

# Chapter 2

## Gradient Descent

### Contents

---

2.1	The algorithm	31
2.2	Vanilla analysis	31
2.3	Bounded gradients: $\mathcal{O}(1/\varepsilon^2)$ steps	33
2.4	Smoothness: $\mathcal{O}(1/\varepsilon)$ steps	34
2.5	Interlude	38
2.6	Strong convexity: $\mathcal{O}(\log(1/\varepsilon))$ steps	39
2.7	Exercises	42

---

## 2.1 The algorithm

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable. We also assume that  $f$  has a global minimum  $\mathbf{x}^*$ , and the goal is to find (an approximation of) it. This usually means that for a given  $\varepsilon > 0$ , we want to find  $\mathbf{x} \in \mathbb{R}^d$  such that

$$f(\mathbf{x}) - f(\mathbf{x}^*) < \varepsilon.$$

In this, we are not making an attempt to get near to  $\mathbf{x}^*$  itself — note that there can be several minima  $\mathbf{x}_1^* \neq \mathbf{x}_2^*$  with  $f(\mathbf{x}_1^*) = f(\mathbf{x}_2^*)$ .

Gradient descent is a very simple iterative algorithm for finding the desired approximation  $\mathbf{x}$ , under suitable conditions that we will get to. Gradient descent computes a sequence  $\mathbf{x}_0, \mathbf{x}_1, \dots$  of vectors such that  $\mathbf{x}_0$  is arbitrary, and

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma \nabla f(\mathbf{x}_t), \quad t \geq 0. \quad (2.1)$$

Here,  $\gamma$  is a fixed *stepsize*, but it may also make sense to have  $\gamma$  depend on  $t$ . For now,  $\gamma$  is fixed. As the vector  $-\nabla f(\mathbf{x}_t)$  points into a direction of descent of  $f$  at  $\mathbf{x}_t$ , the idea is to move a little bit into this direction and then iterate. We hope that after not too many iterations  $t$ ,  $f(\mathbf{x}_t) - f(\mathbf{x}^*) < \varepsilon$ ; see Figure 2.1 for an example.

The choice of  $\gamma$  is critical for the performance. If  $\gamma$  is too small, the process might take too long, and if  $\gamma$  is too large, we are in danger of overshooting. It is not clear at this point whether there is a “right” stepsize.

## 2.2 Vanilla analysis

Let  $\mathbf{x}_t$  be some iterate in the sequence (2.1). We do have an inequality that bounds  $f(\mathbf{x}_t) - f(\mathbf{x}^*)$ , namely the one saying that the graph of  $f$  lies above all its tangent hyperplanes; indeed, applying (1.2) with  $\mathbf{x} = \mathbf{x}_t$ ,  $\mathbf{y} = \mathbf{x}^*$  and reshuffling terms, we obtain

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*). \quad (2.2)$$

By definition of gradient descent (2.1),  $\nabla f(\mathbf{x}_t) = (\mathbf{x}_t - \mathbf{x}_{t+1})/\gamma$ , hence

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{1}{\gamma} (\mathbf{x}_t - \mathbf{x}_{t+1})^\top (\mathbf{x}_t - \mathbf{x}^*). \quad (2.3)$$

