

Near-Optimal Active Learning of Multi-Output Gaussian Processes

Yehong Zhang[†] and Trong Nghia Hoang^{*} and Kian Hsiang Low[†] and Mohan Kankanhalli[†]

Department of Computer Science[†], Interactive Digital Media Institute^{*}
National University of Singapore, Republic of Singapore
{yehong, lowkh, mohan}@comp.nus.edu.sg[†], idmhtn@nus.edu.sg^{*}

Abstract

This paper¹ addresses the problem of active learning of a *multi-output Gaussian process* (MOGP) model representing multiple types of coexisting correlated environmental phenomena. In contrast to existing works, our active learning problem involves selecting not just the most informative sampling locations to be observed but also the types of measurements at each selected location for minimizing the predictive uncertainty (i.e., posterior joint entropy) of a target phenomenon of interest given a sampling budget. Unfortunately, such an entropy criterion scales poorly in the numbers of candidate sampling locations and selected observations when optimized. To resolve this issue, we first exploit a structure common to sparse MOGP models for deriving a novel active learning criterion. Then, we exploit a relaxed form of submodularity property of our new criterion for devising a polynomial-time approximation algorithm that guarantees a constant-factor approximation of that achieved by the optimal set of selected observations. Empirical evaluation on real-world datasets shows that our proposed approach outperforms existing algorithms for active learning of MOGP and single-output GP models.

1 Introduction

For many budget-constrained environmental sensing and monitoring applications in the real world, active learning/sensing is an attractive, frugal alternative to passive high-resolution (hence, prohibitively costly) sampling of the spatially varying target phenomenon of interest. Different from the latter, active learning aims to select and gather the *most informative* observations for modeling and predicting the spatially varying phenomenon given some sampling budget constraints (e.g., quantity of deployed sensors, energy consumption, mission time).

In practice, the target phenomenon often coexists and correlates well with some auxiliary type(s) of phenomena whose measurements may be more spatially correlated, less noisy (e.g., due to higher-quality sensors), and/or less tedious to sample (e.g., due to greater availability/quantity, higher sampling rate, and/or lower sampling cost of deployed sensors of these type(s)) and can consequently be

exploited for improving its prediction. For example, to monitor soil pollution by some heavy metal (e.g., Cadmium), its complex and time-consuming extraction from soil samples can be alleviated by supplementing its prediction with correlated auxiliary types of soil measurements (e.g., pH) that are easier to sample (Goovaerts 1997). Similarly, to monitor algal bloom in the coastal ocean, plankton abundance correlates well with auxiliary types of ocean measurements (e.g., chlorophyll a, temperature, and salinity) that can be sampled more readily. Other examples of real-world applications include remote sensing, traffic monitoring (Chen et al. 2012; Chen, Low, and Tan 2013; Chen et al. 2015), monitoring of groundwater and indoor environmental quality (Xu et al. 2014), and precision agriculture (Webster and Oliver 2007), among others. All of the above practical examples motivate the need to design and develop an active learning algorithm that selects not just the *most informative* sampling locations to be observed but also the types of measurements (i.e., target and/or auxiliary) at each selected location for minimizing the predictive uncertainty of unobserved areas of a target phenomenon given a sampling budget, which is the focus of our work here².

To achieve this, we model all types of coexisting phenomena (i.e., target and auxiliary) jointly as a *multi-output Gaussian process* (MOGP) (Álvarez and Lawrence 2011; Bonilla, Chai, and Williams 2008; Teh and Seeger 2005), which allows the spatial correlation structure of each type of phenomenon and the cross-correlation structure between different types of phenomena to be formally characterized. More importantly, unlike the non-probabilistic multivariate regression methods, the probabilistic MOGP regression model allows the predictive uncertainty of the target phenomenon (as well as the auxiliary phenomena) to be formally quantified (e.g., based on entropy or mutual information criterion) and consequently exploited for deriving the active learning criterion.

To the best of our knowledge, this paper is the first to present an efficient algorithm for active learning of a MOGP

²Our work here differs from multivariate spatial sampling algorithms (Bueso et al. 1999; Le, Sun, and Zidek 2003) that aim to improve the prediction of *all* types of coexisting phenomena, for which existing active learning algorithms for sampling measurements only from the target phenomenon can be extended and applied straightforwardly, as detailed in Section 3.

¹This paper is an extended version (with proofs) of (Zhang et al. 2016).

model. We consider utilizing the entropy criterion to measure the predictive uncertainty of a target phenomenon, which is widely used for active learning of a single-output GP model. Unfortunately, for the MOGP model, such a criterion scales poorly in the number of candidate sampling locations of the target phenomenon (Section 3) and even more so in the number of selected observations (i.e., sampling budget) when optimized (Section 4). To resolve this scalability issue, we first exploit a structure common to a unifying framework of sparse MOGP models (Section 2) for deriving a novel active learning criterion (Section 3). Then, we define a relaxed notion of submodularity³ called ϵ -submodularity and exploit the ϵ -submodularity property of our new criterion for devising a polynomial-time approximation algorithm that guarantees a constant-factor approximation of that achieved by the optimal set of selected observations (Section 4). We empirically evaluate the performance of our proposed algorithm using three real-world datasets (Section 5).

2 Modeling Coexisting Phenomena with Multi-Output Gaussian Process (MOGP)

Convolved MOGP (CMOGP) Regression. A number of MOGP models such as co-kriging (Webster and Oliver 2007), parameter sharing (Skolidis 2012), and *linear model of coregionalization* (LMC) (Teh and Seeger 2005; Bonilla, Chai, and Williams 2008) have been proposed to handle multiple types of correlated outputs. A generalization of LMC called the *convolved MOGP* (CMOGP) model has been empirically demonstrated in (Álvarez and Lawrence 2011) to outperform the others and will be the MOGP model of our choice due to its approximation whose structure can be exploited for deriving our active learning criterion and in turn an efficient approximation algorithm, as detailed later.

Let M types of coexisting phenomena be defined to vary as a realization of a CMOGP over a domain corresponding to a set D of sampling locations such that each location $x \in D$ is associated with noisy realized (random) output measurement $y_{\langle x, i \rangle}$ ($Y_{\langle x, i \rangle}$) if x is observed (unobserved) for type i for $i = 1, \dots, M$. Let $D_i^+ \triangleq \{\langle x, i \rangle\}_{x \in D}$ and $D^+ \triangleq \bigcup_{i=1}^M D_i^+$. Then, measurement $Y_{\langle x, i \rangle}$ of type i is defined as a convolution between a smoothing kernel $K_i(x)$ and a latent measurement function $L(x)$ ⁴ corrupted by an additive noise $\varepsilon_i \sim \mathcal{N}(0, \sigma_{n_i}^2)$ with noise variance $\sigma_{n_i}^2$:

$$Y_{\langle x, i \rangle} \triangleq \int_{x' \in D} K_i(x - x') L(x') dx' + \varepsilon_i.$$

As shown in (Álvarez and Lawrence 2011), if $\{L(x)\}_{x \in D}$ is a GP, then $\{Y_{\langle x, i \rangle}\}_{\langle x, i \rangle \in D^+}$ is also a GP, that is, every finite subset of $\{Y_{\langle x, i \rangle}\}_{\langle x, i \rangle \in D^+}$ follows a multivariate

³The original notion of submodularity has been used in (Krause and Golovin 2014; Krause, Singh, and Guestrin 2008) to theoretically guarantee the performance of their algorithms for active learning of a *single-output* GP model.

⁴To ease exposition, we consider a single latent function. Note, however, that multiple latent functions can be used to improve the fidelity of modeling, as shown in (Álvarez and Lawrence 2011). More importantly, our proposed algorithm and theoretical results remain valid with multiple latent functions.

Gaussian distribution. Such a GP is fully specified by its *prior* mean $\mu_{\langle x, i \rangle} \triangleq \mathbb{E}[Y_{\langle x, i \rangle}]$ and covariance $\sigma_{\langle x, i \rangle \langle x', j \rangle} \triangleq \text{cov}[Y_{\langle x, i \rangle}, Y_{\langle x', j \rangle}]$ for all $\langle x, i \rangle, \langle x', j \rangle \in D^+$, the latter of which characterizes the spatial correlation structure for each type of phenomenon (i.e., $i = j$) and the cross-correlation structure between different types of phenomena (i.e., $i \neq j$). Specifically, let $\{L(x)\}_{x \in D}$ be a GP with prior covariance $\sigma_{xx'} \triangleq \mathcal{N}(x - x' | \underline{0}, P_0^{-1})$ and $K_i(x) \triangleq \sigma_{s_i} \mathcal{N}(x | \underline{0}, P_i^{-1})$ where $\sigma_{s_i}^2$ is the signal variance controlling the intensity of measurements of type i , P_0 and P_i are diagonal precision matrices controlling, respectively, the degrees of correlation between latent measurements and cross-correlation between latent and type i measurements, and $\underline{0}$ denotes a column vector comprising components of value 0. Then,

$$\sigma_{\langle x, i \rangle \langle x', j \rangle} = \sigma_{s_i} \sigma_{s_j} \mathcal{N}(x - x' | \underline{0}, P_0^{-1} + P_i^{-1} + P_j^{-1}) + \delta_{xx'}^{ij} \sigma_{n_i}^2 \quad (1)$$

where $\delta_{xx'}^{ij}$ is a Kronecker delta of value 1 if $i = j$ and $x = x'$, and 0 otherwise.

Supposing a column vector y_X of realized measurements is available for some set $X \triangleq \bigcup_{i=1}^M X_i$ of tuples of observed locations and their corresponding measurement types where $X_i \subset D_i^+$, a CMOGP regression model can exploit these observations to provide a Gaussian predictive distribution $\mathcal{N}(\mu_{Z|X}, \Sigma_{ZZ|X})$ of the measurements for any set $Z \subseteq D^+ \setminus X$ of tuples of unobserved locations and their corresponding measurement types with the following *posterior* mean vector and covariance matrix:

$$\begin{aligned} \mu_{Z|X} &\triangleq \mu_Z + \Sigma_{ZX} \Sigma_{XX}^{-1} (y_X - \mu_X) \\ \Sigma_{ZZ|X} &\triangleq \Sigma_{ZZ} - \Sigma_{ZX} \Sigma_{XX}^{-1} \Sigma_{XZ} \end{aligned} \quad (2)$$

where $\Sigma_{AA'} \triangleq (\sigma_{\langle x, i \rangle \langle x', j \rangle})_{\langle x, i \rangle \in A, \langle x', j \rangle \in A'}$ and $\mu_A \triangleq (\mu_{\langle x, i \rangle})_{\langle x, i \rangle \in A}^\top$ for any $A, A' \subseteq D^+$.

