

BAYESCG AS AN UNCERTAINTY AWARE VERSION OF CG*

TIM W. REID[†], ILSE C. F. IPSEN[†], JON COCKAYNE[‡], AND CHRIS J. OATES[§]

Abstract. The Bayesian Conjugate Gradient method (BayesCG) is a probabilistic generalization of the Conjugate Gradient method (CG) for solving linear systems with real symmetric positive definite coefficient matrices. Our CG-based implementation of BayesCG under a structure-exploiting prior distribution represents an ‘uncertainty-aware’ version of CG. Its output consists of CG iterates and posterior covariances that can be propagated to subsequent computations. The covariances have low-rank and are maintained in factored form. This allows easy generation of accurate samples to probe uncertainty in downstream computations. Numerical experiments confirm the effectiveness of the low-rank posterior covariances.

Key words. Symmetric positive semi-definite matrix, Krylov space method, Gaussian probability distribution, Bayesian inference, covariance matrix, mean, Moore-Penrose inverse, projectors in semi-definite inner products

AMS subject classifications. 65F10, 62F15, 65F50, 15A06, 15A10

1. Introduction. The solution of linear systems

$$(1.1) \quad \mathbf{Ax}_* = \mathbf{b},$$

with symmetric positive definite coefficient matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is an important problem in computational science and engineering. For large and sparse matrices \mathbf{A} , the preferred solver is the Conjugate Gradient method (CG) [26, 31]. This is a Krylov subspace method that, starting from a user-specified initial guess \mathbf{x}_0 , produces iterates \mathbf{x}_m that, the user hopes, ultimately converge to the solution \mathbf{x}_* . In practice, CG is terminated early, once the residual $\|\mathbf{b} - \mathbf{Ax}_m\|$ is sufficiently small in some norm. Early termination introduces a source of uncertainty since the solution \mathbf{x}_* has not been exactly computed.

We seek to create an ‘uncertainty aware’ version of CG that models the uncertainty in our knowledge of \mathbf{x}_* due to early termination. From the UQ perspective, this represents an instance of model discrepancy with epistemic uncertainties. Our motivation is to understand how the accuracy of the CG output \mathbf{x}_m affects downstream computations in a *computational pipeline* [12, Section 5], [25], that is, sequences of computations where the output of one computation is the input to another [7, 23, 40, 43, 44]. Traditional normwise CG error estimates are inadequate, because subsequent computations may not be able to make effective use of them. In contrast, a probabilistic model of the uncertainty, in the form of a distribution, can be propagated so that downstream computations can sample from the distribution to probe the effect of uncertainty on their own computations.

This is the mission of *probabilistic numerics*¹: Modelling the uncertainty in de-

*Submitted to the editors

Funding: The work was supported in part by NSF grant DMS-1745654 (TWR, ICFI), NSF grant DMS-1760374 and DOE grant DE-SC0022085 (ICFI), and the Lloyd’s Register Foundation Programme on Data Centric Engineering at the Alan Turing Institute (CJO).

[†]Department of Mathematics, North Carolina State University, Raleigh, NC 27695-8205, USA, (twreid@alumni.ncsu.edu, ipsen@ncsu.edu)

[‡]Department of Mathematical Sciences, University of Southampton, Southampton SO17 1BJ, UK (jon.cockayne@soton.ac.uk)

[§]School of Mathematics, Statistics and Physics, Newcastle University, Newcastle upon Tyne NE1 7RU, UK (chris.oates@ncl.ac.uk)

¹<https://www.probabilistic-numerics.org/>

37 deterministic computations with a probabilistic treatment of the errors [25, 42]. The
 38 origins of probabilistic numerics can be traced back to Poincaré [42], while a rigor-
 39 ous modern perspective is established in [12]. Probabilistic numerical methods have
 40 been developed for Bayesian optimization [38], subsequently applied to hyperparam-
 41 eter optimization in machine learning [46]; numerical integration [4, 14, 29], sparse
 42 Cholesky decompositions [45], and solution of ordinary and partial differential equa-
 43 tions [8, 34, 41, 52].

44 In the context of linear solvers, probabilistic solvers posit a *prior distribution* rep-
 45 resenting initial epistemic uncertainty about a quantity of interest, which can be the
 46 solution [1, 7, 9, 53] or the matrix inverse [1, 2, 24]. They then condition on the finite
 47 amount of information obtained during m iterations to produce a *posterior distribu-*
 48 *tion* that reflects the reduced uncertainty [9, Section 1.2], [42]. The interpretation of
 49 CG as a probabilistic solver was pioneered in the context of optimization [24], followed
 50 by the development of the *Bayesian Conjugate Gradient method (BayesCG)* [9] as a
 51 general purpose solver in statistics. However, current versions of BayesCG have two
 52 drawbacks: they are computationally expensive; and their posterior distributions do
 53 not model the uncertainty accurately.

54 **1.1. Contributions and outline.** We propose an efficient uncertainty-aware
 55 CG implementation in the form of BayesCG (Algorithm 3.1), and establish its proper
 56 foundation within probabilistic numerics (sections 2 and 3).

57 We design a new *Krylov prior* distribution for BayesCG, which is motivated by
 58 the *Krylov subspace prior* [9, section 4.1], which is a *non-singular* structured prior
 59 based on Krylov spaces, whose posterior distributions are expensive and not always
 60 meaningful. In contrast, our new Krylov prior is generally singular, depends on quan-
 61 tities computed by CG, and produces low-rank posteriors that lend themselves to
 62 efficient sampling in downstream computations. We proceed in two steps.

- 63 1. Extension of BayesCG to singular prior covariances (section 2).

64 We show that under reasonable assumptions, the theoretical and computa-
 65 tional properties of BayesCG from [9] extend to prior covariances that are
 66 singular. This extension to singular priors paves the way for an efficient
 67 BayesCG implementation that produces meaningful posteriors. Auxiliary re-
 68 sults and technical proofs are postponed to the end (Appendices A and B).

- 69 2. Introduction of the new Krylov prior and its properties (section 3).

70 This singular prior covariance exploits structure and adapts to BayesCG,
 71 with posteriors whose means are identical to the corresponding CG iterates,
 72 and whose covariances describe a realistic level of uncertainty. The posterior
 73 covariances are maintained in factored form, and are therefore highly accurate
 74 and easy to approximate, as confirmed by numerical experiments (section 4).

75 **1.2. Notation.** Bold uppercase letters, like \mathbf{A} , represent matrices, with \mathbf{I} den-
 76 noting the identity. The Moore-Penrose inverse of \mathbf{A} is \mathbf{A}^\dagger . Bold lowercase letters,
 77 like \mathbf{x}_* , represent vectors; italic lowercase letters, like α , scalars; and italic uppercase
 78 letters, like X_0 , random variables. A multivariate Gaussian distribution with mean \mathbf{x}
 79 and covariance Σ is denoted by $\mathcal{N}(\mathbf{x}, \Sigma)$, and $X \sim \mathcal{N}(\mathbf{x}, \Sigma)$ is a Gaussian random
 80 variable. We assume exact arithmetic throughout the theoretical sections 2 and 3.

81 **2. Introduction to BayesCG with singular priors.** We extend the appli-
 82 cability of BayesCG from definite to semi-definite prior covariances, and discuss the
 83 theory (section 2.1), recursive computation of posterior distributions (section 2.2),
 84 and choices for prior distributions (section 2.3).

85 **2.1. Theoretical properties of BayesCG under singular priors.** We derive
 86 expressions for the BayesCG posterior means and covariances under singular priors
 87 (Theorem 2.1), express the posteriors in terms of projectors (Theorem 2.4), and estab-
 88 lish the optimality of the posterior means (Theorem 2.6). The proofs are analogous
 89 to earlier proofs for non-singular priors in [1, 9], and relegated to Appendix A and
 90 the supplement.

91 BayesCG computes posterior distributions $\mathcal{N}(\mathbf{x}_m, \Sigma_m)$ by conditioning the prior
 92 $\mathcal{N}(\mathbf{x}_0, \Sigma_0)$ on information from $m \leq n$ linearly independent search directions \mathbf{S}_m .
 93 Specifically, the posterior is the distribution of the random variable $X \sim \mathcal{N}(\mathbf{x}_0, \Sigma_0)$
 94 conditioned on the random variable $Y = \mathbf{S}_m^T \mathbf{A} X$ taking the value $\mathbf{S}_m^T \mathbf{A} \mathbf{x}_*$. The
 95 conditioning relies on two properties of Gaussian distributions:

- 96 (i) *Stability*: linear transformations of Gaussians remain Gaussian [39, Section 1.2].
 97 (ii) *Conjugacy*: posteriors from Gaussian priors conditioned under linear information
 98 remain Gaussian [51, Theorem 6.20].

99 We start with the extension of BayesCG to singular priors.

100 **THEOREM 2.1** (Extension of [9, Proposition 1]). *Let $\mathcal{N}(\mathbf{x}_0, \Sigma_0)$ be a prior with*
 101 *a symmetric positive semi-definite covariance $\Sigma_0 \in \mathbb{R}^{n \times n}$. Let $m \leq \text{rank}(\Sigma_0)$, and*
 102 *let the matrix of search directions $\mathbf{S}_m \equiv [\mathbf{s}_1 \ \cdots \ \mathbf{s}_m] \in \mathbb{R}^{n \times m}$ have linearly inde-*
 103 *pendent columns so that $\Lambda_m \equiv \mathbf{S}_m^T \mathbf{A} \Sigma_0 \mathbf{A} \mathbf{S}_m$ is non-singular. Then the BayesCG*
 104 *posterior $\mathcal{N}(\mathbf{x}_m, \Sigma_m)$ has mean and covariance*

$$105 \quad (2.1) \quad \mathbf{x}_m = \mathbf{x}_0 + \Sigma_0 \mathbf{A} \mathbf{S}_m \Lambda_m^{-1} \mathbf{S}_m^T (\mathbf{b} - \mathbf{A} \mathbf{x}_0)$$

$$106 \quad (2.2) \quad \Sigma_m = \Sigma_0 - \Sigma_0 \mathbf{A} \mathbf{S}_m \Lambda_m^{-1} \mathbf{S}_m^T \mathbf{A} \Sigma_0.$$

108 *Proof.* See supplement. □

109 **REMARK 2.2.** *Theorem 2.1 requires the existence of search directions that produce*
 110 *a nonsingular Λ_m , and the purpose this theorem is to derive an expression for how*
 111 *to compute the posterior distribution resulting from any valid set of search directions.*
 112 *Section 2.2 presents the recursive computation of search directions that make Λ_m non-*
 113 *singular, while the supplement presents an example of a non-recursive construction.*

114 Next we derive explicit expressions for the posterior covariances in terms of or-
 115 thogonal projectors onto $\text{range}(\Sigma_0 \mathbf{A} \mathbf{S}_m)$. To this end we exploit the close relation
 116 between Gaussian conditioning and orthogonal projections [1, Section 3]; and gener-
 117 alize the notion of projector [48, page 111] to semi-definite inner products to allow for
 118 singular priors Σ_0 ,

119 **DEFINITION 2.3** ([28, section 0.6.1]). *Let $\mathbf{B} \in \mathbb{R}^{n \times n}$ be symmetric positive semi-*
 120 *definite, and $\mathbf{P} \in \mathbb{R}^{n \times n}$. If $\mathbf{P}^2 = \mathbf{P}$ and $(\mathbf{B}\mathbf{P})^T = \mathbf{B}\mathbf{P}$, then \mathbf{P} is a \mathbf{B} -orthogonal*
 121 *projector, with $(\mathbf{I} - \mathbf{P})^T \mathbf{B} \mathbf{P} = \mathbf{0}$.*

122 Now we are ready to express the posterior distributions in Theorem 2.1 in terms
 123 of Σ_0^\dagger -orthogonal projectors.

124 **THEOREM 2.4** (Extension of [10, Proposition 3]). *Under the assumptions of*
 125 *Theorem 2.1*

$$126 \quad (2.3) \quad \mathbf{P}_m \equiv \Sigma_0 \mathbf{A} \mathbf{S}_m \Lambda_m^{-1} \mathbf{S}_m^T \mathbf{A} \Sigma_0 \Sigma_0^\dagger$$

127 *is a Σ_0^\dagger -orthogonal projector onto $K_m \equiv \text{range}(\Sigma_0 \mathbf{A} \mathbf{S}_m)$.*

If additionally $\mathbf{x}_* - \mathbf{x}_0 \in \text{range}(\boldsymbol{\Sigma}_0)$, then the posterior satisfies

$$\begin{aligned}\mathbf{x}_m &= (\mathbf{I} - \mathbf{P}_m)\mathbf{x}_0 + \mathbf{P}_m\mathbf{x}_* \\ \boldsymbol{\Sigma}_m &= (\mathbf{I} - \mathbf{P}_m)\boldsymbol{\Sigma}_0, \quad \mathbf{P}_m\boldsymbol{\Sigma}_m = \mathbf{0}.\end{aligned}$$

Proof. See Appendix A. \square

Theorem 2.4 expresses the posterior mean \mathbf{x}_m as the sum of two projections: the projection of the solution \mathbf{x}_* onto $\text{range}(\mathbf{P}_m)$, and the projection of the prior mean \mathbf{x}_0 onto the complementary space $\text{range}(\mathbf{P}_m)^\perp$. As for the posterior covariance $\boldsymbol{\Sigma}_m$, it is the projection of the prior covariance $\boldsymbol{\Sigma}_0$ onto the complementary space $\text{range}(\mathbf{P}_m)^\perp$.

REMARK 2.5. *Theorem 2.4 implies that $\mathbf{P}_m\mathbf{x}_m = \mathbf{P}_m\mathbf{x}_*$ and $\mathbf{P}_m\boldsymbol{\Sigma}_m\mathbf{P}_m^T = \mathbf{0}$. As a consequence, if $X \sim \mathcal{N}(\mathbf{x}_m, \boldsymbol{\Sigma}_m)$, then the distribution of $\mathbf{P}_m(X - \mathbf{x}_*)$ is Gaussian with mean $\mathbf{P}_m\mathbf{x}_m - \mathbf{P}_m\mathbf{x}_* = \mathbf{0}$ and covariance $\mathbf{P}_m\boldsymbol{\Sigma}_m\mathbf{P}_m^T = \mathbf{0}$. Thus, within $\text{range}(\mathbf{P}_m)$, there is no uncertainty in our knowledge of \mathbf{x}_* . We can interpret the posterior as a conjecture about the unknown location of \mathbf{x}_* in the complementary subspace $\text{range}(\mathbf{P}_m)^\perp$.*

Theorem 2.4 implies the following optimality for the posterior mean: It is the vector closest to the solution \mathbf{x}_* in the affine space $\mathbf{x}_0 + K_m$, with K_m as in Theorem 2.1.

THEOREM 2.6 (Extension of [1, Proposition 4]). *Under all the assumptions of Theorem 2.4, the posterior mean satisfies*

$$(2.4) \quad \mathbf{x}_m = \arg \min_{\mathbf{x} \in \mathbf{x}_0 + K_m} (\mathbf{x}_* - \mathbf{x})^T \boldsymbol{\Sigma}_0^\dagger (\mathbf{x}_* - \mathbf{x}).$$

Additionally, $(\mathbf{x}_* - \mathbf{x}_m)^T \boldsymbol{\Sigma}_0^\dagger (\mathbf{x}_* - \mathbf{x}_m) = 0$ if and only if $\mathbf{x}_m = \mathbf{x}_*$.

