
Bias-Free Scalable Gaussian Processes via Randomized Truncations

Andres Potapczynski^{*1} Luhuan Wu^{*2} Dan Biderman^{*1} Geoff Pleiss¹ John P. Cunningham^{1,2}

Abstract

Scalable Gaussian Process methods are computationally attractive, yet introduce modeling biases that require rigorous study. This paper analyzes two common techniques: early truncated conjugate gradients (CG) and random Fourier features (RFF). We find that both methods introduce a systematic bias on the learned hyperparameters: CG tends to underfit while RFF tends to overfit. We address these issues using randomized truncation estimators that eliminate bias in exchange for increased variance. In the case of RFF, we show that the bias-to-variance conversion is indeed a trade-off: the additional variance proves detrimental to optimization. However, in the case of CG, our unbiased learning procedure meaningfully outperforms its biased counterpart with minimal additional computation. Our code is available at <https://github.com/cunningham-lab/RTGPS>.

1. Introduction

Gaussian Processes (GP) are popular and expressive non-parametric models, and considerable effort has gone into alleviating their cubic runtime complexity. Notable successes include inducing point methods (e.g. Snelson & Ghahramani, 2006; Titsias, 2009; Hensman et al., 2013), finite-basis expansions (e.g. Rahimi & Recht, 2008; Mutn̄y & Krause, 2018; Wilson et al., 2020; Loper et al., 2020), nearest neighbor truncations (e.g. Datta et al., 2016; Katzfuss et al., 2021), and iterative numerical methods (e.g. Cunningham et al., 2008; Cutajar et al., 2016; Gardner et al., 2018). Common to these techniques is the classic *speed-bias tradeoff*: coarser GP approximations afford faster but more biased solutions that in turn affect both the model’s predictions and learned hyperparameters. While a few papers analyze the

bias of inducing point methods (Bauer et al., 2016; Burt et al., 2019), the biases of other approximation techniques, and their subsequent impact on learned GP models, have not been rigorously studied.

Here we scrutinize the biases of two popular techniques – random Fourier features (RFF) (Rahimi & Recht, 2008) and conjugate gradients (CG) (e.g. Cunningham et al., 2008; Cutajar et al., 2016; Gardner et al., 2018). These methods are notable due to their popularity and because they allow dynamic control of the speed-bias tradeoff: at any model evaluation, the user can adjust the number of CG iterations or RFF features to a desired level of approximation accuracy. In practice, it is common to truncate these methods to a fixed number of iterations/features that is deemed adequate. However, such truncation will stop short of an exact (machine precision) solution and potentially lead to biased optimization outcomes.

We provide a novel theoretical analysis of the biases resulting from RFF and CG on the GP log marginal likelihood objective. Specifically, we prove that CG is biased towards hyperparameters that underfit the data, while RFF is biased towards overfitting. In addition to yielding suboptimal hyperparameters, these biases hurt posterior predictions, regardless of the inference method used at test-time. Perhaps surprisingly, this effect is not subtle, as we will demonstrate. Our analysis suggests there is value in debiasing GP learning with CG and RFF. To do so, we turn to recent work that shows the merits of exchanging the speed-bias tradeoff for a *speed-variance tradeoff* (Beatson & Adams, 2019; Chen et al., 2019; Luo et al., 2020; Oktay et al., 2020). These works all introduce a randomization procedure that reweights elements of a fast truncated estimator, eliminating its bias at the cost of increasing its variance.

We thus develop bias-free versions of GP learning with CG and RFF using randomized truncation estimators. In short, we randomly truncate the number of CG iterations and RFF features, while reweighting intermediate solutions to maintain unbiasedness. Our variant of CG uses the **Russian Roulette** estimator (Kahn, 1955), while our variant of RFF uses the **Single Sample** estimator of Lyne et al. (2015). We believe our **RR-CG** and **SS-RFF** methods to be the first to produce unbiased estimators of the GP log marginal likelihood with $< \mathcal{O}(N^3)$ computation.

^{*}Equal contribution (randomly ordered). ¹Zuckerman Institute, Columbia University ²Statistics Department, Columbia University. Correspondence to: Andres Potapczynski <ap3635@columbia.edu>, Luhuan Wu <lw2827@columbia.edu>, Dan Biderman <db3236@columbia.edu>.

Finally, through extensive empirical evaluation, we find our methods and their biased counterparts indeed constitute a bias-variance tradeoff. Both RR-CG and SS-RFF are unbiased, recovering nearly the same optimum as the exact GP method, while GP trained with CG and RFF often converge to solutions with worse likelihood. We note that bias elimination is not always practical. For SS-RFF, the optimization is slow, due to the large auxiliary variance needed to counteract the slowly decaying bias of RFF. On the other hand, RR-CG incurs a minimal variance penalty, likely due to the favorable convergence properties of CG. In a wide range of benchmark datasets, RR-CG demonstrates similar or better predictive performance compared to CG using the same expected computational time.

To summarize, this work offers three main contributions:

- theoretical analysis of the bias of CG- and RFF-based GP approximation methods (§3)
- RR-CG and SS-RFF: bias-free versions of these popular GP approximation methods (§4)
- results demonstrating the value of our RR-CG and SS-RFF methods (§3 and 5).

2. Background

We consider observed data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ for $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$, and the standard GP model:

$$f(\cdot) \sim \mathcal{GP}(\mu(\cdot), k(\cdot, \cdot)),$$

$$\mathbf{y}_i = f(\mathbf{x}_i) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

where $k(\cdot, \cdot)$ is the covariance kernel, $\mu(\cdot)$ is set to zero without loss of generality, and hyperparameters are collected into the vector θ , which is optimized as:

$$\begin{aligned} \theta &= \arg \min_{\theta} \mathcal{L}(\theta) \\ \mathcal{L}(\theta) &= -\log p(\mathbf{y} | \mathbf{X}; \theta) \\ &= \frac{1}{2} \left[\underbrace{\log |\mathbf{K}_{\mathbf{X}\mathbf{X}}|}_{\text{model complexity}} + \underbrace{\mathbf{y}^\top \mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}}_{\text{data fit}} + N \log 2\pi \right] \end{aligned} \quad (1)$$

where $\mathbf{K}_{\mathbf{X}\mathbf{X}} \in \mathbb{R}^{N \times N}$ is the Gram matrix of all data points with diagonal observational noise:

$$\mathbf{K}_{\mathbf{X}\mathbf{X}}[i, j] = k(\mathbf{x}_i, \mathbf{x}_j) + \sigma^2 \mathbb{1}_{i=j}.$$

Following standard practice, we optimize θ with gradients:

$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{1}{2} \text{tr} \left[\mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \frac{\partial \mathbf{K}_{\mathbf{X}\mathbf{X}}}{\partial \theta} \right] - \mathbf{y}^\top \mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \frac{\partial \mathbf{K}_{\mathbf{X}\mathbf{X}}}{\partial \theta} \mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}. \quad (2)$$

Three terms thus dominate the computational complexity: $\mathbf{y}^\top \mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$, $\log |\mathbf{K}_{\mathbf{X}\mathbf{X}}|$, and $\text{tr} \left\{ \mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \frac{\partial \mathbf{K}_{\mathbf{X}\mathbf{X}}}{\partial \theta} \right\}$. The common

approach to computing this triad is the Cholesky factorization, requiring $\mathcal{O}(N^3)$ time and $\mathcal{O}(N^2)$ space.

Extensive literature has accelerated the inference and hyperparameter learning of GP. Two very popular strategies are using *conjugate gradients* (Cunningham et al., 2008; Cutajar et al., 2016; Gardner et al., 2018; Wang et al., 2019) to approximate the linear solves in Eq. (2), and *random Fourier features* (Rahimi & Recht, 2008), which constructs a randomized finite-basis approximation of the kernel.

2.1. Conjugate Gradients

To apply conjugate gradients to GP learning, we begin by replacing the gradient in Eq. (2) with a stochastic estimate (Cutajar et al., 2016; Gardner et al., 2018):

$$\frac{\partial \mathcal{L}}{\partial \theta} \approx \frac{1}{2} \left[\mathbf{z}^\top \mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \frac{\partial \mathbf{K}_{\mathbf{X}\mathbf{X}}}{\partial \theta} \mathbf{z} - \mathbf{y}^\top \mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \frac{\partial \mathbf{K}_{\mathbf{X}\mathbf{X}}}{\partial \theta} \mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y} \right], \quad (3)$$

where \mathbf{z} is a random variable such that $\mathbb{E}[\mathbf{z}] = \mathbf{0}$ and $\mathbb{E}[\mathbf{z}\mathbf{z}^\top] = \mathbf{I}$. Note that the first term constitutes a stochastic estimate of the trace term (Hutchinson, 1989). Thus, stochastic optimization of Eq. (1) can be reduced to computing the linear solves $\mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$ and $\mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{z}$.

Conjugate gradients (CG) (Hestenes et al., 1952) is an iterative algorithm for solving positive definite linear systems $\mathbf{A}^{-1} \mathbf{b}$. It consists of a three-term recurrence, where each new term requires only a matrix-vector multiplication with \mathbf{A} . More formally, each CG iteration computes a new term of the following summation:

$$\mathbf{A}^{-1} \mathbf{b} = \sum_{i=1}^N \gamma_i \mathbf{d}_i, \quad (4)$$

where the γ_i are coefficients and the \mathbf{d}_i are conjugate search directions (Golub & Van Loan, 2012). N iterations of CG produce all N summation terms and recover the exact solution. In practice, exact convergence may require more than N iterations due to inaccuracies of floating point arithmetic. However, the summation converges exponentially, and so $J \ll N$ iterations may suffice to achieve high accuracy.

CG is an appealing method for computing $\mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$ and $\mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{z}$ due to its computational complexity and its potential for GPU-accelerated matrix products. J iterations takes at most $\mathcal{O}(JN^2)$ time and $\mathcal{O}(N)$ space if the matrix-vector products are performed in a map-reduce fashion (Wang et al., 2019). However, ill-conditioned kernel matrices hinder the convergence rate (Cutajar et al., 2016), and so the J^{th} CG iteration may yield an inaccurate approximation of Eq. (3).

2.2. Random Fourier Features

Rahimi & Recht (2008) introduce a randomized finite-basis approximation to stationary kernels:

$$k(\mathbf{x}, \mathbf{x}^\theta) = k(\mathbf{x} - \mathbf{x}^\theta) \approx \left(\mathbf{x} \right)^\top \left(\mathbf{x}^\theta \right) \quad (5)$$

where $(\mathbf{x}) \in \mathbb{R}^J$ and $J \ll N$. The RFF approximation relies on Bochner’s theorem (Bochner et al., 1959): letting $\tau = \mathbf{x} - \mathbf{x}^\ell$, all stationary kernels $k(\tau)$ on \mathbb{R}^d can be exactly expressed as the Fourier dual of a nonnegative measure $P(\ell)$: $k(\tau) = \int P(\ell) \exp(i\ell^\top \tau) d\ell$. A Monte Carlo approximation of this Fourier transform yields:

$$k(\tau) \approx \frac{2}{J} \sum_{j=1}^{\lfloor J/2 \rfloor} \exp(i\ell_j^\top \tau), \quad \ell_j \sim P(\ell),$$

which simplifies to a finite-basis approximation:

$$\mathbf{K}_{\mathbf{X}\mathbf{X}} \approx [(\mathbf{x}_1) \dots (\mathbf{x}_n)] [(\mathbf{x}_1) \dots (\mathbf{x}_n)]^\top, \\ (\mathbf{x}) = [\cos \ell_j^\top \mathbf{x}, \sin \ell_j^\top \mathbf{x}]_{j=1}^{J/2}, \quad \ell_j \sim P(\ell).$$

For many common kernels, $P(\ell)$ can be computed in closed-form (e.g. RBF kernels have zero-mean Gaussian spectral densities). The approximated log likelihood can be computed in $\mathcal{O}(J^3 + N)$ time and $\mathcal{O}(JN)$ space using the Woodbury inversion lemma and the matrix determinant lemma, respectively. The number of random features $J/2$ is a user choice, with typical values between 100-1000. More features lead to more accurate kernel approximations.

