

**An efficient scaled spectral preconditioner for sequences of
symmetric positive definite linear systems**

YOUSSEF DIOUANE, SELIME GÜROL, OUSSAMA MOUHTAL AND DOMINIQUE
ORBAN

Technical Report TR/PA/24/134

1 **AN EFFICIENT SCALED SPECTRAL PRECONDITIONER FOR**
2 **SEQUENCES OF SYMMETRIC POSITIVE DEFINITE LINEAR**
3 **SYSTEMS***

4 Y. DIOUANE[†], S.GÜROL[‡], O.MOUHTAL^{‡†}, AND D.ORBAN[†]

5 **Abstract.** We explore a *scaled* spectral preconditioner for the efficient solution of sequences
6 of symmetric and positive-definite linear systems. We design the scaled preconditioner not only as
7 an approximation of the inverse of the linear system but also with consideration of its use within
8 the conjugate gradient (CG) method. We propose three different strategies for selecting a scaling
9 parameter, which aims to position the eigenvalues of the preconditioned matrix in a way that reduces
10 the energy norm of the error, the quantity that CG monotonically decreases at each iteration. Our
11 focus is on accelerating convergence especially in the early iterations, which is particularly important
12 when CG is truncated due to computational cost constraints. Numerical experiments provide in
13 data assimilation confirm that the *scaled* spectral preconditioner can significantly improve early CG
14 convergence with negligible computational cost.

15 **Key words.** Sequence of linear systems, conjugate gradient method, deflated CG, spectral
16 preconditioner, convergence rate, data assimilation.

17 **MSC codes.** 68Q25, 68R10, 68U05

18 **1. Introduction.** Efficiently solving sequences of symmetric positive-definite
19 (SPD) linear systems

20 (1.1) $A^{(j)}x^{(j)} = b^{(j)}, \quad j = 1, 2, \dots$

21 is crucial in various inverse problems of computational science and engineering. For
22 instance, in data assimilation [4, 15], where one aims to solve a large-scale weighted
23 regularized nonlinear least-squares problem via the truncated Gauss-Newton algo-
24 rithm (GN) [10, 20], each iteration involves solving a linear least-squares subproblem.
25 The latter may be formulated as a large SPD linear system, typically solved using
26 the preconditioned conjugate-gradient method (PCG). Since consecutive systems do
27 not differ significantly, recycling Krylov subspace information has been explored and
28 proven to be effective [6, 17, 11, 19].

29 One way of recycling Krylov subspace information involves leveraging search di-
30 rections obtained from PCG on earlier systems to construct a limited-memory quasi-
31 Newton preconditioner (LMP) [17, 19]. This preconditioner, built solely from PCG
32 information, does not require explicit knowledge of any matrix in the sequence, mak-
33 ing it particularly suitable for applications where only matrix-vector products are
34 available, which is the case of data assimilation. [11] generalizes this limited-memory
35 preconditioner, and introduces specific variants when used with eigen- or Ritz pairs.

36 They focused on a first-level preconditioner, capable of clustering most eigenvalues
37 at 1 with few outliers, is already available for the first linear system in sequence.
38 Then, they used LMP as a second-level preconditioner to improve the efficiency of
39 the first. The goal of the LMP is to capture directions in a low-dimensional subspace
40 that the first-level preconditioner may miss, and use them to improve convergence of
41 PCG. When $A^{(j)} = A$ for all j , spectral analysis of the preconditioned matrix when

*Submitted to the editors DATE.

Funding: This work was funded by French National Programme LEFE/INSU.

[†]GERAD and Department of Mathematics and Industrial Engineering, Polytechnique Montréal.
(youssef.diouane@polymtl.ca, dominique.orban@polymtl.ca).

[‡]CERFACS / CECI CNRS UMR 5318, Toulouse, France. (gurol@cerfacs.fr, mouhtal@cerfacs.fr).

used with k pairs has shown that LMP can cluster at least k eigenvalues at 1, and that the eigenvalues of the preconditioned matrix interlace with those of the original matrix [11]. The efficiency of this approach has been demonstrated in a real-life data assimilation applications [11, 24].

We focus on improving the performance of the *spectral LMP* [7, 11], which is built by using eigenpairs of $A^{(j)}$. The spectral LMP shares the same formulation as the abstract balancing domain decomposition method [18] and is equivalent to deflation-based preconditioning when used with a specific initial point [24].

When designing preconditioners for PCG, the primary focus in the literature is mostly on A and the significance of the initial guess is overlooked. Although the importance of the initial guess is mentioned, its impact on the choice of a preconditioner is not well studied. Favorable eigenvalue distributions are also highlighted in terms of clustering, but there is little emphasis on the position of the clusters. The performance of the preconditioner is also measured in terms of the total number of iterations to converge, with little focus on the convergence in the early iterations. When PCG is truncated before convergence due to computational budget or when used as a solver within a optimization method like GN, the effect of the preconditioner on the early convergence of PCG is also crucial. In this paper, we aim to explore those overlooked aspects to design a good preconditioner. We not only aim to improve convergence by reducing the total number of iterations but also ensure that, from the very first iteration, the preconditioned iterates outperform those of the original system. In doing so, we specifically focus on strategically positioning the eigenvalues captured by the LMP, in that the energy norm of the error at each iteration of CG is reduced.

The paper is organized as follows. In [Section 2](#) we start by introducing the necessary notation. In [Section 3](#), we review PCG and its convergence properties. We then discuss the characteristics of an efficient preconditioner that can be applied to (1.1). [Section 4](#) is our main contribution. We define the *scaled* spectral preconditioner and discuss its properties. Next, we outline three key approaches for selecting the scaling parameter, which influences the positioning of the eigenvalue cluster, to reduce total number of iterations and enhance convergence in the early iterations. In [Section 5](#), we provide numerical experiments using the Lorenz 95 reference model from data assimilation to validate theoretical results. Finally, conclusions and perspectives are discussed in [Section 6](#).

2. Notation. The matrix $A \in \mathbb{R}^{n \times n}$ is always SPD. Its spectral radius is $\rho(A)$. Its spectral decomposition is $A = S\Lambda S^\top$ with $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$, $\lambda_1 \geq \dots \geq \lambda_n > 0$, and $S = [s_1 \ \dots \ s_n]$ orthogonal. Its i -th eigenvalue is $\nu_i(A)$. Its range space is $\mathcal{R}(A)$. The A -norm, or *energy norm*, of vector x is $\|x\|_A = \sqrt{x^\top A x}$. The spectral norm is $\|\cdot\|_2$.

3. Background.

3.1. CG algorithm. The Conjugate Gradient (CG) method [13] is the workhorse for $Ax = b$ with SPD $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$. If $x_0 \in \mathbb{R}^n$ is an initial guess and $r_0 = b - Ax_0$ is the initial residual, then at every step $\ell = 1, 2, \dots, n$, CG produces a unique approximation [22, p.176]

$$(3.1) \quad x_\ell \in x_0 + \mathcal{K}_\ell(A, r_0) \quad \text{such that} \quad r_\ell \perp \mathcal{K}_\ell(A, r_0),$$

which is equivalent [22, p.126] to

$$(3.2) \quad \|x^* - x_\ell\|_A = \min_{x \in x_0 + \mathcal{K}_\ell(A, r_0)} \|x^* - x\|_A,$$

88 where x^* is the exact solution, $\mathcal{K}_\ell(A, r_0) := \text{span}\{r_0, Ar_0, \dots, A^{\ell-1}r_0\}$ is the ℓ -th
 89 Krylov subspace generated by A and r_0 . In exact arithmetic, the method terminates
 90 in at most μ iterations, where μ is the grade of r_0 with respect to A , i.e., the maximum
 91 dimension of the Krylov subspace generated by A and r_0 [22]. The most popular
 92 and computationally efficient variant of (3.1) is the original formulation of [13], that
 93 recursively updates coupled 2-term recurrences for $x_{\ell+1}$, $r_{\ell+1}$, and the search direction
 94 $p_{\ell+1}$. Algorithm 3.1 states the complete algorithm. A common stopping criterion is
 95 based on sufficient decrease of the relative residual norm. However, in practical data
 96 assimilation implementations, a fixed number of iterations is used as stopping criterion
 97 due to computational budget constraints. CG is presented alongside its companion
 98 formulation, Algorithm 3.2, to be detailed in Subsection 3.3.
 99

Algorithm 3.1 CG

```

1:  $r_0 = b - Ax_0$ 
2:
3:  $\rho_0 = r_0^\top r_0$ 
4:  $p_0 = r_0$ 
5: for  $\ell = 0, 1, \dots$  do
6:    $q_\ell = Ap_\ell$ 
7:    $\alpha_\ell = \rho_\ell / (q_\ell^\top p_\ell)$ 
8:    $x_{\ell+1} = x_\ell + \alpha_\ell p_\ell$ 
9:    $r_{\ell+1} = r_\ell - \alpha_\ell q_\ell$ 
10:
11:    $\rho_{\ell+1} = r_{\ell+1}^\top r_{\ell+1}$ 
12:    $\beta_{\ell+1} = \rho_{\ell+1} / \rho_\ell$ 
13:    $p_{\ell+1} = r_{\ell+1} + \beta_{\ell+1} p_\ell$ 
14: end for
    
```

Algorithm 3.2 PCG

```

1:  $\hat{r}_0 = b - A\hat{x}_0$ 
2:  $z_0 = F\hat{r}_0$ 
3:  $\hat{\rho}_0 = \hat{r}_0^\top z_0$ 
4:  $\hat{p}_0 = z_0$ 
5: for  $\ell = 0, 1, \dots$  do
6:    $\hat{q}_\ell = A\hat{p}_\ell$ 
7:    $\hat{\alpha}_\ell = \hat{\rho}_\ell / (\hat{q}_\ell^\top \hat{p}_\ell)$ 
8:    $\hat{x}_{\ell+1} = \hat{x}_\ell + \hat{\alpha}_\ell \hat{p}_\ell$ 
9:    $\hat{r}_{\ell+1} = \hat{r}_\ell - \hat{\alpha}_\ell \hat{q}_\ell$ 
10:   $z_{\ell+1} = F\hat{r}_{\ell+1}$ 
11:   $\hat{\rho}_{\ell+1} = \hat{r}_{\ell+1}^\top z_{\ell+1}$ 
12:   $\hat{\beta}_{\ell+1} = \hat{\rho}_{\ell+1} / \hat{\rho}_\ell$ 
13:   $\hat{p}_{\ell+1} = z_{\ell+1} + \hat{\beta}_{\ell+1} \hat{p}_\ell$ 
14: end for
    
```

100 **3.2. Convergence properties of CG.** The approximation x_ℓ uniquely deter-
 101 mined by (3.1) minimizes the error in the energy norm:

$$\text{102 (3.3)} \quad \|x^* - x_\ell\|_A^2 = \min_{p \in \mathbb{P}_\ell(0)} \|p(A)(x^* - x_0)\|_A^2 = \min_{p \in \mathbb{P}_\ell(0)} \sum_{i=1}^n p(\lambda_i)^2 \frac{\eta_i^2}{\lambda_i},$$

103 where $\eta_i = s_i^\top r_0$ and $\mathbb{P}_\ell(0)$ is the set of polynomials of degree at most ℓ with value 1
 104 at zero [22, p.193]. Thus, at each iteration, CG solves a certain weighted polynomial
 105 approximation problem over the discrete set $\{\lambda_1, \dots, \lambda_n\}$. Moreover, if $z_1^{(\ell)}, \dots, z_\ell^{(\ell)}$
 106 are the ℓ roots of the solution p_ℓ^* to (3.3),

$$\text{107 (3.4)} \quad \|x^* - x_\ell\|_A^2 = \sum_{i=1}^n p_\ell^*(\lambda_i)^2 \frac{\eta_i^2}{\lambda_i} = \sum_{i=1}^n \prod_{j=1}^{\ell} \left(1 - \frac{\lambda_i}{z_j^{(\ell)}}\right)^2 \frac{\eta_i^2}{\lambda_i}.$$

108 The $z_j^{(\ell)}$ are the *Ritz values* [5]. From (3.4), if $z_j^{(\ell)}$ is close to a λ_i , we expect a
 109 significant reduction in the error in energy norm. Based on the above, [5] explains the
 110 rate of convergence of CG in terms of the convergence of the Ritz values to eigenvalues
 111 of A . Assuming that $\lambda_1, \dots, \lambda_n$ take on the r distinct values ρ_1, \dots, ρ_r , CG converges
 112 in at most r iterations [20, Theorem 5.4].

