

---

# Computation-Aware Robust Gaussian Processes

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Gaussian Processes (GPs) are flexible nonparametric statistical models equipped  
2 with principled uncertainty quantification for both noise and model uncertainty.  
3 However, their cubic inference complexity requires them to be combined with  
4 approximation techniques when applied to large datasets. Recent work demon-  
5 strated that such approximations introduce an additional source of uncertainty,  
6 *computational uncertainty*, and that the latter could be quantified, leading to the  
7 *computation-aware* GP, also known as IterGP. In this short communication, we  
8 demonstrate that IterGP is not “robust”, in the sense that a quantity of interest,  
9 the posterior influence function, is not bounded. Subsequently, drawing inspira-  
10 tion from recent work on Robust Conjugate GPs, we introduce a novel class of  
11 GPs: IterRCGPs. We carry out a number of theoretical analyses, demonstrating  
12 the robustness of IterRCGPs among other things.

## 13 1 Introduction

14 Gaussian Processes (GPs, [Rasmussen and Williams \(2006\)](#)) are a class of probabilistic models en-  
15 joying many properties such as universal approximation or closed-form computations. Due to their  
16 principled uncertainty quantification, they are becoming increasingly popular when applied in high-  
17 stakes domains like medical datasets ([Cheng et al., 2019](#); [Chen et al., 2023](#)) or used as a surrogate  
18 model in Bayesian Optimization ([Garnett, 2023](#)). This being said, GPs suffer from a cubic infer-  
19 ence complexity, hindering their use on large datasets. As a remedy, approximation techniques  
20 like Sparse Variational Gaussian Processes ([Titsias, 2009](#)) or the Nyström approximation are often  
21 used ([Williams and Seeger, 2000](#); [Wild et al., 2023](#)).

22 These approximations introduce bias in uncertainty quantification, which, as recently demonstrated,  
23 can be quantified and combined with mathematical uncertainty, leading to the development of  
24 *computation-aware* GPs ([Wenger et al., 2022](#)), also known as IterGPs. This combined uncertainty  
25 is shown to be the correct measure for capturing overall uncertainty, as limited computation intro-  
26 duces computational error. While this analysis applies to standard GPs, many practical applications  
27 require variations, e.g., to deal with heteroscedasticity or outliers.

28 Recent work by [Altamirano et al. \(2024\)](#) introduced the robust conjugate GP (RCGP), which unifies  
29 three classes of GPs. RCGP retains conjugacy, enabling a closed-form posterior while exhibiting  
30 a robustness property. However, like standard GPs, RCGP faces significant inference complexity,  
31 necessitating approximation methods such as sparse variational RCGP, and therefore suggesting the  
32 use of the framework developed by [Wenger et al. \(2022\)](#).

33 **Contributions.** Our work can be seen as bridging the gap between computation-aware GPs and  
34 Robust Conjugate GPs. As such, our contributions are mainly theoretical and can be summarized as  
35 follows:

- 36 • We present IterRCGP, a novel computation-aware Gaussian Process (GP) framework that  
37 extends IterGPs by accommodating a broader range of observation noise models.
- 38 • We demonstrate that IterRCGP inherits the robustness properties characteristic of RCGP.
- 39 • We establish that IterRCGP exhibits convergence behavior and worst-case errors analogous  
40 to IterGP.

## 41 2 Preliminaries

42 We first introduce notations for GP regression (Rasmussen and Williams, 2006). Let  $\mathcal{D} =$   
43  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  be a dataset, with  $(\mathbf{x}_j, y_j) \in \mathbb{R}^d \times \mathbb{R}$  such that  $y_j = f(\mathbf{x}_j) + \epsilon$  and  
44  $\epsilon \sim \mathcal{N}(0, \sigma_{\text{noise}}^2)$  for all  $j$ . The latent function  $f$  is modeled with a GP prior:

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')). \quad (1)$$

45 This defines a distribution over functions  $f$  whose mean is  $\mathbb{E}[f(\mathbf{x})] = m(\mathbf{x})$  and covariance  
46  $\text{cov}[f(\mathbf{x}), f(\mathbf{x}')] = k(\mathbf{x}, \mathbf{x}')$ .  $k$  is a kernel function measuring the similarity between in-  
47 puts. For any finite-dimensional collection of inputs  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , the function values  $\mathbf{f} =$   
48  $[f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)]^\top \in \mathbb{R}^n$  follow a multivariate normal distribution  $\mathbf{f} \sim \mathcal{N}(\mathbf{m}, \mathbf{K})$ , where  
49  $\mathbf{m} = [m(\mathbf{x}_1), \dots, m(\mathbf{x}_n)]^\top$  and  $\mathbf{K} \in \mathbb{R}^{n \times n} = [k(\mathbf{x}_j, \mathbf{x}_l)]_{1 \leq j, l \leq n}$  is the kernel matrix.

50 Given  $\mathcal{D}$ , the posterior predictive distribution  $p(f(\mathbf{x}) \mid \mathcal{D})$  is Gaussian for all  $\mathbf{x}$  with mean  $\mu_*(\mathbf{x})$   
51 and variance  $k_*(\mathbf{x}, \mathbf{x})$ , such that

$$\begin{aligned} \mu_*(\mathbf{x}) &= m(\mathbf{x}) + \mathbf{k}_\mathbf{x}^\top (\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{I})^{-1} (\mathbf{y} - \mathbf{m}), \\ k_*(\mathbf{x}, \mathbf{x}) &= k(\mathbf{x}, \mathbf{x}) - \mathbf{k}_\mathbf{x}^\top (\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{I})^{-1} \mathbf{k}_\mathbf{x}, \end{aligned}$$

52 where  $\mathbf{y} = [y_1, \dots, y_n] \in \mathbb{R}^n$  and  $\mathbf{k}_\mathbf{x} = [k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n)]^\top \in \mathbb{R}^n$ .

53 Next, we introduce an extension of GPs: Robust Conjugate Gaussian Processes (RCGPs).

54 **Robust conjugate Gaussian process.** RCGP follows the generalized Bayesian inference frame-  
55 work, substituting the classical likelihood with the loss function  $L_n^w$  (Altamirano et al., 2024) defined  
56 as

$$L_n^w(\mathbf{f}, \mathbf{x}, \mathbf{y}) = \frac{1}{n} \left( \sum_{j=1}^n w^2(\mathbf{x}_j, y_j) s_{\text{model}}^2(\mathbf{x}_j, y_j) + 2 \nabla_y [w^2(\mathbf{x}_j, y_j) s_{\text{model}}(\mathbf{x}_j, y_j)] \right), \quad (2)$$

57 where  $s_{\text{model}}(\mathbf{x}, y) = \sigma_{\text{noise}}^{-2} (f(\mathbf{x}) - y)$ ,  $\sigma_{\text{noise}}^2 > 0$ . The core component of  $L_n^w$  is the weighting  
58 function  $w$ , which depends on  $\mathbf{x}$  and  $y$ . Altamirano et al. (2024)[Table 1] provides three weighting  
59 functions corresponding to homoscedastic, heteroscedastic, and outliers-robust GPs. Building on  
60  $L_n^w$ , the authors further define the RCGP’s predictive posterior distribution  $p^w(f(\mathbf{x}) \mid \mathcal{D})$  as follows:

$$\hat{\mu}_*(\mathbf{x}) = m(\mathbf{x}) + \mathbf{k}_\mathbf{x}^\top \overbrace{(\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{J}_\mathbf{w})^{-1} (\mathbf{y} - \mathbf{m}_\mathbf{w})}^{\hat{\mathbf{v}}} \quad (3)$$

$$\hat{k}_*(\mathbf{x}, \mathbf{x}) = k(\mathbf{x}, \mathbf{x}) - \mathbf{k}_\mathbf{x}^\top \tilde{\mathbf{K}}^{-1} \mathbf{k}_\mathbf{x} \quad (4)$$

61 for  $\mathbf{w} = (w(\mathbf{x}_1, y_1), \dots, w(\mathbf{x}_n, y_n))^\top$ ,  $\tilde{\mathbf{K}} = \mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{J}_\mathbf{w}$ ,  $\mathbf{m}_\mathbf{w} = \mathbf{m} + \sigma_{\text{noise}}^2 \nabla_y \log(\mathbf{w}^2)$ , and  
62  $\mathbf{J}_\mathbf{w} = \text{diag}(\frac{\sigma_{\text{noise}}^2}{2} \mathbf{w}^{-2})$ . A key advantage of RCGP is its robustness to outliers and non-Gaussian  
63 errors. While vanilla GPs exhibit an unbounded posterior influence function, RCGP, under certain  
64 conditions, maintains a bounded posterior influence function (Altamirano et al., 2024)[Proposition  
65 3.2].

