

## Line search methods

Given a starting point  $\mathbf{x}_k$  and a direction  $\mathbf{d}_k$ , we still need to decide on  $\alpha_k$ , i.e., how far to move. With  $\mathbf{x}_k$  and  $\mathbf{d}_k$  fixed, we can think of the remaining problem as a one-dimensional optimization problem where we would like to choose  $\alpha$  to minimize (or at least reduce)

$$\phi(\alpha) = f(\mathbf{x}_k + \alpha\mathbf{d}_k).$$

Note that we don't necessarily need to find the true minimum – we aren't even sure that we are moving in the right direction at this point – but we would generally still like to make as much progress as possible before calculating a new direction  $\mathbf{d}_{k+1}$ . There are many methods for doing this, here are three:

### Fixed step size

We can just use a constant step size  $\alpha_k = \alpha$ . This will work if the step size is small enough, but usually this results in using more iterations than necessary. This is actually a very commonly used approach since if your problem is small enough, this may not matter.

### Exact line search

Another approach is to solve the one-dimensional optimization program

$$\underset{\alpha \geq 0}{\text{minimize}} \phi(\alpha).$$

There are a variety of strategies you could take here (e.g., applying a bisection search or some similar one-dimensional optimization strategy) to try to solve this problem. This is typically not worth

the trouble. However, there are certain instances (e.g., least squares and other unconstrained convex quadratic programs) when it can be solved analytically, in which case it is generally a good idea.

**Example: Minimizing a quadratic function** Suppose we wish to solve the optimization problem

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{x}^T \mathbf{b}.$$

For example, this optimization problem arises in the context of solving least squares problems. Suppose that we have selected a step direction  $\mathbf{d}_k$ . In this case

$$\phi(\alpha) = \frac{1}{2} (\mathbf{x}_k + \alpha \mathbf{d}_k)^T \mathbf{Q} (\mathbf{x}_k + \alpha \mathbf{d}_k) - (\mathbf{x}_k + \alpha \mathbf{d}_k)^T \mathbf{b}.$$

This is a quadratic function of  $\alpha$ , and thus we can compute the optimal step size by finding the  $\alpha$  such that  $\phi'(\alpha) = 0$ . By expanding out the quadratic term, it is easy to show that

$$\phi'(\alpha) = \alpha \mathbf{d}_k^T \mathbf{Q} \mathbf{d}_k + \mathbf{d}_k^T \mathbf{Q} \mathbf{x}_k - \mathbf{d}_k^T \mathbf{b}.$$

Setting this equal to zero and solving for  $\alpha$  yields the step size

$$\alpha_k = \frac{\mathbf{d}_k^T (\mathbf{b} - \mathbf{Q} \mathbf{x}_k)}{\mathbf{d}_k^T \mathbf{Q} \mathbf{d}_k}.$$

## Backtracking

Exact line search is generally not worth the trouble, but the problem with a fixed step size is that we cannot guarantee convergence of  $\alpha$  is too large, but when  $\alpha$  is too small we may not make much progress on

each iteration. A popular strategy is to do a rudimentary search for an  $\alpha$  that gives us sufficient progress as measured by the inequality

$$f(\mathbf{x}_k) - f(\mathbf{x}_k + \alpha \mathbf{d}_k) \geq -c_1 \alpha \langle \mathbf{d}_k, \nabla f(\mathbf{x}_k) \rangle,$$

where  $c_1 \in (0, 1)$ . This is known as the **Armijo condition**. For  $\alpha$  satisfying the inequality we have that the reduction in  $f$  is proportional to both the step length  $\alpha$  and the directional derivative in the direction  $\mathbf{d}_k$ .