113 Using (3.3) and maximizing over the values $p(\lambda_i)$ [22, p.194] leads to

$$114 \quad (3.5) \quad \frac{\|x^* - x_\ell\|_A}{\|x^* - x_0\|_A} \leq \min_{p \in \mathbb{P}_\ell(0)} \max_{1 \leq i \leq n} |p(\lambda_i)|.$$

115 By replacing $\{\lambda_1, \dots, \lambda_n\}$ with the interval $[\lambda_1, \lambda_n]$ and using Chebyshev polynomials,
116 we obtain an upper bound [22, p.194]:

$$117 \quad (3.6) \quad \frac{\|x^* - x_\ell\|_A}{\|x^* - x_0\|_A} \leq 2 \left(\frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1} \right)^\ell,$$

118 where $\kappa(A) := \lambda_1/\lambda_n$ is the condition number of A . While (3.5) and (3.6) provide the
119 worst-case behavior of CG [12], the convergence properties may vary significantly from
120 the worst case for a specific initial approximation. Note also that upper bounds (3.5)
121 and (3.6) only depend on A , and not on r_0 . Though (3.6) relates the convergence
122 behavior of CG to $\kappa(A)$, one should be careful as convergence is also influenced by
123 the clustering of the eigenvalues and their positioning [2, 3].

124 **3.3. Properties of a good preconditioner.** In many practical applications,
125 a preconditioner is essential for accelerating the convergence of CG [1, 25]. Assume
126 that a preconditioner $F = UU^\top \in \mathbb{R}^{n \times n}$ is available in a factored form, where U is
127 SPD, and consider the system with split preconditioner

$$128 \quad (3.7) \quad U^\top AUy = U^\top b,$$

129 whose matrix is also SPD. System (3.7) can then be solved with CG. The latter
130 updates estimate y_ℓ that can be used to recover $\hat{x}_\ell := Uy_\ell$. Algorithm 3.2, the pre-
131 conditioned conjugate gradients method, is equivalent to the procedure just described,
132 but only involves solves with F and does not assume knowledge of U [8, p.532]. PCG
133 updates \hat{x}_ℓ directly.

PCG looks for an approximate solution in the Krylov subspace

$$x_0 + UK_\ell(U^\top AU, U^\top r_0),$$

134 and as in (3.3), it minimizes the energy norm,

$$135 \quad (3.8) \quad \|x^* - \hat{x}_\ell\|_A = \min_{q \in \mathbb{P}_\ell(0)} \|Uq(U^\top AU)U^{-1}(x^* - x_0)\|_A.$$

136 Although there is no general method for building a good preconditioner [1, 25],
137 leveraging the convergence properties of CG on (3.8) often leads to the following
138 criteria: (i) F should approximate the inverse of A , (ii) F should be cheap to apply,
139 (iii) $\kappa(U^\top AU)$ should be smaller than $\kappa(A)$, and (iv) $U^\top AU$ should have a more
140 favorable distribution of eigenvalues than A . Note that, all four criteria only focus on
141 A and overlook the significance of the initial guess.

142 **3.4. Preconditioning for a sequence of linear systems.** In the context
143 of (1.1), it is common to use a first level preconditioner, $F^{(1)}$, for the initial linear
144 system, $A^{(1)}x^{(1)} = b^{(1)}$. The selection of the first-level preconditioner depends on the
145 problem and may take into account both the physics of the problem and the algebraic
146 structure of $A^{(1)}$ [1, 25, 21]. To further accelerate convergence of an iterative method
147 such as PCG on subsequent linear systems $A^{(j+1)}x^{(j+1)} = b^{(j+1)}$, one can perform

148 a low-rank update of the most-recent preconditioner, $F^{(j)}$, leveraging information
 149 obtained from solving $A^{(j)}x^{(j)} = b^{(j)}$ [17, 11].

150 One common choice of low-rank update is to use the (approximate) spectrum
 151 of $A^{(j)}$ [6, 7, 11]. The main idea is to capture the eigenvalues not captured by the
 152 first-level preconditioner, and cluster them to a positive quantity, typically around 1.

153 In this paper, we will consider the case where only the right-hand side is changing
 154 over the sequence of the linear systems, i.e., $A^{(j)} = A$ for all j . Perturbation analysis
 155 with respect to A will be presented in a forthcoming paper.

156 **4. A scaled spectral preconditioner.** We focus on the scaled spectral precon-
 157 ditioner, known in the literature as the deflating preconditioner [7] or spectral Limited
 158 Memory Preconditioner (LMP) [11], which is defined using a scaling parameter that
 159 determines the positioning of the cluster. We will provide several strategies for the
 160 choice of the scaling parameter, which has a significant impact on the convergence of
 161 PCG.

162 Let us assume that k largest eigenvalues of A , i.e. $\{\lambda_i\}_{i=1}^k$, are available. We
 163 define the spectral preconditioner

$$164 \quad (4.1) \quad F_\theta := I_n + \sum_{i=1}^k \left(\frac{\theta}{\lambda_i} - 1 \right) s_i s_i^\top = I_n + S_k (\theta \Lambda_k^{-1} - I_k) S_k^\top = S \begin{bmatrix} \theta \Lambda_k^{-1} & \\ & I_{n-k} \end{bmatrix} S^\top,$$

165 where $S_k := [s_1 \ \cdots \ s_k]$ and $\Lambda_k := \text{diag}(\lambda_1, \dots, \lambda_k)$. The factor of $F_\theta = U_\theta^2$ is

$$166 \quad (4.2) \quad U_\theta = U_\theta^\top := I_n + \sum_{i=1}^k \left(\sqrt{\frac{\theta}{\lambda_i}} - 1 \right) s_i s_i^\top = S \begin{bmatrix} \sqrt{\theta} \Lambda_k^{-\frac{1}{2}} & \\ & I_{n-k} \end{bmatrix} S^\top.$$

167 Preconditioner F_θ clusters $\lambda_1, \dots, \lambda_k$ around θ , and leaves the rest of the spectrum
 168 untouched, i.e.,

$$169 \quad (4.3) \quad U_\theta A U_\theta = S \begin{bmatrix} \theta I_k & \\ & \bar{\Lambda}_k \end{bmatrix} S^\top = \theta S_k S_k^\top + \bar{S}_k \bar{\Lambda}_k \bar{S}_k^\top,$$

170 where $\bar{S}_k := [s_{k+1} \ \cdots \ s_n]$ and $\bar{\Lambda}_k := \text{diag}(\lambda_{k+1}, \dots, \lambda_n)$. As in (3.8), PCG mini-
 171 mizes

$$172 \quad \|x^* - \hat{x}_\ell(\theta)\|_A = \min_{q \in \mathbb{P}_\ell(0)} \|U_\theta q (U_\theta A U_\theta) U_\theta^{-1} (x^* - x_0)\|_A \\ 173 \quad (4.4) \quad = \min_{q \in \mathbb{P}_\ell(0)} \|q (U_\theta A U_\theta) (x^* - x_0)\|_A,$$

174 where we used $U_\theta q (U_\theta A U_\theta) U_\theta^{-1} = U_\theta U_\theta^{-1} q (U_\theta A U_\theta) = q (U_\theta A U_\theta)$. Using (3.3) in the
 175 context of (4.4), we obtain the following result.

176 **THEOREM 4.1.** *Let $\hat{x}_\ell(\theta)$ be generated at iteration ℓ of Algorithm 3.2 applied to*
 177 *$Ax = b$ with preconditioner (4.1). Then,*

$$178 \quad (4.5) \quad \|x^* - \hat{x}_\ell(\theta)\|_A^2 = \min_{q \in \mathbb{P}_\ell(0)} \sum_{i=1}^k \frac{\eta_i^2}{\lambda_i} q(\theta)^2 + \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} q(\lambda_i)^2,$$

179 where $\eta_i = s_i^\top r_0$ is the i -th component of the initial residual in the basis S .

180 *Proof.* Given (4.3), we have for any polynomial q ,

$$181 \quad q(U_\theta A U_\theta) = Sq \left(\begin{bmatrix} \theta I_k & \\ & \bar{\Lambda}_k \end{bmatrix} \right) S^\top.$$

182 Since $x^* - x_0 = A^{-1}r_0 = S\Lambda^{-1}S^\top r_0$,

$$183 \quad (4.6) \quad q(U_\theta A U_\theta)(x^* - x_0) = Sq \left(\begin{bmatrix} \theta I_k & \\ & \bar{\Lambda}_k \end{bmatrix} \right) \Lambda^{-1}S^\top r_0.$$

184 Substituting (4.6) into (4.4), we obtain the result. \square

185 The scaled LMP (4.2) is typically used with $\theta = 1$. This choice is operational in
186 numerical weather forecast [6, 24]. In the next subsections, we explore various choices
187 for θ aiming to improve convergence properties of PCG.

188 **4.1. On the choice of the scaling parameter.** The scaling parameter θ ,
189 which defines the position of the cluster, is often set to 1 [6, 7, 11]. This choice is
190 motivated by several factors, such as the eigenvalue distribution of A , the behavior of
191 the first-level preconditioner, and the convergence behavior of PCG.

