
Lecture 4

Iterative algorithms for solving a problem of the form

$$\text{minimize } f(\mathbf{x}) \quad (4.1)$$

on \mathbb{R}^d generate a sequence of vectors $\mathbf{x}_0, \mathbf{x}_1, \dots$ in the hope that this sequence converges to a (local or global) minimizer \mathbf{x}^* of (4.1). In this lecture we study what it means for a sequence to converge, and how to quantify the speed of convergence. The concepts are applied to gradient descent.

4.1 Convergence of iterative methods

A sequence of vectors $\{\mathbf{x}_k\}$ in \mathbb{R}^d , $k \geq 0$, *converges* to a vector \mathbf{x}^* with respect to a norm $\|\cdot\|$ as $k \rightarrow \infty$, written $\mathbf{x}_k \rightarrow \mathbf{x}^*$, if the sequence of numbers $\|\mathbf{x}_k - \mathbf{x}^*\|$ converges to zero. More formally, if for every $\varepsilon > 0$ there exists an index N such that for all $n \geq N$,

$$\|\mathbf{x}_n - \mathbf{x}^*\| < \varepsilon.$$

Iterative algorithms will rarely find the exact solution to a problem like (4.1), so we will usually be happy to find a solution that differs from the true one by at most some specified accuracy.

Definition 4.1. A sequence of vectors $\{\mathbf{x}_k\}$, $k \geq 0$, is said to converge to \mathbf{x}^*

- (a) linearly (or Q-linear, Q for quotient), if there exist an $r \in (0, 1)$ such that for sufficiently large k ,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq r \|\mathbf{x}_k - \mathbf{x}^*\|.$$

- (b) superlinearly, if

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = 0,$$

- (c) with order p , if there exists a constant $M > 0$, such that for sufficiently large k ,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq M \|\mathbf{x}_k - \mathbf{x}^*\|^p.$$

The case $p = 2$ is called *quadratic convergence*.

Of course, as mentioned earlier, these definitions depend on the choice of a norm. It can be shown that quadratic convergence implies superlinear convergence, and superlinear convergence implies linear convergence.

Example 4.2. Consider the sequence of numbers $x_k = 1/2^{r^k}$ for some $r > 1$. Clearly, $x_k \rightarrow x^* = 0$ as $k \rightarrow \infty$. Moreover,

$$x_{k+1} = \frac{1}{2^{r^{k+1}}} = \frac{1}{2^{r^k r}} = \left(\frac{1}{2^{r^k}}\right)^r = x_k^r,$$

which shows that the sequence has rate of convergence r .

4.2 Convergence of gradient descent

In this section, a norm $\|\cdot\|$ will refer to the 2-norm, unless otherwise stated. We now study the convergence of gradient descent for the least squares problem

$$\text{minimize } f(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2, \quad (4.1)$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m \geq n$ of full rank. As we have seen in Lecture 3, the gradient descent method is the procedure

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{r}_k,$$

where the step length and the *residual* are given by

$$\alpha_k = \frac{\|\mathbf{r}_k\|^2}{\|\mathbf{A}\mathbf{r}_k\|^2}, \quad \mathbf{r}_k = \mathbf{A}^\top(\mathbf{b} - \mathbf{A}\mathbf{x}_k) = -\nabla f(\mathbf{x}_k).$$

At the minimizer, the residual is

$$\mathbf{r} = -\nabla f(\mathbf{x}^*) = \mathbf{A}^\top(\mathbf{b} - \mathbf{A}\mathbf{x}^*) = 0, \quad (4.2)$$

and as the sequence \mathbf{x}_k converges to \mathbf{x}^* , the norms of the residuals converge to 0. Conversely, the residual is related to the difference $\mathbf{x}_k - \mathbf{x}^*$ by

$$\mathbf{r}_k = \mathbf{A}^\top(\mathbf{b} - \mathbf{A}\mathbf{x}_k) = \mathbf{A}^\top(\mathbf{b} - \mathbf{A}\mathbf{x}_k - (\mathbf{b} - \mathbf{A}\mathbf{x}^*)) = \mathbf{A}^\top \mathbf{A}(\mathbf{x}_k - \mathbf{x}^*), \quad (4.3)$$

where we used the “intelligent zero” (4.2). Therefore

$$\|\mathbf{x}_k - \mathbf{x}^*\| = \|(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{r}\| \leq \|(\mathbf{A}^\top \mathbf{A})^{-1}\| \|\mathbf{r}_k\|,$$

where $\|\mathbf{B}\| = \max_{\mathbf{x} \neq 0} \|\mathbf{B}\mathbf{x}\|/\|\mathbf{x}\|$ is the operator norm of a matrix \mathbf{B} with respect to the 2-norm. Consequently, if the sequence $\|\mathbf{r}_k\|$ converges to zero, so does the sequence $\|\mathbf{x}_k - \mathbf{x}^*\|$. A reasonable criterium to stop the algorithm is therefore when the residual norm $\|\mathbf{r}_k\|$ is below a predefined tolerance ε .