Sparse CMOGP Regression. A limitation of the CMOGP model is its poor scalability in the number $|X|$ of observations: Computing its Gaussian predictive distribution (2) requires inverting Σ_{XX} , which incurs $\mathcal{O}(|X|^3)$ time. To improve its scalability, a unifying framework of sparse CMOGP regression models such as the deterministic training conditional, fully independent training conditional, and *partially independent training conditional* (PITC) approximations (Álvarez and Lawrence 2011) exploit a vector $L_U \triangleq (L(x))_{x \in U}^\top$ of inducing measurements for some small set $U \subset D$ of inducing locations (i.e., $|U| \ll |D|$) to approximate each measurement $Y_{\langle x, i \rangle}$:

$$Y_{\langle x, i \rangle} \approx \int_{x' \in D} K_i(x - x') \mathbb{E}[L(x') | L_U] dx' + \varepsilon_i.$$

They also share two structural properties that can be exploited for deriving our active learning criterion and in turn an efficient approximation algorithm: **(P1)** Measurements of different types (i.e., $Y_{D_i^+}$ and $Y_{D_j^+}$ for $i \neq j$) are conditionally independent given L_U , and **(P2)** Y_X and Y_Z are conditionally independent given L_U . PITC will be used as the sparse CMOGP regression model in our work here since the others in the unifying framework impose further assumptions. With the above structural properties, PITC can utilize

the observations to provide a Gaussian predictive distribution $\mathcal{N}(\mu_{Z|X}^{\text{PITC}}, \Sigma_{ZZ|X}^{\text{PITC}})$ where

$$\begin{aligned}\mu_{Z|X}^{\text{PITC}} &\triangleq \mu_Z + \Gamma_{ZX}(\Gamma_{XX} + \Lambda_X)^{-1}(y_X - \mu_X) \\ \Sigma_{ZZ|X}^{\text{PITC}} &\triangleq \Gamma_{ZZ} + \Lambda_Z - \Gamma_{ZX}(\Gamma_{XX} + \Lambda_X)^{-1}\Gamma_{XZ}\end{aligned}\quad (3)$$

such that $\Gamma_{AA'} \triangleq \Sigma_{AU}\Sigma_{UU}^{-1}\Sigma_{UA'}$ for any $A, A' \subseteq D^+$, Σ_{AU} (Σ_{UU}) is a covariance matrix with covariance components $\sigma_{\langle x, i \rangle x'} = \sigma_{s_i} \mathcal{N}(x - x' | 0, P_0^{-1} + P_i^{-1})$ for all $\langle x, i \rangle \in A$ and $x' \in U$ ($\sigma_{xx'}$ for all $x, x' \in U$), $\Sigma_{UA'}$ is the transpose of $\Sigma_{A'U}$, and Λ_A is a block-diagonal matrix constructed from the M diagonal blocks of $\Sigma_{AA|U} \triangleq \Sigma_{AA} - \Gamma_{AA}$ for any $A \subseteq D^+$, each of which is a matrix $\Sigma_{A_i A_i | U}$ for $i = 1, \dots, M$ where $A_i \subseteq D_i^+$ and $A \triangleq \bigcup_{i=1}^M A_i$. Note that computing (3) does not require the inducing locations U to be observed. Also, the covariance matrix Σ_{XX} in (2) is approximated by a reduced-rank matrix Γ_{XX} summed with the resulting sparsified residual matrix Λ_X . So, by using the matrix inversion lemma to invert the approximated covariance matrix $\Gamma_{XX} + \Lambda_X$ and applying some algebraic manipulations, computing (3) incurs $\mathcal{O}(|X|(|U|^2 + (|X|/M)^2))$ time (Álvarez and Lawrence 2011) in the case of $|U| \leq |X|$ and evenly distributed observations among all M types.

3 Active Learning of CMOGP

Recall from Section 1 that the entropy criterion can be used to measure the predictive uncertainty of the unobserved areas of a target phenomenon. Using the CMOGP model (Section 2), the Gaussian posterior joint entropy (i.e., predictive uncertainty) of the measurements Y_Z for any set $Z \subseteq D^+ \setminus X$ of tuples of unobserved locations and their corresponding measurement types can be expressed in terms of its posterior covariance matrix $\Sigma_{ZZ|X}$ (2) which is independent of the realized measurements y_X :

$$H(Y_Z | Y_X) \triangleq \frac{1}{2} \log(2\pi e)^{|Z|} |\Sigma_{ZZ|X}|.$$

Let index t denote the type of measurements of the target phenomenon⁵. Then, active learning of a CMOGP model involves selecting an optimal set $X^* \triangleq \bigcup_{i=1}^M X_i^*$ of N tuples (i.e., sampling budget) of sampling locations and their corresponding measurement types to be observed that minimize the posterior joint entropy of type t measurements at the remaining unobserved locations of the target phenomenon:

$$X^* \triangleq \arg \min_{X: |X|=N} H(Y_{V_t \setminus X_t} | Y_X) \quad (4)$$

where $V_t \subset D_t^+$ is a finite set of tuples of candidate sampling locations of the target phenomenon and their corresponding measurement type t available to be selected for observation. However, evaluating the $H(Y_{V_t \setminus X_t} | Y_X)$ term in (4) incurs $\mathcal{O}(|V_t|^3 + N^3)$ time, which is prohibitively expensive when the target phenomenon is spanned by a large number $|V_t|$ of candidate sampling locations. If auxiliary types of phenomena are missing or ignored (i.e.,

⁵Our proposed algorithm can be extended to handle multiple types of target phenomena, as demonstrated in Section 5.

$M = 1$), then such a computational difficulty can be eased instead by solving the well-known *maximum entropy sampling* (MES) problem (Shewry and Wynn 1987): $X_t^* = \arg \max_{X_t: |X_t|=N} H(Y_{X_t})$ which can be proven to be equivalent to (4) by using the chain rule for entropy $H(Y_{V_t}) = H(Y_{X_t}) + H(Y_{V_t \setminus X_t} | Y_{X_t})$ and noting that $H(Y_{V_t})$ is a constant. Evaluating the $H(Y_{X_t})$ term in MES incurs $\mathcal{O}(|X_t|^3)$ time, which is independent of $|V_t|$. Such an equivalence result can in fact be extended and applied to minimizing the predictive uncertainty of *all* M types of coexisting phenomena, as exploited by multivariate spatial sampling algorithms (Bueso et al. 1999; Le, Sun, and Zidek 2003):

$$\arg \max_{X: |X|=N} H(Y_X) = \arg \min_{X: |X|=N} H(Y_{V \setminus X} | Y_X), \quad (5)$$

where $V \triangleq \bigcup_{i=1}^M V_i$ and V_i is defined in a similar manner to V_t but for measurement type $i \neq t$. This equivalence result (5) also follows from the chain rule for entropy $H(Y_V) = H(Y_X) + H(Y_{V \setminus X} | Y_X)$ and the fact that $H(Y_V)$ is a constant. Unfortunately, it is not straightforward to derive such an equivalence result for our active learning problem (4) in which a target phenomenon of interest coexists with auxiliary types of phenomena (i.e., $M > 1$): If we consider maximizing $H(Y_X)$ or $H(Y_{X_t})$, then it is no longer equivalent to minimizing $H(Y_{V_t \setminus X_t} | Y_X)$ (4) as the sum of the two entropy terms is not necessarily a constant.

Exploiting Sparse CMOGP Model Structure. We derive a new equivalence result by considering instead a constant entropy $H(Y_{V_t} | L_U)$ that is conditioned on the inducing measurements L_U used in sparse CMOGP regression models (Section 2). Then, by using the chain rule for entropy and structural property **P2** shared by sparse CMOGP regression models in the unifying framework (Álvarez and Lawrence 2011) described in Section 2, (4) can be proven (see Appendix A) to be equivalent to

$$X^* \triangleq \arg \max_{X: |X|=N} H(Y_{X_t} | L_U) - I(L_U; Y_{V_t \setminus X_t} | Y_X) \quad (6)$$

where

$$I(L_U; Y_{V_t \setminus X_t} | Y_X) \triangleq H(L_U | Y_X) - H(L_U | Y_{X \cup V_t \setminus X_t}) \quad (7)$$

is the conditional mutual information between L_U and $Y_{V_t \setminus X_t}$ given Y_X . Our novel active learning criterion in (6) exhibits an interesting exploration-exploitation trade-off: The inducing measurements L_U can be viewed as latent structure of the sparse CMOGP model to induce conditional independence properties **P1** and **P2**. So, on one hand, maximizing the $H(Y_{X_t} | L_U)$ term aims to select tuples X_t of locations with the most uncertain measurements Y_{X_t} of the target phenomenon and their corresponding type t to be observed given the latent model structure L_U (i.e., exploitation). On the other hand, minimizing the $I(L_U; Y_{V_t \setminus X_t} | Y_X)$ term (7) aims to select tuples X to be observed (i.e., possibly of measurement types $i \neq t$) so as to rely less on measurements $Y_{V_t \setminus X_t}$ of type t at the remaining unobserved locations of the target phenomenon to infer latent model structure L_U (i.e., exploration) since $Y_{V_t \setminus X_t}$ won't be sampled.

Supposing $|U| \leq |V_t|$, evaluating our new active learning criterion in (6) incurs $\mathcal{O}(|U|^3 + N^3)$ time for every $X \subset V$ and a *one-off* cost of $\mathcal{O}(|V_t|^3)$ time (Appendix B). In contrast, computing the original criterion in (4) requires $\mathcal{O}(|V_t|^3 + N^3)$ time for every $X \subset V$, which is more costly, especially when the number N of selected observations is much less than the number $|V_t|$ of candidate sampling locations of the target phenomenon due to, for example, a tight sampling budget or a large sampling domain that usually occurs in practice. The trick to achieving such a computational advantage can be inherited by our approximation algorithm to be described next.

4 Approximation Algorithm

Our novel active learning criterion in (6), when optimized, still suffers from poor scalability in the number N of selected observations (i.e., sampling budget) like the old criterion in (4) because it involves evaluating a prohibitively large number of candidate selections of sampling locations and their corresponding measurement types (i.e., exponential in N). However, unlike the old criterion, it is possible to devise an efficient approximation algorithm with a theoretical performance guarantee to optimize our new criterion, which is the main contribution of our work in this paper.

The key idea of our proposed approximation algorithm is to greedily select the next tuple of sampling location and its corresponding measurement type to be observed that maximally increases our criterion in (6), and iterate this till N tuples are selected for observation. Specifically, let

$$F(X) \triangleq H(Y_{X_t}|L_U) - I(L_U; Y_{V_t \setminus X_t}|Y_X) + I(L_U; Y_{V_t}) \quad (8)$$

denote our active learning criterion in (6) augmented by a positive constant $I(L_U; Y_{V_t})$ to make $F(X)$ non-negative. Such an additive constant $I(L_U; Y_{V_t})$ is simply a technical necessity for proving the performance guarantee and does not affect the outcome of the optimal selection (i.e., $X^* = \arg \max_{X: |X|=N} F(X)$). Then, our approximation algorithm greedily selects the next tuple $\langle x, i \rangle$ of sampling location x and its corresponding measurement type i that maximizes $F(X \cup \{\langle x, i \rangle\}) - F(X)$:

$$\begin{aligned} \langle x, i \rangle^+ &\triangleq \arg \max_{\langle x, i \rangle \in V \setminus X} F(X \cup \{\langle x, i \rangle\}) - F(X) \\ &= \arg \max_{\langle x, i \rangle \in V \setminus X} H(Y_{\langle x, i \rangle}|Y_X) - \delta_i H(Y_{\langle x, i \rangle}|Y_{X \cup V_t \setminus X_t}) \end{aligned} \quad (9)$$

where δ_i is a Kronecker delta of value 0 if $i = t$, and 1 otherwise. The derivation of (9) is in Appendix C. Our algorithm updates $X \leftarrow X \cup \{\langle x, i \rangle^+\}$ and iterates the greedy selection (9) and the update till $|X| = N$ (i.e., sampling budget is depleted). The intuition to understanding (9) is that our algorithm has to choose between observing a sampling location with the most uncertain measurement (i.e., $H(Y_{\langle x, t \rangle}|Y_X)$) of the target phenomenon (i.e., type t) vs. that for an auxiliary type $i \neq t$ inducing the largest reduction in predictive uncertainty of the measurements at the remaining unobserved locations of the target phenomenon since $H(Y_{\langle x, i \rangle}|Y_X) - H(Y_{\langle x, i \rangle}|Y_{X \cup V_t \setminus X_t}) = I(Y_{\langle x, i \rangle}; Y_{V_t \setminus X_t}|Y_X) = H(Y_{V_t \setminus X_t}|Y_X) - H(Y_{V_t \setminus X_t}|Y_{X \cup \{\langle x, i \rangle\}})$.

