

Chapter 2

This lecture covers smooth continuous optimization and provides an introduction to convex optimization. For this purpose, we require several basic definitions such as gradients, Hessian matrices, etc. This chapter also "scratches" the surface of properties of optimization functions: Taylor expansion is reviewed and types of stationary points are introduced. Several special conditions that benefits optimization, including Lipschitz continuity, Lipschitz gradient continuity and convexity are also covered.

The main algorithm for this chapter will be gradient descent (GD), as well as projected GD. Additionally, these notes explain convergence rates. We will see how further global assumptions lead to improved convergence guarantees.

Lipschitz conditions | Gradient Descent

As discussed in Chapter 1, this course mostly covers general smooth optimization, where the objective function can be pictured as a continuous curve in high-dimensions. However, there are other important classes of optimization problems, not covered in this course, as shown in the following figure. Some of them are typically explored as a special topic, for example discrete optimization and integer programming. However, the stress of this course is not restricted only on convex functions, but on more general but smooth functions.

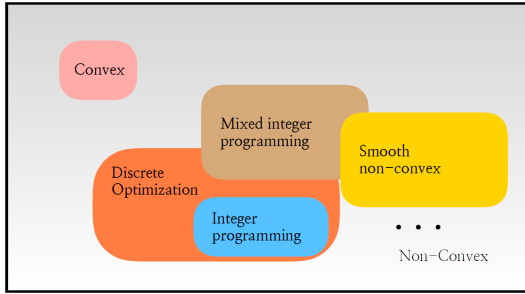


Fig. 3. Landscape of optimization

Derivatives, gradients and Hessians. Algorithms and heuristics in optimization often involve derivatives as a means of approaching an optimal solution. In short, the derivative tells you the direction and magnitude of steepest ascent (or descent) as a linear function.

Definition 1. (First-order Derivative) The derivative of a univariate function $f : \mathbb{R} \rightarrow \mathbb{R}$ at a point x is defined as:

$$\frac{\partial f}{\partial x} = f'(x) = \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon) - f(x)}{\epsilon}.$$

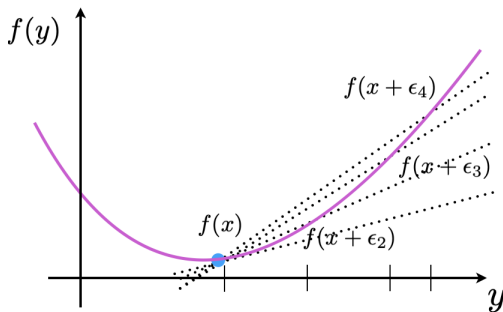


Fig. 5. Graphical illustration of first-order derivative

The derivative of f represents the slope f at a very local area near point x ; i.e., it gives information of how much f changes within a very small area.

This in turn introduces the second-order derivative, which is recursively defined as the derivative of the derivative, which depicts how rapidly the derivative changes.

Definition 2. (Second-order Derivative) The second-order derivative of a univariate function $f : \mathbb{R} \rightarrow \mathbb{R}$ at a point x is defined as:

$$\frac{\partial^2 f}{\partial x^2} = f''(x) = \lim_{\epsilon \rightarrow 0} \frac{f'(x + \epsilon) - f'(x)}{\epsilon}.$$

The second-order derivative represents the *local curvature* of f : how the slope of the function changes.

Moving to higher-dimensions, we introduce the notion of the gradient.

Definition 3. The gradient of a multivariate function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_p} \end{bmatrix} \in \mathbb{R}^p$$

where

$$\begin{aligned} \frac{\partial f}{\partial x_i} &= \lim_{\epsilon \rightarrow 0} \frac{f(\dots, x_{i-1}, x_i + \epsilon, x_{i+1}, \dots) - f(\dots, x_{i-1}, x_i, x_{i+1}, \dots)}{\epsilon} \\ &= \frac{f(x + \epsilon e_i) - f(x)}{\epsilon} \end{aligned}$$

The following fact relates linear forms of the gradient (i.e., inner product) to one-dimensional derivative, evaluated at zero point.

