# STEEPEST DESCENT AND CONJUGATE GRADIENT METHODS WITH VARIABLE PRECONDITIONING[∗]

ANDREW V. KNYAZEV[†] AND ILYA LASHUK[‡]

**Abstract.** We show that the conjugate gradient method with variable preconditioning in certain situations cannot give any improvement compared to the steepest descent method for solving a linear system with a symmetric positive definite (SPD) matrix of coefficients. We assume that the preconditioner is SPD on each step, and that the condition number of the preconditioned system matrix is bounded from above by a constant independent of the step number. Our proof is geometric and is based on the simple fact that a nonzero vector multiplied by all SPD matrices with a condition number bounded by a constant generates a circular cone.

**Key words.** Steepest descent, conjugate gradient method, variable preconditioning, preconditioner, condition number

**1. Introduction.** The steepest descent (SD) and conjugate gradient (CG) methods are well-known iterative procedures to obtain approximate solutions $x_k$ of a linear system $Ax = b$ with symmetric and positive definite (SPD) matrix A. To accelerate the convergence of the error $e_k = x - x_k$ to zero in the SD and CG methods preconditioning is often used, i.e. on every iteration $k$ an operator $B_k$, called the preconditioner, is introduced, possibly different for each iteration. A general algorithm of the preconditioned SD or CG (PSD or PCG respectively) that we consider in this paper can be presented in the following way, e.g., [2]: given $A$, $b$, $\{B_k\}$, $\{m_k\}$, $x_0$, for $k = 0, 1, \ldots$:

$$(1.1) \quad \begin{aligned} r_k &= b - Ax_k, \\ s_k &= B_k^{-1} r_k, \\ p_k &= s_k - \sum_{l=k-m_k}^{k-1} \frac{(As_k, p_l)}{(Ap_l, p_l)} p_l, \\ x_{k+1} &= x_k + \frac{(r_k, p_k)}{(Ap_k, p_k)} p_k, \end{aligned}$$

where

$$(1.2) \quad 0 \le m_k \le k \text{ and } m_{k+1} \le m_k + 1.$$

The latter condition appears in [6] and ensures that the formula for $p_k$ in (1.1) performs the standard Gram–Schmidt $A$-orthogonalizations to previous search directions. The choice $m_k = k$, $\min\{k, 1\}$, or 0 corresponds to full orthogonalization, the standard PCG, or PSD methods, respectively. The full orthogonalization performs explicit $A$-orthogonalizations to all previous search directions. It is well known that the $A$-orthogonalization terms with $l < k - 1$ in the sum vanish in exact arithmetic if the preconditioner is SPD and fixed.

It is also well known that if the preconditioner is SPD and fixed, $B_k = B = B^* > 0$, the preconditioned method (1.1) using this preconditioner can be viewed as an unpreconditioned method applied to the preconditioned system $B^{-1}Ax = B^{-1}b$ in the $B$-based inner product $(x, y)_B = (Bx, y)$. This implies that the theory obtained for unpreconditioned methods remains valid for preconditioned methods. The situation changes, however, if different preconditioners $B_k$ are used on each iteration of the PCG method.

The present paper concerns the behavior of method (1.1), where the preconditioner $B_k$ varies from step to step, but remains SPD on each step and the spectral condition number

$$\kappa\left(B_k^{-1}A\right) = \frac{\lambda_{\max}\left(B_k^{-1}A\right)}{\lambda_{\min}\left(B_k^{-1}A\right)}$$

is bounded from above by some constant $\kappa_{\max}$ independent of the step number $k$. We note that the matrix $B_k^{-1}A$ is SPD with respect to, e.g., $B_k$ inner product, so its eigenvalues are real positive.