It is also interesting to figure out whether our approximation algorithm may avoid selecting tuples of a certain auxiliary type $i \neq t$ and formally analyze the conditions under which it will do so, as elucidated in the following result:

Proposition 1 *Let $V_{-t} \triangleq \bigcup_{i \neq t} V_i$, $X_{-t} \triangleq \bigcup_{i \neq t} X_i$, $\rho_i \triangleq \sigma_{s_i}^2 / \sigma_{n_i}^2$, and $R(\langle x, i \rangle, V_t \setminus X_t) \triangleq \sum_{\langle x', t \rangle \in V_t \setminus X_t} \mathcal{N}(x - x' | \underline{0}, P_0^{-1} + P_i^{-1} + P_t^{-1})^2$. Assuming absence of suppressor variables, $H(Y_{\langle x, i \rangle}|Y_X) - H(Y_{\langle x, i \rangle}|Y_{X \cup V_t \setminus X_t}) \leq 0.5 \log(1 + 4\rho_t \rho_i R(\langle x, i \rangle, V_t \setminus X_t))$ for any $\langle x, i \rangle \in V_{-t} \setminus X_{-t}$.*

Its proof (Appendix D) relies on the following assumption of the absence of suppressor variables which holds in many practical cases (Das and Kempe 2008): Conditioning does not make $Y_{\langle x, i \rangle}$ and $Y_{\langle x', t \rangle}$ more correlated for any $\langle x, i \rangle \in V_{-t} \setminus X_{-t}$ and $\langle x', t \rangle \in V_t \setminus X_t$. Proposition 1 reveals that when the signal-to-noise ratio ρ_i of auxiliary type i is low (e.g., poor-quality measurements due to high noise) and/or the cross correlation (1) between measurements of the target phenomenon and auxiliary type i is small due to low $\sigma_{s_i}^2 \sigma_{s_t}^2 R(\langle x, i \rangle, V_t \setminus X_t)$, our greedy criterion in (9) returns a small value, hence causing our algorithm to avoid selecting tuples of auxiliary type i .

Theorem 1 (Time Complexity) *Our approximation algorithm incurs $\mathcal{O}(N(|V||U|^2 + N^3) + |V_t|^3)$ time.*

Its proof is in Appendix E. So, our approximation algorithm only incurs quartic time in the number N of selected observations and cubic time in the number $|V_t|$ of candidate sampling locations of the target phenomenon.

Performance Guarantee. To theoretically guarantee the performance of our approximation algorithm, we will first motivate the need to define a relaxed notion of *submodularity*. A submodular set function exhibits a natural diminishing returns property: When adding an element to its input set, the increment in its function value decreases with a larger input set. To maximize a nondecreasing and submodular set function, the work of Nemhauser, Wolsey, and Fisher (1978) has proposed a greedy algorithm guaranteeing a $(1 - 1/e)$ -factor approximation of that achieved by the optimal input set.

The main difficulty in proving the submodularity of $F(X)$ (8) lies in its mutual information term being conditioned on X . Some works (Krause and Guestrin 2005; Renner and Maurer 2002) have shown the submodularity of such conditional mutual information by imposing conditional independence assumptions (e.g., Markov chain). In practice, these strong assumptions (e.g., $Y_A \perp Y_{\langle x, t \rangle} | Y_{V_t \setminus X_t}$ for any $A \subseteq X$ and $\langle x, i \rangle \in V_{-t} \setminus X_{-t}$) severely violate the correlation structure of multiple types of coexisting phenomena and are an overkill: The correlation structure can in fact be preserved to a fair extent by relaxing these assumptions, which consequently entails a relaxed form of submodularity of $F(X)$ (8); a performance guarantee similar to that of Nemhauser, Wolsey, and Fisher (1978) can then be derived for our approximation algorithm.

Definition 1 *A function $G : 2^B \rightarrow \mathbb{R}$ is ϵ -submodular if $G(A' \cup \{a\}) - G(A') \leq G(A \cup \{a\}) - G(A) + \epsilon$ for any $A \subseteq A' \subseteq B$ and $a \in B \setminus A'$.*

Lemma 1 Let $\sigma_{n^*}^2 \triangleq \min_{i \in \{1, \dots, M\}} \sigma_{n_i}^2$. Given $\epsilon_1 \geq 0$, if

$$\sum_{\langle x, i \rangle \langle x, i \rangle | \tilde{X} \cup V_t \setminus X_t}^{\text{PITC}} - \sum_{\langle x, i \rangle \langle x, i \rangle | X \cup V_t \setminus X_t}^{\text{PITC}} \leq \epsilon_1 \quad (10)$$

for any $\tilde{X} \subseteq X$ and $\langle x, i \rangle \in V_t \setminus X_t$, then $F(X)$ is ϵ -submodular where $\epsilon = 0.5 \log(1 + \epsilon_1/\sigma_{n^*}^2)$.

Its proof is in Appendix F. Note that (10) relaxes the above example of conditional independence assumption (i.e., assuming $\epsilon_1 = 0$) to one which allows $\epsilon_1 > 0$. In practice, ϵ_1 is expected to be small: Since further conditioning monotonically decreases a posterior variance (Xu et al. 2014), an expected large set $V_t \setminus X_t$ of tuples of remaining unobserved locations of the target phenomenon tends to be informative enough to make $\sum_{\langle x, i \rangle \langle x, i \rangle | \tilde{X} \cup V_t \setminus X_t}^{\text{PITC}}$ small and hence the non-negative variance reduction term and ϵ_1 in (10) small.

Furthermore, (10) with a given small ϵ_1 can be realized by controlling the discretization of the domain of candidate sampling locations. For example, by refining the discretization of V_t (i.e., increasing $|V_t|$), the variance reduction term in (10) decreases because it has been shown in (Das and Kempe 2008) to be submodular in many practical cases. We give another example in Lemma 2 to realize (10) by controlling the discretization such that every pair of selected observations are sufficiently far apart.

It is easy to derive that $F(\emptyset) = 0$. The ‘‘information never hurts’’ bound for entropy (Cover and Thomas 1991) entails a nondecreasing $F(X)$: $F(X \cup \{\langle x, i \rangle\}) - F(X) = H(Y_{\langle x, i \rangle} | Y_X) - \delta_i H(Y_{\langle x, i \rangle} | Y_{X \cup V_t \setminus X_t}) \geq H(Y_{\langle x, i \rangle} | Y_X) - H(Y_{\langle x, i \rangle} | Y_{X \cup V_t \setminus X_t}) \geq 0$. The first inequality requires $\sigma_{n^*}^2 \geq (2\pi e)^{-1}$ so that $H(Y_{\langle x, i \rangle} | Y_A) = 0.5 \log 2\pi e \sum_{\langle x, i \rangle \langle x, i \rangle | A}^{\text{PITC}} \geq 0.5 \log 2\pi e \sigma_{n^*}^2 \geq 0$,⁶ which is reasonable in practice due to ubiquitous noise. Combining this result with Lemma 1 yields the performance guarantee:

Theorem 2 Given $\epsilon_1 \geq 0$, if (10) holds, then our approximation algorithm is guaranteed to select X s.t. $F(X) \geq (1 - 1/e)(F(X^*) - N\epsilon)$ where $\epsilon = 0.5 \log(1 + \epsilon_1/\sigma_{n^*}^2)$.

Its proof (Appendix G) is similar to that of the well-known result of Nemhauser, Wolsey, and Fisher (1978) except for exploiting ϵ -submodularity of $F(X)$ in Lemma 1 instead of submodularity.

Finally, we present a discretization scheme that satisfies (10): Let ω be the smallest discretization width of V_i for $i = 1, \dots, M$. Construct a new set $V^- \subset V$ of tuples of candidate sampling locations and their corresponding measurement types such that every pair of tuples are at least a distance of $p\omega$ apart for some $p > 0$; each candidate location thus has only one corresponding type. Such a construction V^- constrains our algorithm to select observations sparsely across the spatial domain so that any $\langle x, i \rangle \in V_t \setminus X_t$ has sufficiently many neighboring tuples of remaining unobserved locations of the target phenomenon from $V_t \setminus X_t$ to keep $\sum_{\langle x, i \rangle \langle x, i \rangle | \tilde{X} \cup V_t \setminus X_t}^{\text{PITC}}$ small and hence the variance reduction term and ϵ_1 in (10) small. Our previous theoretical results still hold if V^- is used instead of V . The result below gives the minimum value of p to satisfy (10):

⁶ $\sum_{\langle x, i \rangle \langle x, i \rangle | A}^{\text{PITC}} \geq \sigma_{n^*}^2$ is proven in Lemma 3 in Appendix D.

Lemma 2 Let $\sigma_{s^*}^2 \triangleq \max_{i \in \{1, \dots, M\}} \sigma_{s_i}^2$, ℓ be the largest first diagonal component of $P_0^{-1} + P_i^{-1} + P_j^{-1}$ for all $i, j = 1, \dots, M$, and $\xi \triangleq \exp(-\omega^2/(2\ell))$. Given $\epsilon_1 > 0$ and assuming absence of suppressor variables, if

$$p^2 > \log \left\{ \frac{1}{2\sigma_{s^*}^2} \min \left(\frac{\sigma_{n^*}^2}{N}, \frac{1}{2} \left(\sqrt{\epsilon_1^2 + \frac{4\epsilon_1 \sigma_{n^*}^2}{N}} - \epsilon_1 \right) \right) \right\} / \log \xi,$$

then (10) holds. See Appendix H for its proof.

5 Experiments and Discussion

This section evaluates the predictive performance of our approximation algorithm (m-Greedy) empirically on three real-world datasets: (a) *Jura* dataset (Goovaerts 1997) contains concentrations of 7 heavy metals collected at 359 locations in a Swiss Jura region; (b) *Gilgai* dataset (Webster 1977) contains electrical conductivity and chloride content generated from a line transect survey of 365 locations of Gilgai territory in New South Wales, Australia; and (c) *indoor environmental quality* (IEQ) dataset (Bodik et al. 2004) contains temperature ($^\circ\text{F}$) and light (Lux) readings taken by 43 temperature sensors and 41 light sensors deployed in the Intel Berkeley Research lab. The sampling locations for the Jura and IEQ datasets are shown in Appendix I.

The performance of m-Greedy is compared to that of the (a) maximum variance/entropy (m-Var) algorithm which greedily selects the next location and its corresponding measurement type with maximum posterior variance/entropy in each iteration; and (b) greedy maximum entropy (s-Var) (Shewry and Wynn 1987) and mutual information (s-MI) (Krause, Singh, and Guestrin 2008) sampling algorithms for gathering observations *only* from the target phenomenon.

For all experiments, k-means is used to select inducing locations U by clustering all possible locations available to be selected for observation into $|U|$ clusters such that each cluster center corresponds to an element of U . The hyper-parameters (i.e., $\sigma_{s_i}^2$, $\sigma_{n_i}^2$, P_0 and P_i for $i = 1, \dots, M$) of MOGP and single-output GP models are learned using the data via maximum likelihood estimation (Álvarez and Lawrence 2011). For each dataset, observations (i.e., 100 for Jura and Gilgai datasets and 10 for IEQ dataset) of type t are randomly selected to form the test set T ; the tuples of candidate sampling locations and corresponding type t therefore become less than that of auxiliary types. The *root mean squared error* (RMSE) metric $\sqrt{|T|^{-1} \sum_{x \in T} (y_{\langle x, t \rangle} - \mu_{\langle x, t \rangle | X})^2}$ is used to evaluate the performance of the tested algorithms. All experimental results are averaged over 50 random test sets. For a fair comparison, the measurements of all types are normalized before using them for training, prediction, and active learning.

Jura Dataset. Three types of correlated lg-Cd, Ni, and lg-Zn measurements are used in this experiment; we take the log of Cd and Zn measurements to remove their strong skewness, as proposed as a standard statistical practice in (Webster and Oliver 2007). The measurement types with the smallest and largest signal-to-noise ratios (respectively, lg-Cd and Ni; see Appendix J) are each set as type t .

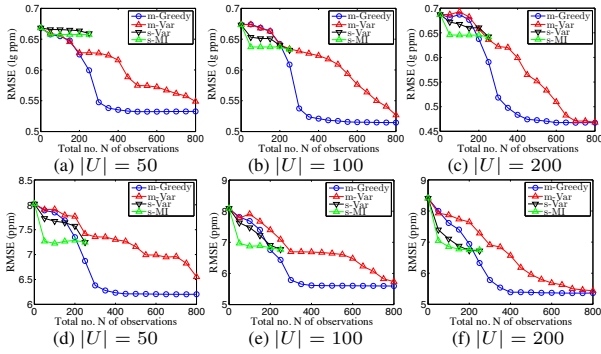


Figure 1: Graphs of RMSEs vs. no. N of observations with (a-c) lg-Cd and (d-f) Ni as type t and varying no. $|U| = 50, 100, 200$ of inducing locations for Jura dataset.