Claim 1. Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a differentiable function. For two points $x, y \in \mathbb{R}^p$ and for scalar γ , we have:

$$\nabla f(x)^\top y = \left. \frac{\partial f(x + \gamma y)}{\partial \gamma} \right|_{\gamma=0}.$$

Similarly, we define derivative for multivariate vector function.

Definition 4. The Jacobian of a multivariate vector function $f : \mathbb{R}^p \rightarrow \mathbb{R}^m$ is given by:

$$Df(x) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_p} \\ \vdots & \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_p} \end{bmatrix} \in \mathbb{R}^{m \times p}$$

Taking the Jacobian of the gradient yields the Hessian or second-order information of f :

Definition 5. The Hessian of a multivariate function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_p} \\ \vdots & \vdots & & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_p} & \frac{\partial^2 f}{\partial x_2 \partial x_p} & \dots & \frac{\partial^2 f}{\partial x_p^2} \end{bmatrix}$$

The Hessian matrix of a continuous function is symmetric. The Hessian matrix also provides information about the curvature of the function f . E.g., for a point x^* , when $\nabla^2 f(x^*) \succ 0$, then x^* is (at least) a strict local minimizer of f . Alternatively, when $\nabla^2 f(x^*) \prec 0$, then x^* is a strict local maximizer of f . See the next figure for some geometric interpretation.

We can relate quadratic forms with the Hessian matrix to one-dimensional derivatives.

Claim 2. Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a twice-differentiable function. Let $x, y \in \mathbb{R}^p$ and γ a scalar. Then:

$$\langle \nabla^2 f(x + \gamma y) \cdot y, y \rangle = \frac{\partial^2 f(x + \gamma y)}{\partial \gamma^2}.$$

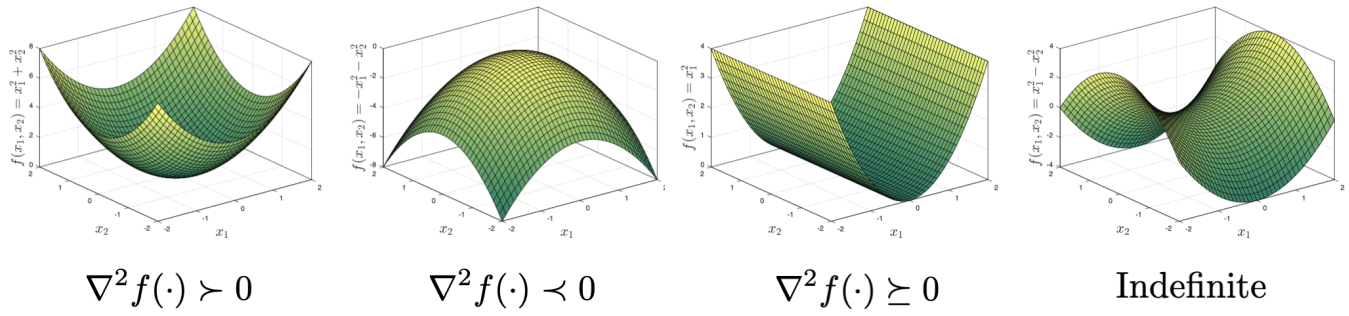


Fig. 6. How Hessian looks like around interesting points of a two-dimensional function f (z-axis).

Taylor expansion of a function f . Are there any intuitive ways of approximating the behavior of a function? Yes, the *Taylor expansion* of the function may be used to approximate the function locally.

Definition 6. (Taylor Series) Assuming that f is n -differentiable, then the Taylor series of f centered at x_0 is

$$\begin{aligned}
 T_\alpha(x) &= \sum_{k=0}^{\infty} \frac{f^{(k)}(\alpha)(x-\alpha)^k}{k!} \\
 &= \frac{f(\alpha)}{0!} + \frac{f'(\alpha)}{1!}(x-\alpha) + \frac{f''(\alpha)}{2!}(x-\alpha)^2 + \dots
 \end{aligned}$$

The k -th order Taylor approximation is the above series truncated at the k^{th} term in the sum.