192 We investigate clustering the eigenvalues at a general $\theta > 0$, which, compared
193 with the conventional choice of 1, results in enhanced convergence of PCG. It is
194 important to note that the notion of “better convergence” may vary across different
195 applications. For instance, in some applications, one may require high accuracy, in
196 which case, a better convergence may be defined as a lower number of iterations. In
197 other applications, we may want to get an approximate solution quickly, which requires
198 to improve the convergence especially in the early iterations. In this case, there is
199 no guarantee that the early preconditioned iterates will provide a better reduction
200 in the energy norm compared to the unpreconditioned iterates (Subsection 4.2). For
201 certain applications, such as numerical weather forecast, where PCG is stopped before
202 reaching convergence due to computational budget, early convergence properties could
203 be of critical importance. As a first direction, we will focus on the following question:

204 *Is there $\theta > 0$ such that for any x_0 ,*

$$205 \quad (4.7) \quad \|x^* - \hat{x}_\ell(\theta)\|_A \leq \|x^* - x_\ell\|_A, \quad \ell = 1, \dots, n?$$

206 To accelerate early convergence, we will investigate optimal choices for θ with respect
207 to the error in the energy norm at the first iteration of PCG, i.e.,

$$208 \quad \min_{\theta} \Phi(\theta) := \|x^* - \hat{x}_1(\theta)\|_A^2.$$

209 We focus solely on the first iteration as it allows us to derive the optimal value of θ
210 in closed form.

211 On the other hand, for PCG, it is well known that removing eigenvalues causing
212 convergence delay can improve the convergence rate significantly [6, 11]. This can be
213 done by using deflation techniques, in which the aim is to “hide” (problematic) parts of
214 the spectrum of A from PCG, so that the convergence rate of PCG is improved [14, 23].
215 Finally, our focus will be also on answering the question

216 *Can we choose $\theta > 0$ such that for any x_0 , PCG generates iterates*
217 *close to those of deflation techniques?*

218 **4.2. θ providing lower error in energy norm.** In general, although scaled
 219 spectral preconditioning is expected to help reduce the number of iterations required
 220 to achieve convergence, (4.7) may not hold for any choice of $\theta > 0$ and all iterations
 221 ℓ as given by the following proposition.

222 **PROPOSITION 4.2.** *Let x_1 be the first CG iterate when solving $Ax = b$. Let $\hat{x}_1(\theta)$
 223 be generated at the first iteration of [Algorithm 3.2](#) applied to $Ax = b$ with preconditioner
 224 (4.1). Let x_0 be such that $\eta_i^2 = \lambda_i$ for $i = k, k + 1$ and $\eta_i = 0$ otherwise.
 225 Then,*

$$226 \quad \|x^* - \hat{x}_1(\theta)\|_A^2 \leq \|x^* - x_1\|_A^2 \iff \frac{\lambda_{k+1}^2}{\lambda_k} \leq \theta \leq \lambda_k.$$

227 *Proof.* For $\ell = 1$, (3.4) yields $\|x^* - x_1\|_A^2 = p_1^*(\lambda_k)^2 + p_1^*(\lambda_{k+1})^2$, where

$$228 \quad p_1^*(\lambda) = 1 - \frac{r_0^\top r_0}{r_0^\top A r_0} \lambda = 1 - \frac{\lambda_k + \lambda_{k+1}}{\lambda_k^2 + \lambda_{k+1}^2} \lambda.$$

229 Similarly, (4.5) gives $\|x^* - \hat{x}_1(\theta)\|_A^2 = q_{1,\theta}^*(\theta)^2 + q_{1,\theta}^*(\lambda_{k+1})^2$, where

$$230 \quad q_{1,\theta}^*(\lambda) = 1 - \frac{r_0^\top F_\theta r_0}{r_0^\top F_\theta A F_\theta r_0} \lambda = 1 - \frac{\theta + \lambda_{k+1}}{\theta^2 + \lambda_{k+1}^2} \lambda$$

231 is the polynomial that realizes the minimum. Using these relations, we obtain

$$232 \quad \|x^* - x_1\|_A^2 = \left(1 - \frac{\lambda_k + \lambda_{k+1}}{\lambda_k^2 + \lambda_{k+1}^2} \lambda_k\right)^2 + \left(1 - \frac{\lambda_k + \lambda_{k+1}}{\lambda_k^2 + \lambda_{k+1}^2} \lambda_{k+1}\right)^2 = \frac{(\lambda_k - \lambda_{k+1})^2}{\lambda_k^2 + \lambda_{k+1}^2}$$

233 and

$$234 \quad \|x^* - \hat{x}_1(\theta)\|_A^2 = \left(1 - \frac{\theta + \lambda_{k+1}}{\theta^2 + \lambda_{k+1}^2} \theta\right)^2 + \left(1 - \frac{\theta + \lambda_{k+1}}{\theta^2 + \lambda_{k+1}^2} \lambda_{k+1}\right)^2 = \frac{(\theta - \lambda_{k+1})^2}{\theta^2 + \lambda_{k+1}^2}.$$

235 Hence,

$$236 \quad \frac{(\theta - \lambda_{k+1})^2}{\theta^2 + \lambda_{k+1}^2} \leq \frac{(\lambda_k - \lambda_{k+1})^2}{\lambda_k^2 + \lambda_{k+1}^2} \iff \frac{\lambda_{k+1}^2}{\lambda_k} \leq \theta \leq \lambda_k. \quad \square$$

237 **Proposition 4.2** shows that (4.7) is not satisfied for all $\theta > 0$. If $\theta > 0$ lies outside
 238 of $[\lambda_{k+1}^2/\lambda_k, \lambda_k]$, then $\|x^* - \hat{x}_1(\theta)\|_A > \|x^* - x_1\|_A$ for x_0 as defined in [Proposition 4.2](#).

239 In what comes next, we focus on the properties of θ such that (4.7) is guaranteed
 240 for all iterations ℓ , and for any given x_0 . An intuitive approach is to identify a range
 241 of θ values where the eigenvalue ratios of the preconditioned matrix are less than or
 242 equal to those of the unpreconditioned matrix, as noted in [12, Lemma 1]. The next
 243 lemma shows that this property holds for $\theta \in [\lambda_{k+1}, \lambda_k]$, and for such choice, there
 244 exists a polynomial that promotes favorable PCG convergence.

245 **LEMMA 4.3.** *Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$, $\ell \in \{1, \dots, n\}$, and $k \in \{1, \dots, \ell\}$. For
 246 any $\theta \in [\lambda_{k+1}, \lambda_k]$, and any polynomial p of degree ℓ such that $p(0) = 1$ and whose
 247 roots all lie in $[\lambda_n, \lambda_1]$, there exists a polynomial q of degree ℓ such that $q(0) = 1$ and*

$$248 \quad |q(\theta)| \leq |p(\lambda_i)|, \quad i = 1, \dots, k$$

$$249 \quad |q(\lambda_i)| \leq |p(\lambda_i)|, \quad i = k + 1, \dots, n.$$

250 *Proof.* Let us denote $(\mu_j)_{1 \leq j \leq \ell}$ the roots of the polynomial p given in decreasing
 251 order, so $p(\lambda) = \prod_{i=1}^{\ell} \left(1 - \frac{\lambda}{\mu_i}\right)$ for any $\lambda \geq 0$. Then, three cases may occur:

Case 1: For all $j \in \{1, \dots, \ell\}$, $\mu_j < \theta$, we choose $q(\lambda) = p(\lambda)$, then simply we
 have for $i \in \{k+1, \dots, n\}$, $|q(\lambda_i)| = |p(\lambda_i)|$. For $i \in \{1, \dots, k\}$, using the property
 that $\mu_j < \theta \leq \lambda_i$, we obtain

$$1 - \frac{\lambda_i}{\mu_j} \leq 1 - \frac{\theta}{\mu_j} \leq 0.$$

252 Thus, we have $|1 - \frac{\theta}{\mu_j}| \leq |1 - \frac{\lambda_i}{\mu_j}|$, and consequently $|q(\theta)| \leq |p(\lambda_i)|$.

Case 2: If for all $j \in \{1, \dots, \ell\}$, $\theta \leq \mu_j$, we choose $q(\lambda) = \prod_{j=1}^{\ell} \left(1 - \frac{\lambda}{\theta}\right) =$
 $\left(1 - \frac{\lambda}{\theta}\right)^{\ell}$. Then simply for $i \in \{1, \dots, k\}$, $|q(\theta)| = 0 \leq |p(\lambda_i)|$. For $i \in \{k+1, \dots, n\}$,
 using the property $\lambda_{k+1} \leq \theta \leq \mu_j$, we obtain

$$0 \leq 1 - \frac{\lambda_i}{\lambda_{k+1}} \leq 1 - \frac{\lambda_i}{\theta} \leq 1 - \frac{\lambda_i}{\mu_j}.$$

253 Therefore, for $i = k+1, \dots, n$, $|q(\lambda_i)| \leq |p(\lambda_i)|$.

Case 3: let $s \in \{1, \dots, \ell-1\}$ such that for $j = 1, \dots, s$, $\theta \leq \mu_j \leq \lambda_1$, and for
 $j = s+1, \dots, \ell$, $\lambda_n \leq \mu_j < \theta$. Let's denote

$$q(\lambda) = \prod_{j=1}^s \left(1 - \frac{\lambda}{\theta}\right) \prod_{j=s+1}^{\ell} \left(1 - \frac{\lambda}{\mu_j}\right) = \left(1 - \frac{\lambda}{\theta}\right)^s \prod_{j=s+1}^{\ell} \left(1 - \frac{\lambda}{\mu_j}\right).$$

We have $q(\theta) = 0$, so $|q(\theta)| \leq |p(\lambda_i)|$ for $i \in \{1, \dots, k\}$. For $i \in \{k+1, \dots, n\}$ and
 $j \in \{1, \dots, s\}$, we have

$$0 \leq 1 - \frac{\lambda_i}{\lambda_{k+1}} \leq 1 - \frac{\lambda_i}{\theta} \leq 1 - \frac{\lambda_i}{\mu_j},$$

254 because $\lambda_{k+1} \leq \theta \leq \mu_j$. Therefore, for $i = k+1, \dots, n$, $|q(\lambda_i)| \leq |p(\lambda_i)|$. \square

255 Now, we can present a result that enables comparing the error in energy norm between
 256 the preconditioned system given by (3.7) and the unpreconditioned system, $Ax = b$.