2.3. Unbiased Randomized Truncation

We will now briefly introduce Randomized Truncation Estimators, which are the primary tool we use to unbiased the CG and RFF log marginal likelihood estimates. At a high level, assume that we wish to estimate some quantity ψ that can be expressed as a (potentially-infinite) series:

$$\psi = \sum_{j=1}^H \Delta_j, \quad H \in \mathbb{N} \cup \{\infty\}.$$

Here and in the following sections, Δ_j can either be random or deterministic. To avoid the expensive evaluation of the full summation, a randomized truncation estimator chooses a random term $J \in \{1, \dots, H\}$ with probability mass function $P(J) = P(\mathcal{J} = J)$ after which to truncate computation. In the following, we introduce two means of deriving unbiased estimators by upweighting the summation terms.

The Russian Roulette estimator (Kahn, 1955) obtains an unbiased estimator $\bar{\psi}_J$ by truncating the sum after $J \sim P(J)$ terms and dividing the surviving terms by their survival probabilities:

$$\bar{\psi}_J = \sum_{j=1}^{\mathcal{J}} \frac{\Delta_j}{P(\mathcal{J} \geq j)} = \sum_{j=1}^{\mathcal{J}} \frac{1_{\mathcal{J} \geq j}}{P(\mathcal{J} \geq j)} \Delta_j, \quad (6)$$

and, $E[\bar{\psi}_J] = \sum_{i=1}^H \Delta_i = \psi$. (See appendix for further derivation.) The choice of $P(J)$ determines both the computational efficiency and the variance of $\bar{\psi}_J$. A thin-tailed $P(J)$ will often truncate sums after only a few terms ($J \ll H$).

However, tail events ($J \approx H$) are upweighted inversely to their low survival probability, and so thin-tailed truncation distributions may lead to high variance.

The Single Sample estimator (Lyne et al., 2015) implements an alternative reweighting scheme. After drawing $J \sim P(J)$, it computes a single summation term Δ_J , which it upweights by $1/P(J)$:

$$\bar{\psi}_J = \frac{\Delta_J}{P(J)} = \sum_{j=1}^{\mathcal{J}} \frac{1_{\mathcal{J}=j}}{P(\mathcal{J}=j)} \Delta_j. \quad (7)$$

This procedure is unbiased, and it amounts to estimating ψ using a single importance-weighted sample from $P(J)$ (see appendix). Again, $P(J)$ controls the speed/variance trade-off. We refer the reader to (Beatson & Adams, 2019) for a detailed comparison of these two estimators. We emphasize that both estimators remain unbiased even if Δ_j is a random variable, as long as it is independent from the random truncation integer J .

3. GP Learning with CG and RFF is Biased

Here we prove that early truncated CG and RFF provide biased approximations to the terms comprising the GP log marginal likelihood (Eq. 1). We also derive the bias decay rates for each method. We then empirically demonstrate these biases and show they affect the hyperparameters learned through optimization. Remarkably, we find that the above biases are *highly systematic*: CG-based GP learning favors underfitting hyperparameters while RFF-based learning favors overfitting hyperparameters.

3.1. CG Biases GP Towards Underfitting

In the GP literature, CG has often been considered an “exact” method for computing the log marginal likelihood (Cunningham et al., 2008; Cutajar et al., 2016), as the iterations are only truncated after reaching a pre-specified residual error threshold (e.g. 10^{-10}). However, as CG is applied to ever-larger kernel matrices it is common to truncate the CG iterations before reaching this convergence (Wang et al., 2019). While this accelerates the hyperparameter learning process, the resulting solves and gradients can no longer be considered “exact.” In what follows, we show that the early-truncated CG optimization objective is not only approximate but also systematically biased towards underfitting.

To analyze the early-truncation bias, we adopt the analysis of Gardner et al. (2018) that recovers the GP log marginal likelihood (Eq. 1) from the stochastic gradient’s CG estimates of $\hat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$ and $\hat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{z}$ (Eq. 3). Recall the two terms in the log marginal likelihood are the “data fit” term $\mathbf{y}^\top \hat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$ and the “model complexity” term $\log |\hat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}|$. The first term falls directly out of the CG estimate of $\hat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$, while a

Figure 1. CG (left) systematically overestimates $\log j \kappa_{xx}^{-1} j$ and underestimates $\kappa_{xx}^{-1} y$ whereas RFF (right) does the opposite. The dashed orange line shows the exact values computed by Cholesky. Our unbiased methods, RR-CG (left) and SS-RFF (right), recover the true $\log j \kappa_{xx}^{-1} j$ and $\kappa_{xx}^{-1} y$ values. For these two methods, the x-axis indicates the expected number of iterations/features.

If $J < N$, CG underestimates the inverse quadratic term and overestimates the log determinant in expectation:

$$u_J = y^T \kappa_{xx}^{-1} y; \quad E_z[v_J] = \log j \kappa_{xx}^{-1} j; \quad (9)$$

The biases of both terms decay at a rate of $O(C^{-2J})$, where C is a constant that depends on the conditioning of κ_{xx} .

Proof sketch. The direction of the biases can be proved using a connection between CG and numeric quadrature. u_J and v_J are exactly equal to the point Gauss quadrature approximation of $y^T \kappa_{xx}^{-1} y$ and $\int_0^1 \log T_z^{(J)} e_1$.

Figure 2. Kernel lengthscale values learned by optimizing (biased) CG and RFF log marginal likelihood approximations. CG overestimates the optimal kernel lengthscale whereas RFF underestimates it. We plot the divergence (in log-ratio scale) between the learned and true lengthscales as a function of the number of CG iterations (left) and of the number of RFF samples (right).

represented as Riemann-Stieltjes integrals. The sign of the CG approximation bias follows from the standard Gauss quadrature error bound, which is negative for $\kappa_{xx}^{-1} y$ and positive for $\log j \kappa_{xx}^{-1} j$. The convergence rates are from standard bounds on CG (Golub & Van Loan, 2012) and the analysis of Ubaru et al. (2017). See appendix for a full proof.

stochastic estimate of $\log j \kappa_{xx}^{-1} j$ can be obtained through the byproducts of CG's computation for $\kappa_{xx}^{-1} z$. Gardner et al. (2018) show that the CG coefficients in Eq. (4) can be manipulated to produce a partial tridiagonalization of $T_z^{(J)} = Q_z^{(J)T} \kappa_{xx} Q_z^{(J)}$; where $T_z^{(J)} \in \mathbb{R}^{J \times J}$ is tridiagonal. $T_z^{(J)}$ can compute the Stochastic Lanczos Quadrature estimate of $\kappa_{xx}^{-1} z$ (Ubaru et al., 2017; Dong et al., 2017):

$$\log j \kappa_{xx}^{-1} j = E \left[z^T (\log T_z^{(J)}) z \right] + \int_0^1 \log T_z^{(J)} e_1; \quad (8)$$

where $\log(\cdot)$ is the matrix logarithm and e_1 is the first row of the identity matrix. The following theorem analyzes the bias of these $\kappa_{xx}^{-1} y$ and $\log j \kappa_{xx}^{-1} j$ estimates:

Theorem 1. Let u_J and v_J be the estimates of $\kappa_{xx}^{-1} y$ and $\log j \kappa_{xx}^{-1} j$ respectively after J iterations of CG; i.e.:

$$u_J = y^T \prod_{i=1}^J d_i^{-1}; \quad v_J = \int_0^1 \log T_z^{(J)} e_1;$$

Fig. 1 confirms our theoretical analysis and demonstrates the systematic biases of CG. We plot the log marginal likelihood terms for a subset of the PoleTele UCI dataset, varying the number of CG iterations J used to produce the estimates. Compared against the exact terms (computed with Cholesky), we see an overestimation of $\log j \kappa_{xx}^{-1} j$ and an underestimation of $\kappa_{xx}^{-1} y$. These biases are most prominent when using few CG iterations.

We turn to study the effect of CG learning on hyperparameters. Since the log marginal likelihood is a nonconvex function of θ , it is not possible to directly prove how the bias affects θ . Nevertheless, we know intuitively that underestimating $\kappa_{xx}^{-1} y$ de-prioritizes model fit while overestimating $\log j \kappa_{xx}^{-1} j$ over-penalizes model complexity. Thus, the learned hyperparameters will likely underfit the data. Underfitting may manifest in an overestimation of the learned lengthscale as low values of θ increase the flexibility and the complexity of the model. This hypothesis is empirically confirmed in Fig. 2 (left panel). We train a GP regression model on a toy dataset $y = \sin(5x) + \epsilon$

and $\mathbf{N} \sim \mathcal{N}(0; 0.01)$. We fix all hyperparameters other than the lengthscale, which is learned using both CG-based optimization and (exact) Cholesky-based optimization. The overestimation of decays with the number of CG iterations.

3.2. RFF Biases GP Towards Overfitting

Previous work has studied the accuracy of RFF's approximation to the entries of the Gram matrix $\mathbf{K}(x; x')$ (Rahimi & Recht, 2008; Sutherland & Schneider, 2015). However, to the best of our knowledge there has been little analysis of nonlinear functions of this approximate Gram matrix, such as $\mathbf{K}_{XX}^{-1}y$ and $\log \det \mathbf{K}_{XX}$ appearing in the GP objective. Interestingly, we find that RFF systematically biases these terms:

Theorem 2. Let \mathbf{K}_J be the RFF approximation with $m=2$ random features. In expectation, \mathbf{K}_J overestimates the inverse quadratic and underestimates the log determinant:

$$\mathbb{E}_{P(\mathcal{J})} \mathbf{y}^T \mathbf{K}_J^{-1} \mathbf{y} > \mathbf{y}^T \mathbf{K}_{XX}^{-1} \mathbf{y} \quad (10)$$

$$\mathbb{E}_{P(\mathcal{J})} \log \det \mathbf{K}_J < \log \det \mathbf{K}_{XX} \quad (11)$$

The biases of both terms decay at a rate $\mathcal{O}(1/J)$.

Proof sketch. The direction of the biases is a straightforward application of Jensen's inequality, noting that \mathbf{K}_{XX}^{-1} is a convex function and $\log \det \mathbf{K}_{XX}$ is a concave function. The magnitude of the bias is derived from a second-order Taylor expansion that closely resembles the analysis of Nowozin (2018). See appendix for full proof.

