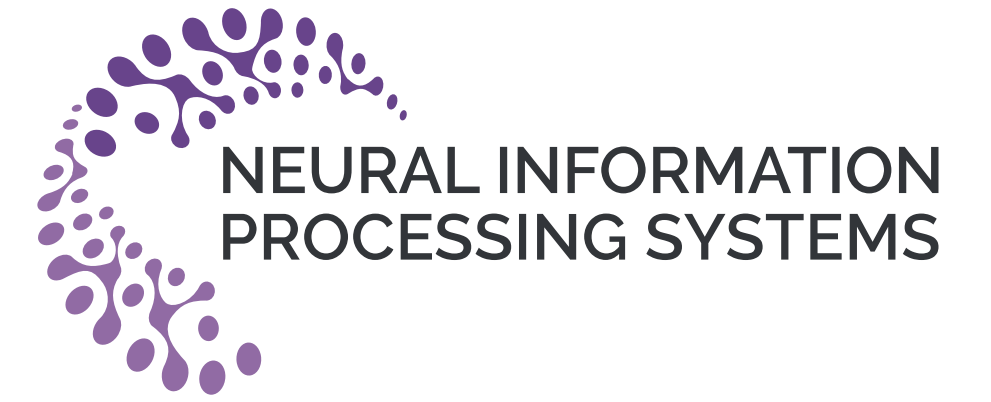


Probabilistic Linear Solvers for Machine Learning

Jonathan Wenger, Philipp Hennig

University of Tübingen and Max Planck Institute for Intelligent Systems, Tübingen, Germany



Abstract: Linear systems are the bedrock of virtually all numerical computation. Machine learning poses specific challenges for the solution of such systems due to their scale, characteristic structure, stochasticity and the central role of uncertainty in the field. Unifying earlier work we propose a class of probabilistic linear solvers which jointly infer the matrix, its inverse and the solution from matrix-vector product observations. This class emerges from a fundamental set of desiderata which constrains the space of possible algorithms and recovers the method of conjugate gradients under certain conditions. We demonstrate how to incorporate prior spectral information in order to calibrate uncertainty and experimentally showcase the potential of such solvers for machine learning.

