

## 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  $u_J$  and  $v_J$  – CG’s estimates for  $\mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$  and  $\log |\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}|$  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  $\mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$  and  $\mathbf{z}^\top (\log \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}) \mathbf{z}$  as Riemann-Stieltjes integrals.** First, we note that  $\mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$  and  $\mathbf{z}^\top (\log \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}) \mathbf{z}$  (our stochastic trace estimate of  $\log |\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}|$ ) can both be expressed by the quadratic form  $\mathbf{w}^\top f(\mathbf{A}) \mathbf{w}$ , where  $f(\mathbf{A})$  denotes a matrix function of the matrix  $\mathbf{A}$ , and  $\mathbf{w}$  is a vector. Letting  $\mathbf{A} = \mathbf{P} \mathbf{\Lambda} \mathbf{P}^\top$  be the eigendecomposition of  $\mathbf{A}$ , we can write this quadratic form as

$$\mathbf{w}^\top f(\mathbf{A}) \mathbf{w} = \mathbf{w}^\top \mathbf{P} f(\mathbf{\Lambda}) \mathbf{P}^\top \mathbf{w} = \|\mathbf{w}\|^2 \sum_{i=1}^N f(\lambda_i) \mu_i^2,$$

where  $\lambda_{\min} \leq \lambda_i \leq \lambda_{\max}$  are the diagonal elements of  $\mathbf{\Lambda}$  (i.e. the eigenvalues) – ordered from smallest to largest, and  $\mu_i$  are the components of  $\mathbf{P}^\top \mathbf{w} / \|\mathbf{w}\|$ . The summation in the above equation can be expressed as a Riemann-Stieltjes integral:

$$\begin{aligned} I[f] &:= \|\mathbf{w}\|^2 \sum_{i=1}^N f(\lambda_i) \mu_i^2 \\ &= \|\mathbf{w}\|^2 \int_{\lambda_{\min}}^{\lambda_{\max}} f(t) d\mu_{\mathbf{A}}(t), \end{aligned} \quad (\text{S1})$$

where the measure  $\mu_{\mathbf{A}}(t)$  is a piecewise constant function

$$\mu_{\mathbf{A}}(t) = \begin{cases} 0, & \text{if } t < \lambda_{\min}, \\ \sum_{j=1}^i \mu_j^2, & \text{if } \lambda_i \leq t < \lambda_{i-1} \\ \sum_{j=1}^N \mu_j^2, & \text{if } \lambda_{\max} \leq t. \end{cases} \quad (\text{S2})$$

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  $\mathbf{A}$  via the partial tridiagonalization  $\mathbf{T}_{\mathbf{w}}^{(J)} = \mathbf{Q}_{\mathbf{w}}^{(J)\top} \mathbf{A} \mathbf{Q}_{\mathbf{w}}^{(J)}$ .  $\mathbf{Q}_{\mathbf{w}}^{(J)} \in \mathbb{R}^{N \times J}$  is an orthonormal matrix with  $\mathbf{w} / \|\mathbf{w}\|$  as its first column, and  $\mathbf{T}_{\mathbf{w}}^{(J)}$  is a  $J \times J$  tridiagonal matrix. Briefly, the columns of  $\mathbf{Q}^{(J)}$  matrices are computed by performing Gram-Schmidt orthogonalization on the Krylov subspace  $[\mathbf{w}, \mathbf{A}\mathbf{w}, \mathbf{A}^2\mathbf{w}, \dots, \mathbf{A}^{J-1}\mathbf{w}]$ , and storing the orthogonalization constants in  $\mathbf{T}_{\mathbf{w}}^{(J)}$ .

The Lanczos algorithm is commonly used to estimate quadratic forms:

$$\begin{aligned} \mathbf{w}^\top f(\mathbf{A}) \mathbf{w} &\approx \mathbf{w}^\top \mathbf{Q}_{\mathbf{w}}^{(J)} \mathbf{f}(\mathbf{T}_{\mathbf{w}}^{(J)}) \mathbf{Q}_{\mathbf{w}}^{(J)\top} \mathbf{w} \\ &= \|\mathbf{w}\|^2 \mathbf{e}_1^\top f(\mathbf{T}_{\mathbf{w}}^{(J)}) \mathbf{e}_1, \end{aligned} \quad (\text{S3})$$

where  $\mathbf{e}_1$  is the first unit vector. Note that Eq. (S3) holds because the columns of  $\mathbf{Q}_{\mathbf{w}}^{(J)}$  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  $J$ -term Gauss quadrature rule applied to the Riemann-Stieltjes integral in Eq. (S1). We can thus analyze Lanczos estimates of  $\mathbf{w}^\top f(\mathbf{A}) \mathbf{w}$  using standard Gauss quadrature error analysis.

