t = f(\mathbf{x}_t) + \eta_t$ . To be precise, we say that an algorithm is *PAC- $\delta$*  for the explicit (respectively implicit) level set problem if it stops at a time  $T_\delta$  which is measurable with respect to the filtration  $(\mathcal{F}_t)_{t \geq 1}$  and returns  $G_\alpha$  (and in the implicit setting returns  $G_\epsilon$ ) with probability at least  $1 - \delta$ . If  $f(\mathbf{x})$  is very close to the threshold, it may take an enormous number of samples to determine whether it is above or below the threshold, so in practice we introduce a  $\tilde{\beta} \geq 0$  tolerance that ensures that any learner has a finite sample complexity (see theorems) and allows for misclassification of points very near to the threshold. But in the discussion that follows, assume that  $f(x)$  is bounded away from the threshold.

Our approach is based on modeling  $f$  in a Reproducing Kernel Hilbert Space (RKHS)  $\mathcal{H}$ . Let  $\phi : \mathbb{R}^d \mapsto \mathcal{H}$  be the “feature map” associated with the RKHS. Since

$|f(\mathbf{x})| \leq B$  for all  $\mathbf{x} \in \mathcal{X}$ , there exists a  $\theta_* \in \mathcal{H}$  and a scalar  $h \geq 0$  such that  $\max_{\mathbf{x} \in \mathcal{X}} |f(\mathbf{x}) - \langle \theta_*, \phi(\mathbf{x}) \rangle_{\mathcal{H}}| \leq h$ . When  $h = 0$ ,  $f \in \mathcal{H}$ , and in general we allow  $h \geq 0$  (typically small) in the interest of generality. Our sample complexity bounds will depend on  $h$  and  $\|\theta_*\|_{\mathcal{H}}$  which we denote  $\|\theta_*\|$ . If  $h$  is small, then  $f$  is well approximated as a linear function of the feature maps  $\phi(\mathbf{x})$ . We refer to the case when  $h > 0$  as being *misspecified* and otherwise when  $h = 0$  as being *well-specified*. This class of functions is frequently used for level-set estimation because it is often sufficiently rich to model real-world functions but also contains enough structure to quantify the uncertainty of generalizing a learned function to unmeasured locations. One note of departure from the existing literature is that we do not assume the unknown function is precisely captured by a function in an RKHS, only that it is well approximated by one (i.e., the misspecified setting). In the discussion that follows, we additionally assume  $|\mathcal{X}| < \infty$  for simplicity since in practice given an arbitrary bounded domain we can replace  $\mathcal{X}$  with a finite cover.

## 2 RELATED WORKS

The level-set estimation problem naturally connects to several related ideas in Bayesian optimization and multi-armed bandits. In the former setting, methods tend to sample greedily according to an acquisition function that seeks to minimize the uncertainty of the learner about the level set. The first work on level set estimation that employed the use of Gaussian processes and introduced the Straddle heuristic is due to [Bryan et al. \(2005\)](#). These ideas were further developed in [Gotovos \(2013\)](#) which proposed the LSE and LSE-imp algorithms for explicit and implicit level set respectively. They provide a theoretical guarantee on the sample complexities of LSE and LSE-imp, and as we will show below, our sample complexity is always at least as good as their stated bounds. [Bogunovic et al. \(2016\)](#) further connected Bayesian optimization with level set estimation and considered the setting of heteroscedastic noise. The work of [Shekhar and Javidi \(2019\)](#) focuses on the level-set problem in a continuous domain, and provides an algorithm that maintains a notion of uncertainty over regions, providing a potentially improved computational complexity, along with tighter sample complexity bounds compared to LSE for certain kernels and smoothness assumptions. The work of [Zanette et al. \(2018\)](#) reposes level-set estimation as a classification problem and introduces a novel acquisition function. [Iwazaki et al. \(2020\)](#) extends the work of [Zanette et al. \(2018\)](#) to improve model robustness in quality control applications. [Bogunovic \(2019\)](#); [Vakili et al. \(2021\)](#) demonstrate frequentist guarantees for Gaussian process algorithms. ([Bect et al., 2012](#);

<sup>1</sup>For ease of exposition, we assume  $f(\mathbf{x}_*) \geq 0$ . This is easily removed by taking  $\epsilon < 0$  if  $f(\mathbf{x}_*) < 0$ .

Azzimonti et al., 2021) employ a sequential experimental design approaches for estimating failure probability given a threshold form a density that is expensive to evaluate. (Chevalier et al., 2014) proposes a kriging-based approach for the same problem. This line of work is also related to Gaussian Process Bandits, namely the GPUCB algorithm and improved variants (Srinivas et al., 2009; Chowdhury and Gopalan, 2017; Valko et al., 2013). Ha et al. (2020) introduces a Bayesian Neural Network approach for active level set estimation using Monte Carlo dropout techniques. Table 1 in the appendix summarizes the results we are aware of in the Gaussian process setting.

In the multi-armed and linear bandit setting, the explicit level set estimation problem is related to threshold bandits where one seeks to find all arms above an explicit threshold (Locatelli et al., 2016; Jamieson and Jain, 2018; Degenne et al., 2020). The approach of Degenne et al. (2020), would provide an asymptotically optimal algorithm in the linear setting, however we are not aware of any other works that provide an optimal finite-time guarantee. The implicit level set problem in the standard multi-armed bandit setting is equivalent to the *multiplicative all- $\epsilon$*  problem introduced by Mason et al. (2020). Algorithm 2 recovers the sample complexities of the instance-optimal (ST)<sup>2</sup> algorithm given there. Finally, our experimental design techniques are inspired by Soare et al. (2014); Fiez et al. (2019), and especially the recent work of Camilleri et al. (2021) that introduces the RIPS estimator which we use to perform experimental design in an RKHS.

### 3 EXPLICIT LEVEL SET ESTIMATION

In recent years, adaptive experimental design has arisen as a popular paradigm for active learning in structured settings, for example in linear bandits and RKHS (Soare et al., 2014; Fiez et al., 2019; Camilleri et al., 2021), and we adapt these ideas for the level set problem. To motivate this paradigm, in the following example we focus on the well-specified linear case where  $\phi(\mathbf{x}) = \mathbf{x}, \beta = 0, h = 0$  where we recall  $h$  denotes the misspecification and  $\tilde{\beta}$  denotes the error tolerance as defined in Section 1.1. Imagine we have access to a collection of  $n$ -measurements  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  and let  $\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^d} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \theta)^2$  be the least squares estimator. Standard results show that with probability greater than  $1 - \delta$ , we have for all  $\mathbf{x} \in \mathcal{X}$  simultaneously

$$|\mathbf{x}^\top (\hat{\theta} - \theta_*)| \leq \|\mathbf{x}\|_{(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top)^{-1}} \sqrt{\frac{2 \log(2|\mathcal{X}|/\delta)}{n}},$$

where the additional factor of  $|\mathcal{X}|$  in the logarithm arises from a union bound over  $\mathcal{X}$ . In particular, if our

data is chosen so that for each arm  $\mathbf{x} \in \mathcal{X}$

$$|\mathbf{x}^\top \theta_* - \alpha| > \|\mathbf{x}\|_{(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top)^{-1}} \sqrt{\frac{2 \log(2|\mathcal{X}|/\delta)}{n}}, \quad (1)$$

we see that for any  $\mathbf{x}$  such that  $\mathbf{x}^\top \theta_* > \alpha$ ,

$$\mathbf{x}^\top \hat{\theta} = \mathbf{x}^\top \theta_* + \mathbf{x}^\top (\hat{\theta} - \theta_*) > \mathbf{x}^\top \theta_* - |\mathbf{x}^\top \theta_* - \alpha| > \alpha.$$

The first inequality stems from equation (1) where we have sampled such that the error  $\mathbf{x}^\top (\hat{\theta} - \theta_*)$  is less than the margin to the threshold  $|\mathbf{x}^\top \theta_* - \alpha|$ . Hence, if  $\mathbf{x}^\top \theta_* > \alpha$  then  $\mathbf{x}^\top \hat{\theta} > \alpha$ . This same argument may be repeated for  $\mathbf{x} : \mathbf{x}^\top \theta_* < \alpha$ . Therefore  $\{\mathbf{x} : \mathbf{x}^\top \hat{\theta} > \alpha\} = \{\mathbf{x} : \mathbf{x}^\top \theta_* > \alpha\} = G_\alpha$ , i.e. we have a high probability guarantee that we return the correct set of arms above the threshold. Letting  $\lambda_x = n_x/n$  be the proportion of times we sample  $\mathbf{x} \in \mathcal{X}$ , we see that equation (1) is equivalent to

$$n \geq \max_{\mathbf{x} \in \mathcal{X}} \frac{\|\mathbf{x}\|^2_{(\sum_{\mathbf{x} \in \mathcal{X}} \lambda_x \mathbf{x} \mathbf{x}^\top)^{-1}}}{(\theta_*^\top \mathbf{x} - \alpha)^2}. \quad (2)$$

In particular, this implies that to achieve a good sample complexity we can minimize the right side of this expression over all possible distributions  $\lambda \in \Delta_{\mathcal{X}}$  where  $\Delta_{\mathcal{X}} = \{\lambda \in \mathbb{R}^{|\mathcal{X}|} : \sum_{\mathbf{x} \in \mathcal{X}} \lambda_x = 1, \lambda_x \geq 0 \forall \mathbf{x}\}$ . Indeed as the following theorem shows, this gives a lower bound on this problem.

**Theorem 3.1.** *Assume  $\eta_t \stackrel{iid}{\sim} \mathcal{N}(0, 1) \forall t$ . In the well-specified linear setting when  $\phi(\mathbf{x}) = \mathbf{x}$  and  $f(\mathbf{x}) = \theta_*^\top \mathbf{x}$ , for any  $\delta > 0$ , any PAC- $\delta$  algorithm with stopping time  $T_\delta$  that returns the set  $G_\alpha$  with probability at least  $1 - \delta$  must satisfy*

$$\frac{\mathbb{E}[T_\delta]}{\log(1/2.4\delta)} \geq 2 \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X}} \frac{\|\mathbf{x}\|_{A(\lambda)}^2}{(\theta_*^\top \mathbf{x} - \alpha)^2}$$

where  $A(\lambda) := \sum_{\mathbf{x} \in \mathcal{X}} \lambda_x \mathbf{x} \mathbf{x}^\top$ .

**Remark.** We prove this result for completeness in the appendix using ideas from Fiez et al. (2019). A similar result has appeared previously in the Appendix of Degenne et al. (2020) which also shows its tightness.

As a concrete interpretation of the lower bound, consider the case where  $\mathbf{x}_i = \mathbf{e}_i$ , the  $i^{\text{th}}$  standard basis vector. Then the mean of arm  $i$  is  $\theta_*^\top \mathbf{e}_i = [\theta_*]_i$ , the  $i^{\text{th}}$  entry of  $\theta_*$ . This setting removes all structure by making the mean of each point independent of the others, and we may solve the optimization in Theorem 3.1 in closed form. Namely, the fraction of samples given to arm  $i$ , denoted  $\lambda^{(i)} \propto ([\theta_*]_i - \alpha)^{-2} \log(1/\delta)$  is proportional to its inverse gap squared. This leads to a lower bound of  $\mathbb{E}[\tau_\delta] \geq \sum_{i=1}^n ([\theta_*]_i - \alpha)^{-2} \log(1/\delta)$  with matches the known lower bounds from (Jamieson and Jain, 2018; Locatelli et al., 2016) which are specific to this setting.

We now operationalize this lower bound to provide an algorithm for level set estimation that has a nearly matching upper bound. In the following sections, we will explain our algorithm and the adaptations necessary to handle the general setting of the RKHS.

### 3.1 Algorithm

Motivated by this lower bound, we now provide an experimental design approach in the general case. In this setting, we recall the feature map  $\phi : \mathbb{R}^d \mapsto \mathcal{H}$  and  $h \geq 0$  represents the possibly nonzero misspecification level. Despite these changes, the same intuition from the linear case in Theorem 3.1 applies. We have a set of vectors  $\phi(\mathbf{x})_1, \dots, \phi(\mathbf{x}_n) \in \mathcal{H}$  and an unknown parameter vector  $\theta_* \in \mathcal{H}$  such that  $f(\mathbf{x}) \approx \theta_*^\top \phi(\mathbf{x})$ . Ideally, we would sample according to a distribution  $\lambda_*$  that achieves the minimum in the lower bound in Theorem 3.1, however this is not possible since  $\lambda_*$  depends on the a priori unknown  $\theta_*$ . Instead, we approximate this distribution by solving a series of designs based on the information we have thus far. Furthermore, we allow for a tolerance  $\tilde{\beta} \geq 0$  reflecting the fact that depending on the setting, practitioners may be satisfied with an approximate solution if it requires fewer samples to learn.

Our approach, MELK (Misspecified Explicit Level set via Kernelization), for the generalized RKHS setting is given in Algorithm 1. MELK proceeds in phases. To keep track of the points it has identified so far, MELK maintains two sets: 1)  $\hat{G}_t$  is the set of all points that up to round  $t$  have been declared as being in  $G_\alpha$  by MELK, that is  $f(\mathbf{x}) > \alpha$ . 2)  $\hat{B}_t$  is the set of all points declared as being in  $G_\alpha^c$ . The remaining, uncertain points are *active* and in the set  $\mathcal{A}_t$ . Motivated by the lower bound from the linear setting, it then computes the experimental design:  $\lambda_t = \arg \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{A}_t} \|\phi(\mathbf{x})\|_{A^{(\gamma)}(\lambda)}^2$  with  $A^{(\gamma)}(\lambda) := \sum_{\mathbf{x} \in \mathcal{X}} \lambda_{\mathbf{x}} \phi(\mathbf{x}) \phi(\mathbf{x})^\top + \gamma I$  where  $\gamma$  is a necessary regularization in the kernelized (infinite-dimensional) setting. Indeed, the number of samples taken in each round equals  $N_t \approx \min_{\lambda} \max_{\mathbf{x} \in \mathcal{A}_t} \frac{\|\phi(\mathbf{x})\|_{A^{(\gamma)}(\lambda)}^2}{(2^{-t})^2}$  from  $\lambda_t$ . This guarantees that at the end of the round,  $\mathcal{A}_{t+1} \subset \{\mathbf{x} \in \mathcal{X} : |\theta_*^\top \mathbf{x} - \alpha| \leq 2^{-(t+1)}\}$  and, we can interpret our design as an approximation to the lower bound on the points that are remaining. MELK declares that  $\mathbf{x} \in G_\alpha$  if  $\hat{\theta}^\top \phi(\mathbf{x}) - 2^{-t} \gtrsim \alpha$  and adds  $\mathbf{x}$  to the set  $\hat{G}_t$ . Similarly, MELK adds  $\mathbf{x}$  to declares  $\mathbf{x} \in G_\alpha^c$  and adds  $\mathbf{x}$  to  $\hat{B}_t$  if  $\hat{\theta}^\top \phi(\mathbf{x}) + 2^{-t} \lesssim \alpha$ . Finally, MELK terminates when either all arms have been added to the sets  $\hat{G}_t$  or  $\hat{B}_t$  or when  $t \gtrsim \log_2(1/\tilde{\beta})$  and it has achieved the practitioner’s desired tolerance of  $\tilde{\beta}$ .

MELK leverages a Robust Inverse Propensity Scoring

(RIPS) estimator introduced in Camilleri et al. (2021) and reviewed in Appendix C. Previous works in linear bandits have utilized rounding procedures for sampling followed by ordinary least squares that are not applicable in the infinite dimensional setting. Instead, the RIPS estimator appeals to an inverse propensity score estimator plus robust mean estimation. We state the guarantee of the RIPS estimator below<sup>2</sup>.

**Theorem 3.2** (Theorem 1, (Camilleri et al., 2021)). *Consider the model  $y = \langle \phi(\mathbf{x}), \theta_* \rangle_{\mathcal{H}} + \zeta_{\mathbf{x}} + \eta$  for misspecification  $|\zeta_{\mathbf{x}}| \leq h$  where it is assumed that  $|\langle \phi(\mathbf{x}), \theta_* \rangle_{\mathcal{H}} + \zeta_{\mathbf{x}}| \leq B$ ,  $\mathbb{E}[\eta] = 0$ , and  $\mathbb{E}[\eta^2] \leq \sigma^2$ . Fix any finite sets  $\mathcal{X} \subset \mathbb{R}^d$  and  $\mathcal{V} \subset \mathcal{H}$ , feature map  $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ , number of samples  $\tau$ , regularization  $\gamma > 0$ , and distribution  $\lambda \in \Delta_{\mathcal{X}}$ . If  $\tau \geq 2 \log(|\mathcal{V}|/\delta)$  then with probability at least  $1 - \delta$ , RIPS returns  $\hat{\theta}$  satisfying*

$$\max_{\mathbf{v} \in \mathcal{V}} \frac{|\langle \hat{\theta}, \mathbf{v} \rangle - \langle \theta_*, \mathbf{v} \rangle|}{\|\mathbf{v}\|_{A^{(\gamma)}(\lambda)}^{-1}} \leq 2\sqrt{\gamma} \|\theta_*\| + 2h + 4\sqrt{\frac{(B^2 + \sigma^2)}{\tau} \log\left(\frac{2|\mathcal{V}|}{\delta}\right)}.$$

**Computational Considerations.** We note briefly that while we state the optimal design in terms of the potentially infinite dimensional  $\phi(\mathbf{x})$  for clarity, we never explicitly compute  $\phi(\mathbf{x})$  and instead resort to the kernel trick (see Appendix G). Furthermore the design can be computed using first order optimization methods, such as Frank-Wolfe (Lattimore and Szepesvári, 2020; Todd, 2016). The total computational cost of each design is  $\text{poly}(|\mathcal{X}|)$ . Though these designs can be expensive to compute, this is done very rarely by the algorithm. In particular, for  $T$  total samples drawn by MELK, the design is computed  $O(\log_2(T))$  times leading to an overall computational cost of  $O(\text{poly}(|\mathcal{X}|) \log_2(T))$  for computing the design. By contrast, any algorithm that computes an acquisition function at every sample suffers computational complexity  $\Omega(T)$  for the design. Furthermore, for Gaussian process approaches, the added cost of computing posterior means and variances leads to an overall computational cost of either  $\Omega(\text{poly}(|\mathcal{X}|)T)$  or  $\Omega(|\mathcal{X}| \text{poly}(T))$  depending on implementation for computing acquisition functions. We focus on the complexity of computing the design and acquisition functions as this is frequently the core computational bottleneck of algorithms for level set estimation and the complexity of drawing samples is usually negligible by comparison. Hence, when many samples are drawn, MELK can be significantly more efficient than past approaches.

<sup>2</sup>We carefully detail RIPS from (Camilleri et al., 2021) as it is important for understanding the behavior of the algorithms we present, but the experimental designs we propose are not consequences of that work.



---

**Algorithm 1** MELK: Misspecified Explicit Level set via Kernelization
 

---

**Require:** Arms  $\mathcal{X}$ ,  $\phi$ ,  $\sigma \geq 0$ ,  $\delta > 0$ ,  $\gamma \geq 0$ , threshold  $\alpha$ , tolerance  $\tilde{\beta}$

- 1:  $t \leftarrow 1$ ,  $\hat{G}_1 \rightarrow \emptyset$ ,  $\hat{B}_1 \leftarrow \emptyset$ ,  $\mathcal{A}_1 \leftarrow \mathcal{X}$
  - 2: **while**  $|\hat{G}_t \cup \hat{B}_t| < |\mathcal{X}|$  and  $t \leq \lceil \log_2(4/\tilde{\beta}) \rceil$  **do**
  - 3:      $\delta_t \leftarrow \delta/2t^2$
  - 4:     Let  $\lambda_t \in \Delta_{\mathcal{X}}$  minimize  $g(\lambda; \mathcal{A}_t; \gamma)$  where
 
$$g(\lambda; \mathcal{V}; \gamma) := \max_{\mathbf{x} \in \mathcal{V}} \|\phi(\mathbf{x})\|_{A^{(\gamma)}(\lambda)}^2$$
  - 5:      $q_t \leftarrow 16 \cdot 2^{2t} g(\lambda_t; \mathcal{A}_t; \gamma) (B^2 + \sigma^2) \log(2t^2 |\mathcal{X}|^2 / \delta)$
  - 6:
  - 7:     Set  $N_t \leftarrow \lceil \max\{q_t, 2 \log(|\mathcal{X}|/\delta)\} \rceil$  and sample  $x_1, \dots, x_{N_t}$  observing noisy function values  $y_1, \dots, y_{N_t}$  according to  $\lambda_t$ .
  - 8:      $\hat{\theta}_t \leftarrow \text{RIPS}(\mathcal{A}_t, \{\mathcal{A}^{(\gamma)}(\lambda_t)^{-1} \phi(x_i) y_i\}_{i=1}^{N_t})$ , Alg 3 in Appendix C
  - 9:     **for**  $\mathbf{x} \in \mathcal{A}_t$  **do**
  - 10:         **if**  $\hat{\theta}_t^T \phi(\mathbf{x}) < \alpha - 2 \cdot 2^{-t}$  **then**
  - 11:              $\hat{B}_{t+1} \leftarrow \mathbf{x}$
  - 12:              $\mathcal{A}_{t+1} \leftarrow \mathcal{A}_t \setminus \{\mathbf{x}\}$
  - 13:         **else if**  $\hat{\theta}_t^T \phi(\mathbf{x}) > \alpha + 2 \cdot 2^{-t}$  **then**
  - 14:              $\hat{G}_{t+1} \leftarrow \hat{G}_t \cup \{\mathbf{x}\}$
  - 15:              $\mathcal{A}_{t+1} \leftarrow \mathcal{A}_t \setminus \{\mathbf{x}\}$
  - 16:      $t \leftarrow t+1$
- return**  $\hat{\mathcal{R}} := \mathcal{X} \setminus \hat{B}_t$
- 

### 3.2 Optimal Sample Complexity for Explicit Level Set Estimation

Next we state MELK's complexity, deferring constants and doubly logarithmic factors to the appendix for readability.

