

# III. Unconstrained Optimization of Smooth Functions

# Unconstrained optimization

We will start our discussion about solving convex optimization programs by considering the unconstrained case. Our template problem is

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x}), \tag{1}$$

where  $f$  is convex. While we state this problem as a search over all of  $\mathbb{R}^N$ , almost everything we say here can be applied to minimizing a convex function over an *open* set<sup>1</sup>.

We will first discuss two fundamental results. First, for any convex  $f$ , we will give conditions under which a minimizer to (1) exists, and show that if  $\mathbf{x}^*$  is a local minimizer of (1), then it is also a global minimizer. Second, under the conditions that  $f(\mathbf{x})$  is convex and differentiable, we will show that  $\mathbf{x}^*$  is a minimizer of (1) if and only if the derivative is equal to zero:

$$\mathbf{x}^* \text{ is a global minimizer} \iff \nabla f(\mathbf{x}^*) = \mathbf{0}.$$

Something similar to this gradient condition is also true for non-differentiable (but still convex)  $f$ ; we will explore this a little later in the course.

---

<sup>1</sup>Intuition: Open sets don't have boundaries, closed sets do. The entire point of the study of constrained optimization, which we will get to next, comes down to treating the fact that the solution can (and probably is) on the boundary of your constraint set.

## Existence of minimizers

We begin by recalling a fundamental result in analysis: the *Weierstrass extreme value theorem*. If  $f(\mathbf{x})$  is a continuous functional on a compact<sup>2</sup> set  $\mathcal{K} \subset \mathbb{R}^N$ , then it attains its minimum value at least once. That is,

$$\underset{\mathbf{x} \in \mathcal{K}}{\text{minimize}} \quad f(\mathbf{x})$$

has a minimizer in  $\mathcal{K}$  — there exists an  $\mathbf{x}^* \in \mathcal{K}$  such that  $f(\mathbf{x}^*) \leq f(\mathbf{x})$  for all  $\mathbf{x} \in \mathcal{K}$ . For a proof of this, see just about any introductory text on analysis. The same is also true for  $f$  achieving its maximum value in  $\mathcal{K}$ .

In the unconstrained setting, we are interested in

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x}),$$

where  $f$  is convex. Simple examples illustrate that the minimum does not necessarily have to be achieved for any  $\mathbf{x}$ . That is, there is no  $\mathbf{x}^*$  such that  $f(\mathbf{x}^*) \leq f(\mathbf{x})$  for all  $\mathbf{x} \in \mathbb{R}^N$ . For example,  $f(x) = e^{-x}$  does not have a minimizer on the real line, even though it is as convex and as smooth as can be.

There is, however, a class of functions for which we can guarantee a global minimizer in the unconstrained setting. If the sublevel sets of  $f$ ,

$$\mathcal{S}_\alpha(f) = \{\mathbf{x} \in \mathbb{R}^N : f(\mathbf{x}) \leq \alpha\}$$

are *compact* (again, this means closed and bounded), then there will be at least one global minimizer. This should be easy to see — just

---

<sup>2</sup>In  $\mathbb{R}^N$ , saying a set is compact is the same as saying that it is closed and bounded.

choose  $\alpha$  such that  $\mathcal{S}_\alpha(f)$  is non-empty, then

$$\underset{\mathbf{x} \in \mathcal{S}_\alpha(f)}{\text{minimize}} \quad f(\mathbf{x})$$

has a minimizer (by the extreme value theorem), and this also clearly corresponds to a minimizer of  $f$  over  $\mathbb{R}^N$ . If  $f$  is continuous (which all convex functions with  $\text{dom } f = \mathbb{R}^N$  are), then having compact sublevel sets is the same as being *coercive*: for every sequence  $\{\mathbf{x}_k\} \subset \mathbb{R}^N$  with  $\|\mathbf{x}_k\|_2 \rightarrow \infty$ , we have  $f(\mathbf{x}_k) \rightarrow \infty$  as well. (I will let you prove that at home.)

