

§6.5. Descent methods and conjugate gradients (CG)

→ Descent methods, and the improved descent method (or the CG method) is a powerful technique for solving

$$\boxed{AU = f},$$

when A is SPD (symmetric positive definite), or a SND (symmetric negative definite).

→ We begin considering the basic method of descent method first and extend the concept to CG.

→ The descent method is based on solving a minimization problem of the fun $\phi: \mathbb{R}^m \rightarrow \mathbb{R}$ defined by

$$\boxed{\phi(U) = \frac{1}{2}U^T A U - U^T f} \quad \dots \quad (1)$$

Ex. $m=2$, $U = (u_1, u_2)$, $A = \begin{bmatrix} a_{11} & a_{12} \\ a_{12} & a_{22} \end{bmatrix} \rightarrow$ symmetric

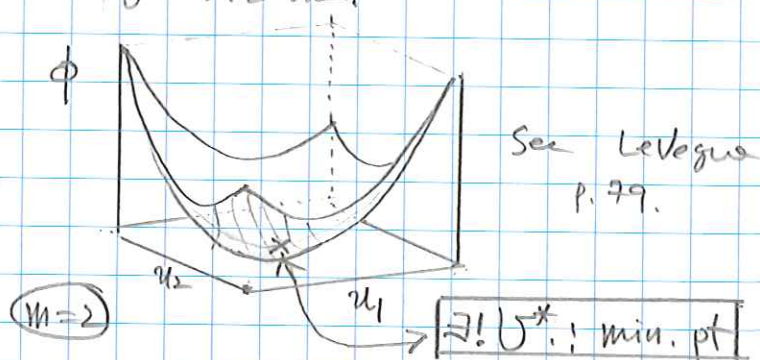
$$\begin{aligned} \phi(U) &= \phi(u_1, u_2) \\ &= \frac{1}{2} \begin{bmatrix} u_1 & u_2 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} \\ a_{12} & a_{22} \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} - \begin{bmatrix} u_1 & u_2 \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \end{aligned}$$

$$= \frac{1}{2} (a_{11} u_1^2 + 2a_{12} u_1 u_2 + a_{22} u_2^2) - (u_1 f_1 + u_2 f_2)$$

→ a quadratic fun of u_1 & u_2 .

⇒ Since A is SPD:

∃! U^* : a min. pt. that minimized ϕ .



⇒ To find $U^* = (u_1^*, u_2^*)$, we look for

$$\begin{cases} 0 = \frac{\partial \phi}{\partial u_1} = a_{11}u_1^* + a_{12}u_2^* - f_1 \\ 0 = \frac{\partial \phi}{\partial u_2} = a_{21}u_1^* + a_{22}u_2^* - f_2 \end{cases} \dots \textcircled{2}$$

⇒ $\textcircled{2}$ is equivalent to $AU^* = f$:

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} u_1^* \\ u_2^* \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \end{bmatrix}$$

⇒ Therefore, we conclude that if we find U^* for the minimization problem of ϕ , then we solve the linear system $AU = f$.

Prk. In general, for any m , the above is also true &

for
$$\begin{cases} U \in \mathbb{R}^m, \\ A \in \mathbb{R}^{m \times m} : \text{SPD}, \end{cases}$$

We seek for U^* s.t.

$$\boxed{0 = \nabla \phi(U^*) = AU^* - f}$$

Prk. If A is indefinite (neither positive nor negative definite), then the eigenvalues of A are not all of the same sign, then $\exists U^*$ s.t. $\nabla \phi(U^*) = 0$, but U^* is not a min or max pt, rather a saddle pt.

→ Any iteration method to find U^* will fail in this case.

→ Now our goal is to minimize $\phi(U)$, or find U^* s.t

$$\boxed{\nabla \phi(U^*) = 0}, \quad \phi(U) = \frac{1}{2} U^T A U - U^T f.$$

In order to solve $AU = f$.

→ The question on "how" to find U^* is then the next important topic;

→ Depending on the choice of a search direction from any initial pt. $U = (u_1, u_2)$ on $\phi(U)$, to $U^* = (u_1^*, u_2^*)$, one can achieve slow or fast convergence.

→ Two ways to choose the k th search direction S_k :

- (i) $S_k = -\nabla \phi(U_k)$; the negative gradient
⇒ the method of steepest descent
- (ii) $S_k = p_k (\neq -\nabla \phi(U_k))$, where p_k is A-conjugate to all previous search directions, $p_j, j=0, 1, \dots, k-1$, i.e.