### Equivalence between CG and the Lanczos algorithm.

We will now show an equivalence between our estimates  $u_J \approx \mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$  and  $v_J \approx \mathbf{z}^\top (\log \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}) \mathbf{z}$  and Lanczos algorithm approximations. Note that we have already established  $v_J \approx \mathbf{z}^\top (\log \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}) \mathbf{z}$  as a Lanczos algorithm approximation in Eq. (8):

$$\mathbf{z}^\top = \|\mathbf{z}\|^2 \mathbf{e}_1^\top (\log \mathbf{T}_{\mathbf{z}}^{(J)}) \mathbf{e}_1.$$

For  $u_J \approx \mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$ , we exploit a connection between CG and Lanczos (Golub & Van Loan, 2012, Ch. 11.3):

$$\sum_{j=1}^J \gamma_j \mathbf{d}_j = \|\mathbf{y}\| \mathbf{Q}_{\mathbf{y}}^{(J)} (\mathbf{T}_{\mathbf{y}}^{(J)})^{-1} \mathbf{e}_1, \quad (\text{S4})$$

where  $\sum_{j=1}^J \gamma_j \mathbf{d}_j$  (see Eq. 4) is the  $J^{\text{th}}$  CG approximation to  $\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$ . Multiplying Eq. (S4) by  $\mathbf{y}^\top$ , we have

$$\begin{aligned} u_J &= \mathbf{y}^\top \left( \sum_{j=1}^J \gamma_j \mathbf{d}_j \right) \\ &= \|\mathbf{y}\| \mathbf{y}^\top \mathbf{Q}_{\mathbf{y}}^{(J)} (\mathbf{T}_{\mathbf{y}}^{(J)})^{-1} \mathbf{e}_1 \\ &= \|\mathbf{y}\|^2 \mathbf{e}_1^\top (\mathbf{T}_{\mathbf{y}}^{(J)})^{-1} \mathbf{e}_1. \end{aligned}$$

Therefore, the CG approximations of  $\mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$  and  $\mathbf{z}^\top (\log \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}) \mathbf{z}$  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:

$$\mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y} = \|\mathbf{y}\|^2 \int_{\lambda_{\min}}^{\lambda_{\max}} t^{-1} d\mu_{\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}}(t), \quad (\text{S5})$$

$$\mathbf{z}^\top (\log \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}) \mathbf{z} = \|\mathbf{z}\|^2 \int_{\lambda_{\min}}^{\lambda_{\max}} \log(t) d\mu_{\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}}(t), \quad (\text{S6})$$

where the measure  $\mu_{\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}}(t)$  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. If  $L_G^{(J)}[f]$  is the  $J$ -term Gaussian quadrature approximation of the Riemann-Stieltjes integral  $I[f]$  (see Eq. S1), the error can be exactly expressed as:

$$I[f] - L_G^{(J)}[f] = (\gamma_1 \cdots \gamma_{J-1})^2 \frac{f^{(2J)}(\eta)}{(2J)}, \quad (\text{S7})$$

for some  $\eta \in [\lambda_{\min}, \lambda_{\max}]$ , where  $f^{(2J)}$  is the  $2J$ th derivative of  $f$  and  $\{\gamma_i\}$  are some quantities that depend on the spectrum of the matrix  $\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}$  (Golub & Meurant, 2009, Ch. 6).

Turning back to  $u_J$ , the corresponding function  $f(t) = t^{-1}$  has even derivatives of the form  $f^{(2J)} = (2J)! t^{-(2J+1)} > 0, \forall J, \forall t \in (\lambda_{\min}, \lambda_{\max})$ . Replacing  $I[f]$  with the integral defined by Eq. (S5), we have that  $\mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y} - u_J \geq 0$ , which proves that CG underestimates  $\mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$ . Similarly for  $v_J$ , the corresponding  $f(t) = \log t$  has even derivatives of the form  $f^{(2J)} = -(2J-1)! t^{-2J} < 0, \forall J, \forall t \in (\lambda_{\min}, \lambda_{\max})$ . Replacing  $I[f]$  with the integral defined by Eq. (S6), we have that  $\mathbf{z}^\top (\log \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}) \mathbf{z} - v_J \leq 0$ .

