

Numerical Methods II

Convergence of Gradient and Conjugate Gradient Methods

Marcel Oliver

1 Notation

We apply the Gradient and the Conjugate Gradient methods to the problem of finding the minimum of the quadratic functional

$$\Phi(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x} - \mathbf{x}^T \mathbf{b}, \quad (1)$$

where $A \in \mathbb{R}^{n \times n}$ is symmetric and positive definite. Recall that both methods are *descent methods*. They construct a minimizing sequence

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k \quad (2)$$

where \mathbf{d}_k is a descent direction, i.e. $\mathbf{d}_k^T \nabla \Phi(\mathbf{x}_k) < 0$, so that

$$\Phi(\mathbf{x}_k + \alpha \mathbf{d}_k) < \Phi(\mathbf{x}_k) \quad (3)$$

for small positive values of α . For quadratic functionals (1), the value of α that minimizes Φ along the line through \mathbf{x}_k in the direction \mathbf{d}_k is easily found to be

$$\alpha = \frac{\mathbf{d}_k^T \mathbf{r}_k}{\mathbf{d}_k^T A \mathbf{d}_k}. \quad (4)$$

We denote the error in the k -th step by

$$\mathbf{e}_k = \mathbf{x} - \mathbf{x}_k, \quad (5)$$

\mathbf{x} being the location of the true minimum, and define the *residual*

$$\mathbf{r}_k = A \mathbf{e}_k = \mathbf{b} - A \mathbf{x}_k. \quad (6)$$

The last equality is true since the location of the minimum is the solution of $A \mathbf{x} = \mathbf{b}$. See the notes from Numerical Methods I for details.

Finally, we introduce the A -norm

$$\|\mathbf{x}\|_A^2 \equiv \mathbf{x}^T A \mathbf{x}. \quad (7)$$

2 Convergence of the Gradient Method

In the gradient method we always walk down the direction of steepest descent, i.e.

$$\mathbf{d}_k = -\nabla\Phi(\mathbf{x}_k) = \mathbf{r}_k. \quad (8)$$

A direct computation of the error norm shows that

$$\begin{aligned} \|\mathbf{e}_{k+1}\|_A^2 &= \|\mathbf{x} - (\mathbf{x}_k + \alpha_k \mathbf{d}_k)\|_A^2 \\ &= \|\mathbf{e}_k - \alpha_k \mathbf{d}_k\|_A^2 \\ &= \|\mathbf{e}_k\|_A^2 - 2\alpha_k \mathbf{e}_k^T A \mathbf{d}_k + \alpha_k^2 \|\mathbf{d}_k\|_A^2 \\ &= \|\mathbf{e}_k\|_A^2 - 2 \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T A \mathbf{r}_k} \mathbf{e}_k^T A \mathbf{r}_k + \frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{(\mathbf{r}_k^T A \mathbf{r}_k)^2} \mathbf{r}_k^T A \mathbf{r}_k \\ &= \|\mathbf{e}_k\|_A^2 \left(1 - \frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{\mathbf{r}_k^T A \mathbf{r}_k \mathbf{r}_k^T A^{-1} \mathbf{r}_k} \right). \end{aligned} \quad (9)$$

To proceed further, we need the following result.

Lemma 1 (Kantorovich inequality). *Let $A \in \mathbb{R}^{n \times n}$ be symmetric and positive definite, and let $0 < \lambda_{\min} < \lambda_{\max}$ denote its smallest and largest eigenvalue, respectively. Then*

$$\min_{\mathbf{y} \neq 0} \frac{(\mathbf{y}^T \mathbf{y})^2}{\mathbf{y}^T A \mathbf{y} \mathbf{y}^T A^{-1} \mathbf{y}} = \frac{4 \lambda_{\min} \lambda_{\max}}{(\lambda_{\min} + \lambda_{\max})^2}. \quad (10)$$

Applying the Kantorovich inequality to the right side of (9) and noting that

$$1 - \frac{4 \lambda_{\min} \lambda_{\max}}{(\lambda_{\min} + \lambda_{\max})^2} = \frac{(\lambda_{\min} - \lambda_{\max})^2}{(\lambda_{\min} + \lambda_{\max})^2}, \quad (11)$$

we find that

$$\|\mathbf{e}_{k+1}\|_A \leq \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}} \|\mathbf{e}_k\|_A \equiv \frac{\kappa - 1}{\kappa + 1} \|\mathbf{e}_k\|_A, \quad (12)$$

where $\kappa = \lambda_{\max}/\lambda_{\min}$ is the so-called spectral condition number of A .