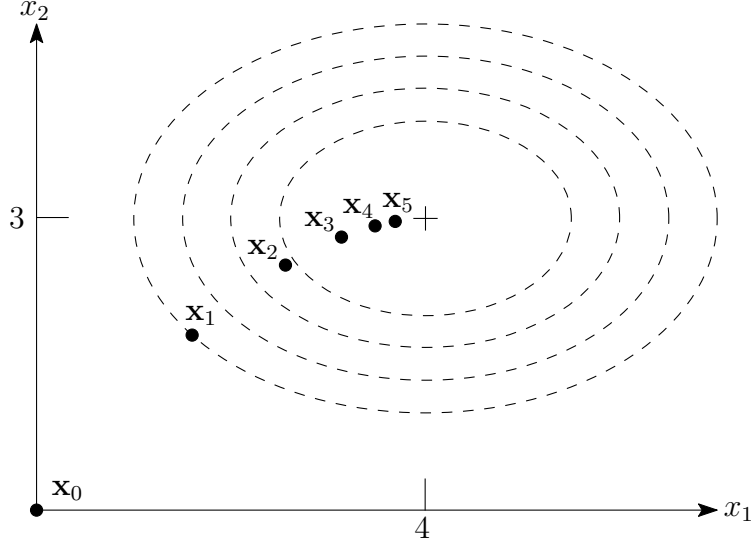


Figure 2.1: Example run of gradient descent on the quadratic function  $f(x_1, x_2) = 2(x_1 - 4)^2 + 3(x_2 - 3)^2$  with global minimum  $(4, 3)$ ; we have chosen  $\mathbf{x}_0 = (0, 0)$ ,  $\gamma = 0.1$ ; dashed lines represent level sets of  $f$  (points of constant  $f$ -value)

Now we apply (somewhat out of the blue, but this will clear up in the next step) the basic vector equation  $2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2$  to obtain

$$\begin{aligned} f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) \\ &= \frac{1}{2\gamma} (\gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) \end{aligned} \quad (2.4)$$

again by using the definition (2.1) of gradient descent. Next we sum this up over some initial values of  $t$ , so that the latter two terms in the bracket cancel in a telescoping sum.

$$\begin{aligned} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_T - \mathbf{x}^*\|^2) \\ &\leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \end{aligned} \quad (2.5)$$

This gives us an upper bound for the *average* error  $f(\mathbf{x}_t) - f(\mathbf{x}^*)$ ,  $t = 0, \dots, T - 1$ , hence in particular for the error incurred by the iterate with the smallest function value. The last iterate is not necessarily the best one: gradient descent with fixed stepsize  $\gamma$  will in general also make steps that overshoot and actually increase the function value; see Exercise 12(i).

The question is of course: is this bound any good? In general, the answer is no. A dependence on  $\|\mathbf{x}_0 - \mathbf{x}^*\|$  is to be expected (the further we start from  $\mathbf{x}^*$ , the longer we will take); the dependence on the squared gradients is more of an issue, and if we cannot control them, we cannot say much.

## 2.3 Bounded gradients: $\mathcal{O}(1/\varepsilon^2)$ steps

Here is the cheapest “solution” to squeeze something out of the vanilla analysis: let us simply assume that all gradients of  $f$  are bounded in norm. This rules out many interesting functions, though, since functions with bounded gradients only have at most linear growth. Equivalently, such functions are Lipschitz continuous over  $\mathbb{R}^d$ . But for example,  $f(x) = x^2$  (a supermodel in the world of convex functions) already doesn’t qualify, as  $\nabla f(x) = 2x$ —and this is unbounded as  $x$  tends to infinity. But let’s care about supermodels later.

**Theorem 2.1.** *Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable with a global minimum  $\mathbf{x}^*$ ; furthermore, suppose that  $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$  and  $\|\nabla f(\mathbf{x})\| \leq L$  for all  $\mathbf{x}$ . Choosing the stepsize*

$$\gamma := \frac{R}{L\sqrt{T}},$$

*gradient descent (2.1) yields*

$$\frac{1}{T} \sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{RL}{\sqrt{T}}.$$

*Proof.* This is a simple calculation on top of (2.5): after plugging in the bounds  $R$  and  $L$ , we want to choose  $\gamma$  such that

$$q(\gamma) = \frac{L^2 T \gamma}{2} + \frac{R^2}{2\gamma}$$

is minimized. Setting the derivative to zero yields the above value of  $\gamma$ , and  $q(R/(L\sqrt{T})) = RL\sqrt{T}$ . Dividing by  $T$ , the result follows.  $\square$

This means that in order to achieve  $\min_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \varepsilon$ , we need  $\mathcal{O}(1/\varepsilon^2)$  many iterations, considering  $R$  and  $L$  as constants. This is not particularly good when it comes to concrete numbers (think of desired error  $\varepsilon = 10^{-6}$  when  $R, L$  are somewhat larger). On the other hand, the number of steps does not depend on  $d$ , the dimension of the space. This is very important since we often optimize in high-dimensional spaces. Of course,  $R$  and  $L$  may depend on  $d$ , but in many relevant cases, this dependence is mild.