For the following, we use the norm

$$\|\mathbf{x}\|_A^2 := \langle \mathbf{A}\mathbf{x}, \mathbf{A}\mathbf{x} \rangle = \|\mathbf{A}\mathbf{x}\|^2.$$

We have the bounds

$$\|\mathbf{A}^\dagger\|^{-1}\|\mathbf{x}\| \leq \|\mathbf{x}\|_A \leq \|\mathbf{A}\|\|\mathbf{x}\|,$$

where $\mathbf{A}^\dagger = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top$ is the *Moore-Penrose pseudoinverse* of \mathbf{A} . The important consequence of these inequalities is that convergence with respect to the $\|\cdot\|_A$ norm is equivalent to convergence with respect to the 2-norm.

A simple but nevertheless very useful observation is the following.

Lemma 4.1. *The difference between the function value at some $\mathbf{x} \in \mathbb{R}^d$ and the optimal value is*

$$f(\mathbf{x}) - f(\mathbf{x}^*) = \frac{1}{2} \|\mathbf{x} - \mathbf{x}^*\|_A^2.$$

Proof. This proceeds by straight-forward calculation:

$$\begin{aligned} f(\mathbf{x}^*) - f(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{x}^*\|_A^2 &= \frac{1}{2} \|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|^2 - \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 + \frac{1}{2} \|\mathbf{A}(\mathbf{x} - \mathbf{x}^*)\|^2 \\ &= \frac{1}{2} \|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|^2 - \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 + \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b} - (\mathbf{A}\mathbf{x}^* - \mathbf{b})\|^2 \\ &= \|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|^2 - \langle \mathbf{A}\mathbf{x}^* - \mathbf{b}, (\mathbf{A}\mathbf{x} - \mathbf{b}) \rangle \\ &= \langle \mathbf{A}\mathbf{x}^* - \mathbf{b}, \mathbf{A}\mathbf{x}^* - \mathbf{b} - (\mathbf{A}\mathbf{x} - \mathbf{b}) \rangle \\ &= \langle \mathbf{A}\mathbf{x}^* - \mathbf{b}, \mathbf{A}(\mathbf{x}^* - \mathbf{x}) \rangle \\ &= \langle \mathbf{A}^\top (\mathbf{A}\mathbf{x}^* - \mathbf{b}), \mathbf{x}^* - \mathbf{x} \rangle = \langle \nabla f(\mathbf{x}^*), \mathbf{x}^* - \mathbf{x} \rangle = 0, \end{aligned}$$

where the last equality follows from (4.2). \square

We next give a bound on the number of iterations that are guaranteed to bring the residual below a certain value. We then state (without proof) a more precise convergence bound in terms of the condition number, that shows that gradient descent for linear least squares has linear convergence.

Theorem 4.2. *Let $\varepsilon > 0$ be given and assume*

$$N > \frac{2\|\mathbf{A}\|^2}{\varepsilon^2} \|\mathbf{x}_0 - \mathbf{x}^*\|_A^2$$

for a starting point \mathbf{x}_0 , where \mathbf{x}^ is the minimizer of $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$. If $\{\mathbf{x}_k\}_{k \geq 0}$ is the sequence of vectors generated by gradient descent for this function, then*

$$\min_{0 \leq k \leq N} \|\mathbf{r}_k\| < \varepsilon.$$

We first derive a bound on the decrease of the function value at each step.

Lemma 4.3. *For the function f as in (4.1),*

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) = \frac{1}{2} \frac{\|\mathbf{r}_k\|^4}{\|\mathbf{A}\mathbf{r}_k\|^2} \geq \frac{1}{2} \frac{\|\mathbf{r}_k\|^2}{\|\mathbf{A}\|^2},$$

where $\|\mathbf{A}\| := \max_{\mathbf{x} \neq \mathbf{0}} \|\mathbf{A}\mathbf{x}\|/\|\mathbf{x}\|$ is the operator norm of \mathbf{A} with respect to the 2-norm.