257 **THEOREM 4.4.** *Let $(x_{\ell})_{\ell \in \{1, \dots, n\}}$ and $\hat{x}_{\ell}(\theta)_{\ell \in \{1, \dots, n\}}$ be the sequences generated by*
 258 *CG and PCG with F_{θ} with $\theta \in [\lambda_{k+1}, \lambda_k]$, respectively, when solving $Ax = b$. Assume*
 259 *that $\hat{x}_0(\theta) = x_0$. Then, for all $\ell = 1, \dots, n$, $\|x^* - \hat{x}_{\ell}(\theta)\|_A \leq \|x^* - x_{\ell}\|_A$.*

260 *Proof.* Let $\ell \in \{1, \dots, n\}$. From (3.4),

261 (4.8)
$$\|x^* - x_{\ell}\|_A^2 = \min_{p \in \mathbb{P}_{\ell}(0)} \|p_{\ell}(A)(x^* - x_0)\|_A^2 = \sum_{i=1}^n \frac{\eta_i^2}{\lambda_i} p_{\ell}^*(\lambda_i)^2,$$

262 where η_i represents the components of the initial residual $r_0 = b - Ax_0$ in the
 263 eigenspace of A . Applying Lemma 4.3 to p_{ℓ}^* , there exists a polynomial q of degree ℓ
 264 with $q(0) = 1$ such that

265
$$|q(\theta)| \leq |p_{\ell}^*(\lambda_i)|, \quad i \in \{1, \dots, k\}$$

 266
$$|q(\lambda_i)| \leq |p_{\ell}^*(\lambda_i)|, \quad i \in \{k+1, \dots, n\}.$$

267 Applying these inequalities to (4.8) yields

$$\begin{aligned}
 268 \quad \|x^* - x_\ell\|_A^2 &= \sum_{i=1}^n \frac{\eta_i^2}{\lambda_i} p_\ell^*(\lambda_i)^2 \geq \sum_{i=1}^k \frac{\eta_i^2}{\lambda_i} q(\theta)^2 + \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} q(\lambda_i)^2 \\
 269 \quad &\geq \min_{q \in \mathbb{P}_\ell(0)} \left(\sum_{i=1}^k \frac{\eta_i^2}{\lambda_i} q(\theta)^2 + \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} q(\lambda_i)^2 \right) = \|x^* - \hat{x}_\ell(\theta)\|_A^2. \quad \square
 \end{aligned}$$

270 **Theorem 4.4** offers a range of choices for θ . Next, we discuss the practical and
 271 theoretical choices from this range. Let us remind that to construct the spectral
 272 LMP (4.2), we are given k eigenpairs. As a result, one practical choice is $\theta = \lambda_k$.
 273 This idea is summarized in the following corollary.

274 **COROLLARY 4.5.** *Let $\theta = \lambda_k$. Then, $\|x^* - \hat{x}_\ell(\lambda_k)\|_A \leq \|x^* - x_\ell\|_A$ for any x_0
 275 and for all $\ell \in \{1, \dots, n\}$.*

276 The next theorem shows that increasing k results in improved convergence.

277 **THEOREM 4.6.** *Let $1 < k_1 \leq k_2 < n$ and $\theta_{k_1} \in [\lambda_{k_1+1}, \lambda_{k_1}]$, $\theta_{k_2} \in [\lambda_{k_2+1}, \lambda_{k_2}]$
 278 with, $\theta_{k_2} \leq \theta_{k_1}$. Let $(\hat{x}_\ell(\theta_{k_1}))_{\ell \in \{1, \dots, n\}}$, $(\hat{x}_\ell(\theta_{k_2}))_{\ell \in \{1, \dots, n\}}$ be the sequences obtained
 279 from PCG iterates when solving $Ax = b$ using $F_{\theta_{k_1}}$ and $F_{\theta_{k_2}}$ respectively with an
 280 arbitrary initial guess x_0 . Then, for all $\ell \in \{1, \dots, n\}$, one has:*

$$281 \quad \|x^* - \hat{x}_\ell(\theta_{k_2})\|_A \leq \|x^* - \hat{x}_\ell(\theta_{k_1})\|_A.$$

282 *Proof.* The eigenvalues of the preconditioned matrix using $F_{\theta_{k_1}}$ and $F_{\theta_{k_2}}$ are given
 283 in decreasing order respectively as

$$284 \quad \rho_i = \begin{cases} \theta_{k_1} & i \in \{1, \dots, k_1\} \\ \lambda_i & \text{otherwise,} \end{cases} \quad \text{and} \quad \tilde{\rho}_i = \begin{cases} \theta_{k_2} & i \in \{1, \dots, k_2\} \\ \lambda_i & \text{otherwise.} \end{cases}$$

285 As $k_1 < k_2$, it follows that $\tilde{\rho}_{k_2} \leq \rho_{k_1} = \theta_{k_1}$. Therefore, $\tilde{\rho}_i$ can be expressed as a
 286 function of ρ_i as

$$287 \quad \tilde{\rho}_i = \begin{cases} \theta_{k_2} \in [\rho_{k_2+1}, \rho_{k_2}] & i \in \{1, \dots, k_2\} \\ \rho_i & \text{otherwise.} \end{cases}$$

288 Using **Lemma 4.3**, for the polynomial $q_{\ell, \theta_{k_1}}^*$, there exists a polynomial q of degree ℓ
 289 with $q(0) = 1$, such that for $i \in \{1, \dots, n\}$,

$$\begin{aligned}
 290 \quad |q(\theta_{k_2})| &\leq |q_{\ell, \theta_{k_1}}^*(\rho_i)|, \quad i \in \{1, \dots, k_2\} \\
 291 \quad |q(\rho_i)| &\leq |q_{\ell, \theta_{k_1}}^*(\rho_i)|, \quad i \in \{k_2 + 1, \dots, n\}
 \end{aligned}$$

292

293 Applying this result to (4.5) yields that

$$\begin{aligned}
 294 \quad \|x^* - \hat{x}_\ell(\theta_{k_1})\|_A^2 &= \sum_{i=1}^n \frac{\eta_i^2}{\lambda_i} q_{\ell, \theta_{k_1}}^*(\rho_i)^2 \\
 295 \quad &\geq \sum_{i=1}^{k_2} \frac{\eta_i^2}{\lambda_i} q(\theta_{k_2})^2 + \sum_{i=k_2+1}^n \frac{\eta_i^2}{\lambda_i} q(\rho_i)^2 \\
 296 \quad &\geq \min_{q \in \mathbb{P}_\ell(0)} \left(\sum_{i=1}^{k_2} \frac{\eta_i^2}{\lambda_i} q(\theta_{k_2})^2 + \sum_{i=k_2+1}^n \frac{\eta_i^2}{\lambda_i} q(\lambda_i)^2 \right) = \|x^* - \hat{x}_\ell(\theta_{k_2})\|_A^2. \quad \square
 \end{aligned}$$

297 One can see that $k_1 < k_2 \implies \theta_{k_2} \leq \theta_{k_1}$, since λ_i are in decreasing order. In addition,
 298 when $k_1 = k_2$, [Theorem 4.6](#) shows that λ_{k_1+1} is the best choice in $[\lambda_{k_1+1}, \lambda_{k_1}]$ in terms
 299 of reducing the error with respect to the unpreconditioned system.

300 **4.3. Optimal choice for θ with respect to the initial residual.** Our ob-
 301 jective is to determine the value of θ that minimizes the energy norm of the error at
 302 the initial iterate. This will provide us with the optimal reduction at the first iterate,
 303

$$304 \quad (4.9) \quad \theta_r \in \arg \min_{\theta > 0} \Phi(\theta) := \|x^* - \hat{x}_1(\theta)\|_A^2.$$

305 The expression for θ_r is stated in the following theorem.

306 **THEOREM 4.7.** *Let $r_0 = b - Ax_0$. The unique $\lambda_n \leq \theta_r \leq \lambda_{k+1}$ satisfying (4.9) is*

$$307 \quad (4.10) \quad \theta_r := \frac{\sum_{i=k+1}^n \lambda_i \eta_i^2}{\sum_{i=k+1}^n \eta_i^2} = \frac{r_0^\top Ar_0 - r_0 S_k \Lambda_k S_k^\top r_0}{r_0^\top r_0 - r_0^\top S_k S_k^\top r_0}.$$

308 *Proof.* First, [Theorem 4.1](#) implies

$$309 \quad (4.11) \quad \|x^* - \hat{x}_1(\theta)\|_A^2 = \sum_{i=1}^k \frac{\eta_i^2}{\lambda_i} q_{1,\theta}^*(\theta)^2 + \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} q_{1,\theta}^*(\lambda_i)^2$$

310 where $\eta_i = s_i^\top r_0$ and $q_{1,\theta}^*(\lambda) = 1 - \frac{r_0^\top F_\theta r_0}{r_0^\top F_\theta A F_\theta r_0} \lambda$. Using (4.1), we obtain

$$311 \quad (4.12) \quad r_0^\top F_\theta r_0 = \theta \sum_{i=1}^k \frac{\eta_i^2}{\lambda_i} + \sum_{i=k+1}^n \eta_i^2 \quad \text{and} \quad r_0^\top F_\theta A F_\theta r_0 = \theta^2 \sum_{i=1}^k \frac{\eta_i^2}{\lambda_i} + \sum_{i=k+1}^n \lambda_i \eta_i^2.$$

312 Then, for all $\theta > 0$, $\Phi(\theta)$ simplifies to

$$313 \quad \Phi(\theta) = a_1 \left(\frac{a_2 \theta - a_3}{a_1 \theta^2 + a_3} \right)^2 + \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} \left(1 - \frac{a_1 \theta + a_2}{a_1 \theta^2 + a_3} \lambda_i \right)^2,$$

314 where $a_1 = \sum_{i=1}^k \frac{\eta_i^2}{\lambda_i}$, $a_2 = \sum_{i=k+1}^n \eta_i^2$ and $a_3 = \sum_{i=k+1}^n \lambda_i \eta_i^2$. The derivative of Φ is

$$315 \quad \Phi'(\theta) = \frac{2a_1}{(a_1 \theta^2 + a_3)^3} (a_2 \theta - a_3) (a^2 \theta^3 + a_1 a_2 \theta^2 + a_1 a_3 \theta + a_2 a_3).$$

316 Since $\Phi'(\theta) < 0$ on $]0, \frac{a_3}{a_2}[$ and $\Phi'(\theta) > 0$ on $]\frac{a_3}{a_2}, +\infty[$, then $\frac{a_3}{a_2}$ is the global minimizer
 317 of Φ on \mathbb{R}_+^* and is unique. Hence,

$$318 \quad \theta_r = \arg \min_{\theta > 0} \Phi(\theta) = \frac{a_3}{a_2} = \frac{\sum_{i=k+1}^n \lambda_i \eta_i^2}{\sum_{i=k+1}^n \eta_i^2}.$$

319 Moreover,

$$320 \quad \lambda_n = \frac{\sum_{i=k+1}^n \lambda_n \eta_i^2}{\sum_{i=k+1}^n \eta_i^2} \leq \theta_r \leq \frac{\sum_{i=k+1}^n \lambda_i \eta_i^2}{\sum_{i=k+1}^n \eta_i^2} = \lambda_{k+1}.$$

321 The expression for θ_r can be rewritten in terms of S_k , Λ_k , and r_0 as follows:

$$322 \quad \theta_r = \frac{\sum_{i=1}^n \lambda_i \eta_i^2 - \sum_{i=1}^k \lambda_i \eta_i^2}{\sum_{i=1}^n \eta_i^2 - \sum_{i=1}^k \eta_i^2} = \frac{r_0^\top Ar_0 - r_0 S_k \Lambda_k S_k^\top r_0}{r_0^\top r_0 - r_0^\top S_k S_k^\top r_0}. \quad \square$$

Note that θ_r can be interpreted as the center of mass for the remaining part of the spectrum in which the weights are determined by η_i^2 , i.e.

$$\sum_{i=k+1}^n \eta_i^2 (\theta_r - \lambda_i) = 0.$$

323 Let us now look at the first iterate,

$$324 \quad (4.13) \quad \hat{x}_1(\theta_r) = x_0 + \frac{r_0^\top F_{\theta_r} r_0}{r_0^\top F_{\theta_r} A F_{\theta_r} r_0} F_{\theta_r} r_0,$$