Proof of the Kantorovich inequality. The Kantorovich inequality (10) is invariant under rescaling of \mathbf{y} . It is therefore sufficient to prove it for arbitrary unit vectors. Let $\mathbf{v}_1, \dots, \mathbf{v}_n$ be an orthonormal basis of \mathbb{R}^n consisting of eigenvectors of A with corresponding eigenvalues $\lambda_1, \dots, \lambda_n$. We can write

$$\mathbf{y} = \sum_{i=1}^n y_i \mathbf{v}_i \quad (13)$$

so that

$$\|\mathbf{y}\|_2^2 = \sum_{i=1}^n y_i^2 = 1, \quad (14)$$

$$\mathbf{y}^T A \mathbf{y} = \sum_{i=1}^n y_i^2 \lambda_i, \quad (15)$$

$$\mathbf{y}^T A^{-1} \mathbf{y} = \sum_{i=1}^n y_i^2 \lambda_i^{-1}. \quad (16)$$

Thus, proving the Kantorovich inequality reduces to solving the constraint optimization problem *find the maximum of*

$$g(\mathbf{y}) = \mathbf{y}^T A \mathbf{y} \mathbf{y}^T A^{-1} \mathbf{y} \quad (17)$$

under the constraint

$$h(\mathbf{y}) = \mathbf{y}^T \mathbf{y} - 1 = 0. \quad (18)$$

Such maximum must necessarily satisfy

$$\nabla g(\mathbf{y}) = \mu \nabla h(\mathbf{y}), \quad (19)$$

where μ is the Lagrange multiplier. By direct computation,

$$\frac{\partial g}{\partial y_i} = 2 y_i \lambda_i \mathbf{y}^T A^{-1} \mathbf{y} + 2 y_i \lambda_i^{-1} \mathbf{y}^T A \mathbf{y}, \quad (20)$$

$$\frac{\partial h}{\partial y_i} = 2 y_i, \quad (21)$$

so that (19) reads

$$y_i \lambda_i \mathbf{y}^T A^{-1} \mathbf{y} + y_i \lambda_i^{-1} \mathbf{y}^T A \mathbf{y} = \mu y_i, \quad (22)$$

or

$$(\lambda_i^2 \mathbf{y}^T A^{-1} \mathbf{y} - \lambda_i \mu + \mathbf{y}^T A \mathbf{y}) \frac{y_i}{\lambda_i} = 0. \quad (23)$$

The expression in parenthesis is a quadratic equation in λ_i , and can only be zero for at most two distinct eigenvalues. Therefore, there can be at most two non-zero components y_i and y_j . (If an eigenvalue occurs with multiplicity larger than one, $g(\mathbf{y})$ depends only on the norm of the projection of \mathbf{y} onto the corresponding eigenspace, so the statement remains valid.) Dividing (22) by y_i and equating the expressions for indices i and j , we obtain

$$\lambda_i \left(\frac{y_i^2}{\lambda_i} + \frac{y_j^2}{\lambda_j} \right) + \frac{y_i^2 \lambda_i + y_j^2 \lambda_j}{\lambda_i} = \lambda_j \left(\frac{y_i^2}{\lambda_i} + \frac{y_j^2}{\lambda_j} \right) + \frac{y_i^2 \lambda_i + y_j^2 \lambda_j}{\lambda_j}, \quad (24)$$

or

$$(y_j^2 - y_i^2) \left(\frac{\lambda_i}{\lambda_j} + \frac{\lambda_j}{\lambda_i} - 2 \right) = (y_j^2 - y_i^2) \frac{(\lambda_i - \lambda_j)^2}{\lambda_i \lambda_j}. \quad (25)$$

