
Supplement to *Generalization Errors and Learning Curves for Regression with Multi-task Gaussian Processes*

Kian Ming A. Chai

School of Informatics, University of Edinburgh,
10 Crichton Street, Edinburgh EH8 9AB, UK
k.m.a.chai@ed.ac.uk

S.1 Introduction

We provide material supplementary to the the main paper. Much of this will be detailed proofs for the propositions and corollaries. To distinguish from the equations in the main paper, the equations here are prefixed with ‘‘S.’’.

S.2 Proof for Proposition 5

In this section, we give the proof for Proposition 5 in the main text.

S.2.1 Proof for Proposition 5a

Recall from (3) that $\sigma_T^2(\rho)$ is given by

$$\sigma_T^2(\rho) \stackrel{\text{def}}{=} k_{**} - \begin{pmatrix} \mathbf{k}_{T*}^x \\ \rho \mathbf{k}_{S*}^x \end{pmatrix}^\top \begin{pmatrix} K_{TT}^x + \sigma_n^2 I & \rho K_{TS}^x \\ \rho K_{ST}^x & K_{SS}^x + \sigma_n^2 I \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{k}_{T*}^x \\ \rho \mathbf{k}_{S*}^x \end{pmatrix} \quad (\text{S.1})$$

To perform the matrix inverse in the above equation, we use the following formula for inverting block matrices:

Theorem S.1. *Banachiewicz inversion formula (see e.g., [1]).*

$$\begin{aligned} \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}^{-1} &= \begin{pmatrix} A_{11}^{-1} + A_{11}^{-1} A_{12} C^{-1} A_{21} A_{11}^{-1} & -A_{11}^{-1} A_{12} C^{-1} \\ -C^{-1} A_{21} A_{11}^{-1} & C^{-1} \end{pmatrix} \\ &= \begin{pmatrix} A_{11}^{-1} & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} -A_{11}^{-1} A_{12} \\ I \end{pmatrix} C^{-1} \begin{pmatrix} -A_{21} A_{11}^{-1} & I \end{pmatrix}, \end{aligned}$$

where

$$C \stackrel{\text{def}}{=} A_{22} - A_{21} A_{11}^{-1} A_{12}.$$

The role of C in the above theorem is played by $A(\rho)$ defined as

$$A(\rho) \stackrel{\text{def}}{=} K_{SS}^x + \sigma_n^2 I - \rho^2 K_{ST}^x (K_{TT}^x + \sigma_n^2 I)^{-1} K_{TS}^x. \quad (\text{S.2})$$

In addition, we let

$$\begin{aligned} \mathbf{v}(\rho) &\stackrel{\text{def}}{=} \begin{pmatrix} -\rho K_{ST}^x (K_{TT}^x + \sigma_n^2 I)^{-1} & I \end{pmatrix} \begin{pmatrix} \mathbf{k}_{T*}^x \\ \rho \mathbf{k}_{S*}^x \end{pmatrix} \\ &= -\rho K_{ST}^x (K_{TT}^x + \sigma_n^2 I)^{-1} \mathbf{k}_{T*}^x + \rho \mathbf{k}_{S*}^x \\ &= \rho (\mathbf{k}_{S*}^x - K_{ST}^x (K_{TT}^x + \sigma_n^2 I)^{-1} \mathbf{k}_{T*}^x). \end{aligned} \quad (\text{S.3})$$

Then

$$\sigma_T^2(\rho) = k_{**} - \begin{pmatrix} \mathbf{k}_{T*}^x \\ \rho \mathbf{k}_{S*}^x \end{pmatrix}^\top \begin{pmatrix} (K_{TT}^x + \sigma_n^2 I)^{-1} & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \mathbf{k}_{T*}^x \\ \rho \mathbf{k}_{S*}^x \end{pmatrix} - \mathbf{v}(\rho)^\top [A(\rho)]^{-1} \mathbf{v}(\rho) \quad (\text{S.4})$$

$$= k_{**} - (\mathbf{k}_{T*}^x)^\top (K_{TT}^x + \sigma_n^2 I)^{-1} \mathbf{k}_{T*}^x - \mathbf{v}(\rho)^\top [A(\rho)]^{-1} \mathbf{v}(\rho) \quad (\text{S.5})$$

We can identify $\mathbf{v}(\rho) = \rho \mathbf{v}(1)$. If we also write \mathbf{v}_1 for $\mathbf{v}(1)$, then

$$\sigma_T^2(\rho) = k_{**} - (\mathbf{k}_{T*}^x)^\top (K_{TT}^x + \sigma_n^2 I)^{-1} \mathbf{k}_{T*}^x - \rho^2 \mathbf{v}_1^\top [A(\rho)]^{-1} \mathbf{v}_1 \quad (\text{S.6})$$

By substituting 0 for ρ above, the first two terms on the right of the equation can be identified with $\sigma_T^2(0)$. Thus

$$\sigma_T^2(\rho) = \sigma_T^2(0) - \rho^2 \mathbf{v}_1^\top [A(\rho)]^{-1} \mathbf{v}_1. \quad (\text{S.7})$$