Def. A-conjugate: $\boxed{p_k^T (A p_j) = 0}, \quad j=0, 1, \dots, k-1.$

⇒ Conjugate gradient (CG) method.

Note, let's denote " k " to be our subindex for the k th iteration from now.

→ Once we set our search direction, S_k , then we find a new pt. U_{k+1} from U_k iteratively!

$$\boxed{U_{k+1} = U_k + \alpha_k S_k} \quad \dots \textcircled{3} \quad k=0, 1, \dots$$

where α_k is a scalar, satisfying!

$$\boxed{\min_{\alpha \in \mathbb{R}} \phi(U_{k+1}) = \min_{\alpha \in \mathbb{R}} \phi(U_k + \alpha_k S_k)} \quad \dots \textcircled{4}$$

→ Here, $\begin{cases} \text{(i)} & \alpha_k > 0 \rightarrow \text{not yet converged} \\ \text{(ii)} & \alpha_k = 0 \rightarrow \text{converged.} \end{cases}$

→ To solve $\textcircled{4}$, we compute the derivative of ϕ w.r.t. α and set it to zero.

→ To take such derivative, let's first write ϕ as a continuous fun of α (i.e., removing the discrete "k" for now);

$$\begin{aligned} \phi(U + \alpha S) &= \frac{1}{2} (U + \alpha S)^T A (U + \alpha S) - (U + \alpha S)^T f \\ &= \left(\frac{1}{2} U^T A U - U^T f \right) + \alpha (S^T A U - S^T f) + \frac{1}{2} \alpha^2 S^T A S \end{aligned}$$

$$\rightarrow \frac{d}{d\alpha} \phi(U + \alpha S) = S^T A U - S^T f + \alpha S^T A S = 0$$

$$\rightarrow \boxed{\alpha = \frac{S^T (f - AU)}{S^T A S}} \quad \dots \textcircled{5}$$

→ Note here that, in iteration, we get

$$\nabla \phi(U_k) = AU_k - f \equiv -r_k, \text{ and we set} \quad \left. \right\} \text{9/27/16}$$

$$\text{Def } \boxed{r_k = f - AU_k} \text{ : the residual vector } \dots \textcircled{6}$$

(A) Steepest Descent

→ For the steepest descent method, we set

$$\boxed{s_k = -\nabla\phi(U_k)} \quad ; \quad \text{the direction that always points in the direction of most rapid decrease of } \phi$$
$$= r_k \quad ; \quad \textcircled{7}$$

$$\rightarrow \left\{ \begin{array}{l} \textcircled{5} + \textcircled{6} + \textcircled{7} : \quad \alpha_k = \frac{r_k^T r_k}{r_k^T A r_k} \quad \dots \textcircled{8} \\ \textcircled{3} : \quad U_{k+1} = U_k + \alpha_k r_k \end{array} \right.$$

→ Note in $\textcircled{8}$, there are two matrix-vector multiplications

(i) $r_k = f - \underline{A}U_k$, &

(ii) $\underline{A}r_k$

→ We can reduce them into one by replacing $\underline{A}U_k$ with $\underline{A}r_{k-1}$

$$\begin{aligned} r_k &= f - \underline{A}U_k \\ &= f - \underline{A}(U_{k-1} + \alpha_{k-1} r_{k-1}) \\ &= \overset{\downarrow}{r_k} - \alpha_{k-1} \underline{A}r_{k-1} \quad \dots \textcircled{9} \end{aligned}$$

→ From $\textcircled{9}$, we removed the need to compute $\underline{A}U_k$, but just kept $\underline{A}r_k$.

Algorithm

Choose a guess U_0

$$r_0 = f - AU_0 \quad (\text{initial residual})$$

for $k=1, 2, \dots$

$$W_{k-1} = AR_{k-1} \quad (\text{the only matrix-vector multiplication})$$

$$\alpha_{k-1} = \frac{r_{k-1}^T r_{k-1}}{r_{k-1}^T W_{k-1}} \quad (\text{compute search parameter})$$

$$U_k = U_{k-1} + \alpha_{k-1} r_{k-1} \quad (\text{update soln})$$

$$r_k = r_{k-1} - \alpha_{k-1} W_{k-1} \quad (\text{compute new residual}) \\ = \text{new search direction}$$