Here, f is assumed to be differentiable as many times as we like. The Taylor expansion gives a good (local) estimate of the function. E.g., when we keep only the first two terms, it is a linear approximation of the function near α .

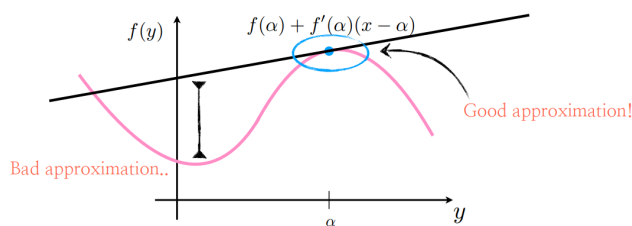


Fig. 4. The first-order Taylor expansion provides a good estimation of the function near the point α , but easily drifts away when we move a little bit away from it.

When keeping the first three terms, we obtain a quadratic approximation of f ; see the next figure.

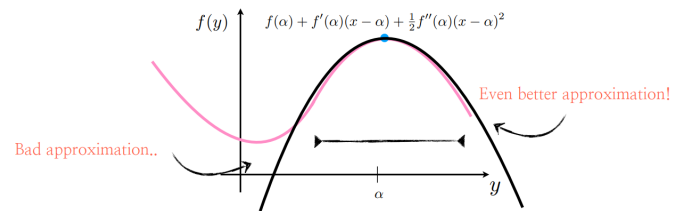


Fig. 7. The second-order Taylor expansion estimates a function better near point α .

Keeping more terms provides a more accurate approximation. For univariate function, this is attainable. However, the complexity increases significantly in high-order Taylor expansion of multivariate functions.

Definition 7. The Taylor expansion of a multivariate function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ at point $\alpha \in \mathbb{R}^p$ is

$$f(x) \approx \frac{f(\alpha)}{0!} + \langle \nabla f(\alpha), x - \alpha \rangle + \frac{1}{2} \langle \nabla^2 f(\alpha)(x - \alpha), (x - \alpha) \rangle + \dots$$

This is a natural generalization of the unidimensional version. For a first-order Taylor expansion approximation, we obtain:

$$f(x) \approx f(\alpha) + \langle \nabla f(\alpha), x - \alpha \rangle, \quad \alpha \in \mathbb{R}^p,$$

while for a second-order one, we obtain:

$$f(x) \approx f(\alpha) + \langle \nabla f(\alpha), x - \alpha \rangle + \frac{1}{2} \langle \nabla^2 f(\alpha)(x - \alpha), x - \alpha \rangle, \quad \alpha \in \mathbb{R}^p,$$

In particular, Taylor's expansion implies the following:

Lemma 1. Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a differentiable function. Let two points $x, y \in \mathbb{R}^p$. Then:

$$f(y) = f(x) + \langle \nabla f(x), y - x \rangle + \int_0^1 (1 - \gamma) \frac{\partial^2 f(x + \gamma(y - x))}{\partial \gamma^2} d\gamma$$

These approximations give us a local approximation of a function, making it possible for us to find good minimizers of the function, which is exactly what we need in iterative optimization methods. I.e., instead of optimizing directly $\min_x f(x)$, we first compute e.g., its second-order approximation around a point x_0 :

$$\min_x \left\{ f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{1}{2} \langle \nabla^2 f(x_0)(x - x_0), x - x_0 \rangle \right\},$$

which in turn is just a minimization of a quadratic function:

$$\min_x \left\{ h^\top x + \frac{1}{2} x^\top H x \right\}$$

for some $h \in \mathbb{R}^p$ and $H \in \mathbb{R}^{p \times p}$. And, solving quadratic problems is a type of optimization that we know how to complete.

Optima. It is never hard to spot the minimum of a function, whenever it can be drawn on a paper. For a computer though, this is like a grid search. Moreover, unfortunately, real problems are usually multidimensional that we cannot even draw it on a paper, and a multidimensional grid is computational prohibitive. Consequently, we have no idea of the global shape of the function and rely on limited local information to search for the minimum. We want to call this as *agnostic optimization*. See Figure 8 and its solution in Figure 9.

