

Meta-Learning Hypothesis Spaces

Parnian Kassraie

*ETH Zurich
Switzerland*

PKASSRAIE@ETHZ.CH

Jonas Rothfuss

*ETH Zurich
Switzerland*

JONAS.ROTHFUSS@INF.ETHZ.CH

Andreas Krause

*ETH Zurich
Switzerland*

KRAUSEA@ETHZ.CH

Abstract

Obtaining reliable, adaptive confidence sets for prediction functions (hypotheses) is a central challenge in sequential decision-making tasks, such as bandits and model-based reinforcement learning. These confidence sets typically rely on prior assumptions on the hypothesis space, e.g., the *known* kernel of a Reproducing Kernel Hilbert Space (RKHS). Hand-designing such kernels is error prone, and misspecification may lead to poor or unsafe performance. In this work, we propose to *meta-learn a kernel from offline data* (META-KEL). For the case where the unknown kernel is a combination of known base kernels, we develop an estimator based on structured sparsity. Under mild conditions, we guarantee that our estimated RKHS yields valid confidence sets that, with increasing amounts of offline data, become *as tight as those given the true unknown kernel*. We demonstrate our approach on the kernelized bandit problem (a.k.a. Bayesian optimization), where we establish regret bounds competitive with those given the true kernel.

Keywords: Meta-Learning, Confidence Bounds, Sequential Decision-making

1. Introduction

A number of well-studied machine learning problems such as bandits, Bayesian optimization (BO) and model-based reinforcement learning are characterized by an agent that sequentially interacts with an unknown, responsive system. Throughout the interaction, the agent’s goal is to maximize the reward based on an unknown underlying function f . Common to such sequential decision-making problems is an exploration-exploitation trade-off. That is, the agent needs to optimize its reward while learning more about the unknown function f . Confidence sets capture and quantify the uncertainty of the learner about f . Thus, they are an integral tool for directing exploration towards areas of high uncertainty and balancing it against exploitation. Moreover, in safety-critical applications, confidence sets are used to reason about the safety of actions. Thus, they are central to efficiency and safety of exploration. In theoretical analysis of sequential decision-making algorithms, a common assumption is that f resides in an RKHS with a *known* kernel function. This assumption allows for the construction of the confidence sets. In practice, however, the true kernel is unknown and needs to be hand-crafted based on the problem instance.

This is a delicate task, since the hand-crafted hypothesis space has to contain the unknown target function f . If this is not the case, the learner may be over-confident and converge to a sub-optimal policy. At the same time, we want the chosen hypothesis space to be as small so that the variance of the associated learner is low and the agent converges quickly. This constitutes a dilemma, where we need to trade off efficiency with a potential loss in consistency.

We approach this dilemma in a data-driven manner. Many applications of sequential decision-making, such as hyper-parameter tuning with BO or online nonlinear control, are of repetitive nature. Often, there is available data from similar but not identical tasks which have been solved before. Therefore, we propose to *meta-learn the kernel function, and thus the RKHS, from offline meta-data*. Our method, *Meta-Kernel Learning (META-KEL)*, works with a generic (i.e., not necessarily i.i.d.) data model and may be applied to a variety of sequential decision-making tasks.

We formally analyze the problem when the true kernel is a combination of known base kernels. We prove that the solution to META-KEL corresponds to an RKHS which contains the true function space (Theorem 4). Further, the meta-learned kernel has a sparse structure (Theorem 6) which reduces the variance of the resulting learner, and makes the learner more efficient for solving the downstream sequential decision-making problem. With mild assumptions on the data, we show that the meta-learned kernel yields confidence bands matching the ones given oracle knowledge of the true kernel, as more samples of the meta-data are provided (Theorem 5). These provably reliable confidence estimates constitute the key contribution of our work and distinguishes META-KEL from prior attempts, a summary of which is given in Appendix A.1.

To demonstrate how META-KEL can be applied to a sequential task, in Appendix A.3 we analyze a Bayesian optimization algorithm when it uses the meta-learned kernel and compare it to the same algorithm when it has knowledge of the true kernel, i.e., the oracle algorithm. By increasing size of the meta-learning data, the regret bound we obtain approaches the rate of the oracle (Theorem 7).

2. Problem Statement

Consider a sequential decision-making problem, where the agent repeatedly interacts with an environment and makes observations $\mathbf{y}_t = f(\mathbf{x}_t) + \varepsilon_t$ of an unknown function $f : \mathcal{X} \rightarrow \mathbb{R}$ residing in an RKHS \mathcal{H}_{k^*} that corresponds to an *unknown* kernel function k^* .¹ We further assume that the function has a bounded kernel norm $\|f\|_{k^*} \leq B$ and that the domain $\mathcal{X} \subset \mathbb{R}^{d_0}$ is compact. The observation noise ε_t are i.i.d. samples from a zero-mean sub-Gaussian distribution with variance proxy σ^2 . At every step t , the chosen input \mathbf{x}_t only depends on the history up to step t , denoted by the random sequence $H_{t-1} = \{(\mathbf{x}_\tau, y_\tau) : 1 \leq \tau \leq t-1\}$.

1. Appendix B.1 presents a compact refresher on the RKHS.

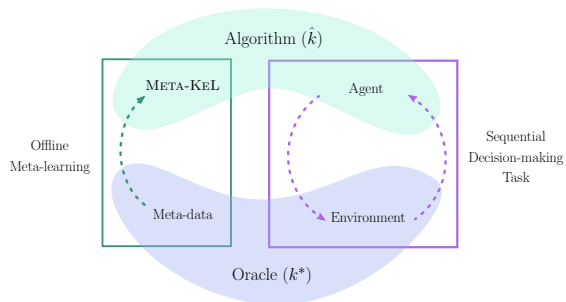


Figure 1: Overview of the described framework with k^* as the true kernel function and \hat{k} as the solution to META-KEL.

Depending on the application, \mathbf{y}_t can serve different purposes: For instance, it can describe the stochastic reward model of a bandit problem, or the transition dynamics of an RL agent.

For solving such problems, a central prerequisite for numerous algorithms are *confidence sets* for $f(\mathbf{x})$ based on the history H_{t-1} to balance exploration and exploitation at any step t . For any $\mathbf{x} \in \mathcal{X}$, the set $\mathcal{C}_{t-1}(\mathbf{x})$ defines an interval to which $f(\mathbf{x})$ belongs with high probability such that $\mathbb{P}(\forall \mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \in \mathcal{C}_{t-1}(\mathbf{x})) \geq 1 - \delta$. The center of this interval reflects the agent’s current knowledge, relevant for exploitation, and the width corresponds to the uncertainty, guiding further exploration. When the true kernel is known, an approach commonly used in the kernelized bandit literature

$$\mathcal{C}_{t-1}(k; \mathbf{x}) = [\mu_{t-1}(k; \mathbf{x}) - \nu_t \sigma_{t-1}(k; \mathbf{x}), \mu_{t-1}(k; \mathbf{x}) + \nu_t \sigma_{t-1}(k; \mathbf{x})] \quad (1)$$

where the exploration coefficient ν_t depends on the desired confidence level $1 - \delta$, and may be set based on the objective of the decision-making task. The functions μ_{t-1} and σ_{t-1} set the center and width of the confidence set as

$$\begin{aligned} \mu_{t-1}(k; \mathbf{x}) &= \mathbf{k}_{t-1}^T(\mathbf{x})(\mathbf{K}_{t-1} + \bar{\sigma}^2 \mathbf{I})^{-1} \mathbf{y}_{t-1} \\ \sigma_{t-1}^2(k; \mathbf{x}) &= k(\mathbf{x}, \mathbf{x}) - \mathbf{k}_{t-1}^T(\mathbf{x})(\mathbf{K}_{t-1} + \bar{\sigma}^2 \mathbf{I})^{-1} \mathbf{k}_{t-1}(\mathbf{x}) \end{aligned} \quad (2)$$

where $\bar{\sigma}$ is a constant, $\mathbf{y}_{t-1} = [y_\tau]_{\tau < t}$ is the vector of observed values, $\mathbf{k}_{t-1}(\mathbf{x}) = [k(\mathbf{x}, \mathbf{x}_\tau)]_{\tau < t}$, and $\mathbf{K}_{t-1} = [k(\mathbf{x}_i, \mathbf{x}_j)]_{i, j < t}$ is the kernel matrix. Hence working with the right kernel function plays an integral role in constructing well-specified sets. Since, in practice, the true kernel k^* is not known by the learner, most approaches use a hand-designed kernel that suits the problem instance at hand or conservatively pick an expressive kernel that constructs a rich RKHS which is very likely to contain f .

Addressing the issue of selecting a correct and yet efficient kernel, we pursue a data-driven approach and *meta-learn a kernel that provably yields valid confidence intervals*. This guarantee is valid regardless of how the meta-data is gathered, as long as it satisfies some basic conditions discussed later in Assumptions 1 and 3. We consider an offline collection of datasets $\mathcal{D}_{n,m} = \{(\mathbf{x}_{s,i}, y_{s,i})_{i \leq n}\}_{s \leq m}$ from m possibly non-i.i.d. tasks, each with a sample size n . Suppose, for each task s , labels are generated by $y_{s,i} = f_s(\mathbf{x}_{s,i}) + \varepsilon_{s,i}$ for $i \leq n$, where $\varepsilon_{s,i}$ are zero-mean i.i.d. sub-Gaussian noise with variance proxy σ^2 . We assume the tasks are related by the fact that all $f_s : \mathcal{X} \rightarrow \mathbb{R}$ come from the same function class \mathcal{H}_{k^*} and have a bounded RKHS norm $\|f_s\|_{k^*} \leq B$. We do not make any assumptions on the policy based on which the points $\mathbf{x}_{s,i}$ are chosen.

Assumptions Our analysis requires some assumptions on the kernel function. In particular, we assume that k^* is a finite combination of known base kernels,

$$k^*(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^p \eta_j^* k_j(\mathbf{x}, \mathbf{x}'), \quad (3)$$

where the weight vector $\boldsymbol{\eta}^* \geq 0$ is *unknown*. Without loss of generality, we assume that k^* and the base kernels are all *normalized*, i.e., $\|\boldsymbol{\eta}^*\|_1 \leq 1$ and $k_j(\mathbf{x}, \mathbf{x}') \leq 1$ for all $1 \leq j \leq p$ and $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$. The weight vector $\boldsymbol{\eta}^*$ is potentially *sparse*, since not all the candidate kernels k_j actively contribute to the construction of k^* . We use $J_{k^*} = \{1 \leq j \leq p : \eta_j^* \neq 0\}$ to refer to the group of base kernels that are present in k^* . The sparse construction of k^* imposes favorable structure on the meta-data, which essentially allows us to meta-model-select the hypothesis space and recover the true sparsity pattern denoted by J_{k^*} . We further

assume that each k_j has a d_j -dimensional feature map, i.e., $k_j(\mathbf{x}, \mathbf{x}') = \phi_j^T(\mathbf{x})\phi_j(\mathbf{x}')$, where $\phi_j \in \mathbb{R}^{d_j}$. For the scope of this paper, we assume that $d_{\max} < \infty$, where $d_{\max} := \max_{j \leq p} d_j$.

Let $\phi(\mathbf{x})$ denote the d -dimensional feature map for k^* where $d = \sum_{j=1}^p d_j$ and

$$\phi(\mathbf{x}) = \left(\sqrt{\eta_1^*} \phi_1^T(\mathbf{x}), \dots, \sqrt{\eta_p^*} \phi_p^T(\mathbf{x}) \right)^T.$$

For each task s , the $f_s \in \mathcal{H}_{k^*}$, and, by the Mercer's theorem may be decomposed as

$$f_s(\mathbf{x}) = \phi^T(\mathbf{x})\beta_s^* = \sum_{j=1}^p \sqrt{\eta_j^*} \phi_j^T(\mathbf{x})\beta_s^{*(j)}, \quad (4)$$

where $\beta_s^* \in \mathbb{R}^d$ is the coefficients vector of task s and $\beta_s^{*(j)} \in \mathbb{R}^{d_j}$ is the sub-vector corresponding to kernel k_j . It is not possible to meta-select a base kernel k_j which has not contributed to the generation of the meta-data. Therefore, if a base kernel is active in the construction of \mathcal{H}_{k^*} , it is only natural to assume that there is some task in the meta-data which reflects this presence. More formally, we assume that, for any $j \in J_{k^*}$, there exists some $s \leq m$ where $\beta_s^{*(j)} \neq 0$. We define $\beta^* = (\beta_1^{*T}, \dots, \beta_m^{*T})^T \in \mathbb{R}^{md}$ as the concatenated coefficients vector for all tasks. To refer to the group of coefficients that correspond to kernel k_j across all tasks, we use $\beta^{*(j)} = ((\beta_1^{*(j)})^T, \dots, (\beta_m^{*(j)})^T)^T \in \mathbb{R}^{md_j}$. Our next assumption guarantees that the meta-learning problem is not ill-posed.

