

Preconditioning for Scalable Gaussian Process Hyperparameter Optimization



Jonathan Wenger^{1,2,3} Geoff Pleiss³ Philipp Hennig^{1,2} John Cunningham³ Jacob Gardner⁴

¹ University of Tübingen

² Max Planck Institute for Intelligent Systems, Tübingen

³ Columbia University

⁴ University of Pennsylvania



MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS



COLUMBIA | Zuckerman Institute
MORTIMER B. ZUCKERMAN MIND BRAIN BEHAVIOR INSTITUTE

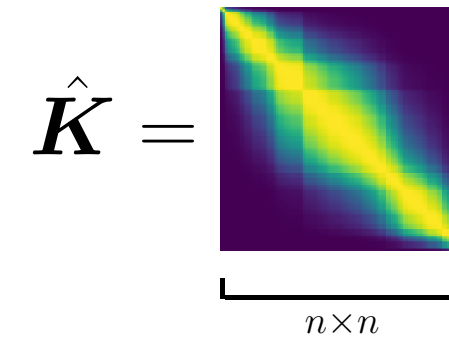


Goal: Scalable GP Hyperparam. Optim.

Need to: Evaluate log-marginal likelihood \mathcal{L} and its derivative $\frac{\partial}{\partial \theta} \mathcal{L}$ repeatedly.

Challenge: Costly $\mathcal{O}(n^3)$ operations with the kernel matrix.

- ▷ linear solves $\hat{\mathbf{K}}^{-1}(\cdot)$
- ▷ matrix traces $\log \det(\hat{\mathbf{K}}) = \text{tr}(\log(\hat{\mathbf{K}}))$ and $\text{tr}(\hat{\mathbf{K}}^{-1} \frac{\partial \hat{\mathbf{K}}}{\partial \theta})$



Known: Reducable to Matrix-Vector Mult.

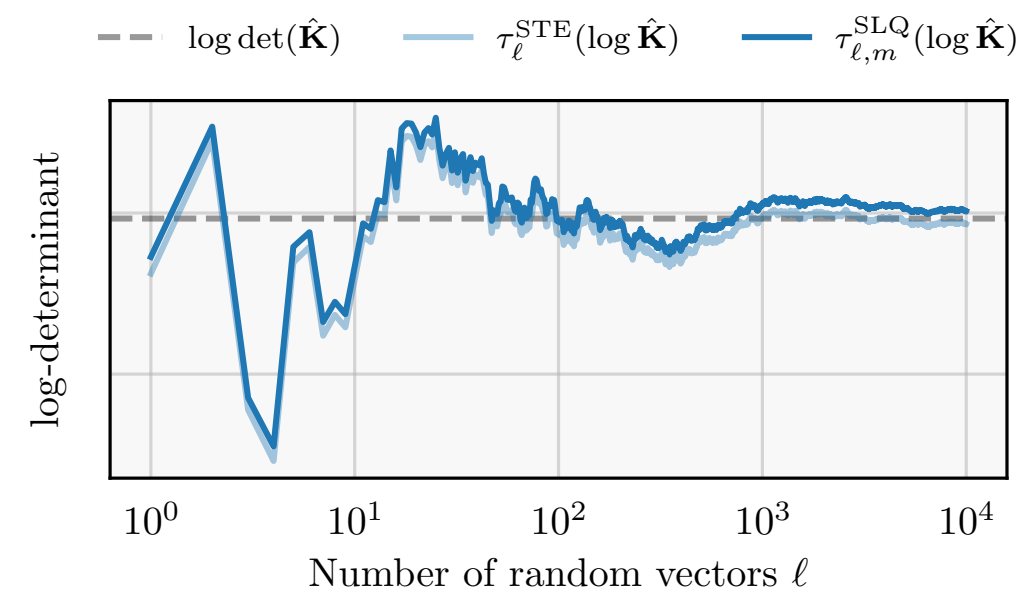
Linear Solves: Preconditioned CG **Matrix Traces:** Stochastic trace estimator

Great because ...