Now, observe that $K_{ST}^x (K_{TT}^x + \sigma_n^2 I)^{-1} K_{TS}^x$ is positive semi-definite, since we can factorize it into the form XX^\top for some matrix X . Proceeding from this, we have

$$\begin{aligned} & K_{ST}^x (K_{TT}^x + \sigma_n^2 I)^{-1} K_{TS}^x \succeq 0 \\ \iff & (1 - \rho^2) K_{ST}^x (K_{TT}^x + \sigma_n^2 I)^{-1} K_{TS}^x \succeq 0 \quad \text{since } \rho^2 \in [0, 1] \\ \iff & -\rho^2 K_{ST}^x (K_{TT}^x + \sigma_n^2 I)^{-1} K_{TS}^x \succeq -K_{ST}^x (K_{TT}^x + \sigma_n^2 I)^{-1} K_{TS}^x \\ \iff & K_{SS}^x + \sigma_n^2 I - \rho^2 K_{ST}^x (K_{TT}^x + \sigma_n^2 I)^{-1} K_{TS}^x \succeq K_{SS}^x + \sigma_n^2 I - K_{ST}^x (K_{TT}^x + \sigma_n^2 I)^{-1} K_{TS}^x \\ \text{i.e.,} & \quad A(\rho^2) \succeq A(1) \\ \iff & \quad [A(\rho^2)]^{-1} \preceq [A(1)]^{-1} \\ \iff & \quad \mathbf{v}_1^\top [A(\rho^2)]^{-1} \mathbf{v} \leq \mathbf{v}_1^\top [A(1)]^{-1} \mathbf{v} \\ \iff & \quad \sigma_T^2(0) - \rho^2 \mathbf{v}_1^\top [A(\rho^2)]^{-1} \mathbf{v} \geq \sigma_T^2(0) - \rho^2 \mathbf{v}_1^\top [A(1)]^{-1} \mathbf{v} \\ \text{i.e.,} & \quad \sigma_T^2(\rho^2) \geq \sigma_T^2(0) - \rho^2 \mathbf{v}_1^\top [A(1)]^{-1} \mathbf{v} \end{aligned}$$

We write A_1 for $A(1)$. To complete the proof, we use the identity $\mathbf{v}_1^\top A_1^{-1} \mathbf{v}_1 = \sigma_T^2(0) - \sigma_T^2(1)$, which is obtained by substituting 1 for ρ into (S.7). Further re-grouping of terms leads to the result. \square

Remark The expression for $\sigma_T^2(\rho)$ given by (S.7) can also be obtained by repeated conditioning. Consider the following covariance matrix between the query f_*^T , the noisy observations \mathbf{y}_S^S at X_S for task S , and the noisy observations \mathbf{y}_T^T at X_T for task T ,

$$\mathbb{C} \begin{pmatrix} f_*^T \\ \mathbf{y}_S^S \\ \mathbf{y}_T^T \end{pmatrix} = \begin{pmatrix} k_{**} & \rho (\mathbf{k}_{S*}^x)^\top & (\mathbf{k}_{T*}^x)^\top \\ \rho \mathbf{k}_{S*}^x & K_{SS}^x + \sigma_n^2 I & \rho K_{ST}^x \\ \mathbf{k}_{T*}^x & \rho K_{TS}^x & K_{TT}^x + \sigma_n^2 I \end{pmatrix}. \quad (\text{S.8})$$

By conditioning on \mathbf{y}_T^T , we obtain

$$\mathbb{C} \begin{pmatrix} f_*^T \\ \mathbf{y}_S^S \end{pmatrix} \Big| \mathbf{y}_T^T = \begin{pmatrix} \sigma_T(0) & (\mathbf{v}(\rho))^\top \\ \mathbf{v}(\rho) & A(\rho) \end{pmatrix}. \quad (\text{S.9})$$

Conditioning subsequently on \mathbf{y}_S^S gives

$$\sigma_T^2(\rho) \stackrel{\text{def}}{=} \mathbb{C}(f_*^T | \mathbf{y}_T^T, \mathbf{y}_S^S) = \sigma_T^2(0) - \mathbf{v}(\rho)^\top [A(\rho)]^{-1} \mathbf{v}(\rho). \quad (\text{S.10})$$

To complete, we simply write $\rho \mathbf{v}_1$ for $\mathbf{v}(\rho)$.

S.2.2 Proof for Proposition 5b

Recall that the exact posterior variance is

$$\sigma_T^2(\rho) = \sigma_T^2(0) - \rho^2 \mathbf{v}_1^\top [A(\rho)]^{-1} \mathbf{v}_1, \quad (\text{S.11})$$

i.e., (S.7). Denote the lower bound by $\underline{\sigma}_T^2(\rho)$. Then from the proof for Proposition 5a, we have

$$\underline{\sigma}_T^2(\rho) = \sigma_T^2(0) - \rho^2 \mathbf{v}_1^T A_1^{-1} \mathbf{v}_1. \quad (\text{S.12})$$

Define the gap between the exact posterior variance and its lower bound as

$$\begin{aligned} g(\rho^2) &\stackrel{\text{def}}{=} \sigma_T^2(\rho) - \underline{\sigma}_T^2(\rho) \\ &= -\rho^2 \mathbf{v}_1^T [A(\rho)]^{-1} \mathbf{v}_1 + \rho^2 \mathbf{v}_1^T A_1^{-1} \mathbf{v}_1. \end{aligned} \quad (\text{S.13})$$