## Local minima are also global minima

Let's nail down a precise statement here right from the start:

Let  $f(\mathbf{x})$  be convex function on  $\mathbb{R}^N$ , and suppose  $\mathbf{x}^*$  is a local minimizer of  $f$  in that there exists an  $\epsilon > 0$  such that

$$f(\mathbf{x}^*) \leq f(\mathbf{x}) \quad \text{for all } \|\mathbf{x} - \mathbf{x}^*\|_2 \leq \epsilon.$$

Then  $\mathbf{x}^*$  is also a global minimizer:  $f(\mathbf{x}^*) \leq f(\mathbf{x})$  for all  $\mathbf{x} \in \mathbb{R}^N$ .

To prove this, suppose that there was an  $\hat{\mathbf{x}} \neq \mathbf{x}^*$  such that  $f(\hat{\mathbf{x}}) < f(\mathbf{x}^*)$ . Then by the convexity of  $f$ ,

$$\begin{aligned} f(\mathbf{x}^* + \theta(\hat{\mathbf{x}} - \mathbf{x}^*)) &\leq (1 - \theta)f(\mathbf{x}^*) + \theta f(\hat{\mathbf{x}}) \\ &< f(\mathbf{x}^*) \quad \text{for all } 0 \leq \theta \leq 1. \end{aligned}$$

But choosing a small enough value of  $\theta$  puts  $\mathbf{x}^* + \theta(\hat{\mathbf{x}} - \mathbf{x}^*)$  in the neighborhood where  $\mathbf{x}^*$  is supposed to be a local min. Specifically,

if we take  $\theta < \epsilon / \|\hat{\mathbf{x}} - \mathbf{x}^*\|_2$ , then the inequality above directly contradicts the assertion that  $\mathbf{x}^*$  is a local minimizer. Thus no such  $\hat{\mathbf{x}}$  can exist.

Our final result in this section gives a sufficient (but definitely not necessary) condition for the minimizer to be unique.

Let  $f$  be strictly convex on  $\mathbb{R}^N$ . If  $f$  has a global minimizer, then it is unique.

This is again easy to argue by contradiction. Let  $\mathbf{x}^*$  be a global minimizer, and suppose that there existed an  $\hat{\mathbf{x}} \neq \mathbf{x}^*$  with  $f(\hat{\mathbf{x}}) = f(\mathbf{x}^*)$ . But then there would be many  $\mathbf{x}$  which achieve smaller values, as for all  $0 < \theta < 1$ ,

$$\begin{aligned} f(\theta \mathbf{x}^* + (1 - \theta) \hat{\mathbf{x}}) &< \theta f(\mathbf{x}^*) + (1 - \theta) f(\hat{\mathbf{x}}) \\ &= f(\mathbf{x}^*). \end{aligned}$$

As this would contradict the assertion that  $\mathbf{x}^*$  is a global minimizer, no such  $\hat{\mathbf{x}}$  can exist.

We close this section by re-emphasizing that the entire discussion above would stay the same if we replaced  $\text{minimize}_{\mathbf{x} \in \mathbb{R}^N} f(\mathbf{x})$  with  $\text{minimize}_{\mathbf{x} \in \mathcal{U}} f(\mathbf{x})$  for any open set  $\mathcal{U} \subset \mathbb{R}^N$ .

# Optimality conditions for unconstrained optimization

How do we know when we have a minimizer of a convex function on our hands? What is our “certificate of optimality”? For the time being, we will assume that  $f$  is differentiable (that  $\nabla f(\mathbf{x})$  exists at every point we are considering). There are a comparable set of results for non-smooth  $f$  that we will discuss in last segment of the course.

In this course we have already mentioned the following, but have never actually discussed exactly why it is true.

Let  $f$  be a convex differentiable function on  $\mathbb{R}^N$ . Then  $\mathbf{x}^*$  solves

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x})$$

if and only if  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ .

The proof of this relies on the critical fact that we can decrease  $f$  by moving in a direction which has an obtuse angle with the gradient.