**The convergence rates** of  $u_J$  and  $v_J$  follow from Eq. 4.4 of Ubaru et al. (2017). Let  $\rho = (\sqrt{\kappa} + 1)/(\sqrt{\kappa} - 1)$ , where  $\kappa$  is the condition number of  $\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}$  (i.e. the ratio of its maximum and minimum singular values). Then

$$\begin{aligned} \|\mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y} - u_J\| &\leq C \rho^{-2J}, \\ \|\mathbf{z}^\top (\log \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}) \mathbf{z} - v_J\| &\leq C \rho^{-2J}, \end{aligned}$$

where  $C$  is a constant that depends on the extremal eigenvalues of  $\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}$ .  $\square$

## 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  $\mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$  term we have that

$$\begin{aligned} \mathbb{E}_{\mathbb{P}(\boldsymbol{\omega})} \left[ \mathbf{y}^\top \widehat{\mathbf{K}}_J^{-1} \mathbf{y} \right] &= \mathbf{y}^\top \mathbb{E}_{\mathbb{P}(\boldsymbol{\omega})} \left[ \widehat{\mathbf{K}}_J^{-1} \right] \mathbf{y} \\ &\geq \mathbf{y}^\top \left( \mathbb{E}_{\mathbb{P}(\boldsymbol{\omega})} \left[ \widehat{\mathbf{K}}_J \right] \right)^{-1} \mathbf{y} \\ &= \mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y} \end{aligned}$$

A similar procedure, but in the opposite direction applies to the  $\log |\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}|$ .

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 residual  $(\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}} - \widetilde{\mathbf{K}}_J)$  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  $\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}$ . Below is a formal argument.

For a fixed  $\boldsymbol{\omega}$  we can write

$$\begin{aligned} \mathbf{y}^\top \widehat{\mathbf{K}}_J^{-1} \mathbf{y} &= \mathbf{y}^\top \left( \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}} - \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}} + \widetilde{\mathbf{K}}_J \right)^{-1} \mathbf{y} \\ &= \mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \left( \mathbf{I} - \mathbf{I} + \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \right)^{-1} \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \mathbf{y} \\ &= \mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \left( \mathbf{I} - \left( \mathbf{I} - \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \right) \right)^{-1} \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \mathbf{y} \end{aligned} \quad (\text{S8})$$

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

$$\begin{aligned} &\left( \mathbf{I} - \left( \mathbf{I} - \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \right) \right)^{-1} \\ &= \sum_{t=0}^{\infty} \left( \mathbf{I} - \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \right)^t \end{aligned}$$

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

$$\begin{aligned} \mathbb{E}_{\mathbb{P}(\boldsymbol{\omega})} \left[ \mathbf{y}^\top \widehat{\mathbf{K}}_J^{-1} \mathbf{y} \right] &= \mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y} \\ &+ \sum_{t=2}^{\infty} \mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \mathbb{E}_{\mathbb{P}(\boldsymbol{\omega})} \left[ \left( \mathbf{I} - \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \right)^t \right] \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \mathbf{y} \end{aligned} \quad (\text{S9})$$

where the first term of the series ( $t = 0$ ) is the data-fit term being approximated and the second term of the series cancels out ( $t = 1$ ). Following the analysis in (Nowozin, 2018) we translate the central moments of the random variable  $\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2}$  to the central moments of  $\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \left( \phi(\boldsymbol{\omega}) \phi(\boldsymbol{\omega})^\top + \sigma^2 \mathbf{I} \right) \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2}$  denoted as  $\mu_i$  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  $J$ , see (Angelova, 2012)). We have that

$$\mathbb{E}_{\mathbb{P}(\boldsymbol{\omega})} \left[ \left( \mathbf{I} - \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \right)^2 \right] = \frac{\mu_2}{J} \quad (\text{S10})$$

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

$$\mathbb{E}_{\mathbb{P}(\boldsymbol{\omega})} \left[ \mathbf{y}^\top \widehat{\mathbf{K}}_J^{-1} \mathbf{y} \right] - \mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y} = \mathcal{O}(1/J)$$

For the model complexity term  $\log |\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}|$  we follow a similar procedure as before. For a fixed  $\omega$

$$\begin{aligned} \log |\widetilde{\mathbf{K}}_J| &= \log |\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}} - \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}} + \widetilde{\mathbf{K}}_J| \\ &= \log |\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}| - \log \left| \mathbf{I} - \left( \mathbf{I} - \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \right) \right| \end{aligned}$$

Focusing on the last term we have that

$$\begin{aligned} &\log \left| \mathbf{I} - \left( \mathbf{I} - \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \right) \right| \\ &= \text{tr} \left( \log \left( \mathbf{I} - \left( \mathbf{I} - \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \right) \right) \right) \end{aligned}$$