If $\|r_k\| < \epsilon$, $\epsilon = \text{tolerance}$, then
stop

else
continue

end if

end for

Performance

(i) slow

(ii) See LeVeque (Figures 4.3, 4.4)

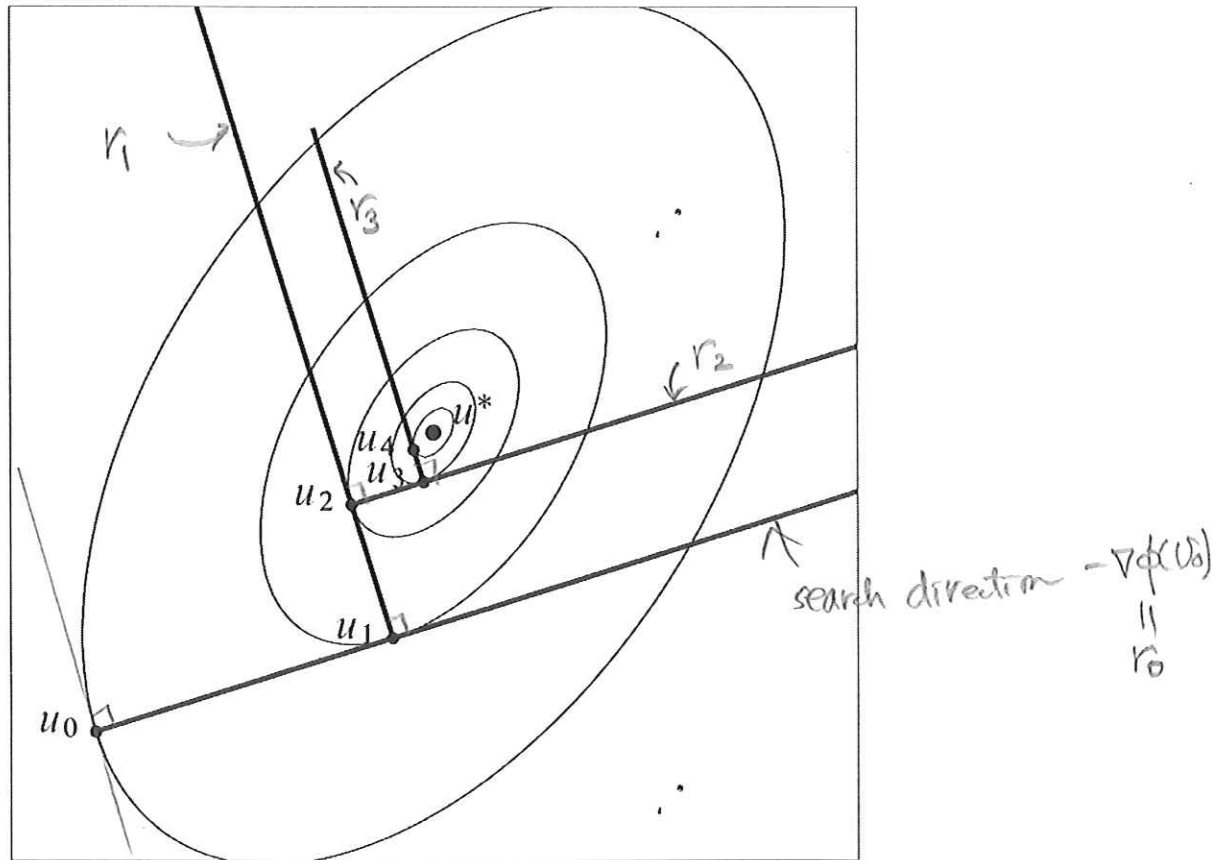


Figure 4.3. Several iterates of the method of steepest descent in the case $m = 2$. The concentric ellipses are level sets of $\phi(u)$.

Principle r_0, r_1, r_2, \dots the search directions are always orthogonal to the contour lines.

→ We move along the direction of the gradient to the pt. where $\phi(u)$ is minimized along the line.

→ $r_0 \perp r_1, r_1 \perp r_2, \dots, r_{k-1} \perp r_k$, etc.