**Theorem 3.3.** Fix  $\delta > 0$ , threshold  $\alpha > 0$ , tolerance  $\tilde{\beta}$ , and regularization  $\gamma \geq 0$ . Define  $\Delta_{\min}(\alpha) := \min_{\mathbf{x} \in \mathcal{X}} |\phi(\mathbf{x})^T \theta_* - \alpha|$ . Define also

$$\bar{\beta}(\alpha) = \min \left\{ \beta > 0 : 4(\sqrt{\gamma} \|\theta_*\| + h) \times \left( 2 + \sqrt{\min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X} : |\phi(\mathbf{x})^T \theta_* - \alpha| \leq \beta} \|\phi(\mathbf{x})\|_{A^{(\gamma)}(\lambda)}^2} \right) \leq \beta \right\}.$$

With probability at least  $1 - \delta$ , MELK returns a set  $\hat{\mathcal{R}}$  at time  $T_{\delta}$  such that

$$\hat{\mathcal{R}} \supseteq \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \geq \alpha + \bar{\beta}(\alpha)\}$$

$$\text{and } \hat{\mathcal{R}} \subseteq \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \geq \alpha - \tilde{\beta} - \bar{\beta}(\alpha)\}$$

and for any  $\alpha, \tilde{\beta}$  such that  $\max(\Delta_{\min}(\alpha), \tilde{\beta}) \geq \bar{\beta}(\alpha)$

$$T_{\delta} \leq (B^2 + \sigma^2) \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X}} \frac{\|\phi(\mathbf{x})\|_{A^{(\gamma)}(\lambda)}^2}{\max\{(\phi(\mathbf{x})^T \theta_* - \alpha)^2, \tilde{\beta}^2\}} \times \log((\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \log(|\mathcal{X}| \delta^{-1}).$$

We now contextualize the result of our theorem. In the well specified setting with  $\phi(\mathbf{x}) = \mathbf{x}$ ,  $h = 0$ ,  $\tilde{\beta} = 0$ , and  $\gamma = 0$  MELK will terminate and return  $G_{\alpha}$  in a time

$$T_{\delta} \lesssim (B^2 + \sigma^2) \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X}} \frac{\|\mathbf{x}\|_{A(\lambda)}^2}{(\mathbf{x}^T \theta_* - \alpha)^2} \log(\Delta_{\min}^{-1}) \log\left(\frac{|\mathcal{X}|}{\delta}\right)$$

samples which nearly matches the rate suggested by the linear lower bound in Theorem 3.1. The added factor of  $\log(|\mathcal{X}|)$  stems from a union bound, while the dependence on  $\log(\Delta_{\min}^{-1})$  is an additional overhead incurred as MELK builds up an estimate of the optimal sample allocation over rounds. We visualize this estimation process in Figure 1 in the experiments.

In the more general misspecified setting when  $h > 0$ , we cannot expect to return  $G_{\alpha}$  exactly and  $\bar{\beta}(\alpha)$  characterizes the limit of how well one can estimate  $f(\mathbf{x})$ . Hence,  $\mathbf{x}$ 's with gaps smaller than  $\bar{\beta}(\alpha)$  cannot reliably be detected by MELK. To better understand this quantity, note that for any  $\gamma' \in \mathbb{R}$  if we run MELK with  $\gamma = \gamma'/T$ , Lemma 2 of Camilleri et al. (2021) can be used to show that  $\bar{\beta}(\alpha) \lesssim (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{\Gamma_T}$  where  $\Gamma_T := \sup_{\lambda \in \Delta_{\mathcal{X}}} \log \det(TA^{(0)}(\lambda) + \gamma'I)$  is the *maximum information gain* as defined by Srinivas et al. (2009); Gotovos (2013); Bogunovic et al. (2016). Additionally, it can be shown that  $\Gamma_T \leq d_{eff}$ , where  $d_{eff}$  is the *effective dimension* of  $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n) \in \mathcal{H}$  as defined in Alaoui and Mahoney (2014); Derezhinski et al. (2020). In particular, to ensure that MELK correctly identifies all points that are at least some gap  $\Delta > h$  away from the threshold, then we can choose  $\gamma$  so that  $\Delta > (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{\Gamma_T}$ . In practice we find that  $\gamma = 1/T$  works well. Finally, the user may additionally set a tolerance  $\tilde{\beta} > 0$ . In this case, we err on the side of potentially returning extra arms that are not in  $G_{\alpha}$  and show that the returned set  $\hat{\mathcal{R}}$  contains all  $\mathbf{x}$  such that  $f(\mathbf{x}) > \alpha + \bar{\beta}(\alpha)$  and none such that  $f(\mathbf{x}) < \alpha - \tilde{\beta} - \bar{\beta}(\alpha)$ . If however, a more selective criteria is desired, the following remark characterizes the output if  $\hat{G}_t$  is returned instead.

**Remark.** If MELK instead returns  $\hat{\mathcal{R}} = \hat{G}_t$  then with probability at least  $1 - \delta$   $\hat{\mathcal{R}} \supseteq \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \geq \alpha + \tilde{\beta} + \bar{\beta}(\alpha)\}$  and  $\hat{\mathcal{R}} \subseteq \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \geq \alpha - \bar{\beta}(\alpha)\}$ .

**Contrast with Existing Approaches.** The experimental design based sampling approach is a departure from past work on level set estimation. As opposed to constructing an acquisition function and then bounding the sample complexity of the resulting algorithm as past works have done, we instead begin with an oracle sampling scheme that arises from a lower bound and attempt to design a practical sampling scheme that matches it as more data is collected. In what follows, we compare the guarantees of MELK to the prior art such as Gotovos (2013); Shekhar and Javidi (2019); Bogunovic et al. (2016). As a technical point, we note that

these past results are specialized to the Gaussian process setting where a prior on  $f$  is known. By contrast, our work makes no assumption of a prior distribution. [Bogunovic \(2019\)](#); [Vakili et al. \(2021\)](#) achieve similar guarantees for the frequentist setting. Ignoring these technicalities, our results are tighter than what were previously known.

The past state of the art sample complexities all guarantee that algorithms terminate at the smallest time  $T$  satisfying  $T \gtrsim \Gamma_T \Delta_{\min}(\alpha)^{-2}$  up to log factors (cf. Thm 1 of ([Gotovos, 2013](#)), Cor. 3.1 of ([Bogunovic et al., 2016](#)), Thm 1 of ([Shekhar and Javidi, 2019](#)), etc.). If we run MELK with  $\gamma = \gamma'/T$  then

$$\begin{aligned} & \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x}} \frac{\|\phi(\mathbf{x})\|_{A(\gamma)(\lambda)}^2}{(\phi(\mathbf{x})^T \theta_* - \alpha)^2} \\ & \leq \min_{\lambda \in \Delta_{\mathcal{X}}} \frac{\max_{\mathbf{x}} \|\phi(\mathbf{x})\|_{(A(\lambda) + \gamma I)}^2}{\min_{\mathbf{x}} (\phi(\mathbf{x})^T \theta_* - \alpha)^2} \leq 3\Gamma_T \Delta_{\min}(\alpha)^{-2} \end{aligned}$$

where the final inequality follows from Lemma 2 of [Camilleri et al. \(2021\)](#) and the definition of  $\Delta_{\min}(\alpha)$ .

**Remark.** Combining the above analysis with the result of Theorem 3.3 highlights that MELK likewise terminates at or before a time  $T$  satisfying  $T \gtrsim \Gamma_T \Delta_{\min}(\alpha)^{-2}$ , though it may stop long before this as the above bound employing  $\Gamma_T$  is only tight in the pathological case when  $|\phi(\mathbf{x})^T \theta_* - \alpha| = \Delta_{\min}(\alpha) \forall \mathbf{x} \in \mathcal{X}$ .

**Remark.** The lower bounds of [Scarlett et al. \(2017\)](#); [Cai and Scarlett \(2021\)](#) show that a dependence  $\Omega(\sqrt{\Gamma_T})$  is necessary in the worst case for functions living in an RKHS. Hence, MELK is instance optimal in the linear regime by Theorem 3.1 and at least minimax optimal in general.

## 4 IMPLICIT LEVEL SET ESTIMATION

In the *implicit* level-set problem, for an  $\epsilon \geq 0$  we seek to identify the set  $G_\epsilon = \{\mathbf{x} : f(\mathbf{x}) > (1 - \epsilon)f(\mathbf{x}_*)\}$ . Note that unlike the explicit setting where the threshold  $\alpha$  was a given input to the algorithm, now the equivalent notion of a threshold value  $\alpha$  is equal to  $(1 - \epsilon)f(\mathbf{x}_*)$ , an unknown quantity since it relies on knowledge of the unknown function  $f$ . A naive strategy would be to attempt estimate  $(1 - \epsilon)f(\mathbf{x}_*)$  directly and then apply explicit level-set estimation techniques using this estimated threshold value. Indeed, this is precisely the strategy of past works ([Mason et al., 2020](#); [Gotovos, 2013](#)). Perhaps surprisingly however, it turns out that estimating the threshold is unnecessary and potentially wasteful. Towards developing lower bound to guide an experimental design, we begin with a simple but powerful observation.

**Lemma 4.1.**  $\mathbf{x} \in G_\epsilon \iff \forall \mathbf{x}' \in \mathcal{X} : f(\mathbf{x}) \geq (1 - \epsilon)f(\mathbf{x}')$ . Conversely,  $\mathbf{x} \in G_\epsilon^c \iff \exists \mathbf{x}' : f(\mathbf{x}) < (1 - \epsilon)f(\mathbf{x}')$ .

*Proof.*

$$\begin{aligned} \mathbf{x} \in G_\epsilon & \iff \nexists \mathbf{x}' : (1 - \epsilon)f(\mathbf{x}') > f(\mathbf{x}) \\ & \iff \forall \mathbf{x}' : (1 - \epsilon)f(\mathbf{x}') \leq f(\mathbf{x}) \end{aligned}$$

where the second equivalence holds by definition since  $\mathbf{x}_*$  maximizes  $(1 - \epsilon)f(\mathbf{x}')$  and we have that  $f(\mathbf{x}) > (1 - \epsilon)f(\mathbf{x}_*)$  for any  $\mathbf{x} \in G_\epsilon$ . The statement for  $\mathbf{x} \in G_\epsilon^c$  holds via the negation  $\square$

The following corollary specializes the previous lemma to the well specified case.

**Corollary 4.1.1.** *In the well specified setting where  $h = 0$ ,*

$$\mathbf{x} \in G_\epsilon \iff \forall \mathbf{x}' \in \mathcal{X} : \theta_*^\top (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')) \geq 0$$

and conversely,

$$\mathbf{x} \in G_\epsilon^c \iff \exists \mathbf{x}' : \theta_*^\top (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')) < 0.$$

This lemma highlights that to determine if  $\mathbf{x} \in G_\epsilon$ , one need only check if

$$\theta_*^\top (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')) > 0 \text{ for all } \mathbf{x}' \in \mathcal{X}.$$

In particular, this does not require any estimate of the threshold  $(1 - \epsilon)f(\mathbf{x}_*)$ . Instead, it is only necessary to maintain estimates of ordered pairs of points  $(\mathbf{x}, \mathbf{x}')$  without searching for  $\mathbf{x}_*$  directly. Next, to guide our algorithm design we look to an information-theoretic lower bound.

**Theorem 4.2.** *In the well-specified linear setting when  $\phi(\mathbf{x}) = \mathbf{x}$  and  $f(\mathbf{x}) = \theta_*^\top \mathbf{x}$ , for any  $\delta > 0$ , any algorithm that returns the set  $G_\epsilon$  with probability at least  $1 - \delta$  must satisfy*

$$\begin{aligned} \frac{\mathbb{E}[T_\delta]}{\log(1/2.4\delta)} & \geq 2 \min_{\lambda \in \Delta_{\mathcal{X}}} \max \left\{ \max_{\mathbf{z} \in G_\epsilon} \max_{\mathbf{x}' \in \mathcal{X}} \frac{\|\mathbf{z} - (1 - \epsilon)\mathbf{x}'\|_{A(\lambda)}^2}{(\theta_*^\top (\mathbf{z} - (1 - \epsilon)\mathbf{x}'))^2}, \right. \\ & \left. \max_{\mathbf{x} \in G_\epsilon^c} \min_{\mathbf{x}' \in \mathcal{X}} \frac{\|\mathbf{x} - (1 - \epsilon)\mathbf{x}'\|_{A(\lambda)}^2}{(\theta_*^\top (\mathbf{x} - (1 - \epsilon)\mathbf{x}'))^2} \right\} \end{aligned}$$

where  $T_\delta$  denotes the random stopping time.

Notably, the directions  $\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')$  naturally arise in the lower bound. This suggests an optimal sampling distribution  $\lambda^*$  that achieves the minimum of the inequality in 4.2. As was the case in explicit level set estimation, this sampling distribution also depends on the unknown  $\theta_*$ .

## 4.1 Algorithm

Motivated by the lower bound, we propose Algorithm 2 called MILK (Misspecified Implicit Level set via Kernelization) which proceeds in phases where we attempt to progressively match the optimal distribution from the lower bound as was done by MELK for the explicit setting. The key difference, however is that MILK instead computes a design to optimally estimate  $\theta_*^\top(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))$  rather than  $\theta_*^\top\phi(\mathbf{x})$  as in MELK. Given active set  $\mathcal{A} \subset \mathcal{X} \times \mathcal{X}$  of pairs define,

$$\mathcal{Y}^\epsilon(\mathcal{A}) := \{\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}') : (\mathbf{x}, \mathbf{x}') \in \mathcal{A}\}.$$

The active set in round 1 is initialized as  $\mathcal{A}_1 = \mathcal{X} \times \mathcal{X}$ . MILK keeps track of sets  $\widehat{G}_t \subset \mathcal{X}$  and  $\widehat{B}_t \subset \mathcal{X}$  of arms it believes to be in  $G_\epsilon$  and  $G_\epsilon^c$  and makes use of the RIPS procedure to robustly estimate means. As the algorithm proceeds, in each round  $t$  an optimal design is computed over remaining difference vectors in  $\mathcal{Y}^\epsilon(\mathcal{A}_t)$  and the number of samples  $N_t$  is sufficient to ensure that  $|(\theta_* - \widehat{\theta})^\top(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))| \leq 2^{-t+1}$ . Then for every arm that has not been added to  $\widehat{G}_t$  or  $\widehat{B}_t$ , MILK does the following:

$$\text{if } \exists \mathbf{x}' : \widehat{\theta}^\top((\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')) < 2^{-t}$$

then  $\mathbf{x}$  is added to  $\widehat{B}_t$ . In our proof, we show this condition occurs if and only if there exists a  $\mathbf{x}'$  such that  $\theta_*^\top(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')) < 0$ . If this occurs, all pairs of the form  $(\mathbf{x}, \mathbf{x}')$  or  $(\mathbf{x}', \mathbf{x})$ ,  $\mathbf{x}' \in \mathcal{X}$  are removed from  $\mathcal{A}_t$ <sup>3</sup>. Semantically, if MILK can ensure that  $\mathbf{x}$  is not in  $G_\epsilon$ , then  $\mathbf{x}$  is never sampled again. Otherwise, for any  $\mathbf{x}'$  if  $\widehat{\theta}^\top(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')) > 2^{-t}$ , the single pair  $(\mathbf{x}, \mathbf{x}')$  is removed from  $\mathcal{A}_t$ . An arm  $\mathbf{x}$  is only ever added to  $\widehat{G}_t$  if  $\{(\mathbf{x}, \mathbf{x}'), \mathbf{x}' \in \mathcal{X}\} \cap \mathcal{A}_t = \emptyset$  which occurs when

$$\forall \mathbf{x}' : \exists t' \text{ such that } \widehat{\theta}_{t'}^\top((\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')) > 2^{-t'}.$$

In our proof, we show that this occurs if and only if  $\theta_*^\top(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')) > 0$  for all  $\mathbf{x}' \in \mathcal{X}$  which is both necessary and sufficient by Lemma 4.1. Note that even if  $\mathbf{x}$  has been added to  $\widehat{G}_t$  implying that all pairs  $(\mathbf{x}, \mathbf{x}')$  have been removed from  $\mathcal{A}_t$ ,  $\mathbf{x}$  may be present in other pairs  $(\mathbf{x}', \mathbf{x})$  which can be necessary to determine if  $\mathbf{x}' \in G_\epsilon$ . Finally, the algorithm terminates when either every arm has been added to either  $\widehat{G}_t$  or  $\widehat{B}_t$  or it has reached a round  $t \gtrsim \log_2(1/\tilde{\beta})$  when the desired tolerance  $\tilde{\beta}$  is achieved.

## 4.2 Theoretical Guarantees

Next we state MILK's complexity, again deferring constants and doubly logarithmic factors to the appendix for readability.

<sup>3</sup>We assume that pairs are ordered, i.e.  $(\mathbf{x}, \mathbf{x}') \neq (\mathbf{x}', \mathbf{x})$  for  $\mathbf{x} \neq \mathbf{x}'$ .

---

### Algorithm 2 MILK: Misspecified Implicit Level set via Kernelization

---

**Require:** Arms  $\mathcal{X}$ ,  $\phi$ ,  $\delta > 0$ ,  $\epsilon > 0$ ,  $\gamma \geq 0$ , tolerance  $\tilde{\beta}$

- 1:  $t \leftarrow 1$ ,  $\widehat{G}_1 \rightarrow \emptyset$ ,  $\widehat{B}_1 \leftarrow \emptyset$ ,  $\mathcal{A}_1 \leftarrow \{(\mathbf{x}, \mathbf{x}'), \mathbf{x}, \mathbf{x}' \in \mathcal{X}\}$
  - 2: **while**  $|\widehat{G}_t \cup \widehat{B}_t| < |\mathcal{X}|$  and  $t \leq \lceil \log_2(4/\tilde{\beta}) \rceil$  **do**
  - 3:      $\delta_t \leftarrow \delta/2t^2$
  - 4:     Let  $\lambda_t \in \Delta_{\mathcal{X}}$  minimize  $g(\lambda; \mathcal{A}_t; \gamma)$  where
 
$$g(\lambda, \mathcal{V}; \gamma) := \max_{(\mathbf{x}, \mathbf{x}') \in \mathcal{V}} \|\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')\|_{A^{(\gamma)}(\lambda)^{-1}}^2$$
  - 5:      $q_t \leftarrow 16 \cdot 2^{2t} g(\lambda_t; \mathcal{A}_t; \gamma) (B^2 + \sigma^2) \log(2t^2 |\mathcal{X}|^2 / \delta)$
  - 6:
  - 7:     Set  $N_t \leftarrow \lceil \max\{q_t, 2 \log(|\mathcal{X}|/\delta)\} \rceil$  and sample  $x_1, \dots, x_{N_t}$  observing noisy function values  $y_1, \dots, y_{N_t}$  according to  $\lambda_t$ .
  - 8:      $\widehat{\theta}_t \leftarrow \text{RIPS}(\mathcal{Y}^\epsilon(\mathcal{A}_t), \{A^{(\gamma)}(\lambda_t)^{-1} \phi(x_i) y_i\}_{i=1}^{N_t})$
  - 9:     **for**  $(\mathbf{x}, \mathbf{x}') \in \mathcal{A}_t$  **do**
  - 10:         **if**  $\widehat{\theta}_t^\top(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')) < -2 \cdot 2^{-t}$  **then**
  - 11:              $\widehat{B}_{t+1} \leftarrow \mathbf{x}$
  - 12:              $\mathbf{x}$ -pairs  $\leftarrow \{(\mathbf{x}, \mathbf{x}') \text{ and } (\mathbf{x}', \mathbf{x}) | \mathbf{x}' \in \mathcal{X}\}$
  - 13:              $\mathcal{A}_{t+1} \leftarrow \mathcal{A}_t \setminus \mathbf{x}$ -pairs
  - 14:         **else if**  $\widehat{\theta}_t^\top(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')) > 2 \cdot 2^{-t}$
  - 15:              $\mathcal{A}_{t+1} \leftarrow \mathcal{A}_t \setminus \{(\mathbf{x}, \mathbf{x}')\}$
  - 16:         **if**  $\{(\mathbf{x}, \mathbf{x}') | \mathbf{x}' \in \mathcal{X}\} \cap \mathcal{A}_t = \emptyset$  **then**
  - 17:              $\widehat{G}_{t+1} \leftarrow \widehat{G}_t \cup \{\mathbf{x}\}$
  - 18:      $t \leftarrow t + 1$
  - 19:     **return**  $\widehat{\mathcal{R}} := \mathcal{X} \setminus \widehat{B}_t$
- 

**Theorem 4.3.** Fix  $\delta > 0$ ,  $\epsilon > 0$ , tolerance  $\tilde{\beta}$ , and regularization  $\gamma > 0$ . Define  $\Delta_{\min}(\epsilon) = \min_{\mathbf{x}} |\theta_*^\top(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))|$ . Define also

$$\bar{\beta}(\epsilon) = \min_{\beta > 0} \left\{ 4(\sqrt{\gamma} \|\theta_*\| + h) \left( 2 + \sqrt{\min_{\lambda \in \Delta_{\mathcal{X}}} \nu(\lambda, \beta)} \right) \leq \beta \right\},$$

$$\nu(\lambda, \beta) := \max_{\substack{(\mathbf{x}, \mathbf{x}') \in \mathcal{X} \times \mathcal{X} \\ |\theta_*^\top(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))| \leq \beta}} \|\theta_*^\top(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))\|_{A^{(\gamma)}(\lambda)^{-1}}^2.$$

With probability  $1 - \delta$ , MILK returns a set  $\widehat{\mathcal{R}}$  at a time  $T_\delta$  such that

$$\widehat{\mathcal{R}} \supseteq \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \geq (1 - \epsilon)f(\mathbf{x}_*) + \bar{\beta}(\epsilon)\} \text{ and}$$

$$\widehat{\mathcal{R}} \subseteq \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \geq (1 - \epsilon)f(\mathbf{x}_*) - \tilde{\beta} - \bar{\beta}(\epsilon)\}$$

and for any  $\epsilon, \tilde{\beta}$  such that  $\max(\Delta_{\min}(\epsilon), \tilde{\beta}) \geq \bar{\beta}(\epsilon)$

$$T_\delta \leq (B^2 + \sigma^2) H^{\text{MILK}}(\theta_*) \log_2((\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \log\left(\frac{|\mathcal{X}|}{\delta}\right)$$

$$\text{for } H^{\text{MILK}}(\theta_*) = \min_{\lambda \in \Delta_{\mathcal{X}}} \left\{ H_\lambda^{\text{MILK}-G_\epsilon}(\theta_*) \vee H_\lambda^{\text{MILK}-G_\epsilon^c}(\theta_*) \right\},$$

where

$$H_\lambda^{\text{MILK}-G_\epsilon}(\theta_*) := \max_{\mathbf{x} \in G_\epsilon} \max_{\mathbf{x}' \in \mathcal{X}} \frac{\|\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')\|_{A^{(\gamma)}(\lambda)^{-1}}^2}{\max\{(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^\top \theta_*\}^2, \tilde{\beta}^2},$$

$$\text{and } H_\lambda^{\text{MILK-}G_\epsilon^c}(\theta_*) := \frac{\max_{\mathbf{x} \in G_\epsilon^c} \max_{\mathbf{x}'} \frac{\|\phi(\mathbf{x}) - (1-\epsilon)\phi(\mathbf{x}')\|_{A^{(\gamma)(\lambda)}^{-1}}^2}{\max\{((\phi(\mathbf{x}) - (1-\epsilon)\phi(\mathbf{x}'))^\top \theta_*)^2, \tilde{\beta}^2\}}}{}$$

The statement of Theorem 4.3 for MILK is similar that of 3.3 for MELK. In the well specified case when  $\tilde{\beta} = 0$ , MILK returns  $G_\epsilon$  exactly at a time  $T_\delta$  that satisfies

$$T_\delta \lesssim (B^2 + \sigma^2) H^{\text{MILK}}(\theta_*) \log_2(\Delta_{\min}(\epsilon)) \log(|\mathcal{X}|\delta^{-1})$$

In this case, however,  $H^{\text{MILK}}(\theta_*)$  is a maximum of two different complexity terms.  $H_\lambda^{\text{MILK-}G_\epsilon}$  represents the complexity of identifying all  $\mathbf{x} \in G_\epsilon$ . Similarly,  $H_\lambda^{\text{MILK-}G_\epsilon^c}$  represents the complexity of identifying all  $\mathbf{x} \in G_\epsilon^c$ . Similar to the explicit setting, in the misspecified case when  $h > 0$ ,  $\tilde{\beta}(\epsilon)$  similarly represents the limit of how well we can estimate  $f(\mathbf{x})$  for any  $\mathbf{x} \in \mathcal{X}$  and  $\tilde{\beta}$  allows for an additional tolerance such that MILK detects all  $\mathbf{x}$  for which  $f(\mathbf{x}) > (1-\epsilon)f(\mathbf{x}_*) + \tilde{\beta}(\epsilon)$  and none worse than  $f(\mathbf{x}) < (1-\epsilon)f(\mathbf{x}_*) - \tilde{\beta}(\epsilon) - \tilde{\beta}$ . The following remark addresses the setting where MILK returns  $\hat{G}_t$  instead.