Again, Fig. 1 confirms the systematic biases of RFF, which decay at a rate proportional to the number of features, as predicted by Thm. 2. Hence, RFF should affect the learned hyperparameters in a manner opposite to CG. Overestimating $\mathbf{y}^T \mathbf{K}_{XX}^{-1} \mathbf{y}$ emphasizes data fitting while underestimating $\log \det \mathbf{K}_{XX}$ reduces the model complexity penalty, overall resulting in overfitting behavior. Following the intuition presented in Sec. 3.1, we expect the lengthscale to be underestimated, as empirically confirmed by Fig. 2 (right panel). The figure also illustrates the slow decay of the RFF bias.

4. Bias-free Scalable Gaussian Processes

We debias the estimates of both the GP training objective in Eq. (1) and its gradient in Eq. (2) (as approximated by CG and RFF) using unbiased randomized truncation estimators. To see how such estimators apply to GP hyperparameter learning, we note that both CG and RFF recover the true log marginal likelihood (or an unbiased estimate thereof) in their limits:

Observation 1. CG recovers the exact log marginal likeli-

hood in expectation in at most iterations:

$$\mathbf{y}^T \mathbf{K}_{XX}^{-1} \mathbf{y} = \mathbf{y}^T \sum_{j=1}^P \mathbf{d}_j \mathbf{d}_j^T \mathbf{y} \quad (12)$$

$$\log \det \mathbf{K}_{XX} = \mathbb{E}_z \sum_{j=1}^P \log \|\mathbf{z} + \mathbf{e}_j\|_2^2 \quad (13)$$

By the law of large numbers, RFF converges almost surely to the exact log marginal likelihood as the number of random features goes to infinity:

$$\mathbf{y}^T \mathbf{K}_{XX}^{-1} \mathbf{y} = \lim_{J \rightarrow \infty} \mathbf{y}^T \mathbf{K}_J^{-1} \mathbf{y} \quad (14)$$

$$\log \det \mathbf{K}_{XX} = \lim_{J \rightarrow \infty} \log \det \mathbf{K}_J \quad (15)$$

To maintain the scalability of CG and RFFs while eliminating bias, we express the log marginal likelihood terms in Eqs. (12) to (15) as summations amenable to randomized truncation. We then apply the Russian Roulette and Single Sample estimators of Sec. 2.3 to avoid computing all summation terms while obtaining the same result in expectation.

4.1. Russian Roulette-Truncated CG (RR-CG)

The stochastic gradient in Eq. (3) requires performing two solves: $\mathbf{K}_{XX}^{-1} \mathbf{y}$ and $\mathbf{K}_{XX}^{-1} \mathbf{z}$. Using the summation formulation of CG (Eq. 4), we can write these two solves as series:

$$\mathbf{K}_{XX}^{-1} \mathbf{y} = \sum_{i=1}^P \mathbf{d}_i \mathbf{d}_i^T \mathbf{y}; \quad \mathbf{K}_{XX}^{-1} \mathbf{z} = \sum_{i=1}^P \mathbf{d}_i \mathbf{d}_i^T \mathbf{z}$$

where each CG iteration computes a new term of the summation. By applying the Russian Roulette estimator from Eq. (6), we obtain the following unbiased estimates:

$$\mathbf{K}_{XX}^{-1} \mathbf{y} \approx \sum_{j=1}^J \frac{P}{j} \mathbf{d}_j \mathbf{d}_j^T \mathbf{y} = P \sum_{j=1}^J \frac{1}{j} \mathbf{d}_j \mathbf{d}_j^T \mathbf{y}; \quad \mathbf{K}_{XX}^{-1} \mathbf{z} \approx \sum_{j=1}^J \frac{P^0}{j} \mathbf{d}_j \mathbf{d}_j^T \mathbf{z} = P^0 \sum_{j=1}^J \frac{1}{j} \mathbf{d}_j \mathbf{d}_j^T \mathbf{z} \quad (16)$$

These unbiased solves produce an unbiased optimization gradient in Eq. (3); we refer to this approach as Russian Roulette CG (RR-CG). With the appropriate truncation distribution $P(\mathcal{J})$, this estimate affords the same computational complexity of standard CG without its bias.

We must compute two independent estimates of $\mathbf{K}_{XX}^{-1} \mathbf{y}$ with different $J = P(\mathcal{J})$ in order for the $\mathbf{y}^T \mathbf{K}_{XX}^{-1} \mathbf{y}$ term in Eq. (3) to be unbiased. Thus, the unbiased gradient requires 3 calls to RR-CG, as opposed to the 2 CG calls needed for the biased gradient. Nevertheless, RR-CG has the same $\mathcal{O}(JN^2)$ complexity as standard CG – and the additional solve can be computed in parallel with the others.

We can also use the Russian Roulette estimator to compute the log marginal likelihood itself, though this is not strictly necessary for gradient-based optimization. (See appendix.)

Choosing the truncation distribution. Since the Russian Roulette estimator is unbiased for any choice of $P(\mathcal{J})$,

we wish to choose a truncation distribution that balances computational cost and variance. In Fig. 1 we plot the empirical mean of the computational cost and variance of the RR-CG estimator using 10^4 samples from an exponential distribution. We find that RR-CG produces unbiased estimates of the $\log \mathbb{E}[\mathbb{K}_{XX}^{-1} y]$ and $\log \mathbb{E}[\mathbb{K}_{XX}^{-1} j]$ terms that are indistinguishable from the exact values computed with Cholesky. Reducing the expected truncation iteration (the x-axis) increases the standard error of empirical means, demonstrating the speed-variance trade-off. We summarize our estimates and choices of distribution as follows:

Theorem 3. The approximation $\log \mathbb{E}[\mathbb{K}_{XX}^{-1} j]$ and $\log \mathbb{E}[\mathbb{K}_{XX}^{-1} y]$ using RR-CG decays at a rate of $O(C^{-2J})$. Therefore the truncation distribution that maximizes the ROE is $P(J) / C^{-2J}$, where C is a constant that depends on the conditioning of \mathbb{K}_{XX} . The expected computation and variance of $P(J)$ is finite.

Proof sketch Beatson & Adams (2019) show that the truncation distribution that maximizes the ROE is proportional to the rate of decay of our approximation divided by its computational cost. The error of CG decays as $O(C^{-2J})$, and the cost of each summation term is constant with respect to J . In practice, we vary the exponential decaying rate to control the expectation and variance of $P(J)$. To further reduce the variance of our estimator, we set a minimum number of CG iterations to be computed, as in (Luo et al., 2020). See appendix for full proof.

In practice, however, we do not have access to \mathbb{K}_{XX} since computing the conditioning of \mathbb{K}_{XX} is impractical. Yet, we can change the base of the exponential and add a temperature parameter to rescale the function and control the rate of decay of the truncation distribution as a sensible alternative. Thus, we follow a more general exponential decay distribution:

$$P(J) / e^{-J}; \quad J = J_{\min}; \quad ; N \quad (17)$$

where J_{\min} is the minimum truncation number. By varying the values of β and J_{\min} , we can control the expectation and standard deviation of $P(J)$. In practice we found that having the standard deviation value between 10 and 20 achieves stable GP learning process, which can be obtained by tuning β between 0.05 and 0.1 for sufficiently large datasets (e.g. $N = 500$). We also noticed that the method is not sensitive to these choices of hyperparameters and that they work well across all the experiments. The expected truncation number can be further tuned by varying β . We emphasize that these choices impact the speed-variance tradeoff. By setting a large J_{\min} we decrease the speed by requiring more baseline computations but also decrease the variance (since the minimum approximations have the largest deviations from the ground truth).

4.2. Single Sample-Truncated RFF (SS-RFF)

Denoting \mathbb{K}_j as the kernel matrix estimated by random Fourier features, we can write $\log \mathbb{E}[\mathbb{K}_{XX}^{-1} j]$ as the following telescoping series:

$$\log \mathbb{E}[\mathbb{K}_{XX}^{-1} j] = \log j \mathbb{K}_{1j} + \sum_{j=2}^{N-1} \log j \mathbb{K}_{jj} - \log j \mathbb{K}_{j-1j} \quad (18)$$

$$+ \log j \mathbb{K}_{XX}^{-1} j - \log j \mathbb{K}_{N-2-1j} \\ = \log j \mathbb{K}_{1j} + \sum_{j=2}^{N-2} \log j \mathbb{K}_{jj} - \log j \mathbb{K}_{j-1j} \quad (19)$$

where j is defined as $\log j \mathbb{K}_{jj} - \log j \mathbb{K}_{j-1j}$ for all $j < N-2$, and $\log j \mathbb{K}_{N-2-1j}$ is defined as $\log j \mathbb{K}_{XX}^{-1} j - \log j \mathbb{K}_{N-2-1j}$. Note that each j is now a random variable, since it depends on the random Fourier frequencies. Crucially, we only include $N-2$ terms in the series so that no term requires more than $O(N^3)$ computation in expectation. (For any $N = 2$, \mathbb{K}_j is a full-rank matrix and thus is as computationally expensive as the true \mathbb{K}_{XX} matrix.) We construct a similar telescoping series for $\log \mathbb{E}[\mathbb{K}_{XX}^{-1} y]$.

As with Eq. (16), we approximate the series in Eq. (19) with a randomized truncation estimator, this time using the Single Sample estimator (7):

$$\log \mathbb{E}[\mathbb{K}_{XX}^{-1} j] \approx \log j \mathbb{K}_{1j} + \sum_{j=2}^J \log j \mathbb{K}_{jj} - \log j \mathbb{K}_{j-1j} = P(J) \quad (20)$$

where J is drawn from the truncation distribution $P(J)$ with support over $\{2, 3, \dots, N-2\}$. Note that the Single Sample estimator requires computing $3 \log$ determinants ($\log j \mathbb{K}_{1j}$, $\log j \mathbb{K}_{J-1j}$, and $\log j \mathbb{K}_{Jj}$) for a total of $O(J^3 + NJ^2)$ computations and $O(NJ)$ memory. This is asymptotically the same requirement as standard RFF. The Russian Roulette estimator, on the other hand, incurs a computational cost of $O(NJ^3 + J^4)$ as it requires computing $\log j \mathbb{K}_{1j}$ through $\log j \mathbb{K}_{Jj}$ which quickly becomes impractical for large J . A similar Single Sample estimator constructs an unbiased estimate of $\log \mathbb{E}[\mathbb{K}_{XX}^{-1} y]$. Backpropagating through these Single Sample estimates produces unbiased estimates of the log marginal likelihood gradient in Eq. (2).

¹Solely minimizing variance is not appealing, as this is achieved by $P(J) = \delta_{J=N}$ which has no computational savings. Choosing the truncation distribution. For the Single Sample estimator we do not have to optimize the ROE since

minimizing the variance of this estimator does not result in a degenerate distribution.

Theorem 4. The truncation distribution that minimizes the variance of the SS-RFF estimators $\log j \mathbb{R}_{XX} j$ and $y^T \mathbb{R}_{XX}^{-1} y$ is $P(J) / 1=J$. The expected variance and computation of $P(J)$ is finite.

Proof sketch. The minimum variance distribution can be found by solving a constrained optimization problem. In practice, we can further decrease the variance of our estimator by fixing a minimum value of RFF features to be used in Eq. (19) and by increasing the step size $\alpha(N)$ between the elements at each $j = \log j \mathbb{R}_{c,j} - \log j \mathbb{R}_{c(j-1),j}$. See appendix for full proof.