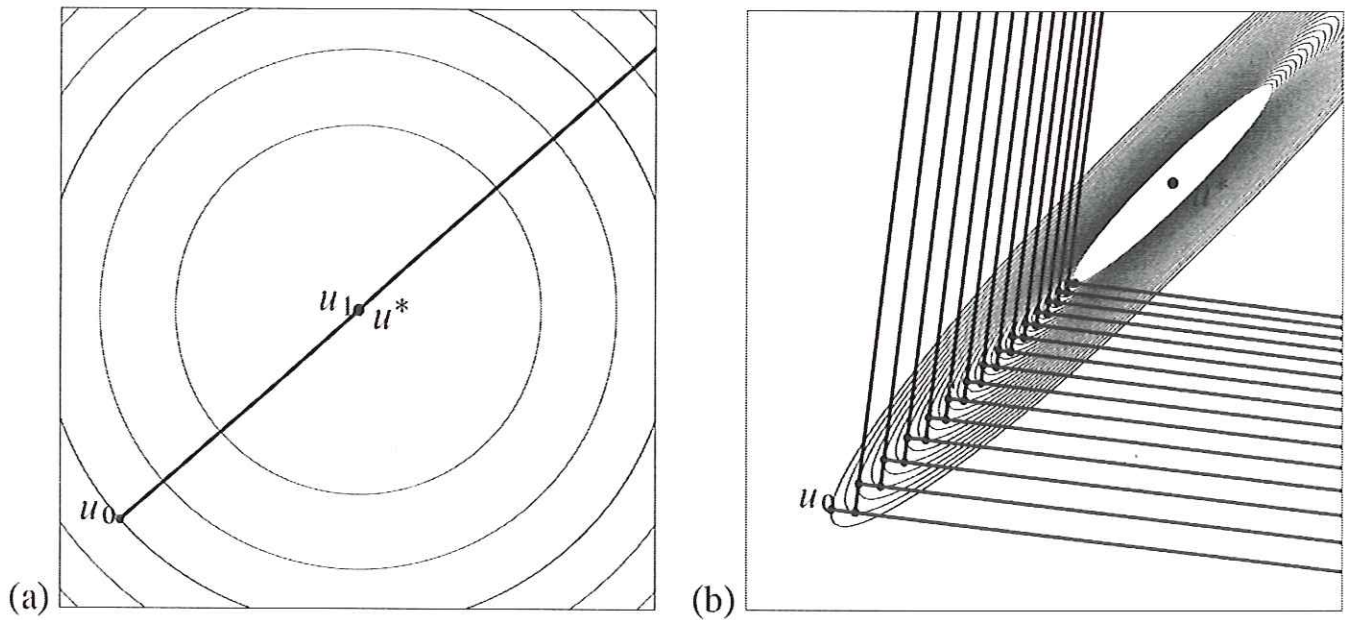


Figure 4.4. (a) If A is a scalar multiple of the identity, then the level sets of $\phi(u)$ are circular and steepest descent converges in one iteration from any initial guess u_0 . (b) If the level sets of $\phi(u)$ are far from circular, then steepest descent may converge slowly.

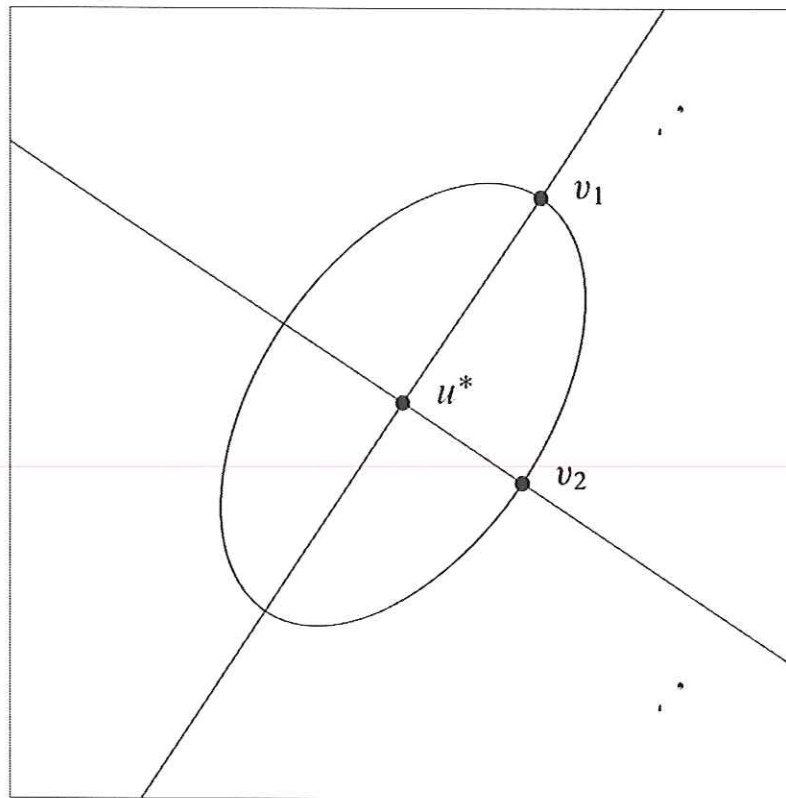


Figure 4.5. The major and minor axes of the elliptical level set of $\phi(u)$ point in the directions of the eigenvectors of A .