Let  $f$  be a function on  $\mathbb{R}^N$  that is differentiable at  $\mathbf{x}$ , and let  $\mathbf{d} \in \mathbb{R}^N$  be a vector obeying  $\langle \mathbf{d}, \nabla f(\mathbf{x}) \rangle < 0$ . Then for small enough  $t > 0$ ,

$$f(\mathbf{x} + t\mathbf{d}) < f(\mathbf{x}).$$

We call such a  $\mathbf{d}$  a **descent direction** from  $\mathbf{x}$ . Similarly, if  $\langle \mathbf{d}, \nabla f(\mathbf{x}) \rangle > 0$ , then for small enough  $t > 0$ ,

$$f(\mathbf{x} + t\mathbf{d}) > f(\mathbf{x}).$$

We call such a  $\mathbf{d}$  an **ascent direction** from  $\mathbf{x}$ .

This fundamental fact is a direct consequence of the Taylor theorem (see the Technical Details section below): for any  $\mathbf{u} \in \mathbb{R}^N$ ,

$$f(\mathbf{x} + \mathbf{u}) = f(\mathbf{x}) + \langle \mathbf{u}, \nabla f(\mathbf{x}) \rangle + h(\mathbf{u}) \|\mathbf{u}\|_2,$$

where  $h(\mathbf{u}) : \mathbb{R}^N \rightarrow \mathbb{R}$  is some function satisfying  $h(\mathbf{u}) \rightarrow 0$  as  $\mathbf{u} \rightarrow \mathbf{0}$ . Taking  $\mathbf{u} = t\mathbf{d}$ , we have

$$f(\mathbf{x} + t\mathbf{d}) = f(\mathbf{x}) + t(\langle \mathbf{d}, \nabla f(\mathbf{x}) \rangle + h(t\mathbf{d})\|\mathbf{d}\|_2).$$

For  $t > 0$  small enough, we can make  $|h(t\mathbf{d})| \cdot \|\mathbf{d}\|_2 < |\langle \mathbf{d}, \nabla f(\mathbf{x}) \rangle|$ , and so the term inside the parentheses above is negative if  $\langle \mathbf{d}, \nabla f(\mathbf{x}) \rangle$  is negative, and it is positive if  $\langle \mathbf{d}, \nabla f(\mathbf{x}) \rangle$  is positive.

At a particular point  $\mathbf{x}^*$ , the only way we can make  $\langle \mathbf{d}, \nabla f(\mathbf{x}^*) \rangle \geq 0$  for all choices of  $\mathbf{d}$  is if  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ . So clearly

$$\mathbf{x}^* \text{ is a minimizer} \quad \Rightarrow \quad \nabla f(\mathbf{x}^*) = \mathbf{0}.$$

On the other hand, if  $f$  is convex, then

$$f(\mathbf{x}^* + t\mathbf{d}) \geq f(\mathbf{x}^*) + t\langle \mathbf{d}, \nabla f(\mathbf{x}^*) \rangle,$$

for all  $t \in \mathbb{R}$  and choices of  $\mathbf{d} \in \mathbb{R}^N$ . This now makes it clear that

$$\nabla f(\mathbf{x}^*) = \mathbf{0} \quad \Rightarrow \quad \mathbf{x}^* \text{ is a minimizer.}$$

Again, for everything we have said in this section, you can use any open domain  $\mathcal{U}$  in place of  $\mathbb{R}^N$ .

# Algorithms for unconstrained minimization

We want to solve

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x}),$$

where  $f$  is convex. In this section, in addition to assuming that  $f$  is differentiable (so its gradient exists at every point), we will also assume that it is **smooth** (we will consider a few different definitions of “smooth” — qualitatively, this just means that the gradient changes in a controlled manner as we move from point to point).<sup>3</sup>

Since  $f$  is convex, a necessary and sufficient condition for  $\mathbf{x}^*$  to be a minimizer is that the gradient vanishes:

$$\nabla f(\mathbf{x}^*) = \mathbf{0}.$$

As noted above, it is not a given that such a  $\mathbf{x}^*$  exists, but in this section we will assume that  $f$  does have (at least one) minimizer, and we will denote the optimal value as  $p^* = f(\mathbf{x}^*)$ .