## 66 3 Computation-aware RCGPs

67 In the same spirit of Wenger et al. (2022), we treat the representer weights  $\hat{\mathbf{v}}$  introduced in Equation 3  
68 as a random variable with the prior  $p(\hat{\mathbf{v}}) = \mathcal{N}(\hat{\mathbf{v}}; \mathbf{0}, \tilde{\mathbf{K}}^{-1})$ . We then update  $p(\hat{\mathbf{v}})$  by iteratively

69 applying the tractable matrix-vector multiplication. For a particular iteration  $i \in \{0, \dots, n\}$ , we  
 70 have the current belief distribution  $p(\hat{\mathbf{v}}) = \mathcal{N}(\hat{\mathbf{v}}; \tilde{\mathbf{v}}_i, \tilde{\Sigma}_i)$  where

$$\tilde{\mathbf{v}}_i = \tilde{\mathbf{v}}_{i-1} + \tilde{\Sigma}_{i-1} \tilde{\mathbf{K}} \mathbf{s}_i (\mathbf{s}_i^\top \tilde{\mathbf{K}} \tilde{\Sigma}_{i-1} \tilde{\mathbf{K}} \mathbf{s}_i)^{-1} \tilde{\alpha}_i = \tilde{\mathbf{C}}_i (\mathbf{y} - \mathbf{m}_w) \quad (5)$$

$$\tilde{\Sigma}_i = \tilde{\Sigma}_{i-1} - \tilde{\Sigma}_{i-1} \tilde{\mathbf{K}} \mathbf{s}_i (\mathbf{s}_i^\top \tilde{\mathbf{K}} \tilde{\Sigma}_{i-1} \tilde{\mathbf{K}} \mathbf{s}_i)^{-1} \mathbf{s}_i^\top \tilde{\mathbf{K}} \tilde{\Sigma}_{i-1} \quad (6)$$

$$\tilde{\alpha}_i = \mathbf{s}_i^\top \underbrace{\tilde{\mathbf{K}}(\hat{\mathbf{v}} - \tilde{\mathbf{v}}_{i-1})}_{\mathbf{r}_{i-1}} \quad (7)$$

$$\tilde{\mathbf{C}}_i = \tilde{\mathbf{K}}^{-1} - \tilde{\Sigma}_i \quad (8)$$

71 Here,  $\mathbf{s}_i$  denotes the policy corresponding to a specific approximation method (Wenger *et al.*,  
 72 2022)[Table 1]. This policy serves as the projection of the residual  $\mathbf{r}_{i-1}$  results in  $\alpha_i$ . The belief  
 73 regarding the representer weights encodes the computational error as an added source of uncertainty,  
 74 which can be integrated with the inherent uncertainty of the mathematical posterior.

75 We obtain the predictive posterior of IterRCGP by integrating out the representer weights:  
 76  $p(f(\mathbf{x})|\mathcal{D}) = \int p(f(\mathbf{x})|\hat{\mathbf{v}})p(\hat{\mathbf{v}})d\hat{\mathbf{v}} = \mathcal{N}(f; \hat{\mu}_i(\mathbf{x}), \hat{k}_i(\mathbf{x}, \mathbf{x}))$  where

$$\hat{\mu}_i(\mathbf{x}) = m(\mathbf{x}) + \mathbf{k}_*^\top \tilde{\mathbf{v}}_i \quad (9)$$

$$\hat{k}_i(\mathbf{x}, \mathbf{x}) = k(\mathbf{x}, \mathbf{x}) - \mathbf{k}_x^\top \tilde{\mathbf{K}}^{-1} \mathbf{k}_x + \underbrace{\mathbf{k}_x^\top \tilde{\Sigma}_i \mathbf{k}_x}_{k_i^{\text{comp}}(\mathbf{x}, \mathbf{x})} = k(\mathbf{x}, \mathbf{x}) - \underbrace{\mathbf{k}_x^\top \tilde{\mathbf{C}}_i \mathbf{k}_x}_{\text{combined uncertainty}} \quad (10)$$

77 IterRCGP follows [Algorithm 1] from Wenger *et al.* (2022) to compute an estimate weights  $\tilde{\mathbf{v}}_i$  and  
 78 the rank- $i$  precision matrix approximation  $\tilde{\mathbf{C}}_i$ .

## 79 4 Theoretical results

80 In this section, we present the theoretical properties of IterRCGP, building upon the IterGP frame-  
 81 work and the RCGP class. Our theoretical analysis primarily aims to establish the following key  
 82 results:

- 83 • Robustness property of IterGP and IterRCGP (Proposition 1).
- 84 • Convergence of IterRCGP’s posterior mean in reproducing kernel Hilbert space (RKHS)  
 85 norm (Proposition 2) and pointwise (Corollary 4).
- 86 • Combined uncertainty of IterRCGP is a tight worst-case bound on the relative distance  
 87 to all potential latent functions shifted by the function  $\mathbf{m}_w$  consistent with computational  
 88 observations, similar to its IterGP counterpart (Proposition 3).

89 We establish the robustness properties of IterGP and IterRCGP using the Posterior Influence Func-  
 90 tion (PIF) as the robustness criterion. Appendix 1 provides a detailed definition of PIF. The propo-  
 91 sition presented below is closely related to Altamirano *et al.* (2024)[Proposition 3.2].

92 **Proposition 1.** (*Robustness property*)  
 93 Suppose  $f \sim \mathcal{GP}(m, k)$ ,  $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma_{\text{noise}}^2 \mathbf{I})$  and let  $C'_k \in \mathbb{R}; k = 1, 2, 3$  be constants independent  
 94 of  $y_m^c$ . For any given iteration  $i \in \{0, \dots, n\}$ , IterGP regression has the PIF

$$\text{PIF}_{\text{IterGP}}(y_m^c, \mathcal{D}, i) = C'_1 (y_m - y_m^c)^2 \quad (11)$$

95 which is not robust:  $\text{PIF}_{\text{IterGP}}(y_m^c, \mathcal{D}, i) \rightarrow \infty$  as  $|y_m^c| \rightarrow \infty$ . In contrast, for the IterRCGP with  
 96  $\sup_{\mathbf{x}, y} w(\mathbf{x}, y) < \infty$ ,

$$\text{PIF}_{\text{IterRCGP}}(y_m^c, \mathcal{D}, i) = C'_2 (w(x_n, y_n^c)^2 y_n^c)^2 + C'_3 \quad (12)$$

97 Therefore, if  $\sup_{\mathbf{x}, y} y w(\mathbf{x}, y)^2 < \infty$ , IterRCGP regression is robust since  
 98  $\sup_{y_m^c} |\text{PIF}_{\text{IterRCGP}}(y_m^c, \mathcal{D}, i)| < \infty$ .

99 The proposition demonstrates that IterGP and IterRCGP inherit the same robustness properties as  
 100 their respective counterparts, GP and RCGP. Specifically, the condition  $\sup_{\mathbf{x}, y} w(\mathbf{x}, y) < \infty$  ensures  
 101 each observation has a finite weight, which is the key factor underpinning robustness.

102 The following proposition is analogous to [Theorem 1] in Wenger *et al.* (2022).

103 **Proposition 2.** (Convergence in RKHS norm of the robust posterior mean approximation)  
 104 Let  $\mathcal{H}_k$  be the RKHS w.r.t. kernel  $k$ ,  $\sigma_{\text{noise}}^2 > 0$  and let  $\hat{\boldsymbol{\mu}}_* - \mathbf{m} \in \mathcal{H}_k$  be the unique solution to  
 105 following minimization problem

$$\operatorname{argmin}_{f \in \mathcal{H}_k} L_n^w(\mathbf{f}, \mathbf{x}, \mathbf{y}) + \frac{1}{2n} \|\mathbf{f}\|_{\mathcal{H}_k}^2 \quad (13)$$

106 which is equivalent to the mathematical RCGP mean posterior shifted by prior mean  $\mathbf{m}$ . Then for  
 107  $i \in \{0, \dots, n\}$  the IterRCGP posterior mean  $\hat{\boldsymbol{\mu}}_i$  satisfies:

$$\|\hat{\boldsymbol{\mu}}_* - \hat{\boldsymbol{\mu}}_i\|_{\mathcal{H}_k} \leq \hat{\rho}(i) c(\mathbf{J}_w) \|\hat{\boldsymbol{\mu}}_* - \mathbf{m}\|_{\mathcal{H}_k} \quad (14)$$

108 where  $\hat{\rho}$  is the relative bound errors corresponding to the number of iterations  $i$  and the constant  
 109  $c(\mathbf{J}_w) = \sqrt{1 + \frac{\lambda_{\max}(\mathbf{J}_w)}{\lambda_{\min}(\mathbf{K})}} \rightarrow 1$  as  $\lambda_{\max}(\mathbf{J}_w) \rightarrow 0$ .

110 Appendix B provides more details about the relative bound errors. Proposition 2 provides a bound  
 111 on the RKHS-norm error between the posterior mean of IterRCGP and the mathematical posterior  
 112 mean of RCGP.

113 The final proposition parallels [Theorem 2] in Wenger *et al.* (2022), demonstrating that the combined  
 114 uncertainty is a tight bound for all functions  $g$  that could have yielded the same computational  
 115 outcomes.

116 **Proposition 3.** (Combined and computational uncertainty as worst-case errors)

117 Let  $\sigma_{\text{noise}}^2 \geq 0$  and  $\hat{k}_i(\cdot, \cdot) = \hat{k}_*(\cdot, \cdot) + k_i^{\text{comp.}}(\cdot, \cdot)$  be the combined uncertainty of IterRCGP. Fur-  
 118 thermore, let  $\mathbf{g} = [g(\mathbf{x}_1), \dots, g(\mathbf{x}_n)] \in \mathbb{R}^n$ . Then, for any new  $\mathbf{x} \in \mathcal{X}$  we have

$$\sup_{\|g - m_w\|_{\mathcal{H}_{k\sigma_w}} \leq 1} \underbrace{g(\mathbf{x}) - \hat{\mu}^g(\mathbf{x})}_{\text{math. err.}} + \underbrace{\hat{\mu}^g(\mathbf{x}) - \hat{\mu}_i^g(\mathbf{x})}_{\text{comp. err.}} = \sqrt{\hat{k}_i(\mathbf{x}, \mathbf{x}) + \sigma_{\text{noise}}^2} \quad (15)$$

$$\sup_{\|g - m_w\|_{\mathcal{H}_{k\sigma_w}} \leq 1} \underbrace{\hat{\mu}^g(\mathbf{x}) - \hat{\mu}_i^g(\mathbf{x})}_{\text{comp. err.}} = \sqrt{k_i^{\text{comp.}}(\mathbf{x}, \mathbf{x})} \quad (16)$$

119 where  $\hat{\mu}^g(\cdot) = k(\cdot, \mathbf{X}) \tilde{\mathbf{K}}^{-1}(\mathbf{g} - \mathbf{m}_w)$  is the RCGP's posterior and  $\hat{\mu}_i^g(\cdot) = k(\cdot, \mathbf{X}) \tilde{\mathbf{C}}_i(\mathbf{g} - \mathbf{m}_w)$   
 120 IterRCGP's posterior mean for the latent function  $g$  and the function  $m_w$  lies in  $\mathcal{H}_{k\sigma_w}$ .