Ignoring the first term, which is negative, gives

$$\begin{aligned} g(\rho^2) &\leq \rho^2 \mathbf{v}_1^T A_1^{-1} \mathbf{v}_1 \\ &= \rho^2 [\sigma_T^2(0) - \sigma_T^2(1)]. \quad \square \end{aligned}$$

S.2.3 Proof for Proposition 5c

We rewrite the gap (S.13) between the exact posterior variance and its lower bound as

$$g(\rho^2) = \rho^2 \mathbf{v}_1^T [A_1^{-1} - A(\rho)]^{-1} \mathbf{v}_1. \quad (\text{S.14})$$

Next, express $A(\rho)$ as

$$A(\rho) = A_1 + (1 - \rho^2) K_{ST}^x (K_{TT}^x + \sigma_n^2 I)^{-1} K_{TS}^x, \quad (\text{S.15})$$

so that we can use the Woodbury identity to expand its inverse in S.14, and write

$$g(\rho^2) = \mathbf{v}_1^T A_1^{-1} K_{ST}^x [B(\rho^2)]^{-1} K_{TS}^x A_1^{-1} \mathbf{v}_1, \quad (\text{S.16})$$

where

$$B(\rho^2) \stackrel{\text{def}}{=} D(\rho^2) + \frac{1}{\rho^2} K_{TS}^x A_1^{-1} K_{ST}^x \quad (\text{S.17})$$

$$D(\rho^2) \stackrel{\text{def}}{=} \frac{1}{\rho^2(1 - \rho^2)} (K_{TT}^x + \sigma_n^2 I). \quad (\text{S.18})$$

Notice that the dependence of $g(\rho^2)$ on ρ^2 is only through $B(\rho^2)$. We differentiate $g(\rho^2)$ with respect to ρ^2 :

$$\frac{dg}{d\rho^2} = \mathbf{v}_1^T A_1^{-1} K_{ST}^x [B(\rho^2)]^{-1} C(\rho^2) [B(\rho^2)]^{-1} K_{TS}^x A_1^{-1} \mathbf{v}_1, \quad (\text{S.19})$$

$$\text{where } C(\rho^2) \stackrel{\text{def}}{=} -\frac{dB}{d\rho^2} = \frac{1 - 2\rho^2}{\rho^2(1 - \rho^2)} D + \frac{1}{\rho^4} K_{TS}^x A_1^{-1} K_{ST}^x = \frac{1}{\rho^2} B - \frac{1}{1 - \rho^2} D. \quad (\text{S.20})$$

Substituting the last expression for $C(\rho^2)$ back into (S.19) gives

$$\frac{dg}{d\rho^2} = \frac{1}{\rho^2} g - \frac{1}{1 - \rho^2} h, \quad (\text{S.21})$$

$$\text{where } h(\rho^2) \stackrel{\text{def}}{=} \mathbf{v}_1^T A_1^{-1} K_{ST}^x [B(\rho^2)]^{-1} D(\rho^2) [B(\rho^2)]^{-1} K_{TS}^x A_1^{-1} \mathbf{v}_1. \quad (\text{S.22})$$

Put inequality $D(\rho^2) \preceq B(\rho^2)$ into $h(\rho^2)$ to give

$$\begin{aligned} h(\rho^2) &\leq \mathbf{v}_1^T A_1^{-1} K_{ST}^x [B(\rho^2)]^{-1} B(\rho^2) [B(\rho^2)]^{-1} K_{TS}^x A_1^{-1} \mathbf{v}_1 \\ &= \mathbf{v}_1^T A_1^{-1} K_{ST}^x [B(\rho^2)]^{-1} K_{TS}^x A_1^{-1} \mathbf{v}_1 \\ &= g(\rho^2). \end{aligned} \quad (\text{S.23})$$

Putting the above inequality into (S.21) leads to

$$\frac{dg}{d\rho^2} \geq f(\rho^2) g(\rho^2), \quad \text{where } f(\rho^2) \stackrel{\text{def}}{=} \frac{1}{\rho^2} - \frac{1}{1 - \rho^2}. \quad (\text{S.24})$$

Since $\underline{\sigma}_T^2(\rho)$ is a lower bound, $g(\rho^2) \geq 0$ (also see the quadratic form in (S.16)). In addition, the multiplicative factor $f(\rho^2)$ is positive for $\rho^2 \in [0, 1/2[$, zero at $\rho^2 = 1/2$, and negative for $\rho^2 \in]1/2, 1]$. Thus $dg/d\rho^2 > 0$ for $\rho^2 \in [0, 1/2[$. Therefore g is monotonically increasing within $\rho^2 \in [0, 1/2[$, and its maximum value must be at $\hat{\rho}^2 \geq 1/2$. \square

S.3 Proof for Proposition 7

Recall that we seek an upper bound $\bar{\sigma}_n^2$ for σ_n^2 such that $\Delta(\rho, \sigma_n^2, \bar{\sigma}_n^2) \leq 0$ for all test locations, where

$$\Delta(\rho, \sigma_n^2, s^2) \stackrel{\text{def}}{=} (\mathbf{k}_*^x)^\top [(\Sigma(1, \sigma_n^2, s^2))^{-1} - (\Sigma(\rho, \sigma_n^2, \sigma_n^2))^{-1}] \mathbf{k}_*^x. \quad (\text{S.25})$$