The main result of this paper is that under our assumptions the preconditioned method (1.1) with (1.2) turns into the PSD method with the worst possible convergence rate on every iteration by properly choosing the preconditioners $B_k$. Thus one can only guarantee the convergence rate for the method (1.1) with (1.2) just the same as for the PSD method, $m_k = 0$, obtained by Kantorovich [4]:

(1.3)
$$\frac{\|e_{k+1}\|_A}{\|e_k\|_A} \le \frac{\kappa_{\max} - 1}{\kappa_{\max} + 1}.$$

Our proof is geometric and is based on the simple fact that a nonzero vector multiplied by all SPD matrices with a condition number bounded by a constant generates a circular cone.

Different aspects of the PCG methods with variable preconditioning are considered, e.g., in [1, 2], where rather general nonlinear preconditioning is introduced, and in [3, 6] that mainly deal with the case when preconditioner $B_k$ on each iteration $k$ approximates a fixed SPD operator $B$. In [1, 2, 6], convergence estimates for the PCG method are proved that are weaker than (1.3). For recent results and other aspects of variable preconditioning, see [7–10] and references there. No attempts are apparently made in the literature to obtain a result similar to that of the present paper.

The rest of the paper is organized as follows. In Section 2 we prove that a nonzero vector multiplied by all SPD matrices with a condition number bounded by a constant generates a circular cone. Basic properties of the method (1.1), most importantly, the local optimality, are derived in Section 3. In Section 4 we apply our results from Section 2 about the cone to obtain a new proof of estimate (1.3). In Section 5 we analyze the convergence of the PCG method with variable preconditioning and prove our main result, described above.

**2. Circular cones represent sets of SPD matrices with varying condition numbers.** For any pair of real non-zero vectors $x$ and $y$ we define the angle between $x$ and $y$ in the usual way as

$$\angle(x, y) = \arccos\left(\frac{(x, y)}{\|x\|\,\|y\|}\right) \in [0, \pi].$$

The following statement is inspired by a ball theorem from [5].

THEOREM 2.1. *The set $\{Cx\}$, where $x$ is a fixed nonzero real vector and $C$ runs through all SPD matrices with condition number $\kappa(C)$ bounded from above by some $\kappa_{\max}$, is a pointed circular cone, specifically,*

$$\{Cx, \quad C = C^* > 0, \kappa(C) \leq \kappa_{\max}\} = \left\{y: \quad \sin \angle(x, y) \leq \frac{\kappa_{\max} - 1}{\kappa_{\max} + 1}\right\}.$$

Theorem 2.1 can be proved by constructing our cone as the smallest pointed cone that includes the ball considered in [5] and using the corresponding result of [5], but we provide an direct proof here based on the following two lemmas. The first lemma is simple and states that the set in question cannot be larger than the cone:

LEMMA 2.2. *Let $x$ be a non-zero real vector, let $C$ be SPD with spectral condition number $\kappa(C)$. Then*

$$\sin \angle(x, Cx) \leq \frac{\kappa(C) - 1}{\kappa(C) + 1}$$

*Proof.* Denote $y = Cx$. We have $(x, Cx) = (y, C^{-1}y) > 0$, since $C$ is SPD, so $y \neq 0$ and $\angle(x, y) < \frac{\pi}{2}$. A positive scaling of $C$ and thus of $y$ is obviously irrelevant, so let us choose $y$ to be the orthogonal projection of $x$ onto 1-dimensional subspace spanned by the original $y$. Then from elementary 2D geometry it follows that $\|y - x\| = \|x\| \sin \angle(x, y)$. The orthogonal projection of a vector onto a subspace is the best approximation to the vector from the subspace, thus

$$\|x\| \sin \angle(x, y) = \|y - x\| \leq \|sy - x\| = \|sCx - x\| \leq \|sC - I\| \|x\|$$

for any scalar $s$, where $I$ is the identity. Taking $s = 2/(\lambda_{\max}(C) + \lambda_{\min}(C))$, where $\lambda_{\min}(C)$ and $\lambda_{\max}(C)$ are the minimal and maximal eigenvalues of $C$, respectively, we get $\|sC - I\| = (\kappa_{\max} - 1)/(\kappa_{\max} + 1)$. $\square$

The second lemma implies that every point in the cone can be represented as $Cx$ for some SPD matrix $C$ with $\kappa(C) \leq \kappa_{\max}$.