For the experiments we started with 500 features and also tried various step sizes $\alpha \in \{1, 10, 100\}$. The variance of the estimator decreases as we increase the probability of features and taking long steps $\alpha = 100$, the variance of the estimator still requires taking several steps before converging, making SS-RFF computationally impractical.

Toy problem. Similar to RR-CG, in Fig. 1 we plot the empirical mean of the SS-RFF estimator using 100 samples. We find that SS-RFF produces unbiased estimates of the $y^T \mathbb{R}_{XX}^{-1} y$ and $\log j \mathbb{R}_{XX} j$ terms. However, these estimates have a higher variance when compared to the estimates of RR-CG. Reducing the expected truncation iterations $E(J)$ (x-axis) increases the standard error of the empirical means, demonstrating the speed-variance trade-off.

4.3. Analysis of the Bias-free Methods

Randomized truncations and conjugate gradients have existed for many decades (Hestenes et al., 1952; Kahn, 1955), but have rarely been used in conjunction. Filippone & Engler (2015) proposed a closely related method, combining randomized early-truncation with CG to draw GP posterior samples over covariance parameters. We differ by providing the first theoretical analysis of the biases on GP learning incurred by CG and RFF, and proceed to remedy these biases using the Russian Roulette and the Simple Sampling estimators.

To some extent, randomized truncation methods are antithetical to the original intention of CG: producing deterministic and nearly exact solves. For large-scale applications, where early truncation is necessary for computational tractability, the ability to trade bias for variance is beneficial. This fact is especially true in the context of GP learning, where the bias of early truncation is systematic and cannot be simply explained away as numerical imprecision.

Randomized truncation estimates are often used to estimate

Figure 3. Optimization landscape of a GP with two hyperparameters. The SS-RFF and RR-CG models converge to similar hyperparameter values that are nearly optimal, while the RFF and CG models converge to suboptimal solutions. In addition, the stochastic effect of the randomized truncation is visible in the trajectories of RR-CG and SS-RFF. Moreover, CG (and RR-CG) models truncate after 20 iterations (in expectation); RFF (and SS-RFF) models use 700 features (in expectation).

in finite series, where it is challenging to design truncation distributions with finite expected computation and/or variance. We avoid such issues since CG and the telescoping RFF summations are both finite.

5. Results

First, we show that our bias-free methods recover nearly the same hyperparameters as exact methods (i.e. Cholesky-based optimization), whereas models that use CG and RFF converge to suboptimal hyperparameters. Since RR-CG and SS-RFF eliminate bias at the cost of increased variance, we then demonstrate the optimization convergence rate and draw conclusions on our methods' applicability. Finally, we compare models optimized with RR-CG against a host of approximate GP methods across a wide range of UCI datasets (Asuncion & Newman, 2007). All experiments are implemented in GPyTorch (Gardner et al., 2018)

Optimization trajectories of bias-free GP. Fig. 3 displays the optimization landscape – the log marginal likelihood of the PoleTele dataset – as a function of (RBF) kernel lengthscales and noise σ^2 . As expected, an exact GP (optimized using Cholesky, Fig. 3 upper left) recovers the optimum. Notably, the GP trained with standard CG and RFF converges to suboptimal hyperparameters (upper right/lower right). RR-CG and SS-RFF models (trained with 20 iterations and 700 features in expectation, respectively) successfully eliminate this bias, and recover nearly

Figure 4. The GP optimization objective for models trained with RFF and SS-RFF. (PoleTele dataset, RBF kernel, Adam optimizer). RFF models converge to sub-optimal log marginal likelihoods. SS-RFF models converge to (near) optimum values, yet require more than 100 as many optimization steps.

Figure 5. The GP optimization objective for models trained with CG and RR-CG. (Bike dataset, RBF kernel, Adam optimizer). RR-CG models converge to optimal solutions, while the (biased) CG models diverge. Increasing the expected truncation of RR-CG only slightly improves optimization convergence; models converge in < 100 steps of Adam.

the same parameters as the exact model (upper center/lower left). These plots also show the speed-variance tradeoff of randomized truncation. SS-RFF and RR-CG have noisy optimization trajectories due to auxiliary truncation variance.

Convergence of GP hyperparameter optimization.

Figs. 4 and 5 plot the exact GP log marginal likelihood of the parameters learned by each method during optimization. Each trajectory corresponds to a RBF-kernel GP trained on the PoleTele dataset (Fig. 4) and the Bike dataset (Fig. 5).

Fig. 4 shows that RFF models converge to solutions with worse log likelihoods, and more RFF features slow the rate of optimization. Additionally, we see the cost of the auxiliary variance needed to debias SS-RFF: while SS-RFF models achieve better optima than their biased counterparts, they take 2-3 orders of magnitude longer to converge, despite using a truncation distribution that minimizes variance. We thus conclude that SS-RFF has too much variance to be practical for GP hyperparameter learning.

Fig. 5 on the other hand shows that RR-CG is minimally affected by its auxiliary variance. The GP trained with RR-CG converges in roughly 100 iterations, nearly matching Cholesky-based optimization. Decreasing the expected truncation value from $E[J] = 40$ to 20 slightly slows this convergence. We note that the bias induced by standard CG can be especially detrimental to GP learning. On this dataset, the biased models deviate from their Cholesky counterparts and eventually diverge away from the optimum.

Predictive performance of bias-free GP. Lastly, we compare the predictive performance of GPs that use RR-CG, CG, and Cholesky for hyperparameter optimization. We emphasize that the RR-CG and CG methods only make use of early truncation approximations during training. At test time, we compute the predictive posterior by running CG to a tolerance of 10^{-4} , which we believe can be considered “exact” for practical intents and purposes. Additionally, we include four other (biased) scalable GP approximations methods as baseline: RFF, Stochastic Variational Gaussian Processes (SVGP) (Hensman et al., 2013), generalized Product of Expert Gaussian Processes (POE) (Cao & Fleet, 2014; Deisenroth & Ng, 2015), and stochastic gradient-based Gaussian Processes (sgGP) (Chen et al., 2020). We note that the RFF, SVGP, and sgGP methods introduce bias and variance, as these methods rely on randomization and approximation.

We use CG with $J = 100$ iterations, and RR-CG with $E[J] = 100$ expected iterations; both methods use the preconditioner of Gardner et al. (2018). All RFF models use 700 random features. For SVGP, we use 2024 inducing points and minibatches of size 2024 as in (Wang et al., 2019). The POE models are comprised of GP experts that are each trained on 2024 data point subsets. For sgGP, the subsampled datasets are constructed by selecting a random point $x; y$ and its 15 nearest neighbors as in (Chen et al., 2020). Each dataset is randomly split to 64% training, 16% validation and 20% testing sets. All kernels are RBF with a separate lengthscale per dimension. See appendix for more details.

We report prediction accuracy (RMSE) and negative log likelihood (NLL) in Fig. 6 (see appendix for full tables on predictive performance and training time). We make two key observations: (i) RR-CG meaningfully debiases CG.

When the bias of CG is not detrimental to optimization (e.g. CG with 100 iterations is close to convergence for the Elevators dataset), RR-CG has similar performance. However, when the CG bias is more significant (e.g. the KEGG dataset), the bias-free RR-CG improves the GP predictive RMSE and NLL. We also include a figure displaying the predictive performance of RR-CG and CG with increasing number of (expected) CG iterations in appendix (ii). RR-CG recovers the same optimum as the “ground-truth” method

Figure 6. Root-mean-square-error (RMSE) and negative log likelihood (NLL) of GP trained with CG (light purple), RR-CG (dark purple) and various approximate methods (grey). Dashed red lines indicates Cholesky-based GP performance (when applicable). Results are averaged over 3 dataset splits. Missing RFF and sgGP results correspond to (very high) outlier NLL / RMSE values. In almost all experiments, GP learning with RR-CG achieves similar or better performance compared to that with CG at the same computational cost.

(i.e. Cholesky) does, as indicated by the red-dashed line in training, we find that variance is almost always preferable to bias. Models trained with RR-CG achieve better performance than those trained with standard CG, and tend to recover the hyperparameters learned with exact methods.

While RR-CG obtains the lowest RMSE on all but 2 datasets, we note that the (biased) GP approximations sometimes achieve lower NLL. For example, SVGP has a lower NLL than that of RR-CG on the Bike dataset, despite having a higher RMSE. We emphasize that this is not a failing of RR-CG inference. The SVGP NLL is even better than that of the exact (Cholesky) GP, suggesting a potential model misspecification for this particular dataset. Since SVGP overestimates the observational noise (Bauer et al., 2016), it may obtain a better NLL when outliers are abundant.

Though we cannot compare against the Cholesky posterior on larger datasets, we hypothesize that the NLL/RMSE discrepancy on these datasets is due to a similar modeling issue. We note that the RR-CG algorithm is not limited to GP applications. Future work should explore applying RR-CG to other optimization problems with large-scale solves.

6. Conclusion

We prove that CG and RFF introduce systematic biases to the GP log marginal likelihood objective: CG-based training will favor underfitting models, while RFF-based training will promote overfitting. Modifying these methods with randomized truncation converts these biases into variance, enabling unbiased stochastic optimization. Our results show that this bias-to-variance exchange indeed constitutes a trade-off. The convergence of SS-RFF is impractically slow, likely due to the truncation variance needed to eliminate RFF's slowly-decaying bias. However, for CG-based

Acknowledgements

This work was supported by the Simons Foundation, McKnight Foundation, the Grossman Center, and the Gatsby Charitable Trust.

References

- Angelova, J. A. On moments of samples mean and variance. In *Int. J. Pure Appl. Math* pp. 79:67–85, 2012.
- Asuncion, A. and Newman, D. Uci machine learning repository, 2007.