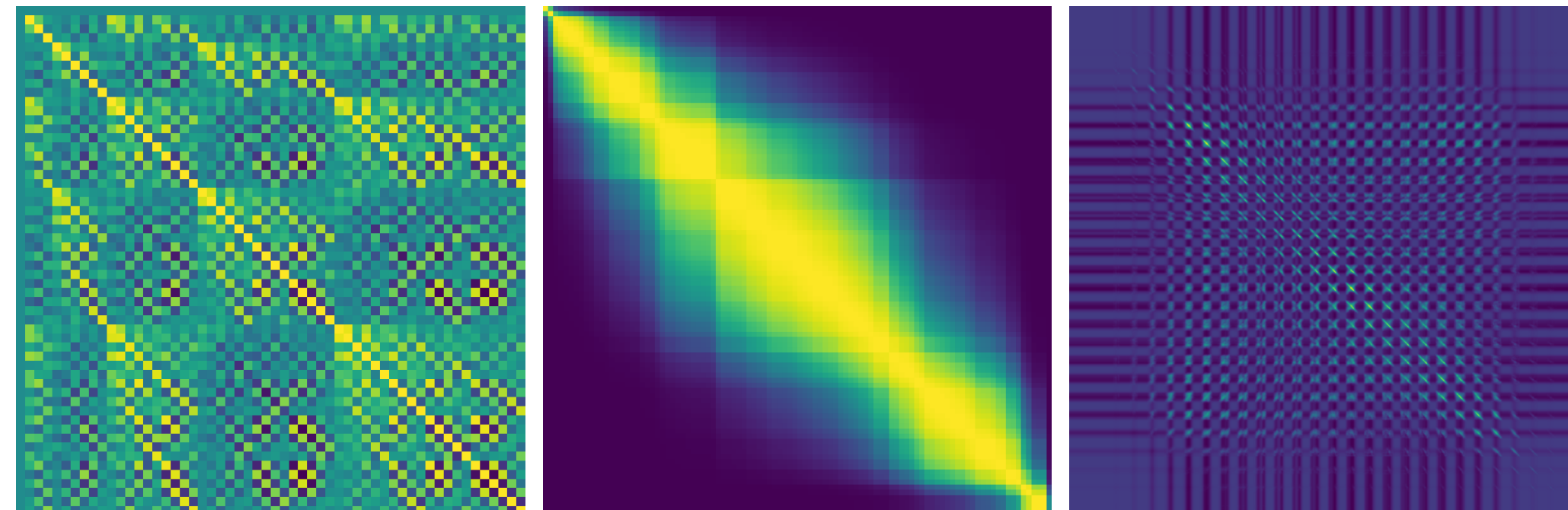
Contribution

- ▷ Unified probabilistic linear solver estimating \mathbf{A} , \mathbf{A}^{-1} and the solution \mathbf{x}_*
- ▷ Incorporation of prior information on the matrix \mathbf{A} or its inverse \mathbf{A}^{-1}
- ▷ Covariance class which enables uncertainty calibration via prior spectral information and recovers the conjugate gradient method

Linear Systems in Machine Learning

$$\mathbf{A}\mathbf{x}_* = \mathbf{b}$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ symmetric positive definite.



(a) General lin. model $\mathbf{X}^T \mathbf{X}$ (b) Kernel matrix $k_{\text{RBF}}(\mathbf{X}, \mathbf{X})$ (c) Hessian matrix $\nabla^2 f(\theta)$

ML-specific Systems

- ▷ large-scale
- ▷ charact. structure
- ▷ generative information
- ▷ noise

Applications: regression, kernel methods, Kalman filtering, Gaussian processes, spectral graph theory, differential equations, (stochastic) second-order optimization

Probabilistic Linear Solvers

Probabilistic linear solvers (PLS) [1, 2, 3] iteratively build a model for the linear operator \mathbf{A} , its inverse $\mathbf{H} = \mathbf{A}^{-1}$ or the solution \mathbf{x}_* , represented by random variables \mathbf{A}, \mathbf{H} or \mathbf{x} . In the framework of *probabilistic numerics* [4, 5] such solvers are Bayesian agents performing *inference* via linear *observations* $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_k] \in \mathbb{R}^{n \times k}$ resulting from *actions* $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_k] \in \mathbb{R}^{n \times k}$ given by a *policy* $\pi(\mathbf{s} \mid \mathbf{A}, \mathbf{H}, \mathbf{x}, \mathbf{A}, \mathbf{b})$. For a matrix-variate prior $p(\mathbf{A})$ or $p(\mathbf{H})$ our solver computes posterior beliefs over the matrix, its inverse and the solution.

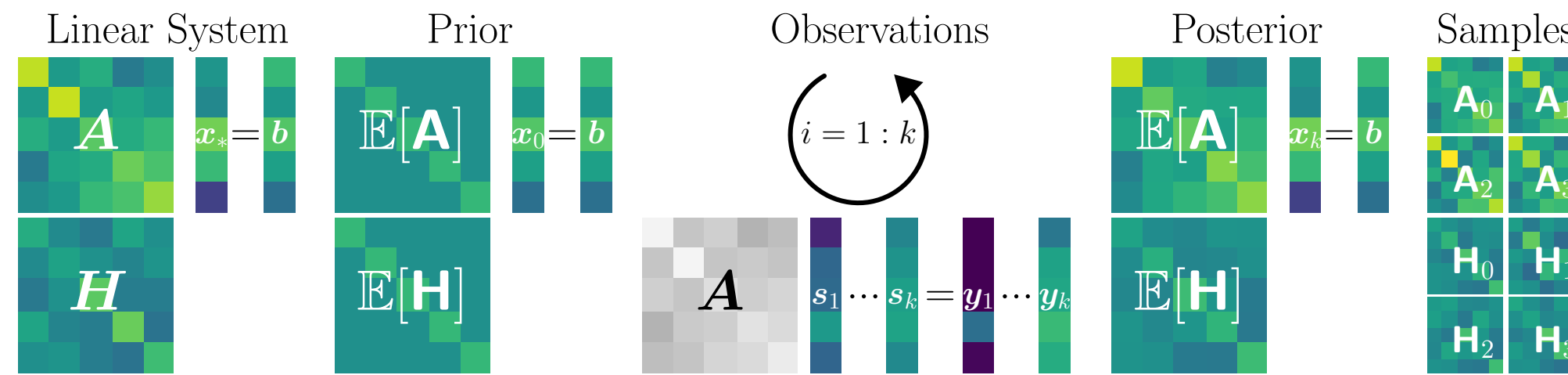


Figure: Illustration of a probabilistic linear solver.