What happens if we don't know  $R$  and/or  $L$ ? An idea is to "guess"  $R$  and  $L$ , run gradient descent with  $T$  and  $\gamma$  resulting from the guess, check whether the result has absolute error at most  $\varepsilon$ , and repeat with a different guess otherwise. This fails, however, since in order to compute the absolute error, we need to know  $f(\mathbf{x}^*)$  which we typically don't. But Exercise 13 asks you to show that knowing  $R$  is sufficient.

We conclude this section by remarking that bounded gradients are actually equivalent to *Lipschitz continuity* of  $f$ .

**Lemma 2.2** (Exercise 14). *Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be differentiable,  $L \in \mathbb{R}_+$ . Then the following two statements are equivalent.*

- (i)  $\|\nabla f(\mathbf{x})\| \leq L$  for all  $\mathbf{x} \in \mathbb{R}^d$ .
- (ii)  $|f(\mathbf{x}) - f(\mathbf{y})| \leq L\|\mathbf{x} - \mathbf{y}\|$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ .

## 2.4 Smoothness: $\mathcal{O}(1/\varepsilon)$ steps

Our workhorse in the vanilla analysis was the first-order characterization of convexity: for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ , we have

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \quad (2.6)$$

Next we want to require that  $f$  is not "too convex", intuitively meaning that the curvature of the bowl is bounded.

**Definition 2.3.** *Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable,  $L \in \mathbb{R}_+$ .  $f$  is called smooth (with parameter  $L$ ) if*

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d. \quad (2.7)$$

Recall that (2.6) says that for any  $\mathbf{x}$ , the graph of  $f$  is above its tangential hyperplane at  $(\mathbf{x}, f(\mathbf{x}))$ . In contrast, (2.7) says that for any  $\mathbf{x}$ , the graph of  $f$  is below a not-too-steep tangential paraboloid at  $(\mathbf{x}, f(\mathbf{x}))$ ; see Figure 2.2.

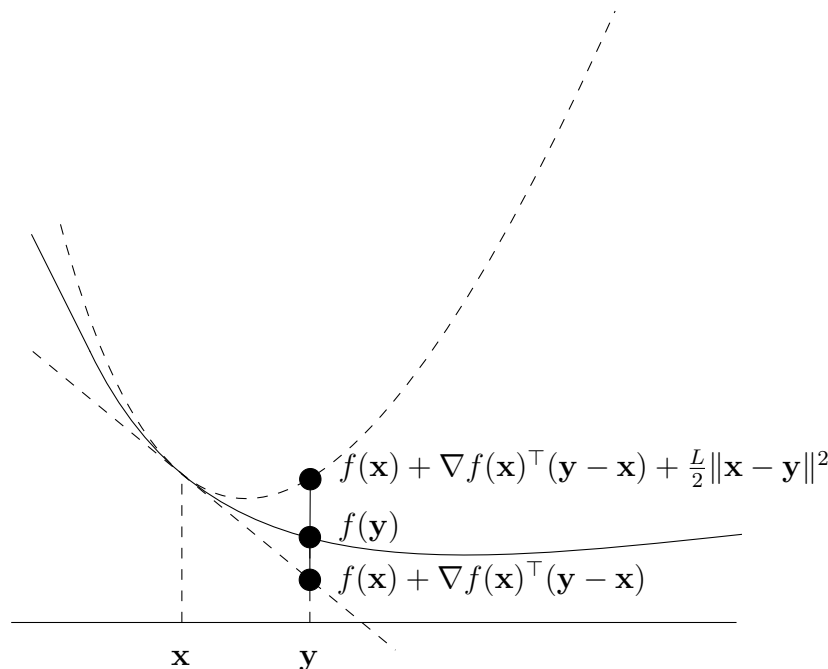


Figure 2.2: A smooth convex function

Let us discuss some cases. If  $L = 0$ , (2.6) and (2.7) together require that

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top(\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d,$$