121 The consequence of Proposition 3 is then formalized through the following corollary:

122 **Corollary 4.** (Pointwise convergence of robust posterior mean)

123 Assume the conditions of Proposition 3 hold and assume the latent function  $g \in \mathcal{H}_{k\sigma_w}$ . Let  $\hat{\boldsymbol{\mu}}$  be the  
 124 corresponding mathematical RCGP posterior mean and  $\hat{\boldsymbol{\mu}}_i$  the IterRCGP posterior mean. It holds  
 125 that

$$\frac{|g(\mathbf{x}) - \hat{\mu}_i(\mathbf{x})|}{\|g\|_{\mathcal{H}_{k\sigma_w}}} \leq \sqrt{\hat{k}_i(\mathbf{x}, \mathbf{x}) + \sigma_{\text{noise}}^2} \quad (17)$$

$$\frac{\hat{\mu}(\mathbf{x}) - \hat{\mu}_i(\mathbf{x})}{\|g\|_{\mathcal{H}_{k\sigma_w}}} \leq \sqrt{k_i^{\text{comp.}}(\mathbf{x}, \mathbf{x})} \quad (18)$$

## 126 5 Conclusion

127 In this paper, we demonstrated that computation-aware GPs as presented by Wenger *et al.* (2022)  
 128 lack robustness in the PIF sense. Subsequently, we introduced Iter RCGPs, a novel class of provably  
 129 robust computation-aware GPs. Since our work mainly involves theoretical analyses, our immediate  
 130 perspective is to run numerical experiments using synthetic and real-world datasets. Next, one  
 131 interesting avenue for applying Iter RCGPs is that of Bayesian Optimization (BO), a domain where  
 132 uncertainty quantification is key to coming up with good exploration policies.

133 Indeed, the issue of refined uncertainty quantification has recently gained attention in BO. One ap-  
 134 proach addresses this by jointly optimizing the selection of the optimal data point along with the  
 135 SVGP parameters and the locations of the inducing points (Maus *et al.*, 2024). Another study incor-  
 136 porates conformal prediction into BO by leveraging the conformal Bayes posterior and proposing  
 137 generalized versions of the corresponding BO acquisition functions (Stanton *et al.*, 2023).

138 **References**

- 139 Altamirano, M., Briol, F.-X., and Knoblauch, J. (2024). Robust and conjugate gaussian process  
140 regression. In *The 41st International Conference on Machine Learning*.
- 141 Chen, Y., Prati, A., Montgomery, J., and Garnett, R. (2023). A multi-task gaussian process model for  
142 inferring time-varying treatment effects in panel data. In *Proceedings of The 26th International  
143 Conference on Artificial Intelligence and Statistics*.
- 144 Cheng, L., Ramchandran, S., Vatanen, T., Lietzén, N., Lahesmaa, R., Vehtari, A., and Lähdesmäki,  
145 H. (2019). An additive gaussian process regression model for interpretable non-parametric anal-  
146 ysis of longitudinal data. *Nature Communications*.
- 147 Garnett, R. (2023). *Bayesian Optimization*. Cambridge University Press.
- 148 Kanagawa, M., Hennig, P., Sejdinovic, D., and Sriperumbudur, B. K. (2018). Gaussian processes and  
149 kernel methods: A review on connections and equivalences. *arXiv preprint arXiv:1807.02582*.
- 150 Maus, N., Kim, K., Pleiss, G., Eriksson, D., Cunningham, J. P., and Gardner, J. R. (2024).  
151 Approximation-aware bayesian optimization.
- 152 Rasmussen, C. and Williams, C. (2006). *Gaussian Processes for Machine Learning*. MIT Press.
- 153 Schölkopf, B., Herbrich, R., and Smola, A. J. (2001). A generalized representer theorem. In *Inter-  
154 national conference on computational learning theory*, pages 416–426. Springer.
- 155 Stanton, S., Maddox, W., and Wilson, A. G. (2023). Bayesian optimization with conformal predic-  
156 tion sets. In *International Conference on Artificial Intelligence and Statistics*.
- 157 Titsias, M. (2009). Variational learning of inducing variables in sparse gaussian processes. In  
158 *Artificial intelligence and statistics*.
- 159 Wenger, J., Pleiss, G., Pförtner, M., Hennig, P., and Cunningham, J. P. (2022). Posterior and compu-  
160 tational uncertainty in gaussian processes. In *Advances in Neural Information Processing Systems*.
- 161 Wild, V., Kanagawa, M., and Sejdinovic, D. (2023). Connections and equivalences between the  
162 nyström method and sparse variational gaussian processes.
- 163 Williams, C. and Seeger, M. (2000). Using the nyström method to speed up kernel machines.  
164 *Advances in neural information processing systems*.

165 **A Proof of Proposition 1**

166 **Posterior influence function.** Given the dataset  $\mathcal{D} = \{(\mathbf{x}_j, y_j)\}_{j=1}^n$ , we define the contamination  
 167 of  $\mathcal{D}$  indexed by  $m \in \{1, \dots, n\}$  as  $\mathcal{D}_m^c = (\mathcal{D} \setminus (\mathbf{x}_m, y_m)) \cup (\mathbf{x}_m, y_m^c)$ . PIF in general, aims  
 168 to measure the impact of  $y_m^c$  on inference through the divergence between the contaminated and  
 169 uncontaminated posteriors  $p(\mathbf{f}|\mathcal{D}_m^c)$  and  $p(\mathbf{f}|\mathcal{D})$ :

$$\text{PIF}(y_m^c, \mathcal{D}) = \text{KL}(p(\mathbf{f}|\mathcal{D}) \| p(\mathbf{f}|\mathcal{D}_m^c)) \quad (\text{S1})$$

170 where we call a posterior robust if  $\sup_{y \in \mathcal{Y}} |\text{PIF}(y_m^c, \mathcal{D})| < \infty$ .

171 We then establish the following lemma to prove Proposition 1.

172 **Lemma 5.** For an arbitrary matrix  $\hat{\mathbf{S}} \in \mathbb{R}^{m \times n}$  and positive semidefinite matrix  $\hat{\mathbf{B}} \in \mathbb{R}^{n \times n}$ , we  
 173 have that

$$\text{Tr}((\hat{\mathbf{S}}\hat{\mathbf{B}}\hat{\mathbf{S}}^\top)^{-1}) = \hat{\mathbf{S}}^{+\top} \hat{\mathbf{B}}^{-1/2} \hat{\mathbf{G}} \hat{\mathbf{B}}^{-1/2} \hat{\mathbf{S}}^+ \quad (\text{S2})$$

174 where we define  $\hat{\mathbf{G}} = \mathbf{I} - \hat{\mathbf{B}}^{-1/2}(\mathbf{I} - \hat{\mathbf{S}}^+\hat{\mathbf{S}})(\hat{\mathbf{B}}^{-1/2}(\mathbf{I} - \hat{\mathbf{S}}^+\hat{\mathbf{S}}))^+$  and  $^+$  denotes the Moore-Penrose  
 175 inverse.

176 *Proof:*

177

178 The whole proof is derived from an answer to a question posted on the [Mathematics Stack Exchange](#)  
 179 [Forums](#), which we write here for conciseness.

180 Denote  $\hat{\mathbf{O}} = \mathbf{I} - \hat{\mathbf{S}}^+\hat{\mathbf{S}}$  and  $\mathbf{H}(\alpha) = (\hat{\mathbf{S}}(\alpha\mathbf{I} + \hat{\mathbf{B}}^{-1})^{-1}\hat{\mathbf{S}}^\top)^{-1}$ . Note that

$$(\hat{\mathbf{S}}\hat{\mathbf{B}}\hat{\mathbf{S}}^\top)^{-1} = \lim_{\alpha \rightarrow 0} \mathbf{H}(\alpha) \quad (\text{S3})$$

181 By applying Woodbury matrix identity, we can rewrite  $\mathbf{H}(\alpha)$  as follows:

$$\mathbf{H}(\alpha) = \left( \frac{1}{\alpha} \hat{\mathbf{S}}\hat{\mathbf{S}}^\top - \frac{1}{\alpha} \hat{\mathbf{S}}\hat{\mathbf{B}}^{-1/2} \left( \mathbf{I} + \frac{1}{\alpha} \hat{\mathbf{B}}^{-1} \right)^{-1} \frac{1}{\alpha} \hat{\mathbf{B}}^{-1/2} \hat{\mathbf{S}}^\top \right)^{-1} \quad (\text{S4})$$