LEMMA 2.3. *Let $x$ and $y$ be non-zero real vectors, such that $\angle(x, y) \in \left[0, \frac{\pi}{2}\right)$. Then there exists SPD matrix $C$, such that $Cx = y$ and*

$$\frac{\kappa(C) - 1}{\kappa(C) + 1} = \sin \angle(x, y).$$

*Proof.* Denote $\alpha = \angle(x, y)$. A positive scaling of vector $y$ is irrelevant, so as in the previous proof we choose $y$ to be the orthogonal projection of $x$ onto 1-dimensional subspace spanned by the original $y$, then $\|y - x\| = (\sin \alpha) \|x\|$, so the vectors $y - x$ and $(\sin \alpha)x$ are of the same length. This implies that there exists a Housholder reflection $H$ such that $H((\sin \alpha)x) = y - x$, cf. [5], so $(I + (\sin \alpha)H)x = y$. We define $C = I + (\sin \alpha)H$ to get $Cx = y$. Any Housholder reflection is symmetric and has only two distinct eigenvalues $\pm 1$, so $C$ is also symmetric and has only two distinct positive eigenvalues $1 \pm \sin \alpha$, as $\alpha \in [0, \pi/2)$, and we conclude that $C > 0$ and $\kappa(C) = (1 + \sin \alpha)/(1 - \sin \alpha)$. $\square$

**3. Local optimality of the method with variable preconditioning.** Here we discuss some basic properties of the general method (1.1) with (1.2). The starting point is the following known, e.g., [1, 2, 6], $A$-orthogonality relations:

(3.1) $$(e_{k+1}, p_j)_A = 0, \quad j = k - m_k, \ldots, k.$$

We now use (3.1) to prove the local optimality of the method (1.1) with (1.2).

THEOREM 3.1. *The A-norm of the error, $\|e_{k+1}\|_A$, in the general method (1.1) with (1.2) is bounded from above by the A-norm of the error of one step of the PSD method, where $m_k = 0$, using $x_k$ as the initial guess and $B_k$ as the preconditioner, i.e.,*

$$\|e_{k+1}\|_A \leq \min_\alpha \|e_k - \alpha s_k\|_A.$$

*Proof.* By (1.1),

$$s_k = p_k + \sum_{l=k-m_k}^{k-1} \frac{(As_k, p_l)}{(Ap_l, p_l)} p_l,$$

so in addition to (3.1), we have $(e_{k+1}, s_k)_A = 0$. Errors on two subsequent iterations in (1.1) are related as

$$e_{k+1} = e_k - \frac{(e_k, p_k)_A}{(p_k, p_k)_A} p_k,$$

which follows directly from (1.1) and implies that

$$e_{k+1} \in e_k + \operatorname{span}\{s_k, p_{k-m_k}, \dots, p_{k-1}\}.$$

Putting this together with $A$-orthogonality relations (3.1), we deduce that

$$\|e_{k+1}\|_A = \min_{p \in \operatorname{span}\{s_k, p_{k-m_k}, \dots, p_{k-1}\}} \|e_k - p\|_A.$$

As a corollary, we obtain the statement of the theorem. □

**4. Convergence rate estimate for variable preconditioning.** The classical [4] convergence rate bound (1.3) for the PSD method is "local" in a sense that it relates the $A$-norm of the error on two subsequent iterations and doesn't depend on previous iterations, thus, it remains valid when the preconditioner $B_k$ changes from iteration to iteration, but the condition number $\kappa\left(B_k^{-1}A\right)$ is bounded from above by some constant $\kappa_{\max}$ independent of $k$. The goal of this section is to give a new simple proof of the estimate (1.3) for the PSD method, based on our cone Theorem 2.1, and to extend this proof to cover the general method (1.1) with (1.2).