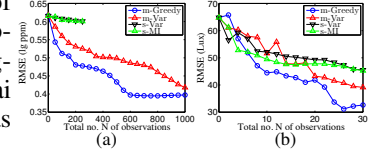
Figs. 1a-c and 1d-f show, respectively, results of the tested algorithms with lg-Cd and Ni as type t . It can be observed that the RMSE of m-Greedy decreases more rapidly than that of m-Var, especially when observations of auxiliary types are selected after about $N = 200$. This is because our algorithm selects observations of auxiliary types that induce the largest reduction in predictive uncertainty of the measurements at the remaining unobserved locations of the target phenomenon (Section 4). In contrast, m-Var may select observations that reduce the predictive uncertainty of auxiliary types of phenomena, which does not directly achieve the aim of our active learning problem. With increasing $|U|$, both m-Greedy and m-Var reach smaller RMSEs, but m-Greedy can achieve this faster with much less observations. As shown in Figs. 1a-f, m-Greedy performs much better than s-Var and s-MI, which means observations of correlated auxiliary types can indeed be used to improve the prediction of the target phenomenon. Finally, by comparing the results between Figs. 1a-c and 1d-f, the RMSE of m-Greedy with Ni as type t decreases faster than that with lg-Cd as type t , especially in the beginning (i.e., $N \leq 200$) due to higher-quality Ni measurements (i.e., larger signal-to-noise ratio).

Gilgai Dataset. In this experiment, the lg-Cl contents at depth 0-10cm and 30-40cm are used jointly as two types of target phenomena while the log of electrical conductivity, which is easier to measure at these depths, is used as the auxiliary type. Fig. 2a shows results of the average RMSE over the two lg-Cl types with $|U| = 100$. Similar to the results of the Jura dataset, with two types of target phenomena, the RMSE of m-Greedy still decreases more rapidly with increasing N than that of m-Var and achieves a much smaller RMSE than that of s-Var and s-MI; the results of s-Var and s-MI are also averaged over two independent single-output GP predictions of lg-Cl content at the two depths.

IEQ Dataset. Fig. 2b shows results with light as type t and $|U| = 40$. The observations are similar to that of the Jura and Gilgai datasets: RMSE of m-Greedy decreases faster than that of the other algorithms. More importantly, with the same number of observations, m-Greedy achieves much smaller RMSE than s-Var and s-MI that can sample only from the target phenomenon. This is because m-Greedy selects observations of the auxiliary type (i.e., temperature) that are less

noisy ($\sigma_{n_i}^2 = 0.13$) than that of light ($\sigma_{n_t}^2 = 0.23$), which demonstrates its advantage over s-Var and s-MI when type t measurements are noisy (e.g., due to poor-quality sensors).

Figure 2: Graphs of RMSEs vs. no. N of observations with (a) lg-Cl as types t for Gilgai dataset and (b) light as type t for IEQ dataset.



6 Related Work

Existing works on active learning with multiple output measurement types are not driven by the MOGP model and have not formally characterized the cross-correlation structure between different types of phenomena: Some spatial sampling algorithms (Bueso, Angulo, and Alonso 1998; Angulo and Bueso 2001) have simply modeled the auxiliary phenomenon as a noisy perturbation of the target phenomenon that is assumed to be latent, which differs from our work here. *Multi-task active learning* (MTAL) and *active transfer learning* (ATL) algorithms have considered the prediction of each type of phenomenon as one task and used the auxiliary tasks to help learn the target task. But, the MTAL algorithm of Zhang (2010) requires relations between different classification tasks to be manually specified, which is highly non-trivial to achieve in practice and not applicable to MOGP regression. Some ATL and active learning algorithms (Roth and Small 2006; Zhao et al. 2013) have used active learning criteria (e.g., margin-based criterion) specific to their classification or recommendation tasks that cannot be readily tailored to MOGP regression.

7 Conclusion

This paper describes a novel efficient algorithm for active learning of a MOGP model. To resolve the issue of poor scalability in optimizing the conventional entropy criterion, we exploit a structure common to a unifying framework of sparse MOGP models for deriving a novel active learning criterion (6). Then, we exploit the ϵ -submodularity property of our new criterion (Lemma 1) for devising a polynomial-time approximation algorithm (9) that guarantees a constant-factor approximation of that achieved by the optimal set of selected observations (Theorem 2). Empirical evaluation on three real-world datasets shows that our approximation algorithm m-Greedy outperforms existing algorithms for active learning of MOGP and single-output GP models, especially when measurements of the target phenomenon are more noisy than that of the auxiliary types. For our future work, we plan to extend our approach by generalizing non-myopic active learning (Cao, Low, and Dolan 2013; Hoang et al. 2014; Ling, Low, and Jaillot 2016; Low, Dolan, and Khosla 2009; Low, Dolan, and Khosla 2008; Low, Dolan, and Khosla 2011) of single-output GPs to that of MOGPs and improving its scalability to big data through parallelization (Chen et al. 2013; Low et al. 2015), online learning (Xu et al. 2014), and stochastic variational inference (Hoang, Hoang, and Low 2015).

Acknowledgments. This research was carried out at the

SeSaMe Centre. It is supported by Singapore NRF under its IRC@SG Funding Initiative and administered by IDMPO.

References

- [Álvarez and Lawrence 2011] Álvarez, M. A., and Lawrence, N. D. 2011. Computationally efficient convolved multiple output Gaussian processes. *JMLR* 12:1459–1500.
- [Angulo and Bueso 2001] Angulo, J. M., and Bueso, M. C. 2001. Random perturbation methods applied to multivariate spatial sampling design. *Environmetrics* 12(7):631–646.
- [Bodík et al. 2004] Bodík, P.; Guestrin, C.; Hong, W.; Madden, S.; Paskin, M.; and Thibaux, R. 2004. <http://www.select.cs.cmu.edu/data/labapp3/index.html>.
- [Bonilla, Chai, and Williams 2008] Bonilla, E. V.; Chai, K. M. A.; and Williams, C. K. 2008. Multi-task Gaussian process prediction. In *Proc. NIPS*.
- [Bueso, Angulo, and Alonso 1998] Bueso, M. C.; Angulo, J. M.; and Alonso, F. J. 1998. A state-space model approach to optimum spatial sampling design based on entropy. *Environmental and Ecological Statistics* 5(1):29–44.
- [Bueso et al. 1999] Bueso, M. C.; Angulo, J. M.; Cruz-Sanjulián, J.; and García-Aróstegui, J. L. 1999. Optimal spatial sampling design in a multivariate framework. *Math. Geology* 31(5):507–525.
- [Cao, Low, and Dolan 2013] Cao, N.; Low, K. H.; and Dolan, J. M. 2013. Multi-robot informative path planning for active sensing of environmental phenomena: A tale of two algorithms. In *Proc. AAMAS*.
- [Chen et al. 2012] Chen, J.; Low, K. H.; Tan, C. K.-Y.; Oran, A.; Jaillet, P.; Dolan, J. M.; and Sukhatme, G. S. 2012. Decentralized data fusion and active sensing with mobile sensors for modeling and predicting spatiotemporal traffic phenomena. In *Proc. UAI*, 163–173.
- [Chen et al. 2013] Chen, J.; Cao, N.; Low, K. H.; Ouyang, R.; Tan, C. K.-Y.; and Jaillet, P. 2013. Parallel Gaussian process regression with low-rank covariance matrix approximations. In *Proc. UAI*, 152–161.
- [Chen et al. 2015] Chen, J.; Low, K. H.; Jaillet, P.; and Yao, Y. 2015. Gaussian process decentralized data fusion and active sensing for spatiotemporal traffic modeling and prediction in mobility-on-demand systems. *IEEE Trans. Autom. Sci. Eng.* 12:901–921.
- [Chen, Low, and Tan 2013] Chen, J.; Low, K. H.; and Tan, C. K.-Y. 2013. Gaussian process-based decentralized data fusion and active sensing for mobility-on-demand system. In *Proc. RSS*.
- [Cover and Thomas 1991] Cover, T., and Thomas, J. 1991. *Elements of Information Theory*. John Wiley & Sons, Inc.
- [Das and Kempe 2008] Das, A., and Kempe, D. 2008. Algorithms for subset selection in linear regression. In *Proc. STOC*, 45–54.
- [Golub and Van Loan 1996] Golub, G. H., and Van Loan, C.-F. 1996. *Matrix Computations*. Johns Hopkins Univ. Press, 3rd edition.
- [Goovaerts 1997] Goovaerts, P. 1997. *Geostatistics for Natural Resources Evaluation*. Oxford Univ. Press.
- [Hoang et al. 2014] Hoang, T. N.; Low, K. H.; Jaillet, P.; and Kankanhalli, M. 2014. Nonmyopic ϵ -Bayes-optimal active learning of Gaussian processes. In *Proc. ICML*, 739–747.
- [Hoang, Hoang, and Low 2015] Hoang, T. N.; Hoang, Q. M.; and Low, K. H. 2015. A unifying framework of anytime sparse Gaussian process regression models with stochastic variational inference for big data. In *Proc. ICML*, 569–578.
- [Krause and Golovin 2014] Krause, A., and Golovin, D. 2014. Submodular function maximization. In Bordeaux, L.; Hamadi, Y.; and Kohli, P., eds., *Tractability: Practical Approaches to Hard Problems*. Cambridge Univ. Press. 71–104.
- [Krause and Guestrin 2005] Krause, A., and Guestrin, C. 2005. Near-optimal nonmyopic value of information in graphical models. In *Proc. UAI*.
- [Krause, Singh, and Guestrin 2008] Krause, A.; Singh, A.; and Guestrin, C. 2008. Near-optimal sensor placements in Gaussian processes: Theory, efficient algorithms and empirical studies. *JMLR* 9:235–284.
- [Le, Sun, and Zidek 2003] Le, N. D.; Sun, L.; and Zidek, J. V. 2003. Designing networks for monitoring multivariate environmental fields using data with monotone pattern. Technical Report #2003-5, Statistical and Applied Mathematical Sciences Institute, Research Triangle Park, NC.
- [Ling, Low, and Jaillet 2016] Ling, C. K.; Low, K. H.; and Jaillet, P. 2016. Gaussian process planning with Lipschitz continuous reward functions: Towards unifying Bayesian optimization, active learning, and beyond. In *Proc. AAAI*.
- [Low et al. 2015] Low, K. H.; Yu, J.; Chen, J.; and Jaillet, P. 2015. Parallel Gaussian process regression for big data: Low-rank representation meets Markov approximation. In *Proc. AAAI*.
- [Low, Dolan, and Khosla 2008] Low, K. H.; Dolan, J. M.; and Khosla, P. 2008. Adaptive multi-robot wide-area exploration and mapping. In *Proc. AAMAS*, 23–30.
- [Low, Dolan, and Khosla 2009] Low, K. H.; Dolan, J. M.; and Khosla, P. 2009. Information-theoretic approach to efficient adaptive path planning for mobile robotic environmental sensing. In *Proc. ICAPS*.
- [Low, Dolan, and Khosla 2011] Low, K. H.; Dolan, J. M.; and Khosla, P. 2011. Active Markov information-theoretic path planning for robotic environmental sensing. In *Proc. AAMAS*, 753–760.
- [Nemhauser, Wolsey, and Fisher 1978] Nemhauser, G. L.; Wolsey, L. A.; and Fisher, M. L. 1978. An analysis of approximations for maximizing submodular set functions - I. *Mathematical Programming* 14(1):265–294.
- [Petersen and Pedersen 2012] Petersen, K. B., and Pedersen, M. S. 2012. *The Matrix Cookbook*.
- [Renner and Maurer 2002] Renner, R., and Maurer, U. 2002. About the mutual (conditional) information. In *Proc. IEEE ISIT*.