**Remark:** If the algorithm instead returns  $\hat{\mathcal{R}} = \hat{G}_t$ , then with probability at least  $1 - \delta$

$$\begin{aligned} \hat{\mathcal{R}} &\supseteq \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \geq (1-\epsilon)f(\mathbf{x}_*) + \tilde{\beta} + \tilde{\beta}(\epsilon)\} \text{ and} \\ \hat{\mathcal{R}} &\subseteq \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \geq (1-\epsilon)f(\mathbf{x}_*) - \tilde{\beta}(\epsilon)\}. \end{aligned}$$

### Comparison with the Lower Bound

The complexity term  $H^{\text{MILK}}(\theta_*)$  naturally breaks into two terms.  $H^{\text{MILK-}G_\epsilon}(\theta_*)$  represents the complexity of finding arms in  $G_\epsilon$  and it matches a corresponding term in the lower bound.  $H^{\text{MILK-}G_\epsilon^c}(\theta_*)$  represents the complexity of removing arms in  $G_\epsilon^c$  but is slightly different than the term in the lower bound. As a consequence of Theorem 4.1 of Mason et al. (2020) however, one can show the term given in the lower bound for  $\mathbf{x} \in G_\epsilon^c$  is not achievable except asymptotically as  $\delta \rightarrow 0$  in general. Instead, the problem of implicit level set estimation reduces to the problem of all  $\epsilon$ -good arm identification in multi-armed bandits studied by Mason et al. (2020) when  $\phi(\mathbf{x}) = \mathbf{x}$ ,  $h = 0$ , and  $\mathbf{x}_i = e_i$ . We show in the appendix that MILK's sample complexity matches the optimal finite time rate up to logarithmic factors as shown in Mason et al. (2020).

### Contrast with Existing Results

As was shown in the explicit setting, we can show that the sample complexity bound in Theorem 4.3 improves on the current state of the art. Take  $\gamma = \gamma'/T$  for any  $\gamma' \in \mathbb{R}$ . Then we may bound  $H^{\text{MILK-}G_\epsilon}(\theta_*)$  as

$$\min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x}, \mathbf{x}'} \left\{ \frac{\|\phi(\mathbf{x}) - (1-\epsilon)\phi(\mathbf{x}')\|_{A^{(\gamma)(\lambda)}^{-1}}^2}{((\phi(\mathbf{x}) - (1-\epsilon)\phi(\mathbf{x}'))^\top \theta_*)^2} \right\}$$

$$\begin{aligned} &\stackrel{(a)}{\leq} 4 \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x}, \mathbf{x}'} \left\{ \frac{(1-\epsilon)^2 \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_{A^{(\gamma)(\lambda)}^{-1}}^2}{((\phi(\mathbf{x}) - (1-\epsilon)\phi(\mathbf{x}'))^\top \theta_*)^2} \right. \\ &\quad \left. \vee \frac{\epsilon^2 \|\phi(\mathbf{x})\|_{A^{(\gamma)(\lambda)}^{-1}}^2}{((\phi(\mathbf{x}) - (1-\epsilon)\phi(\mathbf{x}'))^\top \theta_*)^2} \right\} \\ &\stackrel{(b)}{\leq} 4 \frac{(1+\epsilon)^2}{\Delta_{\min}(\epsilon)^2} \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x}, \mathbf{x}'} \left\{ \|\phi(\mathbf{x}') - \phi(\mathbf{x})\|_{A^{(\gamma)(\lambda)}^{-1}}^2 \right. \\ &\quad \left. \vee \|\phi(\mathbf{x})\|_{(A^{(\gamma)(\lambda)})^{-1}}^2 \right\} \\ &\leq \frac{8(1+\epsilon)^2}{\Delta_{\min}(\epsilon)^2} \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x}} \|\phi(\mathbf{x})\|_{A^{(\gamma)(\lambda)}^{-1}}^2 \\ &\stackrel{(c)}{\leq} \frac{12(1+\epsilon)^2}{\Delta_{\min}(\epsilon)^2} \Gamma_T \end{aligned}$$

where (a) follows by the triangle inequality, (b) by definition of  $\Delta_{\min}(\epsilon)$  and (c) follows by Lemma 2 of Camilleri et al. (2021). A similar computation follows for  $H^{\text{MILK-}G_\epsilon^c}(\theta_*)$ . Hence, MILK is at most  $O(\Gamma_T \Delta_{\min}^{-2})$  though it can be much tighter as inequality (b) is tight only in the worst case when all gaps are equal. In particular, the result of Theorem 4.3 is tighter than Theorem 2 of Gotovos (2013).

## 5 EXPERIMENTS

In this section, we compare our algorithms to existing baselines in the literature. Additional details of these methods and our experiments are in the Appendix.

**Warm-Up: Optimal Sampling.** In Figure 1 we illustrate the sampling behavior of MELK. We let  $\mathcal{X} = \{(\frac{i}{30}, \frac{j}{30})\}_{i,j=1}^{30}$  and considered the squared exponential kernel  $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2/2\ell^2)$  with parameter  $\ell = 0.1$ . We also chose  $\theta_* \sim \mathcal{N}(0, I_{900})$  and show a contour plot of  $f(\mathbf{x}) = \theta_*^\top \phi(\mathbf{x})$ . The black curve represents the boundary of the  $\alpha = 0$  level set. We plot the sample allocations as the algorithm progresses (taking  $\gamma = 0$ ). The initial distribution is mostly uniform with several sampling modes. In later rounds, the points nearest to the boundary of the level set, given by the black curve are sampled, and eventually, only the points with the smallest gaps (the most difficult regions) receive samples. As the number of samples in round  $t$  is proportional to  $2^{2t}$ , we compute the sum of the designs weighted by the  $2^{2t}$  to show the overall sampling design. Additionally, we plot the asymptotic allocation suggested by Theorem 3.1, namely  $\lambda_* = \arg \min_{\lambda} \max_{\mathbf{x} \in \mathcal{X}} \|\phi(\mathbf{x})\|_{A^{(\gamma)(\lambda)}^{-1}}^2 / (\theta_*^\top \phi(\mathbf{x}) - \alpha)^2$ . In particular, the weighted sum of the designs taken by MELK is nearly identical to  $\lambda_*$ .

**Gaussian Process Level Set Estimation.** For our main empirical evaluation, we focused on the Gaussian Process setting for the explicit level set problem. In the explicit level-set case we compare to LSE (Gotovos,



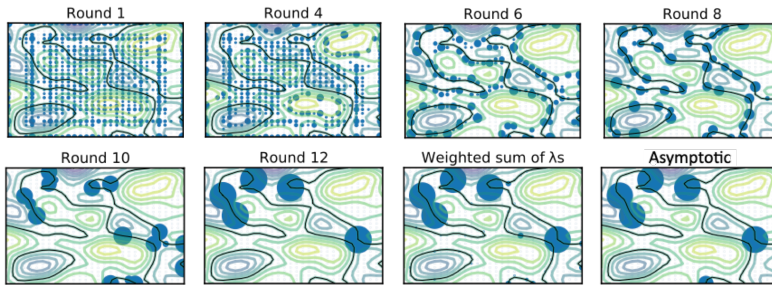


Figure 1: Allocations across rounds for a function  $f(x, y)$  with a threshold of  $\alpha = 0$  shown in black.

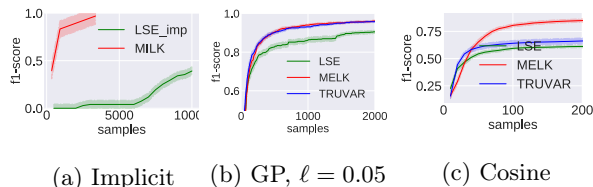


Figure 2: Performance of MELK and MILK versus Gaussian process baseline.

2013) and TruVar (Bogunovic et al., 2016). We drew a function  $f : [0, 1] \rightarrow \mathbb{R}$  from the Gaussian process  $\mathcal{N}(0, k(\mathbf{x}, \mathbf{x}'))$  where the kernel is a squared exponential kernel with parameter  $\ell = .05$  and  $[0, 1]$  was uniformly discretized into 200 points. We assumed that the noise variance was  $\sigma^2 = 1$  (high noise) and the threshold was chosen so that 10% of the function values were above it. In this setting, we implement a batched version of MELK that draws a fixed batch size of samples each round (namely 10) and then recomputes the design. This reflects the practical constraint that experimenters may wish to collect a fixed number of samples at a time rather than a potentially growing amount. To provide a fair comparison to the GP-based methods, we computed a posterior distribution on  $f$  in each round. For each point we replaced our theoretically justified confidence intervals in the RKHS setting with confidence intervals arising from the posterior, namely  $\hat{\mu}_t(\mathbf{x}) \pm \beta^{1/2} \hat{\sigma}_t(\mathbf{x})$  where  $\hat{\mu}_t, \hat{\sigma}_t$  are the posterior mean and standard deviations respectively. As in past works, we take  $\beta^{1/2} = 3$  as theoretically justified choices of  $\beta$  (eg. Theorem 1 of (Srinivas et al., 2009)) tend to be overly conservative. We also took  $\gamma$  dropping like  $1/i$  on the  $i$ -th round we computed the design. We ran 25 repetitions drawing a new choice of  $f$  each run. Figure 2b shows the average F1 score of the set of points each algorithm declares to be in  $G_\alpha$  respectively with bars denoting 1 standard error. Our algorithm performs very similarly to TruVar - an algorithm whose acquisition function samples in a way to reduce the average variance, unlike our method which tries to reduce the maximum variance.

Our second comparison is in Figure 2c: we took  $f(x) = \cos(8\pi x)$ ,  $\ell = .1$ ,  $\sigma = .2$  (low noise regime) and chose the threshold so that 30% of points were above it. We

then considered 700 points uniformly in  $[0, 1]$ . In the appendix, we vary the underlying parameters of  $\ell, \sigma^2$  to demonstrate the performance of these algorithms in different regimes.

**Linear Implicit Case.** We additionally compare against LSE-imp in the linear setting where  $\phi(\mathbf{x}) = \mathbf{x}$  on a benchmark example from the linear bandits literature designed to test the effectiveness of adaptive sampling algorithms (Soare et al., 2014). For  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ , we take  $\mathbf{x}_1 = \mathbf{x}_* = \theta_* = e_1$  and  $\mathbf{x}_2 = e_2$ . The remaining  $\mathbf{x}_3, \dots, \mathbf{x}_n$  are set so that their first two coordinates are  $\cos(\pi/4(1 + \xi))e_1$  and  $\sin(\pi/4(1 + \xi))e_2$  for  $\xi \sim \text{Unif}(-.2, .2)$ . We set the threshold  $\alpha = 0.5$ ,  $n = 100$ , and  $d = 25$ . Though it is far below  $\alpha$ , sampling arm  $\mathbf{x}_2$  provides the most information about which arms exceed the threshold. In this setting, we ran both algorithms with the exact confidence intervals as specified by their respective theoretical guarantees leading to large sample complexities, and we include further details in the appendix. Indeed, we see in 2a that MILK outperforms LSE-imp.

## 6 CONCLUSION

In this work, we provide the first instance optimal algorithms for explicit and implicit level set estimation and provide theoretical and empirical justification for our algorithms. In Appendix A we further explore the potential impacts and limitations of this work.

References

- Alaoui, A. E. and Mahoney, M. W. (2014). Fast randomized kernel methods with statistical guarantees. *arXiv preprint arXiv:1411.0306*.
- Allen-Zhu, Z., Li, Y., Singh, A., and Wang, Y. (2017). Near-optimal design of experiments via regret minimization. In *International Conference on Machine Learning*, pages 126–135. PMLR.
- Azzimonti, D., Ginsbourger, D., Chevalier, C., Bect, J., and Richet, Y. (2021). Adaptive design of experiments for conservative estimation of excursion sets. *Technometrics*, 63(1):13–26.
- Bect, J., Ginsbourger, D., Li, L., Picheny, V., and Vazquez, E. (2012). Sequential design of computer experiments for the estimation of a probability of failure. *Statistics and Computing*, 22(3):773–793.
- Bogunovic, I. (2019). Robust adaptive decision making: Bayesian optimization and beyond. Technical report, EPFL.
- Bogunovic, I., Scarlett, J., Krause, A., and Cevher, V. (2016). Truncated variance reduction: A unified approach to bayesian optimization and level-set estimation. *arXiv preprint arXiv:1610.07379*.
- Bryan, B., Schneider, J., Nichol, R., Miller, C. J., Genovese, C. R., and Wasserman, L. (2005). Active learning for identifying function threshold boundaries. In *NIPS*, pages 163–170. Citeseer.
- Cai, X. and Scarlett, J. (2021). On lower bounds for standard and robust gaussian process bandit optimization. In *International Conference on Machine Learning*, pages 1216–1226. PMLR.
- Camilleri, R., Jamieson, K., and Katz-Samuels, J. (2021). High-dimensional experimental design and kernel bandits. In *International Conference on Machine Learning*, pages 1227–1237. PMLR.
- Chevalier, C., Bect, J., Ginsbourger, D., Vazquez, E., Picheny, V., and Richet, Y. (2014). Fast parallel kriging-based stepwise uncertainty reduction with application to the identification of an excursion set. *Technometrics*, 56(4):455–465.
- Chowdhury, S. R. and Gopalan, A. (2017). On kernelized multi-armed bandits. In *International Conference on Machine Learning*, pages 844–853. PMLR.
- Degenne, R., Ménard, P., Shang, X., and Valko, M. (2020). Gamification of pure exploration for linear bandits. In *International Conference on Machine Learning*, pages 2432–2442. PMLR.
- Derezinski, M., Liang, F., and Mahoney, M. (2020). Bayesian experimental design using regularized determinantal point processes. In *International Conference on Artificial Intelligence and Statistics*, pages 3197–3207. PMLR.
- Fiez, T., Jain, L., Jamieson, K., and Ratliff, L. (2019). Sequential experimental design for transductive linear bandits. *arXiv preprint arXiv:1906.08399*.
- Gotovos, A. (2013). Active learning for level set estimation. Master’s thesis, Eidgenössische Technische Hochschule Zürich, Department of Computer Science,.
- Ha, H., Gupta, S., Rana, S., and Venkatesh, S. (2020). High dimensional level set estimation with bayesian neural network. *arXiv preprint arXiv:2012.09973*.
- Iwazaki, S., Inatsu, Y., and Takeuchi, I. (2020). Bayesian experimental design for finding reliable level set under input uncertainty. *IEEE Access*, 8:203982–203993.
- Jamieson, K. G. and Jain, L. (2018). A bandit approach to sequential experimental design with false discovery control. *Advances in Neural Information Processing Systems*, 31:3660–3670.
- Jun, K.-S., Jain, L., Mason, B., and Nassif, H. (2020). Improved confidence bounds for the linear logistic model and applications to linear bandits. *arXiv preprint arXiv:2011.11222*.
- Katz-Samuels, J., Jain, L., Karnin, Z., and Jamieson, K. (2020). An empirical process approach to the union bound: Practical algorithms for combinatorial and linear bandits. *arXiv preprint arXiv:2006.11685*.
- Kaufmann, E., Cappé, O., and Garivier, A. (2016). On the complexity of best-arm identification in multi-armed bandit models. *The Journal of Machine Learning Research*, 17(1):1–42.
- Lattimore, T. and Szepesvári, C. (2020). *Bandit algorithms*. Cambridge University Press.
- Locatelli, A., Gutzeit, M., and Carpentier, A. (2016). An optimal algorithm for the thresholding bandit problem. In *International Conference on Machine Learning*, pages 1690–1698. PMLR.
- Lugosi, G. and Mendelson, S. (2019). Mean estimation and regression under heavy-tailed distributions: A survey. *Foundations of Computational Mathematics*, 19(5):1145–1190.
- Mason, B., Jain, L., Tripathy, A., and Nowak, R. (2020). Finding all  $\epsilon$ -good arms in stochastic bandits. *Advances in Neural Information Processing Systems*, 33.
- Scarlett, J., Bogunovic, I., and Cevher, V. (2017). Lower bounds on regret for noisy gaussian process bandit optimization. In *Conference on Learning Theory*, pages 1723–1742. PMLR.
- Shekhar, S. and Javidi, T. (2019). Multiscale gaussian process level set estimation. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3283–3291. PMLR.

- Soare, M., Lazaric, A., and Munos, R. (2014). Best-arm identification in linear bandits. *arXiv preprint arXiv:1409.6110*.
- Srinivas, N., Krause, A., Kakade, S. M., and Seeger, M. (2009). Gaussian process optimization in the bandit setting: No regret and experimental design. *arXiv preprint arXiv:0912.3995*.
- Todd, M. J. (2016). *Minimum-volume ellipsoids: Theory and algorithms*. SIAM.
- Vakili, S., Bouziani, N., Jalali, S., Bernacchia, A., and Shiu, D.-s. (2021). Optimal order simple regret for gaussian process bandits. *Advances in Neural Information Processing Systems*, 34.
- Valko, M., Korda, N., Munos, R., Flaounas, I., and Cristianini, N. (2013). Finite-time analysis of kernelised contextual bandits. *arXiv preprint arXiv:1309.6869*.
- Zanette, A., Zhang, J., and Kochenderfer, M. J. (2018). Robust super-level set estimation using gaussian processes. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 276–291. Springer.

## Contents

<b>1</b>	<b>INTRODUCTION</b>	<b>1</b>
1.1	Problem Statement . . . . .	2
<b>2</b>	<b>RELATED WORKS</b>	<b>2</b>
<b>3</b>	<b>EXPLICIT LEVEL SET ESTIMATION</b>	<b>3</b>
3.1	Algorithm . . . . .	4
3.2	Optimal Sample Complexity for Explicit Level Set Estimation . . . . .	5
<b>4</b>	<b>IMPLICIT LEVEL SET ESTIMATION</b>	<b>6</b>
4.1	Algorithm . . . . .	7
4.2	Theoretical Guarantees . . . . .	7
<b>5</b>	<b>EXPERIMENTS</b>	<b>8</b>
<b>6</b>	<b>CONCLUSION</b>	<b>9</b>
	<b>Appendices</b>	<b>12</b>
<b>A</b>	<b>Impacts and Limitations</b>	<b>13</b>
<b>B</b>	<b>Summary of Gaussian Processes Approaches for Level Set Estimation</b>	<b>13</b>
<b>C</b>	<b>Robust estimators for function means</b>	<b>14</b>
<b>D</b>	<b>Proofs for Explicit Level Set Estimation</b>	<b>14</b>
D.1	Lower Bound . . . . .	14
D.2	Upper Bound . . . . .	15
<b>E</b>	<b>Proofs for Implicit Level Set Estimation</b>	<b>20</b>
E.1	Lower Bounds . . . . .	20
E.2	Comparison to the lower bound of (Mason et al., 2020) . . . . .	22
E.3	Upper Bound . . . . .	23
<b>F</b>	<b>Additional Experiment Details</b>	<b>29</b>
F.1	Algorithms Implemented . . . . .	29
F.2	Additional Experiments . . . . .	31
<b>G</b>	<b>Reducing Experimental Design in an RKHS to a finite dimensional optimization</b>	<b>33</b>



## A Impacts and Limitations

Active learning uses a design objective to drive a sampling policy. In the simplest cases of active learning, such as regret minimization in standard multiarmed bandits, the relatively simple and unstructured setting leads to simple and easy to interpret sampling rules. For instance, the famed UCB algorithm simply forms confidence widths and pulls the arm with the largest upper bound. The transparency of this sampling rule makes UCB and algorithms like it inherently easy to diagnose and monitor in real time. For past algorithms in level set estimation, the acquisition functions merit easy oversight. By contrast, our work introduces optimal design to the area of level set estimation. As we show in our work, this can lead to improved sample complexity both theoretically and empirically. However, as the sampling distributions are based on a more complicated objective, how the algorithm chooses which data to collect is less immediately obvious or intuitive. This may make detecting issues such as biased sampling harder to detect and guard against, and for any large scale use of these algorithms in the wild, special care should be given to understand which points are being sampled the most and why. Furthermore, a common issue for many active learning approaches, this work included, is the possibility of model mismatch for any assumptions made in the theoretical analysis. While this work removes the need for an assumed prior over the true function  $f$ , other assumptions are still needed for the analysis, such as the function  $f$  not varying in time. If these assumptions are violated, the claims herein need not be true.

Any assumption made in this paper may reasonably be considered a limitation on the work depending on the application domain, though we hope that analytical assumptions may be easily modified to alter the algorithms to the practitioner’s needs. This is true, for instance in the case of all confidence widths we use. Another limitation of this work is computational complexity. The RIPS procedure necessary to compute estimates of individual function values relies on a robust estimator for each  $x \in \mathcal{X}$ . In this work, we leverage the Catoni estimator. While this is efficient for individual  $x$ ’s, as we observed in our experiments, if the set  $\mathcal{X}$  is large, this can become cumbersome. Additionally, how to best optimize the experimental design objectives is an active area of research and must be done carefully. Finally, our algorithms both suffer potentially bad logarithmic terms in the per-round sample complexity, and this can affect the real-world performance of MELK and MILK. The technique of (Katz-Samuels et al., 2020) may be able to avoid this.

## B Summary of Gaussian Processes Approaches for Level Set Estimation

In Table 2, we briefly summarize past algorithmic approaches to level set estimation. In general, past methods center around the design of an *acquisition function* which at each time  $t$  tells the algorithm which point to go sample. By contrast, the algorithms in this paper both use experimental design to select batches of samples to go gather at one time.

Algorithm	Acquisition Function	Theoretical guarantee
Straddle	$\arg \max_i u_i(t) - \tau \wedge \tau - \ell_i(t)$	None, $u_i(t)$ and $\ell_i(t)$ are set as $1.96 \cdot \sigma_{t-1}$ .
LSE	$\arg \max_i u_i(t) - \tau \wedge \tau - \ell_i(t)$	$\eta$ -approximate solution in $T \lesssim \frac{\gamma_t \log(n/\delta)}{\eta^2}$
TruVar	$\arg \min_{x_i} \sum_{x_j} \sigma_{t-1 x_i}^2(x_j)$	$\eta$ -approximate solution in $T \lesssim \frac{\Gamma_t \log(n/\delta)}{\eta^2}$
RMILE	$\arg \max_{x_i} \{\mathbb{E} \sum_{x_j} (\mathbb{P}_{GP x_i}(f(x_j) > \tau) - \mathbb{P}_{GP}(f(x_j) > \tau)), \sigma^2(x_i)\}$	can be shown to be similar to A-optimality, no complexity guarantee
MELK	G-optimal design	Matching upper and lower bounds in the linear case.