Therefore, $y_i^2 = y_j^2 = \frac{1}{2}$ unless $\lambda_i = \lambda_j$ and, by direct substitution into (17), candidates for the maximum value are

$$\frac{1}{4}(\lambda_i + \lambda_j) \left(\frac{1}{\lambda_i} + \frac{1}{\lambda_j} \right) = \frac{1}{4} \left(\frac{\lambda_i}{\lambda_j} + 1 \right) \left(\frac{\lambda_j}{\lambda_i} + 1 \right) = \frac{(\lambda_i + \lambda_j)^2}{4 \lambda_i \lambda_j}. \quad (26)$$

The expression in the middle shows that it is increasing in the ratio λ_i/λ_j when $\lambda_i > \lambda_j$; we thus take $\lambda_i = \lambda_{\max}$ and $\lambda_j = \lambda_{\min}$ to complete the proof. \square

3 Convergence of the CG Method

Recall that one of the key properties of the Conjugate Gradient method is that each new iterate is optimal with respect to *all* descent directions from the so-called Krylov subspace

$$\begin{aligned} V_k &= \text{Span}\{\mathbf{d}_1, \dots, \mathbf{d}_{k-1}\} \\ &= \text{Span}\{A^0 \mathbf{r}_1, \dots, A^{k-2} \mathbf{r}_1\}. \end{aligned} \quad (27)$$

For details, see the handout on the derivation of the CG method. According to (6), we can also write

$$V_k = \text{Span}\{A^1 \mathbf{e}_1, \dots, A^{k-1} \mathbf{e}_1\}. \quad (28)$$

Recall also the recursion for the residual,

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k A \mathbf{d}_k \equiv \mathbf{r}_k + A \mathbf{w}_{k+1} \quad (29)$$

where \mathbf{w}_{k+1} is some vector from V_{k+1} . By definition of \mathbf{e}_k , we can also write

$$\mathbf{e}_{k+1} = \mathbf{e}_k + \mathbf{w}_{k+1}. \quad (30)$$

We conclude inductively that $\mathbf{e}_k - \mathbf{e}_1 \in V_k$ for $k = 2, 3, \dots$, and we can write

$$\begin{aligned} \mathbf{e}_k &= \mathbf{e}_1 + \sum_{i=1}^{k-1} \gamma_i A^i \mathbf{e}_1 \\ &\equiv \phi_k(A) \mathbf{e}_1 \end{aligned} \quad (31)$$

where

$$\phi_k(x) = \sum_{i=0}^{k-1} \gamma_i x^i \quad (32)$$

is some polynomial of degree $k-1$ with $\phi_k(0) = 1$.

Since

$$\begin{aligned}
\|\mathbf{e}_k\|_A^2 &= (\mathbf{x} - \mathbf{x}_k)^T A (\mathbf{x} - \mathbf{x}_k) \\
&= \mathbf{x}^T A \mathbf{x} - 2 \mathbf{x}_k^T A \mathbf{x} + \mathbf{x}_k^T A \mathbf{x}_k \\
&= \mathbf{x}^T \mathbf{b} - 2 \mathbf{x}_k^T \mathbf{b} + \mathbf{x}_k^T A \mathbf{x}_k \\
&= 2 \Phi(\mathbf{x}_k) + \mathbf{x}^T \mathbf{b},
\end{aligned} \tag{33}$$

last term on the right being independent of \mathbf{x}_k , optimality of \mathbf{x}_k with respect to a certain subspace is equivalent to optimality of the A -norm of \mathbf{e}_k with respect to the same subspace. In other words, CG is constructed in such a way that the polynomial ϕ_k which appears in (31) is the polynomial that minimizes the A norm of \mathbf{e}_k among all polynomials of the same or lesser degree.

As in the proof of the Kantorovich inequality we express \mathbf{e}_1 in terms of the orthonormal eigenvectors of A , i.e.

$$\mathbf{e}_1 = \sum_{j=1}^n y_j \mathbf{v}_j, \tag{34}$$

so that

$$\mathbf{e}_k = \sum_{j=1}^n y_j \phi_k(\lambda_j) \mathbf{v}_j, \tag{35}$$

$$\|\mathbf{e}_k\|_A^2 = \sum_{j=1}^n y_j^2 \phi_k^2(\lambda_j) \lambda_j. \tag{36}$$