Proof. See Appendix A. \square

Theorems 2.1, 2.4, and 2.6 assume that the search directions are chosen so that $\boldsymbol{\Lambda}_m$ is non-singular. The additional assumption $\mathbf{x}_* - \mathbf{x}_0 \in \text{range}(\boldsymbol{\Sigma}_0)$ in Theorems 2.4 and 2.6 guarantees this nonsingularity for the specific search directions computed by BayesCG, as will be shown in Theorem 2.11.

2.2. Recursive computation of BayesCG posteriors under singular priors. We extend the recursions for posterior distributions under nonsingular prior covariances in [9] to singular ones, and present three results for the efficient implementation of BayesCG: New recursions for the posterior covariances (Theorem 2.7) and the search directions (Theorem 2.8); and a proof that the search directions are well-defined (Theorem 2.11).

The residuals of the posterior means are defined as

$$(2.5) \quad \mathbf{r}_m \equiv \mathbf{b} - \mathbf{A}\mathbf{x}_m, \quad 0 \leq m.$$

THEOREM 2.7 (Extension of Proposition 6 in [9]). *Under the assumptions of Theorem 2.1 if, in addition, the search directions \mathbf{S}_m are $\mathbf{A}\boldsymbol{\Sigma}_0\mathbf{A}$ -orthogonal, then the posterior means and covariances admit the recursions*

$$(2.6) \quad \mathbf{x}_j = \mathbf{x}_{j-1} + \frac{\boldsymbol{\Sigma}_0\mathbf{A}\mathbf{s}_j (\mathbf{s}_j^T \mathbf{r}_{j-1})}{\mathbf{s}_j^T \mathbf{A}\boldsymbol{\Sigma}_0\mathbf{A}\mathbf{s}_j}, \quad 1 \leq j \leq m,$$

and

$$(2.7) \quad \boldsymbol{\Sigma}_j = \boldsymbol{\Sigma}_{j-1} - \frac{\boldsymbol{\Sigma}_0\mathbf{A}\mathbf{s}_j (\boldsymbol{\Sigma}_0\mathbf{A}\mathbf{s}_j)^T}{\mathbf{s}_j^T \mathbf{A}\boldsymbol{\Sigma}_0\mathbf{A}\mathbf{s}_j}, \quad 1 \leq j \leq m.$$

170 *Proof.* See Appendix A. □

171 The denominators $(\mathbf{\Lambda}_m)_{jj} = \mathbf{s}_j^T \mathbf{A} \mathbf{\Sigma}_0 \mathbf{A} \mathbf{s}_j$ in (2.6) and (2.7) are non-zero because
 172 Theorem 2.1 assumes that $\mathbf{\Lambda}_m$ is non-singular.

173 Next is a Lanczos-like recurrence for the $\mathbf{A} \mathbf{\Sigma}_0 \mathbf{A}$ -orthogonal search directions from
 174 [9, Proposition 7].

175 THEOREM 2.8 ([9, Proposition 7] and [11, Proof of Proposition 7, Proposition S4,
 176 and Section S2]). *If the search directions*

$$177 \quad (2.8) \quad \mathbf{s}_1 = \mathbf{r}_0 \neq \mathbf{0}, \quad \mathbf{s}_j = \mathbf{r}_{j-1} - \frac{\mathbf{r}_{j-1}^T \mathbf{r}_{j-1}}{\mathbf{r}_{j-2}^T \mathbf{r}_{j-2}} \mathbf{s}_{j-1}, \quad 2 \leq j \leq m,$$

178 *satisfy the assumptions of Theorem 2.1, then they are an $\mathbf{A} \mathbf{\Sigma}_0 \mathbf{A}$ -orthogonal basis for*
 179 *the Krylov space*

$$180 \quad (2.9) \quad \mathcal{K}_m(\mathbf{A} \mathbf{\Sigma}_0 \mathbf{A}, \mathbf{r}_0) \equiv \text{span}\{\mathbf{r}_0, \mathbf{A} \mathbf{\Sigma}_0 \mathbf{A} \mathbf{r}_0, \dots, (\mathbf{A} \mathbf{\Sigma}_0 \mathbf{A})^{m-1} \mathbf{r}_0\},$$

181 *while the residuals $\mathbf{r}_0, \dots, \mathbf{r}_{m-1}$ are an orthogonal basis for $\mathcal{K}_m(\mathbf{A} \mathbf{\Sigma}_0 \mathbf{A}, \mathbf{r}_0)$.*

182 The maximal number of search directions in (2.8) can be less than n , because
 183 they are a basis for the Krylov subspace $\mathcal{K}_m(\mathbf{A} \mathbf{\Sigma}_0 \mathbf{A}, \mathbf{r}_0)$ whose maximal dimension
 184 can be less than n .

185 DEFINITION 2.9 (Section 2 in [3], Definition 4.2.1 in [31]). *Let $\mathbf{B} \in \mathbb{R}^{n \times n}$ be*
 186 *symmetric positive semi-definite and let $\mathbf{w} \in \mathbb{R}^n$ be a non-zero vector. The grade*
 187 *of \mathbf{w} with respect to \mathbf{B} , or the invariance index for (\mathbf{B}, \mathbf{w}) is the maximal dimension*
 188 *$1 \leq K \leq n$ of the Krylov space,*

$$189 \quad \mathcal{K}_K(\mathbf{B}, \mathbf{w}) = \mathcal{K}_{K+i}(\mathbf{B}, \mathbf{w}), \quad i \geq 1.$$

190 REMARK 2.10. *In Theorem 2.8, if K is the grade of \mathbf{r}_0 with respect to $\mathbf{A} \mathbf{\Sigma}_0 \mathbf{A}$,*
 191 *then $\mathbf{s}_{K+1} = \mathbf{0}$, $\mathbf{r}_K = \mathbf{0}$, while $\mathbf{s}_j \neq \mathbf{0}$ and $\mathbf{r}_{j-1} \neq \mathbf{0}$ for $1 \leq j \leq K$. Additionally,*
 192 *$K \leq \text{rank}(\mathbf{\Sigma}_0)$.*

193 In the following theorem, we show that with the additional assumption that $\mathbf{x}_* -$
 194 $\mathbf{x}_0 \in \text{range}(\mathbf{\Sigma}_0)$, the $\mathbf{A} \mathbf{\Sigma}_0 \mathbf{A}$ -orthogonal search directions from Theorem 2.8 satisfy
 195 the assumptions of Theorem 2.1.

196 THEOREM 2.11. *Let $\mathcal{N}(\mathbf{x}_0, \mathbf{\Sigma}_0)$ be a prior with symmetric positive semi-definite*
 197 *$\mathbf{\Sigma}_0 \in \mathbb{R}^{n \times n}$, K the grade of \mathbf{r}_0 with respect to $\mathbf{A} \mathbf{\Sigma}_0 \mathbf{A}$, and $m \leq K$. If $\mathbf{x}_* - \mathbf{x}_0 \in$*
 198 *$\text{range}(\mathbf{\Sigma}_0)$, then the search directions from Theorem 2.1 produce a nonsingular $\mathbf{\Lambda}_m$,*
 199 *and \mathbf{S}_m is $\mathbf{A} \mathbf{\Sigma}_0 \mathbf{A}$ -orthogonal.*

200 *Proof.* Recursive computation of the BayesCG posteriors requires the search di-
 201 rections $\mathbf{S}_m = [\mathbf{s}_1 \ \dots \ \mathbf{s}_m]$ to be $\mathbf{A} \mathbf{\Sigma}_0 \mathbf{A}$ -orthogonal, so that $\mathbf{\Lambda}_m = \mathbf{S}_m^T \mathbf{A} \mathbf{\Sigma}_0 \mathbf{A} \mathbf{S}_m$
 202 is diagonal [9, Section 2.3]. Furthermore, if $\mathbf{s}_j \notin \ker(\mathbf{\Sigma}_0 \mathbf{A})$, $1 \leq j \leq m$, then $\mathbf{\Lambda}_m$ has
 203 non-zero diagonal elements and is nonsingular.

204 In the following induction proof we show that the search directions are $\mathbf{A} \mathbf{\Sigma}_0 \mathbf{A}$ -
 205 orthogonal and that $\mathbf{s}_i \notin \ker(\mathbf{\Sigma}_0 \mathbf{A})$ and $\mathbf{s}_i \neq \mathbf{0}$, $1 \leq i \leq m$. Since \mathbf{A} and $\mathbf{\Sigma}_0$ are
 206 symmetric, $\ker(\mathbf{\Sigma}_0 \mathbf{A}) = \ker(\mathbf{\Sigma}_0^T \mathbf{A}^T) = \ker((\mathbf{A} \mathbf{\Sigma}_0)^T)$ is the orthogonal complement
 207 of $\text{range}(\mathbf{A} \mathbf{\Sigma}_0)$ in \mathbb{R}^n . Therefore, we can show $\mathbf{s}_i \notin \ker(\mathbf{\Sigma}_0 \mathbf{A})$ by showing $\mathbf{s}_i \in$
 208 $\text{range}(\mathbf{A} \mathbf{\Sigma}_0)$ and $\mathbf{s}_i \neq \mathbf{0}$, $1 \leq i \leq m$.

209 By assumption $m \leq K$, so Remark 2.10 implies $\mathbf{r}_i \neq \mathbf{0}$, $1 \leq i \leq m - 1$.

210 *Induction basis.* The assumption $\mathbf{x}_* - \mathbf{x}_0 \in \text{range}(\boldsymbol{\Sigma}_0)$ implies

$$211 \quad \mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0 = \mathbf{A}(\mathbf{x}_* - \mathbf{x}_0) \in \text{range}(\mathbf{A}\boldsymbol{\Sigma}_0).$$

212 Thus $\mathbf{s}_1 = \mathbf{r}_0 \in \text{range}(\mathbf{A}\boldsymbol{\Sigma}_0)$, and $\mathbf{r}_0 \neq \mathbf{0}$ by assumption. Thus $\mathbf{s}_1 \neq \mathbf{0}$, $\mathbf{s}_1 \notin \ker(\boldsymbol{\Sigma}_0\mathbf{A})$,
213 and $\boldsymbol{\Lambda}_1 = \mathbf{s}_1^T \mathbf{A}\boldsymbol{\Sigma}_0\mathbf{A}\mathbf{s}_1 \neq 0$.

214 *Induction hypothesis.* Assume that $\mathbf{s}_i, \mathbf{r}_i \in \text{range}(\mathbf{A}\boldsymbol{\Sigma}_0)$, $\mathbf{s}_i, \mathbf{r}_i \neq \mathbf{0}$, and $\boldsymbol{\Lambda}_i$ is
215 nonsingular, $1 \leq i \leq m-1$. This, along with Theorem 2.8 implies that $\mathbf{s}_1, \dots, \mathbf{s}_{m-1}$
216 are $\mathbf{A}\boldsymbol{\Sigma}_0\mathbf{A}$ -orthogonal so that $\boldsymbol{\Lambda}_{m-1}$ is a diagonal matrix.

217 *Induction step.* Applying the induction hypothesis $\mathbf{s}_{m-1}, \mathbf{r}_{m-1} \in \text{range}(\mathbf{A}\boldsymbol{\Sigma}_0)$ to
218 (2.8) gives

$$219 \quad (2.10) \quad \mathbf{s}_m = \mathbf{r}_{m-1} - \frac{\mathbf{r}_{m-1}^T \mathbf{r}_{m-1}}{\mathbf{r}_{m-2}^T \mathbf{r}_{m-2}} \mathbf{s}_{m-1}.$$

221 Hence $\mathbf{s}_m \in \text{range}(\mathbf{A}\boldsymbol{\Sigma}_0)$. Multiply (2.10) on the left by \mathbf{r}_{m-1}^T and insert $\mathbf{s}_{m-1}^T \mathbf{r}_{m-1} =$
222 0 from Lemma B.1 into the last summand to get $\mathbf{r}_{m-1}^T \mathbf{s}_m = \mathbf{r}_{m-1}^T \mathbf{r}_{m-1}$, where $\mathbf{r}_{m-1} \neq$
223 0 implies $\mathbf{s}_m \neq \mathbf{0}$. Then $\mathbf{s}_m \in \text{range}(\mathbf{A}\boldsymbol{\Sigma}_0)$ and $\mathbf{s}_m \neq \mathbf{0}$ imply $\mathbf{s}_m \notin \ker(\boldsymbol{\Sigma}_0\mathbf{A})$.

224 The induction hypothesis, Theorem 2.8, and (2.10) imply that the search di-
225 rections $\mathbf{s}_1, \dots, \mathbf{s}_m$ are non-zero and $\mathbf{A}\boldsymbol{\Sigma}_0\mathbf{A}$ -orthogonal. Thus $\boldsymbol{\Lambda}_m$ is nonsingular
226 diagonal, which implies that $\mathbf{s}_i \notin \ker(\boldsymbol{\Sigma}_0\mathbf{A})$, $1 \leq i \leq m$; and with Lemma A.1 that
227 $\mathbf{x}_* - \mathbf{x}_m \in \text{range}(\boldsymbol{\Sigma}_0)$, thus $\mathbf{r}_m = \mathbf{A}(\mathbf{x}_* - \mathbf{x}_m) \in \text{range}(\mathbf{A}\boldsymbol{\Sigma}_0)$. \square

228 **REMARK 2.12.** *The assumption $\mathbf{x}_* - \mathbf{x}_0 \in \text{range}(\boldsymbol{\Sigma}_0)$ in Theorem 2.11, which*
229 *holds automatically if the prior covariance $\boldsymbol{\Sigma}_0$ is nonsingular, is required to guarantee*
230 *the nonsingularity of the diagonal matrices $\boldsymbol{\Lambda}_m$.*

231 *The statistical interpretation of the assumption $\mathbf{x}_* - \mathbf{x}_0 \in \text{range}(\boldsymbol{\Sigma}_0)$ is that the*
232 *solution \mathbf{x}_* must live in the support of the prior, that is, in the subspace of \mathbb{R}^n where*
233 *the probability density function of $\mathcal{N}(\mathbf{x}_0, \boldsymbol{\Sigma}_0)$ is nonzero.*

234 Theorems 2.7, 2.8, and 2.11 form the basis for the BayesCG Algorithm 2.1, which
235 differs from the original BayesCG [9, Algorithm 1] only in the computation of the
236 posterior covariances as a sequence of rank-1 downdates rather than just a single
237 rank- m downdate at the end. Algorithm 2.1 is a Krylov space method; for nonsingular
238 priors $\boldsymbol{\Sigma}_0$ this was established in [9, Section 3], while for singular priors this follows
239 from (2.9) and Theorem 2.6. To show the similarity of BayesCG Algorithm 2.1 to
240 CG, we present the most common implementation of CG in Algorithm 2.2; it is the
241 original version due to Hestenes and Stiefel [26, Section 3].