For the derived bound to be applicable in general, it is necessary that the condition “for all test locations” be equivalent to the condition “for all $\mathbf{k}_*^x \in \mathbb{R}^n$ ”, which is easier to handle; we shall remark on this later. Thus, we start from the requirement that $\Delta(\rho, \sigma_n^2, \bar{\sigma}_n^2) \leq 0$ for all $\mathbf{k}_*^x \in \mathbb{R}^n$:

$$\Delta(\rho, \sigma_n^2, \bar{\sigma}_n^2) \leq 0 \quad \forall \mathbf{k}_*^x \in \mathbb{R}^n \quad (\text{S.26a})$$

$$\Leftrightarrow (\Sigma(1, \sigma_n^2, \bar{\sigma}_n^2))^{-1} - (\Sigma(\rho, \sigma_n^2, \sigma_n^2))^{-1} \preceq 0 \quad (\text{S.26b})$$

$$\Leftrightarrow (\Sigma(1, \sigma_n^2, \bar{\sigma}_n^2))^{-1} \preceq (\Sigma(\rho, \sigma_n^2, \sigma_n^2))^{-1} \quad (\text{S.26c})$$

$$\Leftrightarrow \Sigma(1, \sigma_n^2, \bar{\sigma}_n^2) \succeq \Sigma(\rho, \sigma_n^2, \sigma_n^2) \quad (\text{S.26d})$$

$$\Leftrightarrow \begin{pmatrix} K_{TT}^x & K_{TS}^x \\ K_{ST}^x & K_{SS}^x \end{pmatrix} + \begin{pmatrix} \sigma_n^2 I & 0 \\ 0 & \bar{\sigma}_n^2 I \end{pmatrix} \succeq \begin{pmatrix} K_{TT}^x & K_{TS}^x \\ K_{ST}^x & \rho^{-2} K_{SS}^x \end{pmatrix} + \begin{pmatrix} \sigma_n^2 I & 0 \\ 0 & \rho^{-2} \sigma_n^2 I \end{pmatrix} \quad (\text{S.26e})$$

$$\Leftrightarrow \begin{pmatrix} 0 & 0 \\ 0 & \beta K_{SS}^x \end{pmatrix} \preceq \begin{pmatrix} 0 & 0 \\ 0 & (\bar{\sigma}_n^2 - \rho^{-2} \sigma_n^2) I \end{pmatrix} \quad (\text{S.26f})$$

$$\Leftrightarrow \beta K_{SS}^x \preceq (\bar{\sigma}_n^2 - \rho^{-2} \sigma_n^2) I \quad (\text{S.26g})$$

$$\Leftrightarrow K_{SS}^x \preceq \frac{\bar{\sigma}_n^2 - \rho^{-2} \sigma_n^2}{\beta} I \quad (\text{S.26h})$$

$$\Leftrightarrow \bar{\lambda} \leq \frac{\bar{\sigma}_n^2 - \rho^{-2} \sigma_n^2}{\beta} \quad (\text{S.26i})$$

$$\Leftrightarrow \bar{\sigma}_n^2 \geq \beta \bar{\lambda} + \rho^{-2} \sigma_n^2 = \beta(\bar{\lambda} + \sigma_n^2) + \sigma_n^2 \quad (\text{S.26j})$$

Therefore we have the minimum of the upper bound is

$$\bar{\sigma}_n^2 \stackrel{\text{def}}{=} \beta(\bar{\lambda} + \sigma_n^2). \quad (\text{S.27})$$

The tightness of the bound is evident from the construction of $\bar{\sigma}_n^2$. \square

Remark For the bound to hold in general, we have claimed in the above proof that the condition “for all test locations” must be equivalent to the condition “for all $\mathbf{k}_*^x \in \mathbb{R}^n$ ”. To show this, we shall give a particular example that demands this equivalence. Consider input domain $[-1, 1]$ and $k^x(x, x') = xx'$. We fix x_* . If the observed locations are densely located on $[-1, 1] \setminus \{x_*\}$, then the entries in \mathbf{k}_*^x are densely located on $[-x_*, x_*] \setminus \{x_*^2\}$. Since scaling \mathbf{k}_*^x does not effect the inequality $\Delta(\rho, \sigma_n^2, \bar{\sigma}_n^2) \leq 0$, we have, equivalently, $\mathbf{k}_*^x \in \mathbb{R}^n$.

S.4 Proof for Theorem 12

Theorem 12 upon which Proposition 13 is built depends on the Lemma 4 in Ferrari Trecate et al. [2]. Before we examine a variant of that lemma, some additional notations are needed. Let $\mathbf{y} \stackrel{\text{def}}{=} (y_1 \dots y_n)^\top$ be the n values observed at the set of data locations $X \stackrel{\text{def}}{=} \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$. In addition to the matrix Φ defined in Theorem 12, we introduce an (infinite) vector function $\phi^\top(\mathbf{x}) \stackrel{\text{def}}{=} (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots)^\top$, and an (infinite) diagonal matrix Λ_κ with κ_i on the diagonals.¹ Given the data, consider the space of functions

$$\mathcal{H}_0 \stackrel{\text{def}}{=} \{\phi^\top(\mathbf{x}) L \mathbf{y} \mid L \in \mathbb{R}^{\infty \times n}\}.$$

The aim is to use a function g from \mathcal{H}_0 to estimate the true function f^* drawn from the GP with covariance function $k^x(\cdot, \cdot)$. The quality of this estimation may be evaluated using the following variant of Lemma 4 from [2]:

¹In [2], where the focus is on finite-dimensional approximation of GP, only m major eigenvalues and eigenfunctions are used. However, our focus is to obtain an upper bound on the learning curve, and so we follow [3] and let $m \rightarrow \infty$.