Let P_k denote the vector space of polynomials of degree less or equal to k . The optimality condition of CG can therefore be expressed as

$$\begin{aligned}
\|\mathbf{e}_k\|_A^2 &= \min_{\substack{\phi \in P_{k-1} \\ \phi(0)=1}} \sum_{j=1}^n y_j^2 \phi^2(\lambda_j) \lambda_j \\
&\leq \min_{\substack{\phi \in P_{k-1} \\ \phi(0)=1}} \max_{\lambda \in [\lambda_{\min}, \lambda_{\max}]} \phi^2(\lambda) \sum_{j=1}^n y_j^2 \lambda_j \\
&= \min_{\substack{\phi \in P_{k-1} \\ \phi(0)=1}} \max_{\lambda \in [\lambda_{\min}, \lambda_{\max}]} \phi^2(\lambda) \|\mathbf{e}_1\|_A^2.
\end{aligned} \tag{37}$$

Let T_k denote the Chebychev polynomial of degree k . Then the following is true.

Lemma 2. *If $0 < \lambda_{\min} < \lambda_{\max}$, then*

$$\min_{\substack{\phi \in P_k \\ \phi(0)=1}} \max_{\lambda \in [\lambda_{\min}, \lambda_{\max}]} |\phi(\lambda)| = T_k \left(\frac{\lambda_{\max} + \lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} \right)^{-1} \leq 2 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k, \tag{38}$$

where $\kappa = \lambda_{\max}/\lambda_{\min}$.

Inserting this result into (37), we find that

$$\|e_k^{\text{cg}}\|_A \leq 2 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{k-1} \|e_1\|_A, \quad (39)$$

whereas for the gradient method, from equation (12),

$$\|e_k^{\text{grad}}\|_A \leq \left(\frac{\kappa - 1}{\kappa + 1} \right)^{k-1} \|e_1\|_A. \quad (40)$$

This shows that for poorly conditioned matrices the rate of convergence of the conjugate gradient method is much better than that of the gradient method.

Proof of Lemma 2. Recall that the Chebychev polynomials of order k on the interval $[-1, 1]$ can be written in the form

$$T_k(x) = \cos(k \arccos x). \quad (41)$$

It is clear that $|T_k| \leq 1$ on this interval. Moreover, $T_k(1) = 1$, $T_k(-1) = (-1)^k$, and T_k has $k - 1$ distinct extrema on $(-1, 1)$, alternating between 1 and -1 .

To realize the minimum on the left of (38), we need a polynomial that is uniformly small on the interval $[\lambda_{\min}, \lambda_{\max}]$. Since T_k is of uniform size on $[-1, 1]$, we are tempted to use it as a template for constructing a candidate minimizing polynomial by remapping the interval $[\lambda_{\min}, \lambda_{\max}]$ onto $[-1, 1]$ by a linear affine change of variable ℓ . Setting

$$\phi(x) = \frac{T_k(\ell(x))}{T_k(\ell(0))} \quad (42)$$

will then also satisfy $\phi(0) = 1$. The ansatz $\ell(x) = mx + c$ together with the requirement that $\ell(\lambda_{\min}) = 1$ and $\ell(\lambda_{\max}) = -1$ gives

$$\ell(x) = \frac{\lambda_{\max} + \lambda_{\min} - 2x}{\lambda_{\max} - \lambda_{\min}}. \quad (43)$$

Therefore,

$$\max_{\lambda \in [\lambda_{\min}, \lambda_{\max}]} |\phi(\lambda)| = \frac{1}{|T_k(\ell(0))|} \max_{x \in [-1, 1]} |T_k(x)| = T_k \left(\frac{\lambda_{\max} + \lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} \right)^{-1} \quad (44)$$

To complete the proof of the equality in (38), we need to show that no other polynomial of degree k can yield a bound lower than (44). Assume the contrary, and let ψ denote a polynomial satisfying $\psi(0) = 1$ and

$$\max_{\lambda \in [\lambda_{\min}, \lambda_{\max}]} |\psi(\lambda)| < \max_{\lambda \in [\lambda_{\min}, \lambda_{\max}]} |\phi(\lambda)|. \quad (45)$$

Then $\psi(0) - \phi(0) = 0$, and the graph of ψ must intersect the graph of ϕ exactly k times in the interval $[\lambda_{\min}, \lambda_{\max}]$. This means that $\phi - \psi$ is a polynomial of degree k with $k + 1$ zeroes, and must therefore be identically zero.

The proof of the inequality in (38) is a homework exercise. \square