242 The posterior means in Algorithm 2.1 are closely related to the CG iterates in
243 Algorithm 2.2. In the special case $\boldsymbol{\Sigma}_0 = \mathbf{A}^{-1}$, the BayesCG posterior means are iden-
244 tical to the CG iterates [9, Section 2.3]. The relationship between CG and BayesCG
245 is discussed further in [5, 9, 10, 11, 30], and the results are summarized in the sup-
246 plement.

247 **2.3. Choice of BayesCG prior distribution.** The mean \mathbf{x}_0 in the prior
248 $\mathcal{N}(\mathbf{x}_0, \boldsymbol{\Sigma}_0)$ corresponds to the initial guess in CG, while the covariance $\boldsymbol{\Sigma}_0$ can be
249 any symmetric positive semi-definite matrix that satisfies $\mathbf{x}_* - \mathbf{x}_0 \in \text{range}(\boldsymbol{\Sigma}_0)$. Non-
250 singular priors examined in [9, Section 4.1] include

- 251 \bullet Inverse prior $\boldsymbol{\Sigma}_0 = \mathbf{A}^{-1}$: The posterior means in Algorithm 2.1 are equal to
252 the CG iterates.
- 253 \bullet Natural prior $\boldsymbol{\Sigma}_0 = \mathbf{A}^{-2}$: The posterior means in Algorithm 2.1 converge in
254 a single iteration.

Algorithm 2.1 Bayesian Conjugate Gradient Method (BayesCG)

```

1: Input: spd  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{b} \in \mathbb{R}^n$ ,  $\mathbf{x}_0 \in \mathbb{R}^n$ 
2:     spds  $\Sigma_0 \in \mathbb{R}^{n \times n}$  so that  $\mathbf{x}_* - \mathbf{x}_0 \in \text{range}(\Sigma_0)$ 
3:  $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$  ▷ define initial values
4:  $\mathbf{s}_1 = \mathbf{r}_0$ 
5:  $m = 0$ 
6: while not converged do ▷ iterate through BayesCG Recursions
7:      $m = m + 1$ 
8:      $\alpha_m = (\mathbf{r}_{m-1}^T \mathbf{r}_{m-1}) / (\mathbf{s}_m^T \mathbf{A} \Sigma_0 \mathbf{A} \mathbf{s}_m)$ 
9:      $\mathbf{x}_m = \mathbf{x}_{m-1} + \alpha_m \Sigma_0 \mathbf{A} \mathbf{s}_m$ 
10:     $\Sigma_m = \Sigma_{m-1} - \Sigma_0 \mathbf{A} \mathbf{s}_m (\Sigma_0 \mathbf{A} \mathbf{s}_m)^T / (\mathbf{s}_m^T \mathbf{A} \Sigma_0 \mathbf{A} \mathbf{s}_m)$ 
11:     $\mathbf{r}_m = \mathbf{r}_{m-1} - \alpha_m \mathbf{A} \Sigma_0 \mathbf{A} \mathbf{s}_m$ 
12:     $\beta_m = (\mathbf{r}_m^T \mathbf{r}_m) / (\mathbf{r}_{m-1}^T \mathbf{r}_{m-1})$ 
13:     $\mathbf{s}_{m+1} = \mathbf{r}_m + \beta_m \mathbf{s}_m$ 
14: end while
15: Output:  $\mathbf{x}_m$ ,  $\Sigma_m$ 

```

Algorithm 2.2 Conjugate Gradient Method (CG)

```

1: Input: spd  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{b} \in \mathbb{R}^n$ ,  $\mathbf{x}_0 \in \mathbb{R}^n$ 
2:  $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$  ▷ define initial values
3:  $\mathbf{v}_1 = \mathbf{r}_0$ 
4:  $m = 0$ 
5: while not converged do ▷ iterate through CG Recursions
6:      $m = m + 1$ 
7:      $\gamma_m = (\mathbf{r}_{m-1}^T \mathbf{r}_{m-1}) / (\mathbf{v}_m^T \mathbf{A} \mathbf{v}_m)$ 
8:      $\mathbf{x}_m = \mathbf{x}_{m-1} + \gamma_m \mathbf{v}_m$ 
9:      $\mathbf{r}_m = \mathbf{r}_{m-1} - \gamma_m \mathbf{A} \mathbf{v}_m$ 
10:     $\delta_m = (\mathbf{r}_m^T \mathbf{r}_m) / (\mathbf{r}_{m-1}^T \mathbf{r}_{m-1})$ 
11:     $\mathbf{v}_{m+1} = \mathbf{r}_m + \delta_m \mathbf{v}_m$ 
12: end while
13: Output:  $\mathbf{x}_m$ 

```

- 255 • Identity prior $\Sigma_0 = \mathbf{I}$: The prior is easy to compute, but the posterior means
256 in Algorithm 2.1 converge slowly.
- 257 • Preconditioner prior $\Sigma_0 = (\mathbf{M}^T \mathbf{M})^{-1}$ where $\mathbf{M} \approx \mathbf{A}$: This prior approxi-
258 mates the natural prior.
- 259 • Krylov subspace prior Σ_0 : This prior is defined in terms of a basis for the
260 Krylov space $\mathcal{K}(\mathbf{A}, \mathbf{r}_0)$.

261 Figure 2.1 illustrates the convergence of posterior means and covariances from
262 Algorithm 2.1 under the priors $\Sigma_0 = \mathbf{A}^{-1}$ and $\Sigma_0 = \mathbf{I}$. In both cases the posterior
263 means converge faster than the posterior covariances, suggesting that the covariances
264 are unreasonably pessimistic about the size of the error $\mathbf{x}_* - \mathbf{x}_m$. Section 3.3 presents
265 a detailed discussion of the relation between the trace of the posterior covariance and
266 the error $\mathbf{x}_* - \mathbf{x}_m$ in the posterior means.

267 The example below presents a prior of minimal rank that comprises a maximal
268 amount of information.

269 EXAMPLE 2.13. *If $\mathbf{x}_0 \neq \mathbf{x}_*$, then $\Sigma_0 = (\mathbf{x}_* - \mathbf{x}_0)(\mathbf{x}_* - \mathbf{x}_0)^T$ is is a rank-one*

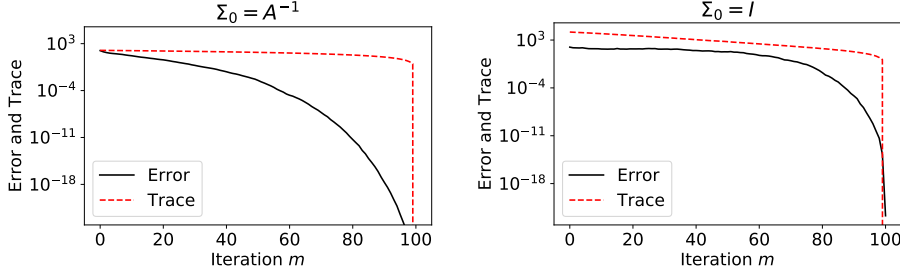


FIGURE 2.1. Convergence of BayesCG Algorithm 2.1 applied to the linear system in section 4.2 under different priors: inverse prior (left panel) and identity prior (right panel). Convergence of the means is displayed as $\|\mathbf{x}_* - \mathbf{x}_m\|_{\mathbf{A}}^2$, while convergence of the covariances is displayed as $\text{trace}(\mathbf{A}\Sigma_m)$.

270 covariance that satisfies $\mathbf{x}_* - \mathbf{x}_0 \in \text{range}(\Sigma_0)$. All rank-one prior covariances for
 271 BayesCG are multiples of this prior.

272 To see this, note that Theorem 2.11 and $\mathbf{A}^{-1}\mathbf{r}_0 = \mathbf{x}_* - \mathbf{x}_0$ imply termination of
 273 Algorithm 2.1 under this prior in a single iteration,

$$274 \quad \mathbf{x}_1 = \mathbf{x}_0 + \frac{1}{\mathbf{r}_0^T \mathbf{A} \underbrace{\mathbf{A}^{-1} \mathbf{r}_0 \mathbf{r}_0^T \mathbf{A}^{-1}}_{\Sigma_0} \mathbf{A} \mathbf{r}_0} \underbrace{\mathbf{A}^{-1} \mathbf{r}_0 \mathbf{r}_0^T \mathbf{A}^{-1}}_{\Sigma_0} \mathbf{A} \mathbf{r}_0 (\mathbf{r}_0^T \mathbf{r}_0) = \mathbf{x}_0 + \mathbf{x}_* - \mathbf{x}_0 = \mathbf{x}_*.$$

275

276 **3. Prior distributions informed by Krylov subspaces.** Motivated by the
 277 ‘Krylov subspace prior’ [9, section 4.1], we introduce a new ‘Krylov prior’ (section 3.1),
 278 derive expressions for the Krylov posteriors (section 3.2), ensure the Krylov posteriors
 279 accurately model uncertainty in \mathbf{x}_* (section 3.3), and develop a practical Krylov
 280 posterior and an efficient implementation of BayesCG as a uncertainty-aware version
 281 of CG (section 3.4).

282 **3.1. General Krylov prior.** We introduce our new Krylov prior (Definition 3.1)
 283 and show that the BayesCG Krylov subspace under the Krylov prior is identical to
 284 the CG Krylov subspace (Lemma 3.2). This Krylov prior is impractical because its
 285 computation amounts to the direct solution of (1.1), however it is the foundation for
 286 the efficient low-rank approximations in section 3.4.

287 The new Krylov prior is defined in terms of the maximal CG Krylov subspace
 288 $\mathcal{K}_K(\mathbf{A}, \mathbf{r}_0)$, where K is the grade of \mathbf{r}_0 with respect to \mathbf{A} (Definition 2.9). The \mathbf{A} -
 289 orthonormal versions of the search directions \mathbf{v}_m in Algorithm 2.2 are

$$290 \quad (3.1) \quad \tilde{\mathbf{v}}_m \equiv \mathbf{v}_m / \sqrt{\mathbf{v}_m^T \mathbf{A} \mathbf{v}_m}, \quad 1 \leq m \leq K.$$

291 As columns of

$$292 \quad (3.2) \quad \mathbf{V} \equiv [\tilde{\mathbf{v}}_1 \quad \cdots \quad \tilde{\mathbf{v}}_K] \in \mathbb{R}^{n \times K} \quad \text{with} \quad \mathbf{V}^T \mathbf{A} \mathbf{V} = \mathbf{I}_K$$

293 they represent an \mathbf{A} -orthonormal basis for $\text{range}(\mathbf{V}) = \mathcal{K}_K(\mathbf{A}, \mathbf{r}_0)$ [26, Theorem 5.1].

294 DEFINITION 3.1. The (general) Krylov prior is $\mathcal{N}(\mathbf{x}_0, \Gamma_0)$, where the mean \mathbf{x}_0 is
 295 an initial guess for \mathbf{x}_* , and the covariance matrix is

$$296 \quad (3.3) \quad \Gamma_0 \equiv \mathbf{V} \Phi \mathbf{V}^T \in \mathbb{R}^{n \times n}$$

297 where \mathbf{V} is as defined in (3.2) and $\mathbf{\Phi} \equiv \text{diag}(\phi_1 \ \phi_2 \ \cdots \ \phi_K) \in \mathbb{R}^{K \times K}$ with $\phi_i > 0$,
 298 $1 \leq i \leq K$. The Krylov prior is ‘general’ because the diagonal elements of $\mathbf{\Phi}$ are
 299 unspecified.

300 The results in this section and in section 3.2 are valid for any choice of posi-
 301 tive diagonal elements in $\mathbf{\Phi}$. A specific choice of diagonal elements is presented in
 302 section 3.3.

303 The Krylov prior covariance has $\text{rank}(\mathbf{\Gamma}_0) = K$ and is singular for $K < n$, hence
 304 the need for singular priors in section 2. Fortunately, $\mathbf{\Gamma}_0$ is a well-defined BayesCG
 305 prior, because it satisfies the crucial condition in Theorem 2.11,

$$306 \quad \mathbf{x}_* - \mathbf{x}_0 \in \mathcal{K}_K(\mathbf{A}, \mathbf{r}_0) = \text{range}(\mathbf{V}) = \text{range}(\mathbf{\Gamma}_0).$$

307 *Intuition.* We give two different interpretations of the decomposition (3.3).

- 308 1. Hermitian eigenvalue problem $\mathbf{A}^{1/2} \mathbf{\Gamma}_0 \mathbf{A}^{1/2} = \mathbf{W} \mathbf{\Phi} \mathbf{W}^T$, where $\mathbf{\Phi}$ contains the
 309 positive eigenvalues, and the eigenvector matrix $\mathbf{W} \equiv \mathbf{A}^{1/2} \mathbf{V}$ has orthonor-
 310 mal columns with $\mathbf{W}^T \mathbf{W} = \mathbf{I}_K$.
- 311 2. Non-Hermitian eigenvalue problem $\mathbf{\Gamma}_0 \mathbf{A} \mathbf{V} = \mathbf{V} \mathbf{\Phi}$ with eigenvalues and eigen-
 312 vectors

$$313 \quad (3.4) \quad \mathbf{\Gamma}_0 \mathbf{A} \tilde{\mathbf{v}}_m = \phi_m \tilde{\mathbf{v}}_m, \quad 1 \leq m \leq K.$$

314 This is the property to be exploited in section 3.2.

315 We show that the BayesCG Krylov subspace under the Krylov prior is identical
 316 to the CG Krylov subspace.

317 LEMMA 3.2. *If $\mathbf{\Gamma}_0$ is the Krylov prior in Definition 3.1, then*

$$318 \quad \mathcal{K}_m(\mathbf{A}, \mathbf{r}_0) = \mathcal{K}_m(\mathbf{A} \mathbf{\Gamma}_0 \mathbf{A}, \mathbf{r}_0), \quad 1 \leq m \leq K.$$

319 *Consequently, K is also the grade of \mathbf{r}_0 with respect to $\mathbf{A} \mathbf{\Gamma}_0 \mathbf{A}$ is K .*

320 *Proof.* An induction proof shows that the Krylov subspaces are the same for the
 321 first K dimensions. Then we prove that the grade of \mathbf{r}_0 with respect to $\mathbf{A} \mathbf{\Sigma} \mathbf{A}$ is K .