Lemma S.2. (cf. [2, Lemma 4]) The generalization error of a function $g \in \mathcal{H}_0$ is

$$\begin{aligned} \epsilon(g \in \mathcal{H}_0, X) &\stackrel{\text{def}}{=} \left\langle (f^*(\mathbf{x}) - g(\mathbf{x}))^2 \right\rangle_{f^*, \mathbf{y}, \mathbf{x}} \\ &= \sum_{i=1}^{\infty} \kappa_i + \text{tr} \left(L \langle \mathbf{y} \mathbf{y}^T \rangle_{\mathbf{y}} L^T \right) - 2 \text{tr} \left(L \left\langle \langle \mathbf{y} f^*(\mathbf{x}) \rangle_{f^*, \mathbf{y}} \phi^T(\mathbf{x}) \right\rangle_{\mathbf{x}} \right) \end{aligned}$$

Proof.

$$\begin{aligned} &\left\langle (f^*(\mathbf{x}) - g(\mathbf{x}))^2 \right\rangle_{f^*, \mathbf{y}, \mathbf{x}} \\ &= \left\langle \left(f^*(\mathbf{x}) - \phi^T(\mathbf{x}) L \mathbf{y} \right)^2 \right\rangle_{f^*, \mathbf{y}, \mathbf{x}} \\ &= \langle f^*(\mathbf{x}) f^*(\mathbf{x}) \rangle_{f^*, \mathbf{x}} + \left\langle \phi^T(\mathbf{x}) L \mathbf{y} \mathbf{y}^T L^T \phi(\mathbf{x}) \right\rangle_{\mathbf{y}, \mathbf{x}} - 2 \left\langle \phi^T(\mathbf{x}) L \langle \mathbf{y} f^*(\mathbf{x}) \rangle_{f^*, \mathbf{y}} \right\rangle_{\mathbf{x}} \\ &= \langle k^{\mathbf{x}}(\mathbf{x}, \mathbf{x}) \rangle_{\mathbf{x}} + \text{tr} \left(L \langle \mathbf{y} \mathbf{y}^T \rangle_{\mathbf{y}} L^T \left\langle \phi(\mathbf{x}) \phi^T(\mathbf{x}) \right\rangle_{\mathbf{x}} \right) - 2 \text{tr} \left(L \left\langle \langle \mathbf{y} f^*(\mathbf{x}) \rangle_{f^*, \mathbf{y}} \phi^T(\mathbf{x}) \right\rangle_{\mathbf{x}} \right) \\ &= \sum_{i=1}^{\infty} \kappa_i + \text{tr} \left(L \langle \mathbf{y} \mathbf{y}^T \rangle_{\mathbf{y}} L^T \right) - 2 \text{tr} \left(L \left\langle \langle \mathbf{y} f^*(\mathbf{x}) \rangle_{f^*, \mathbf{y}} \phi^T(\mathbf{x}) \right\rangle_{\mathbf{x}} \right) \end{aligned}$$

We have used $\left\langle \phi(\mathbf{x}) \phi^T(\mathbf{x}) \right\rangle_{\mathbf{x}} = I$ for the last expression. \square

To proceed further, it is necessary to specify how \mathbf{y} is obtained. For the single-task GP with isotropic noise, and under correct prior specification, each entry in \mathbf{y} is generated via

$$y(\mathbf{x}) \sim \mathcal{N}(f^*(\mathbf{x}), \sigma_n^2). \quad (\text{S.28})$$

This is the setting considered in [2], where it is also shown that minimizing the generalization error with respect to L leads to the GP mean predictor [2, Theorem 5]; see also [4], and [5, Proposition V.1]. The single-task FWO bound on the learning curve of the GP is obtained by minimizing $\langle \epsilon(g \in \mathcal{H}_0, X) \rangle_X$ with respect to g within only the sub-space of functions

$$\mathcal{H}_1 \stackrel{\text{def}}{=} \{ \phi^T(\mathbf{x}) D \Phi^T \mathbf{y} \mid D \text{ is a diagonal matrix} \}.$$

This results in an upper bound on the learning curve. This is because $\mathcal{H}_1 \subseteq \mathcal{H}_0$, so that minimizing $\langle \epsilon(g \in \mathcal{H}_0, X) \rangle_X$ naturally gives predictors that cannot outperform the GP mean predictor. The form of functions in \mathcal{H}_1 is motivated by Projected Bayes Regression [2, Definition 1], wherein the L in \mathcal{H}_0 is constrained to be $M \Phi^T$ for a square matrix M .