Bayesian Inference Framework

Prior $p(\mathbf{A}) = \mathcal{N}(\mathbf{A}; \mathbf{A}_0, \mathbf{W}_0^{\mathbf{A}} \otimes \mathbf{W}_0^{\mathbf{A}})$
 Likelihood $p(\mathbf{Y} \mid \mathbf{A}, \mathbf{S}) = \lim_{\varepsilon \downarrow 0} \mathcal{N}(\mathbf{Y}; \mathbf{A}\mathbf{S}, \varepsilon^2 \mathbf{I} \otimes \mathbf{I}) = \delta(\mathbf{Y} - \mathbf{A}\mathbf{S})$
 Posterior $p(\mathbf{A} \mid \mathbf{S}, \mathbf{Y}) = \mathcal{N}(\mathbf{A}; \mathbf{A}_k, \Sigma_k)$

$$\mathbf{A}_k = \mathbf{A}_0 + \Delta_0^{\mathbf{A}} \mathbf{U}^T + \mathbf{U}(\Delta_0^{\mathbf{A}})^T - \mathbf{U} \mathbf{S}^T \Delta_0^{\mathbf{A}} \mathbf{U}^T$$

$$\Sigma_k = \mathbf{W}_0^{\mathbf{A}} (\mathbf{I}_n - \mathbf{S} \mathbf{U}^T) \otimes \mathbf{W}_0^{\mathbf{A}} (\mathbf{I}_n - \mathbf{S} \mathbf{U}^T)$$

where $\Delta_0^{\mathbf{A}} = \mathbf{Y} - \mathbf{A}_0 \mathbf{S}$ and $\mathbf{U} = \mathbf{W}_0^{\mathbf{A}} \mathbf{S} (\mathbf{S}^T \mathbf{W}_0^{\mathbf{A}} \mathbf{S})^{-1}$. Similarly for inverse model \mathbf{H} .

Algorithm 1: Probabilistic Linear Solver with Uncertainty Calibration

```

1 procedure PROBLINSOLVE( $\mathbf{A}(\cdot), \mathbf{b}, \mathbf{A}, \mathbf{H}$ )           # prior for  $\mathbf{A}$  or  $\mathbf{H}$ 
2    $\mathbf{x}_0 \leftarrow \mathbb{E}[\mathbf{H}]\mathbf{b}$                                # initial guess
3    $\mathbf{r}_0 \leftarrow \mathbf{A}\mathbf{x}_0 - \mathbf{b}$ 
4   while  $\min(\sqrt{\text{tr}(\text{Cov}[\mathbf{x}])}, \|\mathbf{r}_i\|_2) > \max(\delta_{\text{rtol}}\|\mathbf{b}\|_2, \delta_{\text{atol}})$  do
5      $\mathbf{s}_i \leftarrow -\mathbb{E}[\mathbf{H}]\mathbf{r}_{i-1}$                      # compute action via policy
6      $\mathbf{y}_i \leftarrow \mathbf{A}\mathbf{s}_i$                                # make observation
7      $\alpha_i \leftarrow -\mathbf{s}_i^T \mathbf{r}_{i-1} (\mathbf{s}_i^T \mathbf{y}_i)^{-1}$        # optimal step size
8      $\mathbf{x}_i \leftarrow \mathbf{x}_{i-1} + \alpha_i \mathbf{s}_i$                  # update solution estimate
9      $\mathbf{r}_i \leftarrow \mathbf{r}_{i-1} + \alpha_i \mathbf{y}_i$                  # update residual
10     $\mathbf{A} \leftarrow \text{INFER}(\mathbf{A}, \mathbf{s}_i, \mathbf{y}_i)$                # infer posterior distributions
11     $\mathbf{H} \leftarrow \text{INFER}(\mathbf{H}, \mathbf{s}_i, \mathbf{y}_i)$ 
12     $\Phi, \Psi \leftarrow \text{CALIBRATE}(\mathbf{S}, \mathbf{Y})$            # calibrate uncertainty
13     $\mathbf{x} \leftarrow \mathcal{N}(\mathbf{x}_k, \text{Cov}[\mathbf{H}\mathbf{b}])$              # belief over solution
14  return  $(\mathbf{x}, \mathbf{A}, \mathbf{H})$ 