Table 1: Algorithms and theoretical guarantees for explicit LSE

Algorithm	Acquisition function	Theoretical guarantee
LSE-imp	$\arg \max_i \sigma^2(x_i)$	$\eta$ -approximate solution in $T \lesssim \frac{\Gamma_t \log(n/\delta)}{\eta^2}$
MILK	XY optimal design over vectors $\phi(x) - (1 - \epsilon)\phi(x')$	Upper bounds and matching lower for certain cases.

Table 2: Acquisition functions and theoretical guarantees for implicit level set estimation

## C Robust estimators for function means

In order for the algorithm to declare whether points  $\mathbf{x}$  belong in  $G_\alpha$  (or  $G_\epsilon$  in the sequel) or not, we require an estimator of the function values  $f(\mathbf{x})$ . As we have introduced structure by assuming that  $f$  is well approximated by a function  $\theta_*$  in the RKHS  $\mathcal{H}$ , we seek an estimator that leverages this structure to provide accurate estimates of many arms given samples of only a few. As a warmup, in the linear case where  $\phi(\cdot)$  is the identity map, one could form the least squares or regularized least squares estimate of  $\theta_*$  denoted  $\hat{\theta}$  and estimate the mean of any point  $\mathbf{x}$  as  $\hat{\theta}^T \mathbf{x}$ . To sample to estimate  $\theta_*$ , optimal design procedures first compute a design  $\lambda \in \Delta_{\mathcal{X}}$ . Then for a specified number of samples  $N$ , it is common to use an efficient rounding procedure such as (Allen-Zhu et al., 2017) to compute an allocation of the  $N$  samples to the arms  $\mathcal{X}$  such that  $\mathbf{x}_i$  gets roughly  $\lambda_i \cdot N$  samples (Fiez et al., 2019; Jun et al., 2020). Efficient rounding procedures require that  $N = \Omega(d)$ , and while this is a minor assumption in the case of a linear RKHS where  $\phi(\mathbf{x}) = \mathbf{x}$ , in general  $\phi(\mathbf{x})$  may be infinite dimensional, and naive rounding is not possible. Instead of performing rounding given design  $\lambda$ , one may instead sample from  $\lambda$  directly and use inverse propensity scoring (IPS) which avoids bad dimensional factors but can have high variance.

In this work, we leverage the RIPS estimator from (Camilleri et al., 2021) which combines IPS with robust mean estimation and regularization to control variance and is presented in Algorithm 3. RIPS requires a robust mean estimator for its performance and theoretical guarantees. In Theorem 3.2, we state the guarantee of this estimator.

---

### Algorithm 3 RIPS: Robust IPS estimator

---

**Require:** Finite sets  $\mathcal{X} \subset \mathbb{R}^d$  and  $\mathcal{V} \subset \mathcal{H}$ , feature map  $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ , number of samples  $\tau$ , regularization  $\gamma > 0$ , robust mean estimator  $\hat{\mu} : \mathbb{R}^* \rightarrow \mathbb{R}$

$$\lambda := \arg \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{v \in \mathcal{V}} \|v\|_{(A^{(\gamma)}(\lambda))^{-1}}$$

1: Randomly draw  $\tilde{x}_1, \dots, \tilde{x}_\tau$  from  $\mathcal{X}$  according to  $\lambda^*$

2: Set  $W^{(v)} = \hat{\mu}(\{v^\top A^{(\gamma)}(\lambda^*)^{-1} \phi(\tilde{x}_t) \tilde{y}_t\}_{t=1}^\tau)$

**return**  $\hat{\theta} := \arg \min_{\theta} \max_{v \in \mathcal{V}} \frac{|\langle \theta, v \rangle - W^{(v)}|}{\|v\|_{(A^{(\gamma)}(\lambda))^{-1}}}$

---

We next state the complete theoretical guarantee of the RIPS estimator.

**Theorem C.1** (Theorem 1, (Camilleri et al., 2021)). *Consider the model  $y = \langle \phi(\mathbf{x}), \theta^* \rangle_{\mathcal{H}} + \zeta_{\mathbf{x}} + \eta$  for misspecification  $|\zeta_{\mathbf{x}}| \leq h$  where it is assumed that  $|y| \leq B$ ,  $\mathbb{E}[\eta] = 0$ , and  $\mathbb{E}[\eta^2] \leq \sigma^2$ . Fix any finite sets  $\mathcal{X} \subset \mathbb{R}^d$  and  $\mathcal{V} \subset \mathcal{H}$ , feature map  $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ , number of samples  $\tau$  and regularization  $\gamma > 0$ . If the RIPS procedure of Algorithm 3 is run with  $\frac{\delta}{|\mathcal{V}|}$ -robust mean estimator  $\hat{\mu}(\cdot)$  and if  $\tau \geq c_1 \log(|\mathcal{V}|/\delta)$  then with probability at least  $1 - \delta$ , we have*

$$\begin{aligned} \max_{v \in \mathcal{V}} \frac{|W^{(v)} - \langle \theta_*, v \rangle|}{\|v\|_{(A^{(\gamma)}(\lambda))^{-1}}} &\leq \sqrt{\gamma} \|\theta_*\| + h \\ &\quad + c \sqrt{\frac{(B^2 + \sigma^2)}{\tau} \log(2|\mathcal{V}|/\delta)} \end{aligned}$$

Moreover,  $W^{(v)} = \hat{\mu}(\{v^\top A^{(\gamma)}(\lambda)^{-1} \phi(\mathbf{x}_t) y_t\}_{t=1}^\tau)$  can be replaced by  $\langle \hat{\theta}, v \rangle$  by multiplying the RHS by a factor of 2.

For RIPS, we leverage Catoni's estimator (Lugosi and Mendelson, 2019) for which  $c_1 = 2$  and  $c = 4$  suffice.

## D Proofs for Explicit Level Set Estimation

### D.1 Lower Bound

*Proof of Theorem 3.1.* Recall that we have assumed that  $h = 0$  and  $\phi(\mathbf{x}) = \mathbf{x}$ . We begin with a result of (Fiez et al., 2019) that will be useful here.

**Lemma D.1** ((Fiez et al., 2019), Remark 2). *The projection onto the closure of the set  $\{\theta \in \mathbb{R}^d : \theta^T \mathbf{x} < \alpha\}$  under the  $\|\cdot\|_{A(\lambda)}$  norm is given by*

$$\theta_{\mathbf{x}} := \theta - \frac{(\theta^T \mathbf{x} - \alpha)A(\lambda)^{-1}\mathbf{x}}{\|\mathbf{x}\|_{A(\lambda)^{-1}}^2}.$$

By (Kaufmann et al., 2016), we have that the any  $\delta$ -PAC algorithm for all- $\alpha$  requires

$$\min_{\lambda} \frac{KL(1 - \delta, \delta)}{\min_{\theta' \in \text{Alt}(\theta_*)} \|\theta' - \theta_*\|_{A(\lambda)}}$$

where  $\text{Alt}(\theta_*)$  is the set of alternates such that  $G_{\alpha}(\theta_*) \neq G_{\alpha}(\theta')$  for any  $\theta' \in \text{Alt}(\theta_*)$ . The set of alternates may be decomposed as

$$\mathbf{Alt}(\theta_*) = \left( \bigcup_{\mathbf{x} \in G_{\alpha}(\theta_*)} \{\theta' : \mathbf{x} \notin G_{\alpha}(\theta')\} \right) \cup \left( \bigcup_{\mathbf{x} \in G_{\alpha}(\theta_*)^c} \{\theta' : \mathbf{x} \in G_{\alpha}(\theta')\} \right)$$

Note that  $\mathbf{x} \in G_{\alpha}(\theta_*) \iff \theta_*^T \mathbf{x} > \alpha$ . Hence, the set of alternates for any  $\mathbf{x} \in G_{\alpha}(\theta_*)$  such that  $\mathbf{x} \in G_{\alpha}^c(\theta')$  for any  $\theta' \in \mathbf{Alt}(\theta_*)$  is given by

$$A_{\mathbf{x}} := \{\theta \in \mathbb{R}^d : \theta^T \mathbf{x} < \alpha\}.$$

Next note that  $\mathbf{x} \in G_{\alpha}^c(\theta_*) \iff \theta_*^T \mathbf{x} < \alpha$ . Hence, for any  $\mathbf{x} \in G_{\alpha}^c(\theta_*)$  the set of alternates such that  $\mathbf{x} \in G_{\alpha}(\theta')$  for any  $\theta' \in \mathbf{Alt}(\theta_*)$  is given by

$$A_{\mathbf{x}} := \{\theta \in \mathbb{R}^d : \theta^T \mathbf{x} > \alpha\}.$$

Next, we discuss how to project onto  $A_{\mathbf{x}}$ . As this set is open, to be precise, we should take a point in the interior and consider the limit for a sequence approaching the boundary. For brevity, we simply project onto the closure and consider the closures of the  $A_{\mathbf{x}}$  sets. Using the decomposition of  $\mathbf{Alt}(\theta_*)$  we have that

$$\min_{\theta' \in \mathbf{Alt}(\theta_*)} \|\theta' - \theta_*\|_{A(\lambda)} = \min_{\mathbf{x}} \min_{\theta' \in A_{\mathbf{x}}} \|\theta' - \theta_*\|_{A(\lambda)} = \min_{\mathcal{S} \in \{G_{\alpha}, G_{\alpha}^c\}} \min_{\mathbf{x} \in \mathcal{S}} \min_{\theta \in A_{\mathbf{x}}} \|\theta' - \theta_*\|_{A(\lambda)}.$$

For  $\mathbf{x} \in G_{\alpha}(\theta_*)$ , using Lemma D.1 and recalling the definition of the set  $\theta_{\mathbf{x}}$  therein,

$$\min_{\theta' \in A_{\mathbf{x}}} \|\theta' - \theta_*\|_{A(\lambda)} = \min_{\theta' \in \{\theta \in \mathbb{R}^d : \theta^T \mathbf{x} \leq \alpha\}} \|\theta' - \theta_*\|_{A(\lambda)} = \|\theta_{\mathbf{x}} - \theta_*\|_{A(\lambda)}.$$

The statement for points in  $G_{\alpha}^c$  follows identically. Hence,

$$\min_{\theta' \in \mathbf{Alt}(\theta_*)} \|\theta' - \theta_*\|_{A(\lambda)} = \min_{\mathbf{x}} \|\theta_{\mathbf{x}} - \theta_*\|_{A(\lambda)}$$

Note that

$$\|\theta_{\mathbf{x}} - \theta_*\|_{A(\lambda)} = \frac{(\theta_*^T (\mathbf{x}' - \mathbf{x}) - \alpha)^2}{2\|\mathbf{x}\|_{A(\lambda)^{-1}}^2}$$

by Theorem 2 of (Fiez et al., 2019). Hence, any  $\delta$ -PAC algorithm requires at least

$$2 \min_{\lambda} \max_{\mathbf{x}} \frac{\|\mathbf{x}\|_{A(\lambda)^{-1}}^2}{(\theta_*^T \mathbf{x} - \alpha)^2} KL(1 - \delta, \delta)$$

samples in expectation. Noting that the binary entropy  $KL(1 - \delta, \delta) \geq \log(1/2.4\delta)$  completes the proof.  $\square$

## D.2 Upper Bound

Next, we restate Theorem 3.3 that bounds the complexity of MELK.

**Theorem D.2.** *Fix  $\delta > 0$ , threshold  $\alpha > 0$ , tolerance  $\tilde{\beta}$ , and regularization  $\gamma \geq 0$ . Define  $\Delta_{\min}(\alpha) := \min |\phi(\mathbf{x})^T \theta_* - \alpha|$ . Define also*

$$\bar{\beta}(\alpha) = \min\{\beta > 0 : 4(\sqrt{\gamma}\|\theta_*\| + h)(2 + \sqrt{f(\mathcal{X}, \{\phi(\mathbf{x}) | \mathbf{x} \in \mathcal{X}, |\phi(\mathbf{x})^T \theta_* - \alpha| \leq \beta\}; \gamma)}) \leq \beta\}.$$

With probability at least  $1 - \delta$ , *MELK* returns a set  $\widehat{\mathcal{R}} = (\mathcal{X} \setminus \widehat{B}_t)$  at time  $T_\delta$  such that

$$\{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \geq \alpha + \bar{\beta}(\alpha)\} \subseteq \widehat{\mathcal{R}} \subseteq \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \geq \alpha - \tilde{\beta} - \bar{\beta}(\alpha)\}$$

and for any  $\alpha, \tilde{\beta}$  such that  $\max(\Delta_{\min}(\alpha), \tilde{\beta}) \geq \bar{\beta}(\alpha)$

$$T_\delta \leq 256(B^2 + \sigma^2) \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X}} \frac{\|\phi(\mathbf{x})\|_{(A(\lambda) + \gamma I)^{-1}}^2}{\max\{(\phi(\mathbf{x})^T \theta_* - \alpha)^2, \tilde{\beta}^2\}} \log \left( \frac{4|\mathcal{X}|^2 \lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil^2}{\delta} \right) + 2 \log(|\mathcal{X}|/\delta) \lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil$$

Recall the definition of the set  $G_\alpha := \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) > \alpha\}$ .

**Lemma D.3.** For any  $\mathcal{V} \subset \mathcal{X}$  define  $f(\mathcal{X}, \mathcal{V}; \gamma) = \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{v \in \mathcal{V}} \|v\|_{(\sum_{x \in \mathcal{X}} \lambda_x \phi(x)\phi(x)^T + \gamma I)^{-1}}^2$ .

In each round  $t$ , define the event

$$\mathcal{E}_t = \{|\mathbf{x}^T(\widehat{\theta}_t - \theta_*)| \leq 2^{-t} + (\sqrt{\gamma}\|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{A}_t; \gamma)} \forall \mathbf{x} \in \mathcal{A}_t\}$$

Holds  $\mathbb{P}(\bigcup_{t=1}^{\infty} \mathcal{E}_t^c) \leq \delta$ .

*Proof.* Using Theorem 3.2, for any  $\mathbf{x} \in \mathcal{A}_t$  we have that with probability at least  $1 - \delta_t/|\mathcal{X}|^2$

$$\begin{aligned} |\mathbf{x}^T(\widehat{\theta}_t - \theta_*)| &\leq \|\mathbf{x}\|_{(\sum_{x \in \mathcal{X}} \lambda_x \mathbf{x} \mathbf{x}^T + \gamma I)^{-1}} \left( \sqrt{\gamma}\|\theta_*\| + h + c \sqrt{\frac{(B^2 + \sigma^2)}{N_t} \log(2t^2|\mathcal{X}|^2/\delta)} \right) \\ &\leq \sqrt{f(\mathcal{X}, \mathcal{A}_t; \gamma)} \left( \sqrt{\gamma}\|\theta_*\| + h + 2^{-t}/\sqrt{f(\mathcal{X}, \mathcal{A}_t; \gamma)} \right) \\ &\leq 2^{-t} + (\sqrt{\gamma}\|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{A}_t; \gamma)} \end{aligned}$$

Since  $|\mathcal{A}_t| \leq |\mathcal{X}|^2$ ,  $\mathcal{E}_t$  holds for all  $\mathbf{x} \in \mathcal{A}_t$  with probability  $1 - \delta_t$  via a union bound. Taking a second union bound over rounds, we have that

$$\mathbb{P} \left( \bigcup_{t=1}^{\infty} \mathcal{E}_t^c \right) \leq \sum_{t=1}^{\infty} \mathbb{P}(\mathcal{E}_t^c) \leq \sum_{t=1}^{\infty} \delta_t = \sum_{t=1}^{\infty} \frac{\delta}{2t^2} \leq \delta$$

□

Define

$$\bar{t} = \max\{t : (\sqrt{\gamma}\|\theta_*\|_2 + h)(2 + \sqrt{f(\mathcal{X}, \{\mathbf{x} \in \mathcal{X} : |\mathbf{x}^T \theta_* - \alpha| \leq 2^{-t+2}\}; \gamma)}) \leq 2^{-t}\}.$$

As we will see in Lemmas D.6 and D.7,

$$\mathcal{A}_t \subset \{\mathbf{x} \in \mathcal{X} : |\mathbf{x}^T \theta_* - \alpha| \leq 2^{-t+1}\}.$$

Thus for  $t \leq \bar{t}$ , holds on  $\bigcap_t \mathcal{E}_t$  that

$$\forall \mathbf{x} \in \mathcal{A}_t, |\mathbf{x}^T(\widehat{\theta}_t - \theta_*)| \leq 2 \cdot 2^{-t}.$$

**Lemma D.4.** On  $\bigcap_t \mathcal{E}_t$ , when  $t \leq \bar{t}$  holds  $\widehat{G}_t \subset G_\alpha^\phi := \{\mathbf{x} : \phi(\mathbf{x})^T \theta_* > \alpha\}$ .

**Remark:** If  $h = 0$ ,  $G_\alpha^\phi = G_\alpha$ .

*Proof.*

$$\begin{aligned} \mathbf{x} \in \widehat{G}_t &\iff \exists t' \leq t : \phi(\mathbf{x})^T \widehat{\theta}_{t'} \geq \alpha + 2 \cdot 2^{-t'} \\ &\iff \exists t' \leq t : \phi(\mathbf{x})^T(\widehat{\theta}_{t'} - \theta_*) + \phi(\mathbf{x})^T \theta_* \geq \alpha + 2 \cdot 2^{-t'} \\ &\xrightarrow{\bigcap_t \mathcal{E}_t} \phi(\mathbf{x})^T \theta_* > \alpha \\ &\iff \mathbf{x} \in G_\alpha^\phi. \end{aligned}$$

□



**Lemma D.5.** On  $\bigcap_t \mathcal{E}_t$ , when  $t \leq \bar{t}$  holds,  $\widehat{B}_t \subset (G_\alpha^\phi)^c$ .

*Proof.*

$$\begin{aligned} \mathbf{x} \in \widehat{B}_t &\iff \exists t' \leq t : \phi(\mathbf{x})^T \widehat{\theta}_t \leq \alpha - 2 \cdot 2^{-t'} \\ &\iff \exists t' \leq t : \phi(\mathbf{x})^T (\widehat{\theta}_t - \theta_*) + \phi(\mathbf{x})^T \theta_* \leq \alpha - 2 \cdot 2^{-t'} \\ &\stackrel{\bigcap_t \mathcal{E}_t}{\implies} \phi(\mathbf{x})^T \theta_* < \alpha \\ &\iff \mathbf{x} \in (G_\alpha^\phi)^c. \end{aligned}$$

□

**Lemma D.6.** On the event  $\bigcap_t \mathcal{E}_t$ , when  $t \leq \bar{t}$  holds,

$$\mathcal{A}_t \cap G_\alpha^\phi \subset \left\{ \mathbf{x} \in G_\alpha^\phi \mid |\phi(\mathbf{x})^T \theta_* - \alpha| \leq 2^{-t+2} \right\} =: \mathcal{S}_t^{\text{Above}}$$

*Proof.* For any  $\mathbf{x} \in G_\alpha^\phi$  such that  $\phi(\mathbf{x})^T \theta_* > \alpha + 2^{-t+1}$ , if  $t \geq \log(4(\alpha - \phi(\mathbf{x})^T \theta_*)^{-1})$  and  $t \leq \bar{t}$ , then

$$\phi(\mathbf{x})^T \widehat{\theta}_t = \phi(\mathbf{x})^T (\widehat{\theta}_t - \theta_*) + \phi(\mathbf{x})^T \theta_* > -2^{-t+1} + \alpha + 2^{-t+1} = \alpha \geq \alpha$$

which implies that  $\mathbf{x} \in \widehat{G}_t$ . Noting that  $\mathcal{A}_t \cap \widehat{G}_{t-1} = \emptyset$  completes the proof. □

**Lemma D.7.** On the event  $\bigcap_t \mathcal{E}_t$ , when  $t \leq \bar{t}$  holds,

$$\mathcal{A}_t \cap (G_\alpha^\phi)^c \subset \left\{ \mathbf{x} \in (G_\alpha^\phi)^c \mid |\phi(\mathbf{x})^T \theta_* - \alpha| \leq 2^{-t+2} \right\} =: \mathcal{S}_t^{\text{Below}}$$

*Proof.* The proof follows identically as that of Lemma D.6 □

**Remark:** Lemmas D.6 and D.7 jointly imply that  $\mathcal{A}_t \subset \{\mathbf{x} \mid |\phi(\mathbf{x})^T \theta_* - \alpha| \leq 2^{-t+2}\} =: \mathcal{S}_t$  for  $t \leq \bar{t}$ . Furthermore,  $f(\mathcal{X}, \mathcal{A}_t, \gamma) \leq f(\mathcal{X}, \mathcal{S}_t, \gamma)$ .

**Remark:**

The algorithm stops on either of two conditions. On one hand if  $t \geq \lceil \log_2(4/\tilde{\beta}) \rceil =: t_\beta$ , then it has achieved precision  $\tilde{\beta}$  as desired and it terminates. Otherwise, it terminates if  $\widehat{G}_t \cup \widehat{B}_t = \mathcal{X}$ . This occurs when  $\tilde{\beta}$  is very small. Define  $\Delta_{\min}(\alpha) := \min |\phi(\mathbf{x})^T \theta_* - \alpha|$ . Recall

$$\begin{aligned} \bar{t} &= \max\{t : (\sqrt{\gamma}\|\theta_*\|_2 + h)(2 + \sqrt{f(\mathcal{X}, \{\mathbf{x} \in \mathcal{X} : |\phi(\mathbf{x})^T \theta_* - \alpha| \leq 4 \cdot 2^{-t}\}; \gamma)}) \leq 2^{-t}\} \\ &= \max\{t : 4(\sqrt{\gamma}\|\theta_*\|_2 + h)(2 + \sqrt{f(\mathcal{X}, \{\mathbf{x} \in \mathcal{X} : |\phi(\mathbf{x})^T \theta_* - \alpha| \leq 4 \cdot 2^{-t}\}; \gamma)}) \leq 4 \cdot 2^{-t}\} \\ &= -2 + \max\{t : 4(\sqrt{\gamma}\|\theta_*\|_2 + h)(2 + \sqrt{f(\mathcal{X}, \{\mathbf{x} \in \mathcal{X} : |\phi(\mathbf{x})^T \theta_* - \alpha| \leq 2^{-t}\}; \gamma)}) \leq 2^{-t}\} \\ &= -3 - \log_2(\min\{\beta > 0 : 4(\sqrt{\gamma}\|\theta_*\|_2 + h)(2 + \sqrt{f(\mathcal{X}, \{\mathbf{x} \in \mathcal{X} : |\phi(\mathbf{x})^T \theta_* - \alpha| \leq \beta\}; \gamma)}) \leq \beta\}). \end{aligned}$$

This defines

$$\bar{\beta} = \min\{\beta > 0 : 4(\sqrt{\gamma}\|\theta_*\|_2 + h)(2 + \sqrt{f(\mathcal{X}, \{\mathbf{x} \in \mathcal{X} : |\phi(\mathbf{x})^T \theta_* - \alpha| \leq \beta\}; \gamma)}) \leq \beta\}.$$

Let  $t_{\max}$  denote the random variable of the last round before the algorithm terminates. The following Lemmas give a guarantee on the set  $\mathcal{X} \setminus \widehat{B}_t$  at termination.