Theorem 12, our variant of Ferrari Trecate et al.'s Theorem 6 for correlated noise γ^2 , for example as defined by (14), uses

$$y(\mathbf{x}) \sim \mathcal{GP}(f^*(\mathbf{x}), \gamma^2(\mathbf{x}, \mathbf{x}')) \quad (\text{S.29})$$

instead of (S.28). Generating \mathbf{y} in this manner and further restricting g to be from \mathcal{H}_1 leads to the following generalization error using Lemma S.2:

$$\epsilon(g \in \mathcal{H}_1, X) = \sum_{i=1}^{\infty} \kappa_i + \text{tr} (D \Phi^T H \Phi D) - 2 \text{tr} (D \Phi^T \Phi \Lambda_{\kappa}), \quad (\text{S.30})$$

where

$$H_{ij} \stackrel{\text{def}}{=} k^{\mathbf{x}}(\mathbf{x}_i, \mathbf{x}_j) + \gamma^2(\mathbf{x}_i, \mathbf{x}_j). \quad (\text{S.31})$$

The right of (S.30) is quadratic in the diagonal entries of D . Minimizing $\langle \epsilon(g \in \mathcal{H}_1, X) \rangle_X$ with respect to D results in Theorem 12.

S.5 Proof for Proposition 13

We prove the expression for c_i given by (15), which forms the crux of Proposition 13. In order to present the proof, some notations are required. Recall that the cardinality of these sets are $|X| = n$,

$|X_S| = n_S$ and $|X_T| = n_T$, and that $\pi_S \stackrel{\text{def}}{=} n_S/n$. We partition the index set $\mathcal{I} \stackrel{\text{def}}{=} \{1 \dots n\}$ into $\mathcal{I}_T \stackrel{\text{def}}{=} \{1 \dots n_T\}$, and $\mathcal{I}_S \stackrel{\text{def}}{=} \{(n_T + 1) \dots n\}$. We enumerate and order the elements of X so that $X = \{\mathbf{x}_i\}_{i=1}^n$, $X_T = \{\mathbf{x}_i \mid i \in \mathcal{I}_T\}$, and $X_S = \{\mathbf{x}_i \mid i \in \mathcal{I}_S\}$. We shall represent by $\langle \dots \rangle_X$ the expectation over the set X , and write $p(X)dX \stackrel{\text{def}}{=} \prod_{i=1}^n p(\mathbf{x}_i)d\mathbf{x}_i$, where the distributions over the \mathbf{x}_i s are identical; expressions $\langle \dots \rangle_{X_S}$ and $\langle \dots \rangle_{X_T}$, and $p(X_S)dX_S$ and $p(X_T)dX_T$ have similar meanings.

Recall the definition of eigenvalues and eigenfunctions using the integral equation and the orthogonality of eigenfunctions:

$$\int k^x(\mathbf{x}, \mathbf{x}') \phi_i(\mathbf{x}') p(\mathbf{x}') d\mathbf{x}' = \kappa_i \phi_i(\mathbf{x}) \quad \int \phi_i(\mathbf{x}) \phi_i(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = 1. \quad (\text{S.32})$$

Using these two equalities, we can show that

$$\begin{aligned} \int k^x(\mathbf{x}, \mathbf{x}') \phi_i(\mathbf{x}) \phi_i(\mathbf{x}') p(\mathbf{x}) p(\mathbf{x}') d\mathbf{x} d\mathbf{x}' &= \int \left(\int k^x(\mathbf{x}, \mathbf{x}') \phi_i(\mathbf{x}') p(\mathbf{x}') d\mathbf{x}' \right) \phi_i(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\ &= \int (\kappa_i \phi_i(\mathbf{x})) \phi_i(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\ &= \kappa_i \int \phi_i(\mathbf{x}) \phi_i(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\ &= \kappa_i \end{aligned} \quad (\text{S.33})$$

The next two equalities will be used in the main part of the proof:

$$\begin{aligned} \left\langle \sum_{p,q \in \mathcal{I}} \delta_{pq} \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) \right\rangle_X &= \int \sum_{p=1}^n \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_p) p(X) dX \\ &= \sum_{p=1}^n \int \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_p) p(\mathbf{x}_p) d\mathbf{x}_p \\ &= n \int [\phi_i(\mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x} \\ &= n \end{aligned} \quad (\text{S.34})$$

$$\begin{aligned} \left\langle \sum_{p,q \in \mathcal{I}} k^x(\mathbf{x}_p, \mathbf{x}_q) \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) \right\rangle_X &= \int \sum_{p,q \in \mathcal{I}} k^x(\mathbf{x}_p, \mathbf{x}_q) \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) p(X) dX \\ &= \sum_{\substack{p,q \in \mathcal{I} \\ p \neq q}} \int k^x(\mathbf{x}_p, \mathbf{x}_q) \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) p(\mathbf{x}_p) p(\mathbf{x}_q) d\mathbf{x}_p d\mathbf{x}_q \\ &\quad + \sum_{p=1}^n \int k^x(\mathbf{x}_p, \mathbf{x}_p) \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_p) p(\mathbf{x}_p) d\mathbf{x}_p \\ &= \sum_{\substack{p,q \in \mathcal{I} \\ p \neq q}} \kappa_i + \sum_{p=1}^n \int k^x(\mathbf{x}, \mathbf{x}) [\phi_i(\mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x} \quad (*) \\ &= n(n-1)\kappa_i + n \int k^x(\mathbf{x}, \mathbf{x}) [\phi_i(\mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x}, \end{aligned} \quad (\text{S.35})$$

where (S.33) is used in getting to (*). By similar arguments we can show equivalent results when the summations and expectations are taken only over data points in X_S and X_T .