```

Theoretical Properties

- ▷ Conjugate directions method (convergence in $\leq n$ iterations)
- ▷ Recovers the conjugate gradient method [6] given certain priors.
- ▷ Time complexity $\mathcal{O}(kn^2)$
- ▷ Space complexity $\mathcal{O}(kn)$

Covariance Class and Uncertainty Calibration

Problem: Probabilistic Linear Solvers are typically miscalibrated.

Idea: Covariance class which connects \mathbf{A} and \mathbf{H} and enables calibration.

$$\mathbf{W}_0^{\mathbf{A}} = \mathbf{A}\mathbf{S}(\mathbf{S}^T \mathbf{A}\mathbf{S})^{-1} \mathbf{S}^T \mathbf{A} + \mathbf{P}_{\mathbf{S}^\perp} \Phi \mathbf{P}_{\mathbf{S}^\perp},$$

$$\mathbf{W}_0^{\mathbf{H}} = \mathbf{A}_0^{-1} \mathbf{Y}(\mathbf{Y}^T \mathbf{A}_0^{-1} \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{A}_0^{-1} + \mathbf{P}_{\mathbf{Y}^\perp} \Psi \mathbf{P}_{\mathbf{Y}^\perp},$$

Experiments

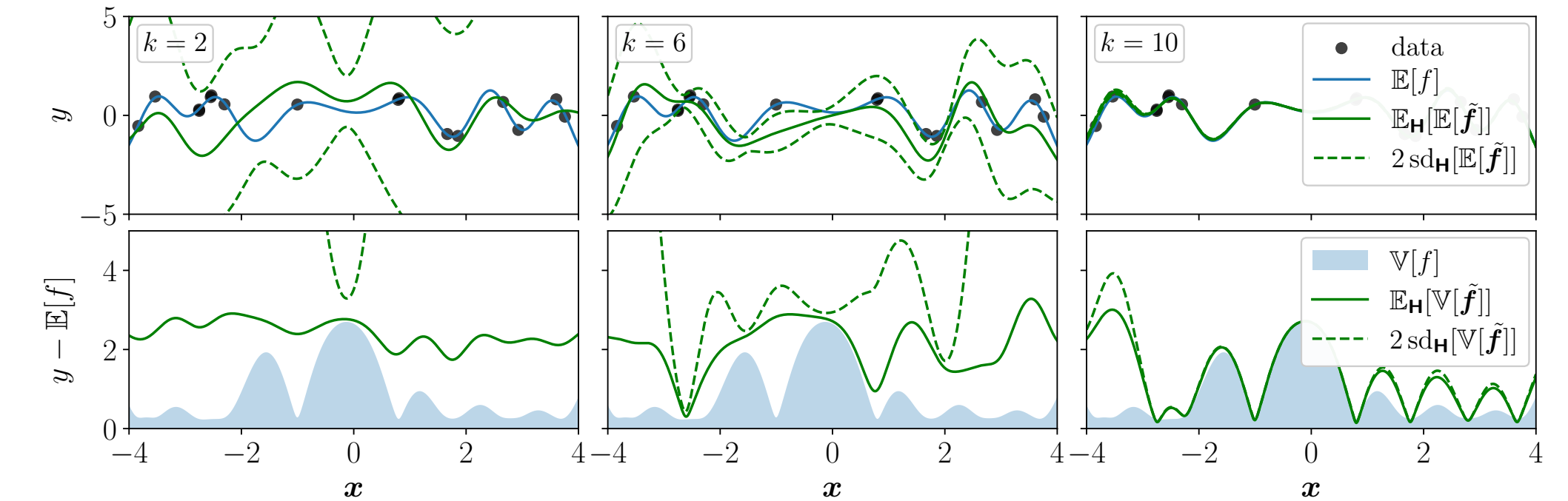


Figure: Numerical uncertainty in GP inference.