Every general-purpose optimization algorithm we will look at in this course is **iterative** — they will all have the basic form:

$k = 0$ ,  $\mathbf{x}^{(0)}$  = initial guess

while not converged

$k = k + 1$

    calculate a direction to move  $\mathbf{d}^{(k)}$

    calculate a step size  $t_k \geq 0$

$\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} + t_k \mathbf{d}^{(k)}$

    check convergence criteria

---

<sup>3</sup>While many problems are smooth, methods for nonsmooth  $f(\mathbf{x})$  are also of great interest, and will be covered later in the course. Nonsmooth methods are not much more involved algorithmically, but they are slightly harder to analyze.



In the coming lectures, we will focus primarily on two methods for computing the direction  $\mathbf{d}^{(k)}$ .

1. **Gradient descent**: we take

$$\mathbf{d}^{(k)} = -\nabla f(\mathbf{x}^{(k-1)}).$$

This is the direction of “steepest descent” (where “steepest” is defined relative to the Euclidean norm). Gradient descent iterations are **cheap**, but typically many iterations are required for convergence.

2. **Newton’s method**: we build a quadratic model around  $\mathbf{x}^{(k)}$  then compute the exact minimizer of this quadratic by solving a system of equations. This corresponds to taking

$$\mathbf{d}^{(k)} = -(\nabla^2 f(\mathbf{x}^{(k-1)}))^{-1} \nabla f(\mathbf{x}^{(k-1)}),$$

that is, the inverse of the Hessian evaluated at  $\mathbf{x}^{(k-1)}$  applied to the gradient evaluated at the same point. Newton iterations tend to be **expensive** (as they require a system solve), but they typically converge in far fewer iterations than gradient descent.

Whichever direction we choose, it should be a **descent direction**;  $\mathbf{d}^{(k)}$  should satisfy

$$\langle \mathbf{d}^{(k)}, \nabla f(\mathbf{x}^{(k-1)}) \rangle \leq 0.$$

Since  $f$  is convex, it is always true that

$$f(\mathbf{x} + t\mathbf{d}) \geq f(\mathbf{x}) + t\langle \mathbf{d}, \nabla f(\mathbf{x}) \rangle,$$

and so to decrease the value of the functional while moving in direction  $\mathbf{d}$ , it is necessary that the inner product above be negative.

## Line search

After the step direction is chosen, we need to compute how far to move. There are many methods for doing this, here are three:

**Exact:** Solve the 1D optimization program

$$\underset{s \geq 0}{\text{minimize}} \quad f(\mathbf{x}^{(k-1)} + s\mathbf{d}^{(k)}).$$

This is typically not worth the trouble, but there are instances (i.e. unconstrained convex quadratic programs) when it can be solved analytically.

**Backtracking:** Start with a step size of  $t = 1$ , then decrease by a factor of  $\beta$  until we are below a certain line.

Fix  $\alpha \in (0, 0.5)$  and  $\beta \in (0, 1)$ .

Given a starting point  $\mathbf{x}$  and direction  $\mathbf{d}$ ,

$t = 1$

repeat

    if  $f(\mathbf{x} + t\mathbf{d}) < f(\mathbf{x}) + \alpha t \langle \mathbf{d}, \nabla f(\mathbf{x}) \rangle$

        converged

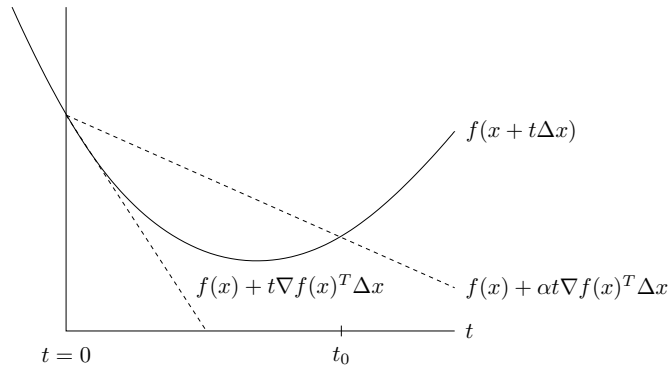
    else

$t = \beta t$

    endif

until converged

Here is a figure (from [BV04, Chapter 9]):



