

1 We thank the reviewers for their insightful feedback! We are encouraged they recognize the importance of probabilistic
2 linear solvers (PLS) for ML (R1, R2, R4) and the need to quantify uncertainty arising from finite computation (R1, R4).

3 **Contributions and Novelty** We are pleased the reviewers appreciate the originality of this work (R2, R3), its method-
4 ological contributions (R1, R2, R3), the stipulated desiderata marking a roadmap for PLS research (R2, R3), the value
5 of the proposed prior class (R2, R4), and the importance of the novel uncertainty calibration procedure (R2, R3). It
6 remedies a primary shortcoming of PLS, enables probabilistic stopping criteria and estimates $\log(\det(\mathbf{A}))$ (see below).
7 We also provide the first practical open-source implementation of a PLS returning distributions over \mathbf{A} , \mathbf{A}^{-1} and
8 \mathbf{x} . This is an important step towards a framework which targets computational resources "to reduce uncertainty as
9 required for a downstream task" (R4). The returned posterior means by the PLS are low-rank approximations to \mathbf{A} and
10 \mathbf{A}^{-1} , which among others find use in kernel methods. The returned covariances provide a bound to the error and their
11 structure can be exploited for novel use cases, e.g. as proposed for probabilistic mesh refinement in Galerkin's method.

12 **Bayesian Interpretation and Prior Class (R1, R3)** The generic inference procedure in Section 2 for a given prior
13 covariance $\mathbf{W}_0^{\mathbf{A}} \otimes \mathbf{W}_0^{\mathbf{A}}$ is Bayesian since it relies on Bayes' theorem. Algorithm 1 performs sequential Bayesian
14 updates for single action - observation pairs $(\mathbf{s}_i, \mathbf{y}_i)$. This can be seen by recognizing that the posterior (see Section 2.1)
15 is of the same form as the prior (for any $1 \leq k \leq n$). Guided by the desiderata in Table 1, we restrict the $n \times n$ DoFs
16 in the prior. This results in the proposed prior class in eq. (3). This prior and our calibration procedure depend on
17 the entire collected 'data' during a run of Algorithm 1. When using the proposed prior class our method is thus not
18 strictly Bayesian in the philosophical sense, but empirical Bayesian (i.e. it uses data to fit hyperparameters of the prior).
19 As this approach is standard in GP regression (where kernel parameters are set by type-II maximum likelihood), we
20 neglected to make this distinction. We will clarify this in the final version. This leaves the question how the algorithm
21 is realizable for the proposed prior (3) given its dependence on future data. The posterior mean in Section 2.1 only
22 depends on $\mathbf{W}_0^{\mathbf{A}} \mathbf{S} = \mathbf{Y}$ not on $\mathbf{W}_0^{\mathbf{A}}$ alone. By eq. (3), this product is given by the previously made observations \mathbf{Y} .
23 Similar reasoning applies for the inverse. Now, the posterior covariances do depend on $\mathbf{W}_0^{\mathbf{A}}$, resp. $\mathbf{W}_0^{\mathbf{H}}$ alone, but
24 during a run of Algorithm 1, we only require $\text{tr}(\text{Cov}[\mathbf{x}])$ for the stopping criterion. We show in Section S4.5 under the
25 assumptions of Theorem 2 how to compute this at any iteration i without access to future actions and observations.

26 **Calibration Procedure (R1, R3)** Calibration ensures that the uncertainty returned by the solver has the right scale,
27 i.e. it bounds the expected (relative) error (see Sections 2.2 and S4.5). Since the policy π only depends on the
28 posterior mean $\mathbb{E}[\mathbf{H}]$ and not the covariance, the hyperparameters Φ, Ψ and thus calibration do *not* change the
29 solution estimate \mathbf{x}_i only its associated covariance at iteration i . While structure in the uncertainty is preserved,
30 miscalibration negatively impacts the probabilistic termination criterion. In our experiments in Table 2 the solver
31 without calibration terminates early, since it is overconfident and thus has larger error than with calibration. *Why not*
32 *return $\mathcal{N}(\hat{\mathbf{x}}^{\text{CG}}, \text{span}(\mathbf{S})^\perp(\text{GP output at } t + 1))$?* (R1): This is similar to Algorithm 1's output assuming Theorem 2
33 holds. However, it omits prior knowledge about the space $\text{span}(\mathbf{S})$ explored by the algorithm (e.g. information about
34 the dominant eigenspectrum). Further, only using the GP prediction at $t + 1$ implies that the algorithm's uncertainty
35 about the action of \mathbf{A} in \mathbf{S}^\perp is of the same order as the next eigenvalue. This ignores any information about eigenvalues
36 $\lambda_{t+2}, \dots, \lambda_n$ contained in Rayleigh quotients and a priori known decay patterns for specific matrix classes.