- [Roth and Small 2006] Roth, D., and Small, K. 2006. Margin-based active learning for structured output spaces. In *Proc. ECML*, 413–424.
- [Shewry and Wynn 1987] Shewry, M. C., and Wynn, H. P. 1987. Maximum entropy sampling. *J. Applied Stat.* 14(2):165–170.
- [Skolidis 2012] Skolidis, G. 2012. *Transfer Learning with Gaussian Processes*. Ph.D. Dissertation, University of Edinburgh.
- [Stewart and Sun 1990] Stewart, G. W., and Sun, J.-G. 1990. *Matrix Perturbation Theory*. Academic Press.
- [Teh and Seeger 2005] Teh, Y. W., and Seeger, M. 2005. Semiparametric latent factor models. In *Proc. AISTATS*, 333–340.
- [Webster and Oliver 2007] Webster, R., and Oliver, M. 2007. *Geostatistics for Environmental Scientists*. John Wiley & Sons, Inc., 2nd edition.
- [Webster 1977] Webster, R. 1977. Spectral analysis of Gilgai soil. *Australian Journal of Soil Research* 15(3):191–204.
- [Xu et al. 2014] Xu, N.; Low, K. H.; Chen, J.; Lim, K. K.; and Ozgul, E. B. 2014. GP-Localize: Persistent mobile robot localization using online sparse Gaussian process observation model. In *Proc. AAAI*, 2585–2592.
- [Zhang et al. 2016] Zhang, Y.; Hoang, T. N.; Low, K. H.; and Kankanhalli, M. 2016. Near-optimal active learning of multi-output Gaussian processes. In *Proc. AAAI*.
- [Zhang 2010] Zhang, Y. 2010. Multi-task active learning with output constraints. In *Proc. AAAI*, 667–672.
- [Zhao et al. 2013] Zhao, L.; Pan, S. J.; Xiang, E. W.; Zhong, E.; Lu, Z.; and Yang, Q. 2013. Active transfer learning for cross-system recommendation. In *Proc. AAAI*, 1205–1211.

A Derivation of Novel Active Learning Criterion (6)

$$\begin{aligned}
& \arg \min_{X:|X|=N} H(Y_{V_t \setminus X_t} | Y_X) \\
&= \arg \max_{X:|X|=N} H(Y_{V_t} | L_U) - H(Y_{V_t \setminus X_t} | Y_X) \\
&= \arg \max_{X:|X|=N} H(Y_{V_t} | L_U) - H(Y_{V_t \setminus X_t} | L_U) + H(Y_{V_t \setminus X_t} | L_U) - \\
&\quad H(Y_{V_t \setminus X_t} | Y_X, L_U) + H(Y_{V_t \setminus X_t} | Y_X, L_U) - H(Y_{V_t \setminus X_t} | Y_X) \\
&= \arg \max_{X:|X|=N} H(Y_{X_t} | L_U, Y_{V_t \setminus X_t}) + I(Y_{V_t \setminus X_t}; Y_X | L_U) - \\
&\quad I(L_U; Y_{V_t \setminus X_t} | Y_X) \\
&= \arg \max_{X:|X|=N} H(Y_{X_t} | L_U) - I(L_U; Y_{V_t \setminus X_t} | Y_X).
\end{aligned}$$

The first equality follows from the fact that $H(Y_{V_t} | L_U)$ is a constant. The third equality is due to the chain rule for entropy $H(Y_{V_t} | L_U) = H(Y_{V_t \setminus X_t} | L_U) + H(Y_{X_t} | L_U, Y_{V_t \setminus X_t})$ as well as the definition of conditional mutual information $I(Y_{V_t \setminus X_t}; Y_X | L_U) \triangleq H(Y_{V_t \setminus X_t} | L_U) - H(Y_{V_t \setminus X_t} | Y_X, L_U)$ and $I(L_U; Y_{V_t \setminus X_t} | Y_X) \triangleq H(Y_{V_t \setminus X_t} | Y_X) - H(Y_{V_t \setminus X_t} | Y_X, L_U)$. The last equality follows from structural property **P2** shared by sparse CMOGP regression models in the unifying framework (Álvarez and Lawrence 2011) described in Section 2, which results in $H(Y_{X_t} | L_U, Y_{V_t \setminus X_t}) = H(Y_{X_t} | L_U)$ and $I(Y_{V_t \setminus X_t}; Y_X | L_U) = 0$.

B Time Complexity of Evaluating Active Learning Criterion in (6)

$$H(Y_{X_t} | L_U) = \frac{1}{2} \log(2\pi e)^{|X_t|} |\Sigma_{X_t X_t | U}|$$

where $\Sigma_{X_t X_t | U} = \Sigma_{X_t X_t} - \Sigma_{X_t U} \Sigma_{U U}^{-1} \Sigma_{U X_t}$ by definition (see last paragraph of Section 2). So, evaluating $H(Y_{X_t} | L_U)$ incurs $\mathcal{O}(|U|^3 + N^3)$ time for every $X \subset V$; this worst-case time complexity occurs when all the tuples in X are of measurement type t (i.e., $X = X_t$).

$$\begin{aligned}
I(L_U; Y_{V_t \setminus X_t} | Y_X) &= H(L_U | Y_X) - H(L_U | Y_{X \cup V_t \setminus X_t}) \\
&= \frac{1}{2} \log \frac{|\Sigma_{U U | X}|}{|\Sigma_{U U | X \cup V_t \setminus X_t}|} \\
&= \frac{1}{2} \log \frac{|\Sigma_{U U | X}|}{|\Sigma_{U U | \bigcup_{i \neq t} X_i \cup V_t}|}
\end{aligned}$$

where

$$\Sigma_{U U | A} = \Sigma_{U U} (\Sigma_{U U} + \Sigma_{U A} \Lambda_A^{-1} \Sigma_{A U})^{-1} \Sigma_{U U}$$

for any $A \subset D^+$, as derived in (Álvarez and Lawrence 2011). Therefore, evaluating $|\Sigma_{U U | X}|$ incurs $\mathcal{O}(|U|^3 + N^3)$ time for every $X \subset V$; this worst-case time complexity occurs when all the tuples in X are of one measurement type.

Let $A \triangleq \bigcup_{i \neq t} X_i \cup V_t$. Then, by the definition of Λ_A (see last paragraph of Section 2),

$$\Sigma_{U A} \Lambda_A^{-1} \Sigma_{A U} = \sum_{i \neq t} \Sigma_{U X_i} \Sigma_{X_i X_i | U}^{-1} \Sigma_{X_i U} + \Sigma_{U V_t} \Sigma_{V_t V_t | U}^{-1} \Sigma_{V_t U}.$$

Evaluating the $\sum_{i \neq t} \Sigma_{U X_i} \Sigma_{X_i X_i | U}^{-1} \Sigma_{X_i U}$ term incurs $\mathcal{O}(|U|^3 + N^3)$ time for every $X \subset V$; this worst-case time

complexity occurs when all the tuples in X are of one measurement type. Note that the $\Sigma_{U V_t} \Sigma_{V_t V_t | U}^{-1} \Sigma_{V_t U}$ term remains the same for every $X \subset V$ (i.e., since it is independent of X) and hence only needs to be computed once in $\mathcal{O}(|V_t|^3)$ time. Therefore, evaluating $|\Sigma_{U U | \bigcup_{i \neq t} X_i \cup V_t}| = |\Sigma_{U U | A}|$ incurs $\mathcal{O}(|U|^3 + N^3)$ time for every $X \subset V$ and a *one-off* cost of $\mathcal{O}(|V_t|^3)$ time. Consequently, evaluating $I(L_U; Y_{V_t \setminus X_t} | Y_X)$ incurs $\mathcal{O}(|U|^3 + N^3)$ time for every $X \subset V$ and a *one-off* cost of $\mathcal{O}(|V_t|^3)$ time.

So, evaluating our active learning criterion in (6) incurs $\mathcal{O}(|U|^3 + N^3)$ time for every $X \subset V$ and a *one-off* cost of $\mathcal{O}(|V_t|^3)$ time.

C Derivation of Greedy Criterion in (9)

If $i = t$, then

$$\begin{aligned}
& F(X \cup \{(x, t)\}) - F(X) \\
&= H(Y_{X \cup \{(x, t)\}} | L_U) - \\
&\quad (H(L_U | Y_{X \cup \{(x, t)\}}) - H(L_U | Y_{X \cup \{(x, t)\} \cup V_t \setminus (X_t \cup \{(x, t)\})}) - \\
&\quad (H(Y_{X_t} | L_U) - (H(L_U | Y_X) - H(L_U | Y_{X \cup V_t \setminus X_t})))) \\
&= H(Y_{X \cup \{(x, t)\}} | L_U) - H(Y_{X_t} | L_U) + \\
&\quad (H(L_U | Y_X) - H(L_U | Y_{X \cup \{(x, t)\}})) \\
&= H(Y_{(x, t)} | Y_{X_t}, L_U) + H(Y_{(x, t)} | Y_X) - H(Y_{(x, t)} | Y_X, L_U) \\
&= H(Y_{(x, t)} | L_U) + H(Y_{(x, t)} | Y_X) - H(Y_{(x, t)} | L_U) \\
&= H(Y_{(x, t)} | Y_X).
\end{aligned} \tag{11}$$

The first equality follows from (6) and (8). The second equality is due to $H(L_U | Y_{X \cup \{(x, t)\} \cup V_t \setminus (X_t \cup \{(x, t)\})}) = H(L_U | Y_{X \cup V_t \setminus X_t})$. The third equality is due to the chain rule for entropy $H(Y_{X_t \cup \{(x, t)\}} | L_U) = H(Y_{X_t} | L_U) + H(Y_{(x, t)} | Y_{X_t}, L_U)$ as well as the definition of conditional mutual information $I(L_U; Y_{(x, t)} | Y_X) \triangleq H(L_U | Y_X) - H(L_U | Y_{X \cup \{(x, t)\}}) = H(Y_{(x, t)} | Y_X) - H(Y_{(x, t)} | Y_X, L_U)$. The second last equality follows from structural property **P2** shared by sparse CMOGP regression models in the unifying framework (Álvarez and Lawrence 2011) described in Section 2.

Otherwise (i.e., $i \neq t$),

$$\begin{aligned}
& F(X \cup \{(x, i)\}) - F(X) \\
&= H(Y_{X \cup \{(x, i)\}} | L_U) - \\
&\quad (H(L_U | Y_{X \cup \{(x, i)\}}) - H(L_U | Y_{X \cup \{(x, i)\} \cup V_t \setminus X_t}) - \\
&\quad (H(Y_{X_t} | L_U) - (H(L_U | Y_X) - H(L_U | Y_{X \cup V_t \setminus X_t})))) \\
&= H(Y_{X \cup \{(x, i)\}} | L_U) - H(Y_{X_t} | L_U) + \\
&\quad (H(L_U | Y_X) - H(L_U | Y_{X \cup \{(x, i)\}})) + \\
&\quad H(L_U | Y_{X \cup V_t \setminus X_t \cup \{(x, i)\}}) - H(L_U | Y_{X \cup V_t \setminus X_t}) \\
&= H(Y_{(x, i)} | Y_X) - H(Y_{(x, i)} | L_U, Y_X) + \\
&\quad H(Y_{(x, i)} | Y_{X \cup V_t \setminus X_t}, L_U) - H(Y_{(x, i)} | Y_{X \cup V_t \setminus X_t}) \\
&= H(Y_{(x, i)} | Y_X) - H(Y_{(x, i)} | L_U) \\
&\quad + H(Y_{(x, i)} | L_U) - H(Y_{(x, i)} | Y_{X \cup V_t \setminus X_t}) \\
&= H(Y_{(x, i)} | Y_X) - H(Y_{(x, i)} | Y_{X \cup V_t \setminus X_t}).
\end{aligned} \tag{12}$$

The first equality follows from (6) and (8). The third equality is due to the definition of conditional mutual information $I(L_U; Y_{(x, i)} | Y_X) \triangleq H(L_U | Y_X) - H(L_U | Y_{X \cup \{(x, i)\}}) = H(Y_{(x, i)} | Y_X) - H(Y_{(x, i)} | L_U, Y_X)$ and $I(L_U; Y_{(x, i)} | Y_{X \cup V_t \setminus X_t}) \triangleq$

$H(L_U|Y_{X \cup V_t} \setminus X_t) - H(L_U|Y_{X \cup V_t} \setminus X_t \cup \{x, i\}) = H(Y_{\langle x, i \rangle} | Y_{X \cup V_t} \setminus X_t) - H(Y_{\langle x, i \rangle} | Y_{X \cup V_t} \setminus X_t, L_U)$. The second last equality follows from structural properties **P1** and **P2** shared by sparse CMOGP regression models in the unifying framework (Álvarez and Lawrence 2011) described in Section 2. Therefore, (9) results.