- ▷ matrix-vector multiplies cost at most $\mathcal{O}(n^2)$
- ▷ no need to store kernel matrix in memory
- ▷ can exploit parallelization and modern hardware (GPUs)

lower time and space
complexity

Problem: Stochastic Trace Estimators



$$\begin{aligned} \text{tr}(f(\hat{\mathbf{K}})) &= n \mathbb{E}[\mathbf{z}_i^\top f(\hat{\mathbf{K}}) \mathbf{z}_i] \\ &\approx \tau_\ell^{\text{STE}}(f(\hat{\mathbf{K}})) = \frac{n}{\ell} \sum_{i=1}^{\ell} \mathbf{z}_i^\top f(\hat{\mathbf{K}}) \mathbf{z}_i \\ &\approx \tau_{\ell,m}^{\text{SLQ}}(f(\hat{\mathbf{K}})) \end{aligned}$$

Bad because ...

- ▷ slow $\mathcal{O}(\ell^{-\frac{1}{2}})$ convergence in number of random vectors
- ▷ adds **noise** into hyperparameter optimization

slows down training

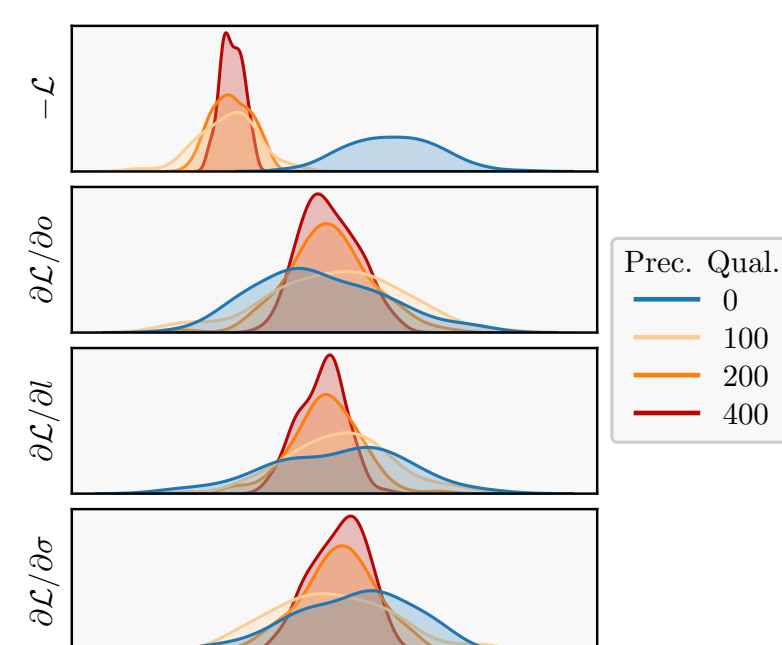
Our work: Precondition Trace Estimators

Insight

Can precondition not only **linear solves** but also **stochastic trace estimators**!

Contributions

- ▷ Preconditioning **reduces variance** of the STE, i.e. accelerates convergence.
- ▷ Theoretical guarantees.
- ▷ Preconditioner choices for given kernels.
- ▷ Up to twelvefold training speedup.



Background: Preconditioning

$$\hat{\mathbf{P}} \approx \hat{\mathbf{K}}$$

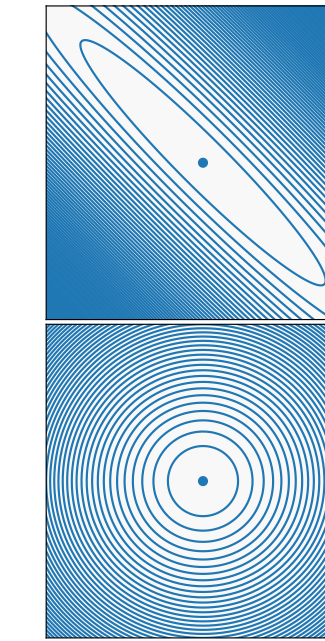
such that $\kappa(\hat{\mathbf{P}}^{-1} \hat{\mathbf{K}}) \ll \kappa(\hat{\mathbf{K}})$ and $\hat{\mathbf{P}}$ is tractable.

- ▷ Computing and storing $\hat{\mathbf{P}}$ is cheap.
- ▷ Linear solves $\mathbf{v} \mapsto \hat{\mathbf{P}}^{-1} \mathbf{v}$ are efficient.
- ▷ Derived properties (determinant, spectrum, ...) known.