Assumption 1 (Group Beta-min Condition) *There exists $c_1 > 0$ s.t. for all $j \in J_{k^*}$ it holds that $\|\beta^{*(j)}\|_2 \geq c_1$.*

3. Meta-learning the Hypothesis Space (Meta-KeL)

In the following section, we present our formulation of the meta-learning problem and analyze the properties of the learned hypothesis space. We meta-learn the kernel by solving the following optimization problem. Then, we set the hypothesis space of the downstream learning algorithm to be the RKHS of the meta-learned kernel.

$$\begin{aligned} \min_{\boldsymbol{\eta}, f_1, \dots, f_m} \quad & \frac{1}{m} \sum_{s=1}^m \left[\frac{1}{n} \sum_{i=1}^n (y_{s,i} - f_s(\mathbf{x}_{s,i}))^2 \right] + \frac{\lambda}{2} \sum_{s=1}^m \|f_s\|_k^2 + \frac{\lambda}{2} \|\boldsymbol{\eta}\|_1 \\ \text{s.t.} \quad & \forall s : f_s \in \mathcal{H}_k, k = \sum_{j=1}^p \eta_j k_j, 0 \leq \boldsymbol{\eta} \end{aligned} \quad (5)$$

We will refer to this problem as *Meta-Kernel Learning* (META-KeL). The first part of the objective is similar to the kernel ridge regression loss, and accounts for how well a series of regularized f_s fit the meta-data. The last term regularizes our choice of the kernel function. We use ℓ_1 -norm regularization for $\boldsymbol{\eta}$ to implicitly perform meta-model-selection. As shown in Theorem 6, the meta-learned kernel will reflect the sparsity pattern of the true kernel. The optimization problem (5) is convex and admits an efficient solution, as explained next.

We first introduce a vectorized formulation of Equation (5). Let $\mathbf{y}_s \in \mathbb{R}^n$ denote the observed values for a task s and $\mathbf{y} = (\mathbf{y}_1^T, \dots, \mathbf{y}_m^T)^T \in \mathbb{R}^{mn}$ the multi-task stacked vector of observations. We then design a multi-task feature matrix. We define Φ to be a $mn \times md$ block-diagonal matrix, where each block s corresponds to $\Phi_s = (\phi(\mathbf{x}_{s,1}), \dots, \phi(\mathbf{x}_{s,n}))^T$,

the $n \times d$ feature matrix of task s . Figure 3 provides an illustration thereof. As shown in Proposition 2, this vectorized design brings forth a parametric equivalent of META-KEl, which happens to be the well-known Group Lasso problem.

Proposition 2 (Solution of Meta-KEl) *Let $k = \sum_j \hat{\eta}_j k_j$ be a solution to Problem (5). Then, for all $1 \leq j \leq p$, it holds that $\hat{\eta}_j = \left\| \hat{\boldsymbol{\beta}}^{(j)} \right\|_2$ with $\hat{\boldsymbol{\beta}} = (\hat{\boldsymbol{\beta}}^{(j)})_{j \leq p}$ as the solution of the following convex optimization problem:*

$$\min_{\boldsymbol{\beta}} \frac{1}{mn} \|\mathbf{y} - \Phi \boldsymbol{\beta}\|_2^2 + \lambda \sum_{j=1}^p \left\| \boldsymbol{\beta}^{(j)} \right\|_2. \quad (6)$$

We show this equivalence by eliminating $\boldsymbol{\eta}$. We use a trick introduced by Bach et al. (2004), which, for $w, v \in \mathbb{R}$ states $2|w| = \min_{v>0} w^2/v + v$. The proof is given in Appendix B.2. Problem (6) can be optimized by any Group Lasso solver. Bach et al. (2011) present a number of coordinate descent algorithms which efficiently find the solution.

Note that Reproducing Kernel Hilbert Spaces are equivalent up to scaling of the kernel function. For $c > 0$, both \mathcal{H}_k and the scaled version \mathcal{H}_{ck} contain the same set of functions. Going from \mathcal{H}_k to \mathcal{H}_{ck} , the RKHS norm of any member f would scale by $1/c$, i.e. $\|f\|_k = c\|f\|_{ck}$. Hence, the norm $\|\hat{\boldsymbol{\eta}}\|_1$ will be irrelevant when meta-learning the function space. This norm can be scaled and still yield the same hypothesis space, only with a scaled operator norm. For consistency of notation, we define \hat{k} as follows. For any two points $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$, set

$$\hat{k}(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^p \frac{\hat{\eta}_j}{c_1} \phi_j^T(\mathbf{x}) \phi_j(\mathbf{x}'), \quad (7)$$

where c_1 is the same constant as in Theorem 1. We denote the set of base kernels active in \hat{k} with $J_{\hat{k}} = \{1 \leq j \leq p : \hat{\eta}_j \neq 0\}$. The meta-learned hypothesis space will then be $\mathcal{H}_{\hat{k}}$.

Properties of the Meta-learned Hypothesis Space The meta-learned \hat{k} can be used as the kernel for a model-based sequential decision making algorithm (see Figure 1). Our goal is to analyze how this choice of kernel affects the success of the algorithm, compared to the oracle algorithm with access to the unknown kernel. To this end, we discuss properties of \hat{k} . For our main result, we require a final technical assumption to ensure that our meta-data is sufficiently diverse. Consider a vector $\mathbf{b} \in \mathbb{R}^{md}$ that adheres to the same group structure as $\boldsymbol{\beta}^*$. For a set of group indices $J \subset \{1, \dots, p\}$, we define $\mathbf{b}_J := (\mathbf{b}^{(j)})_{j \in J}$ as the sub-vector indicated by the groups in J . The following condition ensures that the meta-training data is not degenerate, e.g., no two points $\mathbf{x}_{s,i}$ and $\mathbf{x}_{s,j}$ in the meta-data are identical or too close.

Assumption 3 (Relaxed Sufficient Exploration) *There exists $\kappa = \kappa(s) > 0$ such that*

$$\kappa \leq \min_{J, \mathbf{b}} \frac{\|\Phi \mathbf{b}\|_2}{\sqrt{mn} \|\mathbf{b}_J\|_2}, \quad \text{s.t.} \quad \sum_{j \notin J} \left\| \mathbf{b}^{(j)} \right\|_2 \leq 3 \sum_{j \in J} \left\| \mathbf{b}^{(j)} \right\|_2, \quad \mathbf{b} \neq 0, |J| \leq s.$$

Theorem 4 (Hypothesis Space Recovery) *Set $0 \leq \delta \leq 1$, and choose λ such that,*

$$\lambda \geq \frac{4\sigma}{\sqrt{mn}} \sqrt{1 + \frac{2}{m} \left(\log(2p/\delta) + \sqrt{md_{\max} \log(2p/\delta)} \right)}.$$

If $|J_{k^*}| \leq s$ and Assumption 3 holds with $\kappa(s)$, then $\mathcal{H}_{k^*} \subseteq \mathcal{H}_{\hat{k}}$ and $\|f\|_{\hat{k}} \leq \|f\|_{k^*} \left(1 + \epsilon(n, m) + o(\epsilon(n, m))\right)$, with probability greater than $1 - \delta$, if n and m are large enough to satisfy $\epsilon(n, m) \leq c_1$. The absolute constant c_1 is defined in Theorem 1 and

$$\epsilon(n, m) := \frac{32\sigma s}{\kappa^2(s)\sqrt{mn}} \sqrt{1 + \frac{2}{m} \left(\log(2p/\delta) + \sqrt{md_{\max} \log(2p/\delta)}\right)}.$$

The proof is given in Appendix C.1. Theorem 4 states that, provided enough meta-data, \mathcal{H}_{k^*} is contained in $\mathcal{H}_{\hat{k}}$ with high probability and δ , the probability of failure in recovery, decreases as m and n grow (See Appendix C.3). The theorem implies that the meta-learner benefits more from increasing the number m of meta-data tasks, rather than increasing the sample size n of each task, since $\epsilon(n, m)$ shrinks faster with m compared to n . Note that increasing either m or n will result in convergence and therefore this theorem also holds for the classic offline kernel learning setup when the dataset consists of a single learning task ($m = 1$). Lastly, we can show that provided enough meta-data, if k^* is sparse, then with high probability \hat{k} will also be sparse, i.e. $J_{k^*} = J_{\hat{k}}$. Appendix A.2 formalizes this claim.

4. Sequential Decision-making with Meta-KeL

We now analyze the effect of using \hat{k} as kernel function in the downstream sequential decision-making problem. We adopt the common construction of confidence sets given in Equation (1), and define $\hat{C}_{t-1}(\mathbf{x}) := \mathcal{C}_{t-1}(\hat{k}; \mathbf{x})$. We let $\hat{\mu}_{t-1}(\mathbf{x}) := \mu_{t-1}(\hat{k}; \mathbf{x})$, and $\hat{\sigma}_{t-1}(\mathbf{x}) := \sigma_{t-1}(\hat{k}; \mathbf{x})$, where $\mu_{t-1}(k; \mathbf{x})$ and $\sigma_{t-1}(k; \mathbf{x})$ are as defined in Equation (2) with $\bar{\sigma} = 1 + 2/T$. Theorem 5 shows that for the right choice of ν_t , the set $\hat{C}_{t-1}(\mathbf{x})$ is a valid confidence bound for any $f \in \mathcal{H}_{k^*}$, evaluated at any $\mathbf{x} \in \mathcal{X}$, at any step t , with high probability.

Theorem 5 (Confidence Bounds with Meta-KeL) *Let $f \in \mathcal{H}_{k^*}$ with $\|f\|_{k^*} \leq B$, where k^* is unknown. Under the assumptions of Theorem 4, with probability greater than $1 - \delta$, for all $\mathbf{x} \in \mathcal{X}$ and $1 \leq t \leq T$,*

$$|\hat{\mu}_{t-1}(\mathbf{x}) - f(\mathbf{x})| \leq \nu_t \hat{\sigma}_{t-1}(\mathbf{x}) \left(B \left(1 + \frac{\epsilon(n, m)}{2c_1} \right) + \sigma \sqrt{\hat{d} \log \left(1 + \frac{\bar{\sigma}^{-2}t}{c_1} \right) + 2 + 2 \log(1/\delta)} \right)$$

where $\hat{d} = \sum_{j \in J_{\hat{k}}} d_j$.

The proof is given in Appendix D. As discussed in Section 3, the $\epsilon(n, m)/2c_1$ term shrinks faster than $\mathcal{O}(1/\sqrt{mn})$ and \hat{d} approaches $d^* = \sum_{j \in J_{k^*}} d_j$ at a similar rate. Therefore, Theorem 5 presents a tight confidence bound relative to the case when k^* is known by the agent. Due to Theorem 2 in Chowdhury and Gopalan (2017), the $1 - \delta$ confidence bound would be

$$|\mu_{t-1}(\mathbf{x}) - f(\mathbf{x})| \leq \sigma_{t-1}(\mathbf{x}) \left(B + \sigma \sqrt{d^* \log \left(1 + \bar{\sigma}^{-2}t \right) + 2 + 2 \log(1/\delta)} \right).$$

where the mean and variance functions are defined by $\mu_{t-1}(\mathbf{x}) := \mu_{t-1}(k^*; \mathbf{x})$ and $\sigma_{t-1}(\mathbf{x}) := \sigma_{t-1}(k^*; \mathbf{x})$ with $\bar{\sigma} = 1 + 2/T$. We conclude that the base learner does not require knowledge of the true kernel for constructing confidence sets, as long as there is sufficient meta-data available. Theorem 4 quantifies this notion of sufficiency.