182 Since  $\hat{\mathbf{S}}\hat{\mathbf{S}}^\top$  is invertible, we can apply the Woodbury matrix identity for the second time to obtain

$$\begin{aligned} \mathbf{H}(\alpha) &= \alpha(\hat{\mathbf{S}}\hat{\mathbf{S}}^\top)^{-1} - (\hat{\mathbf{S}}\hat{\mathbf{S}}^\top)^{-1} \hat{\mathbf{S}}\hat{\mathbf{B}}^{-1/2} \\ &\quad \left( -(\mathbf{I} + \frac{1}{\alpha} \hat{\mathbf{B}}^{-1}) + \frac{1}{\alpha} \hat{\mathbf{B}}^{-1/2} \hat{\mathbf{S}}^\top (\hat{\mathbf{S}}\hat{\mathbf{S}}^\top)^{-1} \hat{\mathbf{S}}\hat{\mathbf{B}}^{-1/2} \right)^{-1} \hat{\mathbf{B}}^{-1/2} \hat{\mathbf{S}}^\top (\hat{\mathbf{S}}\hat{\mathbf{S}}^\top)^{-1} \end{aligned} \quad (\text{S5})$$

$$\begin{aligned} &= \alpha(\hat{\mathbf{S}}\hat{\mathbf{S}}^\top)^{-1} + (\hat{\mathbf{S}}\hat{\mathbf{S}}^\top)^{-1} \hat{\mathbf{S}}\hat{\mathbf{B}}^{-1/2} \left( \mathbf{I} + \frac{1}{\alpha} \hat{\mathbf{B}}^{-1/2} (\mathbf{I} - \hat{\mathbf{S}}^\top (\hat{\mathbf{S}}\hat{\mathbf{S}}^\top)^{-1} \hat{\mathbf{S}}) \hat{\mathbf{B}}^{-1/2} \right)^{-1} \\ &\quad \hat{\mathbf{B}}^{-1/2} \hat{\mathbf{S}}^\top (\hat{\mathbf{S}}\hat{\mathbf{S}}^\top)^{-1} \end{aligned} \quad (\text{S6})$$

183 We note that

$$\hat{\mathbf{S}}^\top (\hat{\mathbf{S}}\hat{\mathbf{S}}^\top)^{-1} = \hat{\mathbf{S}}^+ \quad (\text{S7})$$

$$\mathbf{I} - \hat{\mathbf{S}}^\top (\hat{\mathbf{S}}\hat{\mathbf{S}}^\top)^{-1} \hat{\mathbf{S}} = \hat{\mathbf{O}} \quad (\text{S8})$$

184 Then, we rewrite  $\mathbf{H}(\alpha)$  as follows:

$$\mathbf{H}(\alpha) = \alpha(\hat{\mathbf{S}}\hat{\mathbf{S}}^\top)^{-1} + \hat{\mathbf{S}}^{+\top}\hat{\mathbf{B}}^{-1/2} \left( \mathbf{I} + \frac{1}{\alpha}\hat{\mathbf{B}}^{-1/2}\hat{\mathbf{O}}\hat{\mathbf{O}}\hat{\mathbf{B}}^{-1/2} \right)^{-1} \hat{\mathbf{B}}^{-1/2}\hat{\mathbf{S}}^+ \quad (\text{S9})$$

185 Applying the Woodbury matrix identity for the third time provides

$$\mathbf{H}(\alpha) = \alpha(\hat{\mathbf{S}}\hat{\mathbf{S}}^\top)^{-1} + \hat{\mathbf{S}}^{+\top}\hat{\mathbf{B}}^{-1/2}(\mathbf{I} - \hat{\mathbf{B}}^{-1/2}\hat{\mathbf{O}}(\alpha\mathbf{I} + \hat{\mathbf{O}}\hat{\mathbf{B}}^{-1}\hat{\mathbf{O}})^{-1}\hat{\mathbf{O}}\hat{\mathbf{B}}^{-1/2})\hat{\mathbf{B}}^{-1/2}\hat{\mathbf{S}}^+ \quad (\text{S10})$$

186 Since the Moore-Penrose inverse of a matrix  $\mathbf{A}$  is a limit:

$$\mathbf{A}^+ = \lim_{\alpha \rightarrow 0} (\mathbf{A}^\top \mathbf{A} + \alpha \mathbf{I})^{-1} \mathbf{A}^\top = \lim_{\alpha \rightarrow 0} \mathbf{A}^\top (\mathbf{A} \mathbf{A}^\top + \alpha \mathbf{I})^{-1} \quad (\text{S11})$$

187 We can take the limit of  $\mathbf{H}(\alpha)$  as  $\alpha \rightarrow 0$  and apply the limit relation above to obtain the following  
188 result:

$$(\hat{\mathbf{S}}\hat{\mathbf{B}}\hat{\mathbf{S}}^\top)^{-1} = \hat{\mathbf{S}}^{+\top}\hat{\mathbf{B}}^{-1/2} \underbrace{(\mathbf{I} - \hat{\mathbf{B}}^{-1/2}\hat{\mathbf{O}}(\hat{\mathbf{B}}^{-1/2}\hat{\mathbf{O}})^+)}_{\hat{\mathbf{G}}}\hat{\mathbf{B}}^{-1/2}\hat{\mathbf{S}}^+ \quad (\text{S12})$$

189 **PIF for the IterGP.** IterGP regression has the PIF for some constant  $C'_1 \in \mathbb{R}$ .

$$\text{PIF}_{\text{IterGP}}(y_m^c, \mathcal{D}, i) = C'_1 (y_m - y_m^c)^2 \quad (\text{S13})$$

190 and is not robust:  $\text{PIF}_{\text{IterGP}}(y_m^c, \mathcal{D}, i) \rightarrow \infty$  as  $|y_m^c| \rightarrow \infty$ .

191 *Proof:*

192 Let  $p(\mathbf{f}|\mathcal{D}) = \mathcal{N}(\mathbf{f}; \boldsymbol{\mu}_i, \mathbf{K}_i)$  and  $p(\mathbf{f}|\mathcal{D}_m^c) = \mathcal{N}(\mathbf{f}; \boldsymbol{\mu}_i^c, \mathbf{K}_i^c)$  be the uncontaminated and contaminated  
193 computation-aware GP, respectively. Here,

$$\boldsymbol{\mu}_i = \mathbf{m} + \mathbf{K}\mathbf{v}_i \quad (\text{S14})$$

$$\mathbf{K}_i = \mathbf{K}\mathbf{C}_i\sigma_{\text{noise}}^2\mathbf{I}_n \quad (\text{S15})$$

$$\boldsymbol{\mu}_i^c = \mathbf{m} + \mathbf{K}\mathbf{v}_i^c \quad (\text{S16})$$

$$\mathbf{K}_i^c = \mathbf{K}\mathbf{C}_i\sigma_{\text{noise}}^2\mathbf{I}_n \quad (\text{S17})$$

194 Note that both  $\mathbf{K}_i$  and  $\mathbf{K}_i^c$  share the same matrix  $\mathbf{C}_i$ . Then, the PIF has the following form:

$$\text{PIF}_{\text{IterGP}}(y_m^c, \mathcal{D}, i) = \frac{1}{2}(\text{Tr}(\mathbf{K}_i^c\mathbf{K}_i) - n + (\boldsymbol{\mu}_i^c - \boldsymbol{\mu}_i)^\top(\mathbf{K}_i^c)^{-1}(\boldsymbol{\mu}_i^c - \boldsymbol{\mu}_i) + \ln\left(\frac{\det(\mathbf{K}_i^c)}{\det(\mathbf{K}_i)}\right)) \quad (\text{S18})$$

195 Based on [Altamirano et al. \(2024\)](#), the PIF leads to the following form:

$$\text{PIF}_{\text{IterGP}}(y_m^c, \mathcal{D}, i) = \frac{1}{2}((\boldsymbol{\mu}_i^c - \boldsymbol{\mu}_i)^\top(\mathbf{K}_i^c)^{-1}(\boldsymbol{\mu}_i^c - \boldsymbol{\mu}_i)) \quad (\text{S19})$$

196 Notice that the term  $\boldsymbol{\mu}_i^c - \boldsymbol{\mu}_i$  can be written as

$$\boldsymbol{\mu}_i^c - \boldsymbol{\mu}_i = (\mathbf{m} + \mathbf{K}\mathbf{v}_i^c) - (\mathbf{m} + \mathbf{K}\mathbf{v}_i) \quad (\text{S20})$$

$$= \mathbf{K}(\mathbf{v}_i^c - \mathbf{v}_i) \quad (\text{S21})$$

$$= \mathbf{K}(\mathbf{C}_i(\mathbf{y}^c - \mathbf{m}) - \mathbf{C}_i(\mathbf{y} - \mathbf{m})) \quad (\text{S22})$$

$$= \mathbf{K}(\mathbf{C}_i(\mathbf{y}^c - \mathbf{y})) \quad (\text{S23})$$

197 Substituting the RHS of Eq. (S23) to  $\boldsymbol{\mu}_i^c - \boldsymbol{\mu}_i$  in Eq. (S19), we obtain

$$\text{PIF}_{\text{IterGP}}(y_m^c, \mathcal{D}, i) = \frac{1}{2}(\mathbf{C}_i(\mathbf{y}^c - \mathbf{y}))^\top \mathbf{K} (\mathbf{K}\mathbf{C}_i\sigma^2\mathbf{I})^{-1} \mathbf{K}(\mathbf{C}_i(\mathbf{y}^c - \mathbf{y})) \quad (\text{S24})$$

$$= \frac{1}{2}(\mathbf{y}^c - \mathbf{y})^\top \mathbf{C}_i^\top \mathbf{K}\sigma_{\text{noise}}^{-2} \mathbf{I}(\mathbf{y}^c - \mathbf{y}) \quad (\text{S25})$$

198 Note that  $\mathbf{y}$  and  $\mathbf{y}^c$  have only one exception for the  $m$ -th element. Thus, we have

$$\text{PIF}_{\text{IterGP}}(y_m^c, \mathcal{D}, i) = \frac{1}{2}[\mathbf{C}_i^\top \mathbf{K}\sigma^{-2}\mathbf{I}]_{mm}(y_m^c - y_m)^2 \quad (\text{S26})$$

199 **PIF for the IterRCGP.** For the IterRCGP with  $\sup_{\mathbf{x}, y} w(\mathbf{x}, y) < \infty$ , the following holds

$$\text{PIF}_{\text{IterRCGP}}(y_m^c, \mathcal{D}, i) \leq C'_2(w(\mathbf{x}_m, y_m^c)^2 y_m^c)^2 + C'_3 \quad (\text{S27})$$

200 for some constants  $C'_2, C'_3 \in \mathbb{R}$ . Therefore, if  $\sup_{\mathbf{x}, y} y w(\mathbf{x}, y)^2 < \infty$ , the computation-aware  
201 RCGP is robust since  $|\text{PIF}_{\text{IterRCGP}}(y_m^c, \mathcal{D}, i)| < \infty$ .

