

Mathematical Tripos Part II: Michaelmas Term 2023

Numerical Analysis – Lecture 18

4.1 Steepest descent and conjugate gradient methods

For solving $A\mathbf{x} = \mathbf{b}$ with a symmetric positive definite matrix A , we consider iterative methods based on an optimization formulation. Consider the convex quadratic function

$$F(\mathbf{x}) := \frac{1}{2}\langle \mathbf{x}, A\mathbf{x} \rangle - \langle \mathbf{b}, \mathbf{x} \rangle \quad (4.5)$$

where $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^T \mathbf{v}$ is the Euclidean inner product. Note that the global minimizer of F is $\mathbf{x}^* = A^{-1}\mathbf{b}$. Indeed

$$F(\mathbf{x}^* + \mathbf{h}) - F(\mathbf{x}^*) = \langle \mathbf{h}, A\mathbf{x}^* - \mathbf{b} \rangle + \frac{1}{2}\langle \mathbf{h}, A\mathbf{h} \rangle \geq 0$$

for any \mathbf{h} . Observe that F can also be written as

$$F(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}^* - \mathbf{x}\|_A^2 + \text{constant}$$

where $\|\mathbf{y}\|_A := \langle \mathbf{y}, A\mathbf{y} \rangle^{1/2} = \sqrt{\mathbf{y}^T A \mathbf{y}}$ is the A -norm of A . (The constant in the above formulation is a term that does not depend on \mathbf{x} , so it is irrelevant for the purpose of minimizing F , the constant is $\frac{1}{2}\mathbf{b}^T A^{-1}\mathbf{b}$.)

Gradient/Steepest descent The *gradient descent* method for minimizing F has iterates

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla F(\mathbf{x}^{(k)})$$

where $\nabla F(\mathbf{x}^{(k)})$ is the gradient of F at $\mathbf{x}^{(k)}$, and $\alpha_k > 0$ is the step size. For our quadratic function, it is easy to verify that

$$\nabla F(\mathbf{x}^{(k)}) = A\mathbf{x}^{(k)} - \mathbf{b} = -\mathbf{r}^{(k)}$$

where $\mathbf{r}^{(k)} = \mathbf{b} - A\mathbf{x}^{(k)}$ is the residual. There are multiple ways to choose the step size α_k :

- *Constant step-size* $\alpha_k = \alpha$. In this case the iteration takes the form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha(A\mathbf{x}^{(k)} - \mathbf{b}) = (I - \alpha A)\mathbf{x}^{(k)} + \alpha \mathbf{b}$$

which is nothing but a Jacobi-like iteration with $D = \alpha^{-1}I$ (we say Jacobi-like because the diagonal of A is not necessarily equal to $\alpha^{-1}I$). We know from previous lectures that the method converges iff

$$\rho(I - \alpha A) < 1 \iff |1 - \alpha \lambda_i| < 1 \quad \forall \lambda_i \text{ eigenvalues of } A \iff 0 < \alpha < 2/\rho(A).$$

For example, assume the eigenvalues of A are all in $[l, L]$ where $0 < l < L$. Then one can choose $\alpha = 1/L$, and in this case the convergence rate is given by $\rho(I - \frac{1}{L}A) = 1 - l/L$, i.e., the error $\|\mathbf{x}^* - \mathbf{x}^{(k)}\|$ decays like $(1 - l/L)^k$. The quantity $L/l \geq 1$ is known as the condition number of A . We see that, as the condition number grows, the convergence rate becomes worse and worse.

- *Exact line search*. Another way to choose the step size α_k is using line search. Here α_k is chosen so that it achieves the smallest possible value of F along the search direction, i.e., $\alpha_k = \arg \min_{\alpha} F(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$ where $\mathbf{d}^{(k)}$ is the search direction, equal to the negative gradient. Because our function is quadratic, one can get a closed form expression for the optimal α .

Lemma 4.20 Let F be the function defined in (4.5). Let $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{r} = \mathbf{b} - A\mathbf{x}$ be the residual and let $\mathbf{d} \in \mathbb{R}^n$ be a search direction. Then

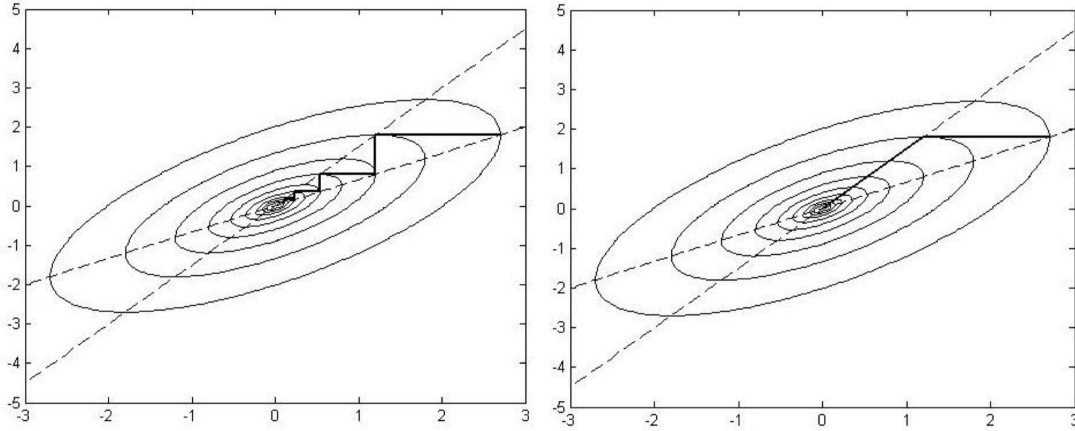
$$\arg \min_{\alpha} F(\mathbf{x} + \alpha \mathbf{d}) = \frac{\langle \mathbf{r}, \mathbf{d} \rangle}{\langle \mathbf{d}, A\mathbf{d} \rangle}. \quad (4.6)$$