References

- Francis Bach, Rodolphe Jenatton, Julien Mairal, and Guillaume Obozinski. Optimization with sparsity-inducing penalties. *arXiv preprint arXiv:1108.0775*, 2011.
- Francis R Bach. Consistency of the group lasso and multiple kernel learning. *Journal of Machine Learning Research*, 9(6), 2008.
- Francis R Bach, Gert RG Lanckriet, and Michael I Jordan. Multiple kernel learning, conic duality, and the smo algorithm. In *Proceedings of the twenty-first international conference on Machine learning*, page 6, 2004.
- Hamsa Bastani and Mohsen Bayati. Online decision making with high-dimensional covariates. *Operations Research*, 68(1):276–294, 2020.
- Soumya Basu, Branislav Kveton, Manzil Zaheer, and Csaba Szepesvári. No regrets for learning the prior in bandits. *Advances in Neural Information Processing Systems*, 34, 2021.
- Felix Berkenkamp, Matteo Turchetta, Angela P Schoellig, and Andreas Krause. Safe model-based reinforcement learning with stability guarantees. *arXiv preprint arXiv:1705.08551*, 2017.
- Felix Berkenkamp, Angela P Schoellig, and Andreas Krause. No-regret bayesian optimization with unknown hyperparameters. *arXiv preprint arXiv:1901.03357*, 2019.
- Ilija Bogunovic and Andreas Krause. Misspecified Gaussian process bandit optimization. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- Craig Boutilier, Chih-wei Hsu, Branislav Kveton, Martin Mladenov, Csaba Szepesvari, and Manzil Zaheer. Differentiable meta-learning of bandit policies. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2020.
- Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge University Press, 2004.
- Peter Bühlmann and Sara Van De Geer. *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media, 2011.
- Florentina Bunea, Johannes Lederer, and Yiyuan She. The group square-root lasso: Theoretical properties and fast algorithms. *IEEE Transactions on Information Theory*, 60(2): 1313–1325, 2013.
- Laurent Cavalier, GK Golubev, Dominique Picard, and AB Tsybakov. Oracle inequalities for inverse problems. *The Annals of Statistics*, 30(3):843–874, 2002.
- Leonardo Cella and Massimiliano Pontil. Multi-task and meta-learning with sparse linear bandits. In *Uncertainty in Artificial Intelligence*, pages 1692–1702. PMLR, 2021.

- Leonardo Cella, Alessandro Lazaric, and Massimiliano Pontil. Meta-learning with stochastic linear bandits. In *International Conference on Machine Learning*, pages 1360–1370. PMLR, 2020.
- Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *International Conference on Machine Learning*, pages 844–853. PMLR, 2017.
- Nello Cristianini, Jaz Kandola, Andre Elisseeff, and John Shawe-Taylor. On kernel target alignment. In *Innovations in machine learning*, pages 205–256. Springer, 2006.
- Sebastian Curi, Felix Berkenkamp, and Andreas Krause. Efficient model-based reinforcement learning through optimistic policy search and planning. *arXiv preprint arXiv:2006.08684*, 2020.
- Theodoros Evgeniou and Massimiliano Pontil. Regularized multi-task learning. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 109–117, 2004.
- Tilmann Gneiting, Fadoua Balabdaoui, and Adrian E Raftery. Probabilistic forecasts, calibration and sharpness. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(2):243–268, 2007.
- Mehmet Gönen and Ethem Alpaydm. Multiple kernel learning algorithms. *The Journal of Machine Learning Research*, 12:2211–2268, 2011.
- Botao Hao, Tor Lattimore, and Mengdi Wang. High-dimensional sparse linear bandits. *arXiv preprint arXiv:2011.04020*, 2020.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *arXiv preprint arXiv:1806.07572*, 2018.
- Sham Kakade, Akshay Krishnamurthy, Kendall Lowrey, Motoya Ohnishi, and Wen Sun. Information theoretic regret bounds for online nonlinear control. *arXiv preprint arXiv:2006.12466*, 2020.
- Marius Kloft, Ulf Brefeld, Sören Sonnenburg, and Alexander Zien. Lp-norm multiple kernel learning. *The Journal of Machine Learning Research*, 12:953–997, 2011.
- Vladimir Koltchinskii and Ming Yuan. Sparse recovery in large ensembles of kernel machines. In *Proceedings of COLT*, volume 69, 2008.
- Branislav Kveton, Martin Mladenov, Chih-Wei Hsu, Manzil Zaheer, Csaba Szepesvari, and Craig Boutilier. Meta-learning bandit policies by gradient ascent. *arXiv preprint arXiv:2006.05094*, 2020.
- Han Liu and Jian Zhang. Estimation consistency of the group lasso and its applications. In *Artificial Intelligence and Statistics*, pages 376–383. PMLR, 2009.
- Karim Lounici, Massimiliano Pontil, Sara Van De Geer, and Alexandre B Tsybakov. Oracle inequalities and optimal inference under group sparsity. *The annals of statistics*, 39(4): 2164–2204, 2011.

- Mathurin Massias, Alexandre Gramfort, and Joseph Salmon. Celer: a fast solver for the lasso with dual extrapolation. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pages 3321–3330, 2018.
- Cheng Soon Ong, Alexander J. Smola, and Robert C. Williamson. Learning the kernel with hyperkernels. *Journal of Machine Learning Research*, 2005.
- Jonas Rothfuss, Vincent Fortuin, Martin Josifoski, and Andreas Krause. Pacoh: Bayes-optimal meta-learning with pac-guarantees. In *International Conference on Machine Learning*, pages 9116–9126. PMLR, 2021a.
- Jonas Rothfuss, Dominique Heyn, Jinfan Chen, and Andreas Krause. Meta-learning reliable priors in the function space. In *Advances in Neural Information Processing Systems*, 2021b.
- Pier Giuseppe Sessa, Ilija Bogunovic, Maryam Kamgarpour, and Andreas Krause. Learning to play sequential games versus unknown opponents. *arXiv preprint arXiv:2007.05271*, 2020.
- Max Simchowitz, Christopher Tosh, Akshay Krishnamurthy, Daniel J Hsu, Thodoris Lykouris, Miro Dudik, and Robert E Schapire. Bayesian decision-making under misspecified priors with applications to meta-learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- Niranjan Srinivas, Andreas Krause, Sham M Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. *arXiv preprint arXiv:0912.3995*, 2009.
- Sattar Vakili, Kia Khezeli, and Victor Picheny. On information gain and regret bounds in gaussian process bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 82–90. PMLR, 2021.
- Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- Zi Wang, Chengtao Li, Stefanie Jegelka, and Pushmeet Kohli. Batched high-dimensional bayesian optimization via structural kernel learning. In *International Conference on Machine Learning*, pages 3656–3664. PMLR, 2017.
- Zi Wang, Beomjoon Kim, and Leslie Pack Kaelbling. Regret bounds for meta bayesian optimization with an unknown gaussian process prior. *arXiv preprint arXiv:1811.09558*, 2018.
- Ziyu Wang and Nando de Freitas. Theoretical analysis of bayesian optimisation with unknown gaussian process hyper-parameters. *arXiv preprint arXiv:1406.7758*, 2014.
- George Wynne, François-Xavier Briol, and Mark Girolami. Convergence guarantees for gaussian process means with misspecified likelihoods and smoothness. *Journal of Machine Learning Research*, 22(123):1–40, 2021.

Appendix A. Details of the Main Text

A.1 Related Work

Numerous sequential decision-making methods rely on confidence sets for uncertainty quantification, e.g., UCB algorithms (Srinivas et al., 2009; Chowdhury and Gopalan, 2017) for Bayesian optimization and bandits, safe exploration and various forms of RL (Berkenkamp et al., 2017; Curi et al., 2020; Kakade et al., 2020; Sessa et al., 2020). Most of these methods assume the true hypothesis space as given. However, in practice, we typically do not know the correct hypothesis space, e.g., in form of a kernel. A body of recent work considers the unknown kernel setting and analyzes the effect of working with a misspecified hypothesis space (Wynne et al., 2021; Simchowitz et al., 2021; Bogunovic and Krause, 2021). Alternatively, Wang and de Freitas (2014) and Berkenkamp et al. (2019) propose to successively expand the hypothesis space throughout the course of BO so that the algorithm remains no-regret in a setting where the kernel lengthscale is unknown.

Our work relates to meta-learning for Bayesian optimization. There is a recent line of algorithms that improve accuracy of base sequential learners via meta-learning, albeit without theoretical guarantees (Rothfuss et al., 2021a,b), or with mild guarantees in special cases (Kveton et al., 2020; Boutilier et al., 2020). There are a number of results on updating bandit priors or policies by meta-learning, under problem settings different than ours. Basu et al. (2021) work with a sequence of multi-armed bandit tasks, and adaptively meta-learn the *mean* of the Gaussian prior used for the next task. Others consider solving a number of structurally similar linear bandit tasks in parallel (Wang et al., 2017; Cella et al., 2020; Cella and Pontil, 2021). They propose how to efficiently update the policy when each learner has access to the data across all tasks. We significantly improve upon the work of Wang et al. (2018) which analyzes the simple regret of the GP-UCB algorithm (Srinivas et al., 2009) for multi-armed and linear bandits, when the mean and variance of the Gaussian prior are unknown, and there is sufficient offline i.i.d. data drawn from the same Gaussian distribution.

Our framework considers structural sparsity at the kernel level, translating to group sparsity for the coefficients vectors, if applied to linear bandits. Thus our work relates to results on sparse linear bandits and Lasso bandits. In this area, Bastani and Bayati (2020) and Hao et al. (2020) give dimension-independent regret bounds for *Explore-Then-Commit* algorithms under certain assumptions over the action set. This work does not consider offline data.

We draw inspiration from the early Multiple Kernel Learning (MKL) literature, which focuses on kernel design for classification with SVMs (Bach et al., 2004; Gönen and Alpaydm, 2011; Kloft et al., 2011; Evgeniou and Pontil, 2004; Cristianini et al., 2006; Ong et al., 2005). In contrast, our key contribution is to derive adaptive confidence bounds from meta-learned kernels for regression, even for non-i.i.d. data. Orthogonal to most prior works, we reduce the kernel learning problem to group Lasso and leverage the properties of the Lasso estimator, in particular seminal results of Lounici et al. (2011) and Bach (2008). Other relevant works on convergence properties of the group Lasso include Koltchinskii and Yuan (2008), Liu and Zhang (2009) and Bunea et al. (2013).

A.2 The Benefit of Structural Sparsity

Consider a conservative hand-picked kernel function

$$k_{\text{full}} = 1/p \sum_{j=1}^p \phi_j^T(\mathbf{x}) \phi_j(\mathbf{x}'), \quad (\text{A.1})$$

which does not use any meta-data and instead incorporates all the considered base kernels. When p and d_{max} are finite, \mathcal{H}_{k^*} is contained in $\mathcal{H}_{k_{\text{full}}}$ and the hand-picked hypothesis space is not misspecified. However, working with an overly large hypothesis space has downsides. Consider using k_{full} to estimate a function $f \in \mathcal{H}_{k^*}$. Then every base kernel, including k_j with $j \notin J_{k^*}$, appears in the construction of the estimator. These terms contribute to the estimation error and increase the variance of the function estimate. This slows down the rate of convergence, compared to the case where only active k_j are present in the kernel function. By meta-learning \hat{k} via META-KeL, we can eliminate irrelevant candidate kernels and produce a structurally sparse hypothesis space. Theorem 6 guarantees this property. Its proof is given in Appendix C.2.

Proposition 6 (Bound on structural sparsity of \hat{k}) *Set $0 < \delta < 1$ and choose λ according to Theorem 4. Let $|J_{k^*}| \leq s$ be the number candidate kernels that contribute to k^* . If Assumption 3 holds with $\kappa(s)$, then with probability greater than $1 - \delta$, the number of kernels active in \hat{k} is bounded by*

$$|J_{\hat{k}}| \leq \frac{4s}{mn\kappa^2(s)}$$

which implies that if $mn > \frac{4s}{p\kappa^2(s)}$, then with the same probability

$$\mathcal{H}_{\hat{k}} \subsetneq \mathcal{H}_{k_{\text{full}}}.$$

Hence, in the presence of enough meta-data, $\mathcal{H}_{\hat{k}}$ is a strict subset of $\mathcal{H}_{k_{\text{full}}}$, and therefore

$$\mathcal{H}_{k^*} \stackrel{\text{w.h.p.}}{\subseteq} \mathcal{H}_{\hat{k}} \stackrel{\text{w.h.p.}}{\subsetneq} \mathcal{H}_{k_{\text{full}}}$$

where the left relation is due to Theorem 4. Figure 2 illustrates the nested sets. We conclude that our meta-learned hypothesis space has favorable properties: it contains the true hypothesis space, and it is sparse in structure, in particular, smaller than the conservative candidate space.

The fact that $\mathcal{H}_{\hat{k}}$ is smaller than $\mathcal{H}_{k_{\text{full}}}$ reduces the complexity of the downstream learning problem and yields faster convergence rates. We provide an example of this effect in Appendix A.3, where we analyze a Bayesian optimization problem, and establish how choosing \hat{k} improves upon k_{full} . Finally, our experiments (e.g. Figure 5) support the claim that in practice the BO algorithm is faster in finding the optimum when it uses the meta-learned kernel.