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

$$\begin{aligned} &\log \left( \mathbf{I} - \left( \mathbf{I} - \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \right) \right) \\ &= \sum_{t=1}^{\infty} \frac{1}{t} \left( \mathbf{I} - \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \right)^t \end{aligned}$$

Again, following the analysis of (Nowozin, 2018), we can express the central moments of the random variable  $\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2}$  to the central moments of  $\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \left( \phi(\omega) \phi(\omega)^\top + \sigma^2 \mathbf{I} \right) \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2}$  denoted as  $\mu_i$  for  $i \geq 2$ . Also, as explained before, the dominant term will be ( $t = 2$ ), therefore we have that thus, using Eq. (S10) we have that

$$\begin{aligned} &\text{tr} \left( \mathbb{E}_{\mathbb{P}(\omega)} \left[ \log \left( \mathbf{I} - \left( \mathbf{I} - \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \widetilde{\mathbf{K}}_J \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1/2} \right) \right) \right] \right) \\ &= \mathcal{O}(1/J) \end{aligned}$$

Therefore, we can conclude that

$$\mathbb{E}_{\mathbb{P}(\omega)} \left[ \log |\widetilde{\mathbf{K}}_J| \right] - \log |\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}| = \mathcal{O}(1/J).$$

□

## 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

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

by computing just the first  $J \ll H$  terms, where  $J \in \{1, \dots, H\}$  is randomly drawn from the truncation distribution  $\mathbb{P}(\mathcal{J} = J)$ . RR and SS estimators, denoted as  $\bar{\psi}$ , offer two strategies for up-weighting the  $J$  computed terms, which, as we will prove below, yield an unbiased estimator  $\mathbb{E}_J[\bar{\psi}] = \psi$ .

### B.1. Unbiasedness of the RR Estimator

The RR estimator of Eq. (6) is unbiased, i.e.,  $\mathbb{E}_J[\bar{\psi}_J] = \psi$ .

*Proof.*

$$\begin{aligned} \mathbb{E}_J[\bar{\psi}_J] &= \mathbb{E}_J \left[ \sum_{j=1}^J \frac{\Delta_j}{\mathbb{P}(\mathcal{J} \geq j)} \right] \\ &= \mathbb{E}_J \left[ \sum_{j=1}^H \frac{\Delta_j}{\mathbb{P}(\mathcal{J} \geq j)} \cdot \mathbb{I}_{j \leq J} \right] \\ &= \sum_{j=1}^H \frac{\Delta_j}{\mathbb{P}(\mathcal{J} \geq j)} \mathbb{E}_J[\mathbb{I}_{j \leq J}] \\ &= \sum_{j=1}^H \frac{\Delta_j}{\mathbb{P}(\mathcal{J} \geq j)} \left[ \sum_{J=1}^H \mathbb{P}(\mathcal{J} = J) \cdot \mathbb{I}_{J \geq j} \right] \\ &= \sum_{j=1}^H \frac{\Delta_j}{\mathbb{P}(\mathcal{J} \geq j)} \mathbb{P}(\mathcal{J} \geq j) = \psi \end{aligned}$$

□

### B.2. Unbiasedness of the SS Estimator

The SS estimator of Eq. (7) is unbiased, i.e.,  $\mathbb{E}_J[\bar{\psi}_J] = \psi$ .

*Proof.*

$$\begin{aligned} \mathbb{E}_J[\bar{\psi}_J] &= \mathbb{E}_J \left[ \sum_{j=1}^H \frac{\Delta_j}{\mathbb{P}(\mathcal{J} = j)} \cdot \mathbb{I}_{Jj} \mathbb{1}_{j=J} \right] \\ &= \sum_{j=1}^H \frac{\Delta_j}{\mathbb{P}(\mathcal{J} = j)} \mathbb{E}_J[\mathbb{I}_{j=J}] \\ &= \sum_{j=1}^H \frac{\Delta_j}{\mathbb{P}(\mathcal{J} = j)} \sum_{J=1}^H \mathbb{P}(\mathcal{J} = J) \cdot \mathbb{I}_{J=j} \\ &= \sum_{j=1}^H \frac{\Delta_j}{\mathbb{P}(\mathcal{J} = j)} \mathbb{P}(\mathcal{J} = j) = \psi \end{aligned}$$