**Figure 9.1** *Backtracking line search.* The curve shows  $f$ , restricted to the line over which we search. The lower dashed line shows the linear extrapolation of  $f$ , and the upper dashed line has a slope a factor of  $\alpha$  smaller. The backtracking condition is that  $f$  lies below the upper dashed line, *i.e.*,  $0 \leq t \leq t_0$ .

The backtracking line search tends to be cheap, and works very well in practice. Typically, we take  $\alpha$  small (to encourage a large step) and  $\beta \in [0.3, 0.8]$ .

**Fixed:** Finally, we can just use a constant step size  $t_k = t$ . This will actually work if the step size is small enough, but usually this results in way too many iterations.

## Convergence of gradient descent

How quickly does gradient descent converge?

I'm glad you asked!

We will look at two different smoothness assumptions on  $f$ , and translate them into convergence rates. In the first case, we assume that  $f$  is twice differentiable (so that its Hessian exists everywhere), and that

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

That is, the eigenvalues of the Hessian (which is always in  $S^N$  for a convex function) are bounded between  $m > 0$  and  $M < \infty$ . We call this assumption **strong convexity** (note that the lower bound by itself means that  $f$  is strictly convex).

Recall that the main consequence of convexity is that we have a way to compute a linear global lower bound at every point:

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle, \quad \text{for all } \mathbf{x}, \mathbf{y}.$$

The main consequence of strong convexity is that we have in addition a quadratic lower and upper bound. By the Taylor theorem, we know that for any  $\mathbf{x}, \mathbf{y}$ ,

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

for some point  $\mathbf{z}$  on the line segment between  $\mathbf{x}$  and  $\mathbf{y}$ . Thus we have the bounds

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

and

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

An immediate consequence of the stronger lower bound in (2) is that we can tell at any point  $\mathbf{x}$  how close to optimal we are. The smallest value the right hand side can take in that inequality is when  $\tilde{\mathbf{y}} = \mathbf{x} - m^{-1} \nabla f(\mathbf{x})$ ; plugging that value into the right hand side yields

$$f(\mathbf{y}) \geq f(\mathbf{x}) - \frac{1}{2m} \|\nabla f(\mathbf{x})\|_2^2,$$

a bound which holds for every  $f(\mathbf{y})$ . In particular, it holds for the optimal value  $p^*$ , and so

$$p^* \geq f(\mathbf{x}) - \frac{1}{2m} \|\nabla f(\mathbf{x})\|_2^2. \quad (4)$$

Thus we get a bound on  $f(\mathbf{x}) - p^*$  just by calculating the norm of the gradient. It also re-iterates that if  $\|\nabla f(\mathbf{x})\|_2$  is small, we are close to the solution.

We can also bound how close  $\mathbf{x}$  is to the minimizer  $\mathbf{x}^*$ . Using (2) with  $\mathbf{y} = \mathbf{x}^*$ ,

$$\begin{aligned} p^* = f(\mathbf{x}^*) &\geq f(\mathbf{x}) + \langle \mathbf{x}^* - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{m}{2} \|\mathbf{x}^* - \mathbf{x}\|_2^2 \\ &\geq f(\mathbf{x}) - \|\mathbf{x}^* - \mathbf{x}\|_2 \|\nabla f(\mathbf{x})\|_2 + \frac{m}{2} \|\mathbf{x}^* - \mathbf{x}\|_2^2. \end{aligned}$$

Then since  $p^* \leq f(\mathbf{x})$ ,

$$-\|\mathbf{x}^* - \mathbf{x}\|_2 \|\nabla f(\mathbf{x})\|_2 + \frac{m}{2} \|\mathbf{x}^* - \mathbf{x}\|_2^2 \leq 0,$$

and so

$$\|\mathbf{x}^* - \mathbf{x}\|_2 \leq \frac{2}{m} \|\nabla f(\mathbf{x})\|_2.$$

These facts are very useful for defining stopping criteria.