Asymptotic approx. error $g(\ell) \rightarrow 0$ of precondition. seq. $\hat{\mathbf{P}}_\ell \rightarrow \hat{\mathbf{K}}$:

$$\kappa(\hat{\mathbf{P}}_\ell^{-1} \hat{\mathbf{K}}) \leq (1 + \mathcal{O}(g(\ell))) \|\hat{\mathbf{K}}\|_F^2$$

Known Use: Accelerate and stabilize linear solves via CG \Rightarrow **bias reduction**

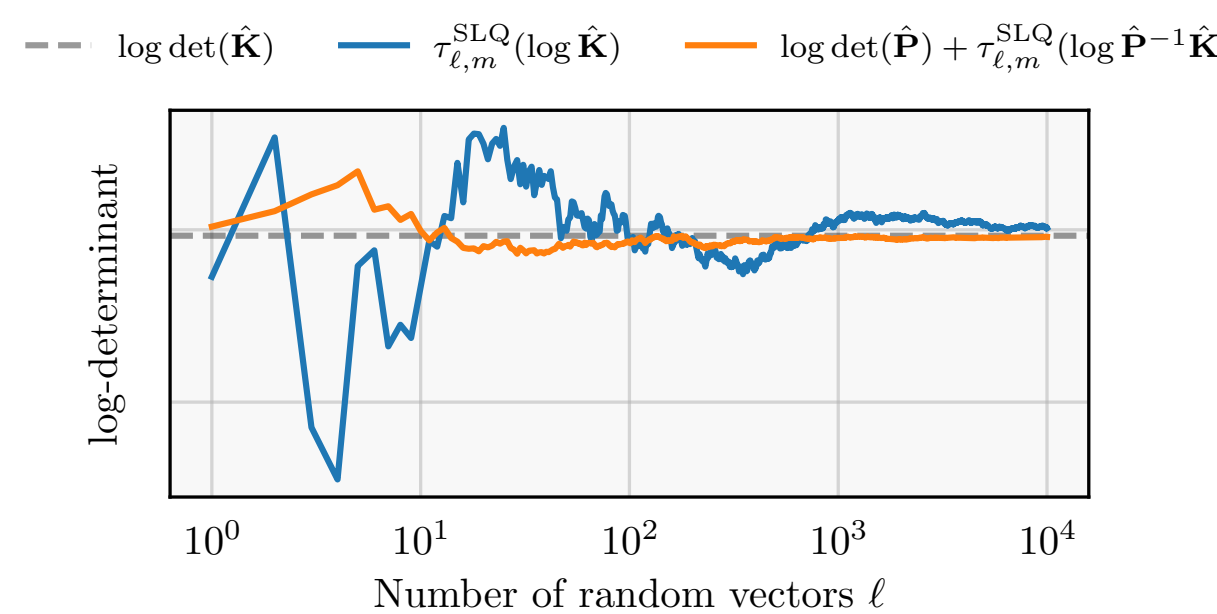


Precond. Log-Determinant Estimation

Idea: Decompose log-determinant into deterministic and stochastic approximation.

$$\log \det(\hat{\mathbf{K}}) = \log \det(\hat{\mathbf{P}}_\ell \hat{\mathbf{P}}_\ell^{-1} \hat{\mathbf{K}}) = \underbrace{\log \det(\hat{\mathbf{P}}_\ell)}_{\text{known}} + \underbrace{\text{tr}(\log(\hat{\mathbf{K}}) - \log(\hat{\mathbf{P}}_\ell))}_{\approx \text{stochastic trace estimate (STE)}}$$

Better preconditioner \Rightarrow smaller stochastic approximation \Rightarrow **variance reduction**



- ▷ Backward pass analogously via automatic differentiation.
- ▷ If we compute a preconditioner for CG, we can simply reuse it at negligible overhead.

If $\hat{\mathbf{P}}_\ell \rightarrow \hat{\mathbf{K}}$ at rate $g(\ell)$, then the STE only requires $\mathcal{O}(\ell^{-\frac{1}{2}} g(\ell))$ random vectors.

Theoretical Results

Probabilistic Error Bounds

Preconditioning not only **reduces bias**, but crucially also **reduces variance**.