We denote the angle between two real nonzero vectors with respect to the $A$-based inner product as

$$\angle_A(x, y) = \arccos\left(\frac{(x, y)_A}{\|x\|_A \|y\|_A}\right) \in [0, \pi]$$

and express the error reduction ratio for the PSD method in terms of the angle with respect to the $A$-based inner product:

LEMMA 4.1. *On every step of the PSD algorithm, (1.1) with $m_k = 0$, the error reduction factor takes the form*

$$\frac{\|e_{k+1}\|_A}{\|e_k\|_A} = \sin(\angle_A(e_k, B_k^{-1}Ae_k)).$$

*Proof.* By (3.1), $(e_{k+1}, p_k)_A = 0$ . Now, for $m_k = 0$, in addition, $p_k = s_k$, so $0 = (e_{k+1}, p_k)_A = (e_{k+1}, s_k)_A = (e_{k+1}, x_{k+1} - x_k)_A$, i.e. the triangle with vertices $x$, $x_k$, $x_{k+1}$ is right-angled in the $A$-inner product, where the hypotenuse is $e_k = x - x_k$. Therefore, $\|e_{k+1}\|_A / \|e_k\|_A = \sin(\angle_A(e_k, x_{k+1} - x_k)) = \sin(\angle_A(e_k, s_k))$, where $s_k = B_k^{-1}(b - Ax_k) = B_k^{-1}Ae_k$ by (1.1). $\square$

Combining the results of Lemmas 2.2 and 4.1 together immediately leads to (1.3) for the PSD method, where $m_k = 0$. Now, taking into account Theorem 3.1, we get

THEOREM 4.2. *The convergence rate bound (1.3) holds for the general method (1.1) with (1.2).*

**5. The convergence rate estimate is sharp.** Here we formulate and prove the main result of the paper that one can only guarantee the convergence rate for method (1.1) with (1.2) just the same as for the PSD method, (1.1) with $m_k = 0$, described by (1.3).

THEOREM 5.1. *For any given SPD matrix $A$, vectors $b$ and $x_0$, and $\kappa_{\max} > 1$, assuming the matrix size larger than the number of iterations, one can choose such a sequence of SPD preconditioners $B_k$, satisfying $\kappa(B_k^{-1}A) \leq \kappa_{\max}$, that the method (1.1) with (1.2) turns into the PSD method, (1.1) with $m_k = 0$, and for each iteration*

$$(5.1) \qquad \frac{\|e_{k+1}\|_A}{\|e_k\|_A} = \frac{\kappa_{\max} - 1}{\kappa_{\max} + 1}.$$

*Proof.* We construct the sequence $B_k$ by induction. First, we choose any vector $q_0$, such that $\sin \angle_A(q_0, e_0) = (\kappa_{\max} - 1)/(\kappa_{\max} + 1)$. According to Lemma 2.3 applied in the $A$-inner product, there exists an $A$-SPD matrix $C_0$ with condition number $\kappa(C_0) = \kappa_{\max}$, such that $C_0 e_0 = q_0$. We define the SPD $B_0 = C_0^{-1}A$, then $\kappa(B_0^{-1}A) = \kappa(C_0) = \kappa_{\max}$. We have $s_k = B_k^{-1}Ae_k$, so such a choice of $B_0$ implies $s_0 = q_0$. Also, we have $p_0 = s_0$, i.e. the first step is always a PSD step, thus, by Lemma 4.1 we have proved (5.1) for $k = 0$. Note that $(e_1, p_0)_A = 0$ by (3.1).