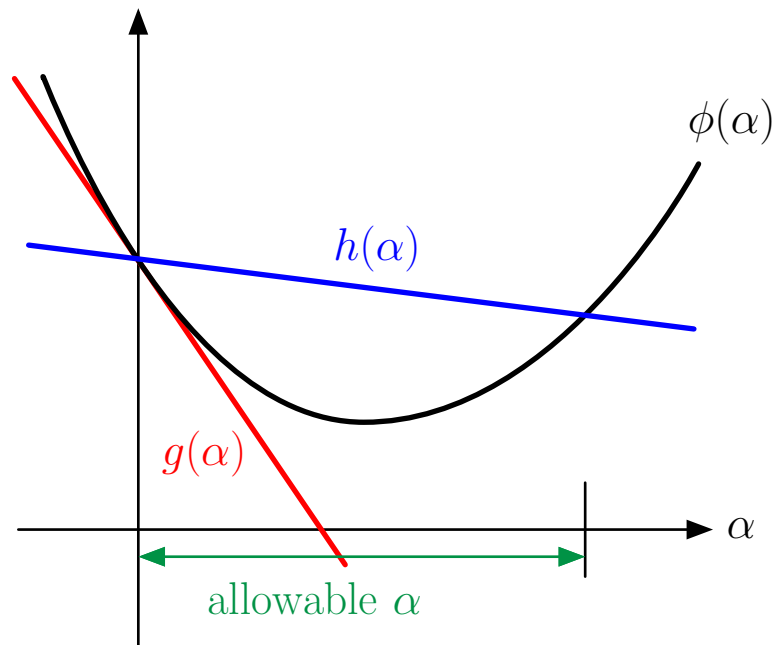
Note that we can equivalently write this condition as

$$\phi(\alpha) \leq h(\alpha) := \phi(0) + c_1 \alpha \phi'(0).$$

Recall that from convexity, we also have that

$$\phi(\alpha) \geq g(\alpha) := \phi(0) + \alpha \phi'(0).$$

Since  $c_1 < 1$ , we always have  $\phi(\alpha) \leq h(\alpha)$  for sufficiently small  $\alpha$ . An example is illustrated below:



We still haven't said anything about how to actually use the Armijo condition to pick  $\alpha$ . Within the set of allowable  $\alpha$  satisfying the condition, the (guaranteed) reduction in  $f$  is proportional to  $\alpha$ , so we would generally like to select  $\alpha$  to be large. Note that the Armijo condition protects us against setting  $\alpha$  to be *too* large, however, it actually does not rule out setting  $\alpha$  to be extremely small, which can be just as much of a problem in practice.

This inspires the following very simple **backtracking** algorithm: start with a large step size of  $\alpha = \bar{\alpha}$ , and then decrease by a factor of  $\rho$  until the Armijo condition is satisfied.

### Backtracking line search

Input:  $\mathbf{x}_k, \mathbf{d}_k, \nabla f(\mathbf{x}_k), \bar{\alpha} > 0, c_1 \in (0, 1),$  and  $\rho \in (0, 1).$

Initialize:  $\alpha = \bar{\alpha}$

**while** Armijo condition not satisfied **do**

$\alpha = \rho\alpha$

**end while**

The backtracking line search tends to be cheap, and works very well in practice. A common choice for  $\bar{\alpha}$  is  $\bar{\alpha} = 1$ , but this can vary somewhat depending on the algorithm. The choice of  $c_1$  can range from extremely small ( $10^{-4}$ , encouraging larger steps) to relatively large (0.3, encouraging smaller steps), and typical values of  $\rho$  range from 0.1, (corresponding to a relatively coarse search) to 0.8 (corresponding to a finer search).

## Convergence of gradient descent

Here we will discuss convergence guarantees for **gradient descent**, i.e., the version of our iterative algorithm where we set

$$\mathbf{d}_k = -\nabla f(\mathbf{x}_k),$$

resulting in the update rule

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k).$$

It is hard to say much about the convergence properties of this approach for *arbitrary* convex functions. However, if  $f$  satisfies certain “regularity conditions”, then we can get very nice guarantees, even for a fixed step size. Here we will look at two different kinds of regularity assumptions on  $f$ , and translate them into convergence rates. Throughout, we will assume that  $f$  is differentiable everywhere.<sup>1</sup>

### Smoothness and strong convexity

We will consider two related kinds of assumptions on  $f$ . One is that  $f$  is **smooth** in a certain sense. Qualitatively, we would just like to assume that the gradient changes in a controlled manner as we move from point to point. Quantitatively, we will assume that  $f$  has a **Lipschitz gradient**. This means that there exists an  $M > 0$  such that

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq M\|\mathbf{x} - \mathbf{y}\|_2, \quad (1)$$

---

<sup>1</sup>Methods for nondifferentiable  $f(\mathbf{x})$  are also of great interest, and will be covered later in the course. These methods are not much more involved algorithmically (although, you obviously will have to replace the gradient with something else), but they are slightly harder to analyze.

for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$ . We will say that such a function is  **$M$ -smooth** or **strongly smooth**.

One can show that  $f$  obeying (1) is actually equivalent to saying that

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{M}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \quad (2)$$

for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$ .