□

### 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

$$\begin{aligned} \mathbb{V}_J(\bar{\psi}) &= \sum_{j=1}^H \frac{\Delta_j^2}{\mathbb{P}(\mathcal{J} = j)^2} \mathbb{V}_J(\mathbb{I}_{\mathcal{J}=j}) \\ &= \sum_{j=1}^H \Delta_j^2 \left( \frac{1 - \mathbb{P}(\mathcal{J} = j)}{\mathbb{P}(\mathcal{J} = j)} \right) \end{aligned}$$

since  $\mathbb{I}_{\mathcal{J}=j}$  is a Bernoulli random variable with probability  $\mathbb{P}(\mathcal{J} = j)$ , 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.

$$\min_p \sum_{j=1}^H \Delta_j^2 \frac{1-p_j}{p_j} \quad \text{s.t.} \quad \sum_{j=1}^H p_j = 1, \quad p_j \geq 0$$

where  $p_j$  is acting as a shorthand of  $\mathbb{P}(\mathcal{J} = j)$ . The Lagrangian of the problem is

$$\mathcal{L}(p, \lambda) = \sum_{j=1}^H \Delta_j^2 \frac{1-p_j}{p_j} + \lambda \left( 1 - \sum_{j=1}^H p_j \right)$$

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

$$\frac{\partial \mathcal{L}}{\partial p_j}(p_j^*, \lambda^*) = - \left( \frac{\Delta_j^2}{p_j^*} \right)^2 + \lambda^* = 0$$

therefore, if we take the ratio of the probabilities with respect to the first we get that  $p_j^* = \frac{\Delta_j}{\Delta_1} p_1^*$ . If we substitute this expression in the equality constraint we get that

$$\mathbb{P}^*(\mathcal{J} = j) = \frac{\Delta_j}{\sum_{i=1}^H \Delta_i} \propto \Delta_j \quad (\text{S11})$$

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  $\Delta_j$  for  $j = 1, \dots, H$ . However, if we possess an estimate or a theoretical bound on rate of decay of each  $\Delta_j$  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  $\psi$  with  $\mathbb{P}(\mathcal{J} = J)$  as our proposal distribution.

$$\psi = \sum_{J=1}^H \Delta_J$$

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

$$= H \sum_{J=1}^H \frac{1}{H} \Delta_J,$$

We now introduce an alternative proposal distribution  $J \sim \mathbb{P}(\mathcal{J} = J)$ :

$$\begin{aligned} &= H \sum_{J=1}^H \frac{\Delta_J}{H \mathbb{P}(\mathcal{J} = J)} \mathbb{P}(\mathcal{J} = J) \\ &= \mathbb{E}_{J \sim \mathbb{P}(\mathcal{J} = J)} \left[ \frac{\Delta_J}{\mathbb{P}(\mathcal{J} = J)} \right], \end{aligned}$$

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  $\mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$  term can be estimated directly from the  $\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$  solve in Eq. (16). The  $\log |\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}|$  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:

$$\begin{aligned} \log |\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}| &\approx \|\mathbf{z}\|^2 \mathbf{e}_1^\top \left( \log \mathbf{T}_{\mathbf{z}}^{(N)} \right) \mathbf{e}_1 \quad (\text{S12}) \\ &= \|\mathbf{z}\|^2 \mathbf{e}_1^\top \left( \log \mathbf{T}_{\mathbf{z}}^{(1)} \right) \mathbf{e}_1 + \sum_{j=2}^N \Delta_j, \end{aligned}$$

where  $\Delta_j = \|\mathbf{z}\|^2 \mathbf{e}_1^\top \left( \log \mathbf{T}_{\mathbf{z}}^{(j)} - \log \mathbf{T}_{\mathbf{z}}^{(j-1)} \right) \mathbf{e}_1$ . We can therefore apply the Russian Roulette estimator with truncation  $\mathbb{P}(J)$  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  $v_J$  which that approximates  $\log |\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}|$  since the procedure is analogous for  $u_J$  which approximates  $\mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$ . Beatson & Adams (2019, Theorem 5.4) state that the optimal RR truncation distribution that maximizes the ROE is

$$\mathbb{P}^*(J \geq j) \propto \sqrt{\frac{\mathbb{E}[\|v_j\|^2]}{c(j) - c(j-1)}}$$

where  $c(j)$  is the computational cost of evaluating  $v_j$  which is the  $j$ -th term being approximated through Russian

Roulette. In this case, the  $v_j$  are nonrandom and single-valued and the cost per iteration is constant with respect to  $j$ . Hence, we can conclude that

$$\mathbb{P}^*(J \geq j) \propto v_j = \mathcal{O}(C^{-2j})$$

which implies that  $\mathbb{P}^*(J \leq j) \propto 1 - \mathcal{O}(C^{-2j})$  since the derivative of an exponential function is of the same form we have that  $\mathbb{P}^*(J = j) \propto \mathcal{O}(C^{-2j})$ .  $\square$

## 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-fit term  $\mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$  and the model complexity term  $\log |\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}|$ . For  $\mathbf{y}^\top \widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$  we create this connection through the Sherman-Morrisson formula and for the  $\log |\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}|$  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.