Second, we make the induction assumption: suppose preconditioners $B_l$ for $l \leq k - 1$ are constructed, such that

$$\frac{\|e_{l+1}\|_A}{\|e_l\|_A} = \frac{\kappa_{\max} - 1}{\kappa_{\max} + 1}$$

and $(e_k, p_l)_A = 0$ hold for all $l \leq k - 1$. The dimension of the space is greater than the total number of iterations by our assumption, so there exists a vector $u_k$, such that $(u_k, p_l)_A = 0$ for $l \leq k - 1$ and $u_k$ and $e_k$ are linearly independent. Then the 2D subspace spanned by $u_k$ and $e_k$ is $A$-orthogonal to $p_l$ for $l \leq k - 1$. In this subspace we can choose a vector $q_k$, such that $\sin \angle_A(q_k, e_k) = (\kappa_{\max} - 1)/(\kappa_{\max} + 1)$. This vector will be obviously $A$-orthogonal to $p_l$, $l \leq k - 1$. Then, applying the same reasoning as for constructing $B_0$, we deduce that there exists an SPD $B_k$ such that $\kappa(B_k^{-1}A) \leq \kappa_{\max}$ and $B_k^{-1}Ae_k = q_k$. With such a choice of $B_k$ we have $s_k = q_k$. Since $q_k = s_k$ is $A$-orthogonal to $p_l$ for all $l \leq k - 1$, it turns out that $p_k = s_k$, no matter how $\{m_k\}$ are chosen. This means that $x_{k+1}$ is obtained from $x_k$ by a steepest descent step. Then we apply Lemma 4.1 and conclude that (5.1) holds. Finally we note, that $(e_{k+1}, p_l)_A = 0$ for all $l \leq k$. Indeed, $(e_{k+1}, p_l)_A = 0$ for all $l \leq k - 1$ since $e_{k+1}$ is linear combination of $e_k$ and $p_k = s_k = q_k$, both $A$-orthogonal to $p_l$ for $l \leq k - 1$. Finally, $(e_{k+1}, p_k)_A = 0$ by (3.1). This completes the construction of $\{B_k\}$ by induction and thus the proof. $\square$

**References.**

[1] O. Axelsson and P. S. Vassilevski. A black box generalized conjugate gradient solver with inner iterations and variable-step preconditioning. *SIAM J. Matrix Anal. Appl.*, 12(4):625–644, 1991. ISSN 0895-4798.

[2] Owe Axelsson. *Iterative solution methods.* Cambridge University Press, Cambridge, 1994. ISBN 0-521-44524-8.

[3] Gene H. Golub and Qiang Ye. Inexact preconditioned conjugate gradient method with inner-outer iteration. *SIAM J. Sci. Comput.*, 21(4):1305–1320 (electronic), 1999/00. ISSN 1064-8275.

[4] L. Kantorovich and G. P. Akilov. *Functional Analysis in Normed Spaces.* Pergamon, NY, 1964.

[5] Klaus Neymeyr. A geometric theory for preconditioned inverse iteration. I. Extrema of the Rayleigh quotient. *Linear Algebra Appl.*, 322(1-3):61–85, 2001. ISSN 0024-3795.

[6] Yvan Notay. Flexible conjugate gradients. *SIAM J. Sci. Comput.*, 22(4):1444–1460 (electronic), 2000. ISSN 1064-8275.

[7] Valeria Simoncini and Daniel B. Szyld. Flexible inner-outer Krylov subspace methods. *SIAM J. Numer. Anal.*, 40(6):2219–2239 (electronic) (2003), 2002. ISSN 0036-1429.

[8] Valeria Simoncini and Daniel B. Szyld. Theory of inexact Krylov subspace methods and applications to scientific computing. *SIAM J. Sci. Comput.*, 25(2): 454–477 (electronic), 2003. ISSN 1064-8275.

[9] Valeria Simoncini and Daniel B. Szyld. On the occurrence of superlinear convergence of exact and inexact Krylov subspace methods. *SIAM Rev.*, 47(2):247–272 (electronic), 2005. ISSN 0036-1445.

[10] Daniel B. Szyld and Judith A. Vogel. FQMR: a flexible quasi-minimal residual method with inexact preconditioning. *SIAM J. Sci. Comput.*, 23(2):363–380 (electronic), 2001. ISSN 1064-8275. Copper Mountain Conference (2000).