(B) Conjugate Gradient Method (CG)

→ The method of steepest descent can be generalized and accelerated if we choose a better search direction S_k , which will accelerate the iteration.

→ As already mentimed, the idea is to take

$$(i) S_k \neq -\nabla\phi(U_k) = r_k, \quad \&$$

$$(ii) S_k = P_k, \quad \text{where } P_k \text{ satisfies } \underline{A}\text{-conjugate}$$

$$P_k^T (A P_j) = 0, \quad j=0, \dots, k-1,$$

which means that

P_k is orthogonal to all previous $A P_j$.

→ Note that A-conjugacy does NOT mean that, for $A \neq I$,

$$(i) P_k \perp P_j, \quad j=0, \dots, k-1, \quad \text{but}$$

$$(ii) P_k \perp A P_j, \quad j=0, \dots, k-1.$$

→ If $A = I$, then this just means the conventional orthogonality of vectors P_k & P_j , $j=0, \dots, k-1$.

→ A-conjugacy is a generalized concept of orthogonality.

Prk. (i) A very bad choice of P_k would be a direction orthogonal to $r_k = -\nabla\phi(U_k)$, which means P_k is tangent to the contour line of ϕ at U_k . (see Fig 4.3)

\Rightarrow In this case, $\phi(U)$ increases along the tangent direction. and the iteration will never converge.

(ii) We therefore want P_k s.t. $P_k^T r_k \neq 0$, then the new point U_{k+1} will be different from U_k , and $\phi(U_{k+1}) < \phi(U_k)$, $\forall k$. (See Fig. 4.6)

(iii) $P_k = -\nabla\phi(U_k)$ will reduce to the steepest descent.

\rightarrow We want better than this & iteration scheme: $U_{k+1} = U_k + \alpha_k P_k$.

(iv) Assuming U^* is the min. pt. that we want to obtain, the question is then "How do we find P_k without knowing U^* ?"

The idea is to observe that, since $P_0 (=r_0)$ is tangent to the level set of ϕ at U_1 ,

$$\Rightarrow P_0 \perp -\nabla\phi(U_1) = r_1 = f - \underline{A}U_1$$

$$\begin{aligned} \Rightarrow 0 &= P_0^T r_1 = P_0^T (f - \underline{A}U_1), & f &= \underline{A}U^* \text{ (exact)} \\ &= P_0^T \underline{A} (U^* - U_1), & U^* - U_1 &= \alpha P_1, \text{ for some } \alpha \neq 0. \\ &= P_0^T \underline{A} \alpha P_1 & \dots & \textcircled{10} \end{aligned}$$

⇒ (I) implies, written in another equivalent form, and after dividing by d_1

$$0 = \underbrace{p_0^T A p_1}_{\in \mathbb{R}} = (p_1^T A p_0)^T = (p_1^T A p_0)^T$$

$$\textcircled{!} \quad \boxed{p_1^T A p_0 = 0} \quad \dots \quad \textcircled{II}$$

⇒ (II) means that p_0 & p_1 are A -conjugate, and if so, using p_1 as a new search direction, the CG iteration converges in only two steps, which is, in general, true for $m=2$.

⇒ If $m > 2$, then we can also show that CG converges in, at most, m iterations.

⇒ If A has n distinct eigenvalues, $A \in \mathbb{R}^m \times \mathbb{R}^m$, then we also see that CG converges in at most n iterations, $n \leq m$.

⇒ However, in practice, CG frequently "converges" to a sufficiently accurate appn to U^* in much less than n iterations.