37 **Importance of Noise (R4)** When referring to noise, we consider matrix-vector products of the form $\mathbf{v} \mapsto (\mathbf{A} + \mathbf{E}_i)\mathbf{v}$,
38 where $\mathbf{E}_i \in \mathbb{R}_{\text{sym}}^{n \times n}$ is Gaussian with zero mean. CG fails to converge in such a setting. While this approach can
39 model floating-point arithmetic, typically in ML settings noise from subsampling dominates. An important example is
40 large-scale empirical risk minimization. Due to memory constraints data needs to be batched and thus only approximate
41 Hessian-vector products $\mathbf{H}_{\text{batch}}\mathbf{v} = (\mathbf{H} + \mathbf{E}_{\text{batch}})\mathbf{v}$ are available. One could use a PLS for Hessian-free optimization
42 in this setting. This results in a trade-off between computing an accurate Hessian by sampling new batches in each
43 iteration of the solver and taking more optimization steps in parameter space. This approach results in an optimizer
44 which interpolates between SGD and Newton's method, depending on the batch size and number of PLS iterations k .

45 **Applications** *Transfer Learning (R1)*: Using a posterior from a related problem as a prior on a new problem has the
46 advantage over only setting $\mathbf{x}_0^{\text{new}} = \mathbf{x}_k^{\text{prev}}$, that uncertainty in already explored directions \mathbf{S} is low. Hence, if the new
47 problem $(\mathbf{A}^{\text{new}}, \mathbf{b}^{\text{new}})$ is similar, the covariance will contract faster. In turn also convergence will be faster (as in subspace
48 recycling). *Kernel matrix inversion (R4)*: We recognize the variety of methods available for Gram matrix inversion in the
49 Gaussian process setting. While a comparison for different priors adapted to the kernel choice vs. a set of inducing point
50 methods is an interesting experiment, this would have exceeded the scope of this paper. *Log-Determinant Estimation*
51 (R3, R4): The PLS can estimate the log-determinant in $\mathcal{O}(n)$ using the proposed ln-Rayleigh regression model for
52 uncertainty calibration via $\ln(\det(\mathbf{A})) = -\sum_{i=1}^n \ln R(\mathbf{A}, \mathbf{s}_i) \approx -(\sum_{i=1}^k \ln R(\mathbf{A}, \mathbf{s}_i) + \sum_{i=k+1}^n \mathbb{E}[\ln R_i | \mathbf{A}, \mathbf{S}])$.
53 *Galerkin's Method (R2, R4)*: When using a PLS as part of Galerkin's method the posterior (predictive) on a refined
54 mesh can be derived analytically (see Proposition S6). We leave comparisons to multi-grid methods for future work.

55 **Other (R3, R4)** *Reorthogonalization (R3)* is a consequence of the policy choice $\mathbf{s}_i = -\mathbb{E}[\mathbf{H}]\mathbf{r}_i$, where $\mathbb{E}[\mathbf{H}]$ depends
56 on all previous search directions as opposed to just \mathbf{s}_{i-1} for (naive) CG. *Unification of PLS theory (R4)*: Bartels et al.
57 [13] demonstrate that the matrix-based view generalizes the solution-based view. We focus on presenting a unified
58 matrix-based framework, which among others, connects the inference perspectives for \mathbf{A} and \mathbf{A}^{-1} in a rigorous way.