**Lemma D.8.** On the event  $\bigcap_{t=1}^\infty \mathcal{E}_t$ , MELK returns a set  $(\mathcal{X} \setminus \widehat{B}_{t_{\max}})$  such that  $\{\mathbf{x} : f(\mathbf{x}) > \alpha + \bar{\beta}(\alpha)\} \subset (\mathcal{X} \setminus \widehat{B}_{t_{\max}})$ .

*Proof.* Take any  $\mathbf{x}$  such that  $f(\mathbf{x}) > \alpha + \bar{\beta}(\alpha)$  and recall that by assumption  $|f(\mathbf{x}) - \phi(\mathbf{x})^T \theta_*| \leq h$  for all  $\mathbf{x} \in \mathcal{X}$ . We consider two cases. In the first case, assume that  $t_{\max} \leq \bar{t}$ . We claim that in this case  $\exists t$  such that  $\mathbf{x} \in \widehat{B}_t$ . We prove this by contradiction. Assume not. Then  $\exists t$  such that

$$\begin{aligned} \widehat{\theta}_t^T \phi(\mathbf{x}) < \alpha - 2^{-t+1} &\iff \phi(\mathbf{x})^T (\widehat{\theta}_t - \theta_*) + \phi(\mathbf{x})^T \theta_* < \alpha - 2^{-t+1} \\ &\stackrel{\mathcal{E}_t}{\implies} -2^{-t} - (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{A}_t; \gamma)} + \phi(\mathbf{x})^T \theta_* < \alpha - 2^{-t+1} \\ &\implies -2^{-t} - (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{S}_t; \gamma)} + \phi(\mathbf{x})^T \theta_* < \alpha - 2^{-t+1} \\ &\implies -(\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{S}_t; \gamma)} + f(\mathbf{x}) - h < \alpha - 2^{-t} \\ &\implies f(\mathbf{x}) < \alpha - 2^{-t} + h + (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{S}_t; \gamma)}. \end{aligned}$$

Recalling that we have assumed that  $f(\mathbf{x}) > \alpha + \bar{\beta}(\alpha)$ . Hence, this implies that

$$\bar{\beta}(\alpha) < -2^{-t} + h + (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{S}_t; \gamma)}.$$

Note that  $\bar{\beta}(\alpha) > 0$ . As we have assumed that,  $t \leq t_{\max} \leq \bar{t}$ , we have that  $2^{-t} \geq (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{S}_t; \gamma)}$  using the definition of  $\bar{t}$ . Hence, we have that

$$h > \bar{\beta}(\alpha) > 4h$$

which is a contradiction where the final inequality follows from the definition of  $\bar{\beta}(\alpha)$  for  $\gamma > 0$ . Hence, in this case we have shown that  $\{\mathbf{x} : f(\mathbf{x}) > \alpha + \bar{\beta}(\alpha)\} \subset (\mathcal{X} \setminus \widehat{B}_{t_{\max}})$ .

In the second case, assume that  $t_{\max} > \bar{t}$  and take  $\mathbf{x}$  such that  $f(\mathbf{x}) > \alpha + \bar{\beta}(\alpha)$ . We claim that  $\mathbf{x} \in \widehat{G}_{\bar{t}}$  and hence  $\mathbf{x} \notin \mathcal{A}_t$  for any  $t > \bar{t}$  and thus is never added to  $\widehat{B}_t$ . This occurs if

$$\begin{aligned} \phi(\mathbf{x})^T \widehat{\theta}_{\bar{t}} > \alpha + 2^{-\bar{t}+1} &\iff \phi(\mathbf{x})^T (\widehat{\theta}_{\bar{t}} - \theta_*) + \phi(\mathbf{x})^T \theta_* > \alpha + 2^{-\bar{t}+1} \\ &\stackrel{\mathcal{E}_{\bar{t}}}{\iff} -2^{-\bar{t}} - (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{A}_{\bar{t}}; \gamma)} + \phi(\mathbf{x})^T \theta_* \geq \alpha + 2^{-\bar{t}+1} \\ &\iff -2^{-\bar{t}} - (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{S}_{\bar{t}}; \gamma)} + \phi(\mathbf{x})^T \theta_* \geq \alpha + 2^{-\bar{t}+1} \\ &\iff \phi(\mathbf{x})^T \theta_* \geq \alpha + 2^{-\bar{t}+1} + 2^{-\bar{t}} + (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{S}_{\bar{t}}; \gamma)} \end{aligned}$$

Recall that  $f(\mathbf{x}) > \alpha + \bar{\beta}(\alpha)$ . Furthermore, we have by the definition of  $\bar{t}$  that

$$2^{-\bar{t}} \geq (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{S}_{\bar{t}}; \gamma)}.$$

Hence, the above is implied by  $\bar{\beta}(\alpha) - h \geq 4 \cdot 2^{-\bar{t}} = 0.5 \bar{\beta}(\alpha)$  where the final equality holds by definition of  $\bar{\beta}(\alpha)$ . Noting that  $\bar{\beta}(\alpha) > 4h$  proves this claim. In summary, we have shown that for any  $\mathbf{x}$  such that  $f(\mathbf{x}) > \alpha + \bar{\beta}(\alpha)$ , if  $t_{\max} \leq \bar{t}$ , then  $\mathbf{x}$  is never added to  $\widehat{B}_t$  and hence is contained in the set  $\mathcal{X} \setminus \widehat{B}_t$  at termination, and if otherwise that  $t_{\max} > \bar{t}$ , then  $\mathbf{x}$  is added to the set  $\widehat{G}_{\bar{t}}$  before round  $\bar{t} + 1$  and hence is removed from the active set and never added to  $\widehat{B}_t$ . Applying this argument to any  $\mathbf{x}$  such that  $f(\mathbf{x}) > \alpha + \bar{\beta}(\alpha)$  completes the proof.  $\square$

**Lemma D.9.** *On the event  $\bigcap_{t=1}^{\infty} \mathcal{E}_t$ , MELK returns a set  $(\mathcal{X} \setminus \widehat{B}_{t_{\max}})$  such that  $(\mathcal{X} \setminus \widehat{B}_{t_{\max}}) \subset \{\mathbf{x} : f(\mathbf{x}) > \alpha - \bar{\beta}(\alpha) - \widetilde{\beta}\}$ .*

*Proof.* Take any  $\mathbf{x}$  such that  $f(\mathbf{x}) < \alpha - \bar{\beta}(\alpha) - \widetilde{\beta}$ . We claim that there exists a  $t \leq t_{\max}$  such that  $\mathbf{x}$  is added to  $\widehat{B}_t$  which implies that  $\mathbf{x} \notin (\mathcal{X} \setminus \widehat{B}_{t_{\max}})$ . Suppose for contradiction that this is not the case. Then for all  $t \leq t_{\max}$ ,

$$\begin{aligned} \widehat{\theta}_t^T \phi(\mathbf{x}) > \alpha - 2^{-t+1} &\iff \phi(\mathbf{x})^T (\widehat{\theta}_t - \theta_*) + \phi(\mathbf{x})^T \theta_* > \alpha - 2^{-t+1} \\ &\stackrel{\mathcal{E}_t}{\implies} 2^{-t} + (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{A}_t; \gamma)} + \phi(\mathbf{x})^T \theta_* > \alpha - 2^{-t+1} \\ &\implies 2^{-t} + (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{S}_t; \gamma)} + \phi(\mathbf{x})^T \theta_* > \alpha - 2^{-t+1} \\ &\implies (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{S}_t; \gamma)} + f(\mathbf{x}) + h > \alpha - 2^{-t+1} - 2^{-t} \\ &\implies f(\mathbf{x}) > \alpha - 2^{-t+1} - 2^{-t} - h - (\sqrt{\gamma} \|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{S}_t; \gamma)}. \end{aligned}$$

Plugging in  $f(\mathbf{x}) < \alpha - \bar{\beta}(\alpha) - \tilde{\beta}$ , the above implies

$$\bar{\beta}(\alpha) + \tilde{\beta} < 2^{-t+1} + 2^{-t} + h + (\sqrt{\gamma}\|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{S}_t; \gamma)} \quad (3)$$

Next, recall that MELK terminates either on the condition that  $t = \lceil \log_2(4/\tilde{\beta}) \rceil$  or that  $\hat{G}_t \cup \hat{B}_t = \mathcal{X}$ . Using this, we brake our analysis into cases.

Case 1:  $t_{\max} = \lceil \log_2(4/\tilde{\beta}) \rceil \leq \bar{t}$ .

In this case, MELK stops due to the  $\tilde{\beta}$  tolerance in a round before  $\bar{t}$ . For  $t \leq \bar{t}$ , we have that  $2^{-t} \geq + (\sqrt{\gamma}\|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{S}_t; \gamma)}$ . Hence, the above implies that

$$\bar{\beta}(\alpha) + \tilde{\beta} < 2^{-t+2} + h.$$

As we have assumed this condition for all  $t \leq t_{\max}$ , we may plug in  $t_{\max}$  which implies

$$\bar{\beta}(\alpha) + \tilde{\beta} < \tilde{\beta} + h.$$

As  $\bar{\beta}(\alpha) > h$ , this is a contradiction. Hence there must exist a  $t$  such that  $\mathbf{x} \in \hat{B}_t$ .

Case 2:  $t_{\max} \leq \bar{t} < \lceil \log_2(4/\tilde{\beta}) \rceil$ .

In this case, MELK terminates before round  $t = \lceil \log_2(4/\tilde{\beta}) \rceil$ . Hence, it does so on the condition that  $\hat{G}_t \cup \hat{B}_t = \mathcal{X}$ . Note that for  $f(\mathbf{x}) < \alpha - \bar{\beta}(\alpha) - \tilde{\beta}$ , we have that  $\mathbf{x} \in (G_\alpha^\phi)^c$  since  $\bar{\beta}(\alpha) > h$  and  $\tilde{\beta} \geq 0$ . If we terminate before round  $\bar{t}$ , we have by Lemma D.5 that  $(G_\alpha^\phi)^c \subset \hat{B}_t$  which implies that  $\mathbf{x} \in \hat{B}_{t_{\max}}$ . This contradicts the assumption that  $\nexists t : \mathbf{x} \in \hat{B}_t$ .

Case 3:  $\bar{t} < t_{\max}$ .

In this case, MELK terminates at a round after  $\bar{t}$ . In this setting, we argue that  $\mathbf{x} \in \hat{B}_{\bar{t}}$ . Recall that for any  $t \leq \bar{t}$ , (3) simplifies to

$$\bar{\beta}(\alpha) + \tilde{\beta} < 2^{-t+2} + h$$

Plugging in  $\bar{t}$ , and noting that  $2^{-\bar{t}+2} = \frac{1}{2}\bar{\beta}(\alpha)$ , the above implies

$$\bar{\beta}(\alpha) + \tilde{\beta} < \frac{1}{2}\bar{\beta}(\alpha) + h.$$

Noting that  $\bar{\beta}(\alpha) > 4h$ , shows that the above is a contradiction. Hence, there exists a  $t \leq \bar{t}$  such that  $\mathbf{x} \in \hat{B}_t$ .

Therefore, in all cases we have shown that for any  $\mathbf{x}$  such that  $f(\mathbf{x}) < \alpha - \bar{\beta}(\alpha) - \tilde{\beta}$ ,  $\mathbf{x} \in \hat{B}_t$ . Therefore, for the returned set  $\mathcal{X} \setminus \hat{B}_{t_{\max}}$ , we have that

$$(\mathcal{X} \setminus \hat{B}_{t_{\max}}) \subset \{\mathbf{x} : f(\mathbf{x}) > \alpha - \bar{\beta}(\alpha) - \tilde{\beta}\}.$$

□

*Proof of Theorem 3.3.* Throughout, assume the high probability event  $\bigcap_T \mathcal{E}_t$ . By Lemmas D.8 and D.9 in conjunction with the high probability event  $\bigcap \mathcal{E}_t$  we have correctness. It remains to control the sample complexity of MELK. Recall that we have assumed that  $\max(\Delta_{\min}(\alpha), \tilde{\beta}) \geq \bar{\beta}(\alpha)$ . This implies that  $\min\{\lceil \log_2(4/\Delta_{\min}(\alpha)) \rceil, \lceil \log_2(4/\tilde{\beta}) \rceil\} \leq \bar{t}$ . Applying Lemmas D.6 and D.7, we have that  $t_{\max} \leq \min\{\lceil \log_2(4/\Delta_{\min}(\alpha)) \rceil, \lceil \log_2(4/\tilde{\beta}) \rceil\} \leq \bar{t}$  and that  $\mathcal{A}_t \subseteq \mathcal{S}_t$  for all rounds  $t$ . Now we proceed by bounding the total number of samples drawn.

$$\tau \leq \sum_{t=1}^{t_{\max}} N_t$$

$$\begin{aligned}
 &\leq \sum_{t=1}^{\min\{\lceil \log_2(4/\Delta_{\min}(\alpha)) \rceil, \lceil \log_2(4/\tilde{\beta}) \rceil\}} N_t \\
 &= \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil} \max\{c_1 \log(|\mathcal{X}|/\delta), c^2 2^{2t} f(\mathcal{A}_t; \gamma) (B^2 + \sigma^2) \log(2t^2 |\mathcal{X}|^2/\delta)\} \\
 &\leq c_1 \log(|\mathcal{X}|/\delta) \lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil + c^2 (B^2 + \sigma^2) \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil} 2^{2t} f(\mathcal{A}_t; \gamma) \cdot \log(2t^2 |\mathcal{X}|^2/\delta) \\
 &= c_1 \log(|\mathcal{X}|/\delta) \lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil + \\
 &\quad c^2 (B^2 + \sigma^2) \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil} 2^{2t} \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{A}_t} \|\mathbf{x}\|^2_{\left(\sum_{\mathbf{z} \in \mathcal{X}} \lambda_t(\mathbf{z}) \phi(\mathbf{z}) \phi(\mathbf{z})^T + \gamma I\right)^{-1}} \cdot \log(2t^2 |\mathcal{X}|^2/\delta) \\
 &\leq c_1 \log(|\mathcal{X}|/\delta) \lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil + \\
 &\quad c^2 (B^2 + \sigma^2) \log\left(\frac{4|\mathcal{X}|^2 \lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil^2}{\delta}\right) \\
 &\quad \cdot \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil} 2^{2t} \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{A}_t} \|\mathbf{x}\|^2_{\left(\sum_{\mathbf{z} \in \mathcal{X}} \lambda_t(\mathbf{z}) \phi(\mathbf{z}) \phi(\mathbf{z})^T + \gamma I\right)^{-1}} \\
 &\stackrel{\mathcal{A}_t \subset \mathcal{S}_t}{\leq} c_1 \log(|\mathcal{X}|/\delta) \lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil + \\
 &\quad c^2 (B^2 + \sigma^2) \log\left(\frac{4|\mathcal{X}|^2 \lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil^2}{\delta}\right) \\
 &\quad \cdot \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil} 2^{2t} \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{S}_t} \|\mathbf{x}\|^2_{\left(\sum_{\mathbf{z} \in \mathcal{X}} \lambda_t(\mathbf{z}) \phi(\mathbf{z}) \phi(\mathbf{z})^T + \gamma I\right)^{-1}}.
 \end{aligned}$$

It remains to control the final summation. To do so, note that

$$\begin{aligned}
 &\frac{1}{\lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil} \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil} 2^{2t} \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{S}_t} \|\mathbf{x}\|^2_{\left(\sum_{\mathbf{z} \in \mathcal{X}} \lambda_t(\mathbf{z}) \phi(\mathbf{z}) \phi(\mathbf{z})^T + \gamma I\right)^{-1}} \\
 &\leq \max_{t \leq \lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil} \min_{\lambda \in \Delta_{\mathcal{X}}} 2^{2t} \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{S}_t} \|\mathbf{x}\|^2_{\left(\sum_{\mathbf{z} \in \mathcal{X}} \lambda_t(\mathbf{z}) \phi(\mathbf{z}) \phi(\mathbf{z})^T + \gamma I\right)^{-1}} \\
 &\leq \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{t \leq \lceil \log_2(4(\Delta_{\min}(\alpha) \vee \tilde{\beta})^{-1}) \rceil} 2^{2t} \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{S}_t} \|\mathbf{x}\|^2_{\left(\sum_{\mathbf{z} \in \mathcal{X}} \lambda_t(\mathbf{z}) \phi(\mathbf{z}) \phi(\mathbf{z})^T + \gamma I\right)^{-1}} \\
 &\leq 16 \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x}} \frac{\|\phi(\mathbf{x})\|^2_{\left(\sum_{\mathbf{z} \in \mathcal{X}} \lambda_t(\mathbf{z}) \phi(\mathbf{z}) \phi(\mathbf{z})^T + \gamma I\right)^{-1}}}{\max\{(\phi(\mathbf{x})^T \theta_* - \alpha)^2, \tilde{\beta}^2\}}
 \end{aligned}$$

Plugging this along with  $c = 4$  and  $c_1 = 2$  for Theorem C.1 from RIPS with the Catoni estimator in completes the proof.  $\square$

## E Proofs for Implicit Level Set Estimation

### E.1 Lower Bounds

*Proof of Theorem 4.2.* Recall that in this setting,  $h = 0$  and  $\phi(\mathbf{x}) = \mathbf{x}$ . By (Kaufmann et al., 2016), we have that the any  $\delta$ -PAC algorithm for all- $\epsilon$  requires

$$\min_{\lambda} \frac{KL(1 - \delta, \delta)}{\min_{\theta' \in \text{Alt}(\theta_*)} \|\theta' - \theta_*\|_{A(\lambda)}}$$

where  $\text{Alt}(\theta_*)$  is the set of alternates such that  $G_\epsilon(\theta_*) \neq G_\epsilon(\theta')$  for any  $\theta' \in \text{Alt}(\theta_*)$ . The set of alternates may be decomposed as

$$\mathbf{Alt}(\theta_*) = \left( \bigcup_{\mathbf{x} \in G_\epsilon(\theta_*)} \{\theta' : \mathbf{x} \notin G_\epsilon(\theta')\} \right) \cup \left( \bigcup_{\mathbf{x} \in G_\epsilon(\theta_*)^c} \{\theta' : \mathbf{x} \in G_\epsilon(\theta')\} \right)$$

By Lemma 4.1,  $\mathbf{x} \in G_\epsilon \iff \forall \mathbf{x}' : \theta_*^T(\mathbf{x} - (1 - \epsilon)\mathbf{x}') > 0$ . Hence, the set of alternates for any  $\mathbf{x} \in G_\epsilon(\theta_*)$  such that  $\mathbf{x} \in G_\epsilon^c(\theta')$  for any  $\theta' \in \mathbf{Alt}(\theta_*)$  is given by

$$A_{\mathbf{x}} := \bigcup_{\mathbf{x}' \in \mathcal{X}} \{\theta \in \mathbb{R}^d : \theta^T(\mathbf{x} - (1 - \epsilon)\mathbf{x}') < 0\}.$$

Furthermore, by Lemma 4.1  $\mathbf{x} \in G_\epsilon^c \iff \exists \mathbf{x}' : \theta_*^T(\mathbf{x} - (1 - \epsilon)\mathbf{x}') < 0$ . Hence, for any  $\mathbf{x} \in G_\epsilon^c(\theta_*)$  the set of alternates such that  $\mathbf{x} \in G_\epsilon(\theta')$  for any  $\theta' \in \mathbf{Alt}(\theta_*)$  is given by

$$A_{\mathbf{x}} := \bigcap_{\mathbf{x}' \in \mathcal{X}} \{\theta \in \mathbb{R}^d : \theta^T(\mathbf{x} - (1 - \epsilon)\mathbf{x}') > 0\}.$$

Next, we discuss how to project onto  $A_{\mathbf{x}}$ . As this set is open, to be precise, we should take a point in the interior and consider the limit for a sequence approaching the boundary. For brevity, we simply project onto the closure and consider the closures of the  $A_{\mathbf{x}}$  sets. Using the decomposition of  $\mathbf{Alt}(\theta_*)$  we have that

$$\min_{\theta' \in \mathbf{Alt}(\theta_*)} \|\theta' - \theta_*\|_{A(\lambda)} = \min_{\mathbf{x}} \min_{\theta' \in A_{\mathbf{x}}} \|\theta' - \theta_*\|_{A(\lambda)} = \min_{\mathcal{S} \in \{G_\epsilon, G_\epsilon^c\}} \min_{\mathbf{x} \in \mathcal{S}} \min_{\theta \in A_{\mathbf{x}}} \|\theta' - \theta_*\|_{A(\lambda)}.$$

Reminiscent of Lemma D.1, we define

$$\theta_{\mathbf{x}, \mathbf{x}'}^\epsilon(\lambda) := \theta_* - [\theta_*^T(\mathbf{x} - (1 - \epsilon)\mathbf{x}')] \frac{\mathcal{A}(\lambda)^{-1}(\mathbf{x} - (1 - \epsilon)\mathbf{x}')}{\|\mathbf{x} - (1 - \epsilon)\mathbf{x}'\|_{\mathcal{A}(\lambda)}^2}.$$

For  $\mathbf{x} \in G_\epsilon(\theta_*)$ , using Lemma D.1,

$$\min_{\theta' \in A_{\mathbf{x}}} \|\theta' - \theta_*\|_{A(\lambda)} = \min_{\theta' \in \bigcup_{\mathbf{x}' \in \mathcal{X}} \{\theta \in \mathbb{R}^d : \theta^T(\mathbf{x} - (1 - \epsilon)\mathbf{x}') < 0\}} \|\theta' - \theta_*\|_{A(\lambda)} = \min_{\mathbf{x}'} \|\theta_{\mathbf{x}, \mathbf{x}'}^\epsilon(\lambda) - \theta_*\|_{A(\lambda)}$$

where the latter equality follows since projecting onto a union of hyperplanes is achieved by the projection onto the closest constituent.