A.3 Meta-KeL for Bayesian Optimization

As an example application, we consider the classic Bayesian optimization problem, but in the case where \mathcal{H}_{k^*} is unknown. This example illustrates how Theorem 5 may be used

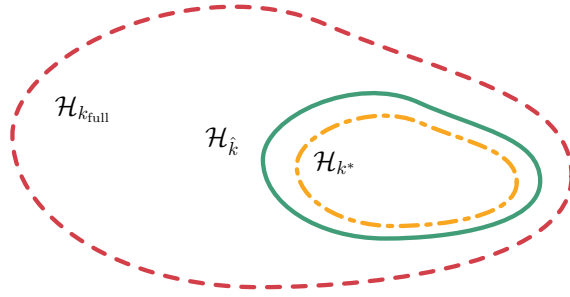


Figure 2: The oracle \mathcal{H}_{k^*} (Eq. 3), the meta-learned $\mathcal{H}_{\hat{k}}$ (Eq. 7) and the hand-picked $\mathcal{H}_{k_{\text{full}}}$ (Eq. A.1) hypothesis spaces (informal)

to prove guarantees for a decision-making algorithm, which uses the meta-learned kernel due to a lack of knowledge of k^* . We follow the setup and BO notation of Srinivas et al. (2009). The agent seeks to maximize an unknown reward function f , sequentially accessed as described in ???. Their goal is to choose actions \mathbf{x}_t which maximize the cumulative reward achieved over T time steps. This is equivalent to minimizing the cumulative regret $R_T = \sum_{t=1}^T [f(\mathbf{x}^*) - f(\mathbf{x}_t)]$, where \mathbf{x}^* is a global maximum of f . Note that if $R_T/T \rightarrow 0$ as $T \rightarrow \infty$ then $\max_{1 \leq t \leq T} f(\mathbf{x}_t) \rightarrow f(\mathbf{x}^*)$, i.e., the learner converges to the optimal value. We will refer to this property as *sublinearity* of the regret. In the spirit of the GP-UCB algorithm (Srinivas et al., 2009), we choose the next point by maximizing the upper confidence bound as determined by Theorem 5

$$\mathbf{x}_t = \arg \max_{\mathbf{x} \in \mathcal{X}} \hat{\mu}_{t-1}(\mathbf{x}) + \nu_t \hat{\sigma}_{t-1}(\mathbf{x}) \quad (\text{A.2})$$

where a suitable choice for ν_t is suggested in Theorem 7.

Corollary 7 (A Regret Bound with Meta-KeL) *Let $\delta \in (0, 1)$. Suppose $f \in \mathcal{H}_{k^*}$ with $\|f\|_{k^*} \leq B$ and that values of f are observed with zero-mean sub-Gaussian noise of variance proxy σ^2 . Then, with probability greater than $1 - \delta$, GP-UCB used together with \hat{k} satisfies*

$$R_T = \mathcal{O} \left(\sqrt{\hat{d} T \log T} \left(B(1 + \epsilon(n, m)) + \sqrt{\hat{d} \log T + \log 1/\delta} \right) \right)$$

provided that the exploration coefficient is set to

$$\nu_t = B(1 + \epsilon(n, m)/2c_1) + \sigma \sqrt{\hat{d} \log(1 + \bar{\sigma}^{-2} t/c_1) + 2 + 2 \log(1/\delta)}.$$

The proof is straightforward. Conditioned on the event that $f \in \mathcal{H}_{\hat{k}}$, we may directly use the regret bound of Chowdhury and Gopalan (2017). Then, by Theorem 4, we calculate the probability of this event (Appendix D.1). The Corollary relies on knowledge of a bound

B on $\|f\|_{k^*}$. However, using techniques of Berkenkamp et al. (2019) it is possible to adapt it even to *unknown* B at increased (but still sublinear) regret.

Theorem 7 shows that GP-UCB using the meta-learned kernel guarantees sublinear regret. We obtain a $\mathcal{O}(\hat{d}B \log T \sqrt{T})$ rate for the regret which is tight compared to the $\mathcal{O}(d^*B \log T \sqrt{T})$ rate satisfied by the oracle. It is insightful to compare this convergence result to a scenario where the hypothesis space is misspecified. For a reward function $f \notin \mathcal{H}_{\hat{k}}$, Bogunovic and Krause (2021) show that the learner will not converge to the global optimum, since the cumulative regret has a lower bound of linear order $\mathcal{O}(T \sqrt{\log T})$. Corollary 7 suggests that by using a sparse kernel we can potentially find the optimal policy faster compared to when the complex kernel k_{full} is used. Recall that $d = \sum_{j=1}^p d_j$, by Theorem 2 of Chowdhury and Gopalan (2017) the regret of GP-UCB used together with k_{full} is bounded by $\mathcal{O}(dpB \log T \sqrt{T})$, since $\|f\|_{k_{\text{full}}} = p\|f\|_{k^*}$. Therefore, using the meta-learned kernel improves the regret bound by a factor of $\hat{d}/(dp)$, implying that the solution may be found faster. The results of our experiments in Figure 5 support this argument.

Note that our approach to guarantee a sublinear regret for GP-UCB without oracle knowledge of k^* naturally generalizes to other sequential decision tasks. In particular, any theoretical result relying on RKHS confidence intervals with a *known* kernel can be immediately extended to use those of the meta-learned kernel.

Appendix B. Details of the main Result

B.1 RKHS Refresher

Here we present a compact reminder of RKHS basics for the sake of completeness and clarifying our notation. We work in a finite-dimensional regime which can also be described by a euclidean vector space. Nevertheless, we use the RKHS notation as it gives a powerful framework and hides away the vector algebra. This section is mainly based on Wainwright (2019). For a positive semi-definite kernel function k over some set $\mathcal{X} \times \mathcal{X}$, the corresponding unique Reproducing Kernel Hilbert Space can be constructed as,

$$\mathcal{H}_k = \left\{ f : \mathcal{X} \rightarrow \mathbb{R} \mid f(\cdot) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \cdot), n \in \mathbb{N}, (\mathbf{x}_i)_{i=1}^n \in \mathcal{X}, \boldsymbol{\alpha} \in \mathbb{R}^n \right\},$$

equipped with the dot product, $\langle f, \bar{f} \rangle_k = \sum_{i,j} \alpha_i \bar{\alpha}_j k(x_i, \bar{x}_j)$. We limit \mathcal{X} to compact sets, and only consider Mercer kernels, i.e., continuous kernel functions that satisfy the Hilbert-Schmidt condition,

$$\int_{\mathcal{X} \times \mathcal{X}} k^2(\mathbf{x}, \mathbf{x}') d\mu(\mathbf{x}) d\mu(\mathbf{x}') \leq \infty$$

where μ is a non-negative measure over \mathcal{X} . Mercer theorem states that under these assumptions, the kernel operator has a sequence of orthonormal eigenfunctions $(\phi_r)_{r \geq 1}$ and non-negative eigenvalues $(\eta_r)_{r \geq 1}$, defined as follows

$$\int_{\mathcal{X}} k(\mathbf{x}, \mathbf{x}') \phi_r(\mathbf{x}') d\mu(\mathbf{x}') = \eta_r \phi_r(\mathbf{x}). \quad (\text{B.1})$$

Moreover, k can be written as their linear combination,

$$k(\mathbf{x}, \mathbf{x}') = \sum_r \eta_r \phi_r(\mathbf{x}) \phi_r(\mathbf{x}').$$

$$\begin{aligned}
\begin{pmatrix} \mathbf{y}_s \\ \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_s \\ \vdots \\ \mathbf{y}_m \end{pmatrix} &= \sum_{j=1}^p \begin{pmatrix} \phi_j^T(\mathbf{x}_{s,1}) \\ \vdots \\ \phi_j^T(\mathbf{x}_{s,n}) \end{pmatrix} \begin{pmatrix} \beta_s^{(j)} \\ \vdots \\ \beta_s^{(j)} \end{pmatrix} + \begin{pmatrix} \epsilon_s \\ \vdots \\ \epsilon_s \\ \vdots \\ \epsilon_m \end{pmatrix} \\
&= \begin{pmatrix} \Phi_1 & & 0 \\ & \ddots & \\ 0 & \Phi_s & 0 \\ & & \ddots \\ & & 0 & \Phi_m \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_s \\ \vdots \\ \beta_m \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_s \\ \vdots \\ \epsilon_m \end{pmatrix} \\
\beta &= \begin{pmatrix} \beta_1^{(1)} \\ \vdots \\ \beta_1^{(j)} \\ \vdots \\ \beta_1^{(p)} \\ \vdots \\ \beta_s^{(1)} \\ \vdots \\ \beta_s^{(j)} \\ \vdots \\ \beta_s^{(p)} \\ \vdots \\ \beta_m^{(1)} \\ \vdots \\ \beta_m^{(j)} \\ \vdots \\ \beta_m^{(p)} \end{pmatrix} \quad \beta^{(j)} = \begin{pmatrix} \beta_1^{(j)} \\ \vdots \\ \beta_s^{(j)} \\ \vdots \\ \beta_m^{(j)} \end{pmatrix}
\end{aligned}$$

Figure 3: Visual guide for the Group Lasso formulation. The red shade shows the group of coefficients which correspond to the effect of kernel k_j . The green shade demonstrates how features from each task come together on the diagonal of the multi-task feature matrix. The coefficients that belong to one task are group together with a green rectangle.

Or as a non-negative combination of base kernels, $k_r(\mathbf{x}, \mathbf{x}') = \phi_r(\mathbf{x})\phi_r(\mathbf{x}')$

$$k(\mathbf{x}, \mathbf{x}') = \sum_r \eta_r k_r(\mathbf{x}, \mathbf{x}').$$

It immediately follows that the unique RKHS corresponding to k takes the form,

$$\mathcal{H}_k = \left\{ f : \mathcal{X} \rightarrow \mathbb{R} \mid f(\cdot) = \sum_{r \geq 1} \beta_r \phi_r(\cdot), \sum_{r: \eta_r \neq 0} \frac{\beta_r^2}{\eta_r} \leq \infty \right\}$$

and the inner product the following form,

$$\langle f, g \rangle_k = \sum_{r: \eta_r \neq 0} \frac{\langle f, \phi_r \rangle \langle g, \phi_r \rangle}{\eta_r},$$

where $\langle \cdot, \cdot \rangle_2$ denotes the inner product in the $L^2(\mathcal{X})$. It is then implied that $\|f\|_k = \sum_{r: \eta_r \neq 0} \beta_r^2 / \eta_r$. Lastly, we define $\phi(\mathbf{x}) = (\sqrt{\eta_r} \phi_r(\mathbf{x}))_{r \geq 1}$ to be the feature map corresponding to k . In this paper we refer to the number of non-zero eigen-values as the dimension of the kernel. Note that under this convention, most kernel functions used in practice, e.g. RBF kernel or the Matérn family, are infinite-dimensional.

B.2 Proof of Proposition 2

By Equation (4) we may write a parametric equivalent of Problem (5) in terms of the feature maps ϕ_j ,

$$\begin{aligned} \min_{\substack{0 \leq \boldsymbol{\eta}, \\ \forall s: \boldsymbol{\beta}_s}} \frac{1}{m} \sum_{s=1}^m \left[\frac{1}{n} \sum_{i=1}^n \left(y_{s,i} - \sum_{j=1}^p \sqrt{\eta_j} \phi_j^T(\mathbf{x}_{s,i}) \boldsymbol{\beta}_s^{(j)} \right)^2 \right] \\ + \frac{\lambda}{2} \sum_{s=1}^m \sum_{j=1}^p \|\boldsymbol{\beta}_s^{(j)}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\eta}\|_1. \end{aligned} \quad (\text{B.2})$$

This problem is jointly convex in $\boldsymbol{\eta}$ and $(\boldsymbol{\beta}_s)_{s \leq m}$, and has an optimal solution (Kloft et al., 2011). Renaming the variable $\boldsymbol{\beta}_s^{(j)} \leftarrow \boldsymbol{\beta}_s^{(j)} / \sqrt{\eta_j}$ gives the following equivalent problem,

$$\min_{\substack{0 \leq \boldsymbol{\eta}, \\ \forall s: \boldsymbol{\beta}_s}} \frac{1}{m} \sum_{s=1}^m \left[\frac{1}{n} \sum_{i=1}^n \left(y_{s,i} - \sum_{j=1}^p \phi_j^T(\mathbf{x}_{s,i}) \boldsymbol{\beta}_s^{(j)} \right)^2 \right] + \frac{\lambda}{2} \sum_{j=1}^p \frac{\sum_{s=1}^m \|\boldsymbol{\beta}_s^{(j)}\|_2^2}{\eta_j} + \frac{\lambda}{2} \|\boldsymbol{\eta}\|_1.$$

There is no constraint connecting the two variables, and renaming $\boldsymbol{\beta}_s^{(j)}$ does not effect the optimization problem with respect to $\boldsymbol{\eta}$. Therefore, $\hat{\boldsymbol{\eta}}$ is also an optima for this problem. Let $(\hat{\boldsymbol{\eta}}, \hat{\boldsymbol{\beta}}_s)$ denote the solution to the problem above. We show that $\hat{\boldsymbol{\eta}}$ has a closed form expression in terms of $\hat{\boldsymbol{\beta}}$. This observation allows us to reduce the joint optimization problem to an equivalent problem which is only over $(\boldsymbol{\beta}_s)$. We use the η -trick introduced in Bach et al. (2004). The authors observe that for any two scalar variables w and v ,

$$|w| = \min_{v \geq 0} \frac{w^2}{2v} + \frac{v}{2},$$

and $\hat{v} = |w|$. Applying this trick with $w = \|\boldsymbol{\beta}^{(j)}\|_2$ and $v = \eta_j$ for all $j \leq p$ gives

$$\min_{\boldsymbol{\eta} \geq 0} \frac{\lambda}{2} \sum_{j=1}^p \frac{\|\boldsymbol{\beta}^{(j)}\|_2^2}{\eta_j} + \frac{\lambda}{2} \|\boldsymbol{\eta}\|_1 = \lambda \sum_{j=1}^p \|\boldsymbol{\beta}^{(j)}\|_2,$$

which results the following equivalent problem

$$\min_{\forall s: \boldsymbol{\beta}_s} \frac{1}{m} \sum_{s=1}^m \left[\frac{1}{n} \sum_{i=1}^n \left(y_{s,i} - \sum_{j=1}^p \phi_j^T(\mathbf{x}_{s,i}) \boldsymbol{\beta}_s^{(j)} \right)^2 \right] + \lambda \sum_{j=1}^p \sqrt{\sum_{s=1}^m \|\boldsymbol{\beta}_s^{(j)}\|_2^2}. \quad (\text{B.3})$$

Note that By definition of $\boldsymbol{\beta}^{(j)}$, the second term satisfies $\sum_{j=1}^p \sqrt{\sum_{s=1}^m \|\boldsymbol{\beta}_s^{(j)}\|_2^2} = \|\boldsymbol{\beta}^{(j)}\|_2$. Finally, by simply using the vectorized notation we get Equation (6), concluding the proof.