meaning that  $f$  is an affine function. A simple calculation shows that our supermodel function  $f(x) = x^2$  is smooth with parameter  $L = 2$ , and the same holds for its  $d$ -dimensional generalization  $f(\mathbf{x}) = \|\mathbf{x}\|^2$  (Exercise 11). The convex function  $f(x) = x^4$  is not smooth: at  $x = 0$ , condition (2.7) reads as

$$y^4 \leq \frac{L}{2}y^2,$$

and there is obviously no  $L$  that works for all  $y$ . In general—and this is the important message here—only functions of asymptotically at most quadratic growth can be smooth. It is tempting to believe that any such

“subquadratic” function is actually smooth, but this is not true. Exercise 12(iii) provides a counterexample.

The operations that we have shown to preserve convexity in Lemma 1.9 also preserve smoothness. This immediately gives us a rich collection of smooth functions.

**Lemma 2.4** (Exercise 15).

- (i) Let  $f_1, f_2, \dots, f_m$  be convex functions that are smooth with parameters  $L_1, L_2, \dots, L_m$ , and let  $\lambda_1, \lambda_2, \dots, \lambda_m \in \mathbb{R}_+$ . Then the convex function  $f := \sum_{i=1}^m \lambda_i f_i$  is smooth with parameter  $\sum_{i=1}^m \lambda_i L_i$ .
- (ii) Let  $f$  be a convex function with  $\text{dom}(f) \subseteq \mathbb{R}^d$  that is smooth with parameter  $L$ , and let  $g : \mathbb{R}^m \rightarrow \mathbb{R}^d$  be an affine function, meaning that  $g(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$ , for some matrix  $A \in \mathbb{R}^{d \times m}$  and some vector  $\mathbf{b} \in \mathbb{R}^d$ . Then the convex function  $f \circ g$  (that maps  $\mathbf{x}$  to  $f(A\mathbf{x} + \mathbf{b})$ ) is smooth with parameter  $L\|A\|^2$ , where

$$\|A\| = \max_{\|\mathbf{x}\|=1} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|}$$

is the 2-norm (or spectral norm) of  $A$ .

**Corollary 2.5.** Let  $f(\mathbf{x}) = \|A\mathbf{x} - \mathbf{b}\|^2$  be a least squares objective. Then  $f$  is smooth with parameter  $L = 2\|A\|^2$ .

We next show that for smooth functions, the vanilla analysis provides a better bound than it does under bounded gradients. In particular, we are able to serve the supermodel  $f(x) = x^2$  now.

**Theorem 2.6.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable with a global minimum  $\mathbf{x}^*$ ; furthermore, suppose that  $f$  is smooth with parameter  $L$  according to (2.7). Choosing

$$\gamma := \frac{1}{L},$$

gradient descent (2.1) with arbitrary  $\mathbf{x}_0$  satisfies the following two properties.

- (i) Function values are monotone decreasing:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2, \quad t \geq 0.$$

(ii)

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

*Proof.* For (i), we directly apply the smoothness condition (2.7) and the definition of gradient descent that yields  $\mathbf{x}_{t+1} - \mathbf{x}_t = -\nabla f(\mathbf{x}_t)/L$ . We compute

$$\begin{aligned} f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{L} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2. \end{aligned}$$

In particular, this lets us now bound the sum of squared gradients after step (2.5) of the vanilla analysis:

$$\frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 \leq \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}_{t+1})) = f(\mathbf{x}_0) - f(\mathbf{x}_T). \quad (2.8)$$

With  $\gamma = 1/L$ , (2.5) then yields

$$\begin{aligned} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\ &\leq f(\mathbf{x}_0) - f(\mathbf{x}_T) + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \end{aligned}$$

equivalently

$$\sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2. \quad (2.9)$$

Hence, by (i),

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{1}{T} \sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

□



This improves over the bounds of Theorem 2.1. Again assuming that  $L$  and  $\|\mathbf{x}_0 - \mathbf{x}^*\|^2$  are constant, we now only need  $\mathcal{O}(1/\varepsilon)$  iterations instead of  $\mathcal{O}(1/\varepsilon^2)$  to achieve  $f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \varepsilon$ . Exercise 16 shows that we do not need to know  $L$  to obtain the same asymptotic runtime.