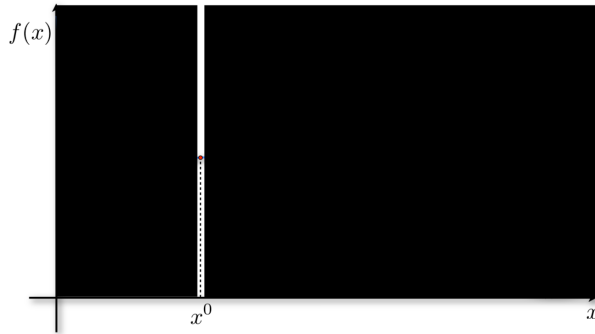


Fig. 8. Agnostic optimization. This point has a 0 derivative. Is it a minimum?

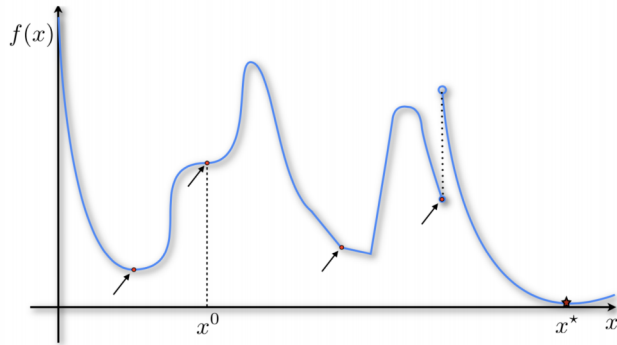


Fig. 9. Agnostic optimization. However, the whole picture is unpredictable. Is it a minimum?

To refer to the optima of a function, we use a set of notations. Without loss of generality, we only discuss minimization.

Definition 8. The global minimizer x^* of a function f satisfies $f(x^*) \leq f(x) \quad \forall x$ in the domain of f .

Definition 9. A local minimizer \hat{x} of a function f satisfies $f(\hat{x}) \leq f(x) \quad \forall x \in \mathcal{N}_{\hat{x}}$,

where $\mathcal{N}_{\hat{x}}$ defines a very small neighborhood around \hat{x} .

We can recognize that a solution is local (global) by the following *necessary* conditions:

- **1st order optimality condition:** $\nabla f(\hat{x}) = 0$.
- **2nd order optimality condition:** $\nabla f(\hat{x}) = 0$ and $\nabla^2 f(\hat{x}) \succeq 0$.

Intuitively, the above states that the function is flat at the point and is convex in its neighborhood. (*Convexity to be defined yet - for the moment think of it as the curve of the function f to look locally as a bowl.*)

Note that these are only necessary conditions, with $f(x) = x^3$ as a simple counterexample at point $x = 0$, which satisfies the two conditions but is not a local minimum.

Lipschitz Conditions

Focus on Figures 8 and 9; there are different points where the gradient is zero, there are points where the gradient is not unique, as well as there are points where the function is discontinuous. Under such generic circumstances, finding the global minima seems to be a difficult task, unless we start making some assumptions about the objective f . Often, many of the objectives f we want to optimize in practice satisfy a form of *Lipschitz continuity*.

Definition 10. A function f is called *Lipschitz continuous*, when

$$|f(x) - f(y)| \leq M \cdot \|x - y\|_2, \quad \forall x, y.$$

This basically means a function should not be too steep, where the steepness is controlled by the constant M . Lipschitz continuity is a stronger than uniform continuity and continuity, and weaker than convexity. Basically, a Lipschitz continuous function may not have abrupt changes.

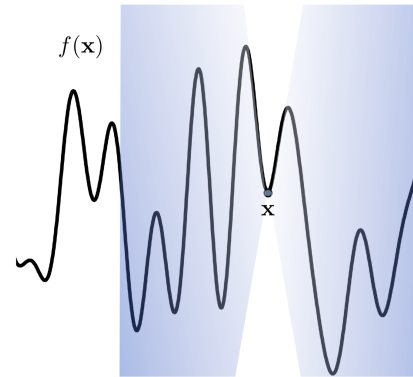


Fig. 10. Illustration of a Lipschitz continuous function, where the width of the cone in white is controlled by M .