This provides some intuition for what kind of structure the Lipschitz gradient condition imposes on  $f$ . Recall that for any convex function, we have that

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle,$$

so if  $f$  is convex, then at any point  $\mathbf{x}$  we can bound  $f$  from *below* by a linear approximation. If  $f$  has a Lipschitz gradient, (2) but we can also bound it from *above* using a quadratic approximation.

In the case that  $f$  is twice differentiable, one can also show that (2) is equivalent to

$$\nabla^2 f(\mathbf{x}) \preceq M\mathbf{I},$$

i.e., that the largest eigenvalue of the Hessian is bounded by  $M$  for all  $\mathbf{x}$ . Note, however, that the Lipschitz gradient condition and the analysis below does not require  $f$  to be twice differentiable.

A closely related assumption that we can make is to assume that  $f$  is **strongly convex** (with strong convexity parameter  $m > 0$ ), meaning that

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{m}{2} \|\mathbf{y} - \mathbf{x}\|_2^2. \quad (3)$$

for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$ . This tells us that not only is  $f$  bounded below by a *linear* approximation (since it is convex), but also by a (nontrivial) convex *quadratic* approximation. Note also that strong convexity implies strict convexity, but strict convexity does not necessarily imply strong convexity.

In the case that  $f$  is twice differentiable, strong convexity is equivalent to

$$\nabla^2 f(\mathbf{x}) \succeq m\mathbf{I}.$$

That is, the eigenvalues of the Hessian are bounded below by  $m > 0$  for all  $\mathbf{x}$ . When combined with the assumption of  $M$ -smoothness, this bounds the conditioning of the Hessian matrix so that its eigenvalues are bounded between  $m > 0$  and  $M < \infty$ . However, again note that strong convexity does not require  $f$  to be twice differentiable.

## Convergence of gradient descent for smooth and strongly convex $f$

If you look back to when we analyzed the convergence of gradient descent for the case where  $f$  was a quadratic function (as in least squares), you may notice that our analysis centered entirely around the identity

$$f(\mathbf{y}) = f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T \mathbf{H}(\mathbf{y} - \mathbf{x}),$$

where in this case  $\mathbf{H} = \nabla^2 f(\mathbf{x})$ . In fact, the main convergence result follows from a pair of inequalities that follow directly from this identity. Specifically, that

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{\lambda_{\max}(\mathbf{H})}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$

and

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{\lambda_{\min}(\mathbf{H})}{2} \|\mathbf{y} - \mathbf{x}\|_2^2.$$

By assuming that  $f$  is  $M$ -smooth and strongly convex, we are essentially making precisely the assumptions that we need to reproduce this analysis. Another perspective is that our assumptions imply that  $f$  is not too different from a quadratic function, and so we have similar convergence guarantees. Specifically, using the same basic approach as before one can show (try this at home!) that for  $\alpha_k = 1/M$  we have

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \left(1 - \frac{m}{M}\right) (f(\mathbf{x}_k) - f(\mathbf{x}^*)).$$

This is an example of *linear convergence*, and using a simple argument from the homework this implies that

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \epsilon,$$

as long as

$$k \geq \frac{\log((f(\mathbf{x}_0) - f(\mathbf{x}^*))/\epsilon)}{\log(1/(1 - m/M))}.$$

## Convergence of gradient descent for smooth $f$

It is also possible to get weaker convergence guarantees without making the assumption that  $f$  is strongly convex. In the technical addendum at the end of these notes, we show that by combining the assumption of  $M$ -smoothness with the definition of convexity and doing some clever manipulations, we can get a guarantee of the form

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{M}{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2.$$



Thus, for  $M$ -smooth functions, we can guarantee that the error is  $O(1/k)$  after  $k$  iterations. Another way to put this is to say that we can guarantee accuracy

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \epsilon$$

as long as

$$k \geq \frac{M}{2\epsilon} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2.$$

Note that if  $\epsilon$  is very small, this says we can expect to need a very large number of iterations. This is **much** slower convergence than what we obtained by assuming strong convexity – it is  $O(\epsilon^{-1})$  versus  $O(\log \epsilon^{-1})$ . As an example, if we wanted to set  $\epsilon = 10^{-6}$ ,  $\epsilon^{-1} = 10^6$  (versus  $\log \epsilon^{-1} \approx 14$ ). Of course, to get the stronger guarantee we had to make a much stronger assumption (strong convexity), which may not always be applicable depending on the objective function you are optimizing.