322 *Induction basis.* Since one-dimensional Krylov subspaces are independent of the
 323 matrix,

$$324 \quad \mathcal{K}_1(\mathbf{A}, \mathbf{r}_0) = \text{span}\{\mathbf{r}_0\} = \mathcal{K}_1(\mathbf{A} \mathbf{\Gamma}_0 \mathbf{A}, \mathbf{r}_0).$$

325 *Induction hypothesis.* Assume that

$$326 \quad \mathcal{K}_i(\mathbf{A}, \mathbf{r}_0) = \mathcal{K}_i(\mathbf{A} \mathbf{\Gamma}_0 \mathbf{A}, \mathbf{r}_0), \quad 1 \leq i \leq m - 1.$$

327 With $\mathbf{V}_{1:m-1} = [\tilde{\mathbf{v}}_1 \ \tilde{\mathbf{v}}_2 \ \cdots \ \tilde{\mathbf{v}}_{m-1}]$ in (3.2) this implies

$$328 \quad (3.5) \quad \text{range}(\mathbf{V}_{1:m-1}) = \mathcal{K}_{m-1}(\mathbf{A}, \mathbf{r}_0) = \mathcal{K}_{m-1}(\mathbf{A} \mathbf{\Gamma}_0 \mathbf{A}, \mathbf{r}_0).$$

329 *Induction step.* From (3.5) follow the expressions for the direct sums,

$$330 \quad (3.6) \quad \mathcal{K}_m(\mathbf{A}, \mathbf{r}_0) = \text{span}\{\mathbf{r}_0\} \oplus \text{range}(\mathbf{A} \mathbf{V}_{1:m-1})$$

$$331 \quad (3.7) \quad \mathcal{K}_m(\mathbf{A} \mathbf{\Gamma}_0 \mathbf{A}, \mathbf{r}_0) = \text{span}\{\mathbf{r}_0\} \oplus \text{range}(\mathbf{A} \mathbf{\Gamma}_0 \mathbf{A} \mathbf{V}_{1:m-1}).$$

332 Then (3.4) and the non-singularity of $\mathbf{\Phi}$ imply

$$333 \quad \text{range}(\mathbf{A} \mathbf{\Gamma}_0 \mathbf{A} \mathbf{V}_{1:m-1}) = \text{range}(\mathbf{A} \mathbf{V}_{1:m-1} \mathbf{\Phi}_{1:m-1}) = \text{range}(\mathbf{A} \mathbf{V}_{1:m-1}).$$

334 Combining this with (3.6) and (3.7) completes the induction,

$$\begin{aligned}
335 \quad \mathcal{K}_m(\mathbf{A}, \mathbf{r}_0) &= \text{span}\{\mathbf{r}_0\} \oplus \text{range}(\mathbf{A}\mathbf{V}_{1:m-1}) \\
336 \quad &= \text{span}\{\mathbf{r}_0\} \oplus \text{range}(\mathbf{A}\mathbf{\Gamma}_0\mathbf{A}\mathbf{V}_{1:m-1}) = \mathcal{K}_m(\mathbf{A}\mathbf{\Gamma}_0\mathbf{A}, \mathbf{r}_0).
\end{aligned}$$

338 *Maximal Krylov space dimension.* If K' is the grade of \mathbf{r}_0 with respect to $\mathbf{A}\mathbf{\Gamma}_0\mathbf{A}$,
339 then the induction implies

$$340 \quad K' \geq \dim(\mathcal{K}_K(\mathbf{A}\mathbf{\Sigma}_0\mathbf{A}, \mathbf{r}_0)) = \dim(\mathcal{K}_K(\mathbf{A}, \mathbf{r}_0)) = K.$$

341 On the other hand, $\text{rank}(\mathbf{A}\mathbf{\Gamma}_0\mathbf{A}) = K$ implies $K' \leq K$. Therefore $K' = K$. \square

342 **3.2. General Krylov posteriors.** We show (Theorem 3.3) that under the
343 Krylov prior, the BayesCG posteriors have means that are identical to the CG it-
344 erates, and covariances that can be factored as in Definition 3.1. This represents the
345 foundation for an efficient implementation of BayesCG (Remark 3.4).

346 Define appropriate submatrices of \mathbf{V} and $\mathbf{\Phi}$,

$$347 \quad (3.8) \quad \mathbf{V}_{i:j} \equiv [\tilde{\mathbf{v}}_i \quad \cdots \quad \tilde{\mathbf{v}}_j], \quad \mathbf{\Phi}_{i:j} \equiv \text{diag}(\phi_i \quad \cdots \quad \phi_j), \quad 1 \leq i < j \leq K.$$

348 In particular, $\mathbf{V} = \mathbf{V}_{1:K}$ and $\mathbf{\Phi} = \mathbf{\Phi}_{1:K}$.

349 **THEOREM 3.3.** *Let $\mathcal{N}(\mathbf{x}_0, \mathbf{\Gamma}_0)$ be the Krylov prior in Definition 3.1, and let*
350 *$\mathcal{N}(\mathbf{x}_m, \mathbf{\Gamma}_m)$ be the posteriors from BayesCG Algorithm 2.1, $1 \leq m \leq K$. Then the*
351 *posterior means \mathbf{x}_m are identical to the corresponding CG iterates in Algorithm 2.2,*
352 *and the posterior covariances can be factored as*

$$353 \quad (3.9) \quad \mathbf{\Gamma}_m = \mathbf{V}_{m+1:K}\mathbf{\Phi}_{m+1:K}(\mathbf{V}_{m+1:K})^T, \quad 1 \leq m < K,$$

354 and $\mathbf{\Gamma}_m = \mathbf{0}$ for $m = K$.

355 *Proof.* We first derive the equality of the posterior means, and then the factor-
356 izations of the covariances.

357 *Posterior means.* The idea is to show equality of the BayesCG posterior means
358 under Krylov and inverse priors since, per the discussion in [9, Section 2.3] and sec-
359 tion 2.3, BayesCG posterior means under the inverse prior are identical to CG iterates.

360 From Theorem 2.1, and the ‘equivalence’ of Algorithm 2.1 under $\mathbf{\Sigma}_0 = \mathbf{A}^{-1}$ and
361 Algorithm 2.2 follows that the BayesCG posterior means under the inverse prior are
362 equal to

$$363 \quad (3.10) \quad \mathbf{x}_m = \mathbf{x}_0 + \mathbf{V}_{1:m}\mathbf{V}_{1:m}^T\mathbf{r}_0.$$

364 Similarly, Theorem 2.1 implies that the BayesCG posterior under the Krylov prior
365 are equal to

$$366 \quad (3.11) \quad \mathbf{x}_m = \mathbf{x}_0 + \mathbf{\Gamma}_0\mathbf{A}\tilde{\mathbf{S}}_m(\tilde{\mathbf{S}}_m^T\mathbf{A}\mathbf{\Gamma}_0\mathbf{A}\tilde{\mathbf{S}}_m)^{-1}\tilde{\mathbf{S}}_m^T\mathbf{r}_0,$$

367 where the columns of $\tilde{\mathbf{S}}_m$ are the search directions from Algorithm 2.1 under the
368 Krylov prior. To show the equality of (3.10) and (3.11), we need to relate $\tilde{\mathbf{S}}_m$ and
369 $\mathbf{V}_{1:m}$ and then include the Krylov prior $\mathbf{\Gamma}_0$.

370 With the submatrices defined as in (3.8) we conclude from (3.2) and Lemma 3.2
371 that

$$372 \quad \text{range}(\tilde{\mathbf{S}}_m) = \mathcal{K}_m(\mathbf{A}\mathbf{\Gamma}_0\mathbf{A}, \mathbf{r}_0) = \text{range}(\mathbf{V}_{1:m}),$$

373 where the columns of $\tilde{\mathbf{S}}_m$ are $\mathbf{A}\Gamma_0\mathbf{A}$ -orthogonal. To show that the columns of $\mathbf{V}_{1:m}$
 374 are also $\mathbf{A}\Gamma_0\mathbf{A}$ -orthogonal, exploit the fact that they are \mathbf{A} -orthonormal and apply
 375 Definition 3.1,

$$376 \quad \mathbf{V}_{1:m}^T \mathbf{A}\Gamma_0 \mathbf{A} \mathbf{V}_{1:m} = \mathbf{V}_{1:m}^T \mathbf{A} \mathbf{V} \Phi \mathbf{V}^T \mathbf{A} \mathbf{V}_{1:m} = \Phi_{1:m},$$

377 which is a diagonal matrix. We have established that the columns of $\tilde{\mathbf{S}}_m$ and $\mathbf{V}_{1:m}$
 378 are $\mathbf{A}\Gamma_0\mathbf{A}$ -orthogonal, with respective leading columns being multiples of \mathbf{r}_0 , thus
 379 are $\mathbf{A}\Gamma_0\mathbf{A}$ -orthogonal bases of $\mathcal{K}_m(\mathbf{A}\Gamma_0\mathbf{A}, \mathbf{r}_0)$. Therefore the columns of $\mathbf{V}_{1:m}$ are
 380 multiples of the columns of $\tilde{\mathbf{S}}_m$. That is

$$381 \quad (3.12) \quad \tilde{\mathbf{S}}_m = \mathbf{V}_{1:m} \Delta$$

382 for some non-singular diagonal matrix $\Delta \in \mathbb{R}^{m \times m}$. Substitute (3.12) into the third
 383 interpretation (3.4) of the Krylov prior,

$$384 \quad \Gamma_0 \mathbf{A} \tilde{\mathbf{S}}_m = \Gamma_0 \mathbf{A} \mathbf{V}_{1:m} \Delta = \mathbf{V}_{1:m} \Phi_{1:m} \Delta$$

385 and this in turn into the second summand of (3.11). Then the non-singularity and
 386 diagonality of both Δ and Φ lead to the simplification

$$387 \quad (3.13) \quad \mathbf{x}_m = \mathbf{x}_0 + \mathbf{V}_{1:m} \Phi_{1:m} \Delta (\Delta \Phi_{1:m} \Delta)^{-1} \Delta \mathbf{V}_{1:m}^T \mathbf{r}_0 = \mathbf{x}_0 + \mathbf{V}_{1:m} \mathbf{V}_{1:m}^T \mathbf{r}_0,$$

389 which is (3.10).

390 *Posterior covariances.* Substituting (3.12) into Theorem 2.1 and simplifying as
 391 in (3.13) gives

$$392 \quad \Gamma_m = \Gamma_0 - \Gamma_0 \mathbf{A} \tilde{\mathbf{S}}_m (\tilde{\mathbf{S}}_m^T \mathbf{A} \Gamma_0 \mathbf{A} \tilde{\mathbf{S}}_m)^{-1} \tilde{\mathbf{S}}_m^T \mathbf{A} \Gamma_0 \\ 393 \quad = \mathbf{V} \Phi \mathbf{V}^T - \mathbf{V}_{1:m} \Phi_{1:m} \mathbf{V}_{1:m}^T = \mathbf{V}_{m+1:K} \Phi_{m+1:K} \mathbf{V}_{m+1:K}^T. \quad \square$$

395 **REMARK 3.4.** *Theorem 3.3 implies that the posteriors from BayesCG under the*
 396 *Krylov prior have means that can be computed with CG, and covariances can be main-*
 397 *tained in factored form without any arithmetic operations. This is the key to the*
 398 *efficient implementation of BayesCG in section 3.4.*

399 **3.3. Krylov posteriors that capture CG convergence.** We present a Krylov
 400 prior with specific diagonal elements (section 3.3.1), discuss the calibration of BayesCG
 401 under this prior (section 3.3.2) and its relation to existing CG error estimation theory
 402 (section 3.3.3).

403 **3.3.1. Specific Krylov prior.** We choose a specific diagonal matrix Φ for the
 404 Krylov prior (Definition 3.6), so that the Krylov posteriors accurately model the un-
 405 certainty in our knowledge of \mathbf{x}_* due to the error $\mathbf{x}_* - \mathbf{x}_m$. We derive error estimates
 406 from samples of the posteriors (Lemma 3.5) and then relate them to CG errors (The-
 407 orem 3.7).

408 Let us start with a general posterior distribution $\mathcal{N}(\mathbf{x}, \Sigma)$. If it indeed accurately
 409 modeled the uncertainty in \mathbf{x}_* due to the approximation error $\mathbf{x}_* - \mathbf{x}$, then we would
 410 expect the difference between samples of $\mathcal{N}(\mathbf{x}, \Sigma)$ and its posterior mean \mathbf{x} to be close
 411 to the actual error,

$$412 \quad (3.14) \quad \mathbb{E} [\|X - \mathbf{x}\|_{\mathbf{A}}^2] = \|\mathbf{x}_* - \mathbf{x}\|_{\mathbf{A}}^2 \quad \text{where } X \sim \mathcal{N}(\mathbf{x}, \Sigma).$$

413 The squared \mathbf{A} -norm error $\|X - \mathbf{x}\|_{\mathbf{A}}^2$ is a quadratic form, whose expected value has
 414 an explicit expression.

415 LEMMA 3.5. *If $X \sim \mathcal{N}(\mathbf{x}, \Sigma)$ is a Gaussian random variable with mean $\mathbf{x} \in \mathbb{R}^n$*
 416 *and symmetric positive semi-definite covariance $\Sigma \in \mathbb{R}^{n \times n}$, then*

$$417 \quad (3.15) \quad \mathbb{E} [\|X - \mathbf{x}\|_{\mathbf{A}}^2] = \text{trace}(\mathbf{A}\Sigma).$$

418 *Proof.* The proof relies on the expected value of a quadratic form in Appendix B.
 419 Set $Z \equiv X - \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ and apply Lemma B.2 to $Z^T \mathbf{A}Z$,

$$420 \quad \mathbb{E} [\|X - \mathbf{x}\|_{\mathbf{A}}^2] = \mathbb{E} [\|Z\|_{\mathbf{A}}^2] = \mathbb{E} [Z^T \mathbf{A}Z] = \text{trace}(\mathbf{A}\Sigma). \quad \square$$

422 Thus, $\text{trace}(\mathbf{A}\Sigma)$ has the potential to be an error indicator. We present a specific
 423 diagonal matrix for the Krylov prior Γ_0 in Definition 3.1, so that its posterior
 424 covariances produce meaningful error estimates $\text{trace}(\mathbf{A}\Gamma_m)$.