A seemingly similar but quite different assumption is that of Lipschitz gradient continuity, where we apply the Lipschitz condition on the gradients of the function.

Definition 11. A function f has *Lipschitz continuous gradients*, when

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|_2, \quad \forall x, y$$

Often, such a function is also called *L-smooth*.

Such a condition forbids sudden changes in the gradient. Using Taylor's expansion, we can prove that

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

which means the function is bounded by a quadratic function.

There are several equivalent characterizations of *Lipschitz gradient continuity* to be aware:

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

$$\nabla^2 f(x) \preceq LI,$$

$$\|\nabla^2 f(x)\|_2 \leq L.$$

Comparison of Lipschitz conditions:

- Lipschitz continuity implies that f should not be too steep.
- Lipschitz gradient continuity implies that changes in the slope of f should not happen suddenly.

Example: Linear regression. In linear regression, the objective $f(x) = \frac{1}{2}\|Ax - b\|^2$ is not Lipschitz continuous—it gets arbitrarily steep when approaching infinity in x —however, it is Lipschitz gradient continuous as

$$\begin{aligned}\|\nabla f(x) - \nabla f(y)\|_2 &= \|A^\top(Ax - b) - A^\top(Ay - b)\|_2 \\ &\leq \|A^\top A\|_2 \cdot \|x - y\|_2,\end{aligned}$$

where $\|A^\top A\|_2$, the largest singular value, serves as the parameter L . This also justifies the equivalent condition:

$$\nabla^2 f(x) \preceq LI.$$

But how can we use the Lipschitz gradient continuity in optimization? A key product of its definition is the inequality:

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

I.e., at a chosen point α , we can upper bound the curve of f (for any y) with a quadratic function, evaluated around α . For a one-dimensional fictional example, one can depict this as in Figure 11.

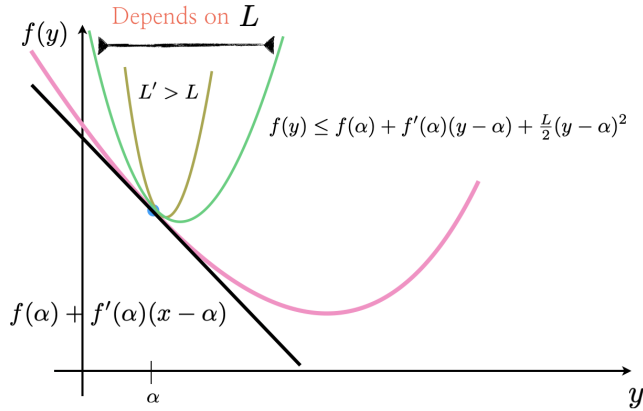


Fig. 11. Illustration of how Lipschitz gradient continuity has algorithmic implications. Here, we want to minimize the one-dimensional $f(y)$ (pink curve). Instead of minimizing f directly—in fact it could be a very complicated function to directly minimize—we will successively construct quadratic (upper-bound) approximations around the current putative solutions, and minimize those approximations. In the figure, we are at point $f(\alpha)$; one can construct the linear local approximation of f around α (black curve); Lipschitz gradient continuity goes further and introduces a quadratic term, “weighted” by the Lipschitz gradient continuity constant L (green curve). Minimizing this quadratic approximation will provide a new point, around which we can form another quadratic approximation, etc. Key observation regarding L is that the larger L is, the steeper the quadratic approximation around the current point is (compare green with khaki curves).