D Proof of Proposition 1

Before proving Proposition 1, the following lemmas are needed:

Lemma 3 For all $X \subset V$ and $\langle x, i \rangle \in V \setminus X$, $\Sigma_{\langle x, i \rangle \langle x, i \rangle | X}^{\text{PITC}} \geq \sigma_{n_i}^2$.

Its proof follows closely to that of Lemma 6 in (Cao, Low, and Dolan 2013).

Lemma 4 Assuming absence of suppressor variables, for all $X \subset V$ and $\langle x, i \rangle, \langle x', j \rangle \in V \setminus X$, $|\Sigma_{\langle x, i \rangle \langle x', j \rangle | X}^{\text{PITC}}| \leq 2|\sigma_{\langle x, i \rangle \langle x', j \rangle}|$.

Proof. If $i = j$, then

$$|\Sigma_{\langle x, i \rangle \langle x', j \rangle}^{\text{PITC}}| = |\sigma_{\langle x, i \rangle \langle x', j \rangle}| \leq 2|\sigma_{\langle x, i \rangle \langle x', j \rangle}|. \quad (13)$$

If $i \neq j$, then

$$\begin{aligned} |\Sigma_{\langle x, i \rangle \langle x', j \rangle}^{\text{PITC}}| &= |\Gamma_{\langle x, i \rangle \langle x', j \rangle}| \\ &= |\sigma_{\langle x, i \rangle \langle x', j \rangle} - \Sigma_{\langle x, i \rangle \langle x', j \rangle | U}| \\ &\leq |\sigma_{\langle x, i \rangle \langle x', j \rangle}| + |\Sigma_{\langle x, i \rangle \langle x', j \rangle | U}| \\ &\leq 2|\sigma_{\langle x, i \rangle \langle x', j \rangle}|. \end{aligned} \quad (14)$$

The first equality is due to (3) while the second equality follows from the definition of $\Gamma_{\langle x, i \rangle \langle x', j \rangle}$ (see last paragraph of Section 2). The last inequality follows from the practical assumption of absence of suppressor variables (Das and Kempe 2008): $|\Sigma_{\langle x, i \rangle \langle x', j \rangle | U}| \leq |\sigma_{\langle x, i \rangle \langle x', j \rangle}|$. Then,

$$|\Sigma_{\langle x, i \rangle \langle x', j \rangle | X}^{\text{PITC}}| \leq |\Sigma_{\langle x, i \rangle \langle x', j \rangle}^{\text{PITC}}| \leq 2|\sigma_{\langle x, i \rangle \langle x', j \rangle}|.$$

The first inequality follows from the practical assumption of absence of suppressor variables (Das and Kempe 2008). The second inequality is due to (13) and (14). \square

Main Proof. Let $B \triangleq V_t \setminus X_t$. Using the spectral theorem, $(\Sigma_{BB|X}^{\text{PITC}})^{-1} = WQW^\top$ where the columns of W are the eigenvectors of $(\Sigma_{BB|X}^{\text{PITC}})^{-1}$ and Q is a diagonal matrix comprising the eigenvalues of $(\Sigma_{BB|X}^{\text{PITC}})^{-1}$. Let $\lambda_{\max}(A)$ and $\lambda_{\min}(A)$ denote, respectively, the maximum and minimum

eigenvalues of matrix A , and $\alpha \triangleq W^\top \Sigma_{B\langle x, i \rangle | X}^{\text{PITC}}$.

$$\begin{aligned} &\Sigma_{\langle x, i \rangle \langle x, i \rangle | X}^{\text{PITC}} - \Sigma_{\langle x, i \rangle \langle x, i \rangle | X \cup V_t \setminus X_t}^{\text{PITC}} \\ &= \Sigma_{\langle x, i \rangle \langle x, i \rangle | X}^{\text{PITC}} - \\ &\quad \left(\Sigma_{\langle x, i \rangle \langle x, i \rangle | X}^{\text{PITC}} - \Sigma_{\langle x, i \rangle B | X}^{\text{PITC}} (\Sigma_{BB|X}^{\text{PITC}})^{-1} \Sigma_{B\langle x, i \rangle | X}^{\text{PITC}} \right) \\ &= \Sigma_{\langle x, i \rangle B | X}^{\text{PITC}} (\Sigma_{BB|X}^{\text{PITC}})^{-1} \Sigma_{B\langle x, i \rangle | X}^{\text{PITC}} \\ &= \Sigma_{\langle x, i \rangle B | X}^{\text{PITC}} WQW^\top \Sigma_{B\langle x, i \rangle | X}^{\text{PITC}} \\ &= \alpha^\top Q \alpha \\ &\leq \lambda_{\max}((\Sigma_{BB|X}^{\text{PITC}})^{-1}) \alpha^\top \alpha \\ &= \frac{\Sigma_{\langle x, i \rangle B | X}^{\text{PITC}} W W^\top \Sigma_{B\langle x, i \rangle | X}^{\text{PITC}}}{\lambda_{\min}(\Sigma_{BB|X}^{\text{PITC}})} \\ &= \frac{\|\Sigma_{\langle x, i \rangle B | X}^{\text{PITC}}\|_2^2}{\lambda_{\min}(\Sigma_{BB|X}^{\text{PITC}})} \\ &= \frac{\sum_{\langle x', t \rangle \in B} |\Sigma_{\langle x, i \rangle \langle x', t \rangle | X}^{\text{PITC}}|^2}{\lambda_{\min}(\Sigma_{BB|X}^{\text{PITC}})} \\ &\leq \frac{\sum_{\langle x', t \rangle \in B} 4|\sigma_{\langle x, i \rangle \langle x', t \rangle}|^2}{\lambda_{\min}(\Sigma_{BB|X}^{\text{PITC}})} \\ &\leq \frac{4\sigma_{s_i}^2 \sigma_{s_t}^2 \sum_{\langle x', t \rangle \in B} \mathcal{N}(x - x' | 0, P_0^{-1} + P_i^{-1} + P_t^{-1})^2}{\sigma_{n_t}^2} \\ &= 4\rho_t \sigma_{s_i}^2 R(\langle x, i \rangle, B). \end{aligned} \quad (15)$$

The first equality is due to the incremental update formula of GP posterior variance (see Appendix C in (Xu et al. 2014)). The first inequality is due to the fact that Q is a diagonal matrix comprising the eigenvalues of $(\Sigma_{BB|X}^{\text{PITC}})^{-1}$. The fifth equality is due to a property of eigenvalues that $\lambda_{\max}(A^{-1}) = 1/\lambda_{\min}(A)$. The sixth equality follows from the fact that $WW^\top = I$. The second inequality follows from Lemma 4. The third inequality is due to (1) and the fact that $\lambda_{\min}(\Sigma_{BB|X}^{\text{PITC}}) = \lambda_{\min}(\Sigma_{BB|X}^{\text{PITC}} - \sigma_{n_t}^2 I + \sigma_{n_t}^2 I) = \lambda_{\min}(\Sigma_{BB|X}^{\text{PITC}} - \sigma_{n_t}^2 I) + \sigma_{n_t}^2 \geq \sigma_{n_t}^2$ since $\lambda_{\min}(\Sigma_{BB|X}^{\text{PITC}} - \sigma_{n_t}^2 I) \geq 0$ (i.e., $\Sigma_{BB|X}^{\text{PITC}} - \sigma_{n_t}^2 I$ is a positive semi-definite matrix). Then,

$$\begin{aligned} &H(Y_{\langle x, i \rangle} | Y_X) - H(Y_{\langle x, i \rangle} | Y_{X \cup V_t} \setminus X_t) \\ &= \frac{1}{2} \log \frac{\Sigma_{\langle x, i \rangle \langle x, i \rangle | X}^{\text{PITC}}}{\Sigma_{\langle x, i \rangle \langle x, i \rangle | X \cup B}^{\text{PITC}}} \\ &\leq \frac{1}{2} \log \frac{\Sigma_{\langle x, i \rangle \langle x, i \rangle | X \cup B}^{\text{PITC}} + 4\rho_t \sigma_{s_i}^2 R(\langle x, i \rangle, B)}{\Sigma_{\langle x, i \rangle \langle x, i \rangle | X \cup B}^{\text{PITC}}} \\ &\leq \frac{1}{2} \log \left(1 + \frac{4\rho_t \sigma_{s_i}^2 R(\langle x, i \rangle, B)}{\sigma_{n_i}^2} \right) \\ &= \frac{1}{2} \log(1 + 4\rho_t \rho_i R(\langle x, i \rangle, B)). \end{aligned}$$

The first inequality is due to (15) while the second inequality follows from Lemma 3.

E Proof of Theorem 1

If $i = t$, then

$$H(Y_{\langle x, t \rangle} | Y_X) = \frac{1}{2} \log(2\pi e) \Sigma_{\langle x, t \rangle \langle x, t \rangle | X}^{\text{PITC}}$$

where $\Sigma_{\langle x,t \rangle \langle x,t \rangle | X}^{\text{PITC}}$ is previously defined in (3). So, evaluating $H(Y_{\langle x,t \rangle} | Y_X)$ incurs $\mathcal{O}(|U|^2)$ time for every $\langle x, t \rangle \in V_t \setminus X_t$ and $\mathcal{O}(|U|^3 + N^3)$ time in each iteration; this worst-case time complexity occurs when all the tuples in X are of one measurement type.

Otherwise (i.e., $i \neq t$),

$$\begin{aligned} & H(Y_{\langle x,i \rangle} | Y_X) - H(Y_{\langle x,i \rangle} | Y_{X \cup V_t \setminus X_t}) \\ &= \frac{1}{2} \log \frac{\Sigma_{\langle x,i \rangle \langle x,i \rangle | X}^{\text{PITC}}}{\Sigma_{\langle x,i \rangle \langle x,i \rangle | X \cup V_t \setminus X_t}^{\text{PITC}}} \\ &= \frac{1}{2} \log \frac{\Sigma_{\langle x,i \rangle \langle x,i \rangle | X}^{\text{PITC}}}{\Sigma_{\langle x,i \rangle \langle x,i \rangle | \bigcup_{i \neq t} X_i \cup V_t}^{\text{PITC}}} \end{aligned}$$

where $\Sigma_{\langle x,i \rangle \langle x,i \rangle | X}^{\text{PITC}}$ and $\Sigma_{\langle x,i \rangle \langle x,i \rangle | \bigcup_{i \neq t} X_i \cup V_t}^{\text{PITC}}$ are previously defined in (3). Therefore, evaluating $\Sigma_{\langle x,i \rangle \langle x,i \rangle | X}^{\text{PITC}}$ incurs $\mathcal{O}(|U|^2)$ time for every $\langle x, i \rangle \in V_t \setminus X_t$ and $\mathcal{O}(|U|^3 + N^3)$ time in each iteration; this worst-case time complexity occurs when all the tuples in X are of one measurement type.