425 DEFINITION 3.6. *The (specific) Krylov prior is $\mathcal{N}(\mathbf{x}_0, \Gamma_0)$, where the mean \mathbf{x}_0 is*
 426 *an initial guess for \mathbf{x}_* , and the covariance matrix is*

$$427 \quad (3.16) \quad \Gamma_0 \equiv \mathbf{V}\Phi\mathbf{V}^T \in \mathbb{R}^{n \times n}$$

428 *where \mathbf{V} is defined in (3.2) and $\Phi \equiv \text{diag}(\phi_1 \ \phi_2 \ \dots \ \phi_K) \in \mathbb{R}^{K \times K}$ has diagonal*
 429 *elements*

$$430 \quad \phi_i = \gamma_i \|\mathbf{r}_{i-1}\|_2^2, \quad 1 \leq i \leq K,$$

432 *where $\gamma_i = \mathbf{r}_{i-1}^T \mathbf{r}_{i-1} / \mathbf{v}_i^T \mathbf{A} \mathbf{v}_i$ are the step sizes in line 7 of CG Algorithm 2.2.*

433 Now we show that the posterior covariances from BayesCG under the specific
 434 Krylov prior reproduce the CG error.

435 THEOREM 3.7. *Let $\mathcal{N}(\mathbf{x}_0, \Gamma_0)$ be the Krylov prior in Definition 3.6, and $\mathcal{N}(\mathbf{x}_m, \Gamma_m)$*
 436 *be the posteriors from BayesCG Algorithm 2.1, $1 \leq m \leq K$. Then*

$$437 \quad \text{trace}(\mathbf{A}\Gamma_m) = \|\mathbf{x}_* - \mathbf{x}_m\|_{\mathbf{A}}^2, \quad 1 \leq m \leq K.$$

438 *Proof.* Apply Lemma 3.5 to the specific Krylov prior in Definition 3.6. From the
 439 cyclic commutativity of the trace and \mathbf{A} -orthonormality of the columns of \mathbf{V} follows

$$440 \quad \text{trace}(\mathbf{A}\Gamma_m) = \text{trace}(\mathbf{A}\mathbf{V}_{m:K}\Phi_{m:K}(\mathbf{V}_{m:K})^T)$$

$$441 \quad (3.17) \quad = \text{trace}((\mathbf{V}_{m:K})^T \mathbf{A} \mathbf{V}_{m:K} \Phi_{m:K}) = \text{trace}(\Phi_{m:K}).$$

443 The diagonal matrix Φ for the specific Krylov prior in Definition 3.6 is chosen so that
 444 $\text{trace}(\Phi_{m:K}) = \|\mathbf{x}_* - \mathbf{x}_m\|_{\mathbf{A}}^2$. Remember that the reduction in the squared \mathbf{A} -norm
 445 error from iteration m to $m + d$ of Algorithm 2.2 equals [26, Theorem 6:1] and [31,
 446 Theorem 5.6.1]

$$447 \quad (3.18) \quad \|\mathbf{x}_* - \mathbf{x}_m\|_{\mathbf{A}}^2 - \|\mathbf{x}_* - \mathbf{x}_{m+d}\|_{\mathbf{A}}^2 = \sum_{i=m+1}^{m+d} \gamma_i \|\mathbf{r}_{i-1}\|_2^2, \quad 0 \leq m < m + d \leq K.$$

448 Setting $d = K - m$ gives $\mathbf{x}_K = \mathbf{x}_*$ and

$$449 \quad \|\mathbf{x}_* - \mathbf{x}_m\|_{\mathbf{A}}^2 = \sum_{i=m+1}^K \gamma_i \|\mathbf{r}_{i-1}\|_2^2, \quad 0 \leq m \leq K.$$

451 Combine this equality with (3.17) to conclude $\phi_i = \gamma_i \|\mathbf{r}_{i-1}\|_2^2$, $1 \leq i \leq K$. \square

452 Thus, the specific Krylov posteriors have covariances that converge at the same
 453 speed as their means.

454 **3.3.2. Calibration of BayesCG under the specific Krylov prior.** A prob-
 455 abilistic numerical linear solver is considered *calibrated* if its posterior distribution
 456 accurately models the uncertainty in \mathbf{x}_* due to the approximation error $\mathbf{x}_* - \mathbf{x}_m$.
 457 Calibration of general probabilistic methods is discussed in [6] and of linear solvers
 458 in [7]. We briefly discuss how Lemma 3.5 and Theorem 3.7 contribute to better
 459 calibration of BayesCG under the specific Krylov prior.

460 Previous probabilistic extensions of CG do not produce posteriors that accurately
 461 model the uncertainty in \mathbf{x}_* [1, Section 6.4], [9, Section 6.1], [53, Section 3]. For
 462 instance, Figure 2.1 illustrates that BayesCG under the priors $\Sigma_0 = \mathbf{A}^{-1}$ and $\Sigma_0 = \mathbf{I}$
 463 has errors $\|\mathbf{x}_* - \mathbf{x}_m\|_{\mathbf{A}}^2$ that converge faster than $\text{trace}(\mathbf{A}\Sigma_m)$. Furthermore, according
 464 to Lemma 3.5, the estimators $\text{trace}(\mathbf{A}\Sigma_m)$ from posterior samples are inaccurate and
 465 do not reflect the true error $\|\mathbf{x}_* - \mathbf{x}_m\|_{\mathbf{A}}^2$. In other words, the posteriors do not
 466 accurately model uncertainty in \mathbf{x}_* .

467 Our approach towards designing posteriors that accurately model the uncertainty
 468 in \mathbf{x}_* relies a judicious choice of the diagonal matrix Φ for the specific Krylov prior,
 469 so that sampling from the posteriors produces accurate error estimates. This can be
 470 viewed as a scaling of the posterior covariance that forces $\text{trace}(\Phi_{m:K}) = \|\mathbf{x}_* - \mathbf{x}_m\|_{\mathbf{A}}^2$.
 471 Alternative approaches for improving posteriors via scaling of the posterior covariances
 472 include [9, Section 4.2], [13, Section 7], and [53, Section 3]

473 Empirical evidence demonstrating that BayesCG under the specific Krylov prior
 474 produces posterior samples with accurate error estimates suggests but does not guar-
 475 antee that it accurately models the uncertainty in \mathbf{x}_* . A rigorous investigation of the
 476 calibration of BayesCG under the specific Krylov prior is the subject of a separate
 477 paper.

478 **3.3.3. Relation to CG error estimation.** The purpose of Lemma 3.5 is to
 479 motivate a choice of Φ so that BayesCG under the specific Krylov prior accurately
 480 models the uncertainty in \mathbf{x}_* due to the approximation error $\mathbf{x}_* - \mathbf{x}_m$.

481 Effective CG error estimation is a well researched area, with most effort focused
 482 on the absolute \mathbf{A} -norm error. One option [49] is to run d additional CG iterations
 483 and apply (3.18) to obtain the underestimate [49, Equation (4.9)],

$$484 \quad (3.19) \quad \sum_{i=m+1}^{m+d} \gamma_i \|\mathbf{r}_{i-1}\|_2^2 \leq \|\mathbf{x}_* - \mathbf{x}_m\|_{\mathbf{A}}^2.$$

485 The rationale is that the error after $m + d$ iterations has become negligible compared
 486 to the error after m iterations, especially in the case of fast convergence. The number
 487 of additional iterations d is usually called the ‘delay’ [37, Section 1], and larger values
 488 of d lead to more accurate error estimates.

489 The estimate (3.19) also coincides with the lower bound from Gaussian quadra-
 490 ture [49, Section 3]. Other lower and upper bounds for the \mathbf{A} -norm error based on
 491 quadrature formulas and tunable with a delay include [17, 18, 19, 35, 36, 37, 49, 50].

492 **3.4. Practical specific Krylov posteriors.** We define low rank approxima-
 493 tions of specific Krylov posterior covariances (Definition 3.8), and present an efficient
 494 CG-based implementation of BayesCG (Algorithm 3.1). It approximates the Krylov
 495 posteriors from delay iterations, thereby avoiding explicit computation of the Krylov
 496 prior, and inherits the fast convergence of CG.

497 The following low-rank approximations are based on the factored form of the
 498 Krylov posteriors in Theorem 3.3 and make use of the submatrices defined in (3.8).

499 DEFINITION 3.8. Let $\mathcal{N}(\mathbf{x}_0, \mathbf{\Gamma}_0)$ be the specific Krylov prior from Definition 3.6
500 with posteriors

$$501 \quad \mathbf{\Gamma}_m = \mathbf{V}_{m+1:K} \mathbf{\Phi}_{m+1:K} (\mathbf{V}_{m+1:K})^T, \quad 1 \leq m < K.$$

502 For $1 \leq d \leq K - m$, extract the leading rank- d submatrices from $\mathbf{V}_{m+1:K}$ and $\mathbf{\Phi}_{m+1:K}$,
503 and define the rank- d approximate Krylov posteriors as $\mathcal{N}(\mathbf{x}_m, \widehat{\mathbf{\Gamma}}_m)$ with

$$504 \quad (3.20) \quad \widehat{\mathbf{\Gamma}}_m \equiv \mathbf{V}_{m+1:m+d} \mathbf{\Phi}_{m+1:m+d} (\mathbf{V}_{m+1:m+d})^T.$$

505 REMARK 3.9. We view (3.20) as approximations of the posteriors resulting from
506 the full-rank prior. Instead, we could also view (3.20) as posteriors from rank- $(m+d)$
507 approximations of the prior $\mathcal{N}(\mathbf{x}_0, \widehat{\mathbf{\Gamma}}_0)$ with $\widehat{\mathbf{\Gamma}}_0 = \mathbf{V}_{1:m+d} \mathbf{\Phi}_{1:m+d} (\mathbf{V}_{1:m+d})^T$. This
508 interpretation of (3.20) is discussed in the supplement. However, from a practical
509 point of view, explicit computation of $\widehat{\mathbf{\Gamma}}_0$ is too expensive and it is not necessary.

510 Following the same argument as Theorem 3.7, one can express the underesti-
511 mate (3.19) for the CG error in terms of the posterior covariance,

$$512 \quad \text{trace}(\mathbf{A} \widehat{\mathbf{\Gamma}}_m) = \sum_{i=m+1}^{m+d} \gamma_i \|\mathbf{r}_{i-1}\|_2^2 \leq \|\mathbf{x}_* - \mathbf{x}_m\|_{\mathbf{A}}^2.$$

513 If the posterior distribution accurately models the uncertainty in the solution, then
514 we expect (3.14) to hold. This means the accuracy of the uncertainty from the ap-
515 proximate Krylov posterior is related to the accuracy of the underestimate (3.19).

516 Algorithm 3.1 represents an efficient computation of BayesCG under rank- d ap-
517 proximate Krylov posteriors, and consists of two loops²:

- 518 1. Run CG until convergence in iteration m and compute the posterior mean
519 \mathbf{x}_m
- 520 2. Run d additional CG iterations and compute the factors $\mathbf{V}_{m+1:m+d}$ and
521 $\mathbf{\Phi}_{m+1:m+d}$ of the rank- d approximate posterior $\widehat{\mathbf{\Gamma}}_m$.

522 *Correctness.* Theorem 3.3 asserts that posteriors of BayesCG under the Krylov
523 prior have means that are identical to CG iterates, and covariances that can be main-
524 tained in factored form involving submatrices of \mathbf{V} and $\mathbf{\Phi}$ from Definition 3.6. The
525 rank d of $\widehat{\mathbf{\Gamma}}_m$ has the same purpose as the ‘delay’ in CG error estimation: a small num-
526 ber of additional iterations to capture the error, and $\text{trace}(\mathbf{A} \widehat{\mathbf{\Gamma}}_m) = \text{trace}(\mathbf{\Phi}_{m+1:m+d})$
527 is equal to the error underestimate (3.19). As a termination criterion one can choose
528 the usual residual norm, or a statistically motivated criterion.

529 *Computational cost.* Algorithm 3.1 performs fewer arithmetic operations than
530 Algorithm 2.1. Specifically, Algorithm 3.1 runs $m+d$ iterations of Algorithm 2.2, and
531 a total of $m+d$ matrix vector products with \mathbf{A} and storage of at most $d+2$ vectors.
532 This is less than Algorithm 2.1, which requires $2m$ matrix vector products with \mathbf{A} , m
533 matrix vector products with $\mathbf{\Sigma}_0$, and storage of $m+2$ vectors.

534 In addition, Algorithm 2.1 requires reorthogonalization to ensure positive semi-
535 definiteness of the posterior covariances [9, Section 6.1]. In contrast, Algorithm 3.1
536 maintains the Krylov posteriors in factored form, thus (i) ensuring symmetric posi-
537 tive semi-definiteness by design; and (ii) reducing the cost of sampling, because the
538 factorizations $\mathbf{\Sigma}_m = \mathbf{F}_m \mathbf{F}_m^T$ are readily available without any computations. The

²The partition of Algorithm 3.1 into two loops is for the purpose expositional clarity. Alternat-
ively, everything could have been merged into a single loop with a conditional.

Algorithm 3.1 BayesCG under rank- d approximations of specific Krylov posterior covariances

```

1: Inputs: spd  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{b} \in \mathbb{R}^n$ ,  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $d \geq 1$ 
2:  $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$  ▷ define initial values
3:  $\mathbf{v}_1 = \mathbf{r}_0$ 
4:  $m = 0$ 
5: while not converged do ▷ CG recursions for posterior means
6:    $m = m + 1$ 
7:    $\eta_m = \mathbf{v}_m^T \mathbf{A} \mathbf{v}_m$ 
8:    $\gamma_m = (\mathbf{r}_{m-1}^T \mathbf{r}_{m-1}) / \eta_m$ 
9:    $\mathbf{x}_m = \mathbf{x}_{m-1} + \gamma_m \mathbf{v}_m$ 
10:   $\mathbf{r}_m = \mathbf{r}_{m-1} - \gamma_m \mathbf{A} \mathbf{v}_m$ 
11:   $\delta_m = (\mathbf{r}_m^T \mathbf{r}_m) / (\mathbf{r}_{m-1}^T \mathbf{r}_{m-1})$ 
12:   $\mathbf{v}_{m+1} = \mathbf{r}_m + \delta_m \mathbf{v}_m$ 
13: end while
14:  $d = \min\{d, K - m\}$  ▷ compute full rank posterior if  $d > K - m$ 
15:  $\mathbf{V}_{m+1:m+d} = \mathbf{0}_{n \times d}$  ▷ define posterior factor matrices
16:  $\Phi_{m+1:m+d} = \mathbf{0}_{d \times d}$ 
17: for  $j = m + 1 : m + d$  do ▷  $d$  additional iterations for posterior covariance
18:    $\eta_j = \mathbf{v}_j^T \mathbf{A} \mathbf{v}_j$ 
19:    $\gamma_j = (\mathbf{r}_{j-1}^T \mathbf{r}_{j-1}) / \eta_j$ 
20:    $\mathbf{V}_j = \mathbf{v}_j / \eta_j$  ▷ store column  $j$  of  $\mathbf{V}$ 
21:    $\Phi_j = \gamma_j \|\mathbf{r}_{j-1}\|_2^2$  ▷ store element  $j$  of  $\Phi$ 
22:    $\mathbf{r}_j = \mathbf{r}_{j-1} - \gamma_j \mathbf{A} \mathbf{v}_j$ 
23:    $\delta_j = (\mathbf{r}_j^T \mathbf{r}_j) / (\mathbf{r}_{j-1}^T \mathbf{r}_{j-1})$ 
24:    $\mathbf{v}_{j+1} = \mathbf{r}_j + \delta_j \mathbf{v}_j$ 
25: end for
26: Output:  $\mathbf{x}_m$ ,  $\mathbf{V}_{m+1:m+d}$ ,  $\Phi_{m+1:m+d}$ 

```

539 last point is important, since the posterior is propagated to subsequent computations
540 which sample from it to probe the effect of the uncertainty in the linear solve. So far,
541 analytical propagation of the posterior has proved elusive, and empirical propagation
542 is our only option.

543 **4. Numerical experiments.** We present numerical experiments to compare (i)
544 Algorithm 3.1 under full or rank- d approximations of specific Krylov posteriors with
545 (ii) Algorithm 2.1 under the inverse prior. After describing the experimental set up
546 (section 4.1), we apply the algorithms to two matrices: a matrix of small dimension
547 (section 4.2), and one of larger dimension (section 4.3).

548 **4.1. Set up of the numerical experiments.** We describe the linear systems
549 in the experiments, reorthogonalization in the algorithms, and sampling from the
550 posterior distributions.³

551 *Linear systems.* We consider two types of symmetric positive-definite linear sys-
552 tems $\mathbf{A}\mathbf{x}_* = \mathbf{b}$: one with a dense matrix \mathbf{A} of dimension $n = 100$, and the other with
553 a sparse preconditioned matrix \mathbf{A} of dimension $n = 11948$. We fix the solution \mathbf{x}_* ,
554 and compute the right hand side from $\mathbf{b} = \mathbf{A}\mathbf{x}_*$.