Lipschitz gradient continuity is also related to Taylor’s expansion. We know from the latter that:

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

for some z . Knowing that for a Lipschitz gradient continuous function we have:

$$\nabla^2 f(x) \preceq LI \Rightarrow \|\nabla^2 f(x)\|_2 \leq \|L \cdot I\|_2 \Rightarrow \|\nabla^2 f(x)\|_2 \leq L.$$

Then,

$$\begin{aligned}\frac{1}{2} \langle \nabla^2 f(z)(y - x), y - x \rangle &\leq \frac{1}{2} \|\nabla^2 f(z)\|_2 \cdot \|y - x\|_2^2 \\ &\leq \frac{1}{2} \|\nabla^2 f(z)\|_2 \cdot \|y - x\|_2^2 \\ &\leq \frac{L}{2} \|y - x\|_2^2.\end{aligned}$$

Combining this with the initial Taylor’s expansion expression, we get:

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

Gradient Descent for Lipschitz continuous gradient f

With Lipschitz gradient continuity, we can establish the convergence of an iterative optimization method, such as gradient descent. In fact, gradient descent can be derived as the method of successively minimizing the quadratic approximations around the current point. (Left for exercise).

Gradient descent is defined as follows:

Definition 12. Let f be a differentiable objective with gradient $\nabla f(\cdot)$. The gradient descent method optimize f iteratively, as in:

$$x_{t+1} = x_t - \eta_t \nabla f(x_t), \quad t = 0, 1, \dots,$$

where x_t is the current estimate, and η_t is the step size or learning rate.

The idea behind gradient descent is simple: given the current point x_t , we can compute the negative gradient $-\nabla f(x_t)$ as direction that f has the steepest slope (locally). (Left for exercise). Following that direction, we carefully select η_t , the step size, to dictate how far on that direction we will move.

While gradient descent is quite simple, there are three actions to be made in order to make it work in practice: *i*) how to choose step size η_t , *ii*) initial point x_0 , and *iii*) when to terminate the algorithm.

—For the step size, several approaches are known, some more practical than others, including:

- $\eta_t = \eta$; i.e., the step size is fixed to a value by the user, and stays fixed for all the iterations;
- $\eta_t = O\left(\frac{c}{t}\right)$ or $\eta_t = O\left(\frac{c}{\sqrt{t}}\right)$ for a constant $c > 0$; i.e., the step size keeps decreasing as we go on with the iterations. It starts aggressively (e.g., for $t = 1$ it can be c), but very fast decreases;
- $\eta_t = \operatorname{argmin}_{\eta} f(x_t - \eta \nabla f(x_t))$; i.e., find the step size that minimizes our objective. This approach makes (computationally) sense only for a narrow set of problems, where solving the above problem has *i*) a closed-form solution, and *ii*) it is not difficult to compute that closed-form solution. In the majority of cases, finding the best η_t is computationally prohibited to perform per iteration, and often it requires the same effort as finding the solution to the original problem.
- Fixed step size procedures, such as the Goldstein-Armijo rule; these are out of the scope of this course, and they are often used in classical numerical analysis (*There is definitely some value studying such rules from a machine learning perspective*).

—For the initial value x_0 , because we know little about the function, we usually start from points that make sense (e.g., unless the data involved in the function definition have abruptly large or small values, starting from $x_0 = 0$ makes sense) or we just pick a random value. How to initialize is (almost) irrelevant for some classes of problems (e.g., convex optimization), but it is extremely important for some more broad

(and in some cases more interesting) class of problems; by important we mean that carefully selecting the starting point either leads to some theory—but in practice several starting points lead to the same performance—or that is required to get good performance in practice.

—For the **termination criterion**, there are various standard criterions, like “killing” execution after T iterations (irrespective of whether we converged or not), checking how much progress we make per iteration through $\|x_{t+1} - x_t\|_2$ or $f(x_{t+1}) - f(x_t)$, or by checking if the norm of the gradient is below a threshold, $\|\nabla f(x_t)\|_2 \leq \varepsilon$.

Performance of gradient descent under smoothness assumptions.