Let $A \triangleq \bigcup_{i \neq t} X_i \cup V_t$. Then, by the definition of Λ_A (see last paragraph of Section 2),

$$\Sigma_{UA} \Lambda_A^{-1} \Sigma_{AU} = \sum_{i \neq t} \Sigma_{UX_i} \Sigma_{X_i X_i | U}^{-1} \Sigma_{X_i U} + \Sigma_{UV_t} \Sigma_{V_t V_t | U}^{-1} \Sigma_{V_t U}.$$

Evaluating the $\sum_{i \neq t} \Sigma_{UX_i} \Sigma_{X_i X_i | U}^{-1} \Sigma_{X_i U}$ term incurs $\mathcal{O}(|U|^3 + N^3)$ time in each iteration; this worst-case time complexity occurs when all the tuples in X are of one measurement type. Note that the $\Sigma_{UV_t} \Sigma_{V_t V_t | U}^{-1} \Sigma_{V_t U}$ term remains the same in each iteration (i.e., since it is independent of X) and hence only needs to be computed once in $\mathcal{O}(|V_t|^3)$ time in our approximation algorithm. As a result, evaluating $\Sigma_{\langle x,i \rangle \langle x,i \rangle | \bigcup_{i \neq t} X_i \cup V_t}^{\text{PITC}} = \Sigma_{\langle x,i \rangle \langle x,i \rangle | A}^{\text{PITC}}$ (specifically, its efficient formulation exploiting $\Sigma_{UA} \Lambda_A^{-1} \Sigma_{AU}$, as shown in (Álvarez and Lawrence 2011)) incurs $\mathcal{O}(|U|^2)$ time for every $\langle x, i \rangle \in V_t \setminus X_t$ and $\mathcal{O}(|U|^3 + N^3)$ time in each iteration, and a *one-off* cost of $\mathcal{O}(|V_t|^3)$ time. Consequently, evaluating $H(Y_{\langle x,i \rangle} | Y_X) - H(Y_{\langle x,i \rangle} | Y_{X \cup V_t \setminus X_t})$ incurs $\mathcal{O}(|U|^2)$ time for every $\langle x, i \rangle \in V_t \setminus X_t$ and $\mathcal{O}(|U|^3 + N^3)$ time in each iteration, and a *one-off* cost of $\mathcal{O}(|V_t|^3)$ time.

Since $|U| \leq |V_t| < |V|$, our approximation algorithm thus incurs $\mathcal{O}(N(|V||U|^2 + N^3) + |V_t|^3)$ time.

F Proof of Lemma 1

To prove that $F(X)$ is ϵ -submodular, we have to show that

$$F(X' \cup \{\langle x, i \rangle\}) - F(X') \leq F(X \cup \{\langle x, i \rangle\}) - F(X) + \epsilon$$

for any $X \subseteq X' \subseteq V$ and $\langle x, i \rangle \in V \setminus X'$. Before doing this, the following lemma is needed:

Lemma 5 *Suppose that $\epsilon_1 \geq 0$ is given. For any $\langle x, i \rangle \in V_t \setminus X'_t$, if $\Sigma_{\langle x,i \rangle \langle x,i \rangle | X \cup V_t \setminus X'_t}^{\text{PITC}} - \Sigma_{\langle x,i \rangle \langle x,i \rangle | X' \cup V_t \setminus X'_t}^{\text{PITC}} \leq \epsilon_1$, then $I(Y_{\langle x,i \rangle}; Y_{V_t \setminus X'_t} | Y_{X'}) \leq I(Y_{\langle x,i \rangle}; Y_{V_t \setminus X'_t} | Y_X) + \epsilon$ where $\epsilon = 0.5 \log(1 + \epsilon_1 / \sigma_{n^*}^2)$.*

Proof. Let $\bar{X} \triangleq X' \setminus X$. Then,

$$\begin{aligned} & I(Y_{\langle x,i \rangle}; Y_{\bar{X}} | Y_{X \cup V_t \setminus X'_t}) \\ &= H(Y_{\langle x,i \rangle} | Y_{X \cup V_t \setminus X'_t}) - H(Y_{\langle x,i \rangle} | Y_{\bar{X} \cup X \cup V_t \setminus X'_t}) \\ &= \frac{1}{2} \log \frac{\Sigma_{\langle x,i \rangle \langle x,i \rangle | X \cup V_t \setminus X'_t}^{\text{PITC}}}{\Sigma_{\langle x,i \rangle \langle x,i \rangle | \bar{X} \cup X \cup V_t \setminus X'_t}^{\text{PITC}}} \\ &\leq \frac{1}{2} \log \frac{\Sigma_{\langle x,i \rangle \langle x,i \rangle | \bar{X} \cup X \cup V_t \setminus X'_t}^{\text{PITC}} + \epsilon_1}{\Sigma_{\langle x,i \rangle \langle x,i \rangle | \bar{X} \cup X \cup V_t \setminus X'_t}^{\text{PITC}}} \\ &= \frac{1}{2} \log \left(1 + \frac{\epsilon_1}{\Sigma_{\langle x,i \rangle \langle x,i \rangle | \bar{X} \cup X \cup V_t \setminus X'_t}^{\text{PITC}}} \right) \\ &\leq \frac{1}{2} \log \left(1 + \frac{\epsilon_1}{\sigma_{n_i}^2} \right) \\ &\leq \frac{1}{2} \log \left(1 + \frac{\epsilon_1}{\sigma_{n^*}^2} \right). \end{aligned} \tag{16}$$

The first inequality is due to the sufficient condition. The second inequality follows from Lemma 3. Then, by the definition of conditional mutual information,

$$\begin{aligned} & I(Y_{\langle x,i \rangle}; Y_{V_t \setminus X'_t} | Y_{\bar{X} \cup X}) + I(Y_{\langle x,i \rangle}; Y_{\bar{X}} | Y_X) \\ &= H(Y_{\langle x,i \rangle} | Y_{\bar{X} \cup X}) - H(Y_{\langle x,i \rangle} | Y_{\bar{X} \cup X \cup V_t \setminus X'_t}) + \\ & \quad H(Y_{\langle x,i \rangle} | Y_X) - H(Y_{\langle x,i \rangle} | Y_{\bar{X} \cup X}) \\ &= H(Y_{\langle x,i \rangle} | Y_X) - H(Y_{\langle x,i \rangle} | Y_{\bar{X} \cup X \cup V_t \setminus X'_t}) \\ &= H(Y_{\langle x,i \rangle} | Y_X) - H(Y_{\langle x,i \rangle} | Y_{X \cup V_t \setminus X'_t}) + \\ & \quad H(Y_{\langle x,i \rangle} | Y_{X \cup V_t \setminus X'_t}) - H(Y_{\langle x,i \rangle} | Y_{\bar{X} \cup X \cup V_t \setminus X'_t}) \\ &= I(Y_{\langle x,i \rangle}; Y_{V_t \setminus X'_t} | Y_X) + I(Y_{\langle x,i \rangle}; Y_{\bar{X}} | Y_{X \cup V_t \setminus X'_t}). \end{aligned}$$

Therefore,

$$\begin{aligned} & I(Y_{\langle x,i \rangle}; Y_{V_t \setminus X'_t} | Y_{X'}) \\ &= I(Y_{\langle x,i \rangle}; Y_{V_t \setminus X'_t} | Y_{\bar{X} \cup X}) \\ &= I(Y_{\langle x,i \rangle}; Y_{V_t \setminus X'_t} | Y_X) + I(Y_{\langle x,i \rangle}; Y_{\bar{X}} | Y_{X \cup V_t \setminus X'_t}) - \\ & \quad I(Y_{\langle x,i \rangle}; Y_{\bar{X}} | Y_X) \\ &\leq I(Y_{\langle x,i \rangle}; Y_{V_t \setminus X'_t} | Y_X) + I(Y_{\langle x,i \rangle}; Y_{\bar{X}} | Y_{X \cup V_t \setminus X'_t}) \\ &\leq I(Y_{\langle x,i \rangle}; Y_{V_t \setminus X'_t} | Y_X) + 0.5 \log \left(1 + \frac{\epsilon_1}{\sigma_{n^*}^2} \right). \end{aligned}$$

The first inequality is due to the fact that conditional mutual information is non-negative. The last inequality follows from (16). \square

Main Proof. To prove that $F(X)$ is ϵ -submodular, we have to show that $H(Y_{\langle x,i \rangle} | Y_{X'}) - \delta_i H(Y_{\langle x,i \rangle} | Y_{X' \cup V_t \setminus X'_t}) \leq H(Y_{\langle x,i \rangle} | Y_X) - \delta_i H(Y_{\langle x,i \rangle} | Y_{X \cup V_t \setminus X_t}) + \epsilon$ for any $X \subseteq X' \subseteq V$ and $\langle x, i \rangle \in V \setminus X'$.

If $i = t$, then $H(Y_{\langle x,i \rangle} | Y_{X'}) \leq H(Y_{\langle x,i \rangle} | Y_X) \leq H(Y_{\langle x,i \rangle} | Y_X) + \epsilon$ for any $\epsilon \geq 0$ due to the ‘‘information never hurts’’ bound for entropy (Cover and Thomas 1991).

Otherwise (i.e., $i \neq t$),

$$\begin{aligned} & H(Y_{\langle x,i \rangle} | Y_{X'}) - H(Y_{\langle x,i \rangle} | Y_{X' \cup V_t \setminus X'_t}) \\ &\leq H(Y_{\langle x,i \rangle} | Y_X) - H(Y_{\langle x,i \rangle} | Y_{X \cup V_t \setminus X_t}) + \epsilon \\ &\leq H(Y_{\langle x,i \rangle} | Y_X) - H(Y_{\langle x,i \rangle} | Y_{X \cup V_t \setminus X_t}) + \epsilon \end{aligned}$$

where $\epsilon = 0.5 \log(1 + \epsilon_1 / \sigma_{n^*}^2)$. The first inequality is due to Lemma 5. The second inequality follows from the ‘‘information never hurts’’ bound for entropy (Cover and Thomas

1991): $H(Y_{\langle x, i \rangle} | Y_{X \cup V_t \setminus X'_t}) \geq H(Y_{\langle x, i \rangle} | Y_{X \cup V_t \setminus X_t})$ since $(V_t \setminus X'_t) \subseteq (V_t \setminus X_t)$.

G Proof of Theorem 2

Our proof here is similar to that of Theorem 1.5 in (Krause and Golovin 2014) which is a generalization of the well-known result of Nemhauser, Wolsey, and Fisher (1978). The key difference is that we exploit ϵ -submodularity of $F(X)$ (i.e., Lemma 1) instead of submodularity, as shown below for completeness.

Let $X^* \triangleq \{\langle x_1, s_1 \rangle^*, \dots, \langle x_N, s_N \rangle^*\}$ be the optimal set of selected observations, X^k be the set of tuples selected by our approximation algorithm in iteration $k = 1, \dots, N$, $X^0 \triangleq \emptyset$, and $\Delta(\langle x, i \rangle | X) \triangleq F(X \cup \{\langle x, i \rangle\}) - F(X)$. Then,

$$\begin{aligned} & F(X^*) \\ & \leq F(X^* \cup X^k) \\ & = F(X^k) + \sum_{j=1}^N \Delta \left(\langle x_j, s_j \rangle^* \left| \bigcup_{r=1}^{j-1} \{\langle x_r, s_r \rangle^*\} \cup X^k \right. \right) \\ & \leq F(X^k) + \sum_{j=1}^N (\Delta(\langle x_j, s_j \rangle^* | X^k) + \epsilon) \\ & \leq F(X^k) + \sum_{j=1}^N (F(X^{k+1}) - F(X^k) + \epsilon) \\ & \leq F(X^k) + N (F(X^{k+1}) - F(X^k) + \epsilon). \end{aligned}$$

The first inequality follows from the nondecreasing property of $F(X)$. The first equality is a straightforward telescoping sum. The second inequality follows from the ϵ -submodularity of $F(X)$, as proven in Lemma 1. The third inequality follows from (9). Then,

$$F(X^*) - F(X^k) \leq N (F(X^{k+1}) - F(X^k) + \epsilon). \quad (17)$$

Let $\zeta_k \triangleq F(X^*) - F(X^k)$. Then, (17) can be rewritten as $\zeta_k \leq N(\zeta_k - \zeta_{k+1} + \epsilon)$ which can be rearranged to yield

$$\zeta_{k+1} \leq \left(1 - \frac{1}{N}\right) \zeta_k + \epsilon. \quad (18)$$

Then, by recursion of (18), it is straightforward to get

$$\zeta_k \leq \left(1 - \frac{1}{N}\right)^k \zeta_0 + N \left(1 - \left(1 - \frac{1}{N}\right)^k\right) \epsilon. \quad (19)$$

Then, by substituting $\zeta_k = F(X^*) - F(X^k)$ and $\zeta_0 = F(X^*) - F(X^0) = F(X^*)$, (19) can be rearranged to

$$\begin{aligned} F(X^k) & \geq \left(1 - \left(1 - \frac{1}{N}\right)^k\right) (F(X^*) - N\epsilon) \\ & \geq (1 - e^{-k/N}) (F(X^*) - N\epsilon). \end{aligned}$$

The second inequality follows from the well-known inequality $e^{-x} \geq 1 - x$. Finally, Theorem 2 is obtained when $k = N$ and $\epsilon = 0.5 \log(1 + \epsilon_1 / \sigma_{n^*}^2)$, as defined in Lemma 1.