Suppose we have a strongly convex function that we minimize using gradient descent with an **exact line search**. We can show that at each iteration, the gap  $f(\mathbf{x}^{(k)}) - p^*$  gets cut down by a fixed factor. We consider a single iteration — we will use  $\mathbf{x}$  to denote the current

point, and  $\mathbf{x}^+ = \mathbf{x} - t_{\text{exact}} \nabla f(\mathbf{x})$  to denote the result of the gradient step. We choose  $t_{\text{exact}}$  by minimizing the following function:

$$\tilde{f}(t) = f(\mathbf{x} - t \nabla f(\mathbf{x})).$$

By strong convexity, we know that

$$\tilde{f}(t) \leq f(\mathbf{x}) - t \|\nabla f(\mathbf{x})\|_2^2 + \frac{Mt^2}{2} \|\nabla f(\mathbf{x})\|_2^2.$$

By definition of  $t_{\text{exact}}$ , we know

$$f(\mathbf{x}^+) = \tilde{f}(t_{\text{exact}}) \leq \tilde{f}(1/M) \leq f(\mathbf{x}) - \frac{1}{2M} \|\nabla f(\mathbf{x})\|_2^2.$$

From (4), we also know that

$$\|\nabla f(\mathbf{x})\|_2^2 \geq 2m(f(\mathbf{x}) - p^*),$$

and so

$$f(\mathbf{x}^+) - p^* \leq f(\mathbf{x}) - p^* - \frac{m}{M}(f(\mathbf{x}) - p^*),$$

which means

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

That is, the gap between the current functional evaluation and the optimal value has been cut down by a factor of  $1 - m/M < 1$ .

Applying this relationship recursively, we see that after  $k$  iterations of gradient descent, we have

$$\frac{f(\mathbf{x}^{(k)}) - p^*}{f(\mathbf{x}^{(0)}) - p^*} \leq \left(1 - \frac{m}{M}\right)^k.$$

Another way to say this is that we can achieve accuracy

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

by taking

$$k \geq \frac{\log(E_0/\epsilon)}{\log(1 - m/M)}, \quad E_0 = f(\mathbf{x}^{(0)}) - p^*,$$

steps.

There are similar results for gradient descent on strongly convex functions using backtracking. They are similar in nature; they have the same linear convergence but with constants that depend on  $\alpha$  and  $\beta$  along with  $m$  and  $M$ . See [BV04, p. 468].

## Lipschitz gradient condition

We can also get (much weaker) convergence results when  $f$  is not strongly convex (or even necessarily twice differentiable), but rather has a **Lipschitz gradient**. This means that there exists an  $L > 0$  such that

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2. \quad (5)$$

The upshot here is that we still have a quadratic upper bound for  $f$  at every point, but not the lower bound.

To be precise, if  $f$  obeys (5), then

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

This follows immediately from the fundamental theorem of calculus<sup>4</sup>:

$$f(\mathbf{y}) - f(\mathbf{x}) = \int_0^1 \langle \mathbf{y} - \mathbf{x}, \nabla f((1-t)\mathbf{x} + t\mathbf{y}) \rangle dt,$$

and so

$$\begin{aligned} f(\mathbf{y}) - f(\mathbf{x}) - \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle &= \int_0^1 \langle \mathbf{y} - \mathbf{x}, \nabla f((1-t)\mathbf{x} + t\mathbf{y}) - \nabla f(\mathbf{x}) \rangle dt \\ &\leq \|\mathbf{y} - \mathbf{x}\|_2 \int_0^1 \|\nabla f((1-t)\mathbf{x} + t\mathbf{y}) - \nabla f(\mathbf{x})\|_2 dt \\ &\leq L \|\mathbf{y} - \mathbf{x}\|_2^2 \int_0^1 t dt \\ &= \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2. \end{aligned}$$