- Bauer, M., van der Wilk, M., and Rasmussen, C. E. Under-
standing probabilistic sparse gaussian process approxi-
mations. In *Advances in Neural Information Processing
Systems*, 2016.
- Beatson, A. and Adams, R. P. Efficient optimization of
loops and limits with randomized telescoping sums. In
International Conference on Machine Learning, 2019.
- Bochner, S. et al. *Lectures on Fourier integrals* volume 42.
Princeton University Press, 1959.
- Burt, D., Rasmussen, C. E., and Van Der Wilk, M. Rates
of convergence for sparse variational gaussian process re-
gression. In *International Conference on Machine Learn-
ing*, pp. 862–871, 2019.
- Cao, Y. and Fleet, D. J. Generalized product of experts
for automatic and principled fusion of gaussian process
predictions. arXiv preprint arXiv:1410.7827, 2014.
- Chen, H., Zheng, L., Al Kontar, R., and Raskutti, G. Stochas-
tic gradient descent in correlated settings: A study on
gaussian processes. *Advances in Neural Information Pro-
cessing Systems*, 33, 2020.
- Chen, R. T. Q., Behrmann, J., Duvenaud, D., and Jacobsen,
J.-H. Residual flows for invertible generative model-
ing. *Advances in Neural Information Processing Systems*
2019.
- Cunningham, J. P., Shenoy, K. V., and Sahani, M. Fast
gaussian process methods for point process intensity esti-
mation. In *International Conference on Machine learning*
pp. 192–199, 2008.
- Cutajar, K., Osborne, M. A., Cunningham, J. P., and Filip-
ppone, M. Preconditioning kernel matrices. *Internat-
ional Conference on Machine Learning*, 2016.
- Datta, A., Banerjee, S., Finley, A. O., and Gelfand, A. E.
Hierarchical nearest-neighbor gaussian process models
for large geostatistical datasets. *Journal of the American
Statistical Association*, 111(514):800–812, 2016.
- Deisenroth, M. and Ng, J. W. Distributed gaussian processes.
In *International Conference on Machine Learning*, pp.
1481–1490. PMLR, 2015.
- Dong, K., Eriksson, D., Nickisch, H., Bindel, D., and Wil-
son, A. G. Scalable log determinants for gaussian process
kernel learning. In *Advances in Neural Information Pro-
cessing Systems*, pp. 6327–6337, 2017.
- Filippone, M. and Engler, R. Enabling scalable stochas-
tic gradient-based inference for gaussian processes by
employing the unbiased linear system solver (ulisse). In
International Conference on Machine Learning, pp. 1015–
1024. PMLR, 2015.
- Gardner, J., Pleiss, G., Weinberger, K. Q., Bindel, D., and
Wilson, A. G. Gpytorch: Blackbox matrix-matrix gaus-
sian process inference with gpu acceleration. *Ad-
vances in Neural Information Processing Systems*
7576–7586, 2018.
- Golub, G. H. and Meurant, G. *Matrices, moments and
quadrature with applications* volume 30. Princeton Uni-
versity Press, 2009.
- Golub, G. H. and Van Loan, C. *Matrix Computations*
The Johns Hopkins University Press, 4th Edition, 2012.
- Hensman, J., Fusi, N., and Lawrence, N. D. Gaussian
processes for big data. *Uncertainty in Arti cial Intelli-
gence* 2013.
- Hestenes, M. R., Stiefel, E., et al. Methods of conjugate
gradients for solving linear systems. *Journal of Research
of the National Bureau of Standards*, 49(1), 1952.
- Hutchinson, M. F. A stochastic estimator of the trace of the
in uence matrix for laplacian smoothing splines. *Com-
munications in Statistics-Simulation and Computation*
18(3):1059–1076, 1989.
- Kahn, H. Use of different monte carlo sampling techniques.
Rand Corporation, 1955.
- Katzfuss, M., Guinness, J., et al. A general framework for
vecchia approximations of gaussian processes. *Stat-
istical Science* 36(1):124–141, 2021.
- Lanczos, C. An iteration method for the solution of the
eigenvalue problem of linear differential and integral op-
erators. *Journal of Research of the National Bureau of
Standards* 45(4), 1950.
- Loper, J., Blei, D., Cunningham, J. P., and Paninski, L.
General linear-time inference for gaussian processes on
one dimensional. arXiv preprint arXiv:2003.05554, 2020.
- Luo, Y., Beatson, A., Norouzi, M., Zhu, J., Duvenaud, D.,
Adams, R. P., and Chen, R. T. Q. Sumo: Unbiased estima-
tion of log marginal probability for latent variable models.
In *International Conference on Learned Representations*
2020.
- Lyne, A.-M., Girolami, M., Atchadé, Y., Strathmann, H.,
Simpson, D., et al. On russian roulette estimates for
bayesian inference with doubly-intractable likelihoods.
Statistical science 30(4):443–467, 2015.
- Mutny, M. and Krause, A. Efficient high dimensional
bayesian optimization with additivity and quadrature
fourier features. *Advances in Neural Information Pro-
cessing Systems*, pp. 9005–9016, 2018.

- Nowozin, S. Debiasing evidence approximations: On importance-weighted autoencoders and jackknife variational inference. In International Conference on Learned Representation, 2018.
- Oktay, D., McGreivy, N., Aduol, J., Beatson, A., and Adams, R. P. Randomized automatic differentiation. arXiv preprint arXiv:2007.10412, 2020.
- Pleiss, G., Gardner, J. R., Weinberger, K. Q., and Wilson, A. G. Constant-time predictive distributions for gaussian processes. In International Conference on Machine Learning, 2018.
- Rahimi, A. and Recht, B. Random features for large-scale kernel machines. In Advances in Neural Information Processing Systems, pp. 1177–1184, 2008.
- Snelson, E. and Ghahramani, Z. Sparse Gaussian processes using pseudo-inputs. In Advances in Neural Information Processing Systems, 2006.
- Sutherland, D. J. and Schneider, J. G. On the error of random fourier features. Conference on Uncertainty in Artificial Intelligence, 2015.
- Titsias, M. K. Variational learning of inducing variables in sparse Gaussian processes. In International Conference on Artificial Intelligence and Statistics, pp. 567–574, 2009.
- Ubaru, S., Chen, J., and Saad, Y. Fast estimation of $\text{tr}(f(A))$ via stochastic lanczos quadrature. SIAM Journal on Matrix Analysis and Applications, 38(4):1075–1099, 2017.
- Wang, K., Pleiss, G., Gardner, J., Tyree, S., Weinberger, K. Q., and Wilson, A. G. Exact gaussian processes on a million data points. In Advances in Neural Information Processing Systems, pp. 14648–14659, 2019.
- Wilson, J., Borovitskiy, V., Terenin, A., Mostowsky, P., and Deisenroth, M. Efficiently sampling functions from gaussian process posteriors. In International Conference on Machine Learning, pp. 10292–10302, 2020.

A. Proof of Biases for CG and RFF

A.1. Proof of Theorem 1

Proof. To prove the bias of the CG log marginal likelihood terms, we rely on a connection between conjugate gradients, the Lanczos algorithm (Lanczos, 1950), and Gauss quadrature. In particular, we demonstrate that $\text{tr}(\log \mathbb{K}_{XX}^{-1})$ and $\text{tr}(\log \mathbb{K}_{XX})$ – CG's estimates of $\int \mathbb{K}_{XX}^{-1} y$ and $\int \mathbb{K}_{XX} j$ terms – are equivalent to Gauss quadrature approximations of two Riemann-Stieltjes integrals. We then use quadrature error analysis to prove the biases of these terms.

Expressing $\int \mathbb{K}_{XX}^{-1} y$ and $\int (\log \mathbb{K}_{XX}) z$ as Riemann-Stieltjes integrals. First, we note that $\int \mathbb{K}_{XX}^{-1} y$ and $\int (\log \mathbb{K}_{XX}) z$ (our stochastic trace estimate $\log \mathbb{K}_{XX} j$) can both be expressed by the quadratic form $w^T f(A)w$, where $f(A)$ denotes a matrix function of the matrix A , and w is a vector. Letting $A = P \Lambda^{-1} P^T$ be the eigendecomposition of A , we can write this quadratic form as

$$w^T f(A)w = w^T P f(\Lambda^{-1}) P^T w = \sum_{i=1}^N k w_i^2 f(\lambda_i^{-1}),$$

where λ_{\min} to λ_{\max} are the diagonal elements of (i.e. the eigenvalues) – ordered from smallest to largest, and w_i are the components of $w = kw$. The summation in the above equation can be expressed as a Riemann-Stieltjes integral:

$$\begin{aligned} \int f(\lambda) d\mu_A(\lambda) &:= \sum_{i=1}^N k w_i^2 f(\lambda_i^{-1}) \\ &= \int_{\lambda_{\min}}^{\lambda_{\max}} f(t) d\mu_A(t); \end{aligned} \quad (S1)$$

where the measure $\mu_A(t)$ is a piecewise constant function

$$\mu_A(t) = \begin{cases} 0 & \text{if } t < \lambda_{\min}; \\ \sum_{j=1}^i k w_j^2 & \text{if } \lambda_i \leq t < \lambda_{i+1}; \\ \sum_{j=1}^N k w_j^2 & \text{if } \lambda_{\max} \leq t. \end{cases} \quad (S2)$$

See (Golub & Meurant, 2009) for more details.

The Lanczos algorithm approximates these integrals with Gauss quadrature. The Lanczos algorithm (Lanczos, 1950), which is briefly described in Sec. 3.1, iteratively expresses a symmetric matrix A via the partial tridiagonalization $T_w^{(J)} = Q_w^{(J)T} A Q_w^{(J)}$. $Q_w^{(J)} \in \mathbb{R}^{N \times J}$ is an orthonormal matrix with $w = kw$ as its first column, and $T_w^{(J)}$ is a $J \times J$ tridiagonal matrix. Briefly, the columns of $Q_w^{(J)}$ matrices are computed by performing Gram-Schmidt orthogonalization on the Krylov subspace $[w; Aw; A^2 w; \dots; A^{J-1} w]$, and storing the orthogonalization constants in $T_w^{(J)}$.

The Lanczos algorithm is commonly used to estimate quadratic forms:

$$\begin{aligned} w^T f(A)w &= w^T Q_w^{(J)T} f(T_w^{(J)}) Q_w^{(J)} w \\ &= \sum_{i=1}^J k w_i^2 f(\tau_i^{(J)}); \end{aligned} \quad (S3)$$

where e_1 is the first unit vector. Note that Eq. (S3) holds because the columns of $Q_w^{(J)}$ are orthonormal.

There is a well-established connection between Eq. (S3) and numeric quadrature (e.g. Golub & Meurant, 2009). More specifically, Eq. (S3) is exactly equivalent to a term Gauss quadrature rule applied to the Riemann-Stieltjes integral in Eq. (S1). We can thus analyze Lanczos estimates of $w^T f(A)w$ using standard Gauss quadrature error analysis.

Equivalence between CG and the Lanczos algorithm. We will now show an equivalence between our estimates $u_J = \int \mathbb{K}_{XX}^{-1} y$ and $v_J = \int (\log \mathbb{K}_{XX}) z$ and Lanczos algorithm approximations. Note that we have already established $v_J = \int (\log \mathbb{K}_{XX}) z$ as a Lanczos algorithm approximation in Eq. (8):

$$z^T = \sum_{i=1}^J k z_i^2 \tau_i^{(J)} e_1;$$

For $u_J = \int \mathbb{K}_{XX}^{-1} y$, we exploit a connection between CG and Lanczos (Golub & Van Loan, 2012, Ch. 11.3):

$$\sum_{j=1}^J d_j = \sum_{j=1}^J k y_j^2 Q_y^{(J)T} (T_y^{(J)})^{-1} e_1; \quad (S4)$$

where $\sum_{j=1}^J d_j$ (see Eq. 4) is the J th CG approximation to $\int \mathbb{K}_{XX}^{-1} y$. Multiplying Eq. (S4) by y^T , we have

$$\begin{aligned} u_J &= y^T \sum_{j=1}^J d_j A \\ &= \sum_{j=1}^J k y_j^2 Q_y^{(J)T} (T_y^{(J)})^{-1} e_1 A \\ &= \sum_{j=1}^J k y_j^2 \tau_j^{(J)} e_1. \end{aligned}$$

Therefore, the CG approximations of $\int \mathbb{K}_{XX}^{-1} y$ and $\int (\log \mathbb{K}_{XX}) z$ are both Lanczos approximations. Putting all the pieces together, we have just shown that u_J and v_J are equivalent to J -term Gauss quadrature rules applied to the following integrals:

$$u_J = \int \mathbb{K}_{XX}^{-1} y = \sum_{j=1}^J k y_j^2 \int_{\lambda_{\min}}^{\lambda_{\max}} t^{-1} d\mu_{\mathbb{K}_{XX}}(t); \quad (S5)$$

$$z^T \log \mathbb{K}_{XX}^{-1} z = \sum_{j=1}^J k z_j^2 \int_{\lambda_{\min}}^{\lambda_{\max}} \log(t) d\mu_{\mathbb{K}_{XX}}(t); \quad (S6)$$

where the measure $\mu_{\mathbb{K}_{XX}}(t)$ is defined by Eq. (S2).