³The Python code used in the numerical experiments can be found at https://github.com/treid5/ProbNumCG_Supp

555 For $n = 100$, the matrix is $\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{Q}^T$ [22, Section 2], where \mathbf{Q} is a random⁴
 556 orthogonal matrix with Haar distribution [47, Section 3], and \mathbf{D} is a diagonal matrix
 557 with eigenvalues [20]

$$558 \quad (4.1) \quad d_{ii} = (10^3)^{(i-1)/99}, \quad 1 \leq i \leq 100.$$

559 The condition number is $\kappa(\mathbf{A}) = 10^3$, and the solution \mathbf{x}_* is sampled from $\mathcal{N}(\mathbf{0}, \mathbf{A}^{-1})$.

560 For $n = 11948$, the matrix $\mathbf{A} = \mathbf{L}^{-1}\mathbf{B}\mathbf{L}^{-T}$ is a sparse preconditioned matrix
 561 where \mathbf{B} is BCSSTK18 from the Harwell-Boeing collection [33], and \mathbf{L} is the incomplete
 562 Cholesky factorization [21, Section 11.1] of the diagonally shifted matrix

$$563 \quad \tilde{\mathbf{B}} = \mathbf{B} + 9.0930 \cdot 10^8 \cdot \text{diag}(\mathbf{B}) \quad \text{with} \quad \max_{1 \leq i \leq n} \left\{ -b_{ii} + \sum_{j \neq i} b_{ij} \right\} = 9.0930 \cdot 10^8.$$

564 The shift forces $\tilde{\mathbf{B}}$ to be diagonally dominant. We compute the factorization of $\tilde{\mathbf{B}}$
 565 with a threshold drop tolerance 10^{-6} to make \mathbf{L} diagonal. The condition number is
 566 $\kappa(\mathbf{A}) \approx 1.57 \cdot 10^6$, and the solution $\mathbf{x}_* = \mathbf{1}$ is the all ones vector.

567 *Reorthogonalization.* Since the posterior covariances in Algorithm 2.1 become
 568 indefinite when the search directions lose orthogonality, reorthogonalization of the
 569 search directions is recommended in every iteration, [9, Section 6.1] and [11, Sec-
 570 tion 4.1]. Following [22, Section 2], we reorthogonalize the residual vectors instead,
 571 as it has the additional advantage of better numerical stability in our experience.
 572 Reorthogonalization consists of classical Gram-Schmidt performed twice because it
 573 is efficient, easy to implement, and produces vectors orthogonal to almost machine
 574 precision [15, 16].

575 *Sampling from the Gaussian distributions.* We exploit the stability of Gaussians,
 576 see section 2.1, to sample from $\mathcal{N}(\mathbf{x}, \Sigma)$ as follows. Let $\Sigma = \mathbf{F}\mathbf{F}^T$ be a factorization of
 577 the covariance with $\mathbf{F} \in \mathbb{R}^{n \times d}$. Sample a standard Gaussian vector⁵ $Z \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$;
 578 multiply it by \mathbf{F} ; and add the mean to obtain $X \equiv \mathbf{x} + \mathbf{F}Z \sim \mathcal{N}(\mathbf{x}, \mathbf{F}\mathbf{F}^T)$.

579 By design, the rank- d approximate Krylov posteriors are maintained in factored
 580 form

$$581 \quad \hat{\Gamma}_m = \mathbf{F}_m \mathbf{F}_m^T \quad \text{where} \quad \mathbf{F}_m \equiv \mathbf{V}_{m+1:m+d} \Phi_{m+1:m+d}^{1/2} \in \mathbb{R}^{n \times d}.$$

582 For all other posteriors Σ_m , we factor the matrix square root [27, Chapter 6] of the
 583 matrix absolute value [27, Chapter 8] of Σ_m ⁶. Factoring the absolute value of Σ_m
 584 enforces positive semi-definiteness of the posteriors which may be lost if BayesCG is
 585 implemented without reorthogonalization.

586 *Convergence.* We display convergence of the mean and covariance with $\|\mathbf{x}_* -$
 587 $\mathbf{x}_m\|_{\mathbf{A}}^2$ and $\text{trace}(\mathbf{A}\Sigma_m)$. In addition, we sample from the posterior, $X \sim \mathcal{N}(\mathbf{x}_m, \Sigma_m)$
 588 and compare the resulting estimate $\|X - \mathbf{x}_m\|_{\mathbf{A}}^2$ to the error $\|\mathbf{x}_* - \mathbf{x}_m\|_{\mathbf{A}}^2$. If the
 589 samples X are accurate estimates, then the posterior distribution is likely to be a
 590 reliable indicator of the uncertainty in the solution \mathbf{x}_* .

⁴The exact random matrix can be reproduced with the python files in our code repository because we specified the random seed.

⁵Most scientific computing packages come with built in functions for sampling from $\mathcal{N}(\mathbf{0}, \mathbf{I})$. In Matlab and Julia the function is `randn` and in Python it is `numpy.random.randn`.

⁶The matrix absolute value of $\mathbf{B} \in \mathbb{R}^{n \times n}$ is $\text{abs}(\mathbf{B}) = (\mathbf{B}^T \mathbf{B})^{1/2}$. If \mathbf{B} is symmetric positive semi-definite, then $\text{abs}(\mathbf{B}) = \mathbf{B}$. Otherwise, the square root of the absolute value is $(\text{abs}(\mathbf{B}))^{1/2} = \mathbf{V}\mathbf{S}^{1/2}\mathbf{V}^T$, where $\mathbf{B} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ is a SVD.

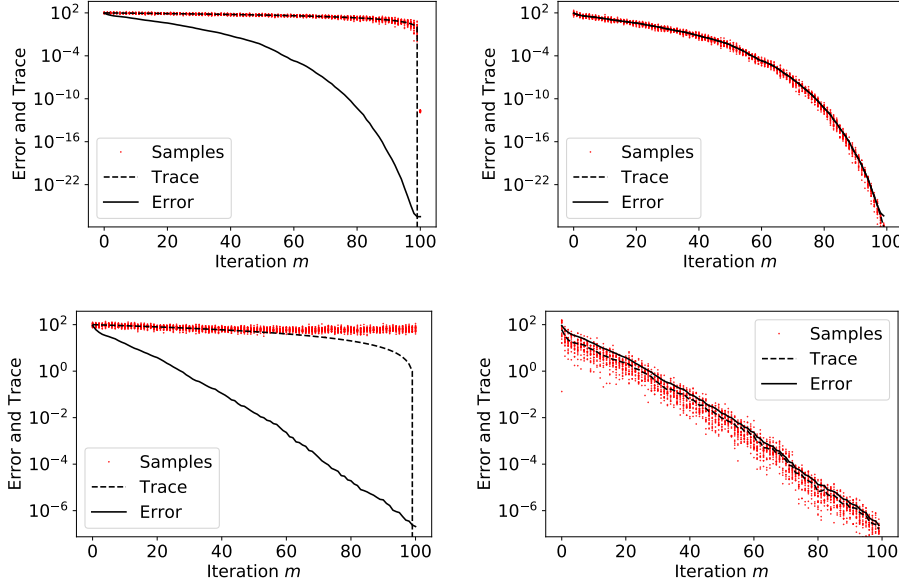


FIGURE 4.1. Error estimates $\|X - \mathbf{x}_m\|_{\mathbf{A}}^2$ and $\text{trace}(\mathbf{A}\Sigma_m)$ from samples $X \sim \mathcal{N}(\mathbf{x}_m, \Sigma_m)$, for the matrix with small dimension $n = 100$. Top row: Algorithm 2.1 with reorthogonalization under the inverse prior (left panel), and Algorithm 3.1 under the full Krylov prior (right panel). Bottom row: Algorithm 2.1 without reorthogonalization under the inverse prior (left panel), and Algorithm 3.1 under the rank-5 approximate Krylov prior (right panel).

591 **4.2. Matrix with small dimension.** We compare Algorithm 2.1 under the
 592 inverse prior, with Algorithm 3.1 under full or rank-5 approximate Krylov posteriors
 593 when applied to the matrix with small dimension $n = 100$.

594 Figure 4.1 illustrates that the posterior means converge at the same speed, regard-
 595 less of reorthogonalization. However, without reorthogonalization, the convergence is
 596 slower.

597 *Algorithm 2.1 under the inverse prior.* The posterior covariances converge more
 598 slowly than the squared errors of the means. Without reorthogonalization, the pos-
 599 terior covariances are indefinite, and the error estimates from the posterior samples
 600 diverge from $\text{trace}(\mathbf{A}\Sigma_m)$ and violate Lemma 3.5. Thus, posteriors from BayesCG
 601 under the inverse prior are not reliable indicators of uncertainty.

602 *Algorithm 3.1 under full or approximate Krylov priors.* The quantity $\text{trace}(\mathbf{A}\Sigma_m)$
 603 equals the error for full rank Krylov posteriors, while it underestimates the error for
 604 rank-5 approximate posteriors. Error estimates from samples of Krylov posteriors
 605 are significantly more accurate than those from the inverse posteriors. Thus, poste-
 606 riors from BayesCG under (approximate) Krylov priors are more reliable indicators
 607 uncertainty.

608 **4.3. Matrix with larger dimension.** We compare Algorithm 3.1 under rank-1
 609 and rank-50 approximate Krylov posteriors, when applied to the matrix with large
 610 dimension $n = 11948$.

611 Figure 4.2 illustrates that the traces of the posterior covariances underestimate
 612 the error. However, the trace of the rank-50 approximate Krylov covariance is more
 613 accurate, because CG error estimates (3.19) are more accurate for larger delays [49,

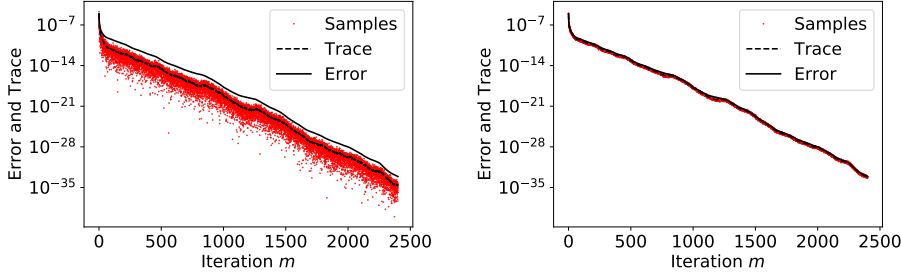


FIGURE 4.2. Error estimates $\|X - \mathbf{x}_m\|_{\mathbf{A}}^2$ and $\text{trace}(\mathbf{A}\Sigma_m)$ from samples $X \sim \mathcal{N}(\mathbf{x}_m, \Sigma_m)$, for the matrix with large dimension $n = 11948$. Left: Algorithm 3.1 under rank-1 approximate Krylov posterior. Right: Algorithm 3.1 under rank-50 approximate Krylov posterior.

614 Section 4]. As expected, error estimates from rank-50 posterior samples are more
 615 tightly concentrated around the true error than those of rank-1 posterior samples.
 616 Thus, BayesCG under higher rank approximate posteriors produces more reliable
 617 indicators of uncertainty.

618 **5. Conclusion.** BayesCG is our ‘uncertainty-aware’ version of CG, that is, a
 619 probabilistic numerical extension of CG that produces a probabilistic model of the
 620 uncertainty about our knowledge of the solution \mathbf{x}_* due to early termination of CG.
 621 Under our Krylov prior, BayesCG produces iterates that are identical to those of
 622 CG (in exact arithmetic), thus converges at the same speed as CG; and its posterior
 623 distributions can be cheaply approximated. Samples from the Krylov posterior and
 624 its low rank approximations produce accurate error estimates, thus represent realistic
 625 indicators of the uncertainty about \mathbf{x}_* .

626 *Future work.* In a forthcoming paper, we focus on the statistical aspects of
 627 BayesCG under the Krylov prior. More specifically, we quantify the approximation
 628 error of low rank approximate Krylov posteriors and investigate the calibration of
 629 BayesCG under low-rank approximate Krylov posteriors.

630 In a separate paper, we assess the effect of CG accuracy in a computational
 631 pipeline in the form of a randomized algorithm for generalized singular value decom-
 632 position [44] with applications to hyper-differential sensitivity analysis [23].

633 Appendix A. Proofs of Theorems 2.4, 2.6 and 2.7.

634 *Proof of Theorem 2.4.* The proof is inspired by the proof of [10, Proposition 3]
 635 for nonsingular Σ_0 . For singular Σ_0 , we replace the inverse by the Moore-Penrose
 636 inverse which satisfies

$$637 \quad (\text{A.1}) \quad \Sigma_0 = \Sigma_0 \Sigma_0^\dagger \Sigma_0.$$

638 The assumption $\mathbf{x}_* - \mathbf{x}_0 \in \text{range}(\Sigma_0)$ implies that there exists $\mathbf{y} \in \mathbb{R}^n$ so that

$$639 \quad (\text{A.2}) \quad \mathbf{x}_* - \mathbf{x}_0 = \Sigma_0 \mathbf{y} = \Sigma_0 \Sigma_0^\dagger \Sigma_0 \mathbf{y} = \Sigma_0 \Sigma_0^\dagger (\mathbf{x}_* - \mathbf{x}_0).$$

640 The proof proceeds in four steps.

641 *Range of \mathbf{P}_m .* On the one hand (2.3) implies

$$642 \quad \text{range}(\mathbf{P}_m) = \text{range} \left(\Sigma_0 \mathbf{A} \Sigma_m \Lambda_m^{-1} \mathbf{S}_m^T \mathbf{A} \Sigma_0 \Sigma_0^\dagger \right) \subset \text{range}(\Sigma_0 \mathbf{A} \Sigma_m).$$

643

644 On the other hand (2.3) and (A.1) imply

$$645 \quad \mathbf{P}_m \boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m = \boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m \boldsymbol{\Lambda}_m^{-1} \overbrace{\mathbf{S}_m^T \mathbf{A} \boldsymbol{\Sigma}_0 \boldsymbol{\Sigma}_0^\dagger \boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m}^{\boldsymbol{\Lambda}_m} = \boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m$$

$$646 \quad \boldsymbol{\Sigma}_0$$

647 so that $\text{range}(\boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m) \subset \text{range}(\mathbf{P}_m)$.

648 Combining the two inclusions gives $\text{range}(\mathbf{P}_m) = K_m \equiv \text{range}(\boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m)$.

649 \mathbf{P}_m is a $\boldsymbol{\Sigma}_0^\dagger$ -orthogonal projector. The above implies

$$650 \quad (\text{A.3}) \quad \mathbf{P}_m^2 = \underbrace{\mathbf{P}_m \boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m}_{\boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m} \boldsymbol{\Lambda}_m^{-1} \mathbf{S}_m^T \mathbf{A} \boldsymbol{\Sigma}_0 \boldsymbol{\Sigma}_0^\dagger = \boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m \boldsymbol{\Lambda}_m^{-1} \mathbf{S}_m^T \mathbf{A} \boldsymbol{\Sigma}_0 \boldsymbol{\Sigma}_0^\dagger = \mathbf{P}_m.$$

$$651 \quad \boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m$$

652 Thus \mathbf{P}_m is a projector. The $\boldsymbol{\Sigma}_0^\dagger$ -orthogonality of \mathbf{P}_m follows from the symmetry of
653 $\boldsymbol{\Sigma}_0^\dagger \mathbf{P}$.