While bounded gradients are equivalent to Lipschitz continuity of  $f$ , smoothness turns out to be equivalent to Lipschitz continuity of  $\nabla f$ .

**Lemma 2.7.** *Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable. The following two statements are equivalent.*

- (i)  $f$  is smooth with parameter  $L$ .
- (ii)  $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ .

A proof can for example be found in the lecture slides of L. Vandenbergh, <http://www.seas.ucla.edu/~vandenbe/236C/lectures/gradient.pdf>.

## 2.5 Interlude

Let us get back to the supermodel  $f(x) = x^2$  (that is smooth with parameter  $L = 2$ , as we observed before). According to Theorem 2.6, gradient descent (2.1) with stepsize  $\gamma = 1/2$  satisfies

$$f(x_t) \leq \frac{1}{t} x_0^2. \quad (2.10)$$

Here we used that the minimizer is  $x^* = 0$ . Let us check how good this bound really is. For our concrete function and concrete stepsize, (2.1) reads as

$$x_{t+1} = x_t - \frac{1}{2} \nabla f(x_t) = x_t - x_t = 0,$$

so we are always done after one step! But we will see in the next section that this is only because the function is particularly beautiful, and on top of that, we have picked the best possible smoothness parameter. To simulate a more realistic situation here, let us assume that we haven't looked at the supermodel too closely and found it to be smooth with parameter  $L = 4$  only (which is a suboptimal but still valid parameter). In this case,  $\gamma = 1/4$  and (2.1) becomes

$$x_{t+1} = x_t - \frac{1}{4} \nabla f(x_t) = x_t - \frac{x_t}{2} = \frac{x_t}{2}.$$

So, we in fact have

$$f(x_t) = f\left(\frac{x_0}{2^t}\right) = \frac{1}{2^{2t}}x_0^2. \quad (2.11)$$

This is still vastly better than the bound of (2.10)! While (2.10) requires  $t \approx x_0^2/\varepsilon$  to achieve  $f(x_t) \leq \varepsilon$ , (2.11) only requires

$$t \approx \frac{1}{2} \log\left(\frac{x_0^2}{\varepsilon}\right),$$

which is an exponential improvement in the number of steps.

## 2.6 Strong convexity: $\mathcal{O}(\log(1/\varepsilon))$ steps

The supermodel function  $f(x) = x^2$  is not only smooth (“not too curved”) but also *strongly convex* (“not too flat”). It will turn out that this is the crucial ingredient that makes gradient descent fast.

**Definition 2.8.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable,  $\mu \in \mathbb{R}_+$ ,  $\mu > 0$ .  $f$  is called *strongly convex* (with parameter  $\mu$ ) if

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d. \quad (2.12)$$

While smoothness according to (2.7) says that for any  $\mathbf{x}$ , the graph of  $f$  is *below* a *not-too-steep* tangential paraboloid at  $(\mathbf{x}, f(\mathbf{x}))$ , strong convexity means that the graph of  $f$  is *above* a *not-too-flat* tangential paraboloid at  $(\mathbf{x}, f(\mathbf{x}))$ . The graph of a smooth *and* strongly convex function is therefore at every point wedged between two paraboloids; see Figure 2.3.

We can also interpret (2.12) as a strengthening of the first-order characterization of convexity. In the form of (2.6) this reads as

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d,$$

and therefore says that every convex function satisfies (2.12) with  $\mu = 0$ .