Theorem (Log-marg. likelihood)

[...] Then with probability $1 - \delta$, the error in the estimate η of the log-marginal likelihood \mathcal{L} satisfies

$$|\eta - \mathcal{L}| \leq \varepsilon_{\text{CG}} + \frac{1}{2}(\varepsilon_{\text{Lanczos}} + \varepsilon_{\text{STE}}) \|\log(\hat{\mathbf{K}})\|_F,$$

where the errors are bounded by

$$\varepsilon_{\text{CG}}(\kappa, m) \leq K_3 \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1} \right)^m \quad (1)$$

$$\varepsilon_{\text{Lanczos}}(\kappa, m) \leq K_1 \left(\frac{\sqrt{2\kappa+1}-1}{\sqrt{2\kappa+1}+1} \right)^{2m} \quad (2)$$

$$\varepsilon_{\text{STE}}(\delta, \ell) \leq C_1 \sqrt{\log(\delta^{-1})} \ell^{-\frac{1}{2}} g(\ell) \quad (3)$$

Theorem (Derivative)

[...] Then with probability $1 - \delta$, the error in the estimate ϕ of the derivative of the log-marginal likelihood $\frac{\partial}{\partial \theta} \mathcal{L}$ satisfies

$$|\phi - \frac{\partial}{\partial \theta} \mathcal{L}| \leq \varepsilon_{\text{CG}} + \frac{1}{2}(\varepsilon_{\text{CG}'} + \varepsilon_{\text{STE}}) \left\| \hat{\mathbf{K}}^{-1} \frac{\partial \hat{\mathbf{K}}}{\partial \theta} \right\|_F$$

where the errors are bounded by

$$\varepsilon_{\text{CG}}(\kappa, m) \leq K_4 \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1} \right)^m \quad (4)$$

$$\varepsilon_{\text{CG}'}(\kappa, m) \leq K_2 \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1} \right)^m \quad (5)$$

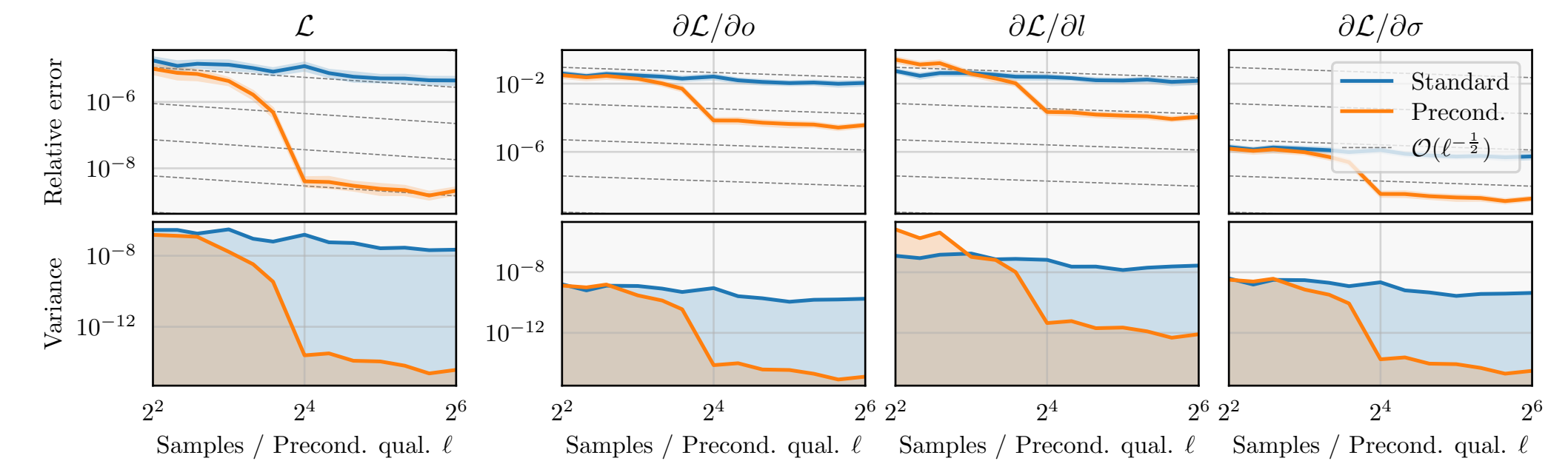
$$\varepsilon_{\text{STE}}(\delta, \ell) \leq C_1 \sqrt{\log(\delta^{-1})} \ell^{-\frac{1}{2}} g(\ell) \quad (6)$$