202 *Proof:*

203 Without loss of generality, we aim to prove the bound for  $m = n$ . We can extend the proof for an  
204 arbitrary  $m \in \{1, \dots, n\}$ . Let  $p^w(\mathbf{f}|\mathcal{D}) = \mathcal{N}(\mathbf{f}; \hat{\boldsymbol{\mu}}_i, \hat{\mathbf{K}}_i)$  and  $p^w(\mathbf{f}|\mathcal{D}_m^c) = \mathcal{N}(\mathbf{f}; \hat{\boldsymbol{\mu}}_i^c, \hat{\mathbf{K}}_i^c)$  be the  
205 uncontaminated and contaminated computation-aware RCGP, respectively. Here,

$$\hat{\boldsymbol{\mu}}_i = \mathbf{m} + \mathbf{K}\tilde{\mathbf{C}}_i\tilde{\mathbf{v}}_i \quad (\text{S28})$$

$$\hat{\mathbf{K}}_i = \mathbf{K}\tilde{\mathbf{C}}_i\sigma_{\text{noise}}^2\mathbf{J}_{\mathbf{w}} \quad (\text{S29})$$

$$\hat{\boldsymbol{\mu}}_i^c = \mathbf{m} + \mathbf{K}\tilde{\mathbf{C}}_i^c\tilde{\mathbf{v}}_i^c \quad (\text{S30})$$

$$\hat{\mathbf{K}}_i^c = \mathbf{K}\tilde{\mathbf{C}}_i^c\sigma_{\text{noise}}^2\mathbf{J}_{\mathbf{w}^c} \quad (\text{S31})$$

206 where  $\mathbf{w}^c = (w(\mathbf{x}_1, y_1), \dots, w(\mathbf{x}_n, y_n^c))^\top$ . The PIF has the following form

$$\text{PIF}_{\text{IterRCGP}}(y_m^c, \mathcal{D}, i) = \frac{1}{2} \left( \underbrace{\text{Tr}((\hat{\mathbf{K}}_i^c)^{-1}\hat{\mathbf{K}}_i) - n}_{(1)} + \underbrace{(\hat{\boldsymbol{\mu}}_i^c - \hat{\boldsymbol{\mu}}_i)^\top (\hat{\mathbf{K}}_i^c)^{-1} (\hat{\boldsymbol{\mu}}_i^c - \hat{\boldsymbol{\mu}}_i)}_{(2)} + \underbrace{\ln \left( \frac{\det(\hat{\mathbf{K}}_i^c)}{\det(\hat{\mathbf{K}}_i)} \right)}_{(3)} \right) \quad (\text{S32})$$

207 We first derive the bound for (1):

$$(1) = \text{Tr}((\hat{\mathbf{K}}_i^c)^{-1}\hat{\mathbf{K}}_i) - n \quad (\text{S33})$$

$$= \text{Tr} \left( (\mathbf{K}\tilde{\mathbf{C}}_i^c\sigma_{\text{noise}}^2\mathbf{J}_{\mathbf{w}^c})^{-1} \mathbf{K}\tilde{\mathbf{C}}_i\sigma_{\text{noise}}^2\mathbf{J}_{\mathbf{w}} \right) - n \quad (\text{S34})$$

$$= \text{Tr}(\sigma_{\text{noise}}^{-2} \mathbf{J}_{\mathbf{w}^c}^{-1} (\tilde{\mathbf{C}}_i^c)^{-1} \tilde{\mathbf{C}}_i\sigma_{\text{noise}}^2\mathbf{J}_{\mathbf{w}}) - n \quad (\text{S35})$$

$$\leq \text{Tr}(\sigma_{\text{noise}}^{-2} \mathbf{J}_{\mathbf{w}^c}^{-1} (\tilde{\mathbf{C}}_i^c)^{-1}) \text{Tr}(\tilde{\mathbf{C}}_i\sigma_{\text{noise}}^2\mathbf{J}_{\mathbf{w}}) - n \quad (\text{S36})$$

$$\leq \text{Tr}(\sigma_{\text{noise}}^{-2} \mathbf{J}_{\mathbf{w}^c}^{-1}) \text{Tr}(\tilde{\mathbf{C}}_i^c)^{-1} \text{Tr}(\tilde{\mathbf{C}}_i\sigma_{\text{noise}}^2\mathbf{J}_{\mathbf{w}}) - n \quad (\text{S37})$$

208 The first and second inequality come from the fact that  $\text{Tr}(\mathbf{A}\mathbf{F}) \leq \text{Tr}(\mathbf{A})\text{Tr}(\mathbf{F})$  for two positive  
209 semidefinite matrices  $\mathbf{A}$  and  $\mathbf{F}$ . Since  $\text{Tr}(\tilde{\mathbf{C}}_i\sigma_{\text{noise}}^2\mathbf{J}_{\mathbf{w}})$  does not contain the contamination term, we



210 can write  $\bar{C}_1 = \text{Tr}(\tilde{\mathbf{C}}_i \sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}})$ . Let  $\mathbf{B} = (\mathbf{S}_i^\top \tilde{\mathbf{K}}^c \mathbf{S}_i)^{-1}$  such that  $\mathbf{C}_i^c = \mathbf{S}_i^\top \mathbf{B} \mathbf{S}_i^\top$ . Observe that  
 211 matrix  $\mathbf{B}$  is positive semidefinite. Thus, we can apply Lemma 5 to obtain the bound of  $\text{Tr}((\tilde{\mathbf{C}}_i^c)^{-1})$ :

$$\text{Tr}((\tilde{\mathbf{C}}_i^c)^{-1}) = \text{Tr}((\mathbf{S}_i^\top \mathbf{B} \mathbf{S}_i^\top)^{-1}) \quad (\text{S38})$$

$$= \text{Tr}(\mathbf{S}_i^{+\top} \mathbf{B}^{-1/2} \mathbf{G} \mathbf{B}^{-1/2} \mathbf{S}_i^+) \quad (\text{S39})$$

$$\leq \text{Tr}(\mathbf{S}_i^+ \mathbf{S}_i^{+\top}) \text{Tr}(\mathbf{B}^{-1/2} \mathbf{B}^{-1/2}) \text{Tr}(\mathbf{G}) \quad (\text{S40})$$

212 where

$$\text{Tr}(\mathbf{G}) = \text{Tr}(\mathbf{I} - \mathbf{B}^{-1/2} (\mathbf{I} - \mathbf{S}_i^+ \mathbf{S}_i) (\mathbf{B}^{-1/2} (\mathbf{I} - \mathbf{S}_i^+ \mathbf{S}_i))^+) \quad (\text{S41})$$

$$= n - \text{Tr}(\mathbf{B}^{-1/2} (\mathbf{I} - \mathbf{S}_i^+ \mathbf{S}_i) (\mathbf{I} - \mathbf{S}_i^+ \mathbf{S}_i)^+ \mathbf{B}^{-1/2}) \quad (\text{S42})$$

$$\leq n - \text{Tr}(\mathbf{B}^{-1/2} \mathbf{B}^{-1/2}) \text{Tr}((\mathbf{I} - \mathbf{S}_i^+ \mathbf{S}_i) (\mathbf{I} - \mathbf{S}_i^+ \mathbf{S}_i)^+) \quad (\text{S43})$$

213 The inequality S40 stems from the trace circular property and the inequality of the product of two  
 214 semidefinite matrices. Note that  $\text{Tr}(\mathbf{G}) \leq n$  since  $\mathbf{B}^{-1/2} \mathbf{B}^{-1/2}$  and  $(\mathbf{I} - \mathbf{S}_i^+ \mathbf{S}_i) (\mathbf{I} - \mathbf{S}_i^+ \mathbf{S}_i)^+$  in  
 215 S43 are positive semidefinite matrices; thus both have non-negative trace value. Therefore, we find  
 216 that

$$\text{Tr}((\tilde{\mathbf{C}}_i^c)^{-1}) \leq n \text{Tr}(\mathbf{S}_i^+ \mathbf{S}_i^{+\top}) \text{Tr}(\mathbf{B}^{-1}) \quad (\text{S44})$$

$$\leq n \text{Tr}(\mathbf{S}_i^+ \mathbf{S}_i^{+\top}) \text{Tr}(\mathbf{S}_i \mathbf{S}_i^\top) \text{Tr}(\tilde{\mathbf{K}}^c) \quad (\text{S45})$$

$$= \bar{C}_2 \text{Tr}(\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c}) \quad (\text{S46})$$

217 where we define  $\bar{C}_2 = n \text{Tr}(\mathbf{S}_i^+ \mathbf{S}_i^{+\top}) \text{Tr}(\mathbf{S}_i \mathbf{S}_i^\top)$ . We then plug S46 into S37 to obtain

$$(1) \leq \text{Tr}(\sigma_{\text{noise}}^{-2} \mathbf{J}_{\mathbf{w}^c}^{-1}) \text{Tr}(\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c}) \bar{C}_1 \bar{C}_2 - n \quad (\text{S47})$$

$$= \left( \sum_{j=1}^n (\sigma_{\text{noise}}^{-2} w^2(\mathbf{x}_j, y_j)) \sum_{k=1}^n (\mathbf{K}_{kk} + \sigma_{\text{noise}}^2 w^{-2}(\mathbf{x}_k, y_k)) \right) \bar{C}_1 \bar{C}_2 - n \quad (\text{S48})$$

$$\leq \left( n^2 \sup_{\mathbf{x}, y} w^2(\mathbf{x}, y) \sup_{\hat{\mathbf{x}}, \hat{y}} w^{-2}(\hat{\mathbf{x}}, \hat{y}) \right) \bar{C}_1 \bar{C}_2 - n = \bar{C}_3 \quad (\text{S49})$$