We now turn to the main part of the proof. The observation noise (co)variance for our upper bound can be expressed as

$$\gamma^2(\mathbf{x}_p, \mathbf{x}_q) \stackrel{\text{def}}{=} \begin{cases} \delta_{pq} \sigma_n^2 & \text{if } p \in \mathcal{I}_T \text{ and } q \in \mathcal{I}_T \\ \beta k^x(\mathbf{x}_p, \mathbf{x}_q) + \rho^{-2} \delta_{pq} \sigma_n^2 & \text{if } p \in \mathcal{I}_S \text{ and } q \in \mathcal{I}_S \\ 0 & \text{otherwise,} \end{cases} \quad (\text{S.36})$$

where δ_{pq} is the Kronecker delta function. In the definition of γ^2 , the first case is when both input locations are for task T , and the second case is when both input locations are for task S . It follows that

$$\begin{aligned}
& \left\langle \sum_{p,q \in \mathcal{I}} \gamma^2(\mathbf{x}_p, \mathbf{x}_q) \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) \right\rangle_X \\
&= \left\langle \sum_{p,q \in \mathcal{I}_T} \delta_{pq} \sigma_n^2 \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) \right\rangle_{X_T} + \left\langle \sum_{p,q \in \mathcal{I}_S} (\beta k^x(\mathbf{x}_p, \mathbf{x}_q) + \rho^{-2} \delta_{pq} \sigma_n^2) \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) \right\rangle_{X_S} \\
&= \sigma_n^2 \left\langle \sum_{p,q \in \mathcal{I}_T} \delta_{pq} \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) \right\rangle_{X_T} + \beta \left\langle \sum_{p,q \in \mathcal{I}_S} k^x(\mathbf{x}_p, \mathbf{x}_q) \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) \right\rangle_{X_S} \\
&\quad + \rho^{-2} \sigma_n^2 \left\langle \sum_{p,q \in \mathcal{I}_S} \delta_{pq} \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) \right\rangle_{X_S} \\
&= n_T \sigma_n^2 + \beta \left[n_S(n_S - 1) \kappa_i + \beta n_S \int k^x(\mathbf{x}, \mathbf{x}) [\phi_i(\mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x} \right] + \rho^{-2} n_S \sigma_n^2 \quad (*) \\
&= \beta n_S(n_S - 1) \kappa_i + \beta n_S \int k^x(\mathbf{x}, \mathbf{x}) [\phi_i(\mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x} + n \sigma_n^2 + \beta n_S \sigma_n^2, \quad (\text{S.37})
\end{aligned}$$

where the X_S and X_T equivalents of (S.34) and (S.35) are applied to get (*). We are now ready to apply Theorem 12 to obtain c_i :

$$\begin{aligned}
c_i &= \frac{1}{n} \langle (\Phi^T H \Phi)_{ii} \rangle_X \\
&= \frac{1}{n} \left\langle \sum_{p,q \in \mathcal{I}} (k^x(\mathbf{x}_p, \mathbf{x}_q) + \gamma^2(\mathbf{x}_p, \mathbf{x}_q)) \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) \right\rangle_X \\
&= \frac{1}{n} \left\langle \sum_{p,q \in \mathcal{I}} k^x(\mathbf{x}_p, \mathbf{x}_q) \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) \right\rangle_X + \frac{1}{n} \left\langle \sum_{p,q \in \mathcal{I}} \gamma^2(\mathbf{x}_p, \mathbf{x}_q) \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) \right\rangle_X \\
&= (n-1) \kappa_i + \int k^x(\mathbf{x}, \mathbf{x}) [\phi_i(\mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x} \quad (*) \\
&\quad + \beta \pi_S (n_S - 1) \kappa_i + \beta \pi_S \int k^x(\mathbf{x}, \mathbf{x}) [\phi_i(\mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x} + \sigma_n^2 + \beta \pi_S \sigma_n^2 \\
&= [(n-1) + \beta n_S (n_S - 1)] \kappa_i + (1 + \beta \pi_S) \int k^x(\mathbf{x}, \mathbf{x}) [\phi_i(\mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x} + (1 + \beta \pi_S) \sigma_n^2 \\
&= [(1 + \beta \pi_S^2) n - (1 + \beta \pi_S)] \kappa_i + (1 + \beta \pi_S) \int k^x(\mathbf{x}, \mathbf{x}) [\phi_i(\mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x} + (1 + \beta \pi_S) \sigma_n^2
\end{aligned}$$

where (S.35) and (S.37) are used to obtain (*). \square

Remark. Although we do not have the proof by Ferrari Trecate et al. [2] for their upper bound on the learning curve for single-task GP with isotropic noise, it is conceivable that some variation of the above proof has been used by them.

S.6 Simulations of learning curve, details

In this section, we give additional details for our simulations of the learning curve.