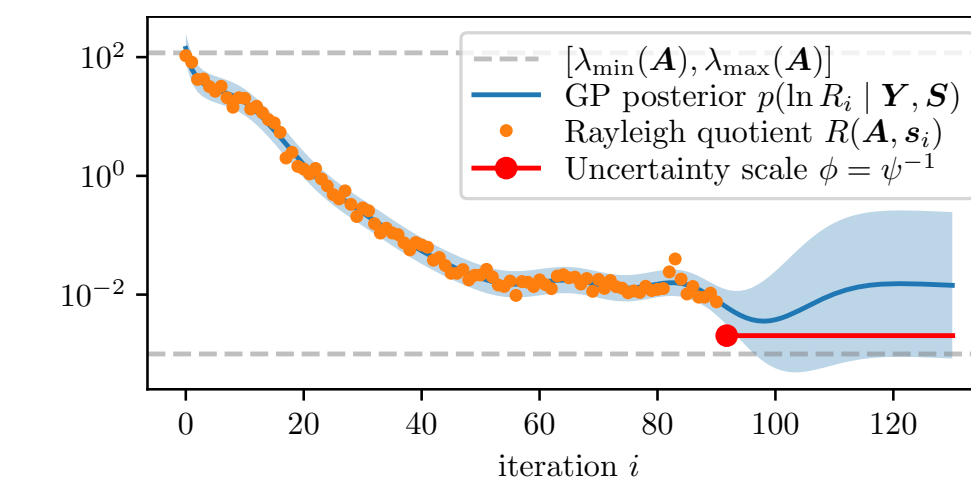


Figure: Rayleigh regression. Calibration via one-dim. GP regression on $\{\ln R(\mathbf{A}, \mathbf{s}_i)\}_{i=1}^k$ for an $n = 1000$ dim. Matérn32 kernel matrix.

Table: Calibration on kernel matrices. For $\bar{w} \approx 0$ the solver is well calibrated, for $\bar{w} \gg 0$ under- and for $\bar{w} \ll 0$ overconfident.

Kernel	n	none	Rayleigh	ε^2	$\bar{\lambda}_{k+1:n}$
Matérn32	10^2	-5.99	-0.24	0.32	0.09
Matérn32	10^3	-1.93	7.53	4.26	4.19
Matérn32	10^4	3.87	17.16	8.48	8.47
Matérn52	10^2	-7.84	-1.01	-0.76	-0.80
Matérn52	10^3	-4.63	1.43	-0.80	-0.81
Matérn52	10^4	-4.34	10.81	0.80	0.80
RBF	10^2	-7.53	-0.70	-0.84	-0.87
RBF	10^3	-4.94	6.60	0.77	0.77
RBF	10^4	0.14	21.32	2.92	2.92

Summary

- ▷ Linear systems in ML exhibit characteristic structure.
- ▷ Limited computational resources induce *numerical uncertainty*.
- ▷ Probabilistic linear solvers
 - ... reinterpret classic linear solvers as Bayesian inference.
 - ... quantify numerical uncertainty when solving linear systems.
 - ... make use of prior generative information.
 - ... allow propagation of information between problems.



ProbNum: Probabilistic numerical methods in Python.

<https://github.com/probabilistic-numerics/probnum>

or alternatively `pip install probnum`.

References

- [1] Philipp Hennig (2015). "Probabilistic Interpretation of Linear Solvers". In: SIAM Journal on Optimization 25(1):pp. 234–260.
- [2] Jon Cockayne et al. (2019). "Bayesian probabilistic numerical methods". In: SIAM Review 61(4):pp. 756–789.
- [3] Simon Bartels et al. (2019). "Probabilistic linear solvers: A unifying view". In: Statistics and Computing 29(6):pp. 1249–1263.
- [4] Philipp Hennig et al. (2015). "Probabilistic numerics and uncertainty in computations". In: Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences 471(2179).
- [5] Chris Oates et al. (2019). "A modern retrospective on probabilistic numerics". In: Statistics and Computing.
- [6] Magnus Rudolph Hestenes et al. (1952). "Methods of conjugate gradients for solving linear systems". In: Journal of Research of the National Bureau of Standards 49.