Applying Gauss quadrature error analysis to u_J and v_J . To analyze the bias of the CG estimates, we make use of standard error Gauss quadrature error analysis. $L_G^{(J)}[f]$ is the J -term Gaussian quadrature approximation of the Riemann-Stieltjes integral $\int [f]$ (see Eq. S1), the error can be exactly expressed as:

$$|f| - L_G^{(J)}[f] = (-1)^{2J} \frac{f^{(2J)}(\xi)}{(2J)!}; \quad (S7)$$

for some $\xi \in [\min; \max]$, where $f^{(2J)}$ is the $2J$ th derivative of f and ξ are some quantities that depend on the spectrum of the matrix K_{XX} (Golub & Meurant, 2009, Ch. 6).

Turning back to u_J , the corresponding function $f(t) = t^{-1}$ has even derivatives of the form $f^{(2J)} = (-1)^J (2J)! t^{-(2J+1)} > 0; 8J; 8t^{-2} (\min; \max)$. Replacing $[f]$ with the integral defined by Eq. (S5), we have that $u_J < 0$, which proves that CG underestimates $K_{XX}^{-1}y$. Similarly for v_J , the corresponding $g(t) = \log t$ has even derivatives of the form $g^{(2J)} = (-1)^J (2J-1)! t^{-2J} < 0; 8J; 8t^{-2} (\min; \max)$. Replacing $[f]$ with the integral defined by Eq. (S6), we have that $(\log K_{XX})z > v_J > 0$.

The convergence rates of u_J and v_J follow from Eq. 4.4 of Ubaru et al. (2017). Let $\kappa = (\rho_{\max}^{-1} + 1) / (\rho_{\min}^{-1} - 1)$, where ρ is the condition number of K_{XX} (i.e. the ratio of its maximum and minimum singular values). Then

$$\begin{aligned} |u_J| &\leq C \kappa^{-2J} \\ |v_J| &\leq C \kappa^{-2J} \end{aligned}$$

where C is a constant that depends on the extremal eigenvalues of K_{XX} . \square

A.2. Proof of Theorem 2

Proof. The inequalities are a result of applying Jensen's inequality to the inverse of a positive definite matrix (which is a convex function) and the log determinant of a positive definite matrix (which is a concave function). For example, for the $y^T K_{XX}^{-1}y$ term we have that

$$\begin{aligned} \mathbb{E}_{P^{(l)}} \left[y^T K_{XX}^{-1}y \right] &= y^T \mathbb{E}_{P^{(l)}} \left[K_{XX}^{-1} \right] y \\ &\geq y^T \mathbb{E}_{P^{(l)}} \left[K_{XX} \right]^{-1} y \\ &= y^T K_{XX}^{-1}y \end{aligned}$$

A similar procedure, but in the opposite direction applies to the $\log |K_{XX}|$.

To estimate the rate of decay of the biases, we rely on two key ideas. The first idea is to translate the central moments

of the kernel matrix being approximated by the random Fourier features into the central moments of these features. This strategy resembles the analysis in (Nowozin, 2018) which uses (Angelova, 2012, Theorem 1). The second idea is to use an approximation to two matrix functions. For the inverse function we use a Neumann series and for the logarithm we use a Taylor series. These series require that the eigenvalues of the approximation residuals $\|K_{XX} - K_J\|$ are less than 1. We will make this assumption as we know that for a large enough J , the kernel approximation will be close to K_{XX} . Below is a formal argument.

For a fixed l , we can write

$$\begin{aligned} y^T K_J^{-1}y &= y^T (K_{XX} - K_J)^{-1}y \\ &= y^T K_{XX}^{-1} \left(I + K_{XX}^{-1}K_J \right)^{-1} y \\ &= y^T K_{XX}^{-1} \left(I - K_{XX}^{-1}K_J + K_{XX}^{-1}K_J K_{XX}^{-1}K_J - \dots \right) y \end{aligned} \quad (S8)$$

this last form allow us to use a Neumann series to expand the inner inverse matrix. Hence, we have that

$$\begin{aligned} &= \sum_{t=0}^{\infty} (-1)^t y^T K_{XX}^{-1} K_J^t K_{XX}^{-1} y \end{aligned}$$

Combining this with Eq. (S8), we have that

$$\begin{aligned} \mathbb{E}_{P^{(l)}} \left[y^T K_J^{-1}y \right] &= y^T K_{XX}^{-1}y \\ &+ \sum_{t=2}^{\infty} (-1)^t y^T K_{XX}^{-1} K_J^t K_{XX}^{-1} y \end{aligned} \quad (S9)$$

where the first term of the series ($t=0$) is the data-term being approximated and the second term of the series cancels out ($= -1$). Following the analysis in (Nowozin, 2018) we translate the central moments of the random variable $K_{XX}^{-1}K_J K_{XX}^{-1}$ to the central moments of $K_{XX}^{-1}K_J K_{XX}^{-1}$ denoted as μ_i for $i \geq 2$. Since the term ($t=2$) will dominate, we will only focus on it (as the others would be of a higher order, in see (Angelova, 2012)). We have that

$$\mathbb{E}_{P^{(l)}} \left[y^T K_{XX}^{-1} K_J^2 K_{XX}^{-1} y \right] = \frac{2}{J} y^T K_{XX}^{-1} y \quad (S10)$$

Therefore, incorporating the previous result into Eq. (S9), allows us to conclude that

$$\mathbb{E}_{P^{(l)}} \left[y^T K_J^{-1}y \right] - y^T K_{XX}^{-1}y = O(1/J)$$

For the model complexity term $\log |K_{XX}^{-1}|$ we follow a similar procedure as before. For a fixed

$$\begin{aligned} \log |K_J| &= \log |K_{XX}^{-1} + R_J| \\ &= \log |K_{XX}^{-1}| + \log |I + K_{XX}^{-1} R_J K_{XX}^{-1}| \end{aligned}$$

Focusing on the last term we have that

$$\begin{aligned} \log |I + K_{XX}^{-1} R_J K_{XX}^{-1}| &= \text{tr} \log |I + K_{XX}^{-1} R_J K_{XX}^{-1}| \\ &= \sum_{t=1}^T \log |I + K_{XX}^{-1} R_J K_{XX}^{-1}|^t \end{aligned}$$

We can rewrite the term in the R.H.S as follows

$$\begin{aligned} \log |I + K_{XX}^{-1} R_J K_{XX}^{-1}| &= \sum_{t=1}^T \frac{1}{t} \text{tr} (K_{XX}^{-1} R_J K_{XX}^{-1})^t \end{aligned}$$

Again, following the analysis of (Nowozin, 2018), we can express the central moments of the random variable $K_{XX}^{-1} R_J K_{XX}^{-1}$ to the central moments of $K_{XX}^{-1} R_J K_{XX}^{-1}$ denoted as μ_i for $i \geq 2$. Also, as explained before, the dominant term will be $(i = 2)$, therefore we have that thus, using Eq. (S10) we have that

$$\begin{aligned} \text{tr} E_{P(\cdot)} \log |I + K_{XX}^{-1} R_J K_{XX}^{-1}| &= O(1/J) \end{aligned}$$

Therefore, we can conclude that

$$E_{P(\cdot)} \log |K_J| = O(1/J)$$

B.1. Unbiasedness of the RR Estimator

The RR estimator of Eq. (6) is unbiased, $E_{P(\cdot)}[\hat{\mu}_J] = \mu$.

Proof.

$$\begin{aligned} E_{P(\cdot)}[\hat{\mu}_J] &= E_{P(\cdot)} \left[\sum_{j=1}^J \frac{X_j^T}{P(J=j)} \right] \\ &= E_{P(\cdot)} \left[\sum_{j=1}^J \frac{X_j^T}{P(J=j)} I_{J=j} \right] \\ &= \sum_{j=1}^J \frac{X_j^T}{P(J=j)} E_{P(\cdot)}[I_{J=j}] \\ &= \sum_{j=1}^J \frac{X_j^T}{P(J=j)} P(J=j) I_{J=j} \\ &= \sum_{j=1}^J X_j^T P(J=j) = \mu \end{aligned}$$

□

B.2. Unbiasedness of the SS Estimator

The SS estimator of Eq. (7) is unbiased, $E_{P(\cdot)}[\hat{\mu}_J] = \mu$.

Proof.

$$\begin{aligned} E_{P(\cdot)}[\hat{\mu}_J] &= E_{P(\cdot)} \left[\sum_{j=1}^J \frac{X_j^T}{P(J=j)} I_{J=j} \right] \\ &= \sum_{j=1}^J \frac{X_j^T}{P(J=j)} E_{P(\cdot)}[I_{J=j}] \\ &= \sum_{j=1}^J \frac{X_j^T}{P(J=j)} P(J=j) I_{J=j} \\ &= \sum_{j=1}^J X_j^T P(J=j) = \mu \end{aligned}$$

□

B. Further Derivations of Randomized Truncation Estimators

Here we prove the unbiasedness of both Russian Roulette (RR) Eq. (6) and Single Sample (SS) Eq. (7) estimators. Recall from Sec. 2.3 that we wish to estimate the expensive

$$\mu = \frac{1}{H} \sum_{j=1}^H f_j; \quad H \geq 2, N[f_j] < \infty;$$

by computing just the first J terms, where $J \leq H$. H_g is randomly drawn from the truncation distribution $P(J = j)$. RR and SS estimators, denoted as $\hat{\mu}_J$ and $\hat{\mu}_J^{\text{SS}}$ respectively, offer two strategies for up-weighting the computed terms, which, as we will prove below, yield an unbiased estimator $E_{P(\cdot)}[\hat{\mu}_J] = \mu$.

B.3. Minimizing the Variance of the SS Estimator

Below we will derive the optimal distribution that minimizes the variance of our SS estimator. Note that for a given truncation distribution, we have that

$$\begin{aligned} V_J &= \sum_{j=1}^J \frac{X_j^T X_j}{P(J=j)^2} V_J(I_{J=j}) \\ &= \sum_{j=1}^J \frac{1}{P(J=j)} P(J=j) X_j^T X_j \end{aligned}$$

since $J = j$ is a Bernoulli random variable with probability $P(J = j)$, we can plug-in its variance and derive the second equality. Hence, to find the truncation distribution that minimizes the variance of the SS estimator we can solve the following constraint optimization problem.

$$\min_p \sum_{j=1}^H \frac{1}{p_j} \quad \text{s.t.} \quad \sum_{j=1}^H p_j = 1; \quad p_j \geq 0$$

where p_j is acting as a shorthand for $P(J = j)$. The Lagrangian of the problem is

$$L(p; \lambda) = \sum_{j=1}^H \frac{1}{p_j} + \lambda \left(\sum_{j=1}^H p_j - 1 \right)$$

where we can ignore the nonnegativity constraints as long as $p_j > 0$ (see solution below). Hence, the first order conditions dictate that

$$\frac{\partial L}{\partial p_j} = -\frac{1}{p_j^2} + \lambda = 0$$

therefore, if we take the ratio of the probabilities with respect to the first we get that $p_j^2 = \frac{1}{p_1^2}$. If we substitute this expression in the equality constraint we get that

$$P^2(J = j) = \frac{1}{H} \frac{1}{j} \quad (S11)$$

We emphasize that this result is a guide for practical choices of the truncation distribution. It is impractical to compute it as it will require evaluating all the p_j for $j = 1; \dots; H$. However, if we possess an estimate or a theoretical bound on rate of decay of each p_j then our unnormalized truncation distribution should also decay at this rate to minimize variance.