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

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

Note the difference between the  $\widetilde{\mathbf{K}}_J$  term used throughout the paper which includes the  $\sigma^2 \mathbf{I}$  against the  $\mathbf{K}_J$  term which only includes the RFF features. Moreover, to reduce clutter we define the following matrix

$$\mathbf{W} := \left( \frac{J}{J+1} \mathbf{K}_J + \sigma^2 \mathbf{I} \right)$$

whose use will become evident throughout the proof.

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

$$\|a\mathbf{A} + b\mathbf{I}\|_2^2 \leq \|\mathbf{A} + b\mathbf{I}\|_2^2$$

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  $a\mathbf{A} + b\mathbf{I} = \mathbf{V}(a\mathbf{D} + b\mathbf{I})\mathbf{V}^\top$ . This implies that the eigenvalues of  $\mathbf{A} + b\mathbf{I}$  are larger than those of  $a\mathbf{A} + b\mathbf{I}$ . An immediate consequence of this remark is that  $\|\mathbf{W}\|_2^2 \leq \|\mathbf{K}_J + \sigma^2 \mathbf{I}\|_2^2$  or that  $\|\mathbf{W}^{-1}\|_2^2 \leq \|\Lambda^{-1} + \sigma^{-2}\|_2^2 \leq \sigma^{-4}$  where  $\Lambda$  contains the positive eigenvalues of  $\mathbf{K}_J$ .

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

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

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

*Proof of Theorem 4.* We start by showing the rate of decay for the  $\Delta_j$  involving the data-fit terms. Note that we can express

$$\begin{aligned} & \mathbf{y}^\top (\mathbf{K}_{J+1} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} \\ &= \mathbf{y}^\top \left( \frac{J}{J+1} \mathbf{K}_J + \sigma^2 \mathbf{I} + \frac{\phi_{J+1} \phi_{J+1}^\top}{J+1} \right)^{-1} \mathbf{y} \quad (\text{S13}) \\ &= \mathbf{y}^\top \left( \mathbf{W} + \frac{\phi_{J+1} \phi_{J+1}^\top}{J+1} \right)^{-1} \mathbf{y} \end{aligned}$$

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

$$\mathbf{W}^{-1} - \frac{(\phi_{J+1}^\top \mathbf{W}^{-1})^\top (\phi_{J+1}^\top \mathbf{W}^{-1})}{(J+1) + \phi_{J+1}^\top \mathbf{W}^{-1} \phi_{J+1}}$$

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

$$\begin{aligned} & \mathbf{y}^\top (\sigma^2 \mathbf{I} + \mathbf{K}_J)^{-1} \mathbf{y} - \mathbf{y}^\top (\sigma^2 \mathbf{I} + \mathbf{K}_{J+1})^{-1} \mathbf{y} \\ & \leq \mathbf{y}^\top \left( \sigma^2 \mathbf{I} + \frac{J}{J+1} \mathbf{K}_J \right)^{-1} \mathbf{y} - \mathbf{y}^\top (\sigma^2 \mathbf{I} + \mathbf{K}_{J+1})^{-1} \mathbf{y} \\ & = \frac{(\phi_{J+1}^\top \mathbf{W}^{-1} \mathbf{y})^2}{(J+1) + \phi_{J+1}^\top \mathbf{W}^{-1} \phi_{J+1}} \\ & \leq \frac{(\phi_{J+1}^\top \mathbf{W}^{-1} \mathbf{y})^2}{J+1} \\ & \leq \frac{\|\phi_{J+1}\|_2^2 \|\mathbf{W}^{-1} \mathbf{y}\|_2^2}{J+1} \\ & \leq \frac{\|\phi_{J+1}\|_2^2 \sigma^{-4} \|\mathbf{y}\|_2^2}{J+1} \end{aligned} \quad (\text{S14})$$

where the first inequality follows from Remark 2, the second inequality occurs since  $\phi_{J+1}^\top \mathbf{W}^{-1} \phi_{J+1} > 0$ , 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

$$\begin{aligned} & \mathbb{E}_{\mathbb{P}(\omega)} \left[ \mathbf{y}^\top (\sigma^2 \mathbf{I} + \mathbf{K}_J)^{-1} \mathbf{y} \right] \\ & - \mathbb{E}_{\mathbb{P}(\omega)} \left[ \mathbf{y}^\top (\sigma^2 \mathbf{I} + \mathbf{K}_{J+1})^{-1} \mathbf{y} \right] \\ & = \mathcal{O}(1/J). \end{aligned}$$

