
Supplementary Material – Scalable Gaussian Process Separation for Kernels with a Non-Stationary Phase

Jan Graßhoff¹ Alexandra Jankowski¹ Philipp Rostalski¹

Overview of Kronecker and Toeplitz Methods

Table 1: Structure exploiting inference and learning methods. All kernels are assumed to be stationary and $\hat{K} = K + \sigma^2 I$.

Kernel and Inputs	Matrix	Linear Solve	Log Determinant
Kernel is separable: $k(\mathbf{x}, \mathbf{z}') = \prod_{d=1}^D k_d(\mathbf{x}^{(d)}, \mathbf{z}'^{(d)})$	$K = \bigotimes_{d=1}^D K_d$ and	noise-free: $K^{-1}\mathbf{y} = \bigotimes_{d=1}^D K_d^{-1}\mathbf{y}$	noise-free: $\log K = \sum_i V_{i,i}$
Inputs on a rectilinear grid: $X = \mathcal{X}_1 \times \dots \times \mathcal{X}_D$	$K = QVQ^\top$	noisy: $\hat{K}^{-1}\mathbf{y} = Q(V + \sigma^2 I)^{-1}Q^\top \mathbf{y}$	noisy: $\log \hat{K} = \sum_i (V_{i,i} + \sigma^2)$
Kernel: $k(x, x')$ Inputs: $x \in \mathbb{R}$ and equispaced	K is Toeplitz	LCG with fast MVMs	(1) circulant approx. (Wilson et al., 2015) (2) stoch. trace estim. (Dong et al., 2017)
Kernel is separable Inputs unstructured	$K \approx WK_UUW^\top$ (Wilson & Nickisch, 2015)	LCG with fast MVMs	(1) scaled eigenvalues (Wilson et al., 2015) (2) stoch. trace estim. (Dong et al., 2017)

Structure Exploitation for Kernels with a Non-Stationary Phase

Table 2: Comparison of the standard SKI (Wilson & Nickisch, 2015) approach using equidistant inducing points U and warpSKI. Inputs may be unstructured (or have partial grid structure). The kernels k (and k_i) are assumed to be stationary and separable. The functions $\phi_i : \mathcal{D} \rightarrow \mathcal{D}_i$ and $\phi : \mathcal{D} \rightarrow \mathcal{D}_1$ with $\mathcal{D}_{\text{in}} \subseteq \mathbb{R}^D$, $\mathcal{D}_i \subseteq \mathbb{R}^D$ are invertible functions. The linear solve $\hat{K}^{-1}\mathbf{y}$ is done by conjugate gradients and the log determinant is approximated using stochastic trace estimation.

Kernel	equidistant U recovers...	warpSKI recovers...
$k(\phi(x), \phi(x'))$ with $x \in \mathbb{R}$	–	Toeplitz structure
$\sum_i k_i(\phi_i(x), \phi_i(x'))$ with $x \in \mathbb{R}$	–	sum over Toeplitz structures
$k(\phi(\mathbf{x}), \phi(\mathbf{x}'))$, $\mathbf{x} \in \mathbb{R}^D$ and ϕ is not an elementwise fnc.	–	Kronecker and Toeplitz structure
$k(\phi(\mathbf{x}), \phi(\mathbf{x}'))$ $\mathbf{x} \in \mathbb{R}^D$ and ϕ is an elementwise fnc.	Kronecker structure	Kronecker and Toeplitz structure
$\sum_i k_i(\phi_i(\mathbf{x}), \phi_i(\mathbf{x}'))$ $\mathbf{x} \in \mathbb{R}^D$ and ϕ_i are elementwise fnc.	sum over Kronecker structures	sum over Kronecker and Toeplitz structures

¹Institute for Electrical Engineering in Medicine, Universität zu Lübeck, Germany. Correspondence to: Jan Graßhoff <j.grasshoff@uni-luebeck.de>.

Experimental Data - Numerical Results

Table 3: Inference runtime (s) for the numerical experiment with n data points and m inducing points. The results are averages over five samples \pm one standard deviation.