325 to better understand the effect of θ_r . Using (4.12) and the value of θ_r ,

$$326 \quad \frac{r_0^\top F_{\theta_r} r_0}{r_0^\top F_{\theta_r} A F_{\theta_r} r_0} = \frac{\sum_{i=k+1}^n \eta_i^2}{\sum_{i=k+1}^n \lambda_i \eta_i^2} = \frac{1}{\theta_r}.$$

327 Therefore, (4.13) simplifies to

$$328 \quad \hat{x}_1(\theta_r) = x_0 + \frac{1}{\theta_r} (\bar{S}_k \bar{S}_k^\top + \theta_r S_k \Lambda_k^{-1} S_k^\top) r_0 = x_0 + S_k \Lambda_k^{-1} S_k^\top r_0 + \frac{1}{\theta_r} \bar{S}_k \bar{S}_k^\top r_0.$$

329 Then, the residual of the first iteration is given by

$$330 \quad (4.14) \quad b - A \hat{x}_1(\theta_r) = r_0 - S_k S_k^\top r_0 - \frac{1}{\theta_r} \bar{S}_k \bar{\Lambda}_k \bar{S}_k^\top r_0 = \bar{S}_k \bar{S}_k^\top r_0 - \frac{1}{\theta_r} \bar{S}_k \bar{\Lambda}_k \bar{S}_k^\top r_0.$$

331 Given (4.14), we conclude that, from the first iteration, we can remove all components
 332 of the residual with respect to S_k , see [Appendix A](#). We now provide an upper bound
 333 for the error in the energy norm for later iterations, $\ell > 1$, beginning with $\hat{x}_1(\theta_r)$.
 334 With this initial point, we ensure that all iterates yield a residual within $\text{Span}(S_k)$.

335 **THEOREM 4.8.** *Let $\hat{x}_\ell(\theta_r)$ be the ℓ -th iterate obtained from PCG when solving*
 336 *$Ax = b$ using the preconditioner F_{θ_r} with an arbitrary initial guess x_0 . Let x_ℓ^{Init} be*
 337 *the ℓ -th iterate generated by CG for solving $Ax = b$ starting from $\hat{x}_1(\theta_r)$ as defined*
 338 *in (4.13). Then, for all $\ell \in \{1, \dots, n\}$, $\|x^* - \hat{x}_{\ell+1}(\theta_r)\|_A \leq \|x^* - x_\ell^{\text{Init}}\|_A$.*

339 *Proof.* From (4.14), the components of $b - A \hat{x}_1(\theta_r)$ in the eigenspace of A are

$$340 \quad 0 \quad (i = 1, \dots, k), \quad \text{and} \quad \eta_i (1 - \lambda_i / \theta_r) \quad (i > k).$$

341 Thus,

$$342 \quad (4.15) \quad \|x - x_\ell^{\text{Init}}\|_A^2 = \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} \left(1 - \frac{\lambda_i}{\theta_r}\right)^2 p_\ell^{*,\text{Init}}(\lambda_i)^2,$$

343 where $p_\ell^{*,\text{Init}}$ is the polynomial that minimizes $p \mapsto \|p(A)(x^* - \hat{x}_1(\theta_r))\|_A^2$ over $\mathbb{P}_\ell(0)$.

344 Define

$$345 \quad \bar{q}(\lambda) = \left(1 - \frac{\lambda}{\theta_r}\right) p_\ell^{*,\text{Init}}(\lambda),$$

346 and note that $\bar{q} \in \mathbb{P}_\ell(0)$. Now we have

$$\begin{aligned}
347 \quad \|x^* - \hat{x}_{\ell+1}(\theta_r)\|_A^2 &= \min_{q \in \mathbb{P}_{\ell+1}(0)} \left(\sum_{i=1}^k \frac{\eta_i^2}{\lambda_i} q(\theta_r)^2 + \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} q(\lambda_i)^2 \right) \\
348 \quad &\leq \sum_{i=1}^k \frac{\eta_i^2}{\lambda_i} \bar{q}(\theta_r)^2 + \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} \bar{q}(\lambda_i)^2 \\
349 \quad &= \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} \left(1 - \frac{\lambda_i}{\theta_r} \right)^2 p_{\ell}^{*,\text{Init}}(\lambda_i)^2 = \|x - x_{\ell}^{\text{Init}}\|_A^2. \quad \square
\end{aligned}$$

350 Note that, one can interpret $\hat{x}_1(\theta_r)$ as the first iteration of CG when solving the
351 unpreconditioned system, starting from $x_0 + S_k \Lambda_k^{-1} S_k^\top r_0$, since the search direction
352 at the first iteration is equal to:

$$353 \quad (4.16) \quad b - A(x_0 + S_k \Lambda_k^{-1} S_k^\top r_0) = b - Ax_0 - S_k S_k^\top r_0 = r_0 - S_k S_k^\top r_0 = \bar{S}_k \bar{S}_k^\top r_0,$$

354 and the step-length α_0 is given as

$$355 \quad \alpha_0 = \frac{1}{\theta_r} = \frac{r_0^\top \bar{S}_k \bar{S}_k^\top r_0}{r_0^\top \bar{S}_k^\top \bar{S}_k A \bar{S}_k \bar{S}_k^\top r_0}.$$

356 This highlights the strong connection between preconditioning, CG with different
357 initial point and deflation techniques [23, 24]. This connection will be explored in
358 detail in the next subsection, providing another choice for the scaling parameter.

359 **4.4. θ as the mid-range between λ_k and λ_n .** We focus now on choosing a
360 scaling parameter θ to obtain approximate iterates to those of deflated CG (see [Algorithm A.1](#)).
361 The deflation technique, with S_k as the deflation subspace, is equivalent
362 to standard CG applied to $Ax = b$ with initial guess

$$363 \quad x_0^{\text{Def}} = x_0 + S_k \Lambda_k^{-1} S_k^\top (b - Ax_0).$$

364 From (4.16), the residual of x_0^{Def} is given as

$$365 \quad b - Ax_0^{\text{Def}} = \bar{S}_k \bar{S}_k^\top r_0.$$

366 One can see that this initial guess gives a residual which is an orthogonal projection
367 of r_0 onto $\text{span}(\bar{S}_k)$, so that the ℓ -th iterate of CG, x_ℓ^{Def} , starting with x_0^{Def} satisfies

$$368 \quad \|x^* - x_\ell^{\text{Def}}\|_A^2 = \min_{q \in \mathbb{P}_\ell(0)} \left(\sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} q(\lambda_i)^2 \right).$$

369 We now provide the main result of this section.

370 **THEOREM 4.9.** *Let $\hat{x}_\ell(\theta)$ be the ℓ -th iterate obtained from PCG iterates when*
371 *solving $Ax = b$ using F_θ starting from an arbitrary initial guess $x_0 \in \mathbb{R}^n$. Let x_ℓ^{Def}*
372 *be the ℓ -th iterate generated with CG when solving $Ax = b$ starting with $x_0^{\text{Def}} =$*
373 *$x_0 + S_k \Lambda_k^{-1} S_k^\top (b - Ax_0)$. Then, in exact arithmetic,*

$$374 \quad (4.17) \quad \left\| x^* - x_{\ell+1}^{\text{Def}} \right\|_A \leq \|x^* - \hat{x}_{\ell+1}(\theta)\|_A \leq \frac{\alpha(\theta)}{\theta} \left\| x^* - x_\ell^{\text{Def}} \right\|_A,$$

375 *with $\alpha(\theta) = \max(|\lambda_{k+1} - \theta|, |\theta - \lambda_n|)$.*

376 *Proof.* Let us start by showing the first inequality. From [Theorem 4.1](#)

$$\begin{aligned}
 377 \quad \|x^* - \hat{x}_{\ell+1}(\theta)\|_A^2 &= \sum_{i=1}^k \frac{\eta_i^2}{\lambda_i} q_{\ell+1,\theta}^*(\theta)^2 + \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} q_{\ell+1,\theta}^*(\lambda_i)^2 \\
 378 \quad &\geq \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} q_{\ell+1,\theta}^*(\lambda_i)^2 \\
 379 \quad &\geq \min_{q \in \mathbb{P}_{\ell+1}(0)} \left(\sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} q(\lambda_i)^2 \right) = \|x^* - x_{\ell+1}^{\text{Def}}\|_A^2.
 \end{aligned}$$

Now, to prove the second inequality, we consider $p_\ell^{*,\text{Def}}$ the polynomial that minimizes $p \mapsto \|p(A)(x^* - x_0^{\text{Def}})\|_A^2$ over $\mathbb{P}_\ell(0)$, i.e.,

$$\|x^* - x_\ell^{\text{Def}}\|_A^2 = \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} p_\ell^{*,\text{Def}}(\lambda_i)^2.$$

380 Consider $\tilde{q}_{\ell+1} \in \mathbb{P}_{\ell+1}(0)$ such as for all $\lambda \in \mathbb{R}$, $\tilde{q}_{\ell+1}(\lambda) = \left(1 - \frac{\lambda}{\theta}\right) p_\ell^{*,\text{Def}}(\lambda)$. Hence,

$$\begin{aligned}
 381 \quad \|x^* - \hat{x}_{\ell+1}(\theta)\|_A^2 &= \sum_{i=1}^k \frac{\eta_i^2}{\lambda_i} q_{\ell+1,\theta}^*(\theta)^2 + \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} q_{\ell+1,\theta}^*(\lambda_i)^2 \\
 382 \quad &\leq \sum_{i=1}^k \frac{\eta_i^2}{\lambda_i} \tilde{q}_{\ell+1}(\theta)^2 + \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} \tilde{q}_{\ell+1}(\lambda_i)^2 \\
 383 \quad &= \sum_{i=k+1}^n \frac{\eta_i^2}{\lambda_i} p_\ell^{\text{Def},*}(\lambda_i) \left(1 - \frac{\lambda_i}{\theta}\right)^2 \\
 384 \quad &\leq \max_{k+1 \leq i \leq n} \left(1 - \frac{\lambda_i}{\theta}\right)^2 \|x^* - x_\ell^{\text{Def}}\|_A^2 = \frac{\alpha(\theta)}{\theta} \|x^* - x_\ell^{\text{Def}}\|_A^2. \quad \square
 \end{aligned}$$

385 Choosing $\theta > 0$ such that $\alpha(\theta)/\theta > 1$ in [\(4.17\)](#) would give a pessimistic upper
 386 bound. For a better bound, we select $\theta > 0$ such that $\alpha(\theta)/\theta \leq 1$, which is equivalent
 387 to impose $\theta \geq \lambda_{k+1}/2$. The value of θ that minimizes $\alpha(\theta)/\theta$ is $\theta^* = (\lambda_{k+1} + \lambda_n)/2$.

388 Given that λ_{k+1} is unknown, and λ_n can be predetermined in various applications,
 389 e.g., in data assimilation problems $\lambda_n = 1$, a practical approach for selecting θ (the
 390 closest to θ^*) is by choosing the average between the λ_k and λ_n , i.e., $\theta_m = (\lambda_k + \lambda_n)/2$,
 391 for which we have $\alpha(\theta_m)/\theta_m = (\lambda_k - \lambda_n)/(\lambda_k + \lambda_n) < 1$. Note that the choice $\theta = \lambda_k$
 392 yields in [\(4.17\)](#) to a worst upper bound compared to θ_m , i.e., $\alpha(\lambda_k)/\lambda_k > \alpha(\theta_m)/\theta_m$.