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

$$|\mathbf{K}_{J+1} + \sigma^2 \mathbf{I}| = \left| \frac{J}{J+1} \mathbf{K}_J + \sigma^2 \mathbf{I} + \frac{\phi_{J+1} \phi_{J+1}^\top}{J+1} \right|$$

then by using the matrix determinant lemma we have that

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

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

$$\begin{aligned} &\mathbb{E}_{\mathbb{P}(\omega)} \left[ \log \left( 1 + \frac{\sigma^{-2}}{J+1} \phi_{J+1}^\top \phi_{J+1} \right) \right] \\ &= \frac{\sigma^{-2} \mathbb{E}_{\mathbb{P}(\omega)} \left[ \phi_{J+1}^\top \phi_{J+1} \right]}{J+1} + \mathcal{O}(1/J^2) \end{aligned}$$

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

$$\begin{aligned} &\mathbb{E}_{\mathbb{P}(\omega)} [\log |\mathbf{K}_{J+1} + \sigma^2 \mathbf{I}|] - \mathbb{E}_{\mathbb{P}(\omega)} [\log |\mathbf{K}_J + \sigma^2 \mathbf{I}|] \\ &= \frac{\sigma^{-2} \mathbb{E}_{\mathbb{P}(\omega)} \left[ \phi_{J+1}^\top \phi_{J+1} \right]}{J+1} + \mathcal{O}(1/J^2) \\ &= \mathcal{O}(1/J) \end{aligned}$$

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

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

□

## 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	<b>.106 ± .000</b>	.150 ± .000	.128 ± .001	.119 ± .002	.112 ± .002
Elevators	10.6K	<b>.360 ± .006</b>	.379 ± .005	<b>.365 ± .001</b>	.376 ± .006	<b>.362 ± .006</b>	<b>.360 ± .006</b>	<b>.360 ± .006</b>
Bike	11.1K	<b>.035 ± .003</b>	.071 ± .002	.063 ± .001	.045 ± .002	1.044 ± .334	<b>.040 ± .005</b>	<b>.035 ± .002</b>
Kin40K	25.6K	—	.248 ± .001	<b>.074 ± .000</b>	.152 ± .001	.081 ± .000	.093 ± .000	.091 ± .001
Protein	25.6K	—	.715 ± .007	<b>.547 ± .001</b>	.664 ± .008	.562 ± .005	<b>.541 ± .008</b>	<b>.541 ± .008</b>
KEGG	31.2K	—	.097 ± .004	.101 ± .001	<b>.088 ± .002</b>	<b>.089 ± .002</b>	.195 ± .064	<b>.087 ± .003</b>
KEGGU	40.7K	—	.125 ± .001	.129 ± .001	.122 ± .001	.123 ± .001	<b>.120 ± .000</b>	<b>.120 ± .000</b>
3DRoad	278K	—	.675 ± .001	.348 ± .001	.439 ± .002	.285 ± .003	.202 ± .003	<b>.114 ± .013</b>

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	<b>.425 ± .013</b>	.469 ± .016	.780 ± .006	<b>.442 ± .015</b>	<b>.426 ± .015</b>	<b>.429 ± .012</b>	<b>.425 ± .013</b>
Bike	11.1K	−.984 ± .018	−.799 ± .010	.099 ± .030	− <b>1.514 ± .024</b>	85.715 ± 49.088	−.971 ± .008	−.982 ± .018
Kin40K	25.6K	—	.464 ± .002	1.407 ± .004	− <b>.410 ± .003</b>	.427 ± .001	.468 ± .003	.449 ± .013
Protein	25.6K	—	1.105 ± .008	1.163 ± .003	1.014 ± .012	.951 ± .004	<b>.934 ± .006</b>	<b>.934 ± .006</b>
KEGG	31.2K	—	−.874 ± .011	6.311 ± 2.323	− <b>1.022 ± .023</b>	− <b>.981 ± .033</b>	.415 ± .653	−.884 ± .009
KEGGU	40.7K	—	−.636 ± .002	4.000 ± .303	− <b>.685 ± .005</b>	−.668 ± .005	−.637 ± .006	−.650 ± .004
3DRoad	278K	—	1.031 ± .002	1.317 ± .004	<b>.601 ± .004</b>	.831 ± .000	<b>.613 ± .010</b>	.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