Appendix C. Proof of Statements in Section 3

For the first three lemmas in this section we follow the technique in Lounici et al. (2011) and occasionally use classical ideas established in Bühlmann and Van De Geer (2011).

Notations and naming conventions When X is a matrix, $\|X\|_2$ and $\|X\|_F$ denote its spectral and Frobenius norm, respectively. Consider the multi-task coefficients vector $\boldsymbol{\beta} \in \mathbb{R}^{md}$, and the sub-vector $\boldsymbol{\beta}^{(j)} \in \mathbb{R}^{md_j}$ which denotes the coefficients corresponding to kernel k_j . Through out this proof, we use the convention “group j ” to refer to the set of indices of $\boldsymbol{\beta}$ which indicate $\boldsymbol{\beta}^{(j)}$. Similarly, we let $\boldsymbol{\Phi}^{(j)}$ denote the $mn \times md_j$ sub-matrix which only has the features coming from group j . Lastly, let $\boldsymbol{\Psi} := \boldsymbol{\Phi}^T \boldsymbol{\Phi} / mn$, then $\boldsymbol{\Psi}^{(j)} = (\boldsymbol{\Phi}^{(j)})^T \boldsymbol{\Phi}^{(j)} / mn$ indicates the $md_j \times md_j$ submatrix that is caused by group j .

C.1 Proof of Theorem 4

Recall the vectorized formulation of the META-KEL loss,

$$\mathcal{L}(\boldsymbol{\beta}) = \frac{1}{mn} \|\mathbf{y} - \boldsymbol{\Phi} \boldsymbol{\beta}\|_2^2 + \lambda \sum_{j=1}^p \|\boldsymbol{\beta}^{(j)}\|_2.$$

Let $\boldsymbol{\varepsilon}_s = (\varepsilon_{s,i})_{i \leq n}$ denote error for task s and $\boldsymbol{\varepsilon} \in \mathbb{R}^{mn}$ the stacked multi-task error vector. Using $\mathbf{y} = \boldsymbol{\Phi} \boldsymbol{\beta} + \boldsymbol{\varepsilon}$, we may decompose the loss into two deterministic and random parts. The term $2\boldsymbol{\varepsilon}^T \boldsymbol{\Phi}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*) / mn$ is the random one and we will refer to as the empirical process. The first typical step in bounding the estimation error of Lasso estimators, is showing that the empirical process, which comes from the noise in observing values of \mathbf{y} , does not play a drastic role. More formally, let $A_j = \{\|(\boldsymbol{\Phi} \boldsymbol{\varepsilon})^{(j)}\|_2 / mn \leq \lambda/4\}$ denote that the event that the image of noise affecting the feature space of $\boldsymbol{\phi}_j$, is dominated by the regularization term of $\boldsymbol{\beta}^{(j)}$. In Lemma 8 we show that $\cap_{j=1}^p A_j$ happens with high probability, if λ is set properly.

Lemma 8 (Regularization term dominates the empirical process) *Set $0 < \delta < 1$. Consider the random event $A = \cap_{j=1}^p A_j$. Then A happens with probability greater than $1 - \delta$, if*

$$\lambda \geq \frac{4\sigma}{\sqrt{mn}} \sqrt{1 + \frac{2}{m} \left(\log(2p/\delta) + \sqrt{md_{\max} \log(2p/\delta)} \right)}.$$

where $d_{\max} = \max_{1 \leq j \leq p} d_j$.

We now show that if the empirical process is controlled by regularization, i.e. if λ is set to be large enough, then $\hat{\boldsymbol{\beta}} = \min \mathcal{L}(\boldsymbol{\beta})$ has favorable properties.

Lemma 9 (Conditional properties of $\hat{\boldsymbol{\beta}}$) *Assume that event A happens. Then for any solution $\hat{\boldsymbol{\beta}}$ of problem 6 and all $\boldsymbol{\beta} \in \mathbb{R}^{md}$, the following hold:*

$$\begin{aligned} \frac{1}{mn} \|\boldsymbol{\Phi}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|_2^2 + \frac{\lambda}{2} \sum_{j=1}^p \|\hat{\boldsymbol{\beta}}^{(j)} - \boldsymbol{\beta}^{*(j)}\|_2 &\leq \|\boldsymbol{\Phi}(\boldsymbol{\beta} - \boldsymbol{\beta}^*)\|_2^2 \\ + 2\lambda \sum_{j: \boldsymbol{\beta}^{(j)} \neq 0} \min \left(\|\boldsymbol{\beta}^{(j)}\|_2, \|\hat{\boldsymbol{\beta}}^{(j)} - \boldsymbol{\beta}^{(j)}\| \right) & \end{aligned} \quad (\text{C.1})$$

$$\left| \left\{ j : \hat{\boldsymbol{\beta}}^{(j)} \neq 0 \right\} \right| \leq \frac{1}{(mn\lambda)^2} \|\boldsymbol{\Phi}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|_2^2 \quad (\text{C.2})$$

Note that by the meta data-generating model (Section 2), $J_{k^*} = \{j : \boldsymbol{\beta}^{*(j)} \neq 0\}$.

Lemma 10 (Complete Variable Screening) Assume $|J_{k^*}| \leq s$ and set $0 \leq \delta \leq 1$. Define

$$\epsilon(n, m) = \frac{32\sigma s}{\kappa^2 \sqrt{mn}} \sqrt{1 + \frac{2}{m} \left(\log(2p/\delta) + \sqrt{md_{\max} \log(2p/\delta)} \right)}$$

Under assumption 3 with $\kappa = \kappa(s)$, if λ is chosen according to lemma 8, then with probability greater than $1 - \delta$

$$\max_{j \in J_{k^*}} \left\| \hat{\beta}^{(j)} - (\beta^*)^{(j)} \right\|_2 \leq \epsilon(n, m) \quad (\text{C.3})$$

and if in addition $\min_{j \in J_{k^*}} \|\beta^{*(j)}\|_2 \geq c_1$, then with the same probability for all $j \in J_{k^*}$

$$\left| \left\| \hat{\beta}^{(j)} \right\|_2 - c_1 \right| \leq \epsilon(n, m). \quad (\text{C.4})$$

We now turn to the first claim made in Theorem 4, and prove that $\mathcal{H}_{k^*} \subseteq \mathcal{H}_{\hat{k}}$. Since $\hat{\eta}_j \geq 0$ and k_i are Mercer, then \hat{k} is also Mercer and corresponds to an RKHS which we have been referring to as $\mathcal{H}_{\hat{k}}$. Consider the RKHS \mathcal{H}_{k^*} , since $|J_{k^*}|$ and each d_j s are finite,

$$\mathcal{H}_{k^*} = \left\{ f : f(\cdot) = \sum_{j \in J_{k^*}} \sqrt{\eta_j^*} \beta_j^T \phi_j(\cdot), \beta \in \mathbb{R}^d, \|\beta_j\| < \infty \right\}$$

Therefore, $f \in \mathcal{H}_{k^*}$ if and only if it is in the finite span of ϕ , defined as

$$\text{FinSpan}(\{\phi_j : j \in J_{k^*}\}) = \left\{ f : f(\cdot) = \sum_{j \in J_{k^*}} \beta_j^T \phi_j(\cdot), \beta \in \mathbb{R}^d, \|\beta_j\| < \infty \right\}$$

Lemma 10 states that $\hat{\beta}_j \geq c_1 - \epsilon(n, m)$ with probability greater than $1 - \delta$ for $j \in J_{k^*}$. Therefore, for any $j \in J_{k^*}$, we get $\hat{\eta}_j \geq 1 - \epsilon(n, m)/c_1$, since we had set $\hat{\eta}_j = \|\beta^{*(j)}\|/c_1$. If $c_1 \geq \epsilon(n, m)$, then $\hat{\eta}_j > 0$ and $j \in J_{\hat{k}}$, implying $J_{k^*} \subset J_{\hat{k}}$. Hence, under the assumptions of the theorem, with probability greater than $1 - \delta$,

$$\mathcal{H}_{k^*} = \text{FinSpan}(\{\phi_j : j \in J_{k^*}\}) \stackrel{\text{w.h.p.}}{\subset} \text{FinSpan}(\{\phi_j : j \in J_{\hat{k}}\}) = \mathcal{H}_{\hat{k}}.$$

The next lemma bounds the \hat{k} -norm of functions contained in \mathcal{H}_{k^*} , concluding the proof for Theorem 4.

Lemma 11 (Bounding the \hat{k} -norm) Set $0 \leq \delta \leq 1$ and choose λ according to Lemma 8 and define $\epsilon(n, m)$ according to Lemma 10. If $|J_{k^*}| \leq s$ and Assumption 3 holds with $\kappa = \kappa(s)$, then under Condition 1, for all $f \in \mathcal{H}_{k^*}$ with $\|f\|_{k^*} \leq B$,

$$\|f\|_{\hat{k}} \leq \left(1 + \frac{\epsilon(n, m)}{2c_1} + o(\epsilon(n, m)) \right) \quad (\text{C.5})$$

C.2 Proof of Proposition 6

Under assumptions of the proposition, event A happens with probability greater than $1 - \delta$. From Equation (C.2) in Lemma 9,

$$|J_{\hat{k}}| \leq \frac{1}{(mn\lambda)^2} \left\| \Phi(\hat{\beta} - \beta^*) \right\|_2^2$$

and by Equation (C.10),

$$\frac{1}{\sqrt{mn}} \left\| \Phi(\hat{\beta} - \beta^*) \right\|_2 \leq \frac{2\lambda\sqrt{s}}{\kappa}$$

which gives

$$|J_{\hat{k}}| \leq \frac{4s}{mn\kappa^2(s)}$$

Therefore, if $mn = \mathcal{O}(s/p)$ then $|J_{\hat{k}}| \leq p = |J_{k_{\text{full}}}|$ with probability greater than $1 - \delta$. and via similar argument as given in the proof of theorem 4,

$$\mathcal{H}_{\hat{k}} = \text{FinSpan}(\{\phi_j : j \in J_{\hat{k}}\}) \stackrel{\text{w.h.p.}}{\subset} \text{FinSpan}(\{\phi_j : j \in J_{k_{\text{full}}}\}) = \mathcal{H}_{k_{\text{full}}}.$$

C.3 Proof of Lemmas used in Section C.1

This section presents the proofs to the helper lemmas introduced before.