393 **4.5. Discussion.** The analysis in this section raises two key questions. The first
 394 is: why use a scaled spectral preconditioner when we know that deflated CG iterations
 395 using the deflated subspace S_k , or using an initial guess as defined in [\(4.13\)](#), produce
 396 better results in exact arithmetic (see [Theorem 4.9](#))? The assumption in this section
 397 is that the eigenpairs used to construct the deflated subspace or the initial guess are
 398 exact, ensuring that components of the initial residual within the eigenspace of S_k are
 399 eliminated. However, when an approximate eigen-spectrum is used, such as the eigen-
 400 spectrum of A is applied to solve a system involving a perturbed matrix, \tilde{A} , the initial
 401 guess may fail to remove the components of the initial residual within the eigenspace

402 of \tilde{A} . For instance, consider the perturbed matrix $\tilde{A} = A + E$, A is modified by a
 403 small perturbation matrix E . This results in the following expression:

$$404 \quad b - \tilde{A}x_0^{\text{Def}} = b - Ax_0^{\text{Def}} + Ex_0^{\text{Def}},$$

405 where the value of $b - Ax_0^{\text{Def}}$ from (4.16) becomes: $b - \tilde{A}x_0^{\text{Def}} = \bar{S}_k \bar{S}_k^\top (b - Ax_0) + Ex_0^{\text{Def}}$.

406 This illustrates that the perturbation E introduces additional components to the
 407 residual, which the initial guess fails to fully eliminate, unlike in the exact case. When
 408 the perturbation exists, we show in numerical experiments that using a scaled spectral
 409 LMP becomes advantageous over deflated CG.

410 The second question is: why not combine the initial guess (4.13) with the scaled
 411 spectral LMP using $\theta = 1$. When the initial guess fails to eliminate components
 412 of the initial residual within the eigenspace of \tilde{A} , these components influence the
 413 convergence of PCG. Their impact on the energy norm of the error can be reduced
 414 by appropriately positioning the largest eigenvalues.

415 **5. Numerical Experiments.** In this section, we illustrate the performance of
 416 the scaled spectral LMP, as defined in (4.2), within the context of a nonlinear weighted
 417 least-squares problem arising in data assimilation, i.e.,

$$418 \quad (5.1) \quad \min_{w_0 \in \mathbb{R}^n} f(w_0) = \min_{w_0 \in \mathbb{R}^n} \frac{1}{2} \|w_0 - w_b\|_{B^{-1}}^2 + \frac{1}{2} \sum_{i=1}^{N_t} \|y_i - \mathcal{H}_i(\mathcal{M}_{t_0, t_i}(w_0))\|_{R_i^{-1}}^2.$$

419 Here, $w_0 = w(t_0)$, is the state at the initial time t_0 , for instance temperature value,
 420 $w_b \in \mathbb{R}^n$ is a priori information at time t_0 and $y_i \in \mathbb{R}^{m_i}$ represents the observation
 421 vector at time t_i for $i = 1, \dots, N_t$. $\mathcal{M}_{t_0, t_i}(\cdot)$ is a nonlinear physical dynamical model
 422 which propagates the state w_0 at time t_0 to the the state w_i at time t_i by solving
 423 the partial differential equations. $\mathcal{H}_i(\cdot)$ maps the state vector w_i to a m_i -dimensional
 424 vector representing the state vector in the observation space. $B \in \mathbb{R}^{n \times n}$, $R_i \in \mathbb{R}^{m_i \times m_i}$
 425 are symmetric positive definite error covariance matrices corresponding to the a priori
 426 and observation model error, respectively.

427 The TGN method [10] is widely used to solve the nonlinear optimization problem
 428 (5.1). At each iteration j of the TGN method, the linearized least-squares approx-
 429 imation to the nonlinear least-squares problem (5.1) is solved. This quadratic
 430 cost function at the j -th iterate is formulated as

$$431 \quad (5.2) \quad Q^{(j)}(s) = \frac{1}{2} \left\| s - (w_b - w_0^{(j)}) \right\|_{B^{-1}}^2 + \frac{1}{2} \sum_{i=1}^{N_t} \|G_i^{(j)} s_i - d_i^{(j)}\|_{R_i^{-1}}^2,$$

432 where $s \in \mathbb{R}^n$, $d_i^{(j)} = y_i - \mathcal{G}_i(w_0^{(j)})$ with $\mathcal{G}_i(w_0^{(j)}) = \mathcal{H}_i(\mathcal{M}_{t_0, t_i}(w_0^{(j)}))$ and $G_i^{(j)}$
 433 represents the Jacobian of \mathcal{G}_i at a given iterate $w_0^{(j)}$. The quadratic cost function (5.2)
 434 is minimized with respect to s which is then used to update the current iterate, i.e.
 435 $w_0^{(j+1)} = w_0^{(j)} + s^{(j)}$, where $s^{(j)}$ is an approximate solution of the problem (5.2). This
 436 process continues till the convergence criterion is met. For large scale problems with
 437 computationally expensive models $\mathcal{M}_{t_0, t_i}(\cdot)$, a limited number of TGN iterations are
 438 applied. The solution to the quadratic problem (5.2) can be found by solving

$$439 \quad (5.3) \quad \left(B^{-1} + (G^{(j)})^\top R^{-1} G^{(j)} \right) s = B^{-1} (w_b - w_0^{(j)}) - (G^{(j)})^\top R^{-1} d^{(j)}.$$

440 where $d^{(j)}$ is a m -dimensional concatenated vector of $d_i^{(j)}$ with $m = \sum_{i=1}^{N_t} m_i$, $G^{(j)} \in$
 441 $\mathbb{R}^{m \times n}$ represents a concatenation of $G_i^{(j)} \in \mathbb{R}^{m_i \times n}$, and $R \in \mathbb{R}^{m \times m}$ is a block diagonal

442 matrix, i.e. $R = \text{diag}(R_1, \dots, R_N)$. The matrix $B^{-1} + (G^{(j)})^\top R^{-1} G^{(j)}$ is SPD,
 443 matrix-vector products with it are accessible only through operators, and n can be
 444 large for data assimilation problems. Hence, CG is widely used to solve such systems.

445 Let us assume that a square root factorization of $B = LL^\top$ is available. The linear
 446 system (5.3) can be then preconditioned by using this *first-level* split preconditioner,

$$447 \quad (5.4) \quad \left(I_n + L^\top (G^{(j)})^\top R^{-1} G^{(j)} L \right) x = L^\top \left(B^{-1}(w_b - w_0^{(j)}) - (G^{(j)})^\top R^{-1} d^{(j)} \right).$$

448 CG at the ℓ -th iteration provides an approximate solution $x_\ell^{(j)}$ which is then used
 449 to obtain an approximate solution of the linear system (5.3), i.e. $s_\ell^{(j)} = Lx_\ell^{(j)}$. In
 450 operational data assimilation problems, in general $m \ll n$. Consequently, the pre-
 451 conditioned matrix $A^{(j)} = I_n + L^\top (G^{(j)})^\top R^{-1} G^{(j)} L$ has $n - m$ eigenvalues clustered
 452 around 1, while the remaining eigenvalues are greater than 1.

453 Since in the context of TGN, a sequence of closely related linear systems is
 454 solved, it is common to update the first-level preconditioner L by using approxi-
 455 mate eigenspectrum of the previous linear system [6, 11]. Let us denote $b^{(j)} :=$
 456 $L^\top \left(B^{-1}(w_b - w_0^{(j)}) - (G^{(j)})^\top R^{-1} d^{(j)} \right)$. For $j = 1$, CG Algorithm 3.1 solves the lin-
 457 ear system $A^{(1)}x = b^{(1)}$, for the variable x . Using the recurrences of CG, we can
 458 easily compute approximate eigenpairs of $A^{(1)}$ (see [22, p.174] for more details).
 459 These pairs can then be used to construct a second-level preconditioner, $U_{\theta_1}^{(1)}$, by
 460 using the formula (4.2). Consequently, $(U_{\theta_1}^{(1)})^2$ is an approximation to the inverse of
 461 the matrix $A^{(1)}$. Then, assuming that $A^{(2)}$ is close to the matrix $A^{(1)}$, for $j = 2$,
 462 CG Algorithm 3.1 is applied to the preconditioned system, $U_{\theta_1}^{(1)} A^{(2)} U_{\theta_1}^{(1)} x = U_{\theta_1}^{(1)} b^{(2)}$.
 463 The approximate solution at ℓ -iterate is obtained from the relation $s_\ell^{(2)} = LU_{\theta_1}^{(1)} x_\ell^{(2)}$.
 464 At the end of the CG, we can obtain approximate eigenpairs of $U_{\theta_1}^{(1)} A^{(2)} U_{\theta_1}^{(1)}$ and use
 465 it to construct a preconditioner for the next linear system. At the j -th outer loop of
 466 TGN, CG is applied to the preconditioned linear system:

$$467 \quad (5.5) \quad \left(U_{\theta_{j-1}}^{(j-1)} \dots U_{\theta_1}^{(1)} A^{(j)} U_{\theta_1}^{(1)} \dots U_{\theta_{j-1}}^{(j-1)} \right) x = U_{\theta_{j-1}}^{(j-1)} \dots U_{\theta_1}^{(1)} b^{(j)},$$

468 and the approximate solution to (5.3) is obtained from $s_\ell^{(j)} = LU_{\theta_{j-1}}^{(j-1)} \dots U_{\theta_1}^{(1)} x_\ell^{(j)}$.

469 **5.1. Setup.** In our numerical experiments, we use the Lorenz-96 [16] model as
 470 the physical dynamical system, $\mathcal{M}_{t_0, t_i}(\cdot)$, which is commonly used as a reference model
 471 in data assimilation. The observation operator $\mathcal{H}(\cdot)$ is defined as a uniform selection
 472 operator, meaning $\mathcal{H}(x)$ extracts a subset of x that is uniformly selected. B is chosen
 473 as a discretized diffusion operator with a standard deviation $\sigma_b = 0.8$ [9]. We consider
 474 $R_1 = R_2 = \sigma_r^2 I_m$ with $\sigma_r = 0.2$. We choose $n = 1000$ and $N_t = 2$, and we consider
 475 two different scenarios, with a different number of observations: (1) *LowObs* with
 476 $m_1 = m_2 = 150$ and (2) *HighObs* with $m_1 = m_2 = 300$. For both cases, 2 outer loops
 477 are performed within TGN. CG is applied to the first linear system $A^{(1)}x = b^{(1)}$ with
 478 100 iterations. Then, approximate largest eigen-pairs of $A^{(1)}$, (S_k, Λ_k) , are computed
 479 and selected based on convergence criteria with a tolerance of $\varepsilon = 10^{-3}$ (See [Section
 480 1.3][24] for further details). With this criteria, the number of selected eigen-pairs is
 481 45 in the *LowObs* case and 26 in the *HighObs* case. Using these pairs, the scaled
 482 LMP, $U_{\theta_1}^{(1)}$, is applied as a preconditioner for $j = 2$. Matrix-vector products with the
 483 preconditioner are carried out via an operator using the selected pairs, meaning the
 484 preconditioner matrix is not explicitly constructed.