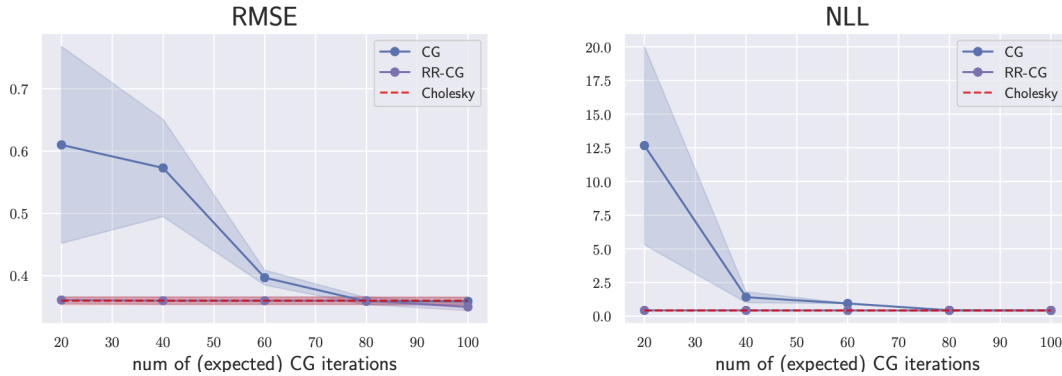


Figure S1. Predictive RMSE (left) and NLL (right) as a function of the number of (expected) CG iterations by optimizing Exact GP with CG (blue solid line) and RR-CG (purple solid line). Lower is better. The red dashed line corresponds to optimizing GP with Cholesky. The results are over three random seeds.

only be run in smaller datasets to be able to compare against Cholesky. We chose Elevators as this dataset requires the largest number of CG iterations before convergence hence making the difference of choosing between CG and RR-CG evident. We emphasize that RR-CG is a better alternative to early-truncated CG as the latter can perform poorly in some cases, whereas the former is robust to the choice of expected truncation number.

**F.2. Convergence of GP hyperparameters for other datasets.**

In this section, we show how SS-RFF and RR-CG converge to the Cholesky solution whereas the biased methods do not. We make this analysis for other datasets than the ones used in the main manuscript.

Using RFF with 1,000 or 1,500 features generates results that clearly diverge from the Cholesky solution and RFF with a 1,000 features generates numerically unstable training. In contrast, SS-RFF with 1,500 features is able to achieve the Cholesky solution after 1,000 iterations. Nonetheless, SS-RFF with 1,000 features also suffers from numerical instability as RFF 1,000 but at a lesser degree.

The RR-CG models converge to optimal solutions, while the (biased) CG models diverge. Only the results for a low number of iterations is plotted since for more than 40 iterations CG and RR-CG are indistinguishable from Cholesky.

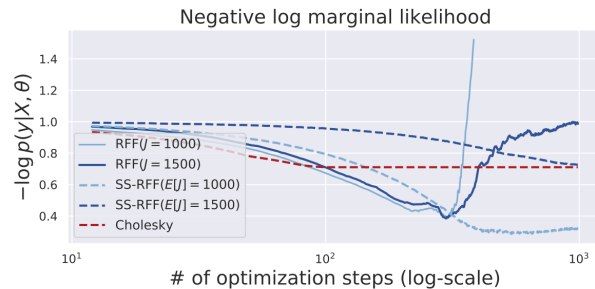


Figure S2. The GP optimization objective for models trained with RFF and SS-RFF. (Bike dataset, RBF kernel, Adam optimizer.) RFF models converge to sub-optimal log marginal likelihoods. SS-RFF models converge to (near) optimum values, yet require more than 100× as many optimization steps.

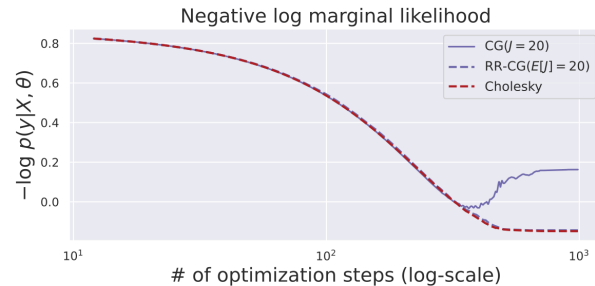


Figure S3. The GP optimization objective for models trained with CG and RR-CG. (PoleTele dataset, RBF kernel, Adam optimizer.) Models converge in < 100 steps of Adam.