Proof. Using the identity $\|\mathbf{x}\|^2 = \langle \mathbf{x}, \mathbf{x} \rangle$ and the bilinearity

$$\|\mathbf{x} + \mathbf{y}\|^2 = \langle \mathbf{x} + \mathbf{y}, \mathbf{x} + \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{x} \rangle + 2\langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{y}, \mathbf{y} \rangle = \|\mathbf{x}\|^2 + 2\langle \mathbf{x}, \mathbf{y} \rangle + \|\mathbf{y}\|^2,$$

we compute

$$\begin{aligned} f(\mathbf{x}_{k+1}) &= f(\mathbf{x}_k + \alpha_k \mathbf{r}_k) = \frac{1}{2} \|\mathbf{A}\mathbf{x}_k + \alpha_k \mathbf{A}\mathbf{r}_k - \mathbf{b}\|^2 \\ &= \frac{1}{2} \|\mathbf{A}\mathbf{x}_k - \mathbf{b}\|^2 + \langle \mathbf{A}\mathbf{x}_k - \mathbf{b}, \alpha_k \mathbf{A}\mathbf{r}_k \rangle + \frac{1}{2} \alpha_k^2 \|\mathbf{A}\mathbf{r}_k\|^2 \\ &= \frac{1}{2} \|\mathbf{A}\mathbf{x}_k - \mathbf{b}\|^2 + \alpha_k \langle \mathbf{A}^\top (\mathbf{A}\mathbf{x}_k - \mathbf{b}), \mathbf{r}_k \rangle + \frac{1}{2} \alpha_k^2 \|\mathbf{A}\mathbf{r}_k\|^2 \\ &= f(\mathbf{x}_k) - \alpha_k \langle \mathbf{r}_k, \mathbf{r}_k \rangle + \frac{1}{2} \alpha_k^2 \|\mathbf{A}\mathbf{r}_k\|^2. \end{aligned} \tag{4.4}$$

Substituting $\alpha_k = \|\mathbf{r}_k\|^2 / \|\mathbf{A}\mathbf{r}_k\|^2$, we get

$$f(\mathbf{x}_{k+1}) = f(\mathbf{x}_k) - \frac{\|\mathbf{r}_k\|^4}{\|\mathbf{A}\mathbf{r}_k\|^2} + \frac{1}{2} \frac{\|\mathbf{r}_k\|^4}{\|\mathbf{A}\mathbf{r}_k\|^2} = f(\mathbf{x}_k) - \frac{1}{2} \frac{\|\mathbf{r}_k\|^4}{\|\mathbf{A}\mathbf{r}_k\|^2}.$$

Since

$$\frac{\|\mathbf{A}\mathbf{r}_k\|}{\|\mathbf{r}_k\|} \leq \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \|\mathbf{A}\|,$$

we get $\|\mathbf{A}\mathbf{r}_k\| \leq \|\mathbf{A}\| \|\mathbf{r}_k\|$, and applying this to (4.4), the inequality

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{1}{2} \frac{\|\mathbf{r}_k\|^2}{\|\mathbf{A}\|^2}.$$

Rearranging the gives the claim. □

Proof of Theorem 4.2. Let

$$g_N := \min_{0 \leq k \leq N} \|\mathbf{r}_k\|.$$

Then since $(N+1)g_N^2 \leq \sum_{k=0}^N \|\mathbf{r}_k\|^2$, and using Lemma 4.3,

$$g_N \leq \sqrt{\frac{1}{N+1} \sum_{k=0}^N \|\mathbf{r}_k\|^2} \leq \frac{\sqrt{2}\|\mathbf{A}\|}{\sqrt{N+1}} \sqrt{f(\mathbf{x}_0) - f(\mathbf{x}^*)} = \frac{\sqrt{2}\|\mathbf{A}\|}{\sqrt{N+1}} \|\mathbf{x}_0 - \mathbf{x}^*\|_{\mathbf{A}}.$$

We know that $g_N < \varepsilon$ if $\sqrt{2\|\mathbf{A}\|^2}\|\mathbf{x}_0 - \mathbf{x}^*\|_{\mathbf{A}}/\sqrt{N+1} < \varepsilon$. Solving the latter for N , we get that

$$N > \frac{2\|\mathbf{A}\|^2}{\varepsilon^2}\|\mathbf{x}_0 - \mathbf{x}^*\|_{\mathbf{A}}^2 - 1.$$

will guarantee a residual smaller than ε . \square

The following theorem shows that the gradient descent method for linear least squares converges linearly with respect to the \mathbf{A} norm. The statement involves the *condition number* of \mathbf{A} ¹. This quantity is defined as

$$\kappa(\mathbf{A}) = \|\mathbf{A}\| \cdot \|\mathbf{A}^\dagger\|.$$

Theorem 4.4. *The error in the $k+1$ -th iterate is bounded by*

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \left(\frac{\kappa^2(\mathbf{A}) - 1}{\kappa^2(\mathbf{A}) + 1} \right) \|\mathbf{x}_k - \mathbf{x}^*\|.$$

In particular, the gradient descent algorithm converges linearly. This Theorem follows from Theorem 4.2 using an inequality known as Kantorovich's inequality.

¹The concept of condition number, introduced by Alan Turing while in Manchester, is one of the most important ideas in numerical analysis, as it is indispensable in studying the performance of numerical algorithms.