For  $\mathbf{x} \in G_\epsilon^c(\theta_*)$  note that  $A_{\mathbf{x}}$  is an intersection of half spaces  $\{\theta \in \mathbb{R}^d : \theta^T(\mathbf{x} - (1 - \epsilon)\mathbf{x}') > 0\}$  for  $\mathbf{x}' \in \mathcal{X}$ . As it is not in general possible to give a closed form expression for projection onto an intersection of convex sets. However, we may at a (possibly very loose) minimum note that the projection onto the union of the hyperplanes is at least as far as the projection onto the furthest hyperplane. Therefore, for any  $\mathbf{x} \in G_\epsilon(\theta_*)^c$ ,

$$\min_{\theta' \in A_{\mathbf{x}}} \|\theta' - \theta_*\|_{A(\lambda)} = \min_{\theta' \in \bigcap_{\mathbf{x}' \in \mathcal{X}} \{\theta \in \mathbb{R}^d : \theta^T(\mathbf{x} - (1 - \epsilon)\mathbf{x}') > 0\}} \|\theta' - \theta_*\|_{A(\lambda)} \leq \max_{\mathbf{x}'} \|\theta_{\mathbf{x}, \mathbf{x}'}^\epsilon(\lambda) - \theta_*\|_{A(\lambda)}$$

Hence we have that

$$\min_{\theta' \in \mathbf{Alt}(\theta_*)} \|\theta' - \theta_*\|_{A(\lambda)} \leq \min \left\{ \min_{\mathbf{x} \in G_\epsilon} \min_{\mathbf{x}'} \|\theta_{\mathbf{x}, \mathbf{x}'}^\epsilon(\lambda) - \theta_*\|_{A(\lambda)}, \min_{\mathbf{x} \in G_\epsilon^c} \max_{\mathbf{x}'} \|\theta_{\mathbf{x}, \mathbf{x}'}^\epsilon(\lambda) - \theta_*\|_{A(\lambda)} \right\}.$$

Note that

$$\|\theta_{\mathbf{x}, \mathbf{x}'}^\epsilon(\lambda) - \theta_*\|_{A(\lambda)} = 2 \frac{(\theta_*^T(\mathbf{x} - (1 - \epsilon)\mathbf{x}'))^2}{\|\mathbf{x} - (1 - \epsilon)\mathbf{x}'\|_{\mathcal{A}(\lambda)}^2}$$

by Theorem 2 of (Fiez et al., 2019). Hence, any  $\delta$ -PAC algorithm requires

$$2 \min_{\lambda} \max \left\{ \max_{\mathbf{x} \in G_\epsilon} \max_{\mathbf{x}'} \frac{\|\mathbf{x} - (1 - \epsilon)\mathbf{x}'\|_{\mathcal{A}(\lambda)}^2}{(\theta_*^T(\mathbf{x} - (1 - \epsilon)\mathbf{x}'))^2}, \max_{\mathbf{x} \in G_\epsilon^c} \min_{\mathbf{x}'} \frac{\|\mathbf{x} - (1 - \epsilon)\mathbf{x}'\|_{\mathcal{A}(\lambda)}^2}{(\theta_*^T(\mathbf{x} - (1 - \epsilon)\mathbf{x}'))^2} \right\} KL(1 - \delta, \delta)$$

samples in expectation. Noting that  $KL(1 - \delta, \delta) \geq \log(1/2.4\delta)$  completes the proof.  $\square$

## E.2 Comparison to the lower bound of (Mason et al., 2020)

Here, we compare the sample complexity given in Theorem 4.3 to the result of Mason et al., (Mason et al., 2020) studying the problem of finding all  $\epsilon$ -good arms in multi-armed bandits. Our setting captures this problem in the special case that  $\phi(\mathbf{x}) = \mathbf{x}$ ,  $\mathbf{x}_i = e_i \in \mathbb{R}^{|\mathcal{X}|}$ ,  $h = 0$ , and  $\tilde{\beta} = 0$ . Additionally, take  $\gamma \rightarrow 0$ . For consistency with the notation of (Mason et al., 2020), let  $\mu_i = f(\mathbf{x}_i)$  and  $|\mathcal{X}| = n$ . In this setting, the problem of implicit level set estimation reduces to identifying the set  $\{i : \mu_i > (1 - \epsilon)\mu_1\}$  where we assume without loss of generality that the means are sorted in descending order such that  $\mu_1 \geq \mu_2 \geq \dots \geq \mu_n$ .

**Lemma E.1.** *The term  $H^{\text{MILK}}(\theta_*) = cH_{(ST)^2}$  for a constant  $c$  where  $H_{(ST)^2}$  is the complexity parameter of the  $(ST)^2$  algorithm from (Mason et al., 2020).*

In particular, (Mason et al., 2020) show in Theorem 4.1 that a complexity of  $H_{(ST)^2}$  is optimal up to logarithmic factors for any fixed  $\delta$  via a moderate confidence bound. This exceeds the lower bound given in Theorem 4.2 specialized to this case. In particular, this highlights that the lower bound given in Theorem 4.2 is not achievable except possibly as  $\delta \rightarrow 0$ . Instead, we show that MILK achieves the optimal non-asymptotic sample complexity for finding all  $\epsilon$ -good arms.

*Proof of Lemma E.1.* First, we recall some notation from (Mason et al., 2020) necessary for this lemma. Let  $\tilde{\alpha}_\epsilon = \min_{i \in G_\epsilon} \mu_i - (1 - \epsilon)\mu_1$  and let  $\tilde{\beta}_\epsilon = \min_{i \in G_\epsilon^c} (1 - \epsilon)\mu_1 - \mu_i$ . For brevity, we let  $k = \arg \min_{i \in G_\epsilon} \mu_i$  and  $k + 1 = \arg \max_{i \in G_\epsilon^c} \mu_i$  where we take  $n > k$ . If this condition does not hold the same argument as below suffices ignoring all terms in  $G_\epsilon^c$ . Hence we have that  $\frac{\mu_k}{1 - \epsilon} = \mu_1 + \frac{\tilde{\alpha}_\epsilon}{1 - \epsilon}$  and  $\frac{\mu_{k+1}}{1 - \epsilon} = \mu_1 - \frac{\tilde{\beta}_\epsilon}{1 - \epsilon}$ . Furthermore, (Mason et al., 2020) restrict to the case of  $\epsilon \in [1/2, 1)$ .

We begin by lower bounding the complexity parameter  $H^{\text{MILK}}(\theta_*)$ . We analyze the two terms given in Theorem 4.3,  $H^{\text{MILK1}}$  and  $H^{\text{MILK2}}$  separately.  $H^{\text{MILK1}}$  reduces to

$$\begin{aligned} \max_{e_i \in G_\epsilon} \max_{e_j} \frac{\|e_j - e_i\|_{A(\lambda)^{-1}}^2}{(\mu_i - (1 - \epsilon)\mu_j)^2} &= \max_{e_i \in G_\epsilon} \max_{e_j} \frac{1/\lambda_i + 1/\lambda_j}{(\mu_i - (1 - \epsilon)\mu_j)^2} \\ &\geq \max \left\{ \max_{e_i \in G_\epsilon} \frac{1/\lambda_i}{(\mu_i - (1 - \epsilon)\mu_1)^2}, \max_{e_j} \frac{1/\lambda_j}{\left(\frac{\mu_k}{1 - \epsilon} - \mu_j\right)^2} \right\} \\ &= \max \left\{ \max_{e_i \in G_\epsilon} \frac{1/\lambda_i}{(\mu_1 - \mu_i - \epsilon)^2}, \max_{e_j} \frac{1/\lambda_j}{\left(\mu_1 + \frac{\tilde{\alpha}_\epsilon}{1 - \epsilon} - \mu_j\right)^2} \right\} \end{aligned}$$

where the final step follows by the definition of  $\tilde{\alpha}_\epsilon$ . The penultimate step follows by first maximizing over  $i \in G_\epsilon$  which introduces a factor of  $\mu_k$ . Then we may multiply the denominator by  $(1 - \epsilon)^2/(1 - \epsilon)^2$  and upper bound  $(1 - \epsilon)^2 \leq 0.25 < 1$  since  $\epsilon \geq 1/2$  to achieve the result.

$H^{\text{MILK2}}$  reduces to

$$\begin{aligned} \max_{e_i \in G_\epsilon^c} \max_{e_j} \frac{\|e_j - e_i\|_{A(\lambda)^{-1}}^2}{((1 - \epsilon)\mu_1 - \mu_i)^2} &= \max_{e_i \in G_\epsilon^c} \max_{e_j} \frac{1/\lambda_i + 1/\lambda_j}{((1 - \epsilon)\mu_1 - \mu_i)^2} \\ &\geq \max \left\{ \max_{e_i \in G_\epsilon^c} \frac{1/\lambda_i}{((1 - \epsilon)\mu_1 - \mu_i)^2}, \max_{e_i \in G_\epsilon^c} \max_{e_j} \frac{1/\lambda_j}{((1 - \epsilon)\mu_1 - \mu_i)^2} \right\} \\ &\geq \max \left\{ \max_{e_i \in G_\epsilon^c} \frac{1/\lambda_i}{((1 - \epsilon)\mu_1 - \mu_i)^2}, \max_{e_j} \frac{1/\lambda_j}{((1 - \epsilon)\mu_1 - \mu_{k+1})^2} \right\} \\ &\geq \max \left\{ \max_{e_i \in G_\epsilon^c} \frac{1/\lambda_i}{((1 - \epsilon)\mu_1 - \mu_i)^2}, \max_{e_j} \frac{1/\lambda_j}{\left(\mu_1 - \frac{\mu_{k+1}}{1 - \epsilon}\right)^2} \right\} \\ &= \max \left\{ \max_{e_i \in G_\epsilon^c} \frac{1/\lambda_i}{((1 - \epsilon)\mu_1 - \mu_i)^2}, \max_{e_j} \frac{1/\lambda_j}{\left(\left(\mu_1 - \frac{\tilde{\beta}_\epsilon}{1 - \epsilon}\right) - \mu_1\right)^2} \right\} \\ &= \max \left\{ \max_{e_i \in G_\epsilon^c} \frac{1/\lambda_i}{((1 - \epsilon)\mu_1 - \mu_i)^2}, \max_{e_j} \frac{1/\lambda_j}{\left(\left(\mu_1 + \frac{\tilde{\beta}_\epsilon}{1 - \epsilon}\right) - \mu_1\right)^2} \right\} \end{aligned}$$



$$\geq \max \left\{ \max_{e_i \in G_\epsilon^c} \frac{1/\lambda_i}{((1-\epsilon)\mu_1 - \mu_i)^2}, \max_{e_j} \frac{1/\lambda_j}{\left(\mu_1 + \frac{\tilde{\beta}_\epsilon}{1-\epsilon} - \mu_j\right)^2} \right\}$$

where the final step follows since  $\mu_1 + \frac{\tilde{\beta}_\epsilon}{1-\epsilon} > \mu_i \forall i$  and  $\mu_j \leq \mu_1$ . The third inequality follows by the same approach as taken for  $H^{\text{MILK1}}$  of multiplying the denominator by  $(1-\epsilon)^2/(1-\epsilon)^2$ .

Hence, we have that

$$H(\theta_*) \geq \min_{\lambda} \max_i \max \left\{ \frac{\frac{1}{\lambda_i}}{((1-\epsilon)\mu_1 - \mu_i)^2}, \frac{\frac{1}{\lambda_i}}{\left(\mu_1 + \frac{\tilde{\alpha}_\epsilon}{1-\epsilon} - \mu_i\right)^2}, \frac{\frac{1}{\lambda_i}}{\left(\mu_1 + \frac{\tilde{\beta}_\epsilon}{1-\epsilon} - \mu_i\right)^2} \right\}.$$

Solving for  $\lambda$  gives

$$H(\theta_*) \geq \sum_{i=1}^n \max \left\{ \frac{1}{((1-\epsilon)\mu_1 - \mu_i)^2}, \frac{1}{\left(\mu_1 + \frac{\tilde{\alpha}_\epsilon}{1-\epsilon} - \mu_i\right)^2}, \frac{1}{\left(\mu_1 + \frac{\tilde{\beta}_\epsilon}{1-\epsilon} - \mu_i\right)^2} \right\} = c_1 \cdot H_{(ST)^2}$$

for a constant  $c_1$ . To upper bound  $H^{\text{MILK}}(\theta_*)$ , we may choose a specific  $\lambda$ . Choosing

$$\lambda_i := \frac{\max\{((1-\epsilon)\mu_1 - \mu_i)^{-2}, (\mu_1 + \frac{\tilde{\alpha}_\epsilon}{1-\epsilon} - \mu_i)^{-2}, (\mu_1 + \frac{\tilde{\beta}_\epsilon}{1-\epsilon} - \mu_i)^{-2}\}}{\sum_j \max\{((1-\epsilon)\mu_1 - \mu_j)^{-2}, (\mu_1 + \frac{\tilde{\alpha}_\epsilon}{1-\epsilon} - \mu_j)^{-2}, (\mu_1 + \frac{\tilde{\beta}_\epsilon}{1-\epsilon} - \mu_j)^{-2}\}},$$

a similar computation shows that  $H^{\text{MILK}}(\theta_*) \leq c_2 H_{(ST)^2}$  for a constant  $c_2$ .  $\square$

### E.3 Upper Bound

First we restate Theorem 4.3 bounding the sample complexity of MILK.

**Theorem E.2.** Fix  $\delta > 0$ , threshold  $\alpha > 0$ , tolerance  $\tilde{\beta}$ , and regularization  $\gamma > 0$ . Define the quantities  $\Delta_{\min}^{\text{Above}}(\epsilon) = \min_{\mathbf{x} \in G_\epsilon} \min_{\mathbf{x}'} \theta_*^\top (\phi(\mathbf{x}) - (1-\epsilon)\phi(\mathbf{x}'))$  and  $\Delta_{\min}^{\text{Below}}(\epsilon) = \min_{\mathbf{x} \in G_\epsilon^c} \max_{\mathbf{x}': (\phi(\mathbf{x}) - (1-\epsilon)\phi(\mathbf{x}'))^\top \theta_* < 0} (\phi(\mathbf{x}) - (1-\epsilon)\phi(\mathbf{x}'))^\top \theta_*$ , and  $\Delta_{\min} = \min\{\Delta_{\min}^{\text{Above}}(\epsilon), \Delta_{\min}^{\text{Below}}(\epsilon)\}$ . Define also

$$\tilde{\beta}(\epsilon) = \min\{\beta > 0 : 4(\sqrt{\gamma}\|\theta_*\| + h)(2 + \sqrt{f(\mathcal{X}, \{\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{X} \times \mathcal{X}) : |\mathbf{y}^\top \theta_*| \leq \beta\}; \gamma)}) \leq \beta\}.$$

With probability  $1 - \delta$ , MILK returns a set  $\hat{\mathcal{R}} = (\mathcal{X} \setminus \hat{B}_t)$  at a time  $T_\delta$  such that

$$\{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \geq (1-\epsilon)f(\mathbf{x}_*) + \tilde{\beta}(\epsilon)\} \subseteq \hat{\mathcal{R}} \subseteq \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \geq (1-\epsilon)f(\mathbf{x}_*) - \tilde{\beta} - \tilde{\beta}(\epsilon)\}$$

and for any  $\alpha, \tilde{\beta}$  such that  $\max(\Delta_{\min}(\epsilon), \tilde{\beta}) \geq \tilde{\beta}(\epsilon)$

$$T_\delta \leq 256(B^2 + \sigma^2)H^{\text{MILK}}(\theta_*) \log \left( \frac{4|\mathcal{X}|^2 \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil^2}{\delta} \right) + 2 \log \left( \frac{|\mathcal{X}|}{\delta} \right) \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil$$

for a sufficiently large constant  $c$  where  $H^{\text{MILK}}(\theta_*) = \min_{\lambda \in \Delta_{\mathcal{X}}} \max\{H_\lambda^{\text{MILK1}}(\theta_*), H_\lambda^{\text{MILK2}}(\theta_*)\}$  and

$$H_\lambda^{\text{MILK1}}(\theta_*) := \max_{\mathbf{x} \in G_\epsilon} \max_{\mathbf{x}' \in \mathcal{X}} \frac{\|\phi(\mathbf{x}) - (1-\epsilon)\phi(\mathbf{x}')\|_{(A(\lambda) + \gamma I)^{-1}}^2}{\max\{((\phi(\mathbf{x}) - (1-\epsilon)\phi(\mathbf{x}'))^\top \theta_*)^2, \tilde{\beta}^2\}}$$

$$H_\lambda^{\text{MILK2}}(\theta_*) := \max_{\mathbf{x} \in G_\epsilon^c} \max_{\mathbf{x}'} \frac{\|\phi(\mathbf{x}) - (1-\epsilon)\phi(\mathbf{x}')\|_{(A(\lambda) + \gamma I)^{-1}}^2}{\max\{((\phi(\mathbf{x}) - (1-\epsilon)\phi(\mathbf{x}_*))^\top \theta_*)^2, \tilde{\beta}^2\}}.$$

Now we show a high probability concentration result that we will use for the remainder of this section.

**Lemma E.3.** For any  $\mathcal{V} \subset \mathcal{Y}^\epsilon(\mathcal{X} \times \mathcal{X})$  define  $f(\mathcal{X}, \mathcal{V}; \gamma) = \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{v} \in \mathcal{V}} \|\mathbf{v}\|_{(\sum_{\mathbf{x} \in \mathcal{X}} \lambda_{\mathbf{x}} \phi(\mathbf{x}) \phi(\mathbf{x})^\top + \gamma I)^{-1}}$ . In each round  $t$ , define the event

$$\mathcal{E}_t = \{|\mathbf{y}^T(\widehat{\theta}_t - \theta_*)| \leq 2^{-t} + (\sqrt{\gamma}\|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{Y}^\epsilon(\mathcal{A}_t)); \gamma)} \forall \mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{A}_t)\}$$

Holds  $\mathbb{P}(\bigcup_{t=1}^{\infty} \mathcal{E}_t^c) \leq \delta$ .

*Proof.* Using Theorem 3.2, for any  $\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{A}_t)$  we have that with probability at least  $1 - \delta_t/|\mathcal{X}|^2$

$$\begin{aligned} |\mathbf{y}^T(\widehat{\theta}_t - \theta_*)| &\leq \|\mathbf{y}\|_{(\sum_{\mathbf{x} \in \mathcal{X}} \lambda_{\mathbf{x}} \phi(\mathbf{x}) \phi(\mathbf{x})^\top + \gamma I)^{-1}} \left( \sqrt{\gamma}\|\theta_*\| + h + c \sqrt{\frac{(B^2 + \sigma^2)}{N_t} \log(2t^2|\mathcal{X}|^2/\delta)} \right) \\ &\leq \sqrt{f(\mathcal{X}, \mathcal{Y}^\epsilon(\mathcal{A}_t); \gamma)} \left( \sqrt{\gamma}\|\theta_*\| + h + 2^{-t} / \sqrt{f(\mathcal{X}, \mathcal{Y}^\epsilon(\mathcal{A}_t); \gamma)} \right) \\ &\leq 2^{-t} + (\sqrt{\gamma}\|\theta_*\| + h) \sqrt{f(\mathcal{X}, \mathcal{Y}^\epsilon(\mathcal{A}_t); \gamma)} \end{aligned}$$

Since  $|\mathcal{Y}^\epsilon(\mathcal{A}_t)| \leq |\mathcal{X}|^2$ ,  $\mathcal{E}_t$  holds for all  $\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{A}_t)$  with probability  $1 - \delta_t$  via a union bound. Taking a second union bound over rounds, we have that

$$\mathbb{P}\left(\bigcup_{t=1}^{\infty} \mathcal{E}_t^c\right) \leq \sum_{t=1}^{\infty} \mathbb{P}(\mathcal{E}_t^c) \leq \sum_{t=1}^{\infty} \delta_t = \sum_{t=1}^{\infty} \frac{\delta}{2t^2} \leq \delta$$

□

Define

$$\bar{t} = \max\{t : (\sqrt{\gamma}\|\theta_*\|_2 + h)(2 + \sqrt{f(\mathcal{X}, \{\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{X} \times \mathcal{X}) : |\mathbf{y}^T \theta_*| \leq 2^{-t+2}\}; \gamma)}) \leq 2^{-t}\}.$$

As we will see in Lemmas E.6 and E.7,

$$\mathcal{Y}^\epsilon(\mathcal{A}_t) \subset \{\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{X} \times \mathcal{X}) : |\mathbf{y}^T \theta_*| \leq 2^{-t+1}\}.$$

Thus for  $t \leq \bar{t}$ , holds on  $\bigcap_t \mathcal{E}_t$  that

$$\forall \mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{A}_t), |\mathbf{y}^T(\widehat{\theta}_t - \theta_*)| \leq 2 \cdot 2^{-t}.$$

**Lemma E.4.** On  $\bigcap_t \mathcal{E}_t$ , when  $t \leq \bar{t}$  holds  $\widehat{G}_t \subset G_\epsilon^\phi := \{\mathbf{x} : (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* > 0 \forall \mathbf{x}' \in \mathcal{X}\}$ .

*Proof.*

$$\begin{aligned} \mathbf{x} \in \widehat{G}_t &\iff \forall \mathbf{x}' \exists t_{x'} \leq \bar{t} : (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \widehat{\theta}_{t_{x'}} \geq 2 \cdot 2^{-t_{x'}} \\ &\iff \forall \mathbf{x}' \exists t_{x'} \leq \bar{t} : (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T (\widehat{\theta}_{t_{x'}} - \theta_*) + (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* \geq 2 \cdot 2^{-t_{x'}} \\ &\stackrel{\bigcap_t \mathcal{E}_t}{\implies} \forall \mathbf{x}' : (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* > 0 \\ &\iff \mathbf{x} \in G_\epsilon^\phi. \end{aligned}$$

□

**Lemma E.5.** On  $\bigcap_t \mathcal{E}_t$ , when  $t \leq \bar{t}$  holds  $\widehat{B}_t \subset (G_\epsilon^\phi)^c$ .

*Proof.*

$$\begin{aligned} \mathbf{x} \in \widehat{B}_t &\iff \exists \mathbf{x}', t_{x'} \leq \bar{t} : (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \widehat{\theta}_t \leq -2 \cdot 2^{-t_{x'}} \\ &\iff \exists \mathbf{x}', t_{x'} \leq \bar{t} : (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T (\widehat{\theta}_t - \theta_*) + (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* \leq -2 \cdot 2^{-t_{x'}} \\ &\stackrel{\bigcap_t \mathcal{E}_t}{\implies} \exists \mathbf{x}' : (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* > \epsilon \\ &\iff \mathbf{x} \in (G_\epsilon^\phi)^c. \end{aligned}$$

□

**Lemma E.6.** On the event  $\bigcap_t \mathcal{E}_t$  for  $t \leq \bar{t}$ ,

$$\{(\mathbf{x}, \mathbf{x}') : (\mathbf{x}, \mathbf{x}'), \mathbf{x} \in G_\epsilon^\phi\} \subset \left\{ (\mathbf{x}, \mathbf{x}') \mid |(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_*| \leq 2^{-t+2} \right\} =: \mathcal{S}_t^{\text{Above}}$$

*Proof.* On  $\bigcap_t \mathcal{E}_t$  for  $t \leq \bar{t}$ , for any  $\mathbf{y} \in \mathcal{A}_t$

$$|\mathbf{y}^T \hat{\theta}_t| \geq |\mathbf{y}^T \theta_*| - |\mathbf{y}^T (\hat{\theta}_t - \theta_*)| \stackrel{\mathcal{E}_t}{\geq} |\mathbf{y}^T \theta_*| - 2 \cdot 2^{-t}.$$

For  $\mathbf{y}$  such that  $|\mathbf{y}^T \theta_*| \geq 2 \cdot 2^{-t+1}$ , the above implies that

$$|\mathbf{y}^T \hat{\theta}_t| \geq 2 \cdot 2^{-t}.$$

By the elimination condition, this implies that  $\mathbf{y}$  is removed from  $\mathcal{A}_t$ . Hence

$$\mathcal{A}_{t+1} \subset \{\mathbf{y} \in \mathcal{Y}(\mathcal{X}) : |\mathbf{y}^T \theta_* - \epsilon| \leq 2 \cdot 2^{-t+1}\}.$$

Specializing this argument to  $\{(\mathbf{x}, \mathbf{x}') : (\mathbf{x}, \mathbf{x}'), \mathbf{x} \in G_\epsilon^\phi\} \subset \mathcal{A}_t$  completes the proof.  $\square$

**Lemma E.7.** On the event  $\bigcap_t \mathcal{E}_t$  for  $t \leq \bar{t}$ ,

$$\{(\mathbf{x}, \mathbf{x}') : (\mathbf{x}, \mathbf{x}'), \mathbf{x} \in (G_\epsilon^\phi)^c\} \subset \left\{ (\mathbf{x}, \mathbf{x}') \mid |(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* - \epsilon| \leq 2^{-t+2} \right. \\ \left. \text{and } \{(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}_*))^T \theta_*\} \geq -2^{-t+2} \right\} =: \mathcal{S}_t^{\text{Below}}$$

*Proof.* The guarantee that  $|(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* - \epsilon| \leq 2^{-t+2}$  for any  $(\mathbf{x}, \mathbf{x}') \in \mathcal{A}_t$  follows by the same argument as Lemma E.6. For the additional statement, that  $(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}_*))^T \theta_* \geq -2^{-t+2}$ , note that if

$$(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}_*))^T \hat{\theta}_t \leq -2^{-t+1}$$

then the pair  $(\mathbf{x}, \mathbf{x}_*)$  is eliminated from  $\mathcal{A}_t$ . If

$$(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}_*))^T \hat{\theta}_t \leq -2^{-t+2},$$

then using this and the event  $\bigcap_t \mathcal{E}_t$

$$(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}_*))^T \hat{\theta}_t = (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}_*))^T (\hat{\theta}_t - \theta_*) + (1 - \epsilon)\phi(\mathbf{x}_*)^T \theta_* \leq -2^{-t+1}.$$

Hence, the only pairs  $(\mathbf{x}, \mathbf{x}_*)$  that remain in  $\mathcal{A}_t$  where  $\mathbf{x}_* \in (G_\epsilon^\phi)^c$  are such that  $(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}_*))^T \theta_* \geq -2^{-t+2}$ . We conclude by noting that the above argument for  $\mathbf{x}_*$  could be repeated for any  $\mathbf{x}'$  such that  $(\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* < 0$ .  $\square$

**Remark:** Lemmas E.6 and E.7 jointly imply that  $\mathcal{A}_t \subset \mathcal{S}_t^{\text{Above}} \cup \mathcal{S}_t^{\text{Below}} =: \mathcal{S}_t$  for  $t \leq \bar{t}$ . Furthermore,  $f(\mathcal{X}, \mathcal{Y}^\epsilon(\mathcal{A}_t), \gamma) \leq f(\mathcal{X}, \mathcal{Y}^\epsilon(\mathcal{S}_t), \gamma)$ .