Proof. The function $F(\mathbf{x} + \alpha \mathbf{d}) = F(\mathbf{x}) - \alpha \langle \mathbf{r}, \mathbf{d} \rangle + \alpha^2 / 2 \langle \mathbf{d}, A \mathbf{d} \rangle$ is quadratic in the single variable α . The minimum is attained at α s.t. $-\langle \mathbf{r}, \mathbf{d} \rangle + \alpha \langle \mathbf{d}, A \mathbf{d} \rangle = 0$ which gives the desired formula. \square

The gradient descent method with exact line search thus takes the form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \frac{\|\mathbf{r}^{(k)}\|_2^2}{\|\mathbf{r}^{(k)}\|_A^2} \mathbf{r}^{(k)},$$

where we used the fact that the gradient direction is $\mathbf{d} = -\nabla F(\mathbf{x}^{(k)}) = \mathbf{r}^{(k)}$. It can be shown that the speed of convergence of the gradient descent with exact line search is, like with the constant step size, $\approx (1 - l/L)^k$ where $0 < l < L$ are the smallest and largest eigenvalues of A . The figure below (left) shows an example of the gradient descent method with exact line search applied to a two-dimensional quadratic function F . Note the zig-zag behaviour of the iterates.



(a) Worst case scenario of steepest descent

(b) Conjugate gradient method applied to the same problem as in (a)

Conjugate directions Let's revisit equation (4.6) for a general direction \mathbf{d} (i.e., not necessarily equal to the negative gradient). Assume $\mathbf{x} = \mathbf{x}^{(k)}$, and let $\mathbf{e}^{(k)} = \mathbf{x}^* - \mathbf{x}^{(k)}$ be the error and $\mathbf{r}^{(k)} = \mathbf{b} - A\mathbf{x}^{(k)} = A\mathbf{e}^{(k)}$ be the residual. Then we can write $\langle \mathbf{r}^{(k)}, \mathbf{d} \rangle = \langle \mathbf{e}^{(k)}, \mathbf{d} \rangle_A$, and so for a general search direction \mathbf{d} with an exact line search, the iterate takes the form $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \frac{\langle \mathbf{e}^{(k)}, \mathbf{d} \rangle_A}{\langle \mathbf{d}, \mathbf{d} \rangle_A} \mathbf{d}$. By subtracting \mathbf{x}^* , the iterates in terms of the error $\mathbf{e}^{(k+1)}$ are given by:

$$\mathbf{e}^{(k+1)} = \mathbf{e}^{(k)} - \frac{\langle \mathbf{e}^{(k)}, \mathbf{d} \rangle_A}{\langle \mathbf{d}, \mathbf{d} \rangle_A} \mathbf{d}. \quad (4.7)$$

Geometrically, this means that $\mathbf{e}^{(k+1)}$ is the projection of $\mathbf{e}^{(k)}$ onto the hyperplane that is A -orthogonal to \mathbf{d} , i.e., we have

$$\langle \mathbf{e}^{(k+1)}, \mathbf{d} \rangle_A = 0 \quad (4.8)$$

Definition 4.21 (Conjugate directions) The vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ are *conjugate* with respect to a symmetric positive definite matrix A if they are nonzero and A -orthogonal: $\langle \mathbf{u}, \mathbf{v} \rangle_A := \langle \mathbf{u}, A\mathbf{v} \rangle = 0$.

The observation above allows us to prove the following important result.

Theorem 4.22 Let $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(n-1)}$ be n nonzero pairwise conjugate directions, and consider the sequence of iterates

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}, \quad \alpha_k = \frac{\langle \mathbf{r}^{(k)}, \mathbf{d}^{(k)} \rangle}{\langle \mathbf{d}^{(k)}, A\mathbf{d}^{(k)} \rangle}.$$

Let $\mathbf{r}^{(k)} = \mathbf{b} - A\mathbf{x}^{(k)}$ be the residual. Then for each $k = 1, \dots, n$, $\mathbf{r}^{(k)}$ is orthogonal to $\text{span}\{\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k-1)}\}$. In particular $\mathbf{r}^{(n)} = 0$.

Proof. Since $\mathbf{r}^{(k)} = A\mathbf{e}^{(k)}$, it suffices to show that $\mathbf{e}^{(k)}$ is A -orthogonal to $\text{span}\{\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k-1)}\}$. The proof is by induction on k . For $k = 0$ there is nothing to prove. Assume the statement is true for $k \geq 0$, and consider the equation (4.7) (with $\mathbf{d} = \mathbf{d}^{(k)}$). From the induction hypothesis, and the fact that the $\mathbf{d}^{(i)}$ are pairwise conjugate directions, we see that $\mathbf{e}^{(k+1)}$ is A -orthogonal to $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k-1)}$. Furthermore, we have already seen in (4.8) that $\langle \mathbf{e}^{(k+1)}, \mathbf{d}^{(k)} \rangle_A = 0$. Thus this shows that $\mathbf{e}^{(k+1)}$ is A -orthogonal to $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k)}$ as desired. \square

So, if a sequence $(\mathbf{d}^{(k)})$ of conjugate directions is at hand, we have an iterative procedure with good approximation properties. In the conjugate gradient method, the (A -orthogonal) basis of conjugate directions is constructed by A -orthogonalization of the sequence of gradients of F at the $\mathbf{x}^{(k)}$; or equivalently the sequence of residuals $\{\mathbf{r}^{(0)}, \dots, \mathbf{r}^{(k)}\}$.