Proof [Proof of Lemma 8] This proof follows a similar treatment of the empirical process to Lemma 3.1 Lounici et al. (2011). Since $\varepsilon_{i,s}$ are i.i.d. zero-mean sub-gaussian variables, we observe that

$$\mathbb{P}(A_j) = \mathbb{P}\left(\left\{\frac{1}{(mn)^2} \varepsilon^T \Phi^{(j)} (\Phi^{(j)})^T \varepsilon \leq \frac{\lambda^2}{16}\right\}\right) = \mathbb{P}\left(\left\{\frac{\sum_{i=1}^{mn} v_i (z_i^2 - 1)}{\sqrt{2}\|\mathbf{v}\|} \leq \alpha\right\}\right)$$

where z_i are i.i.d. sub-gaussian variables with variance proxy 1, v_i denote the eigenvalues of $\Phi^{(j)} (\Phi^{(j)})^T / mn$ and \mathbf{v} is the vector of these eigenvalues. Lastly,

$$\alpha = \frac{\lambda^2 / (16\sigma^2) - \text{Tr}(\Psi)}{\sqrt{2}\|\Psi\|_F}$$

From Equation 27 Cavalier et al. (2002) yields the following inequality,

$$\mathbb{P}(A_j^c) = \mathbb{P}\left(\left|\frac{\sum_{i=1}^{mn} (z_i^2 - 1)v_i}{\sqrt{2}\|\mathbf{v}\|_2}\right| \geq \alpha\right) \leq 2 \exp\left(-\frac{\alpha^2}{2(1 + \sqrt{2}\alpha\|\mathbf{v}\|_\infty/\|\mathbf{v}\|_2)}\right)$$

We choose λ such that the right hand side is bounded by δ/p . By definition of \mathbf{v} we have $\|\mathbf{v}\|_\infty/\|\mathbf{v}\|_2 = \|\Psi\|_2/\|\Psi\|_F$. Then, for A_j^c to happen with probability smaller than δ/p ,

$$\lambda \geq \frac{4\sigma}{\sqrt{mn}} \sqrt{\text{Tr}(\Psi^{(j)}) + 2\|\Psi^{(j)}\|_2 \left(2 \log(2p/\delta) + \sqrt{md_j \log(2p/\delta)}\right)}.$$

Then by union bound, A happens with probability greater than $1 - \delta$ if

$$\lambda \geq \max_j \frac{4\sigma}{\sqrt{mn}} \sqrt{\text{Tr}(\Psi^{(j)}) + 2\|\Psi^{(j)}\|_2 \left(2 \log(2p/\delta) + \sqrt{md_j \log(2p/\delta)}\right)}.$$

Since the base kernels are normalized we may bound the norm and trace of $\Psi^{(j)}$.

$$\text{Tr}(\Psi^{(j)}) = \frac{1}{mn} \sum_{s=1}^m \text{Tr} \left((\Phi_s^{(j)})^T \Phi_s^{(j)} \right) = \frac{1}{mn} \sum_{s=1}^m \sum_{i=1}^n \phi_j^T(\mathbf{x}_{s,i}) \phi_j(\mathbf{x}_{s,i}) = \frac{1}{mn} \sum_{s=1}^m \sum_{i=1}^n k_j(\mathbf{x}_{s,i}, \mathbf{x}_{s,i}) \leq 1$$

Similarly, $\|\Psi^{(j)}\|_2 \leq \max_s \sum_{i=1}^n k_j(\mathbf{x}_{s,i}, \mathbf{x}_{s,i}) \leq 1/m$, and thereby concluding the proof. \blacksquare

Proof [Proof of Lemma 9] For proving this lemma we are essentially only using Cauchy-Schwarz and the Triangle inequality, together with the KKT optimality condition for \mathcal{L} . For any β , since $\hat{\beta}$ is the minimizer of \mathcal{L} and due to the data generating model (Eq. ??), we have

$$\frac{1}{mn} \left\| \Phi(\hat{\beta} - \beta^*) \right\|_2^2 \leq \frac{1}{mn} \left\| \Phi(\beta - \beta^*) \right\|_2^2 + \frac{2}{m} \varepsilon^T \Phi(\hat{\beta} - \beta) + \lambda \sum_{j=1}^p \left(\left\| \beta^{(j)} \right\|_2 - \left\| \hat{\beta}^{(j)} \right\|_2 \right).$$

By Cauchy-Schwarz and the assumption that A happens,

$$\varepsilon^T \Phi(\hat{\beta} - \beta) \leq \sum_{j=1}^p \left\| (\Phi \varepsilon)^{(j)} \right\|_2 \left\| \hat{\beta}^{(j)} - \beta^{(j)} \right\|_2 \leq \frac{mn\lambda}{4} \sum_{j=1}^p \left\| \hat{\beta}^{(j)} - \beta^{(j)} \right\|_2,$$

and thereby,

$$\frac{1}{mn} \left\| \Phi(\hat{\beta} - \beta^*) \right\|_2^2 + \frac{\lambda}{2} \sum_{j=1}^p \left\| \hat{\beta}^{(j)} - \beta^{(j)} \right\|_2 \leq \frac{1}{mn} \left\| \Phi(\beta - \beta^*) \right\|_2^2 + \lambda \sum_{j=1}^p \left(\left\| \hat{\beta}^{(j)} - \beta^{(j)} \right\|_2 + \left\| \beta^{(j)} \right\|_2 - \left\| \hat{\beta}^{(j)} \right\|_2 \right)$$

which gives Inequality C.1. By the KKT optimality conditions for convex losses (Boyd et al., 2004), $\hat{\beta}$ is a minimizer of \mathcal{L} , if and only if $0 \in \partial \mathcal{L}(\hat{\beta})$, where $\partial \mathcal{L}(\hat{\beta})$ denotes the sub-gradient of the loss evaluated at $\hat{\beta}$. Therefore $\hat{\beta}$ satisfies

$$\frac{2}{m} \left(\Phi^T(\mathbf{y} - \Phi \hat{\beta}) \right)^{(j)} = \frac{\lambda \hat{\beta}^{(j)}}{\left\| \hat{\beta}^{(j)} \right\|}, \quad \text{if } \hat{\beta}^{(j)} \neq 0 \quad (\text{C.6})$$

$$\frac{2}{m} \left\| \left(\Phi^T(\mathbf{y} - \Phi \hat{\beta}) \right)^{(j)} \right\|_2 \leq \lambda, \quad \text{if } \hat{\beta}^{(j)} = 0. \quad (\text{C.7})$$

As for Inequality C.2, conditioned on event A , by ?? together with the KKT condition C.6, we obtain that for all $1 \leq j \leq p$ where $\hat{\beta}^{(j)} \neq 0$,

$$\frac{1}{mn} \left\| \left(\Phi^T \Phi(\hat{\beta} - \beta^*) \right)^{(j)} \right\|_2 \geq \lambda.$$

Following the analysis of Lounici et al. we conclude,

$$\begin{aligned} \left| \left\{ j : \hat{\beta}^{(j)} \neq 0 \right\} \right| &\leq \frac{1}{(mn\lambda)^2} \sum_{j: \hat{\beta}^{(j)} \neq 0} \left\| \left(\Phi^T \Phi(\hat{\beta} - \beta^*) \right)^{(j)} \right\|_2^2 \\ &\leq \frac{1}{(mn\lambda)^2} \sum_{j=1}^p \left\| \left(\Phi^T \Phi(\hat{\beta} - \beta^*) \right)^{(j)} \right\|_2^2 \\ &\leq \frac{1}{(mn\lambda)^2} \left\| \Phi^T \Phi(\hat{\beta} - \beta^*) \right\|_2^2 \\ &\leq \frac{1}{(mn\lambda)^2} \left\| \Phi(\hat{\beta} - \beta^*) \right\|_2^2. \end{aligned}$$

Since the kernels k_j are normalized by 1 and $\|\Phi\|_2 \leq \max_j k_j(\mathbf{x}, \mathbf{x}) \leq 1$. \blacksquare

Proof [Proof of Lemma 10] Let $\beta_{J_{k^*}} = (\beta^{(j)})_{j \in J_{k^*}}$ denote the sub coefficient vector that corresponds to all the active groups. Due to Equation (C.1) with $\beta = \beta^*$, conditioned on event A we have

$$\frac{1}{mn} \left\| \Phi(\hat{\beta} - \beta^*) \right\|_2^2 \leq 2\lambda \sum_{j \in J_{k^*}} \left\| \hat{\beta}^{(j)} - \beta^{*(j)} \right\| \leq 2\lambda \sqrt{\sum_{j \in J_{k^*}} s \left\| \hat{\beta}^{(j)} - \beta^{*(j)} \right\|_2^2} = 2\lambda\sqrt{s} \left\| \hat{\beta}_{J_{k^*}} - \beta^*_{J_{k^*}} \right\|_2 \quad (\text{C.8})$$

where the second inequality follows from Cauchy-Schwarz together with the Lemma's assumption $|J_{k^*}| \leq s$. By again using Equation (C.1) we get $\sum_{j=1}^p \left\| \hat{\beta}^{(j)} - \beta^{*(j)} \right\|_2 \leq 4 \sum_{j \in J_{k^*}} \left\| \hat{\beta}^{(j)} - \beta^{*(j)} \right\|_2$, and thereby $\sum_{j \notin J_{k^*}} \left\| \hat{\beta}^{(j)} - \beta^{*(j)} \right\|_2 \leq 3 \sum_{j \in J_{k^*}} \left\| \hat{\beta}^{(j)} - \beta^{*(j)} \right\|_2$. Using assumption 3, this inequality indicates that

$$\left\| \hat{\beta}_{J_{k^*}} - \beta^*_{J_{k^*}} \right\|_2 \leq \frac{1}{\kappa\sqrt{mn}} \left\| \Phi(\hat{\beta} - \beta^*) \right\|_2, \quad (\text{C.9})$$

which together with Equation (C.8) gives,

$$\frac{1}{\sqrt{mn}} \left\| \Phi(\hat{\beta} - \beta^*) \right\|_2 \leq \frac{2\lambda\sqrt{s}}{\kappa} \quad (\text{C.10})$$

The next chain of inequalities proves the first statement of the Lemma (Equation (C.3)). For all $j \in J_{k^*}$,

$$\begin{aligned} \left\| \hat{\beta}^{(j)} - \beta^{*(j)} \right\|_2 &\leq \sum_{j=1}^p \left\| \hat{\beta}^{(j)} - \beta^{*(j)} \right\|_2 \\ &\stackrel{\text{C.1}}{\leq} 4 \sum_{j \in J} \left\| \hat{\beta}^{(j)} - \beta^{*(j)} \right\|_2 \\ &\stackrel{\text{C.8}}{\leq} 4\sqrt{s} \left\| \hat{\beta}_{J_{k^*}} - \beta^*_{J_{k^*}} \right\|_2 \\ &\stackrel{\text{C.9}}{\leq} \frac{4\sqrt{s}}{\kappa\sqrt{mn}} \left\| \Phi(\hat{\beta} - \beta^*) \right\|_2 \\ &\stackrel{\text{C.10}}{\leq} \frac{8\lambda s}{\kappa^2} \end{aligned}$$

Note that the analysis here where carried out conditional on even A . From Lemma 8, we have that A happens with probability greater than $1 - \delta$ if λ is set according to the statement of the Lemma. To conclude the proof, Recall that from the Beta-min condition (Cond. 1) we have $\|\beta^{*(j)}\| \leq c_1$, which together with Equation (C.3) concludes the proof. \blacksquare

Proof [Proof of Lemma 11] Let $f \in \mathcal{H}_{k^*}$, with $\|f\|_k^2 \leq B^2$. Then by construction of \mathcal{H}_{k^*} (see Appendix B.1)

$$f(\mathbf{x}) = \sum_{j=1}^p \sqrt{\eta_j^*} \phi_j^T(\mathbf{x}) (\boldsymbol{\beta})^{(j)}, \quad \|f\|_k^2 = \sum_{j:\eta_j^* \neq 0} \left\| (\boldsymbol{\beta})^{(j)} \right\|_2^2 = \|\boldsymbol{\beta}\|_2^2.$$

We calculate the \hat{k} norm of functions that lie in \mathcal{H}_{k^*} . Let $I = \{1 \leq j \leq p : \eta_j^* \neq 0, \hat{\eta}_j \neq 0\}$.

$$\|f\|_{\hat{k}}^2 = \sum_{j:\hat{\eta}_j \neq 0} \sum_{r=1}^{d_j} \frac{(\langle f, \phi_{j,r} \rangle_2)^2}{\hat{\eta}_j} = \sum_{j:\hat{\eta}_j \neq 0} \sum_{r=1}^{d_j} \frac{(\sqrt{\eta_j^*} \boldsymbol{\beta}_r^{(j)})^2}{\hat{\eta}_j} = \sum_{j \in I} \frac{\eta_j^*}{\hat{\eta}_j} \|\boldsymbol{\beta}^{(j)}\|_2^2 = \sum_{j \in I} \frac{\eta_j^*}{\hat{\eta}_j} \|\boldsymbol{\beta}^{(j)}\|_2^2$$

where $\phi_{j,r}$ denotes the r -th feature in the feature map ϕ_j , similarly $\boldsymbol{\beta}_r^{(j)}$ the r -th element in vector $\boldsymbol{\beta}^{(j)}$, and $\langle \cdot, \cdot \rangle_2$ is the inner product in the $L^2(\mathcal{X})$ space. By applying Cauchy-Schwarz we get

$$\|f\|_{\hat{k}}^2 \leq \sqrt{\sum_{j \in I} \left(\frac{\eta_j^*}{\hat{\eta}_j}\right)^2} \cdot \sqrt{\sum_{j \in I} \|\boldsymbol{\beta}^{(j)}\|_2^4}.$$