**Remark:**

The algorithm stops on either of two conditions. On one hand if  $t \geq \lceil \log_2(4/\tilde{\beta}) \rceil =: t_\beta$ , then it has achieved precision  $\tilde{\beta}$  as desired and it terminates. Otherwise, it terminates if  $\hat{G}_t \cup \hat{B}_t = \mathcal{X}$ . This occurs when  $\tilde{\beta}$  is very small. Define the quantities  $\Delta_{\min}^{\text{Above}}(\epsilon) = \min_{\mathbf{x} \in G_\epsilon^\phi} \min_{\mathbf{x}'} \theta_*^\top (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))$  and  $\Delta_{\min}^{\text{Below}}(\epsilon) = \min_{\mathbf{x} \in G_\epsilon^c} \max_{\mathbf{x}': (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^\top \theta_* < 0} (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^\top \theta_*$ , and  $\Delta_{\min}(\epsilon) = \min\{\Delta_{\min}^{\text{Above}}(\epsilon), \Delta_{\min}^{\text{Below}}(\epsilon)\}$ . Recall

$$\begin{aligned} \bar{t} &= \max\{t : (\sqrt{\gamma}\|\theta_*\|_2 + h)(2 + \sqrt{f(\mathcal{X}, \{\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{X} \times \mathcal{X}) : |\mathbf{y}^T \theta_*| \leq 4 \cdot 2^{-t}\}; \gamma)}) \leq 2^{-t}\} \\ &= \max\{t : 4(\sqrt{\gamma}\|\theta_*\|_2 + h)(2 + \sqrt{f(\mathcal{X}, \{\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{X} \times \mathcal{X}) : |\mathbf{y}^T \theta_*| \leq 4 \cdot 2^{-t}\}; \gamma)}) \leq 4 \cdot 2^{-t}\} \\ &= -2 + \max\{t : 4(\sqrt{\gamma}\|\theta_*\|_2 + h)(2 + \sqrt{f(\mathcal{X}, \{\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{X} \times \mathcal{X}) : |\mathbf{y}^T \theta_*| \leq 2^{-t}\}; \gamma)}) \leq 2^{-t}\} \\ &= -3 + \log_2(\min\{\beta > 0 : 4(\sqrt{\gamma}\|\theta_*\|_2 + h)(2 + \sqrt{f(\mathcal{X}, \{\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{X} \times \mathcal{X}) : |\mathbf{y}^T \theta_*| \leq \beta\}; \gamma)}) \leq \beta\}). \end{aligned}$$

This defines

$$\bar{\beta}(\epsilon) = \min\{\beta > 0 : 4(\sqrt{\gamma}\|\theta_*\|_2 + h)(2 + \sqrt{f(\mathcal{X}, \{\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{X} \times \mathcal{X}) : |\mathbf{y}^T \theta_*| \leq \beta\}}; \gamma)) \leq \beta\}.$$

Let  $t_{\max}$  denote the random variable of the last round before the algorithm terminates. The following Lemmas give a guarantee on the set  $\mathcal{X} \setminus \widehat{B}_t$  at termination.

**Lemma E.8.** *On the event  $\bigcap_{t=1}^{\infty} \mathcal{E}_t$ , MILK returns a set  $(\mathcal{X} \setminus \widehat{B}_{t_{\max}})$  such that  $\{\mathbf{x} : f(\mathbf{x}) > (1 - \epsilon)f(\mathbf{x}_*) + \bar{\beta}(\epsilon)\} \subset (\mathcal{X} \setminus \widehat{B}_{t_{\max}})$ .*

*Proof.* Take any  $\mathbf{x}$  such that  $f(\mathbf{x}) > (1 - \epsilon)f(\mathbf{x}_*) + \bar{\beta}(\alpha)$  and recall that by assumption  $|f(\mathbf{x}) - \phi(\mathbf{x})^T \theta_*| \leq h$  for all  $\mathbf{x} \in \mathcal{X}$ . We consider two cases. In the first case, assume that  $t_{\max} \leq \bar{t}$ . We claim that in this case  $\exists t$  such that  $\mathbf{x} \in \widehat{B}_t$ . We prove this by contradiction. Assume not. Then  $\exists t$  and a  $\mathbf{x}'$  such that

$$\begin{aligned} & \widehat{\theta}_t^T (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')) < -2^{-t+1} \\ & \iff (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T (\widehat{\theta}_t - \theta_*) + (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* < -2^{-t+1} \\ & \stackrel{\mathcal{E}_t, t_{\max} \leq \bar{t}}{\iff} -2^{-t+1} + (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* < -2^{-t+1} \\ & \iff (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* < 0 \\ & \implies f(\mathbf{x}) - (1 - \epsilon)f(\mathbf{x}') < h + (1 - \epsilon)h \end{aligned}$$

Recall that we have assumed that  $f(\mathbf{x}) > (1 - \epsilon)f(\mathbf{x}_*) + \bar{\beta}(\alpha)$  and  $\bar{\beta}(\epsilon) > 4h$  by definition. Hence, this implies that

$$(1 - \epsilon)f(\mathbf{x}_*) - (1 - \epsilon)f(\mathbf{x}') < h + (1 - \epsilon)h - \bar{\beta}(\alpha) < 0$$

which is a contradiction since  $f(\mathbf{x}_*) \geq f(\mathbf{x}')$  by definition. Hence, we have shown in the case that  $t_{\max} \leq \bar{t}$ ,  $\{\mathbf{x} : f(\mathbf{x}) > (1 - \epsilon)f(\mathbf{x}_*) + \bar{\beta}(\epsilon)\} \subset (\mathcal{X} \setminus \widehat{B}_{t_{\max}})$ .

In the second case, assume that  $t_{\max} > \bar{t}$  and take  $\mathbf{x}$  such that  $f(\mathbf{x}) > (1 - \epsilon)f(\mathbf{x}_*) + \bar{\beta}(\alpha)$ . We claim that  $\mathbf{x} \in \widehat{G}_{\bar{t}}$  and hence  $(\mathbf{x}, \mathbf{x}') \notin \mathcal{A}_t$  for any  $t > \bar{t}$  and thus  $\mathbf{x}$  is never added to  $\widehat{B}_t$ . This occurs if for every  $\phi(\mathbf{x}')$

$$\begin{aligned} & (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \widehat{\theta}_{\bar{t}} > 2^{-\bar{t}+1} \\ & \iff (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T (\widehat{\theta}_{\bar{t}} - \theta_*) + (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* > 2^{-\bar{t}+1} \\ & \stackrel{\mathcal{E}_{\bar{t}}}{\iff} -2^{-\bar{t}+1} + (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* \geq 2^{-\bar{t}+1} \\ & \iff (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* \geq 2^{-\bar{t}+2} \\ & \iff (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_* \geq 0.5\bar{\beta}(\epsilon) \\ & \iff f(\mathbf{x}) - (1 - \epsilon)f(\mathbf{x}') \geq 0.5\bar{\beta}(\epsilon) + h + (1 - \epsilon)h \end{aligned}$$

where the penultimate step follows by definition of  $\bar{\beta}(\epsilon)$ . Recall that  $f(\mathbf{x}) > (1 - \epsilon)f(\mathbf{x}_*) + \bar{\beta}(\alpha)$ . Hence, the above is implied by

$$\begin{aligned} & (1 - \epsilon)f(\mathbf{x}_*) + \bar{\beta}(\alpha) - (1 - \epsilon)f(\mathbf{x}') \geq 0.5\bar{\beta}(\epsilon) + h + (1 - \epsilon)h \\ & \iff \bar{\beta}(\epsilon) \geq 0.5\bar{\beta}(\epsilon) + h + (1 - \epsilon)h \end{aligned}$$

where the final step follows by noting that  $f(\mathbf{x}_*) \geq f(\mathbf{x}')$  for any  $\mathbf{x}'$ . The final statement is true since  $\bar{\beta}(\epsilon)$  and thus implies the claim. Therefore, we have shown that  $\mathbf{x} \in \widehat{G}_{\bar{t}}$  and is therefore not added to  $\widehat{B}_t$  in a later round. These two cases together complete the proof.  $\square$

**Lemma E.9.** *On the event  $\bigcap_{t=1}^{\infty} \mathcal{E}_t$ , MILK returns a set  $(\mathcal{X} \setminus \widehat{B}_{t_{\max}})$  such that  $(\mathcal{X} \setminus \widehat{B}_{t_{\max}}) \subset \{\mathbf{x} : f(\mathbf{x}) > (1 - \epsilon)f(\mathbf{x}_*) - \bar{\beta}(\epsilon) - \widetilde{\beta}\}$ .*

*Proof.* Take any  $\mathbf{x}$  such that  $f(\mathbf{x}) < (1 - \epsilon)f(\mathbf{x}_*) - \bar{\beta}(\epsilon) - \widetilde{\beta}$ . We claim that there exists a  $t \leq t_{\max}$  such that  $\mathbf{x}$  is added to  $\widehat{B}_t$  which implies that  $\mathbf{x} \notin (\mathcal{X} \setminus \widehat{B}_{t_{\max}})$ . Suppose for contradiction that this is not the case. Then for all  $t \leq t_{\max}$ ,

$$\widehat{\theta}_t^T (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}_*)) > -2^{-t+1}$$

$$\begin{aligned}
 &\iff (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}_*))^T(\widehat{\theta}_t - \theta_*) + (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}_*))^T\theta_* > -2^{-t+1} \\
 &\stackrel{\mathcal{E}_t}{\implies} 2^{-t} + (\sqrt{\gamma}\|\theta_*\| + h)\sqrt{f(\mathcal{X}, \mathcal{A}_t; \gamma)} + (\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}_*))^T\theta_* > -2^{-t+1} \\
 &\implies (\sqrt{\gamma}\|\theta_*\| + h)\sqrt{f(\mathcal{X}, \mathcal{A}_t; \gamma)} + f(\mathbf{x}) - (1 - \epsilon)f(\mathbf{x}_*) + h + (1 - \epsilon)h > -2^{-t+1} - 2^{-t} \\
 &\implies f(\mathbf{x}) - (1 - \epsilon)f(\mathbf{x}_*) > -2^{-t+1} - 2^{-t} - h - (1 - \epsilon)h - (\sqrt{\gamma}\|\theta_*\| + h)\sqrt{f(\mathcal{X}, \mathcal{S}_t; \gamma)}.
 \end{aligned}$$

Plugging in  $f(\mathbf{x}) < (1 - \epsilon)f(\mathbf{x}_*) - \bar{\beta}(\epsilon) - \tilde{\beta}$ , the above implies

$$\bar{\beta}(\epsilon) + \tilde{\beta} < 2^{-t+1} + 2^{-t} + h + (1 - \epsilon)h + (\sqrt{\gamma}\|\theta_*\| + h)\sqrt{f(\mathcal{X}, \mathcal{S}_t; \gamma)} \quad (4)$$

Next, recall that MILK terminates either on the condition that  $t = \lceil \log_2(4/\tilde{\beta}) \rceil$  or that  $\widehat{G}_t \cup \widehat{B}_t = \mathcal{X}$ . Using this, we break our analysis into cases.

Case 1:  $t_{\max} = \lceil \log_2(4/\tilde{\beta}) \rceil \leq \bar{t}$ .

In this case, MILK stops due to the  $\tilde{\beta}$  tolerance in a round before  $\bar{t}$ . For  $t \leq \bar{t}$ , we have that  $2^{-t} \geq + (\sqrt{\gamma}\|\theta_*\| + h)\sqrt{f(\mathcal{X}, \mathcal{S}_t; \gamma)}$ . Hence, the above implies that

$$\bar{\beta}(\alpha) + \tilde{\beta} < 2^{-t+2} + h + (1 - \epsilon)h.$$

As we have assumed this condition for all  $t \leq t_{\max}$ , we may plug in  $t_{\max}$  which implies

$$\bar{\beta}(\alpha) + \tilde{\beta} < \tilde{\beta} + h + (1 - \epsilon)h.$$

As  $\bar{\beta}(\alpha) > 4h$ , this is a contradiction. Hence there must exist a  $t$  such that  $\mathbf{x} \in \widehat{B}_t$ .

Case 2:  $t_{\max} \leq \bar{t} < \lceil \log_2(4/\tilde{\beta}) \rceil$ .

In this case, MILK terminates before round  $t = \lceil \log_2(4/\tilde{\beta}) \rceil$ . Hence, it does so on the condition that  $\widehat{G}_t \cup \widehat{B}_t = \mathcal{X}$ . Note that for  $f(\mathbf{x}) < \alpha - \bar{\beta}(\alpha) - \tilde{\beta}$ , we have that  $\mathbf{x} \in (G_\alpha^\phi)^c$  since  $\bar{\beta}(\alpha) > h$  and  $\tilde{\beta} \geq 0$ . If we terminate before round  $\bar{t}$ , we have by Lemma E.5 that  $(G_\alpha^\phi)^c \subset \widehat{B}_t$  which implies that  $\mathbf{x} \in \widehat{B}_{t_{\max}}$ . This contradicts the assumption that  $\nexists t : \mathbf{x} \in \widehat{B}_t$ .

Case 3:  $\bar{t} < t_{\max}$ .

In this case, MILK terminates at a round after  $\bar{t}$ . In this setting, we argue that  $\mathbf{x} \in \widehat{B}_{\bar{t}}$ . Recall that for any  $t \leq \bar{t}$ , (4) simplifies to

$$\bar{\beta}(\alpha) + \tilde{\beta} < 2^{-t+2} + h + (1 - \epsilon)h.$$

Plugging in  $\bar{t}$ , and noting that  $2^{-\bar{t}+2} = \frac{1}{2}\bar{\beta}(\alpha)$ , the above implies

$$\bar{\beta}(\alpha) + \tilde{\beta} < \frac{1}{2}\bar{\beta}(\alpha) + h + (1 - \epsilon)h.$$

Noting that  $\bar{\beta}(\alpha) > 4h$ , shows that the above is a contradiction. Hence, there exists a  $t \leq \bar{t}$  such that  $\mathbf{x} \in \widehat{B}_t$ .

Therefore, in all cases we have shown that for any  $\mathbf{x}$  such that  $f(\mathbf{x}) < \alpha - \bar{\beta}(\alpha) - \tilde{\beta}$ ,  $\mathbf{x} \in \widehat{B}_t$ . Therefore, for the returned set  $\mathcal{X} \setminus \widehat{B}_{t_{\max}}$ , we have that

$$(\mathcal{X} \setminus \widehat{B}_{t_{\max}}) \subset \{\mathbf{x} : f(\mathbf{x}) > \alpha - \bar{\beta}(\alpha) - \tilde{\beta}\}.$$

□

*Proof of Theorem 4.3.* Throughout, assume the high probability event  $\bigcap_T \mathcal{E}_t$ . By Lemmas E.8 and E.9 in conjunction with the high probability event  $\bigcap \mathcal{E}_t$  we have correctness. It remains to control the sample complexity of MILK. Recall that we have assumed that  $\max(\Delta_{\min}(\epsilon), \tilde{\beta}) \geq \bar{\beta}(\epsilon)$ . This implies that  $\min\{\lceil \log_2(4/\Delta_{\min}(\epsilon)) \rceil, \lceil \log_2(4/\tilde{\beta}) \rceil\} \leq \bar{t}$ . Applying Lemmas E.6 and E.7, we have that  $t_{\max} \leq$

$\min\{\lceil \log_2(4/\Delta_{\min}(\epsilon)) \rceil, \lceil \log_2(4/\tilde{\beta}) \rceil\} \leq \bar{t}$  and that  $\mathcal{A}_t \subseteq \mathcal{S}_t$  for all rounds  $t$ . Now we proceed by bounding the total number of samples drawn.

$$\begin{aligned}
 \tau &\leq \sum_{t=1}^{t_{\max}} N_t \\
 &\leq \sum_{t=1}^{\min\{\lceil \log_2(4/\Delta_{\min}(\epsilon)) \rceil, \lceil \log_2(4/\tilde{\beta}) \rceil\}} N_t \\
 &= \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil} N_t \\
 &= \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil} \max\{c_1 \log(|\mathcal{X}|/\delta), c^2 2^{2t} f(\mathcal{Y}^\epsilon(\mathcal{A}_t); \gamma)(B^2 + \sigma^2) \log(2t^2 |\mathcal{X}|^2/\delta)\} \\
 &\leq c_1 \log(|\mathcal{X}|/\delta) \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil + c^2 (B^2 + \sigma^2) \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil} 2^{2t} f(\mathcal{Y}^\epsilon(\mathcal{A}_t); \gamma) \cdot \log(2t^2 |\mathcal{X}|^2/\delta) \\
 &= c_1 \log(|\mathcal{X}|/\delta) \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil + \\
 &\quad c^2 (B^2 + \sigma^2) \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil} 2^{2t} \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{A}_t)} \|\mathbf{y}\|_{(A(\lambda) + \gamma I)^{-1}}^2 \cdot \log(2t^2 |\mathcal{X}|^2/\delta) \\
 &\leq c_1 \log(|\mathcal{X}|/\delta) \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil + \\
 &\quad c^2 (B^2 + \sigma^2) \log\left(\frac{4|\mathcal{X}|^2 \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil^2}{\delta}\right) \\
 &\quad \cdot \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil} 2^{2t} \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{A}_t)} \|\mathbf{y}\|_{(A(\lambda) + \gamma I)^{-1}}^2 \\
 &\leq c_1 \log(|\mathcal{X}|/\delta) \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil + \\
 &\quad c^2 (B^2 + \sigma^2) \log\left(\frac{4|\mathcal{X}|^2 \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil^2}{\delta}\right) \\
 &\quad \cdot \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil} 2^{2t} \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{S}_t)} \|\mathbf{y}\|_{(A(\lambda) + \gamma I)^{-1}}^2 \\
 &= c_1 \log(|\mathcal{X}|/\delta) \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil + \\
 &\quad c^2 (B^2 + \sigma^2) \log\left(\frac{4|\mathcal{X}|^2 \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil^2}{\delta}\right) \\
 &\quad \cdot \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil} \min_{\lambda \in \Delta_{\mathcal{X}}} \max\left\{2^{2t} \max_{\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{S}_t^{\text{Above}})} \|\mathbf{y}\|_{(A(\lambda) + \gamma I)^{-1}}^2, 2^{2t} \max_{\mathbf{y} \in \mathcal{Y}^\epsilon(\mathcal{S}_t^{\text{Below}})} \|\mathbf{y}\|_{(A(\lambda) + \gamma I)^{-1}}^2\right\}.
 \end{aligned}$$

where the final equality follows by partitioning  $\mathcal{S}_t = \mathcal{S}_t^{\text{Above}} \cup \mathcal{S}_t^{\text{Below}}$ .

Focusing on this final summation, note that

$$\begin{aligned}
 &\frac{1}{\lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil} \sum_{t=1}^{\lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil} 2^{2t} \min_{\lambda \in \Delta_{\mathcal{X}}} \max\left\{\max_{\mathbf{y} \in \mathcal{S}_t^{\text{Above}}} \|\mathbf{y}\|_{(A(\lambda) + \gamma I)^{-1}}^2, \max_{\mathbf{y} \in \mathcal{S}_t^{\text{Below}}} \|\mathbf{y}\|_{(A(\lambda) + \gamma I)^{-1}}^2\right\} \\
 &\leq \max_{t \leq \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil} \min_{\lambda \in \Delta_{\mathcal{X}}} 2^{2t} \max\left\{\max_{\mathbf{y} \in \mathcal{S}_t^{\text{Above}}} \|\mathbf{y}\|_{(A(\lambda) + \gamma I)^{-1}}^2, \max_{\mathbf{y} \in \mathcal{S}_t^{\text{Below}}} \|\mathbf{y}\|_{(A(\lambda) + \gamma I)^{-1}}^2\right\} \\
 &\leq \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{t \leq \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil} \max\left\{\max_{\mathbf{y} \in \mathcal{S}_t^{\text{Above}}} 2^{2t} \|\mathbf{y}\|_{(A(\lambda) + \gamma I)^{-1}}^2, \max_{\mathbf{y} \in \mathcal{S}_t^{\text{Below}}} 2^{2t} \|\mathbf{y}\|_{(A(\lambda) + \gamma I)^{-1}}^2\right\}
 \end{aligned}$$



$$\begin{aligned}
 &= \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{t \leq \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil} \max \left\{ \max_{(\mathbf{x}, \mathbf{x}') \in \mathcal{S}_t^{\text{above}}} 2^{2t} \|\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')\|_{(A(\lambda) + \gamma I)^{-1}}^2, \right. \\
 &\quad \left. \max_{(\mathbf{x}, \mathbf{x}') \in \mathcal{S}_t^{\text{below}}} 2^{2t} \|\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')\|_{(A(\lambda) + \gamma I)^{-1}}^2 \right\} \\
 &\stackrel{\text{Lemmas E.6, E.7, } \tilde{\beta}}{\leq} 16 \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{t \leq \lceil \log_2(4(\Delta_{\min}(\epsilon) \vee \tilde{\beta})^{-1}) \rceil} \max \left\{ \max_{(\mathbf{x}, \mathbf{x}') \in \mathcal{S}_t^{\text{above}}} \frac{\|\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')\|_{(A(\lambda) + \gamma I)^{-1}}^2}{\max\{((\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_*)^2, \tilde{\beta}^2\}}, \right. \\
 &\quad \left. \max_{(\mathbf{x}, \mathbf{x}') \in \mathcal{S}_t^{\text{below}}} \frac{\|\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')\|_{(A(\lambda) + \gamma I)^{-1}}^2}{\max\{((\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}_*))^T \theta_* - \epsilon)^2, \tilde{\beta}^2\}} \right\} \\
 &\leq 16 \min_{\lambda \in \Delta_{\mathcal{X}}} \max \left\{ \max_{\mathbf{x} \in G_\epsilon} \max_{\mathbf{x}'} \frac{\|\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')\|_{(A(\lambda) + \gamma I)^{-1}}^2}{\max\{((\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}'))^T \theta_*)^2, \tilde{\beta}^2\}}, \right. \\
 &\quad \left. \max_{\mathbf{x} \in G_\epsilon^c} \max_{\mathbf{x}'} \frac{\|\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}')\|_{(A(\lambda) + \gamma I)^{-1}}^2}{\max\{((\phi(\mathbf{x}) - (1 - \epsilon)\phi(\mathbf{x}_*))^T \theta_* - \epsilon)^2, \tilde{\beta}^2\}} \right\}
 \end{aligned}$$

Plugging this in with  $c = 4$  and  $c_1 = 2$  from Theorem C.1 for RIPS with the Catoni estimator completes the proof.  $\square$

## F Additional Experiment Details

In this section we discuss additional experimental details not covered in the main paper. We first give an overview of the algorithms implemented in the following section. All code was written in python and run on a 64 core cluster machine. We have included implementations of all methods and a demo file showing how to call and run the various algorithms.