654 *Posterior mean.* From (2.1), (A.2), and (2.3) follows

$$655 \quad \mathbf{x}_m = \mathbf{x}_0 + \boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m \boldsymbol{\Lambda}_m^{-1} \mathbf{S}_m^T \mathbf{A} (\mathbf{x}_* - \mathbf{x}_0)$$

$$656 \quad = \mathbf{x}_0 + \boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m \boldsymbol{\Lambda}_m^{-1} \mathbf{S}_m^T \mathbf{A} \boldsymbol{\Sigma}_0 \boldsymbol{\Sigma}_0^\dagger (\mathbf{x}_* - \mathbf{x}_0) = (\mathbf{I} - \mathbf{P}_m) \mathbf{x}_0 + \mathbf{P}_m \mathbf{x}_*.$$

658 *Posterior covariance.* From (2.2), (A.1) and (2.3) follows

$$659 \quad \boldsymbol{\Sigma}_m = \boldsymbol{\Sigma}_0 - \boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m \boldsymbol{\Lambda}_m^{-1} \mathbf{S}_m^T \mathbf{A} \boldsymbol{\Sigma}_0$$

$$660 \quad = \boldsymbol{\Sigma}_0 - \boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m \boldsymbol{\Lambda}_m^{-1} \mathbf{S}_m^T \mathbf{A} \boldsymbol{\Sigma}_0 \boldsymbol{\Sigma}_0^\dagger \boldsymbol{\Sigma}_0 = (\mathbf{I} - \mathbf{P}_m) \boldsymbol{\Sigma}_0.$$

662 Multiply $\boldsymbol{\Sigma}_m$ on the left by \mathbf{P}_m and apply (A.3) to obtain $\mathbf{P}_m \boldsymbol{\Sigma}_m = \mathbf{P}_m (\mathbf{I} - \mathbf{P}_m) \boldsymbol{\Sigma}_0 =$
663 $\mathbf{0}$. \square

664 The proof of Theorem 2.6 relies on the next three results related to semi-definite
665 inner product spaces and orthogonal projectors in those spaces.

666 LEMMA A.1. *Under the assumptions of Theorem 2.1, if $\mathbf{x}_* - \mathbf{x}_0 \in \text{range}(\boldsymbol{\Sigma}_0)$,*
667 *then $\mathbf{x}_* - \mathbf{x}_m \in \text{range}(\boldsymbol{\Sigma}_0)$, $1 \leq m \leq n$.*

668 *Proof.* Subtract from \mathbf{x}_* both sides of the posterior mean (2.1),

$$669 \quad \mathbf{x}_* - \mathbf{x}_m = (\mathbf{x}_* - \mathbf{x}_0) - \boldsymbol{\Sigma}_0 \mathbf{A} \mathbf{S}_m \boldsymbol{\Lambda}_m^{-1} \mathbf{S}_m^T \mathbf{A} (\mathbf{x}_* - \mathbf{x}_0), \quad 1 \leq m \leq n.$$

670 The first summand $\mathbf{x}_* - \mathbf{x}_0$ is in $\text{range}(\boldsymbol{\Sigma}_0)$ by assumption, and the second one by
671 design, hence so is the sum. \square

672 LEMMA A.2. *Let $\mathbf{B} \in \mathbb{R}^{n \times n}$ be symmetric positive semi-definite. If $\mathbf{z} \in \text{range}(\mathbf{B})$,*
673 *then $\mathbf{z}^T \mathbf{B} \mathbf{z} = 0$ if and only if $\mathbf{z} = \mathbf{0}$.*

674 *Proof.* Since \mathbf{B} is symmetric positive semi-definite, we can factor $\mathbf{B} \mathbf{B}^T = \mathbf{B}$,
675 where \mathbf{B} has full column rank. Let $\mathbf{w} = \mathbf{B}^T \mathbf{z}$. From $\mathbf{z} \in \text{range}(\mathbf{B}) = \text{range}(\mathbf{B})$, and
676 $\text{range}(\mathbf{B}) = \ker(\mathbf{B}^T)^\perp$ follows that $\mathbf{w} = \mathbf{B}^T \mathbf{z} = \mathbf{0}$ if and only if $\mathbf{z} = \mathbf{0}$. Therefore
677 $\mathbf{w}^T \mathbf{w} = \mathbf{z}^T \mathbf{B} \mathbf{z} = 0$ if and only if $\mathbf{z} = \mathbf{0}$. \square

678 LEMMA A.3. *Let $\mathcal{X} \subseteq \mathbb{R}^n$ be a subspace, $\mathbf{B} \in \mathbb{R}^{n \times n}$ symmetric positive semi-*
679 *definite, and $\mathbf{v} \in \mathbb{R}^n$. If \mathbf{P} is a \mathbf{B} -orthogonal projector onto \mathcal{X} , then*

$$680 \quad \arg \min_{\mathbf{x} \in \mathcal{X}} (\mathbf{v} - \mathbf{x})^T \mathbf{B} (\mathbf{v} - \mathbf{x}) = \{\mathbf{x} \in \mathcal{X} : (\mathbf{x} - \mathbf{P} \mathbf{v})^T \mathbf{B} (\mathbf{x} - \mathbf{P} \mathbf{v}) = 0\}.$$

$$681 \quad \mathbf{x} \in \mathcal{X}$$

682 If additionally $\mathcal{X} \subseteq \text{range}(\mathbf{B})$, then

$$683 \quad \arg \min_{\mathbf{x} \in \mathcal{X}} (\mathbf{v} - \mathbf{x})^T \mathbf{B} (\mathbf{v} - \mathbf{x}) = \mathbf{P}\mathbf{v}.$$

684 *Proof.* After proving the general case, we show that the minimizer is unique if
685 $\mathcal{X} \subseteq \text{range}(\mathbf{B})$.

686 *General case.* Abbreviate the induced semi-norm by $|\mathbf{z}|_{\mathbf{B}}^2 = \mathbf{z}^T \mathbf{B} \mathbf{z}$. Since \mathbf{P} is a
687 projector onto \mathcal{X} , we can write $\mathbf{x} = \mathbf{P}\mathbf{x}$ for $\mathbf{x} \in \mathcal{X}$. Add and subtract $\mathbf{P}\mathbf{v}$ inside the
688 norm to obtain a Pythagoras-like theorem,

$$\begin{aligned} 689 \quad |\mathbf{v} - \mathbf{x}|_{\mathbf{B}}^2 &= |(\mathbf{I} - \mathbf{P})\mathbf{v} + \mathbf{P}(\mathbf{v} - \mathbf{x})|_{\mathbf{B}}^2 \\ 690 \quad &= |(\mathbf{I} - \mathbf{P})\mathbf{v}|_{\mathbf{B}}^2 + |\mathbf{P}(\mathbf{v} - \mathbf{x})|_{\mathbf{B}}^2 + 2\mathbf{v}^T \underbrace{(\mathbf{I} - \mathbf{P})^T \mathbf{B} \mathbf{P}}_{=0} (\mathbf{v} - \mathbf{x}) \\ 691 \quad &= |(\mathbf{I} - \mathbf{P})\mathbf{v}|_{\mathbf{B}}^2 + |\mathbf{P}\mathbf{v} - \mathbf{x}|_{\mathbf{B}}^2. \end{aligned}$$

693 Since the first summand is independent of \mathbf{x} , the minimum is achieved if the second
694 summand is zero.

695 *Uniqueness.* Since \mathbf{P} is a projector onto \mathcal{X} , $\mathbf{P}\mathbf{v} \in \mathcal{X}$. From $\mathcal{X} \subseteq \text{range}(\mathbf{B})$ follows
696 $\mathbf{P}\mathbf{v} \in \text{range}(\mathbf{B})$ and $\mathbf{x} \in \text{range}(\mathbf{B})$. With Lemma A.2 this implies: $|\mathbf{P}\mathbf{v} - \mathbf{x}|_{\mathbf{B}}^2 = 0$
697 only if $\mathbf{P}\mathbf{v} = \mathbf{x}$. \square

698 *Proof of Theorem 2.6.* This is similar to [1, Proof of Proposition 4]. Minimizing
699 (2.4) over the affine space $\mathbf{x}_0 + K_m = \mathbf{x}_0 + \text{range}(\mathbf{\Sigma}_0 \mathbf{A} \mathbf{S}_m)$ is equivalent to shifting
700 by \mathbf{x}_0 and minimizing over K_m ,

$$701 \quad \min_{\mathbf{x} \in \mathbf{x}_0 + K_m} (\mathbf{x}_* - \mathbf{x})^T \mathbf{\Sigma}_0^\dagger (\mathbf{x}_* - \mathbf{x}) = \min_{\mathbf{x} \in K_m} ((\mathbf{x}_* - \mathbf{x}_0) - \mathbf{x})^T \mathbf{\Sigma}_0^\dagger ((\mathbf{x}_* - \mathbf{x}_0) - \mathbf{x}).$$

702 Since $\mathbf{\Sigma}_0$ is symmetric, the $\mathbf{\Sigma}_0^\dagger$ -orthogonal projector \mathbf{P}_m from Theorem 2.4 satisfies
703 $\text{range}(\mathbf{P}_m) = K_m \subseteq \text{range}(\mathbf{\Sigma}_0) = \text{range}(\mathbf{\Sigma}_0^\dagger)$. Therefore, Lemma A.3 implies

$$704 \quad \arg \min_{\mathbf{x} \in K_m} ((\mathbf{x}_* - \mathbf{x}_0) - \mathbf{x})^T \mathbf{\Sigma}_0^\dagger ((\mathbf{x}_* - \mathbf{x}_0) - \mathbf{x}) = \mathbf{P}(\mathbf{x}_* - \mathbf{x}_0).$$

706 From Theorem 2.4 and $K_m = \text{range}(\mathbf{P}_m)$ follows $\mathbf{x}_m - \mathbf{x}_0 = \mathbf{P}_m(\mathbf{x}_* - \mathbf{x}_0) \in K_m$.
707 Thus $\mathbf{x}_m \in \mathbf{x}_0 + K_m$ is the minimizer.

708 The symmetry of $\mathbf{\Sigma}_m$ and Lemmas A.1 and A.2 imply that $(\mathbf{x}_* - \mathbf{x}_m)^T \mathbf{\Sigma}_0^\dagger (\mathbf{x}_* -$
709 $\mathbf{x}_m) = 0$ only if $\mathbf{x}_m = \mathbf{x}_*$. \square

710 *Proof of Theorem 2.7.* Recursion (2.6) was shown in [9, Proposition 6]. The fol-
711 lowing proof for (2.7) is analogous to [11, Proof of Proposition 6]. From (2.2) follows
712 that the posterior covariance at iteration i amounts to a rank- i downdate of the prior,

$$713 \quad \mathbf{\Sigma}_i = \mathbf{\Sigma}_0 - \mathbf{\Sigma}_0 \mathbf{A} \mathbf{S}_i \mathbf{\Lambda}_i^{-1} (\mathbf{\Sigma}_0 \mathbf{A} \mathbf{S}_i)^T, \quad 1 \leq i \leq m.$$

714 Here $\mathbf{\Lambda}_i$ is diagonal due to the $\mathbf{A} \mathbf{\Sigma}_0 \mathbf{A}$ -orthogonality of the search directions, hence a
715 rank- i downdate can be computed as a recursive sequence of i rank-1 downdates,

$$716 \quad \mathbf{\Sigma}_i = \underbrace{\mathbf{\Sigma}_0 - \mathbf{\Sigma}_0 \mathbf{A} \mathbf{S}_{i-1} \mathbf{\Lambda}_{i-1}^{-1} (\mathbf{\Sigma}_0 \mathbf{A} \mathbf{S}_{i-1})^T}_{\mathbf{\Sigma}_{i-1}} - \frac{\mathbf{\Sigma}_0 \mathbf{A} \mathbf{S}_i (\mathbf{\Sigma}_0 \mathbf{A} \mathbf{S}_i)^T}{\mathbf{s}_i^T \mathbf{A} \mathbf{\Sigma}_0 \mathbf{A} \mathbf{S}_i}. \quad \square$$

718 **Appendix B. Auxiliary results.**

719 LEMMA B.1 (Lemma S3 in [11]). *Under the assumptions of Theorem 2.7,*

$$720 \quad \mathbf{s}_j^T \mathbf{r}_i = 0, \quad 1 \leq j \leq i \leq m.$$

721 LEMMA B.2 (Sections 3.2b.1–3.2b.3 in [32]). *Let $Z \sim \mathcal{N}(\mathbf{x}, \Sigma)$ be a Gaussian*
 722 *random variable with mean $\mathbf{x} \in \mathbb{R}^n$ and covariance $\Sigma \in \mathbb{R}^{n \times n}$, and let $\mathbf{B} \in \mathbb{R}^{n \times n}$ be*
 723 *symmetric positive definite. The mean and variance of $Z^T \mathbf{B} Z$ are*

$$724 \quad \mathbb{E}[Z^T \mathbf{B} Z] = \text{trace}(\mathbf{B} \Sigma) + \mathbf{x}^T \mathbf{B} \mathbf{x},$$

$$725 \quad \mathbb{V}[Z^T \mathbf{B} Z] = 2 \text{trace}((\mathbf{B} \Sigma)^2) + 4 \mathbf{x}^T \mathbf{B} \Sigma \mathbf{B} \mathbf{x}.$$

727 **Acknowledgments.** We thank Eric Hallman, Joseph Hart, and the members
 728 of the NCSU Randomized Numerical Analysis RTG for helpful discussions. We are
 729 also most grateful to the reviewers for their recommendations that helped to ensure
 730 mathematical correctness and improve exposition.