Consider the vector $\mathbf{v} = (\|\boldsymbol{\beta}^{(j)}\|_2^2)_{j \in I}$, we observe that $\|\mathbf{v}\|_2 = \sqrt{\sum_{j \in I} \|\boldsymbol{\beta}^{(j)}\|_2^4}$ and $\|\mathbf{v}\|_1 \leq B^2$. Since $\|\cdot\|_2 \leq \|\cdot\|_1$ and due to the assumption $\|\boldsymbol{\eta}^*\|_1 \leq 1$, we obtain

$$\|f\|_{\hat{k}}^2 \leq B^2 \sum_{j \in I} \frac{\eta_j^*}{\hat{\eta}_j} \leq B^2 \|\boldsymbol{\eta}^*\|_1 \max_{j \in I} \frac{1}{\hat{\eta}_j} \leq B^2 \max_{j \in I} \frac{1}{\hat{\eta}_j}. \quad (\text{C.11})$$

It remains to bound $\max_{j \in I} \hat{\eta}_j^{-1}$. We set $\hat{\boldsymbol{\beta}}^{(j)} = \|\hat{\boldsymbol{\beta}}^{(j)}\|_2 / (c_1)$ and under conditions of the theorem, Lemma 10 states that $\hat{\boldsymbol{\beta}}^{(j)} \geq 1 - \epsilon(n, m)/c_1$, for all $j \in J_{k^*}$. Then for members of $I \subset J_{k^*}$,

$$\frac{1}{\hat{\eta}_j} \leq \frac{1}{1 - \epsilon(n, m)/c_1} \leq \left(1 + \frac{\epsilon(n, m)}{c_1} + o(\epsilon(n, m))\right),$$

which implies the following for the norm bound

$$\|f\|_{\hat{k}} \leq B \left(1 + \frac{\epsilon(n, m)}{2c_1} + o(\epsilon(n, m))\right) \quad (\text{C.12})$$

with probability greater than $1 - \delta$. ■

Appendix D. Proof of Statements in Section 4

The following lemma presents a confidence bound for when the learner has oracle knowledge of the true kernel. This lemma plays an integral role in for the proofs in this section.

Lemma 12 (Theorem 2 Chowdhury and Gopalan (2017) for \hat{k}) Let $f \in \mathcal{H}_{\hat{k}}$ for some kernel \hat{k} with a \hat{k} -norm bounded by \hat{B} . Then with probability greater than $1 - \bar{\delta}$, for all $\mathbf{x} \in \mathcal{X}$ and $1 \leq t \leq T$,

$$|\hat{\mu}_{t-1}(\mathbf{x}) - f(\mathbf{x})| \leq \hat{\sigma}_{t-1}(\mathbf{x}) \left(\hat{B} + \sigma \sqrt{2(\hat{\gamma}_{t-1} + 1 + \log(1/\bar{\delta}))} \right) \quad (\text{D.1})$$

where $\hat{\mu}_{t-1}$ and $\hat{\sigma}_{t-1}$ are as defined in Equation 2 with $\bar{\sigma} = 1 + 2/T$.

We skip the proof for this lemma as it is given in Chowdhury and Gopalan (2017), with the same notation. For the kernel \hat{k} we define the *maximum information gain* after $t - 1$ observations as

$$\hat{\gamma}_{t-1} := \max_{[\mathbf{x}_\tau]_{\tau \leq t}} \frac{1}{2} \log \det(\mathbf{I} + \bar{\sigma}^{-2} \hat{\mathbf{K}}_{t-1})$$

This parameter quantifies the speed at which we learn about f , when using the kernel \hat{k} . Note that $\hat{\gamma}_{t-1}$ is independent of any specific realization of H_{t-1} . It only depends on the choice of kernel, the input domain, and the noise variance. The next lemma bounds this parameter.

Lemma 13 (Information Gain Bound) The maximum information gain for \hat{k} after observing t samples satisfies,

$$\hat{\gamma}_t \leq \frac{\hat{d}}{2} \log \left(1 + \frac{\bar{\sigma}^{-2} t}{c_1} \right) = \mathcal{O}(\hat{d} \log t / c_1)$$

where $\hat{d} = \sum_{j \in J_{\hat{k}}} d_j \leq d$.

We now have the main tools for proving Theorem 5.

Proof [Proof of Theorem 5] Assume that $f \in \mathcal{H}_{k^*}$, and that $\|f\|_{k^*} \leq B$. Then by Theorem 4 $f \in \mathcal{H}_{\hat{k}}$, with probability greater than $1 - \delta$. Define \hat{B} as

$$\hat{B} = B \left(1 + \frac{\epsilon(n, m)}{2c_1} + o(\epsilon(n, m)) \right), \quad (\text{D.2})$$

then by Lemma 11, $\|f\|_{\hat{k}} \leq \hat{B}$ with probability greater than $1 - \delta$. We first condition on the event that $f \in \mathcal{H}_{\hat{k}}$ and $\|f\|_{\hat{k}} \leq \hat{B}$. Then Lemma 12 gives the following confidence interval,

$$\mathbb{P} \left(|\hat{\mu}_{t-1}(\mathbf{x}) - f(\mathbf{x})| \leq \hat{\sigma}_{t-1}(\mathbf{x}) \left(\hat{B} + \sigma \sqrt{2(\hat{\gamma}_{t-1} + 1 + \log(1/\bar{\delta}))} \right) \mid f \in \mathcal{H}_{\hat{k}}, \|f\|_{\hat{k}} \leq \hat{B} \right) \geq 1 - \bar{\delta} \quad (\text{D.3})$$

We now remove the conditional event. Let $C_t(\mathbf{x}) := [\hat{\mu}(\mathbf{x}) - \hat{\sigma}(\mathbf{x}), \hat{\mu}(\mathbf{x}) + \hat{\sigma}(\mathbf{x})]$. By the chain rule,

$$\mathbb{P}(f(\mathbf{x}) \in C_t(\mathbf{x})) \geq \mathbb{P}(f(\mathbf{x}) \in C_t(\mathbf{x}) \mid f \in \mathcal{H}_{\hat{k}}) \cdot \mathbb{P}(f \in \mathcal{H}_{\hat{k}}) \geq 1 - \delta - \bar{\delta}$$

Renaming $\delta + \bar{\delta}$ to δ for simplicity, we conclude that with probability greater than $1 - \delta$

$$|\hat{\mu}_{t-1}(\mathbf{x}) - f(\mathbf{x})| \leq \hat{\sigma}_{t-1}(\mathbf{x}) \left(B \left(1 + \frac{\epsilon(n, m)}{2c_1} + o(\epsilon(n, m)) \right) + \sigma \sqrt{2(\hat{\gamma}_{t-1} + 1 + \log(1/\bar{\delta}))} \right).$$

Lastly, Lemma 13 gives an upper bound for $\hat{\gamma}_{t-1}$ which completes the proof. \blacksquare

Proof [Proof of Lemma 13] For this proof, we follow a similar technique as in Vakili et al. (2021). Recall that $\hat{k}(\mathbf{x}, \mathbf{x}') = \sum_{j \in J_{\hat{k}}} \hat{\eta}_j \phi_j^T(\mathbf{x}) \phi_j(\mathbf{x}')$, where $J_{\hat{k}} = \{1 \leq j \leq p : \hat{\eta}_j \neq 0\}$. Define \hat{d} the effective dimension of kernels k_j that correspond to this index set, $\hat{d} = \sum_{j \in J_{\hat{k}}} d_j$. Now consider any arbitrary sequence of inputs $(\mathbf{x}_\tau)_{\tau=1}^t$ and let $\hat{\Phi}_t = \left(\hat{\phi}(\mathbf{x}_1), \dots, \hat{\phi}(\mathbf{x}_t) \right) \in \mathbb{R}^{t \times \hat{d}}$, with $\hat{\phi}(\mathbf{x}) = (\phi_j(\mathbf{x}))_{j \in J_{\hat{k}}}$. Define the $\hat{d} \times \hat{d}$ matrix $\Lambda = \text{diag} \left((\hat{\eta}_j)_{j \in J_{\hat{k}}} \right)$ as the diagonal matrix containing the eigenfunctions of \hat{k} . We have $K_t = \hat{\Phi}_t \Lambda \hat{\Phi}_t^T$. Let $H_t = \Lambda^{1/2} \hat{\Phi}_t^T \hat{\Phi}_t \Lambda^{1/2}$, by the Weinstein-Aronszajn identity,

$$\begin{aligned} \frac{1}{2} \log \det(\mathbf{I} + \bar{\sigma}^{-2} K_t) &= \frac{1}{2} \log \det(\mathbf{I} + \bar{\sigma}^{-2} H_t) \\ &\leq \frac{1}{2} \hat{d} \log (\text{tr}(\mathbf{I} + \bar{\sigma}^{-2} H_t) / \hat{d}) \end{aligned}$$

For positive definite matrices $\mathbf{P} \in \mathbb{R}^{n \times n}$, we have $\log \det \mathbf{P} \leq n \log \text{tr}(\mathbf{P}/n)$. The inequality follows from $\mathbf{I} + \bar{\sigma}^{-2} H_t$ being positive definite, since $\hat{\eta}_j \geq 0$. We may write,

$$\begin{aligned} \frac{1}{2} \log \det(\mathbf{I} + \bar{\sigma}^{-2} K_t) &\leq \frac{1}{2} \hat{d} \log \left(1 + \frac{\bar{\sigma}^{-2}}{\hat{d}} \text{tr}(\Lambda^{1/2} \hat{\Phi}_t^T \hat{\Phi}_t \Lambda^{1/2}) \right) \\ &\leq \frac{1}{2} \hat{d} \log \left(1 + \frac{\bar{\sigma}^{-2}}{\hat{d}} \sum_{\tau=1}^t \text{tr}(\Lambda^{1/2} \hat{\phi}^T(\mathbf{x}_\tau) \hat{\phi}(\mathbf{x}_\tau) \Lambda^{1/2}) \right) \\ &\leq \frac{1}{2} \hat{d} \log \left(1 + \frac{\bar{\sigma}^{-2}}{\hat{d}} \sum_{\tau=1}^t \|\hat{\phi}(\mathbf{x}_\tau) \Lambda^{1/2}\|_2^2 \right) \\ &\leq \frac{1}{2} \hat{d} \log \left(1 + \frac{\bar{\sigma}^{-2}}{\hat{d}} \sum_{\tau=1}^t \sum_{j \in J_{\hat{k}}} \hat{\eta}_j \|\phi_j(\mathbf{x}_\tau)\|_2^2 \right) \\ &\leq \frac{1}{2} \hat{d} \log \left(1 + \frac{\bar{\sigma}^{-2} t}{c_1} \max_j \|\hat{\beta}^{(j)}\|_2 \right) \end{aligned}$$

The next to last inequality holds since $k_j(\mathbf{x}, \mathbf{x}) = \|\phi_j(\mathbf{x})\|_2^2$ is normalized to one and $\hat{\eta}_j = \|\hat{\beta}^{(j)}\|_2 / c_1$. The inequality above holds for any sequence $(\mathbf{x}_\tau)_{\tau \leq t}$,

$$\hat{\gamma}_t = \mathcal{O} \left(\hat{d} \log(t/c_1) \right)$$

\blacksquare

D.1 Proof of Theorem 7

Similar to the previous section, we take advantage of a classic regret bound for the oracle learner and then apply it to our setting.

Lemma 14 (Theorem 3 Chowdhury and Gopalan (2017), for \hat{k}) Set $\delta \in (0, 1)$. If $f \in \mathcal{H}_{\hat{k}}$ with $\|f\|_{\hat{k}} \leq \hat{B}$, then with probability $1 - \delta$, GP-UCB satisfies,

$$R_T = \mathcal{O} \left(\hat{B} \sqrt{T \hat{\gamma}_T} + \sqrt{T \hat{\gamma}_T (\hat{\gamma}_T + \log 1/\delta)} \right)$$

Proof [Proof of Corollary 7] Conditioned on the event that $f \in \mathcal{H}_{\hat{k}}$ and $\|f\|_{\hat{k}} \leq \hat{B}$, Lemma 14 states that

$$\mathbb{P} \left(R_T = \mathcal{O} \left(\hat{B} \sqrt{T \hat{\gamma}_T} + \sqrt{T \hat{\gamma}_T (\hat{\gamma}_T + \log 1/\delta)} \right) \mid f \in \mathcal{H}_{\hat{k}}, \|f\|_{\hat{k}} \leq \hat{B} \right) \geq 1 - \bar{\delta}$$

Then by Lemma 11 and Theorem 4, with probability greater than $1 - \delta - \bar{\delta}$

$$R_T = \mathcal{O} \left(\hat{B} \sqrt{T \hat{\gamma}_T} + \sqrt{T \hat{\gamma}_T (\hat{\gamma}_T + \log 1/\delta)} \right)$$

where \hat{B} is set according to Equation (D.2). Plugging in Lemma 13 to bound the information gain and changing the variable name $\delta + \bar{\delta}$ to δ and concludes the proof. ■

Appendix E. Experiments

In this section, we provide experiments to quantitatively illustrate our theoretical contribution.