**Lemma 2.9** (Exercise 18). *If  $f$  is strongly convex with parameter  $\mu > 0$ , then  $f$  is strictly convex and has a unique global minimum.*

The supermodel  $f(x) = x^2$  is particularly beautiful since it is both smooth and strongly convex with the same parameter  $L = \mu = 2$  (going through the calculations in Exercise 11 again will reveal this). We can easily characterize the class of particularly beautiful functions. These are exactly the ones whose sublevel sets are  $\ell_2$ -balls.

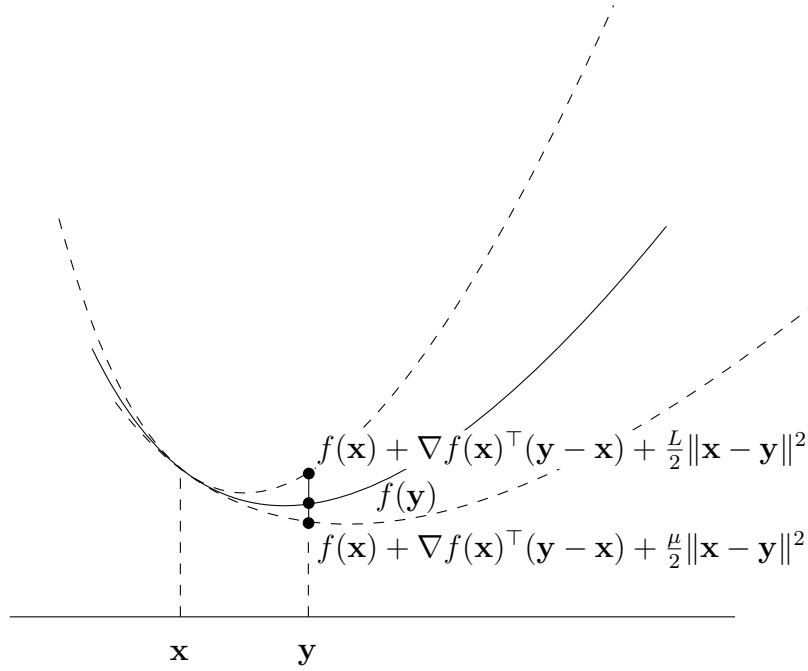


Figure 2.3: A smooth and strongly convex function

**Lemma 2.10** (Exercise 19). Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be strongly convex with parameter  $\mu > 0$  and smooth with parameter  $\mu$ . Prove that  $f$  is of the form

$$f(\mathbf{x}) = \frac{\mu}{2} \|\mathbf{x} - \mathbf{b}\|^2 + c,$$

where  $\mathbf{b} \in \mathbb{R}^d, c \in \mathbb{R}$ .

Once we have a unique global minimum  $\mathbf{x}^*$ , we can attempt to prove that  $\lim_{t \rightarrow \infty} \mathbf{x}_t = \mathbf{x}^*$  in gradient descent. From the vanilla analysis, we already have an inequality that potentially allows us to get started on this, namely (2.4) that we derived from the first-order characterization:

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{1}{2\gamma} (\gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2).$$

If  $f$  is strongly convex, we can start from the strengthening (2.12) instead

to obtain

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{1}{2\gamma} (\gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2. \quad (2.13)$$

Rewriting this yields a bound on  $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2$  in terms of  $\|\mathbf{x}_t - \mathbf{x}^*\|^2$ , along with some “noise” that we still need to take care of:

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq 2\gamma(f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + (1 - \mu\gamma) \|\mathbf{x}_t - \mathbf{x}^*\|^2. \quad (2.14)$$

**Theorem 2.11.** *Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable with a global minimum  $\mathbf{x}^*$ ; furthermore, suppose that  $f$  is smooth with parameter  $L$  according to (2.7) and strongly convex with parameter  $\mu > 0$  according to (2.12). Choosing*

$$\gamma := \frac{1}{L},$$

*gradient descent (2.1) with arbitrary  $\mathbf{x}_0$  satisfies the following two properties.*

(i) *Squared distances to  $\mathbf{x}^*$  are geometrically decreasing:*

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\mu}{L}\right) \|\mathbf{x}_t - \mathbf{x}^*\|^2, \quad t \geq 0.$$