### F.1 Algorithms Implemented

In this section we briefly discuss the algorithms implemented and the hyper-parameters used in the algorithms. The algorithms implemented are as follows:

**Gaussian Process Experiments** For all the algorithms in this section we assumed a GP Prior  $N(0, k(x, x'))$  where  $k(x, x')$  was the RBF kernel given by  $k(x, x') = \exp(-\|x - x'\|^2 / 2\ell^2)$ .

At every time step we build the confidence interval

$$Q_t(\mathbf{x}) := \left[ \mu_{t-1}(\mathbf{x}) \pm \beta_t^{1/2} \sigma_{t-1}(\mathbf{x}) \right]$$

where  $\mu_{t-1}$ , and  $\sigma_{t-1}$  is the posterior mean and variance function over the observed points. For an observation  $\mathbf{y}_t$  at time  $t$  we define  $\mu_t$ , and  $\sigma_t$  as follows:

$$\begin{aligned}
 \mu_t(\mathbf{x}) &:= \mathbf{k}_t(\mathbf{x})^T (\mathbf{K}_t + \sigma^2 \mathbf{I})^{-1} \mathbf{y}_t \\
 k_t(\mathbf{x}, \mathbf{x}') &:= k(\mathbf{x}, \mathbf{x}') - \mathbf{k}_t(\mathbf{x})^T (\mathbf{K}_t + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_t(\mathbf{x}') \\
 \sigma_t^2(\mathbf{x}) &:= k_t(\mathbf{x}, \mathbf{x})
 \end{aligned}$$

where,  $\mathbf{k}_t(\mathbf{x}) = [k(\mathbf{x}_1, \mathbf{x}), \dots, k(\mathbf{x}_t, \mathbf{x})]^T$  and  $\mathbf{K}_t$  is the kernel matrix over the observed points.

1. **LSE:** We implemented the LSE algorithm by (Gotovos, 2013). This algorithm maintains an active set of unclassified points defined as  $U_t$  and the super-level set  $H_t$  and sub-level set  $L_t$ .

At every round LSE selects the most ambiguous point, where the ambiguity is defined as

$$a_t(\mathbf{x}) = \min \{ \max(Q_t(\mathbf{x})) - \alpha, \alpha - \min(Q_t(\mathbf{x})) \}$$

that is, the points LSE is most unsure to classify into  $H_t$  or  $L_t$ . Note that in contrast to this approach MELK follows the optimal allocation over the active set to select the next sample.

2. **TruVar**: We also implemented a modified version of TruVar (Bogunovic et al., 2016) with zero cost and homoscedastic noise. TruVar samples in such a fashion to ensure the maximum decrease of the posterior variance. As above, we maintain a Gaussian Process Posterior and we sample the arm

$$\arg \max_{x \in \mathcal{X}} \sum_{\bar{x} \in \mathcal{A}_t} \sigma_t^2(\bar{x}) - \sum_{\bar{x} \in \mathcal{A}_t} \sigma_{t-1|x}^2(\bar{x})$$

where  $\sigma_{t-1|x}^2(\bar{x})$  is the posterior variance of  $\bar{x}$  if we sample  $x$ .

3. **MELK**: As described in the text, we compute the means and variances of the arms using a Gaussian posterior (identical to above) and eliminate arms when their lower/upper bound is below/above the specified threshold  $\tau$ . We implemented a batched sampling algorithm where we compute the design

$$\min_{\lambda \in \mathcal{X}} \max_{z \in \mathcal{A}_t} \|z\|_{(A(\lambda) + \gamma I)^{-2}}^2$$

ever 10 samples and then sample from it. At the  $i$ -th calculation,  $\gamma = 1/(10 * i)$ . We also use the Frank-Wolfe method to compute the optimal allocation over the active set before every round as described in Section G. We set the step-size of Frank-Wolfe method as 1 and cap the maximum number of iteration to converge for Frank-Wolfe to 500.

### Linear Bandits Examples

Additionally, we also consider comparing algorithms exactly as written using theoretically justified confidence widths in all cases. This presents a challenge as MELK and MILK are designed for the frequentist regime and LSE and TruVar are Bayesian in nature. To level the playing field, we consider all algorithms in the frequentist regime. For this experiment, we focused primarily on comparing MELK to LSE and MILK to LSE-imp. LSE can naturally be adapted to the frequentist setting with the tight RKHS confidence bounds from (Chowdhury and Gopalan, 2017). These bounds scale with the maximum information gain  $\Gamma_T$ . To make the comparison fair, we consider all algorithms in the linear regime where  $\Gamma_T = O(d \log(T))$ . By contrast, for the squared exponential kernel,  $\Gamma_T = O(\log(T)^d)$ , and this leads to overly pessimistic confidence widths preventing a meaningful comparison of the algorithms. Indeed, even for moderate  $d$  such as  $d = 4$ , LSE had confidence widths that were more than an order of magnitude wider for the squared exponential kernel. Hence, we focus on the case of the linear kernel for our experimental comparison where the differences are not so stark. Below, we describe all algorithms in this regime.

LSE follows the same acquisition function described in the previous section. We provide additional details about MELK, MILK, and LSE-imp in this setting.

1. **MELK**: We implement the MELK algorithm as defined in Algorithm 1. Recall that  $|f(x)| \leq B$ , and for the experiments we set  $B = 1$ . We set the confidence parameter  $\delta = 0.1$ , the regularization parameter  $\gamma = 1e - 7$ . Note that we use the original confidence width of  $(B^2 + \sigma^2) \log(2t^2 |\mathcal{X}|^2 / \delta)$  as stated in our algorithm, where  $\sigma^2$  is the noise parameter specific to the environment. We also use the Frank-Wolfe method to compute the optimal allocation over the active set before every round. We set the step-size of Frank-Wolfe method as 0.5 and cap the maximum number of iteration to converge for Frank-Wolfe to 2000.
2. **LSE-imp**: We implement the LSE-Implicit algorithm as stated in (Gotovos, 2013). LSE-Implicit proceeds quite similarly to LSE by constructing the confidence region  $C_t(\mathbf{x})$  (as defined above) and classifying points to the sub-level set  $L_t$  or super-level set  $H_t$ . We set the confidence width as in LSE for calculating the confidence region. Note that LSE-Implicit works in the implicit level set estimation setting and so constructs an estimate of the function maximum to classify points into  $H_t$  or  $L_t$ . It builds an optimistic and pessimistic estimate of the function maximum as

$$f_t^{opt} := \max_{x \in U_t} \max(C_t(x)), \quad f_t^{pes} = \max_{x \in U_t} \min(C_t(x))$$

respectively. A point  $\mathbf{x}$  is classified into  $H_t$  if  $\min(C_t(\mathbf{x})) \geq (1 - \epsilon) f_t^{opt}$  or classified into  $L_t$  if  $\max(C_t(\mathbf{x})) \leq (1 - \epsilon) f_t^{pes}$ . Finally, LSE-Implicit selects the next point with the largest confidence region width, defined as follows:

$$w_t(\mathbf{x}) = \max(C_t(\mathbf{x})) - \min(C_t(\mathbf{x}))$$

such that this leads to more exploration. Again, note that in contrast MILK in Algorithm 2 uses the optimal allocation proportion over the active set to sample the next point.

3. **MILK:** We implement the MILK algorithm as stated in Algorithm 2. Note that MILK proceeds as similarly to MELK but with the allocation calculated over the difference of vectors  $\mathcal{Y}^\epsilon(\mathcal{A})$  over the active set and a different elimination condition depending on  $\epsilon$ . For MILK we set a similar hyper-parameters like MELK. We set the confidence parameter  $\delta = 0.1$ , the regularization parameter  $\gamma = 1e - 7$ , and the confidence width of  $(B^2 + \sigma^2) \log(2t^2|\mathcal{X}|^2/\delta)$ . We use the Frank-Wolfe method to compute the optimal allocation over the active set of points and set the step-size of Frank-Wolfe method as 0.5 and cap the maximum number of iteration to converge for Frank-Wolfe to 2000. Note that we set  $\epsilon$  depending on specific environment setting.

## F.2 Additional Experiments

All experiments were done with 25 repetitions. We consider the  $f1$ -scores on three environments considered below.

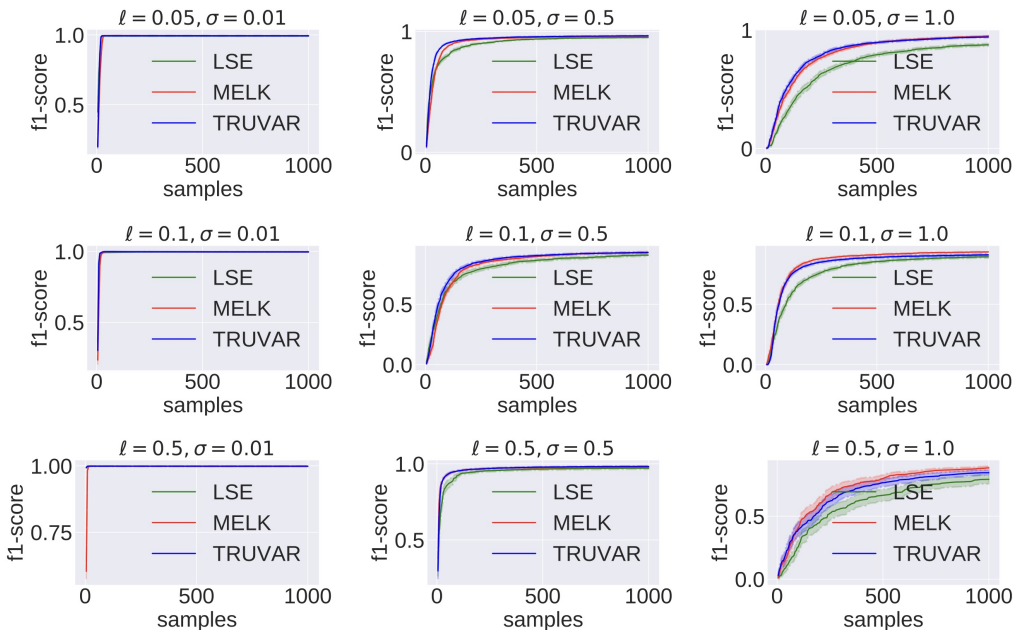


Figure 3:  $f$  drawn randomly from a squared exponential kernel  $N(0, k(\mathbf{x}, \mathbf{x}'))$ .  $\sigma$  denotes the standard deviation of the noise and  $\ell$  denotes the bandwidth of the kernel (i.e.,  $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|/2\ell^2)$ ).

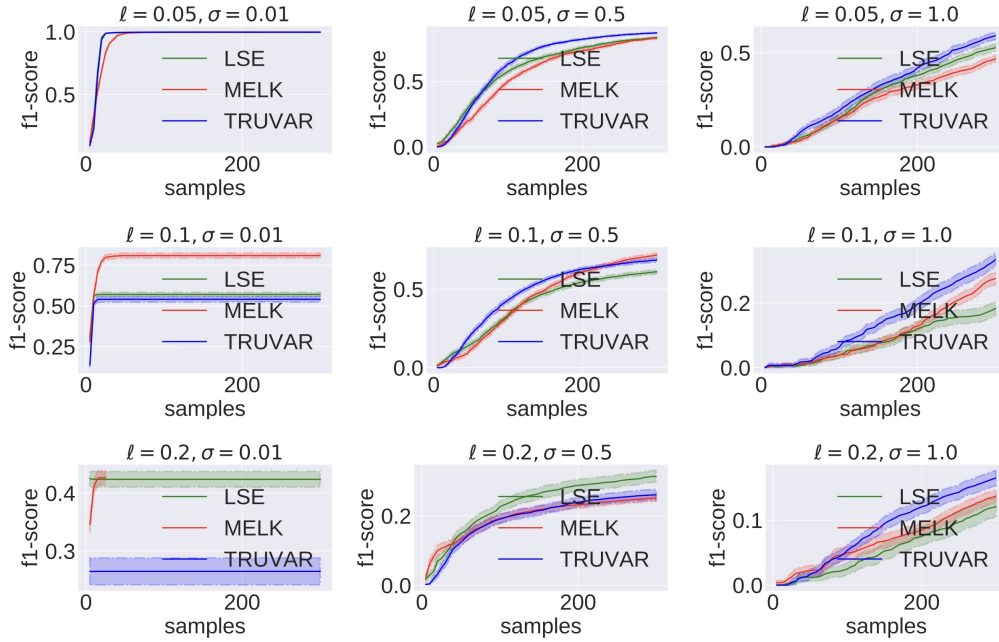


Figure 4:  $f(x) = \cos(8\pi x)$ .  $\sigma$  denotes the standard deviation of the noise and  $\ell$  denotes the bandwidth of the kernel (i.e.,  $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|/2\ell^2)$ ).

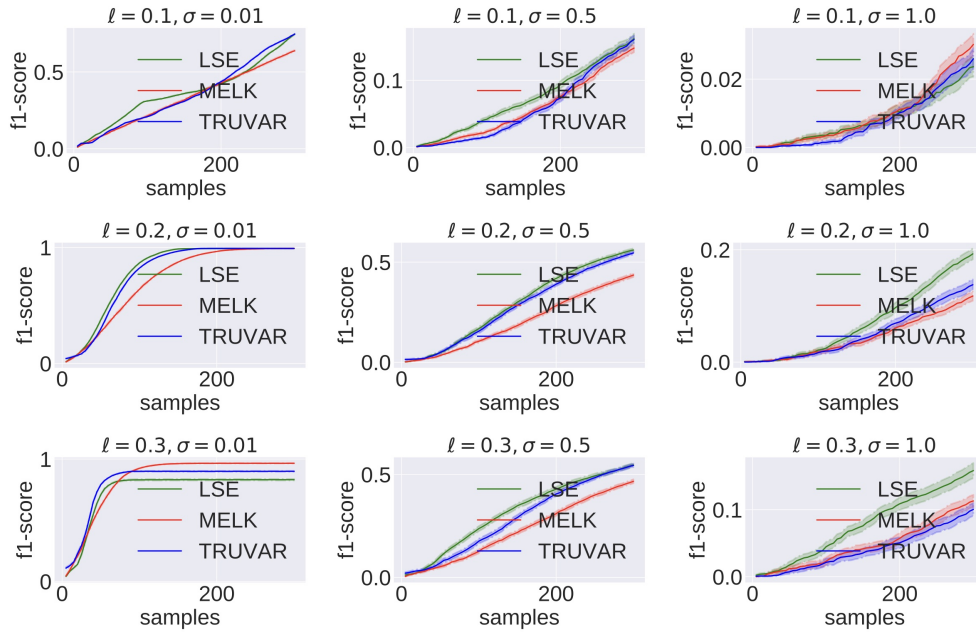


Figure 5:  $f(x, y) = \cos(2\pi x)\sin(2\pi y)$ .  $\sigma$  denotes the standard deviation of the noise and  $\ell$  denotes the bandwidth of the kernel (i.e.,  $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|/2\ell^2)$ ).

### Linear Examples with true confidence widths

Finally we compare the performance of the methods using exact confidence widths.

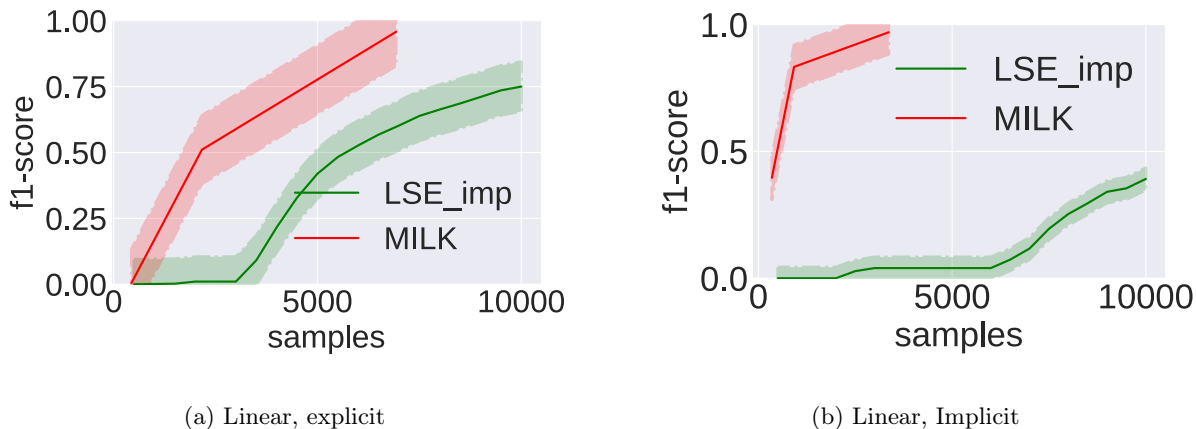


Figure 6: Comparison of algorithms using theoretically justified confidence widths on a linear bandit setting.

For the Linear kernel experiments in Figures 6a and 6b, we run all algorithms with exact confidence intervals as specified by theoretical guarantees and use the theoretical upper bound on information gain  $\gamma_T$  shown in (Srinivas et al., 2009) for the confidence widths from (Valko et al., 2013) needed for LSE. We compare the methods on a benchmark example from the linear bandits literature. For  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ , we take  $\mathbf{x}_1 = \mathbf{x}_* = \theta_* = \mathbf{e}_1$  and  $\mathbf{x}_2 = \mathbf{e}_2$ . The remaining  $\mathbf{x}_3, \dots, \mathbf{x}_n$  are set so that their first two coordinates are  $\cos(\pi/4(1 + \xi))\mathbf{e}_1$  and  $\sin(\pi/4(1 + \xi))\mathbf{e}_2$  for  $\xi \sim \text{Unif}(-.2, .2)$ . We set the threshold  $\alpha = 0.5$ ,  $n = 100$ , and  $d = 25$ . Figure 6a shows that MELK outperforms LSE when both algorithms are run with their exact confidence widths.

In the implicit setting, this example is especially informative and highlights the importance of designing to choose which arms to sample. Though it is far below  $\alpha$ , sampling arm  $\mathbf{x}_2$  provides the most information about which arms exceed the implicit threshold. Indeed, we see in 6b that both MILK greatly outperforms LSE-imp respectively.

## G Reducing Experimental Design in an RKHS to a finite dimensional optimization

In this section we describe the use of the kernel trick and Frank-Wolfe to compute the design

$$f(\lambda) = \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in C} \|\phi(\mathbf{x})\|_{A^{\gamma(\lambda)}-1}$$

where  $C \subset \mathcal{X}$ .

Since this is a convex optimization problem on the finite dimensional simplex  $\Delta_{\mathcal{X}}$  we employ the Frank-Wolfe algorithm. Note that  $\lambda_t$  is at most  $t$ -sparse. The primary challenge is in the computation of the gradient of  $f$ .

---

**Algorithm 4** Frank-Wolfe to minimize  $f$

---

**Require:** Arms  $\mathcal{X}$ , iterations  $T$

- 1:  $\lambda_0 = \mathbf{e}_1$  (first standard basis vector)
  - 2: **for**  $\mathbf{x} \in \mathcal{A}_t$  **do**
  - 3:    $x_t \leftarrow \arg \max_{\mathbf{x} \in \mathcal{X}} \|\phi(\mathbf{x})\|_{A^{\gamma(\lambda)}-1}^2$
  - 4:    $g_t = \nabla_{\lambda_{t-1}} \|\phi(\mathbf{x}_t)\|_{A^{\gamma(\lambda)}-1}^2$
  - 5:    $j_t = \arg \max_{1 \leq j \leq |\mathcal{X}|} e_j^\top g_t$
  - 6:    $\eta_t = \frac{1}{t+2}$
  - 7:    $\lambda_t = (1 - \eta_t)\lambda_{t-1} + \eta_t$
- return**  $\lambda_T$
- 

To do so we leverage a small modification of Lemma 1 of (Camilleri et al., 2021).

**Lemma G.1.** Assume that  $\lambda$  is  $s$ -sparse and (without loss of generality) with its support corresponding to

$x_1, \dots, x_s \in \mathcal{X}$ . Then,

$$\phi(\mathbf{x})^\top A^\gamma(\lambda)^{-1} \phi(\mathbf{y}) = \frac{k(\mathbf{x}, \mathbf{y})}{\gamma} - \frac{1}{\gamma} k_\lambda(\mathbf{x})^\top (K_\lambda + \gamma I_s)^{-1} k_\lambda(\mathbf{y})$$

where  $k_\lambda(\cdot) \in \mathbb{R}^s$  with  $[k_\lambda(\mathbf{x})]_i = \sqrt{\lambda_i} k(\mathbf{x}_i, \mathbf{x})$  for  $i \leq s$  and  $K_\lambda \in \mathbb{R}^{s \times s}$  with  $[K_\lambda]_{i,j} = \sqrt{\lambda_i \lambda_j} k(\mathbf{x}_i, \mathbf{x}_j)$ .

Now, identifying  $\mathcal{X}$  with an indexing of its entries, i.e.  $\mathcal{X} = \{\mathbf{x}^1, \dots, \mathbf{x}^{|\mathcal{X}|}\}$  a computation shows that

$$\mathbf{e}_i^\top [g_t] = -(\phi(\mathbf{x}_t) A^\gamma(\lambda_t)^{-1} \phi(\mathbf{x}^i))^2$$

which can be computed by the above lemma. Note that computationally, the most difficult step is the inversion of a  $t \times t$  matrix at iteration  $t$ . For a small number of iterations ( $< 2000$ ), this is not prohibitive.