Convergence rates for combinations of kernels and preconditioners

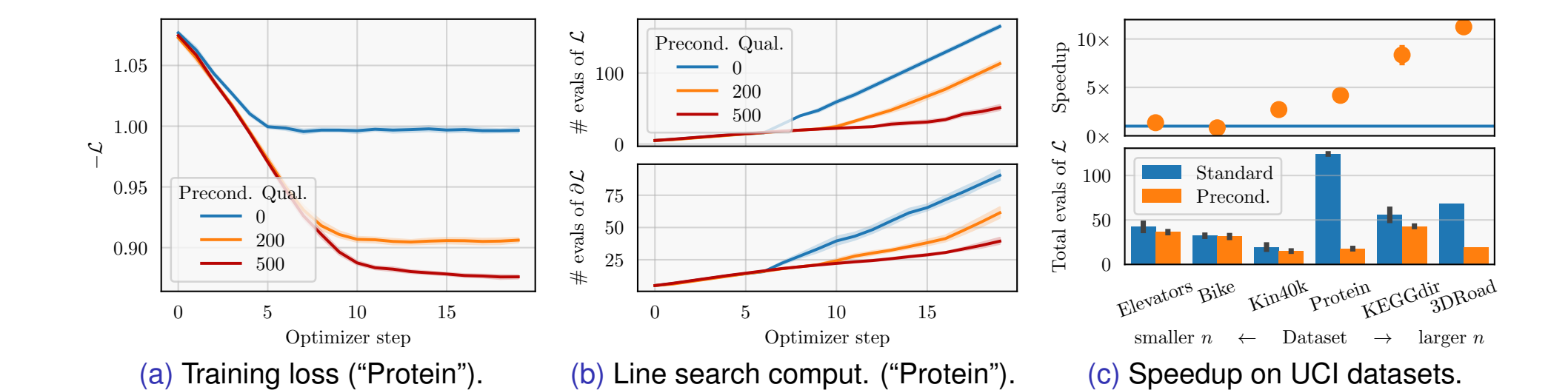
Kernel	d	Preconditioner	$g(\ell)$	Condition
any	\mathbb{N}	none	1	
...
any	\mathbb{N}	RFF	$\ell^{-\frac{1}{2}}$	w/ high probability
RBF	1	partial Cholesky	$\exp(-c\ell)$	for some $c > 0$
RBF	\mathbb{N}	QFF	$\exp(-b\ell^{\frac{1}{d}})$	for some $b > 0$ if $\ell^{\frac{1}{d}} > 2\gamma^{-2}$
Matérn(ν)	\mathbb{N}	partial Cholesky	$\ell^{-(\frac{2\nu}{d}+1)}$	$2\nu \in \mathbb{N}$ and maximin ordering
Matérn(ν)	1	QFF	$\ell^{-(s(\nu)+1)}$	where $s(\nu) \in \mathbb{N}$
mod. Matérn(ν)	\mathbb{N}	QFF	$\ell^{-\frac{s(\nu)+1}{d}}$	where $s(\nu) \in \mathbb{N}$
additive	\mathbb{N}	any	$dg(\ell)$	all summands have rate $g(\ell)$
any	\mathbb{N}	any kernel approx.	$g(\ell)$	\exists uniform convergence bound

Experiments

Preconditioning reduces bias and variance in \mathcal{L} and $\frac{\partial}{\partial \theta} \mathcal{L}$



Preconditioning reduces noise \Rightarrow accelerates hyperparam. optim.



Dataset	n	d	$-\mathcal{L}_{\text{train}} \downarrow$		$-\mathcal{L}_{\text{test}} \downarrow$		RMSE \downarrow		Runtime (s)	
			Standard	Precond.	Standard	Precond.	Standard	Precond.	Standard	Precond.
Elevators	12 449	18	0.4647	0.4377	0.4021	0.4022	0.3484	0.3482	53	39
Bike	13 034	17	-0.9976	-0.9985	-0.9934	-0.9877	0.0446	0.0454	31	37
Kin40k	30 000	8	-0.3339	-0.4332	-0.3141	-0.3135	0.0929	0.0949	187	45
Protein	34 297	9	0.9963	0.9273	0.8869	0.8835	0.5722	0.5577	893	43
KEGGdir	36 620	20	-0.9501	-1.0043	-0.9459	-0.9490	0.0861	0.0864	1450	174
3DRoad	326 155	3	0.7733	0.1284	1.4360	1.1690	0.2982	0.1265	82 200	7306

Paper
on
ARXIV



Implementation
as part of
GPYTORCH