(ii)

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^t \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

*Proof.* For (i), we show that the noise in (2.14) disappears. From Theorem 2.6 (i), we know that

$$f(\mathbf{x}^*) - f(\mathbf{x}_t) \leq f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq -\frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2,$$

and hence the noise can be bounded as follows:

$$\begin{aligned} 2\gamma(f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 &= \frac{2}{L}(f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \frac{1}{L^2} \|\nabla f(\mathbf{x}_t)\|^2 \\ &\leq -\frac{1}{L^2} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{L^2} \|\nabla f(\mathbf{x}_t)\|^2 = 0. \end{aligned}$$

Hence, (2.14) actually yields

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq (1 - \mu\gamma) \|\mathbf{x}_t - \mathbf{x}^*\|^2 = \left(1 - \frac{\mu}{L}\right) \|\mathbf{x}_t - \mathbf{x}^*\|^2.$$

The bound in (ii) follows from smoothness (2.7), using  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  (Lemma 1.13):

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}^*)^\top (\mathbf{x}_t - \mathbf{x}^*) + \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}_t\|^2 = \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}_t\|^2.$$

□

This implies that in order to reach absolute error at most  $\varepsilon$ , we only need  $\mathcal{O}(\log \frac{1}{\varepsilon})$  iterations, where the constant behind the big- $\mathcal{O}$  is roughly  $L/\mu$ .

## 2.7 Exercises

**Exercise 11.** Prove that  $f(\mathbf{x}) = \|\mathbf{x}\|^2$  is smooth with parameter  $L = 2$ .

**Exercise 12.** Consider the function  $f(x) = |x|^{3/2}$ .

- (i) Prove that  $f$  is strictly convex and differentiable, with a unique global minimum  $x^* = 0$ .
- (ii) Prove that for every fixed stepsize  $\gamma$  in gradient descent (2.1) applied to  $f$ , there exists  $x_0$  for which  $f(x_1) > f(x_0)$ .
- (iii) Prove that  $f$  is not smooth.
- (iv) Let  $X \subseteq \mathbb{R}$  be a compact convex set (an interval) such that  $0 \in X$ . Prove that  $f$  is not smooth over  $X$ .

**Exercise 13.** In order to obtain average error at most  $\varepsilon$  in Theorem 2.1, we need to choose iteration number and step size as

$$T \geq \left( \frac{RL}{\varepsilon} \right)^2, \quad \gamma := \frac{R}{L\sqrt{T}}.$$

If  $R$  or  $L$  are unknown, we cannot do this.

But suppose that we know  $R$ . Develop an algorithm that—not knowing  $L$ —finds a vector  $\mathbf{x}$  such that  $f(\mathbf{x}) - f(\mathbf{x}^*) < \varepsilon$ , using at most

$$\mathcal{O} \left( \left( \frac{RL}{\varepsilon} \right)^2 \right)$$

many gradient descent steps!

**Exercise 14.** Prove Lemma 2.2!

**Exercise 15.** Prove Lemma 2.4!

**Exercise 16.** In order to obtain average error at most  $\varepsilon$  in Theorem 2.6, we need to choose

$$\gamma := \frac{1}{L}, \quad T \geq \frac{R^2 L}{2\varepsilon},$$

if  $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$ . If  $L$  is unknown, we cannot do this.

But suppose that we know  $R$ . Develop an algorithm that—not knowing  $L$ —finds a vector  $\mathbf{x}$  such that  $f(\mathbf{x}) - f(\mathbf{x}^*) < \varepsilon$ , using at most

$$\mathcal{O}\left(\frac{R^2 L}{2\varepsilon}\right)$$

many gradient descent steps!

**Exercise 17.** Let  $X = [-a, a] \subseteq \mathbb{R}$ . Prove that  $f(x) = x^4$  is smooth over  $X$  and determine a concrete smoothness parameter  $L$ .

**Exercise 18.** Prove Lemma 2.9!

**Exercise 19.** Prove Lemma 2.10!