## Technical Details: Convergence analysis for $M$ -smooth functions

Here we provide the convergence analysis for gradient descent on  $M$ -smooth functions that are not necessarily strongly convex. From our assumption that  $f$  is  $M$ -smooth, we know that  $f$  satisfies (2), and thus plugging in  $\mathbf{y} = \mathbf{x}_{k+1}$ , we obtain

$$\begin{aligned} f(\mathbf{x}_{k+1}) &\leq f(\mathbf{x}_k) + \left\langle -\frac{1}{M} \nabla f(\mathbf{x}_k), \nabla f(\mathbf{x}_k) \right\rangle + \frac{M}{2} \left\| \frac{1}{M} \nabla f(\mathbf{x}_k) \right\|_2^2 \\ &= f(\mathbf{x}_k) - \frac{1}{M} \|\nabla f(\mathbf{x}_k)\|_2^2 + \frac{1}{2M} \|\nabla f(\mathbf{x}_k)\|_2^2 \\ &= f(\mathbf{x}_k) - \frac{1}{2M} \|\nabla f(\mathbf{x}_k)\|_2^2. \end{aligned} \tag{4}$$

Moreover, by the convexity of  $f$ ,

$$f(\mathbf{x}_k) \leq f(\mathbf{x}^*) + \langle \mathbf{x}_k - \mathbf{x}^*, \nabla f(\mathbf{x}_k) \rangle,$$

where  $\mathbf{x}^*$  is a minimizer of  $f$ , and so we have

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}^*) + \langle \mathbf{x}_k - \mathbf{x}^*, \nabla f(\mathbf{x}_k) \rangle - \frac{1}{2M} \|\nabla f(\mathbf{x}_k)\|_2^2.$$

Substituting  $\nabla f(\mathbf{x}_k) = M(\mathbf{x}_k - \mathbf{x}_{k+1})$  then yields

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq M \langle \mathbf{x}_k - \mathbf{x}^*, \mathbf{x}_k - \mathbf{x}_{k+1} \rangle - \frac{M}{2} \|\mathbf{x}_k - \mathbf{x}_{k+1}\|_2^2. \quad (5)$$

We can re-write this in a slightly more convenient way using the fact that

$$\|\mathbf{a} - \mathbf{b}\|_2^2 = \|\mathbf{a}\|_2^2 - 2\langle \mathbf{a}, \mathbf{b} \rangle + \|\mathbf{b}\|_2^2$$

and thus

$$2\langle \mathbf{a}, \mathbf{b} \rangle - \|\mathbf{b}\|_2^2 = \|\mathbf{a}\|_2^2 - \|\mathbf{a} - \mathbf{b}\|_2^2.$$

Setting  $\mathbf{a} = \mathbf{x}_k - \mathbf{x}^*$  and  $\mathbf{b} = \mathbf{x}_k - \mathbf{x}_{k+1}$  and applying this to (5), we obtain the bound

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \frac{M}{2} (\|\mathbf{x}_k - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|_2^2).$$

This result bounds how far away  $f(\mathbf{x}_{k+1})$  is from the optimal  $f(\mathbf{x}^*)$  in terms (primarily) of the error in the previous iteration:  $\|\mathbf{x}_k - \mathbf{x}^*\|_2^2$ . We can use this result to bound  $f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*)$  in terms of the initial error  $\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$  by a clever argument.

Specifically, this bound holds not only for iteration  $k$ , but for all iterations  $i = 1, \dots, k$ , so we can write down  $k$  inequalities and then sum them up to obtain

$$\sum_{i=1}^k f(\mathbf{x}_i) - f(\mathbf{x}^*) \leq \frac{M}{2} \left( \sum_{i=1}^k \|\mathbf{x}_{i-1} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_i - \mathbf{x}^*\|_2^2 \right).$$

The right-hand side of this inequality is what is called a *telescopic sum*: each successive term in the sum cancels out part of the previous term. Once you write this out, all the terms cancel except for two (one component from the  $i = 1$  term and one from the  $i = k$  term) giving us:

$$\begin{aligned} \sum_{i=1}^k f(\mathbf{x}_i) - f(\mathbf{x}^*) &\leq \frac{M}{2} (\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_k - \mathbf{x}^*\|_2^2) \\ &\leq \frac{M}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2. \end{aligned}$$

Since, as noted above,  $f(\mathbf{x}_i)$  is monotonically decreasing in  $i$ , we also have that

$$k (f(\mathbf{x}_k) - f(\mathbf{x}^*)) \leq \sum_{i=1}^k f(\mathbf{x}_i) - f(\mathbf{x}^*),$$

and thus

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{M}{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2,$$

which is exactly what we wanted to show.