Now, let's consider running gradient descent on such a function with a **fixed step size**  $t \leq 1/L$ . As before, we denote the current iterate as  $\mathbf{x}$  and the next iterate at  $\mathbf{x}^+ = \mathbf{x} - t\nabla f(\mathbf{x})$ . We have

$$\begin{aligned} f(\mathbf{x}^+) &\leq f(\mathbf{x}) - \left(1 - \frac{Lt}{2}\right) t \|\nabla f(\mathbf{x})\|_2^2 \\ &\leq f(\mathbf{x}) - \frac{t}{2} \|\nabla f(\mathbf{x})\|_2^2, \end{aligned}$$

for the range of  $t$  we are considering. By convexity of  $f$ ,

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

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

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

---

<sup>4</sup>The 1D function  $\tilde{f}(t) = f((1-t)\mathbf{x} + t\mathbf{y})$  has derivative  $\tilde{f}'(t) = \langle \mathbf{y} - \mathbf{x}, \nabla f((1-t)\mathbf{x} + t\mathbf{y}) \rangle$ . This is just a simple application of the chain rule.



and then substituting  $\nabla f(\mathbf{x}) = (\mathbf{x} - \mathbf{x}^+)/t$  yields

$$\begin{aligned} f(\mathbf{x}^+) - f(\mathbf{x}^*) &\leq \frac{1}{t} \langle \mathbf{x} - \mathbf{x}^*, \mathbf{x} - \mathbf{x}^+ \rangle - \frac{1}{2t} \|\mathbf{x} - \mathbf{x}^+\|_2^2 \\ &= \frac{1}{2t} (\|\mathbf{x} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^+ - \mathbf{x}^*\|_2^2) . \end{aligned}$$

Summing this difference over  $k$  iterations yields:

$$\begin{aligned} \sum_{i=1}^k f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) &\leq \frac{1}{2t} \left( \sum_{i=1}^k \|\mathbf{x}^{(i-1)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 \right) \\ &= \frac{1}{2t} (\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2^2) \\ &\leq \frac{1}{2t} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 . \end{aligned}$$

Since  $f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*)$  is monotonically decreasing in  $i$ , the  $k^{\text{th}}$  term will be smaller than the average:

$$\begin{aligned} f(\mathbf{x}^{(k)}) - f(\mathbf{x}^*) &\leq \frac{1}{k} \sum_{i=1}^k f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \\ &\leq \frac{1}{2tk} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 . \end{aligned}$$

Note that this convergence guarantee is much slower — it is  $O(1/k)$  in place of  $O(c^k)$  for some  $c < 1$ . This is the price we pay for allowing  $f$  to not be as smooth.

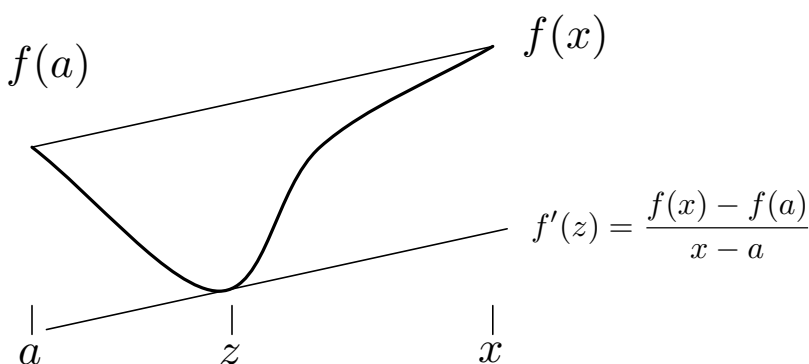
## Technical Details: Taylor's Theorem

The following is an overview of classical results from real analysis. For details and proofs, see just about any text on this subject, including [Rud76, Apo74, Rud86].