B.4. SS estimator as Importance Sampling

Here we derive the SS estimator by importance sampling the quantity $\frac{1}{H} \sum_{j=1}^H \frac{1}{j}$ with $P(J = J)$ as our proposal distribution.

$$\frac{1}{H} \sum_{j=1}^H \frac{1}{j} = \sum_{j=1}^H \frac{1}{j} P(J = j)$$

Next, re-write the summation above as an expectation over the discrete uniform distribution $U[1; H] = \frac{1}{H}; \dots; H$:

$$= \mathbb{E} \left[\frac{1}{J} \right]$$

We now introduce an alternative proposal distribution $P(J = J)$:

$$= \mathbb{E} \left[\frac{1}{J} \frac{P(J = J)}{P(J = J)} \right]$$

Approximating the final expectation using a single Monte Carlo sample results in the SS estimator.

C. Estimating the Marginal Log Likelihood from RR-CG

Recall from Sec. 4.1 that we use the Russian Roulette estimator in conjunction with conjugate gradients to compute an unbiased (stochastic) gradient of the log marginal likelihood. In this section, we briefly describe how to obtain an unbiased estimate of the log marginal likelihood itself.

The $\log \kappa_{XX}^{-1} y$ term can be estimated directly from the $\log \kappa_{XX}^{-1} y$ solve in Eq. (16). The $\log \kappa_{XX}^{-1} y$ term is less straightforward, as the estimate from the CG byproducts (Eq. 8) isn't readily expressed as a summation. Instead, we apply the Russian Roulette estimator to the following telescoping series:

$$\log \kappa_{XX}^{-1} y = \sum_{j=1}^H \left(\log T_z^{(N)} e_1 - \log T_z^{(j)} e_1 \right) \quad (S12)$$

where $j = \log T_z^{(j)} e_1 - \log T_z^{(j-1)} e_1$. We can therefore apply the Russian Roulette estimator with truncation $P(J)$ to Eq. (S12). This telescoping series can be expensive to compute. Since it is unnecessary for gradient-based optimization, we introduce it only as a tool to analyze the RR-CG log marginal likelihood.

D. Optimal Truncation Distributions

D.1. Proof of Theorem 3

Proof. For Thm. 3 we have to combine the results of Thm. 1 with the optimal distribution for RR that is derived in Beatson & Adams (2019). We will only focus on $\log \kappa_{XX}^{-1} y$ which that approximates $\log \kappa_{XX}^{-1} y$ since the procedure is analogous for u_j which approximates $\kappa_{XX}^{-1} y$. Beatson & Adams (2019, Theorem 5.4) state that the optimal RR truncation distribution that maximizes the ROE is

$$P(J = j) = \frac{E[j v_j j^2]}{c(j) - c(j-1)}$$

where $c(j)$ is the computational cost of evaluating $\log \kappa_{XX}^{-1} y$ which is the j -th term being approximated through Russian

Roulette. In this case, the $\mathbf{P}(J = j)$ are nonrandom and single-valued and the cost per iteration is constant with respect to j . Hence, we can conclude that

$$\mathbf{P}(J = j) / v_j = O(C^{-2j})$$

which implies that $\mathbf{P}(J = j) / O(C^{-2j})$ since the derivative of an exponential function is of the same form we have that $\mathbf{P}(J = j) / O(C^{-2j})$. \square

D.2. Proof of Theorem 4

The strategy for Thm. 4 is to relate the kernel approximation using $J+1$ random Fourier features to the kernel approximation using J features for the two components of the log marginal likelihood: the data-term $\log p(\mathbf{y} | \mathbf{X})$ and the model complexity term $\log p(\mathbf{X} | \mathbf{y})$. For $\log p(\mathbf{y} | \mathbf{X})$ we create this connection through the Sherman-Morrison formula and for the $\log p(\mathbf{X} | \mathbf{y})$ we use the matrix determinant lemma. Throughout these proofs, we will require some results regarding positive definite matrices that we add as remarks below. Before stating those remarks, we will first introduce some auxiliary notation.

$$\mathbf{K}_J := \frac{1}{J} \sum_{j=1}^J \mathbf{x}_j \mathbf{x}_j^\top$$

$$(J+1) \mathbf{K}_{J+1} = J \mathbf{K}_J + \mathbf{x}_{J+1} \mathbf{x}_{J+1}^\top$$

Note the difference between the \mathbf{K}_J term used throughout the paper which includes the \mathbf{x}_j against the \mathbf{K}_J term which only includes the RFF features. Moreover, to reduce clutter we define the following matrix

$$\mathbf{W} := \frac{J}{J+1} \mathbf{K}_J + \mathbf{I}$$

whose use will become evident throughout the proof.

Remark 1. Given a positive definite matrix \mathbf{A} we have that for any $\alpha \in (0, 1)$ and any $\beta > 0$

$$k \alpha \mathbf{A} + \beta \mathbf{I} \succeq k \mathbf{A} + \beta \mathbf{I}$$

If we express $\mathbf{A} = \mathbf{V} \mathbf{D} \mathbf{V}^\top$, where \mathbf{V} is an orthogonal matrix and \mathbf{D} is a diagonal matrix containing the positive eigenvalues of \mathbf{A} , then we have that $\alpha \mathbf{A} + \beta \mathbf{I} = \mathbf{V} (\alpha \mathbf{D} + \beta \mathbf{I}) \mathbf{V}^\top$. This implies that the eigenvalues of $\alpha \mathbf{A} + \beta \mathbf{I}$ are larger than those of $\mathbf{A} + \beta \mathbf{I}$. An immediate consequence of this remark is that $\mathbf{W} \succeq \frac{J}{J+1} \mathbf{K}_J + \mathbf{I}$ or that $\mathbf{W}^{-1} \preceq \frac{J+1}{J} (\mathbf{K}_J + \mathbf{I})^{-1}$ where \mathbf{W} contains the positive eigenvalues of \mathbf{K}_J .

Remark 2. Given a positive definite matrix \mathbf{A} we have that for any $\alpha \in (0, 1)$, any $\beta > 0$ and any vector \mathbf{x}

$$\mathbf{x}^\top (\alpha \mathbf{A} + \beta \mathbf{I})^{-1} \mathbf{x} \leq \mathbf{x}^\top (\mathbf{A} + \beta \mathbf{I})^{-1} \mathbf{x} + \beta^{-1} \mathbf{x}^\top \mathbf{x}$$

This remark follows by using the same diagonal decomposition as the one used in the previous remark. Hence, by noting that the eigenvalues of $(\alpha \mathbf{A} + \beta \mathbf{I})^{-1}$ are larger than those of $(\mathbf{A} + \beta \mathbf{I})^{-1}$ and this last matrix has larger eigenvalues than $(\mathbf{A} + \beta \mathbf{I})^{-1}$ then the result follows.

Proof of Theorem 4. We start by showing the rate of decay for the $\mathbf{P}(J = j)$ involving the data-term. Note that we can express

$$\begin{aligned} \mathbf{y}^\top (\mathbf{K}_{J+1} + \mathbf{I})^{-1} \mathbf{y} &= \mathbf{y}^\top \left(\frac{J}{J+1} \mathbf{K}_J + \mathbf{I} + \frac{\mathbf{x}_{J+1} \mathbf{x}_{J+1}^\top}{J+1} \right)^{-1} \mathbf{y} \\ &= \mathbf{y}^\top \left(\mathbf{W} + \frac{\mathbf{x}_{J+1} \mathbf{x}_{J+1}^\top}{J+1} \right)^{-1} \mathbf{y} \end{aligned} \quad (\text{S13})$$

Applying Sherman-Morrison formula to the inverse of the R.H.S results in

$$\mathbf{W}^{-1} \frac{\mathbf{x}_{J+1} \mathbf{W}^{-1} \mathbf{x}_{J+1}}{(J+1) + \mathbf{x}_{J+1} \mathbf{W}^{-1} \mathbf{x}_{J+1}}$$

where the symmetry of \mathbf{W} allows us to express the numerator as above. Substituting the previous result into Eq. (S13) and rearranging terms allows us to conclude that

$$\begin{aligned} \mathbf{y}^\top (\mathbf{K}_{J+1} + \mathbf{I})^{-1} \mathbf{y} &= \mathbf{y}^\top (\mathbf{W} + \frac{\mathbf{x}_{J+1} \mathbf{x}_{J+1}^\top}{J+1})^{-1} \mathbf{y} \\ &= \frac{\mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y}}{(J+1) + \mathbf{x}_{J+1} \mathbf{W}^{-1} \mathbf{x}_{J+1}} \\ &= \frac{\mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y}}{J+1} \frac{1}{1 + \frac{\mathbf{x}_{J+1} \mathbf{W}^{-1} \mathbf{x}_{J+1}}{J+1}} \\ &= \frac{\mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y}}{J+1} \frac{1}{1 + \frac{\mathbf{x}_{J+1} \mathbf{W}^{-1} \mathbf{x}_{J+1}}{J+1}} \end{aligned} \quad (\text{S14})$$

where the first inequality follows from Remark 2, the second inequality occurs since $\mathbf{x}_{J+1} \mathbf{W}^{-1} \mathbf{x}_{J+1} > 0$, the third inequality results from applying Cauchy-Schwarz and the fourth stems from Remark 1. Finally, taking expectations in Eq. (S14) we can conclude that

$$\begin{aligned} \mathbb{E}_{P^{(j)}} \mathbf{y}^\top (\mathbf{K}_{J+1} + \mathbf{I})^{-1} \mathbf{y} &= \mathbb{E}_{P^{(j)}} \mathbf{y}^\top (\mathbf{W} + \frac{\mathbf{x}_{J+1} \mathbf{x}_{J+1}^\top}{J+1})^{-1} \mathbf{y} \\ &= O(1/J) \end{aligned}$$