Points	Inducing Points		
	$m = 9933$	$m = 48824$	$m = 74836$
$n = 10^{2.75}$	0.11 \pm 0.02	0.28 \pm 0.05	0.47 \pm 0.07
$n = 10^3$	0.14 \pm 0.03	0.43 \pm 0.05	0.64 \pm 0.08
$n = 10^{3.5}$	0.34 \pm 0.06	1.10 \pm 0.10	1.61 \pm 0.17
$n = 10^4$	1.95 \pm 0.48	3.80 \pm 0.68	4.76 \pm 0.89
$n = 10^{4.5}$	19.9 \pm 4.4	25.8 \pm 3.7	26.3 \pm 4.7
$n = 10^5$	214 \pm 44	228 \pm 37	241 \pm 49

Table 4: Likelihood evaluation time (s) for the numerical experiment with n data points and m inducing points. The results are averages over five samples \pm one standard deviation.

Points	Inducing Points		
	$m = 9933$	$m = 48824$	$m = 74836$
$n = 10^{2.75}$	1.81 \pm 0.37	4.41 \pm 1.34	8.17 \pm 2.10
$n = 10^3$	1.81 \pm 0.37	5.74 \pm 1.06	10.0 \pm 1.7
$n = 10^{3.5}$	1.92 \pm 0.39	6.15 \pm 1.16	10.9 \pm 1.9
$n = 10^4$	2.35 \pm 0.48	6.67 \pm 0.68	12.0 \pm 2.1
$n = 10^{4.5}$	6.19 \pm 1.45	10.1 \pm 1.7	14.4 \pm 2.8
$n = 10^5$	17.0 \pm 3.7	19.8 \pm 3.6	26.0 \pm 5.2

Table 5: Learning runtime (s) for the numerical experiment with n data points and m inducing points. The results are averages over five samples \pm one standard deviation.

Points	Inducing Points		
	$m = 9933$	$m = 48824$	$m = 74836$
$n = 10^{2.75}$	28.3 \pm 9.3	84.6 \pm 40.5	166 \pm 53
$n = 10^3$	29.0 \pm 5.3	94.3 \pm 34.2	159 \pm 40
$n = 10^{3.5}$	26.6 \pm 16.0	92.0 \pm 20.8	177 \pm 37
$n = 10^4$	55.1 \pm 7.8	117 \pm 35	214 \pm 33
$n = 10^{4.5}$	130 \pm 49	163 \pm 23	252 \pm 88
$n = 10^5$	309 \pm 110	458 \pm 79	398 \pm 153

Table 6: RMSE for the numerical experiment with n data points and m inducing points. The results are averages over five samples \pm one standard deviation.

Points	Inducing Points		
	$m = 9933$	$m = 48824$	$m = 74836$
$n = 10^{2.75}$	0.29 \pm 0.02	0.30 \pm 0.03	0.29 \pm 0.03
$n = 10^3$	0.28 \pm 0.02	0.25 \pm 0.03	0.26 \pm 0.01
$n = 10^{3.5}$	0.20 \pm 0.02	0.19 \pm 0.01	0.19 \pm 0.01
$n = 10^4$	0.17 \pm 0.01	0.16 \pm 0.02	0.16 \pm 0.02
$n = 10^{4.5}$	0.16 \pm 0.03	0.15 \pm 0.02	0.16 \pm 0.02
$n = 10^5$	0.15 \pm 0.01	0.15 \pm 0.02	0.15 \pm 0.02

References

Dong, K., Eriksson, D., Nickisch, H., Bindel, D., and Wilson, A. Scalable log determinants for gaussian process kernel learning. In *Proceedings of the 32Nd International Conference on Neural Information Processing Systems (NeurIPS)*, pp. 6327–6337, 2017.

Wilson, A. G. and Nickisch, H. Kernel interpolation for scalable structured gaussian processes (kiss-gp). In *Proceedings of the 32Nd International Conference on International Conference on Machine Learning (ICML)*, pp. 1775–1784, 2015.

Wilson, A. G., Dann, C., and Nickisch, H. Thoughts on massively scalable gaussian processes. *CoRR*, abs/1511.01870, 2015. URL <http://arxiv.org/abs/1511.01870>.