S.6.1 Continuation of ϵ_T^{avg} in $\pi_S n$

Recall that $\epsilon_T^{\text{avg}}(\rho, \sigma_n^2, \pi_S, n)$ is only defined for values of π_S and n such that $\pi_S n = n_S \in \mathbb{N}_0$. In our simulations, however, we extend the domain to allow $\pi_S n \in \mathbb{R}$, so that smooth curves can be plotted. For the theoretical bounds given by Propositions 11 and 13, this is done by simply using

the respective expressions verbatim. For the experimental bounds, which require sampling over $X \stackrel{\text{def}}{=} X_T \cup X_S$, this is achieved in the manner described next.

For a given π_S , we sample the sizes of the training sets X_T and X_S to satisfy

$$\lfloor \pi_S n \rfloor \leq n_S \leq \lceil \pi_S n \rceil \quad \text{and} \quad \langle n_S \rangle = \pi_S n, \quad (\text{S.38})$$

where the expectation is taken over simulation runs. The first condition ensures that, within each simulation run, the size n_S of X_S is $\pi_S n$ whenever the latter is an integer. The second condition ensures that the ratio n_S/n is consistent with π_S when averaged over multiple simulation runs. For each simulation run, the training set is constructed sequentially by randomly drawing additional training locations. For each new location, we determine its task by using Algorithm 1.

Algorithm 1 Decide the task for a new input

Require: Ratio π_S , required cardinality n of X , and previous cardinality n_S^{old} of X_S .

- 1: **if** $n_S^{\text{old}} < \lfloor \pi_S n \rfloor$ **then**
 - 2: new input is for task S
 - 3: **else if** $n_S^{\text{old}} = \lceil \pi_S n \rceil$ **then**
 - 4: new input is for task T
 - 5: **else**
 - 6: new input is for task S with probability $(\pi_S n - \lfloor \pi_S n \rfloor)$
 {or, equivalently, for task T with probability $1 - (\pi_S n - \lfloor \pi_S n \rfloor)$ }
 - 7: **end if**
-

S.6.2 Analytical averaging over test locations

The simulation study in section 5.3 uses the squared exponential (SE) covariance function with normally distributed inputs. As claimed in section 3.4, the expectation over test locations can be done analytically to obtain the generalization error ϵ_T exactly. In order to do this, we need to be able to compute

$$M_{pq} \stackrel{\text{def}}{=} \int k^x(\mathbf{x}_p, \mathbf{x}_*) k^x(\mathbf{x}_q, \mathbf{x}_*) p(\mathbf{x}_*) d\mathbf{x}_* = \sum_{i=1}^{\infty} \kappa_i^2 \phi_i(\mathbf{x}_p) \phi_i(\mathbf{x}_q) \quad (\text{S.39})$$

for a fixed pair of input locations $(\mathbf{x}_p, \mathbf{x}_q)$. In the one-dimensional case with the SE covariance function and normally distributed inputs, i.e.,

$$k^x(x, x') \stackrel{\text{def}}{=} \exp -\frac{(x - x')^2}{2l^2} \quad \text{and} \quad p(x) \stackrel{\text{def}}{=} \frac{1}{\sqrt{2\pi}\sigma_x} \exp -\frac{x^2}{2\sigma_x^2}, \quad (\text{S.40})$$

we can use the integral expression for M_{pq} to obtain

$$M_{pq} = \frac{l}{\sqrt{2\sigma_x^2 + l^2}} \exp -\frac{\sigma_x^2(x_p - x_q)^2 + l^2(x_p^2 + x_q^2)}{2l^2(2\sigma_x^2 + l^2)}. \quad (\text{S.41})$$

This can be easily generalized to input spaces of higher-dimensions. Note that the infinite sum expression for M_{pq} is not useful in this case, even though analytic expressions for the eigenfunctions are available [6]. This is because the eigenfunctions corresponding to the smaller eigenvalues exhibit larger oscillations around zero in regions where $p(x)$ is low, so that it is hard to determine when to truncate the infinite sum for certain pairs of (x_p, x_q) .

References

- [1] Simo Puntanen and George P. H. Styan. Historical introduction: Issai Schur and the early development of the Schur complement. In Fuzhen Zhang, editor, *The Schur complement and its applications*, Numerical Methods and Algorithms, pages 1–16. Springer, 2005.
- [2] Giancarlo Ferrari Trecate, Christopher K. I. Williams, and Manfred Opper. Finite-dimensional approximation of Gaussian processes. In Michael J. Kearns, Sara A. Solla, and David A. Cohn, editors, *Advances in Neural Information Processing Systems 11*, pages 218–224. The MIT Press, 1999.

- [3] Peter Sollich and Anason Halees. Learning curves for Gaussian process regression: Approximations and bounds. *Neural Computation*, 14(6):1393–1428, 2002.
- [4] George S. Kimeldorf and Grace Wahba. A correspondence between Bayesian estimation on stochastic processes and smoothing by splines. *The Annals of Mathematical Statistics*, 41(2): 495–502, 1970.
- [5] Klaus Ritter. *Average-Case Analysis of Numerical Problems*, volume 1733 of *Lecture Notes in Mathematics*. Springer, 2000.
- [6] Huaiyu Zhu, Christopher K. I. Williams, Richard Rohwer, and Michal Morciniec. Gaussian regression and optimal finite dimensional linear models. In Christopher M. Bishop, editor, *Neural Networks and Machine Learning*, volume 168 of *NATO ASI Series F: Computer and Systems Sciences*, pages 167–184. Springer-Verlag, Berlin, 1998.