You might recall the mean-value theorem from your first calculus class. If  $f : \mathbb{R} \rightarrow \mathbb{R}$  is a differentiable function on the interval  $[a, x]$ , then there is a point inside this interval where the derivative of  $f$  matches the line drawn between  $f(a)$  and  $f(x)$ . More precisely, there exists a  $z \in [a, x]$  such that

$$f'(z) = \frac{f(x) - f(a)}{x - a}.$$

Here is a picture:



We can re-arrange the expression above to say that there is some  $z$  between  $a$  and  $x$  such that

$$f(x) = f(a) + f'(z)(x - a).$$

The mean-value theorem extends to derivatives of higher order; in this case it is known as *Taylor's theorem*. For example, suppose

that  $f$  is twice differentiable on  $[a, x]$ , and that the first derivative  $f'$  is continuous. Then there exists a  $z$  between  $a$  and  $x$  such that

$$f(x) = f(a) + f'(a)(x - a) + \frac{f''(z)}{2}(x - a)^2.$$

In general, if  $f$  is  $k + 1$  times differentiable, and the first  $k$  derivatives are continuous, then there is a point  $z$  between  $a$  and  $x$  such that

$$f(x) = p_{k,a}(x) + \frac{f^{(k+1)}(z)}{k!}(x - a)^{k+1},$$

where  $p_{k,a}(x)$  polynomial formed from the first  $k$  terms of the Taylor series expansion around  $a$ :

$$p_{k,a}(x) = f(a) + f'(a)(x - a) + \frac{f''(a)}{2}(x - a)^2 + \cdots + \frac{f^{(k)}(a)}{k!}(x - a)^k.$$

These results give us a way to quantify the accuracy of the Taylor approximation around a point. For example, if  $f$  is twice differentiable with  $f'$  continuous, then

$$f(x) = f(a) + f'(a)(x - a) + h_1(x)(x - a),$$

for a function  $h_1(x)$  goes to zero as  $x$  goes to  $a$ :

$$\lim_{x \rightarrow a} h_1(x) = 0.$$

In fact, you do not even need two derivatives for this to be true. If  $f$  has a single derivative, then we can find such an  $h_1$ . When  $f$  has two derivatives, then we have an explicit form for  $h_1$ :

$$h_1(x) = \frac{f''(z_x)}{2}(x - a),$$

where  $z_x$  is the point returned by the (generalization of) the mean value theorem for a given  $x$ .

In general, if  $f$  has  $k$  derivatives, then there exists an  $h_k(x)$  with  $\lim_{x \rightarrow a} h_k(x) = 0$  such that

$$f(x) = p_{k,a}(x) + h_k(x)(x - a)^k.$$

All of the results above extend to functions of multiple variables. For example, if  $f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}$  is differentiable, then around any point  $\mathbf{a}$ ,

$$f(\mathbf{x}) = f(\mathbf{a}) + \langle \nabla f(\mathbf{a}), \mathbf{x} - \mathbf{a} \rangle + h_1(\mathbf{x}) \|\mathbf{x} - \mathbf{a}\|_2,$$

where  $h_1(\mathbf{x}) \rightarrow 0$  as  $\mathbf{x}$  approaches  $\mathbf{a}$  from any direction. If  $f(\mathbf{x})$  is twice differentiable and the first derivative is continuous, then there exists  $\mathbf{z}$  on the line between  $\mathbf{a}$  and  $\mathbf{x}$  such that

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

We will use these two particular multidimensional results repeatedly in our analysis throughout the course, referring to them generically as “Taylor’s theorem”.

## References

- [Apo74] T. M. Apostol. *Mathematical Analysis*. Pearson, 2nd edition, 1974.
- [BV04] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.

- [Rud76] W. Rudin. *Principles of Mathematical Analysis*. McGraw-Hill, 1976.
- [Rud86] W. Rudin. *Real and Complex Analysis*. McGraw-Hill, 3rd edition, 1986.