We now move into the rate of decay of the Δ_J terms involving the model complexity terms. Note that we can express

$$\mathbf{K}_{J+1} + \sigma^2 \mathbf{I} = \frac{J}{J+1} \mathbf{K}_J + \sigma^2 \mathbf{I} + \frac{J+1}{J+1} \frac{\mathbf{W}_{J+1}}{J+1}$$

then by using the matrix determinant lemma we have that

$$\begin{aligned} \log \det(\mathbf{K}_{J+1} + \sigma^2 \mathbf{I}) &= \log \det \left(\mathbf{I} + \frac{1}{J+1} \mathbf{W}_{J+1} \mathbf{W}_{J+1}^{-1} \right) + \log \det(\mathbf{K}_J + \sigma^2 \mathbf{I}) \\ &\leq \log \det \left(\mathbf{I} + \frac{1}{J+1} \mathbf{W}_{J+1} \mathbf{W}_{J+1}^{-1} \right) + \log \det(\mathbf{K}_J + \sigma^2 \mathbf{I}) \\ &\leq \log \det \left(\mathbf{I} + \frac{\sigma^2}{J+1} \mathbf{W}_{J+1} \mathbf{W}_{J+1}^{-1} \right) + \log \det(\mathbf{K}_J + \sigma^2 \mathbf{I}) \end{aligned} \quad (\text{S15})$$

where the first inequality follows from Remark 1 and from noting that the determinant is equivalent to the product of the eigenvalues of the matrix. The second inequality follows from Remark 2. Now, we will take the logarithms and the expectations of Eq. (S15). We will first focus on the first term in the RHS. By using a Taylor expansion of the logarithm up to a second term we have that

$$\begin{aligned} \mathbb{E}_{P^{(l)}} \log \det \left(\mathbf{I} + \frac{\sigma^2}{J+1} \mathbf{W}_{J+1} \mathbf{W}_{J+1}^{-1} \right) &= \frac{\sigma^2 \mathbb{E}_{P^{(l)}} \text{tr}(\mathbf{W}_{J+1} \mathbf{W}_{J+1}^{-1})}{J+1} + \mathcal{O}(1/J^2) \\ &= \frac{\sigma^2 \mathbb{E}_{P^{(l)}} \text{tr}(\mathbf{W}_{J+1} \mathbf{W}_{J+1}^{-1})}{J+1} + \mathcal{O}(1/J^2) \end{aligned}$$

therefore, substituting the previous result into Eq. (S15)

$$\begin{aligned} \mathbb{E}_{P^{(l)}} \log \det(\mathbf{K}_{J+1} + \sigma^2 \mathbf{I}) - \mathbb{E}_{P^{(l)}} \log \det(\mathbf{K}_J + \sigma^2 \mathbf{I}) &= \frac{\sigma^2 \mathbb{E}_{P^{(l)}} \text{tr}(\mathbf{W}_{J+1} \mathbf{W}_{J+1}^{-1})}{J+1} + \mathcal{O}(1/J^2) \\ &= \mathcal{O}(1/J) \end{aligned}$$

which concludes the analysis of the rate of decay of the log determinant terms. Thus, since each rate of decay is $\mathcal{O}(1/J)$, then following Eq. (S11) we have that the truncation distribution that minimizes the variance is

$$P(J) \propto \Delta_J = \mathcal{O}(1/J).$$

□

E. Experiment Details.

We provide the experiment details for the predictive performance experiments in Sec. 5.

Experiment setup. To optimize the hyperparameters of the CG, RR-CG and Cholesky models, we use an Adam optimizer with learning rate = 0.01 and a MultiStepLR scheduler dropping the learning rate by a factor of 10 at the 50%, 70% and 90% of the optimization iterations. We run the optimization for 1500, 800 and 300 iterations on small (PoleTele, Elevators and Bike), medium (Kin40K, Protein, KEGG and KEGGU) and large (3DRoad) datasets, respectively. The number of iterations for CG and the expected number of iterations for RR-CG are both set to 100. The latter is achieved by using the truncation distribution from Eq. (17) with $\lambda = 0.05$ and $J_{\min} = 80$.

For CG and RR-CG, we use the rank-5 pivoted Cholesky preconditioner of Gardner et al. (2018). To reduce the number of optimization steps needed for the 3DRoad dataset, we initialize the hyperparameters to those found with the (biased) sgGP method. During evaluation, we compute the predictive mean using 1,000 iterations of CG. Predictive variances are estimated using the rank-100 Lanczos approximation of Pleiss et al. (2018).

For RFF we use 1,000 random Fourier features across all experiments. In terms of optimization, we use an Adam optimizer with learning rate = 0.005 for KEGG, 0.001 for KEGGU and 0.01 for the remaining datasets. We also use a MultiStep scheduler that activates at 85%, 90%, 95% of the optimization iterations with a decay rate of 0.5 and also take a total of 500 optimization iterations on all the datasets.

In all the SVGP models we use 1,024 inducing points, with a full-rank multivariate Gaussian variational distribution. Hyperparameters and variational parameters are jointly optimized for 300 epochs using minibatches of size 1,024. As with the other baselines, we use the an Adam optimizer with a learning rate of 0.01, dropping the learning rate by a factor of 10 after 50%, 70% and 90% of the optimization iterations.

For sgGP, we train using minibatches of 16 data points. As suggested by Chen et al. (2020), the minibatches are constructed by sampling one training data point and selecting its 15 nearest neighbors. To accelerate optimization, we accumulate the gradients of 1,024 minibatches before performing an optimization step (these 1,024 minibatch updates can be performed in parallel, enabling GPU acceleration). We optimize the models for 300 epochs, using the same learning rate and scheduler as with SVGP. During evaluation, we use the same procedure as for CC/RR-CG (1,000 iterations of CG, rank-100 Lanczos variance estimates).

For POE, we divide the training dataset into disjoint subsets, each with $M = 1024$ data points. If N is not divisible by 1024, we pad the final subset with elements from the first subset. We train independent GP with independent hyperparameters on each of the subsets using the same optimization

Dataset	n	RMSE						
		Cholesky	POE	RFF	SVGP	sgGP	CG	RR-CG
PolTele	9.6K	.112 ± .002	.174 ± .006	.106 ± .000	.150 ± .000	.128 ± .001	.119 ± .002	.112 ± .002
Elevators	10.6K	.360 ± .006	.379 ± .005	.365 ± .001	.376 ± .006	.362 ± .006	.360 ± .006	.360 ± .006
Bike	11.1K	.035 ± .003	.071 ± .002	.063 ± .001	.045 ± .002	1.044 ± .334	.040 ± .005	.035 ± .002
Kin40K	25.6K	—	.248 ± .001	.074 ± .000	.152 ± .001	.081 ± .000	.093 ± .000	.091 ± .001
Protein	25.6K	—	.715 ± .007	.547 ± .001	.664 ± .008	.562 ± .005	.541 ± .008	.541 ± .008
KEGG	31.2K	—	.097 ± .004	.101 ± .001	.088 ± .002	.089 ± .002	.195 ± .064	.087 ± .003
KEGGU	40.7K	—	.125 ± .001	.129 ± .001	.122 ± .001	.123 ± .001	.120 ± .000	.120 ± .000
3DRoad	278K	—	.675 ± .001	.348 ± .001	.439 ± .002	.285 ± .003	.202 ± .003	.114 ± .013

Dataset	n	NLL						
		Cholesky	POE	RFF	SVGP	sgGP	CG	RR-CG
PolTele	9.6K	−.464 ± .006	−.252 ± .010	−.159 ± .005	−.442 ± .004	−.480 ± .004	−.354 ± .003	−.458 ± .006
Elevators	10.6K	.425 ± .013	.469 ± .016	.780 ± .006	.442 ± .015	.426 ± .015	.429 ± .012	.425 ± .013
Bike	11.1K	−.984 ± .018	−.799 ± .010	.099 ± .030	−1.514 ± .024	85.715 ± 49.088	−.971 ± .008	−.982 ± .018
Kin40K	25.6K	—	.464 ± .002	1.407 ± .004	−.410 ± .003	.427 ± .001	.468 ± .003	.449 ± .013
Protein	25.6K	—	1.105 ± .008	1.163 ± .003	1.014 ± .012	.951 ± .004	.934 ± .006	.934 ± .006
KEGG	31.2K	—	−.874 ± .011	6.311 ± 2.323	−1.022 ± .023	−.981 ± .033	.415 ± .653	−.884 ± .009
KEGGU	40.7K	—	−.636 ± .002	4.000 ± .303	−.685 ± .005	−.668 ± .005	−.637 ± .006	−.650 ± .004
3DRoad	278K	—	1.031 ± .002	1.317 ± .004	.601 ± .004	.831 ± .000	.613 ± .010	.776 ± .030

Table S1. Root-mean-square-error (RMSE) and negative log-likelihood (NLL) of exact GPs using CG, RRCG and other baselines on UCI regression datasets using a constant prior mean and a RBF kernel with independent lengthscale for each dimension. All trials were averaged over 3 trials with different splits. N and d are the size and dimensionality of the training dataset, respectively.

Dataset	n	Training time (m)						
		Cholesky	POE	RFF	SVGP	sgGP	CG	RR-CG
PolTele	9.6K	22.417 ± .035	1.167 ± .006	.464 ± .002	3.862 ± .018	.629 ± .006	12.793 ± .111	14.968 ± .258
Elevators	10.6K	30.617 ± .016	1.051 ± .008	.503 ± .002	4.236 ± .023	.696 ± .001	14.443 ± .080	16.519 ± .152
Bike	11.1K	34.991 ± .016	1.364 ± .012	.476 ± .014	4.261 ± .023	.691 ± .002	11.634 ± .131	13.669 ± .080
Kin40K	25.6K	—	.930 ± .006	.772 ± .008	9.597 ± .043	1.559 ± .010	11.345 ± .073	12.823 ± .114
Protein	25.6K	—	.851 ± .003	.716 ± .008	11.115 ± .033	1.867 ± .007	10.507 ± .044	12.142 ± .052
KEGG	31.2K	—	1.135 ± .003	.682 ± .002	11.881 ± .058	1.786 ± .009	22.780 ± .021	25.390 ± .089
KEGGU	40.7K	—	1.140 ± .005	.835 ± .001	15.281 ± .070	2.572 ± .029	34.396 ± .026	37.825 ± .064
3DRoad	278K	—	2.176 ± .022	6.089 ± .031	104.164 ± .294	22.615 ± .147	145.396 ± 1.373	158.657 ± .554

Table S2. Total training time (in minutes) of exact GPs using CG, RRCG and other baselines on UCI regression datasets (see the number of optimization iterations for each method in experiment setup). All trials were averaged over 3 trials with different splits.

procedure as Cholesky models. Following Deisenroth & Ng (2015, Eqs. 11 and 12), the posterior distribution for a test input \mathbf{x} is Gaussian with mean $\mu_{\text{POE}}(\mathbf{x})$ and variance $\sigma_{\text{POE}}^2(\mathbf{x})$:

$$\mu_{\text{POE}}(\mathbf{x}) = \frac{\sigma_{\text{POE}}^2(\mathbf{x})}{K} \sum_{k=1}^K \sigma_k^2(\mathbf{x}) \mu_k(\mathbf{x}),$$

$$\sigma_{\text{POE}}^2(\mathbf{x}) = \frac{1}{K} \sum_{k=1}^K \sigma_k^2(\mathbf{x}),$$

where K is the number of independent GP experts and $\mu_k(\cdot)$ and $\sigma_k^2(\cdot)$ are the posterior mean and variance of each expert model.

Full tables. In Table S1, we report the RMSE and NLL numbers which are used to plot Fig. 6 in the paper, and in Table S2 we report the corresponding training time.

F. Additional Experiments

F.1. Predictive Performance with Different CG iterations

Here we include an additional experiment to show that early truncated CG can be detrimental to GP learning while RR-CG remains robust to the expected truncation number.

We conduct GP learning with the CG, RR-CG and Cholesky methods in the Elevators dataset. We vary the number of (expected) iterations for CG and RR-CG from 20 to 100 and plot the corresponding RMSE / NLL in Fig. S1. From the figure, we conclude that: 1) early truncation of the CG algorithm impedes GP optimization and leads to poor predictions. This problem lessens as we increase the number of CG iterations from 20 to 100. 2) The RR-CG model is robust to the expected truncation number, as it keeps comparable RMSE and NLL values to the Cholesky model under different expected truncation numbers. This experiment can