218 Next, we derive the bound for (2). Following Altamirano *et al.* (2024), we have that

$$(2) \leq \lambda_{\max}((\hat{\mathbf{K}}_i^c)^{-1}) \|\hat{\boldsymbol{\mu}}_i^c - \hat{\boldsymbol{\mu}}_i\|_1^2 \quad (\text{S50})$$

219 We expand  $\lambda_{\max}((\hat{\mathbf{K}}_i^c)^{-1})$  and derive the following bound:

$$\lambda_{\max}((\hat{\mathbf{K}}_i^c)^{-1}) = \lambda_{\max}(\sigma^{-2} \mathbf{J}_{\mathbf{w}^c}^{-1} (\tilde{\mathbf{C}}_i^c)^{-1} \mathbf{K}^{-1}) \quad (\text{S51})$$

$$\leq \lambda_{\max}(\sigma_{\text{noise}}^{-2} \mathbf{J}_{\mathbf{w}^c}^{-1}) \lambda_{\max}((\tilde{\mathbf{C}}_i^c)^{-1}) \lambda_{\max}(\mathbf{K}^{-1}) \quad (\text{S52})$$

$$= \lambda_{\max}(\sigma_{\text{noise}}^{-2} \mathbf{J}_{\mathbf{w}^c}^{-1}) \lambda_{\min}(\tilde{\mathbf{C}}_i^c) \lambda_{\max}(\mathbf{K}^{-1}) \quad (\text{S53})$$

$$\leq \lambda_{\max}(\sigma_{\text{noise}}^{-2} \mathbf{J}_{\mathbf{w}^c}^{-1}) \left( \lambda_{\min}((\tilde{\mathbf{K}}^c)^{-1}) \right) \lambda_{\max}(\mathbf{K}^{-1}) \quad (\text{S54})$$

$$\leq \lambda_{\max}(\sigma_{\text{noise}}^{-2} \mathbf{J}_{\mathbf{w}^c}^{-1}) \lambda_{\min}((\tilde{\mathbf{K}}^c)^{-1}) \lambda_{\max}(\mathbf{K}^{-1}) \quad (\text{S55})$$

$$\leq \lambda_{\max}(\sigma_{\text{noise}}^{-2} \mathbf{J}_{\mathbf{w}^c}^{-1}) (\lambda_{\max}(\mathbf{K}) + \lambda_{\max}(\sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c})) \lambda_{\max}(\mathbf{K}^{-1}) \quad (\text{S56})$$

220 The first inequality follows from the maximum eigenvalue of the product of two positive semidefinite  
 221 matrices. The fact that the maximum eigenvalue of a matrix is equal to the minimum eigenvalue  
 222 of the inverse leads to the second equality. Recall that  $\tilde{\mathbf{C}}_i^c = (\tilde{\mathbf{K}}_i^c)^{-1} - \tilde{\Sigma}_i$ . Since  $\tilde{\mathbf{C}}_i^c$ ,  $(\tilde{\mathbf{K}}_i^c)^{-1}$  and  
 223  $\tilde{\Sigma}_i$  are positive semidefinite matrices, the third inequality holds. The fourth inequality stems from  
 224 the equivalence of the maximum eigenvalue and the addition property of the maximum eigenvalue  
 225 of two positive semidefinite matrices.

226 Since  $\mathbf{J}_{\mathbf{w}^c}^{-1} = \text{diag}((\mathbf{w}^c)^2)$ , and  $\sup_{\mathbf{x}, y} w(\mathbf{x}, y) < \infty$ , it holds that  $\lambda_{\max}(\sigma_{\text{noise}}^{-2} \mathbf{J}_{\mathbf{w}^c}^{-1}) = \bar{C}_4 < +\infty$   
 227 and  $\lambda_{\max}(\sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c}) = \bar{C}_5 < +\infty$ , such that

$$\lambda_{\max}((\hat{\mathbf{K}}_i^c)^{-1}) \leq \bar{C}_4(\lambda_{\max}(\mathbf{K}) + \bar{C}_5)\lambda_{\max}(\mathbf{K}^{-1}) = \bar{C}_6 \quad (\text{S57})$$

228 We substitute  $\bar{C}_6$  into (2) to obtain

$$(2) \leq \bar{C}_6 \|\hat{\boldsymbol{\mu}}_i^c - \hat{\boldsymbol{\mu}}_i\|_1^2 \quad (\text{S58})$$

$$= \bar{C}_6 \|(\mathbf{m} + \mathbf{K}\tilde{\mathbf{v}}_i^c) - (\mathbf{m} + \mathbf{K}\tilde{\mathbf{v}}_i)\|_1^2 \quad (\text{S59})$$

$$= \bar{C}_6 \|\mathbf{K}(\tilde{\mathbf{C}}_i^c(\mathbf{y} - \mathbf{m}_{\mathbf{w}^c}) - \tilde{\mathbf{C}}_i(\mathbf{y} - \mathbf{m}_{\mathbf{w}}))\|_1^2 \quad (\text{S60})$$

$$\leq \bar{C}_6 \|\mathbf{K}\|_F \|\tilde{\mathbf{C}}_i^c(\mathbf{y} - \mathbf{m}_{\mathbf{w}^c}) - \tilde{\mathbf{C}}_i(\mathbf{y} - \mathbf{m}_{\mathbf{w}})\|_1^2 \quad (\text{S61})$$

$$\leq \bar{C}_6 \|\mathbf{K}\|_F (\|(\tilde{\mathbf{K}}_i^c)^{-1}(\mathbf{y} - \mathbf{m}_{\mathbf{w}^c}) - (\tilde{\mathbf{K}})^{-1}(\mathbf{y} - \mathbf{m}_{\mathbf{w}})\|_1^2 + \|\tilde{\Sigma}_i^c(\mathbf{y} - \mathbf{m}_{\mathbf{w}^c}) - \tilde{\Sigma}_i(\mathbf{y} - \mathbf{m}_{\mathbf{w}})\|_1^2) \quad (\text{S62})$$

$$\leq q\bar{C}_6 \|\mathbf{K}\|_F (\|(\tilde{\mathbf{K}}_i^c)^{-1}(\mathbf{y} - \mathbf{m}_{\mathbf{w}^c}) - (\tilde{\mathbf{K}})^{-1}(\mathbf{y} - \mathbf{m}_{\mathbf{w}})\|_1^2) \quad (\text{S63})$$

$$= q\bar{C}_6 \|\mathbf{K}\|_F ((\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c})^{-1}(\mathbf{y} - \mathbf{m}_{\mathbf{w}^c}) - (\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}})(\mathbf{y} - \mathbf{m}_{\mathbf{w}}))\|_1^2 \quad (\text{S64})$$

229 for a constant  $q > 0$ . The second equality follows from [Wenger et al. \(2022\)](#)[Eq. (S45)]. The first in-  
 230 equality follows the Cauchy-Schwarz inequality. The second inequality stems from the definition of  
 231  $\tilde{\mathbf{C}}_i$ ,  $\tilde{\mathbf{C}}_i^c$ , and the triangle inequality. Finally, the last inequality holds since  $(\tilde{\mathbf{K}}_i^{-1} - \tilde{\Sigma}_i)$ ,  $\tilde{\mathbf{K}}_i^{-1}$ ,  $\tilde{\Sigma}_i \succeq$   
 232 0.

233 Applying results from [Altamirano et al. \(2024\)](#), we obtain

$$(2) \leq q\bar{C}_6 \|\mathbf{K}\|_F ((\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c})^{-1}(\mathbf{y} - \mathbf{m}_{\mathbf{w}^c}) - (\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}})(\mathbf{y} - \mathbf{m}_{\mathbf{w}}))\|_1^2 \quad (\text{S65})$$

$$\leq q\bar{C}_6 \|\mathbf{K}\|_F 2((\bar{C}_7 + \bar{C}_8)^2 + (\bar{C}_9 + \bar{C}_{10})^2 (w(x_n, y_n^c)^2 y_n^c)^2) \quad (\text{S66})$$

$$\leq \bar{C}_{11} + \bar{C}_{12} (w(x_n, y_n^c)^2 y_n^c)^2 \quad (\text{S67})$$

234 where  $\bar{C}_{11} = q\bar{C}_6 \|\mathbf{K}\|_F 2(\bar{C}_7 + \bar{C}_8)^2$  and  $\bar{C}_{12} = q\bar{C}_6 \|\mathbf{K}\|_F 2(\bar{C}_9 + \bar{C}_{10})^2$ . The terms  
 235  $\bar{C}_7, \bar{C}_8, \bar{C}_9, \bar{C}_{10}$  equal to  $\bar{C}_6, \bar{C}_8, \bar{C}_7, \bar{C}_9$  in [Altamirano et al. \(2024\)](#).