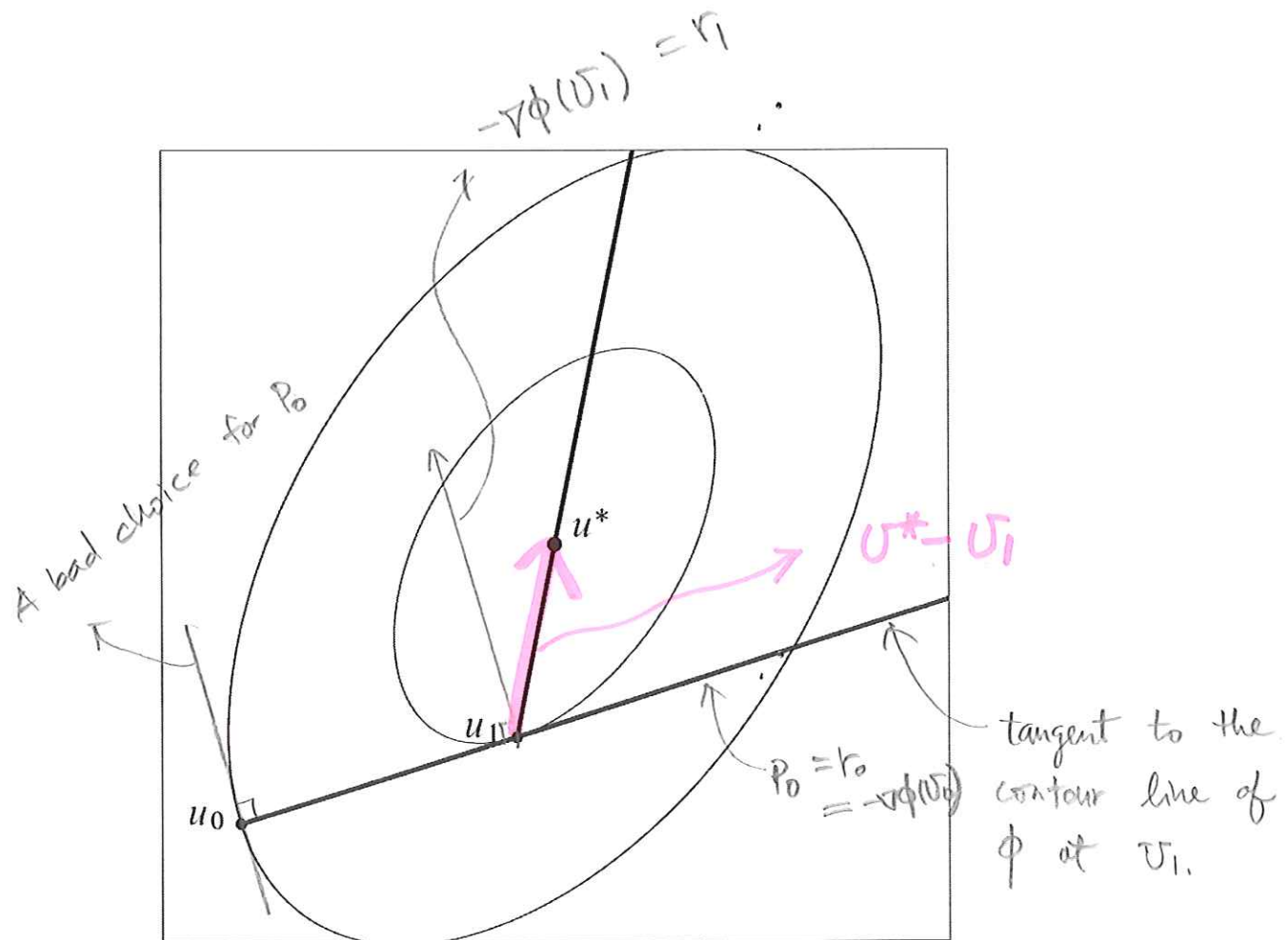


Figure 4.6. The CG algorithm converges in two iterations from any initial guess u_0 in the case $m = 2$. The two search directions used are A -conjugate.

Note that p_0 is tangent to the level set of ϕ at u_1 ,

$$\Rightarrow p_0 \perp -\nabla\phi(u_1) = r_1 = f - Au_1$$

$$\Rightarrow 0 = p_0^T r_1 = p_0^T (f - Au_1), \quad f = Au^*$$

$$= p_0^T A(u^* - u_1), \quad u^* - u_1 = \alpha p_1, \quad \alpha \neq 0$$

$$= p_0^T A \alpha p_1$$

Algorithm

Choose initial guess U_0 (possibly $U_0 = 0$).