Experiment Setup We create a synthetic dataset based on our data model, Equations (??) and (??). We limit the domain to $\mathcal{X} = [-1, 1]$ and use Legendre polynomials P_j as our features ϕ_j . The sequence $(P_j)_{j \geq 1}$ is a natural choice, since it provides an orthonormal basis for $L^2(\mathcal{X})$. Moreover, Legendre polynomials are eigenfunctions to dot-product kernels such as the Neural Tangent Kernel (Jacot et al., 2018). We let $k^*(x, x') = \sum_{j \in J_{k^*}} \eta_j^* P_j(x) P_j(x')$, where J_{k^*} is a random subset of $\{1, \dots, p\}$. Each η_j^* is sampled independently from the standard uniform distribution and the vector η^* is then normalized. Across all experiments, we set $p = 20$ and $s = |J_{k^*}| = 5$. To sample the meta-data $\mathcal{D}_{n,m}$, we choose m independent random subsets of J_{k^*} and generate the functions f_s via Equation (4) where $\beta^{*(j)}$ are drawn from an i.i.d. standard uniform distribution. We then scale the norm $\|f\|_{k^*}$ to $B = 10$. The data for a single task, i.e., $(x_{s,i}, y_{s,i})_{i \leq n}$, is then created by uniformly drawing i.i.d. samples from the domain \mathcal{X} and evaluating f_s at those points. We add Gaussian noise with standard deviation $\sigma = 0.01$ to all data points. Figure 6 in the appendix shows

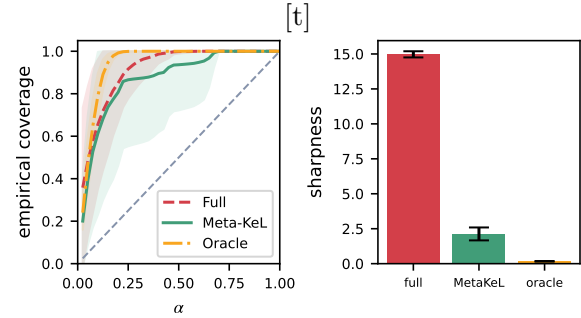


Figure 4: Calibration (left) and sharpness (right) experiment for confidence sets given 4 training samples. Averaged over 50 runs, \hat{k} always gives tight valid confidence intervals.

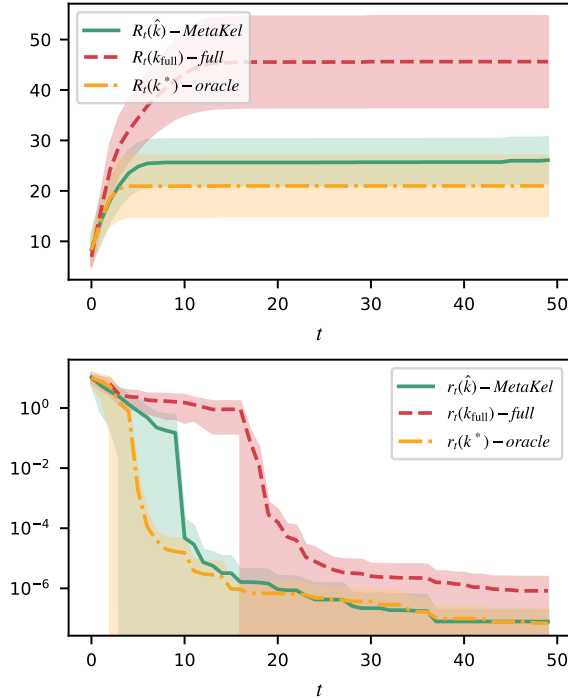


Figure 5: Simple and cumulative regret for GP-UCB. The algorithm converges at a slower pace when using k_{full} .

how random f_s may look like. For all experiments we set $n = m = 50$ unless stated otherwise. To meta-learn \hat{k} , we solve the vectorized META-KEL problem (Eq. 6) over $\beta^{(j)}$ with CELER, a fast solver for the group Lasso (Massias et al., 2018), and then set $\hat{\eta}$ according to Proposition 2. We set $\lambda = 0.03$, such that it satisfies the condition of Theorem 4. As shown in Figure 8 in the appendix, the choice of λ has little effect on the performance of the algorithm.

Confidence Set Experiment We perform calibration and sharpness experiments to assess the meta-learned confidence sets (Gneiting et al., 2007). Figure 4 presents the result. To obtain an α -confidence interval for some $f(x)$ using a kernel k , we assume a $f \sim \text{GP}(0, k)$ prior and calculate the α -quantile of the posterior after observing 4 noisy i.i.d. draws from the function. For each hypothesis, the y -axis of the left plot shows the empirical coverage of the confidence sets, i.e., the fraction of test points contained in the α -confidence intervals for varying levels α . In this plot, if a curve were to fall below the $x = y$ line, it would have implied insufficient coverage and hence over-confident sets. The plot on the right shows the posterior variance averaged across all test points. This quantity, referred to as *sharpness*, reflects the width of the confidence bands. Figure 4 demonstrates that the meta-learned confidence sets are well-calibrated for the entire range of confidence-levels and are tight relative to the true sets. In contrast, k_{full} yields conservative confidence sets, due to considering polynomials P_j that do not contribute to the construction of $f(x)$. We use 1000 test points for calculating the empirical averages. The plot shows the values averaged over 50 runs, where for each the kernel k^* and the data are generated from scratch.

Regret Experiment We verify the performance of GP-UCB when used together with \hat{k} . We generate the random reward function f in a manner similar to f_s of the meta-data. The BO problem is simulated according to ??, and the actions are selected via Equation (A.2). Figure 7 in the appendix shows how this algorithm samples the domain and how the confidence estimates shrink by observing new samples. Keeping the underlying random k^* fixed, we generate 100 random instances of the meta-learning and the BO problem. In Figure 5 we present the average regret and its standard deviation. In these plots, the simple regret of GP-UCB with a kernel k is labeled $r_t(k) = f(x^*) - \max_{\tau \leq t} f(x_\tau)$. Respectively, the cumulative inference regret is $R_t(k) = \sum_{\tau \leq t} f(x^*) - \max_x \mu_{\tau-1}(x)$. The algorithm converges to the optimum using all three kernels. The meta-learned kernel, however, improves upon using k_{full} and results in a performance competitive to when k^* is known by GP-UCB. This behavior empirically confirms Theorem 7.

Appendix F. Additional Details on the Experiments

In this section we present the plots that accompany our experiment design from Section E. Figure 6 illustrates a few samples of the random functions that we optimize over in the experiments. The functions are constructed using the Legendre basis, as explained in the main text.

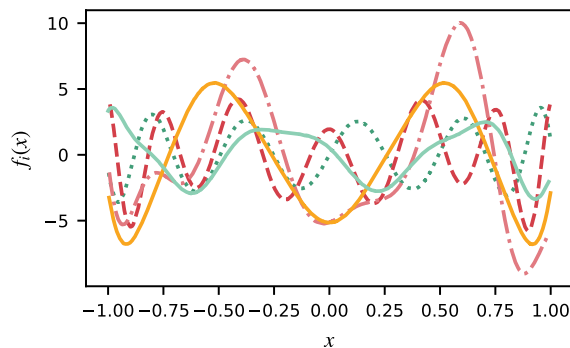


Figure 6: Examples of possible functions f_s for the meta-dataset.

Figure 7 gives an example of a BO problem where we use IGP-UCB together with the meta-learned kernel to find the minimum of a function with as few samples as possible. As the function estimate improves, the confidence sets rapidly shrink and the learner only samples points close to the minimum.

Figure 8 demonstrates that the choice of λ for the META-KEL loss does not have a severe effect on the regret of IGP-UCB. This is only the case if λ satisfies the condition of Theorem 4.

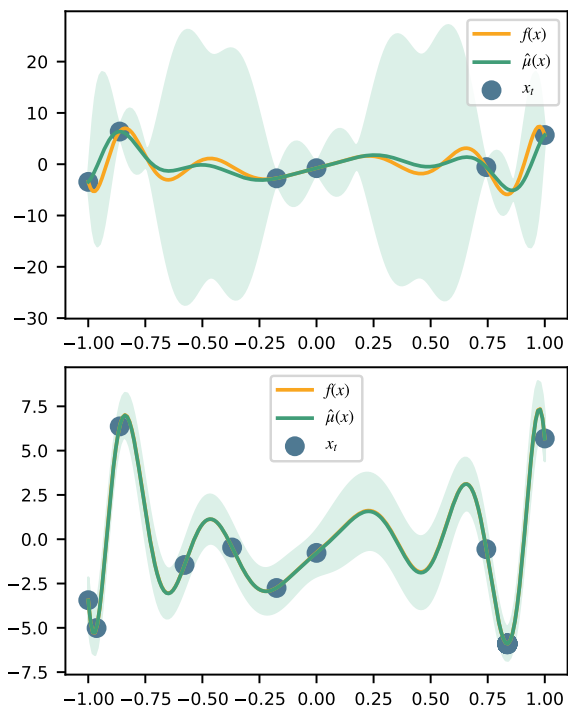


Figure 7: BO (minimization) with META-KEL. Upper plot shows the state at $t = 5$ and the lower plot at $t = 55$.

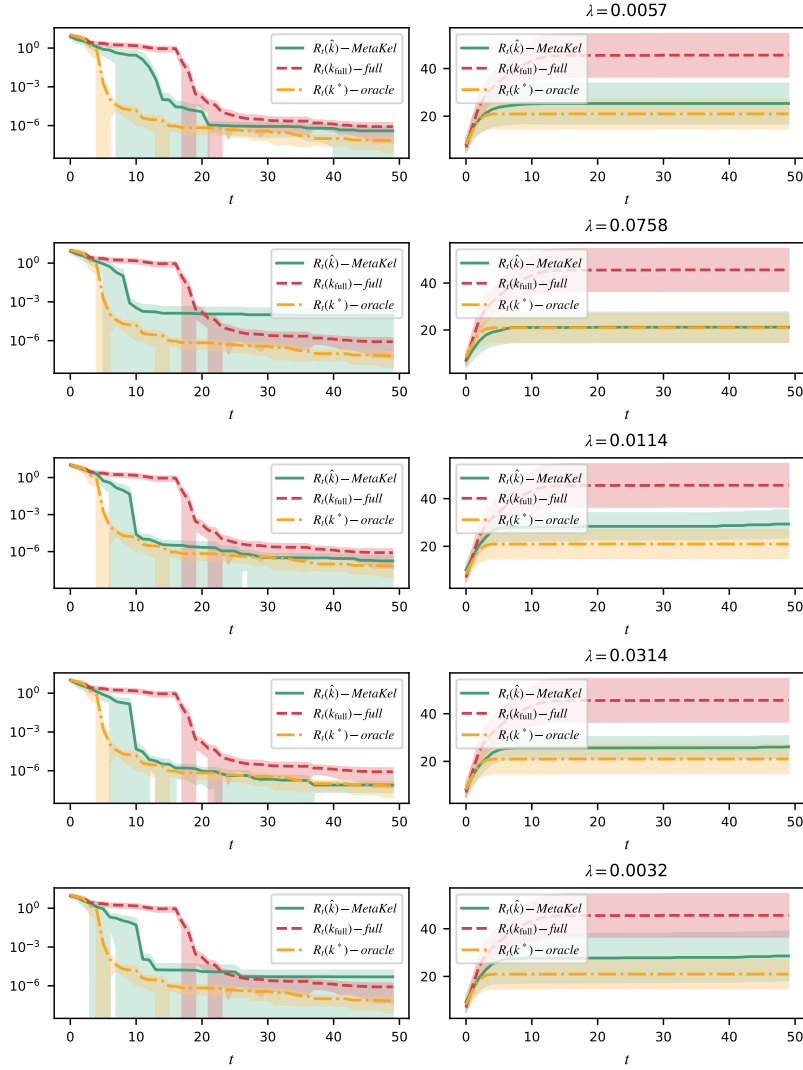


Figure 8: For $m = n = 50$ and $p = 20$, Theorem 4 requires that $\lambda > 0.001$ for recovery to happen with probability greater than $1 - \delta = 0.9$. For λ that satisfies this condition, the particular choice of its value does not effect performance severely.