236 The term (3) can be written as follows:

$$(3) = \ln \left( \frac{\det(\hat{\mathbf{K}}_i^c)}{\det(\hat{\mathbf{K}}_i)} \right) \quad (\text{S68})$$

$$= \ln \left( \frac{\det(\tilde{\mathbf{C}}_i^c \sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c})}{\det(\tilde{\mathbf{C}}_i \sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}})} \right) \quad (\text{S69})$$

$$= \ln(\det(\sigma_{\text{noise}}^{-2} \mathbf{J}_{\mathbf{w}}^{-1} \tilde{\mathbf{C}}_i^{-1}) \det(\tilde{\mathbf{C}}_i^c) \det(\sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c})) \quad (\text{S70})$$

237 Observe that we can write  $\bar{C}_{13} = \ln(\det(\sigma_{\text{noise}}^{-2} \mathbf{J}_{\mathbf{w}}^{-1} \tilde{\mathbf{C}}_i^{-1}))$  since it does not contain the continuation  
 238 term. Furthermore, we obtain

$$(3) = \ln(\bar{C}_{13} \det(\bar{\mathbf{C}}_i^c) \det(\sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c})) \quad (\text{S71})$$

$$\leq \ln(\bar{C}_{13} \det((\tilde{\mathbf{K}}^c)^{-1}) \det(\sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c})) \quad (\text{S72})$$

$$= \ln\left(\bar{C}_{13} \frac{\det(\sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c})}{\det(\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c})}\right) \quad (\text{S73})$$

$$\leq \ln\left(\bar{C}_{13} \frac{\det(\sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c})}{\det(\mathbf{K}) + \det(\sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c})}\right) \quad (\text{S74})$$

239 The first inequality holds since  $((\tilde{\mathbf{K}}_i^c)^{-1} - \tilde{\mathbf{\Sigma}}_i^c)$ ,  $(\tilde{\mathbf{K}}_i^c)^{-1}$ ,  $\tilde{\mathbf{\Sigma}}_i^c \succeq 0$ , so  $\det((\tilde{\mathbf{K}}_i^c)^{-1}) \geq \det(\tilde{\mathbf{\Sigma}}_i^c)$ . The  
 240 last inequality leverages the fact that  $\det(\mathbf{A} + \mathbf{F}) \geq \det(\mathbf{A}) + \det(\mathbf{F})$  for  $\mathbf{A}$  and  $\mathbf{F}$  are positive  
 241 semidefinite matrices. Since  $\det(\mathbf{K})$ ,  $\det(\sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c}) \geq 0$ , we find that

$$\ln\left(\frac{\det(\sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c})}{\det(\mathbf{K}) + \det(\sigma_{\text{noise}}^2 \mathbf{J}_{\mathbf{w}^c})}\right) \leq 1 \quad (\text{S75})$$

242 Leading to the following inequality:

$$(3) \leq \ln(\bar{C}_{13}) = \bar{C}_{14} \quad (\text{S76})$$

243 Finally, putting the three terms together, we obtain the following bound:

$$\text{PIF}_{\text{IterRCGP}}(y_m^c, \mathcal{D}, i) \leq \bar{C}_3 + \bar{C}_{11} + \bar{C}_{12}(w(x_n, y_n^c)^2 y_n^c)^2 + \bar{C}_{14} \quad (\text{S77})$$

$$= C'_2(w(x_n, y_n^c)^2 y_n^c)^2 + C'_3 \quad (\text{S78})$$

244 where  $C'_2 = \bar{C}_{12}$  and  $C'_3 = \bar{C}_3 + \bar{C}_{11} + \bar{C}_{14}$ .

## 245 B Proof of Proposition 2

246 **Unique solution of the empirical-risk minimization problem.** We first show the existence of a  
 247 unique solution to the empirical risk minimization problem corresponding to RCGP. For this pur-  
 248 pose, we set  $\mathbf{m} = \mathbf{0}$ . Following [Altamirano et al. \(2024\)](#) (proof of [Proposition 3.1]), we can rewrite  
 249  $L_n^w$  and formulate the RCGP objective as the following empirical-risk minimization problem:

$$\hat{\mathbf{f}} = \operatorname{argmin}_{\mathbf{f} \in \mathcal{H}_k} \frac{1}{2n} \left( \underbrace{\mathbf{f}^\top \lambda^{-1} \mathbf{J}_{\mathbf{w}}^{-1} \mathbf{f} - 2\mathbf{f}^\top \lambda^{-1} \mathbf{J}_{\mathbf{w}}^{-1} (\mathbf{y} - \mathbf{m}_{\mathbf{w}}) + Q(\mathbf{x}, \mathbf{y}, \lambda)}_{L_n^w} + \|\mathbf{f}\|_{\mathcal{H}_k}^2 \right) \quad (\text{S79})$$

250 where

$$Q(\mathbf{x}, \mathbf{y}, \lambda) = \mathbf{y}^\top \lambda^{-1} \operatorname{diag}(2\lambda^{-1} \mathbf{w}^2) \mathbf{y} - 4\lambda \nabla_{\mathbf{y}} \mathbf{y}^\top \mathbf{w}^2 \quad (\text{S80})$$

251 for  $\lambda > 0$ . We then show the unique solution to [S79](#) through the following lemma:

252

253 **Lemma 6.** *If  $\lambda > 0$  and the kernel  $k$  is invertible, the solution to [S79](#) is a unique, and is given by*

$$\hat{\mathbf{f}}(\mathbf{x}) = \mathbf{k}_{\mathbf{x}} (\mathbf{K} + \lambda \mathbf{J}_{\mathbf{w}})^{-1} (\mathbf{y} - \mathbf{m}_{\mathbf{w}}) = \sum_{j=1}^n \alpha_j k(\mathbf{x}, \mathbf{x}_j), \mathbf{x} \in \mathcal{X} \quad (\text{S81})$$

254 where

$$(\alpha_1, \dots, \alpha_n) = (\mathbf{K} + \lambda \mathbf{J}_w)^{-1}(\mathbf{y} - \mathbf{m}_w) \in \mathbb{R}^n \quad (\text{S82})$$

255 *Proof:*

256 The optimization problem in S79 allows us to apply the representer theorem (Schölkopf *et al.*, 2001).  
257 It implies that the solution of S79 can be written as a weighted sum, i.e.,

$$\hat{\mathbf{f}} = \sum_{j=1}^n \alpha_j k(\cdot, \mathbf{x}_j) \quad (\text{S83})$$

258 for  $\alpha_1, \dots, \alpha_n \in \mathbb{R}$ . Let  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_n]^\top \in \mathbb{R}^n$ . Substituting S83 into S79 provides

$$\operatorname{argmin}_{\boldsymbol{\alpha} \in \mathbb{R}^n} \frac{1}{2n} (\lambda^{-1} \boldsymbol{\alpha}^\top \mathbf{K} \mathbf{J}_w^{-1} \mathbf{K} \boldsymbol{\alpha} - 2\lambda^{-1} \boldsymbol{\alpha}^\top \mathbf{K} \mathbf{J}_w^{-1} (\mathbf{y} - \mathbf{m}_w) + Q(\mathbf{x}, \mathbf{y}, \lambda) + \|\hat{\mathbf{f}}\|_{\mathcal{H}_k}^2) \quad (\text{S84})$$

259 where  $\|\hat{\mathbf{f}}\|_{\mathcal{H}_k}^2 = \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha}$ , following the reproducing property. Taking the differentiation of the  
260 objective w.r.t.  $\boldsymbol{\alpha}$ , setting it equal to zero, and arranging the result yields the following equation:

$$\mathbf{K}(\mathbf{K} + \lambda \mathbf{J}_w) \boldsymbol{\alpha} = \mathbf{K}(\mathbf{y} - \mathbf{m}_w) \quad (\text{S85})$$

261 Since the objective in S84 is a convex function of  $\boldsymbol{\alpha}$ , we find that  $\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{J}_w)^{-1}(\mathbf{y} - \mathbf{m}_w)$   
262 provides the minimum of the objective (S79 and S84). Furthermore, we can verify that  $L_n^w$  is a  
263 convex function w.r.t.  $\mathbf{f}$ . Therefore, we conclude that  $\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{J}_w)^{-1}(\mathbf{y} - \mathbf{m}_w)$  provides  
264 the unique solution to S79. As a remark, Proposition 6 closely connects with [Theorem 3.4] in  
265 Kanagawa *et al.* (2018).  
266

267 **Relative bound errors.** We also provide the equivalence of Proposition 2 in Wenger *et al.* (2022):  
268

269 **Proposition 7.** For any choice of actions a relative bound error  $\hat{\rho}(i)$  s.t.  $\|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_i\|_{\tilde{\mathbf{K}}} \leq \hat{\rho}(i) \|\hat{\mathbf{v}}\|_{\tilde{\mathbf{K}}}$  is  
270 given by

$$\hat{\rho}(i) = (\bar{\mathbf{v}}^\top (\mathbf{I} - \tilde{\mathbf{C}}_i \tilde{\mathbf{K}}) \bar{\mathbf{v}})^{1/2} \leq \lambda_{\max}(\mathbf{I} - \tilde{\mathbf{C}}_i \tilde{\mathbf{K}}) \leq 1 \quad (\text{S86})$$

271 where  $\bar{\mathbf{v}} = \hat{\mathbf{v}} / \|\hat{\mathbf{v}}\|_{\tilde{\mathbf{K}}}$ .

272 The proof is direct since we only need to substitute  $\mathbf{C}_i, \hat{\mathbf{K}}, \mathbf{v}_*$  in Wenger *et al.* (2022) with  $\tilde{\mathbf{C}}_i, \tilde{\mathbf{K}}, \hat{\mathbf{v}}$ ,  
273 respectively.