$$r_0 = f - AU_0 \quad (\text{initial residual})$$

$$p_0 = r_0 \quad (\text{initial search direction})$$

for $k=1, 2, \dots$

$$W_{k-1} = AP_{k-1} \quad (\text{the only matrix-vector multiplication})$$

$$d_{k-1} = \frac{r_{k-1}^T r_{k-1}}{p_{k-1}^T W_{k-1}} \quad (\text{compute search parameter})$$

$$U_k = U_{k-1} + d_{k-1} p_{k-1} \quad (\text{update soln})$$

$$r_k = r_{k-1} - d_{k-1} W_{k-1} \quad (\text{compute new residual} \\ \neq \text{search direction})$$

If $\|r_k\| < \epsilon$, ϵ : tolerance, then

Stop

Else

$$\beta_{k-1} = \frac{r_k^T r_k}{r_{k-1}^T r_{k-1}}$$

(This is A-conjugacy
 $\Rightarrow p_k^T A p_j = 0, j=0, 1, \dots, k-1$)
(\rightarrow See Thm 4.1 LeVeque)

$$p_k = r_k + \beta_{k-1} p_{k-1} \quad (\text{compute new search direction})$$

Endif

Endfor

Compare this algorithm to the steepest descent algorithm presented on page 80. Up through the convergence check it is essentially the same except that the A -conjugate search direction p_{k-1} is used in place of the steepest descent search direction r_{k-1} in several places.

The final two lines in the loop determine the next search direction p_k . This simple choice gives a direction p_k with the required property that p_k is A -conjugate to all the previous search directions p_j for $j = 0, 1, \dots, k-1$. This is part of the following theorem, which is similar to Theorem 38.1 of Trefethen and Bau [91], although there it is assumed that $u_0 = 0$. See also Theorem 2.3.2 in Greenbaum [39].

Theorem 4.1. *The vectors generated in the CG algorithm have the following properties, provided $r_k \neq 0$ (if $r_k = 0$, then we have converged):*

1. p_k is A -conjugate to all the previous search directions, i.e., $p_k^T A p_j = 0$ for $j = 0, 1, \dots, k-1$.
2. The residual r_k is orthogonal to all previous residuals, $r_k^T r_j = 0$ for $j = 0, 1, \dots, k-1$.
3. The following three subspaces of \mathbb{R}^m are identical:

$$\begin{aligned} & \text{span}(p_0, p_1, p_2, \dots, p_{k-1}), \\ & \text{span}(r_0, A r_0, A^2 r_0, \dots, A^{k-1} r_0), \\ & \text{span}(A e_0, A^2 e_0, A^3 e_0, \dots, A^k e_0). \end{aligned} \tag{4.43}$$

The subspace $\mathcal{K}_k = \text{span}(r_0, A r_0, A^2 r_0, \dots, A^{k-1} r_0)$ spanned by the vector r_0 and the first $k-1$ powers of A applied to this vector is called a *Krylov space* of dimension k associated with this vector.

The iterate u_k is formed by adding multiples of the search directions p_j to the initial guess u_0 and hence must lie in the affine spaces $u_0 + \mathcal{K}_k$ (i.e., the vector $u_k - u_0$ is in the linear space \mathcal{K}_k).

We have seen that the CG algorithm can be interpreted as minimizing the function $\phi(u)$ over the space $u_0 + \text{span}(p_0, p_1, \dots, p_{k-1})$ in the k th iteration, and by the theorem above this is equivalent to minimizing $\phi(u)$ over the $u_0 + \mathcal{K}_k$. Many other iterative methods are also based on the idea of solving problems on an expanding sequence of Krylov spaces; see Section 4.4.

4.3.4 Convergence of conjugate gradient

The convergence theory for CG is related to the fact that u_k minimizes $\phi(u)$ over the affine space $u_0 + \mathcal{K}_k$ defined in the previous section. We now show that a certain norm of the error is also minimized over this space, which is useful in deriving estimates about the size of the error and rate of convergence.

Since A is assumed to be SPD, the A -norm defined by

$$\|e\|_A = \sqrt{e^T A e} \tag{4.44}$$

(C) Preconditioned CG (PCG).

→ Although it is a significant improvement over steepest descent, CG can still converge very slowly if

A is ill-conditioned.

⇔ A has a large condition number,

$$\text{cond}(A) \stackrel{\text{def}}{=} \|A\| \cdot \|A^{-1}\|.$$

⇔ A becomes very close to singular and any numerical operations involving

A become very inaccurate, not to mention inverting it.