485 **5.2. Numerical Results.** In this section, we present numerical results only for
 486 the second outer loop ($j = 2$) of the TGN method. We compare the performance of
 487 the methodologies of Table 1 in terms of convergence rate and computational cost.

Method	Description	Initial guess
BPrec	Algorithm 3.1 applied to (5.4)	$x_0 = 0$
sLMP-Base	Algorithm 3.1 applied to (5.5), $\theta_1 = 1$	$x_0 = 0$
Init-sLMP-Base	Algorithm 3.1 applied to (5.5), $\theta_1 = 1$	$x_0 = U_{\theta_1}^{-1} S_k \Lambda_k^{-1} S_k^\top b^{(2)}$
sLMP-λ_k	Algorithm 3.1 applied to (5.5), $\theta_1 = \lambda_k$	$x_0 = 0$
sLMP-θ_r	Algorithm 3.1 applied to (5.5), $\theta_1 = \theta_r$	$x_0 = 0$
sLMP-θ_m	Algorithm 3.1 applied to (5.5), $\theta_1 = (\lambda_k + 1)/2$	$x_0 = 0$
DefCG	Algorithm A.1 applied to (5.4), $W = S_k$	$x_{-1} = 0$

Table 1: Description of methods used in the numerical experiments

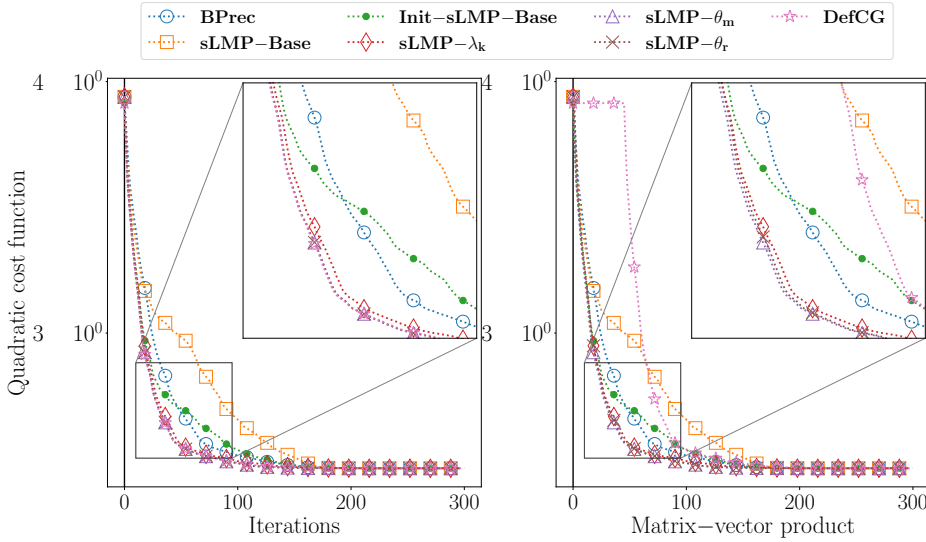


Fig. 1: Quadratic cost function values along all CG iterates (left) and with respect to the number of matrix-vector product with the matrix $A^{(1)}$ and $A^{(2)}$ (right).

488 Note that, for **sLMP- θ_r** we compute θ_r using (4.10) with $r_0 = b^{(2)}$ and $A = A^{(1)}$.
 489 As a result, computation of approximate θ_r requires an extra matrix vector product
 490 with $A^{(1)}$. Figure 1 shows the quadratic cost function values (5.2) and number of
 491 matrix-vector products with $A^{(1)}$ and $A^{(2)}$ along CG iterations.

492 We can easily see that **sLMP-Base** is not necessarily better than **BPrec** espe-
 493 cially in the early iterations. This means that the scaled spectral LMP, clustering the
 494 largest k eigenvalues around 1, might reduce the total number iterations to converge,
 495 however it does not guarantee better convergence for early iterations. The slow con-
 496 vergence of **sLMP-Base** can be partly explained by the fact that perturbations may
 497 cause some eigenvalues to appear near zero, as depicted in Figure 2. When changing
 498 the clustering position from 1 to λ_k by using **sLMP- λ_k** , we can see that the method
 499 performs better than **BPrec**. In this case, however the gap between the cluster and

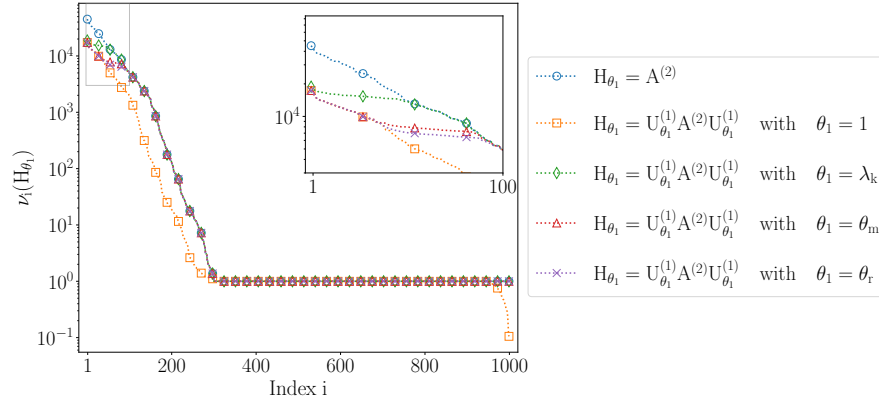


Fig. 2: Spectrum of $U_{\theta_1}^{(1)} A^{(2)} U_{\theta_1}^{(1)}$ for different values of θ_1 on a logarithmic scale. LowObs scenario ($k = 45$).

500 the remaining spectrum as defined in [Theorem 4.9](#), i.e. $\alpha(\theta_1^{(1)})/\theta_1^{(1)}$, can be large.
 501 When clustering around θ_r and θ_m is applied with **sLMP- θ_r** and **sLMP- θ_m** respec-
 502 tively, the value of $\alpha(\theta_1^{(1)})/\theta_1^{(1)}$ reduces for both cases (see [Fig. 2](#)). This improves the
 503 convergence compared to **sLMP- λ_k** as seen from [Figure 1](#).

504 **Init-sLMP-Base** performs better than **sLMP-Base**, i.e. starting from $x_0 =$
 505 $S_k \Lambda_k^{-1} S_k^\top b^{(2)}$ improves performance compared to starting from $x_0 = 0$. This im-
 506 provement arises because the initial residual's components in the eigenbasis of $A^{(2)}$
 507 are reduced. In fact, without any perturbation, these components would be com-
 508 pletely eliminated. Although, the performance is improved with this initial guess,
 509 it can not reach the performance of **DefCG**. This demonstrates that modifying the
 510 initial guess enhances convergence; however, the placement of the eigenvalue cluster-
 511 ing can have an even more significant impact. This is evident from the fact that the
 512 performance of **sLMP- θ_m** and **sLMP- θ_r** are very close to that of **DefCG**.

513 The right panel of [Figure 1](#) shows the values of the quadratic cost function as a
 514 function of the number of matrix-vector products performed with $A^{(j)}$ for $j = 1, 2$
 515 across different methods. Although **DefCG** performs better, it is computationally
 516 expensive as it requires forming the projected matrix $S_k^\top A^{(2)} S_k$. Among the other
 517 techniques, **sLMP- θ_r** requires one additional matrix-vector product with $A^{(1)}$ to com-
 518 pute θ_r . However, as shown in [Figure 1](#), **sLMP- θ_m** and **sLMP- λ_k** do not require
 519 any extra matrix-vector products either $A^{(1)}$ or $A^{(2)}$.

520 These results indicate that the performance of CG, when used with scaled spectral
 521 LMP, can be significantly improved, approaching that of deflated CG, by selecting the
 522 position of the eigenvalue clusters based on CG's convergence properties. The cluster
 523 position is determined by θ , whose computation incurs no additional cost for **sLMP- θ_m**
 524 **sLMP- λ_k** . Conclusions from experiments with *HighObs* are very similar, the
 525 obtained results are depicted in [Figures 3](#) and [4](#) in [Appendix B](#).

526 **6. Conclusion.** We have proposed a *scaled* spectral LMP to accelerate the so-
 527 lution of a sequence of SPD systems $A^{(j)} x^{(j)} = b^{(j)}$ for $j \geq 1$. The *scaled* LMP
 528 incorporates a low-rank update based on k eigenpairs of the matrix A . We have
 529 provided theoretical analysis of the *scaled* spectral LMP when $A^{(j)} = A$. We have
 530 shown that the scaled spectral LMP [\(4.1\)](#) clusters k eigenvalues around the scaling

531 parameter θ , and leaves the rest of the spectrum untouched.

532 We have focused on the choice of θ to ensure that PCG achieves faster convergence,
 533 particularly in the early iterations. In the first approach, we have proposed choosing θ
 534 to guarantee a lower energy norm of the error at each iteration of PCG. In the second
 535 approach, we have obtained an optimum θ in the sense that it minimizes the energy
 536 norm of the error at the first iteration. Our analysis reveals that, with the optimal
 537 θ , the components of the first residual is eliminated from the eigenspace of A , which
 538 aligns with the core principle of deflated CG. Lastly, we have also explored a scaling
 539 parameter that approximates the iterates of deflated CG. We have provided the link
 540 between the deflated CG and PCG with the scaled spectral LMP.

541 We have compared different methods for solving a nonlinear weighted least-
 542 squares problem arising in data assimilation. In our numerical experiments, we used
 543 approximate eigenpairs to construct the scaled spectral LMP. First, we have demon-
 544 strated that selecting θ based on PCG convergence properties significantly accelerates
 545 early convergence compared to the conventional choice of $\theta = 1$. Then, we have shown
 546 that θ values that reduce the spectral gap between θ and the remaining eigenvalues
 547 lead to faster convergence. Additionally, we have compared the scaled spectral LMP
 548 with deflated CG, showing that the scaled spectral LMP produces iterates similar to
 549 deflated CG, but at a negligible computational cost and memory, unlike deflated CG.
 550 These numerical results clearly highlight the importance of selecting the preconditioner
 551 not only as an approximation to the inverse of A , but also with consideration
 552 of its role within PCG. In particular, we have demonstrated the significance of the
 553 placement of clustered eigenvalues, an often overlooked factor in the literature, on the
 554 early convergence of PCG.

555 As the next step, we will provide a detailed theoretical perturbation analysis in a
 556 forthcoming paper. Additionally, we aim to validate the proposed preconditioner in
 557 an operational weather prediction system.

558 **Appendix A. Deflated CG with S_k .** The deflation technique outlined in
 559 [Algorithm A.1](#) is defined for any deflation subspace W , see [\[23\]](#) for more details. The
 560 main idea is to speed-up the CG starting from an initial point such that the initial
 561 residual does not have components in the deflation subspace W and to update the
 562 search directions such that $W^\top Ap_j = 0$. A widely used approach is to choose W as the
 eigenvectors corresponding to the eigenvalues that slows down the CG convergence.