274 **Proof of Proposition 2.** Lemma 6 implies there exists a unique solution to the corresponding RCGP  
275 risk minimization problem. Choosing  $\hat{\rho}(i)$  as described in Proposition 7, we have that  $\|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_i\|_{\tilde{\mathbf{K}}}^2 \leq$   
276  $\hat{\rho}(i) \|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_0\|_{\tilde{\mathbf{K}}}$ , where  $\tilde{\mathbf{v}}_0 = \mathbf{0}$ . Then, for  $i \in \{0, \dots, n\}$  we find that

$$\|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_i\|_{\tilde{\mathbf{K}}}^2 \leq \|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_i\|_{\tilde{\mathbf{K}}}^2 \leq \hat{\rho}^2(i) \|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_0\|_{\tilde{\mathbf{K}}}^2 \quad (\text{S87})$$

$$\leq \hat{\rho}(i)^2 \left( \|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_0\|_{\tilde{\mathbf{K}}}^2 + \frac{\lambda_{\max}(\mathbf{J}_w)}{\lambda_{\min}(\mathbf{K})} \lambda_{\min}(\mathbf{K}) \|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_0\|_2^2 \right) \quad (\text{S88})$$

$$\leq \hat{\rho}(i)^2 \left( \|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_0\|_{\tilde{\mathbf{K}}}^2 + \frac{\lambda_{\max}(\mathbf{J}_w)}{\lambda_{\min}(\mathbf{K})} \|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_0\|_{\tilde{\mathbf{K}}}^2 \right) \quad (\text{S89})$$

$$\leq \hat{\rho}(i)^2 \left( 1 + \frac{\lambda_{\max}(\mathbf{J}_w)}{\lambda_{\min}(\mathbf{K})} \right) \|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_0\|_{\tilde{\mathbf{K}}}^2 \quad (\text{S90})$$

277 The third inequality stems from the definition of  $\mathbf{J}_w$  and the fact that the maximum eigenvalue of  
 278 a diagonal matrix is the largest component of its diagonal. Applying result from [Wenger et al.](#)  
 279 (2022), we have that

$$\|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_i\|_{\mathbf{K}}^2 = \|\hat{\boldsymbol{\mu}}_* - \hat{\boldsymbol{\mu}}_i\|_{\mathcal{H}_k}^2 \quad (\text{S91})$$

280 Combining both results and defining  $c(\mathbf{J}_w) = \left(1 + \frac{\lambda_{\max}(\mathbf{J}_w)}{\lambda_{\min}(\mathbf{K})}\right)$ , we obtain

$$\|\hat{\boldsymbol{\mu}}_* - \hat{\boldsymbol{\mu}}_i\|_{\mathcal{H}_k} = \|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_i\|_{\mathbf{K}} \leq \hat{\rho}(i)c(\mathbf{J}_w)\|\hat{\mathbf{v}} - \tilde{\mathbf{v}}_0\|_{\mathbf{K}} = \hat{\rho}(i)c(\mathbf{J}_w)\|\hat{\boldsymbol{\mu}}_* - \mathbf{m}\|_{\mathcal{H}_k} \quad (\text{S92})$$

### 281 C Proof of Proposition 3

282 Here, we refer to  $\sigma_{\text{noise}}^2$  as  $\sigma^2$  to simplify the notation. Let  $c_j = (\tilde{\mathbf{C}}_i k^{\sigma w}(\mathbf{X}, \mathbf{x}))_j$  for  $j = 1, \dots, n$ ,  
 283 where we define  $k^{\sigma w}(\cdot, \cdot) = k(\cdot, \cdot) + \frac{\sigma^2}{2} \delta_w(\cdot, \cdot)$ , where

$$\delta_w(\mathbf{x}, \mathbf{x}') = \begin{cases} w^{-2}(\mathbf{x}, y) & \mathbf{x} = \mathbf{x}' \text{ and } \mathbf{x} \in \mathcal{D} \\ 2 & \mathbf{x} = \mathbf{x}' \text{ and } \mathbf{x} \notin \mathcal{D} \\ 0 & \mathbf{x} \neq \mathbf{x}' \end{cases} \quad (\text{S93})$$

284 Since  $g, m \in \mathcal{H}_{k^{\sigma w}}$ , it implies that  $g - m \in \mathcal{H}_{k^{\sigma w}}$ . Then, applying [Lemma 3.9] in [Kanagawa et al.](#)  
 285 (2018) provides

$$\left( \sup_{\|g - m_w\|_{\mathcal{H}_{k^{\sigma w}}} \leq 1} g(\mathbf{x}) - \hat{\mu}_i^g(\mathbf{x}) \right)^2 = \left( \sup_{\|g - m_w\|_{\mathcal{H}_{k^{\sigma w}}} \leq 1} g(\mathbf{x}) - \sum_{j=1}^n c_j (g(\mathbf{x}_j) - m_w(\mathbf{x}_j)) \right)^2 \quad (\text{S94})$$

$$= \|k^{\sigma w}(\cdot, \mathbf{x}) - k(\mathbf{x}, \mathbf{X}) \tilde{\mathbf{C}}_i k^{\sigma w}(\mathbf{X}, \cdot)\|_{\mathcal{H}_{k^{\sigma w}}}^2 \quad (\text{S95})$$

$$= \langle k^{\sigma w}(\cdot, \mathbf{x}), k^{\sigma w}(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_k} - 2 \langle k^{\sigma w}(\cdot, \mathbf{x}), k(\mathbf{x}, \mathbf{X}) \tilde{\mathbf{C}}_i k^{\sigma w}(\mathbf{X}, \cdot) \rangle_{\mathcal{H}_k} + \langle k(\mathbf{x}, \mathbf{X}) \tilde{\mathbf{C}}_i k^{\sigma w}(\mathbf{X}, \cdot), k(\mathbf{x}, \mathbf{X}) \tilde{\mathbf{C}}_i k^{\sigma w}(\mathbf{X}, \cdot) \rangle_{\mathcal{H}_k} \quad (\text{S96})$$

286 By reproducing property, we have

$$= k^{\sigma w}(\mathbf{x}, \mathbf{x}) - 2k^{\sigma w}(\mathbf{x}, \mathbf{X}) \tilde{\mathbf{C}}_i k^{\sigma w}(\mathbf{X}, \mathbf{x}) + k(\mathbf{x}, \mathbf{X}) \tilde{\mathbf{C}}_i k^{\sigma w}(\mathbf{X}, \mathbf{X}) \tilde{\mathbf{C}}_i k^{\sigma w}(\mathbf{X}, \mathbf{x}) \quad (\text{S97})$$

287 if  $\mathbf{x} \neq \mathbf{x}_j$  or  $\sigma^2 = 0$ , it holds that  $k^{\sigma w}(\mathbf{x}, \mathbf{X}) = k(\mathbf{x}, \mathbf{X})$ . By definition, we have  $k^{\sigma w}(\mathbf{X}, \mathbf{X}) = \tilde{\mathbf{K}}$   
 288 and by [Wenger et al. \(2022\)](#)[Eq. (S42)], it holds that  $\tilde{\mathbf{C}}_i \tilde{\mathbf{K}} \tilde{\mathbf{C}}_i = \tilde{\mathbf{C}}_i$ . Therefore, we obtain

$$= k(\mathbf{x}, \mathbf{x}) + \sigma^2 - 2k(\mathbf{x}, \mathbf{X}) \tilde{\mathbf{C}}_i k(\mathbf{X}, \mathbf{x}) + k(\mathbf{x}, \mathbf{X}) \tilde{\mathbf{C}}_i \tilde{\mathbf{K}} \tilde{\mathbf{C}}_i k(\mathbf{X}, \mathbf{x}) \quad (\text{S98})$$

$$= k(\mathbf{x}, \mathbf{x}) + \sigma^2 - k(\mathbf{x}, \mathbf{X}) \tilde{\mathbf{C}}_i k(\mathbf{X}, \mathbf{x}) \quad (\text{S99})$$

$$= \hat{k}_i(\mathbf{x}, \mathbf{x}) + \sigma^2 \quad (\text{S100})$$

289 For the last result, we analogously choose  $c_j = ((\tilde{\mathbf{K}}^{-1} - \tilde{\mathbf{C}}_i) k^{\sigma w}(\mathbf{X}, \mathbf{x}))_j$ . Then, we obtain

$$\left( \sup_{\|g-m_w\|_{\mathcal{H}_{k^{\sigma w}}} \leq 1} \hat{\mu}^g(\mathbf{x}) - \hat{\mu}_i^g(\mathbf{x}) \right)^2 = \left( \sup_{\|g-m_w\|_{\mathcal{H}_{k^{\sigma w}}} \leq 1} \sum_{j=0}^n c_j g(\mathbf{x}_j) \right)^2 \quad (\text{S101})$$

$$= \|k(\mathbf{x}, \mathbf{X})(\tilde{\mathbf{K}}^{-1} - \tilde{\mathbf{C}}_i)k^{\sigma w}(\mathbf{X}, \cdot)\|_{\mathcal{H}_{k^{\sigma w}}}^2 \quad (\text{S102})$$

$$= k^{\sigma w}(\mathbf{x}, \mathbf{X})\tilde{\mathbf{K}}^{-1}\tilde{\mathbf{K}}\tilde{\mathbf{K}}^{-1}k^{\sigma w}(\mathbf{X}, \mathbf{x}) - 2k^{\sigma w}(\mathbf{x}, \mathbf{X})\tilde{\mathbf{K}}^{-1}\tilde{\mathbf{K}}\tilde{\mathbf{C}}_i k^{\sigma w}(\mathbf{X}, \mathbf{x}) + k^{\sigma w}(\mathbf{x}, \mathbf{X})\tilde{\mathbf{C}}_i\tilde{\mathbf{K}}\tilde{\mathbf{C}}_i k^{\sigma w}(\mathbf{X}, \mathbf{x}) \quad (\text{S103})$$

$$= k(\mathbf{x}, \mathbf{X})(\tilde{\mathbf{K}}^{-1} - \tilde{\mathbf{C}}_i)k(\mathbf{X}, \mathbf{x}) \quad (\text{S104})$$

$$= k_i^{\text{comp.}}(\mathbf{x}, \mathbf{x}) \quad (\text{S105})$$