H Proof of Lemma 2

Let $B \triangleq \tilde{X} \cup V_t \setminus X_t$ and $A \triangleq X \setminus \tilde{X}$. From the incremental update formula of GP posterior variance (see Appendix C in (Xu et al. 2014)),

$$\begin{aligned} & \Sigma_{\langle x, i \rangle \langle x, i \rangle | B}^{\text{PITC}} - \Sigma_{\langle x, i \rangle \langle x, i \rangle | B \cup A}^{\text{PITC}} \\ & = \Sigma_{\langle x, i \rangle \langle x, i \rangle | B}^{\text{PITC}} - \\ & \quad \left(\Sigma_{\langle x, i \rangle \langle x, i \rangle | B}^{\text{PITC}} - \Sigma_{\langle x, i \rangle A | B}^{\text{PITC}} (\Sigma_{AA | B}^{\text{PITC}})^{-1} \Sigma_{A \langle x, i \rangle | B}^{\text{PITC}} \right) \\ & = \Sigma_{\langle x, i \rangle A | B}^{\text{PITC}} (\Sigma_{AA | B}^{\text{PITC}})^{-1} \Sigma_{A \langle x, i \rangle | B}^{\text{PITC}}. \end{aligned} \quad (20)$$

Let $\Sigma_{AA | B}^{\text{PITC}} \triangleq C + E$ where C is defined as a matrix with the same diagonal components as $\Sigma_{AA | B}^{\text{PITC}}$ and off-diagonal components 0 while E is defined as a matrix with diagonal components 0 and the same off-diagonal components as $\Sigma_{AA | B}^{\text{PITC}}$. Then,

$$\begin{aligned} \|C^{-1}\|_2 & = \lambda_{\max}(C^{-1}) \\ & = \frac{1}{\lambda_{\min}(C)} \\ & = \frac{1}{\min_{\langle x, i \rangle \in A} \Sigma_{\langle x, i \rangle \langle x, i \rangle | B}^{\text{PITC}}} \\ & \leq \frac{1}{\sigma_{n_i}^2} \leq \frac{1}{\sigma_{n^*}^2}. \end{aligned} \quad (21)$$

The first equality is due to a property of matrix norm in Section 10.4.5 in (Petersen and Pedersen 2012). The second equality is due to a property of eigenvalues that $\lambda_{\max}(C^{-1}) = 1/\lambda_{\min}(C)$. The third equality is due to the diagonal property of C . The first inequality is due to Lemma 3.

Matrix E comprises off-diagonal components $\Sigma_{\langle x, i \rangle \langle x', j \rangle | B}^{\text{PITC}}$ for all $\langle x, i \rangle, \langle x', j \rangle \in A$ such that $\langle x, i \rangle \neq \langle x', j \rangle$, each of which has an absolute value not more than $2\sigma_{s^*}^2 \xi^{p^2}$:

$$\begin{aligned} & |\Sigma_{\langle x, i \rangle \langle x', j \rangle | B}^{\text{PITC}}| \\ & \leq 2|\sigma_{s_i} \sigma_{s_j}| |\mathcal{N}(x - x' | 0, P_0^{-1} + P_i^{-1} + P_j^{-1})| \\ & = 2|\sigma_{s_i} \sigma_{s_j}| \exp \left\{ -\frac{1}{2} \sum_{v=1}^d \frac{(x_v - x'_v)^2}{\ell_v^{ij}} \right\} \\ & \leq 2|\sigma_{s_i} \sigma_{s_j}| \exp \left\{ -\frac{(x_1 - x'_1)^2}{2\ell_1^{ij}} \right\} \\ & \leq 2|\sigma_{s_i} \sigma_{s_j}| \exp \left\{ -\frac{p^2 \omega^2}{2\ell} \right\} \\ & = 2|\sigma_{s_i} \sigma_{s_j}| \xi^{p^2} \\ & \leq 2\sigma_{s^*}^2 \xi^{p^2} \end{aligned}$$

where x_v is the v -th component of a d -dimensional location vector x and ℓ_v^{ij} denotes the v -th diagonal component of $P_0^{-1} + P_i^{-1} + P_j^{-1}$. The first inequality follows from Lemma 4. The second equality is due to the precision matrices being diagonal. The third inequality follows from $\ell \triangleq \max_{i, j \in \{1, \dots, M\}} \ell_1^{ij}$ and the fact that the distance between x_1 and x'_1 of any $\langle x, i \rangle, \langle x', j \rangle \in A$ must be at least $p\omega$ due to the construction of V^- . Therefore,

$$\|E\|_2 \leq 2N\sigma_{s^*}^2 \xi^{p^2} \quad (22)$$

due to a property that the 2-norm of a matrix is at most its largest absolute component multiplied by its dimension (Golub and Van Loan 1996).

Similarly, $\Sigma_{\langle x,i \rangle A|B}^{\text{PITC}}$ comprises components $\Sigma_{\langle x,i \rangle \langle x',j \rangle |B}^{\text{PITC}}$ for all $\langle x',j \rangle \in A$, each of which has an absolute value not more than $2\sigma_{s^*}^2 \xi^{p^2}$:

$$|\Sigma_{\langle x,i \rangle \langle x',j \rangle |B}^{\text{PITC}}| \leq 2|\sigma_{\langle x,i \rangle \langle x',j \rangle}| \leq 2\sigma_{s^*}^2 \xi^{p^2}. \quad (23)$$

Now,

$$\begin{aligned} & \Sigma_{\langle x,i \rangle A|B}^{\text{PITC}}(C+E)^{-1}\Sigma_{A\langle x,i \rangle |B}^{\text{PITC}} - \Sigma_{\langle x,i \rangle A|B}^{\text{PITC}}C^{-1}\Sigma_{A\langle x,i \rangle |B}^{\text{PITC}} \\ &= \Sigma_{\langle x,i \rangle A|B}^{\text{PITC}} \left\{ (C+E)^{-1} - C^{-1} \right\} \Sigma_{A\langle x,i \rangle |B}^{\text{PITC}} \\ &\leq \|\Sigma_{\langle x,i \rangle A|B}^{\text{PITC}}\|_2^2 \|(C+E)^{-1} - C^{-1}\|_2 \\ &\leq \sum_{\langle x',j \rangle \in A} |\Sigma_{\langle x,i \rangle \langle x',j \rangle |B}^{\text{PITC}}|^2 \frac{\|C^{-1}\|_2 \|E\|_2}{\|C^{-1}\|_2 - \|E\|_2} \\ &\leq 4N\sigma_{s^*}^4 \xi^{2p^2} \frac{\|C^{-1}\|_2 \|E\|_2}{\|C^{-1}\|_2 - \|E\|_2}. \end{aligned} \quad (24)$$

The first inequality is due to Cauchy-Schwarz inequality and submultiplicativity of the matrix norm (Stewart and Sun 1990). The second inequality follows from an important result in the perturbation theory of matrix inverses (in particular, Theorem III.2.5 in (Stewart and Sun 1990)). It requires the assumption $\|C^{-1}E\|_2 < 1$. Using (21), (22), and the matrix norm property in Section 10.4.2 in (Petersen and Pedersen 2012), this assumption can be satisfied by

$$\|C^{-1}E\|_2 \leq \|C^{-1}\|_2 \|E\|_2 \leq \frac{2N\sigma_{s^*}^2 \xi^{p^2}}{\sigma_{n^*}^2} < 1.$$

Then,

$$p^2 > \log \left(\frac{\sigma_{n^*}^2}{2N\sigma_{s^*}^2} \right) / \log \xi. \quad (25)$$

The last inequality in (24) is due to (23). Then, from both (20) and (24),

$$\begin{aligned} & \Sigma_{\langle x,i \rangle \langle x,i \rangle |B}^{\text{PITC}} - \Sigma_{\langle x,i \rangle \langle x,i \rangle |B \cup A}^{\text{PITC}} \\ &= \Sigma_{\langle x,i \rangle A|B}^{\text{PITC}}(C+E)^{-1}\Sigma_{A\langle x,i \rangle |B}^{\text{PITC}} \\ &\leq \Sigma_{\langle x,i \rangle A|B}^{\text{PITC}}C^{-1}\Sigma_{A\langle x,i \rangle |B}^{\text{PITC}} + 4N\sigma_{s^*}^4 \xi^{2p^2} \frac{\|C^{-1}\|_2 \|E\|_2}{\|C^{-1}\|_2 - \|E\|_2} \\ &\leq \|\Sigma_{\langle x,i \rangle A|B}^{\text{PITC}}\|_2^2 \|C^{-1}\|_2 + 4N\sigma_{s^*}^4 \xi^{2p^2} \frac{\|C^{-1}\|_2 \|E\|_2}{\|C^{-1}\|_2 - \|E\|_2} \\ &\leq 4N\sigma_{s^*}^4 \xi^{2p^2} \|C^{-1}\|_2 + 4N\sigma_{s^*}^4 \xi^{2p^2} \frac{\|C^{-1}\|_2 \|E\|_2}{\|C^{-1}\|_2 - \|E\|_2} \\ &= 4N\sigma_{s^*}^4 \xi^{2p^2} \|C^{-1}\|_2 \left(1 + \frac{\|E\|_2}{\|C^{-1}\|_2 - \|E\|_2} \right) \\ &= \frac{4N\sigma_{s^*}^4 \xi^{2p^2}}{\|C^{-1}\|_2 - \|E\|_2} \\ &\leq \frac{4N\sigma_{s^*}^4 \xi^{2p^2}}{\sigma_{n^*}^2 - 2N\sigma_{s^*}^2 \xi^{p^2}}. \end{aligned}$$

The first inequality is due to (24). The second inequality is due to Cauchy-Schwarz inequality. The third inequality is due to (23). The last inequality follows from (21) and (22).

To satisfy (10) in Lemma 1, let

$$\frac{4N\sigma_{s^*}^4 \xi^{2p^2}}{\sigma_{n^*}^2 - 2N\sigma_{s^*}^2 \xi^{p^2}} \leq \epsilon_1.$$

Then,

$$p^2 \geq \log \left\{ \frac{1}{4\sigma_{s^*}^2} \left(\sqrt{\epsilon_1^2 + \frac{4\epsilon_1\sigma_{n^*}^2}{N}} - \epsilon_1 \right) \right\} / \log \xi. \quad (26)$$

Finally, from both (25) and (26), Lemma 2 results.

I Jura and IEQ Datasets

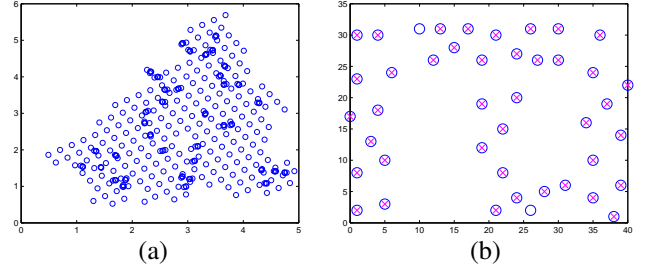


Figure 3: Sampling locations for the (a) Jura (km) and (b) IEQ (m) datasets where ‘o’ and ‘x’ denote locations of temperature and light sensors, respectively.

J Signal-to-Noise Ratios for Jura Dataset

	Ig-Cd	Ni	Ig-Zn
$\sigma_{s_i}^2$	2.2204	8.8280	2.3198
$\sigma_{n_i}^2$	0.0853	0.1130	0.0596
ρ_i	26.0305	78.1239	38.9228

Table 1: Signal-to-noise ratios ρ_i of Ig-Cd, Ni, and Ig-Zn measurements for Jura dataset with $|U| = 100$.