(ex) $r_k = f - AU_k$

↑
this becomes inaccurate

(Note: $\text{Cond}(A) \geq 1$, $\text{cond}(I) = 1$)

→ In this case, we precondition A by implicitly using

$M^{-1}A$, instead of A , where

(i) M is a matrix for which $Mz_k = r_k$ is easily solved, with $r_k = f - AU_k$: residual, &

(ii) $M^{-1}A$ is relatively well-conditioned.

Given $AU_k = f$,
we consider $M^{-1}AU_k = M^{-1}f$
instead.

→ Typically, since CG only works for symmetric matrices,

and since $M^{-1}A$ is not guaranteed to be symmetric for CG, we in fact need to use

$\underline{\underline{L}}^{-1} \underline{\underline{A}} \underline{\underline{L}}$ instead of $\underline{\underline{M}}^{-1} \underline{\underline{A}}$, where $\underline{\underline{M}} = \underline{\underline{L}} \underline{\underline{L}}^T$.

→ However, fortunately, it turns out that the PCG algorithm can be suitably rearranged so that only $\underline{\underline{M}}$ is used and the corresponding matrix $\underline{\underline{L}}$ is not explicitly required. (See LeVeque)

→ The resulting PCG algorithm is essentially the same as the one for CG, but we solve the system

$$\underline{\underline{M}} z_k = r_k \quad \text{for} \quad z_k = \underline{\underline{M}}^{-1} r_k$$
$$= \underline{\underline{M}}^{-1} (f - \underline{\underline{A}} U_k)$$
$$= \underline{\underline{M}}^{-1} f - \underbrace{\underline{\underline{M}}^{-1} \underline{\underline{A}}}_{\uparrow} U_k$$

this is better than $\underline{\underline{A}} U_k$.

and use z_k in place of r_k in the CG algorithm.

Algorithm: PCG

choose initial guess U_0

$$r_0 = f - \underline{A}U_0$$

$$\text{Solve } \underline{M}z_0 = r_0 \text{ for } z_0 \quad (\Leftrightarrow z_0 = \underline{M}^{-1}r_0)$$

$$p_0 = z_0$$

for $k=1, 2, \dots$

$$w_{k-1} = \underline{A}p_{k-1} \quad (= \underline{A}z_{k-1} = \underline{A}\underline{M}^{-1}r_{k-1})$$

$$\alpha_{k-1} = \frac{r_{k-1}^T r_{k-1}}{p_{k-1}^T w_{k-1}}$$

$$U_k = U_{k-1} + \alpha_{k-1} p_{k-1}$$

$$r_k = r_{k-1} - \alpha_{k-1} w_{k-1}$$

If $\|r_k\| < \varepsilon$, ε : tolerance, then

stop

Else

$$\text{Solve } \underline{M}z_k = r_k \text{ for } z_k \quad \Leftrightarrow (z_k = \underline{M}^{-1}r_k)$$

$$p_{k-1} = \frac{z_k^T r_k}{r_k^T r_{k-1}}$$

$$p_k = z_k + p_{k-1} p_{k-1}$$

Endif

Endfor

→ The choice of an appropriate preconditioner \underline{M} :

(i) depends on the trade-off between the gain in the convergence rate & the increased cost per iteration due to \underline{M} .

(ii) is a big active area of research.

→ The most commonly used types of preconditioners :

(i) Diagonal (or Jacobi) : $\underline{M} = \text{diag}(\underline{A})$

(ii) Block diagonal (or Block Jacobi) :

$$I = \{i \mid i=1, \dots, m\}$$

$$I = I_1 \cup I_2 \cup \dots \cup I_s, \quad I_k : \text{disjoint each other}$$

$$\underline{M} = (m_{ij})_{i,j} \begin{cases} a_{ij}, & \text{if } i, j \in I_k, \text{ for some } k, \\ 0, & \text{if } i \in I_p, j \in I_q, p \neq q. \end{cases}$$

(iii) SSOR (symmetric SOR) :

$$\underline{A} = \underline{L} + \underline{D} + \underline{L}^T,$$

$$\underline{M} = (\underline{D} + \underline{L}) \underline{D}^{-1} (\underline{D} + \underline{L})^T.$$

→ We can show that $\text{cond}(\underline{M}^{-1}\underline{A}) = \mathcal{O}(\sqrt{\text{cond}(\underline{A})})$

(iv) incomplete Cholesky factorization,

(v) polynomial : \underline{M}^{-1} is taken to be a poly in \underline{A} that approximates \underline{A}^{-1} .

(vi) Approximate inverse.