Algorithm A.1 Deflated-CG

- 1: Choose k linearly independent vectors w_1, w_2, \dots, w_k .
 - 2: Define $W = [w_1, w_2, \dots, w_k]$, and choose x_{-1} .
 - 3: Set $x_0^{\text{Def}} = x_{-1} + W(W^\top AW)^{-1}W^\top r_{-1}$, where $r_{-1} = b - Ax_{-1}$. $W^\top r_0 = 0$
 - 4: Set $p_0 = r_0 - W(W^\top AW)^{-1}W^\top Ar_0$. $W^\top Ap_0 = 0$
 - 5: **for** $j = 1, 2, \dots$ **do**
 - 6: $\alpha_{j-1} = r_{j-1}^\top r_{j-1} / (p_{j-1}^\top Ap_{j-1})$
 - 7: $x_j^{\text{Def}} = x_{j-1}^{\text{Def}} + \alpha_{j-1} p_{j-1}$
 - 8: $r_j = r_{j-1} - \alpha_{j-1} Ap_{j-1}$ $W^\top r_j = 0$
 - 9: $\beta_{j-1} = r_j^\top r_j / (r_{j-1}^\top r_{j-1})$
 - 10: $p_j = \beta_{j-1} p_{j-1} + r_j - W(W^\top AW)^{-1}W^\top Ar_j$ $W^\top Ap_j = 0$
 - 11: **end for**
-

563

564

If we choose $W = S_k$, and using the fact that $S_k^\top AS_k = \Lambda_k$ and $AS_k = S_k \Lambda_k$, we

565 can achieve the following simplifications:

- 566 • $x_0^{\text{Def}} = x_{-1} + S_k \Lambda_k^{-1} S_k^\top r_{-1}$,
- 567 • $p_0 = r_0 - S_k S_k^\top r_0$.
- 568 • $p_j = \beta_{j-1} p_{j-1} + r_j - S_k S_k^\top r_j$.

569 LEMMA A.1. *The residual r_j and the direction p_j are orthogonal to $\text{span}(S_k)$.*

570 *Proof.* We proceed by induction. For $j = 0$, $r_0 = r_{-1} - S_k S_k^\top r_{-1}$, from which
 571 it follows that $S_k^\top r_0 = 0$. As a consequence, $S_k^\top p_0 = 0$. Assume that r_j and p_j are
 572 orthogonal to $\text{span}(S_k)$ for j . We have $r_{j+1} = r_j - \alpha_j A p_j$. From [23, Proposition 3.3],
 573 replacing W by S_k , we have $S_k^\top A p_j = 0$. Since $p_j, r_j \perp \text{span}(S_k)$ by assumption, it
 574 follows that $r_{j+1} \perp \text{span}(S_k)$. For $p_{j+1} = \beta_j p_j + r_{j+1} - S_k S_k^\top r_{j+1} = \beta_j p_j + r_{j+1}$, we
 575 get $p_{j+1} \perp \text{Span}(S_k)$ since $S_k^\top r_{j+1} = 0$ as shown and $p_j \perp \text{Span}(S_k)$ by assumption. \square

576 From Lemma A.1, it follows that $p_j = \beta_{j-1} p_{j-1} + r_j - S_k S_k^\top r_j = \beta_{j-1} p_{j-1} + r_j$.
 577 With these simplifications, it is clear that in exact arithmetic, deflated CG, when
 578 used with the deflated subspace consisting of a set of eigenvectors of A , generates
 579 iterates equivalent to those generated by using the initial guess x_0^{Def} in standard CG.

580 Appendix B. Results for the *HighObs* scenario.

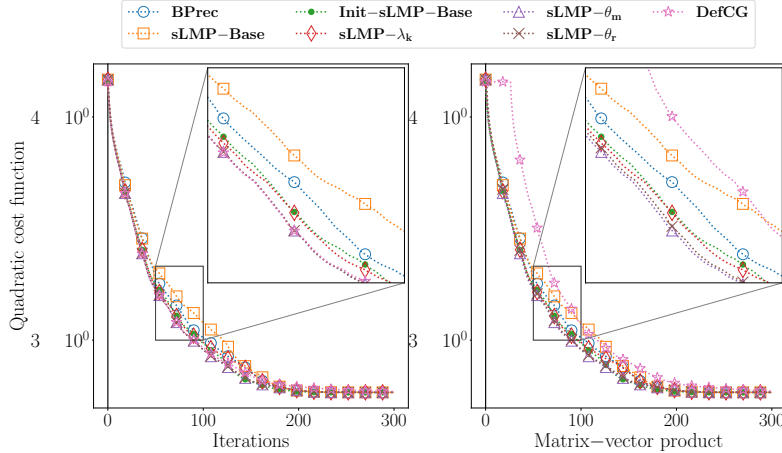


Fig. 3: Quadratic cost function values along all CG iterates and with respect to the number of matrix-vector product for the *HighObs* scenario.

581

REFERENCES

- 582 [1] M. BENZI, *Preconditioning techniques for large linear systems: A survey*, Journal of Computa-
 583 tional Physics, 182 (2002), pp. 418–477.
- 584 [2] E. CARSON, J. LIESEN, AND Z. STRAKOŠ, *Towards understanding cg and gmres through exam-
 585 ples*, Linear Algebra and its Applications, (2024).
- 586 [3] E. CARSON AND Z. STRAKOŠ, *On the cost of iterative computations*, Philosophical Transactions
 587 of the Royal Society A, 378 (2020), p. 20190050.
- 588 [4] R. DALEY, *Atmospheric data analysis*, Cambridge University Press, 1991.
- 589 [5] V. DER SLUIS AND V. DER VORST., *The rate of convergence of conjugate gradients*, Numerische
 590 Mathematik, 48 (1986), pp. 543–560.
- 591 [6] M. FISHER, J. NOCEDAL, Y. TRÉMOLET, AND S. J. WRIGHT, *Data assimilation in weather
 592 forecasting: a case study in pde-constrained optimization*, Optimization and Engineering,
 593 10 (2009), pp. 409–426.

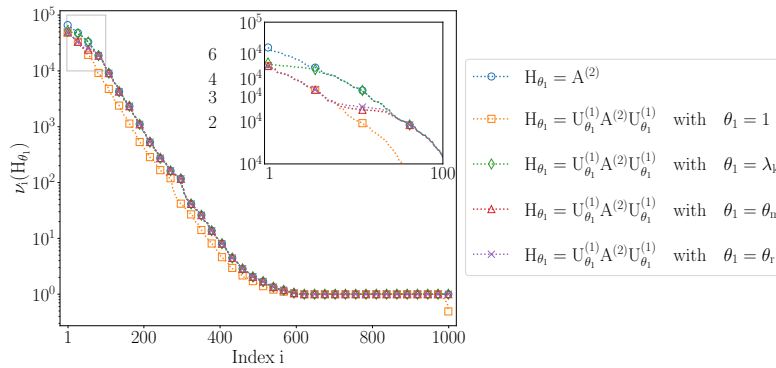


Fig. 4: Spectrum of $U_{\theta_1}^{(1)} A^{(2)} U_{\theta_1}^{(1)}$ with different θ_1 for the *HighObs* scenario ($k = 26$).

- 594 [7] L. GIRAUD AND S. GRATTON, *On the sensitivity of some spectral preconditioners*, SIAM Journal
595 on Matrix Analysis and Applications, 27 (2006), pp. 1089–1105.
- 596 [8] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, The Johns Hopkins University Press,
597 Baltimore, 4th ed., 2013.
- 598 [9] O. GOUX, S. GÜROL, A. T. WEAVER, Y. DIOUANE, AND O. GUILLET, *Impact of correlated
599 observation errors on the conditioning of variational data assimilation problems*, Numerical
600 Linear Algebra with Applications, 31 (2024), p. e2529.
- 601 [10] S. GRATTON, A. S. LAWLESS, AND N. K. NICHOLS, *Approximate gauss–newton methods for
602 nonlinear least squares problems*, SIAM Journal on Optimization, 18 (2007), pp. 106–132.
- 603 [11] S. GRATTON, A. SARTENAER, AND J. TSHIMANGA, *On a class of limited memory preconditioners
604 for large scale linear systems with multiple right-hand sides*, SIAM Journal on Optimiza-
605 tion, 21 (2011), pp. 912–935.
- 606 [12] A. GREENBAUM, *Comparison of splittings used with the conjugate gradient algorithm*, Nu-
607 merische Mathematik, 33 (1979), pp. 181–193.
- 608 [13] M. R. HESTENES AND E. STIEFEL, *Methods of conjugate gradients for solving linear systems*,
609 Journal of research of the National Bureau of Standards, 49 (1952), pp. 409–435.
- 610 [14] K. KAHL AND H. RITTICH, *The deflated conjugate gradient method: Convergence, perturbation
611 and accuracy*, Linear Algebra and its Applications, 515 (2017), pp. 111–129.
- 612 [15] E. KALNAY, *Atmospheric Modeling, Data Assimilation and Predictability*, Cambridge Univer-
613 sity Press, 2002.
- 614 [16] E. N. LORENZ, *Predictability: A problem partly solved*, in Proc. Seminar on predictability, vol. 1,
615 Reading, 1996.
- 616 [17] J. L. MORALES AND J. NOCEDAL, *Automatic preconditioning by limited memory quasi-newton
617 updating*, SIAM Journal on Optimization, 10 (2000), pp. 1079–1096.
- 618 [18] R. NABBEN AND C. VUIK, *A comparison of deflation and the balancing preconditioner*, SIAM
619 Journal on Scientific Computing, 27 (2006), pp. 1742–1759.
- 620 [19] S. G. NASH AND J. NOCEDAL, *A numerical study of the limited memory bfgs method and the
621 truncated-newton method for large scale optimization*, SIAM Journal on Optimization, 1
622 (1991), pp. 358–372.
- 623 [20] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, Springer, New York, NY, USA,
624 2nd ed., 2006.
- 625 [21] J. W. PEARSON AND J. PESTANA, *Preconditioners for Krylov subspace methods: An overview*,
626 GAMM-Mitteilungen, 43 (2020), p. e202000015.
- 627 [22] Y. SAAD, *Iterative methods for sparse linear systems*, PWS Publishing Company, Boston, USA,
628 1996.
- 629 [23] Y. SAAD, M. YEUNG, J. ERHEL, AND F. GUYOMARC’H, *A deflated version of the conjugate
630 gradient algorithm*, SIAM Journal on Scientific Computing, 21 (2000), pp. 1909–1926.
- 631 [24] J. TSHIMANGA, *On a class of limited memory preconditioners for large-scale nonlinear least-
632 squares problems (with application to variational ocean data assimilation)*, PhD thesis,
633 Department of Mathematics, University of Namur, Namur, Belgium, 2007.
- 634 [25] A. J. WATHEN, *Preconditioning*, Acta Numerica, 24 (2015), pp. 329–376.