731

REFERENCES

- 732 [1] S. BARTELS, J. COCKAYNE, I. C. F. IPSEN, AND P. HENNIG, *Probabilistic linear solvers:*
 733 *a unifying view*, Stat. Comput., 29 (2019), pp. 1249–1263, [https://doi.org/10.1007/](https://doi.org/10.1007/s11222-019-09897-7)
 734 [s11222-019-09897-7](https://doi.org/10.1007/s11222-019-09897-7).
- 735 [2] S. BARTELS AND P. HENNIG, *Probabilistic approximate least-squares*, in Proc. 19th Int. Conf.
 736 Artificial Intelligence and Statistics, vol. 51 of Proc. Machine Learning Research, MLR
 737 Press, 2016, pp. 676–684.
- 738 [3] M. BERLJAJA AND S. GÜTTEL, *Generalized rational Krylov decompositions with an application*
 739 *to rational approximation*, SIAM J. Matrix Anal. Appl., 36 (2015), pp. 894–916, <https://doi.org/10.1137/140998081>.
- 740 [4] F.-X. BRIOL, C. J. OATES, M. GIROLAMI, M. A. OSBORNE, AND D. SEJIDINOVIC, *Probabilistic*
 741 *Integration: A Role in Statistical Computation?*, Statist. Sci., 34 (2019), pp. 1 – 22, <https://doi.org/10.1214/18-STS660>.
- 742 [5] D. CALVETTI, *Contributed discussion for “A Bayesian conjugate gradient method”*, Bayesian
 743 Anal., 14 (2019), pp. 937–1012, <https://doi.org/10.1214/19-BA1145>.
- 744 [6] J. COCKAYNE, M. M. GRAHAM, C. J. OATES, AND T. J. SULLIVAN, *Testing whether a learning*
 745 *procedure is calibrated*, 2021, <https://arxiv.org/abs/2012.12670>. arXiv:2012.12670.
- 746 [7] J. COCKAYNE, I. C. F. IPSEN, C. J. OATES, AND T. W. REID, *Probabilistic iterative methods*
 747 *for linear systems*, J. Mach. Learn. Res., 22 (232) (2021), pp. 1–34.
- 748 [8] J. COCKAYNE, C. OATES, T. SULLIVAN, AND M. GIROLAMI, *Probabilistic numerical methods for*
 749 *PDE-constrained Bayesian inverse problems*, AIP Conference Proceedings, 1853 (2017),
 750 p. 060001, <https://doi.org/10.1063/1.4985359>.
- 751 [9] J. COCKAYNE, C. J. OATES, I. C. F. IPSEN, AND M. GIROLAMI, *A Bayesian conjugate gradient*
 752 *method (with discussion)*, Bayesian Anal., 14 (2019), pp. 937–1012, <https://doi.org/10.1214/19-BA1145>. Includes 6 discussions and a rejoinder from the authors.
- 753 [10] J. COCKAYNE, C. J. OATES, I. C. F. IPSEN, AND M. GIROLAMI, *Rejoinder for “A Bayesian*
 754 *conjugate gradient method”*, Bayesian Anal., 14 (2019), pp. 937–1012, <https://doi.org/10.1214/19-BA1145>.
- 755 [11] J. COCKAYNE, C. J. OATES, I. C. F. IPSEN, AND M. GIROLAMI, *Supplementary material for*
 756 *“A Bayesian conjugate-gradient method”*, Bayesian Anal., (2019), [https://doi.org/10.1214/](https://doi.org/10.1214/19-BA1145SUPP)
 757 [19-BA1145SUPP](https://doi.org/10.1214/19-BA1145SUPP).
- 758 [12] J. COCKAYNE, C. J. OATES, T. J. SULLIVAN, AND M. GIROLAMI, *Bayesian probabilistic numeri-*
 759 *cal methods*, SIAM Rev., 61 (2019), pp. 756–789, <https://doi.org/10.1137/17M1139357>.
- 760 [13] V. FANASKOV, *Uncertainty calibration for probabilistic projection methods*, Stat. Comput., 31
 761 (2021), pp. Paper No. 56, 17, <https://doi.org/10.1007/s11222-021-10031-9>, <https://doi.org/10.1007/s11222-021-10031-9>.
- 762 [14] A. GESSNER, O. KANJILAL, AND P. HENNIG, *Integrals over Gaussians under linear domain*
 763 *constraints*, in Proceedings of the Twenty Third International Conference on Artificial
 764 Intelligence and Statistics, S. Chiappa and R. Calandra, eds., vol. 108 of Proceedings
 765 of Machine Learning Research, 2020, pp. 2764–2774, [http://proceedings.mlr.press/v108/](http://proceedings.mlr.press/v108/gessner20a.html)
 766 [gessner20a.html](http://proceedings.mlr.press/v108/gessner20a.html).

- 772 [15] L. GIRAUD, J. LANGOU, M. ROZLOŽNÍK, AND J. VAN DEN ESHOF, *Rounding error analysis of the*
773 *classical Gram-Schmidt orthogonalization process*, Numer. Math., 101 (2005), pp. 87–100,
774 <https://doi.org/10.1007/s00211-005-0615-4>.
- 775 [16] L. GIRAUD, J. LANGOU, AND M. ROZLOZNIK, *The loss of orthogonality in the Gram-Schmidt*
776 *orthogonalization process*, Comput. Math. Appl., 50 (2005), pp. 1069–1075, <https://doi.org/10.1016/j.camwa.2005.08.009>.
- 777 [17] G. H. GOLUB AND G. MEURANT, *Matrices, moments and quadrature*, in Numerical analysis
778 1993 (Dundee, 1993), vol. 303 of Pitman Res. Notes Math. Ser., Longman Sci. Tech.,
779 Harlow, 1994, pp. 105–156.
- 780 [18] G. H. GOLUB AND G. MEURANT, *Matrices, moments and quadrature. II. How to compute the*
781 *norm of the error in iterative methods*, BIT, 37 (1997), pp. 687–705, [https://doi.org/10.](https://doi.org/10.1007/BF02510247)
782 [1007/BF02510247](https://doi.org/10.1007/BF02510247).
- 783 [19] G. H. GOLUB AND Z. STRAKOŠ, *Estimates in quadratic formulas*, Numer. Algorithms, 8 (1994),
784 pp. 241–268, <https://doi.org/10.1007/BF02142693>.
- 785 [20] A. GREENBAUM, *Estimating the attainable accuracy of recursively computed residual meth-*
786 *ods*, SIAM J. Matrix Anal. Appl., 18 (1997), pp. 535–551, [https://doi.org/10.1137/](https://doi.org/10.1137/S0895479895284944)
787 [S0895479895284944](https://doi.org/10.1137/S0895479895284944).
- 788 [21] A. GREENBAUM, *Iterative methods for solving linear systems*, vol. 17 of Frontiers in Applied
789 Mathematics, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA,
790 1997, <https://doi.org/10.1137/1.9781611970937>.
- 791 [22] A. GREENBAUM AND Z. STRAKOŠ, *Predicting the behavior of finite precision Lanczos and*
792 *conjugate gradient computations*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 121–137,
793 <https://doi.org/10.1137/0613011>.
- 794 [23] J. HART, B. VAN BLOEMEN WAANDERS, AND R. HERZOG, *Hyperdifferential sensitivity analysis of*
795 *uncertain parameters in PDE-constrained optimization*, Int. J. for Uncertain. Quantif., 10
796 (2020), pp. 225–248, <https://doi.org/10.1615/Int.J.UncertaintyQuantification.2020032480>.
- 797 [24] P. HENNIG, *Probabilistic interpretation of linear solvers*, SIAM J. Optim., 25 (2015), pp. 234–
798 260, <https://doi.org/10.1137/140955501>.
- 799 [25] P. HENNIG, M. A. OSBORNE, AND M. GIROLAMI, *Probabilistic numerics and uncertainty in*
800 *computations*, Proc. A., 471 (2015), pp. 20150142, 17.
- 801 [26] M. R. HESTENES AND E. STIEFEL, *Methods of conjugate gradients for solving linear systems*, J.
802 Research Nat. Bur. Standards, 49 (1952), pp. 409–436, [https://doi.org/10.6028/jres.049.](https://doi.org/10.6028/jres.049.044)
803 [044](https://doi.org/10.6028/jres.049.044).
- 804 [27] N. J. HIGHAM, *Functions of matrices. Theory and computation*, Society for Industrial
805 and Applied Mathematics (SIAM), Philadelphia, PA, 2008, [https://doi.org/10.1137/1.](https://doi.org/10.1137/1.9780898717778)
806 [9780898717778](https://doi.org/10.1137/1.9780898717778).
- 807 [28] R. A. HORN AND C. R. JOHNSON, *Matrix Analysis*, Cambridge University Press, second ed.,
808 2013.
- 809 [29] T. KARVONEN, C. J. OATES, AND S. SARKKA, *A Bayes-Sard cubature method*, in Advances in
810 Neural Information Processing Systems, S. Bengio, H. Wallach, H. Larochelle, K. Grauman,
811 N. Cesa-Bianchi, and R. Garnett, eds., vol. 31, Curran Associates, Inc., 2018, [https://](https://proceedings.neurips.cc/paper/2018/file/6775a0635c302542da2c32aa19d86be0-Paper.pdf)
812 proceedings.neurips.cc/paper/2018/file/6775a0635c302542da2c32aa19d86be0-Paper.pdf.
- 813 [30] L. LI AND E. X. FANG, *Invited discussion for “A Bayesian conjugate gradient method”*,
814 Bayesian Anal., 14 (2019), pp. 937–1012, <https://doi.org/10.1214/19-BA1145>.
- 815 [31] J. LIESEN AND Z. STRAKOS, *Krylov Subspace Methods: Principles and Analysis*, Oxford Uni-
816 versity Press, 2013.
- 817 [32] A. M. MATHAI AND S. B. PROVOST, *Quadratic forms in random variables: Theory and appli-*
818 *cations*, Dekker, 1992.
- 819 [33] MATRIX MARKET, *BCSSTK18: BCS Structural Engineering Matrices (linear equa-*
820 *tions) R.E. Ginna Nuclear Power Station*, [https://math.nist.gov/MatrixMarket/data/](https://math.nist.gov/MatrixMarket/data/Harwell-Boeing/bcsstruc2/bcsstk18.html)
821 [Harwell-Boeing/bcsstruc2/bcsstk18.html](https://math.nist.gov/MatrixMarket/data/Harwell-Boeing/bcsstruc2/bcsstk18.html).
- 822 [34] T. MATSUDA AND Y. MIYATAKE, *Estimation of ordinary differential equation models with dis-*
823 *cretization error quantification*, SIAM/ASA J. Uncertain. Quantif., 9 (2021), pp. 302–331,
824 <https://doi.org/10.1137/19M1278405>.
- 825 [35] G. MEURANT, *The computation of bounds for the norm of the error in the conjugate gradient*
826 *algorithm*, Numer. Algorithms, 16 (1997), pp. 77–87 (1998), [https://doi.org/10.1023/A:](https://doi.org/10.1023/A:1019178811767)
827 [1019178811767](https://doi.org/10.1023/A:1019178811767). Sparse matrices in industry (Lille, 1997).
- 828 [36] G. MEURANT AND P. TICHÝ, *On computing quadrature-based bounds for the A-norm of the*
829 *error in conjugate gradients*, Numer. Algorithms, 62 (2013), pp. 163–191, [https://doi.org/](https://doi.org/10.1007/s11075-012-9591-9)
830 [10.1007/s11075-012-9591-9](https://doi.org/10.1007/s11075-012-9591-9).
- 831 [37] G. MEURANT AND P. TICHÝ, *Approximating the extreme Ritz values and upper bounds for the*
832 *A-norm of the error in CG*, Numer. Algorithms, 82 (2019), pp. 937–968, <https://doi.org/>
833

- 834 [10.1007/s11075-018-0634-8](https://doi.org/10.1007/s11075-018-0634-8).
- 835 [38] J. MOČKUS, *On Bayesian methods for seeking the extremum*, in Optimization Techniques IFIP
836 Technical Conference Novosibirsk, July 1–7, 1974, G. I. Marchuk, ed., Berlin, Heidelberg,
837 1975, Springer Berlin Heidelberg, pp. 400–404.
- 838 [39] R. J. MUIRHEAD, *Aspects of multivariate statistical theory*, John Wiley & Sons, Inc., New York,
839 1982. Wiley Series in Probability and Mathematical Statistics.
- 840 [40] J. NOCEDAL AND S. J. WRIGHT, *Numerical optimization*, Springer Series in Operations Research
841 and Financial Engineering, Springer, New York, second ed., 2006.
- 842 [41] C. J. OATES, J. COCKAYNE, R. G. AYKROYD, AND M. GIROLAMI, *Bayesian probabilistic numer-*
843 *ical methods in time-dependent state estimation for industrial hydrocyclone equipment*, J.
844 Amer. Statist. Assoc., 114 (2019), pp. 1518–1531, [https://doi.org/10.1080/01621459.2019.](https://doi.org/10.1080/01621459.2019.1574583)
845 [1574583](https://doi.org/10.1080/01621459.2019.1574583).
- 846 [42] C. J. OATES AND T. J. SULLIVAN, *A modern retrospective on probabilistic numerics*, Stat.
847 Comput., 29 (2019), pp. 1335–1351, <https://doi.org/10.1007/s11222-019-09902-z>.
- 848 [43] N. PETRA, H. ZHU, G. STADLER, T. HUGHES, AND O. GHATTAS, *An inexact Gauss-Newton*
849 *method for inversion of basal sliding and rheology parameters in a nonlinear Stokes ice*
850 *sheet model*, J. Glaciology, 58 (2012), p. 889–903, <https://doi.org/10.3189/2012JoG11J182>.
- 851 [44] A. K. SAIBABA, J. HART, AND B. VAN BLOEMEN WAANDERS, *Randomized algorithms for gener-*
852 *alized singular value decomposition with application to sensitivity analysis*, Numer. Linear
853 Algebra Appl., (2021), p. e2364, <https://doi.org/10.1002/nla.2364>.
- 854 [45] F. SCHÄFER, T. J. SULLIVAN, AND H. OWHADI, *Compression, inversion, and approximate*
855 *PCA of dense kernel matrices at near-linear computational complexity*, Multiscale Model.
856 Simul., 19 (2021), pp. 688–730, <https://doi.org/10.1137/19M129526X>.
- 857 [46] J. SNOEK, H. LAROCHELLE, AND R. P. ADAMS, *Practical Bayesian optimization of machine*
858 *learning algorithms*, in Proceedings of the 25th International Conference on Neural In-
859 formation Processing Systems - Volume 2, NIPS’12, Red Hook, NY, USA, 2012, Curran
860 Associates Inc., p. 2951–2959.
- 861 [47] G. W. STEWART, *The efficient generation of random orthogonal matrices with an application*
862 *to condition estimators*, SIAM J. Numer. Anal., 17 (1980), pp. 403–409, [https://doi.org/](https://doi.org/10.1137/0717034)
863 [10.1137/0717034](https://doi.org/10.1137/0717034).
- 864 [48] G. W. STEWART AND J.-G. SUN, *Matrix Perturbation Theory*, Academic Press, San Diego,
865 1990.
- 866 [49] Z. STRAKOŠ AND P. TICHÝ, *On error estimation in the conjugate gradient method and why it*
867 *works in finite precision computations*, Electron. Trans. Numer. Anal., 13 (2002), pp. 56–
868 80.
- 869 [50] Z. STRAKOŠ AND P. TICHÝ, *Error estimation in preconditioned conjugate gradients*, BIT, 45
870 (2005), pp. 789–817, <https://doi.org/10.1007/s10543-005-0032-1>.
- 871 [51] A. M. STUART, *Inverse problems: a Bayesian perspective*, Acta Numer., 19 (2010), pp. 451–559,
872 <https://doi.org/10.1017/S0962492910000061>.
- 873 [52] F. TRONARP, H. KERSTING, S. SÄRKKÄ, AND P. HENNIG, *Probabilistic solutions to ordinary*
874 *differential equations as nonlinear Bayesian filtering: a new perspective*, Stat. Comput.,
875 29 (2019), pp. 1297–1315, <https://doi.org/10.1007/s11222-019-09900-1>.
- 876 [53] J. WENGER AND P. HENNIG, *Probabilistic linear solvers for machine learning*, 2020, [https:](https://arxiv.org/abs/2010.09691)
877 [//arxiv.org/abs/2010.09691](https://arxiv.org/abs/2010.09691). arXiv:2010.09691.