Claim 3. Assume that *i)* f is differentiable, and *ii)* that f has L -Lipschitz continuous gradients. Consider the gradient descent iterate: $x_{t+1} = x_t - \eta_t \nabla f(x_t)$. Then:

$$f(x_{t+1}) \leq f(x_t) - \eta_t \left(1 - \frac{\eta_t L}{2}\right) \cdot \|\nabla f(x_t)\|_2^2.$$

Proof: By using the assumption of Lipschitz gradients, we have:

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{L}{2} \|x_{t+1} - x_t\|_2^2 \\ &= f(x_t) + \langle \nabla f(x_t), x_t - \eta_t \nabla f(x_t) - x_t \rangle \\ &\quad + \frac{L}{2} \|x_t - \eta_t \nabla f(x_t) - x_t\|_2^2 \\ &= f(x_t) - \eta_t \|\nabla f(x_t)\|_2^2 + \frac{\eta_t^2 L}{2} \|\nabla f(x_t)\|_2^2 \\ &= f(x_t) - \eta_t \left(1 - \frac{\eta_t L}{2}\right) \|\nabla f(x_t)\|_2^2 \end{aligned}$$

□

The above result indicates that, *i)* as long as $\eta_t \left(1 - \frac{\eta_t L}{2}\right)$ is positive, by performing gradient descent steps, we decrease the objective value by a non-positive quantity $-\eta_t \left(1 - \frac{\eta_t L}{2}\right) \|\nabla f(x_t)\|_2^2$; *ii)* we can maximize the decrease by maximizing the quantity $\eta_t \left(1 - \frac{\eta_t L}{2}\right)$.

Define $g(\eta) := \eta \left(1 - \frac{\eta L}{2}\right)$. Knowing that $\eta > 0$, we first require $1 - \frac{\eta L}{2} > 0 \Rightarrow \eta < \frac{2}{L}$. Thus, for $0 < \eta < \frac{2}{L}$, we observe that the $g(\eta)$ is maximized when we require the gradient satisfies:

$$g'(\eta) = 0 \Rightarrow 1 - \eta L = 0 \Rightarrow \eta = \frac{1}{L}.$$

For the rest of our theory, we will use $\eta_t = \eta = \frac{1}{L}$. Observe that this step size requires the knowledge of L ; for some objectives this is easy to find (linear regression, logistic regression), but for others it is not.

Claim 4. Gradient descent $x_{t+1} = x_t - \eta_t \nabla f(x_t)$, with $\eta_t = \eta = \frac{1}{L}$, satisfies:

$$f(x_{t+1}) \leq f(x_t) - \frac{1}{2L} \|\nabla f(x_t)\|_2^2$$

Proof: This is true by substituting $\eta_t = \frac{1}{L}$ in the result of Claim 3. □

The above characterize the drop in functions values at the t -th iteration. The idea of convergence is based on the idea of relaxation.

Definition 13. A sequence of real numbers $\{\alpha_t\}_{t=0}^\infty$ is called a *relaxation sequence* if $\alpha_{t+1} \leq \alpha_t$, $t \geq 0$.

Combining all the iterations together, for T iterations, we have:

$$\begin{aligned} f(x_{T+1}) &\leq f(x_T) - \frac{1}{2L} \|\nabla f(x_T)\|_2^2 \\ f(x_T) &\leq f(x_{T-1}) - \frac{1}{2L} \|\nabla f(x_{T-1})\|_2^2 \\ &\vdots \\ f(x_1) &\leq f(x_0) - \frac{1}{2L} \|\nabla f(x_0)\|_2^2 \end{aligned}$$

Summing all these inequalities, and under the observation that $f(x^*) \leq f(x_{T+1})$, we get the following claim.

Claim 5. Over T iterations, gradient descent generates a sequence of points x_1, x_2, \dots , such that:

$$\frac{1}{2L} \sum_{t=0}^T \|\nabla f(x_t)\|_2^2 \leq f(x_0) - f(x^*).$$

First, observe that the right hand side is a constant quantity; it does not depend on the number of iterations. Subsequently, the above result implies that, even if we continue running gradient descent for many iterations, the sum of gradient norms is always bounded by something constant; this indicates that the gradient norms that we add at the very end of the execution has to be small, which further implies convergence to a stationary point.

However, the above say nothing about the convergence rate. For that we have the following claim.

Claim 6. Assume we run gradient descent for T iterations, and we obtain T gradients, $\nabla f(x_t)$, for $t \in \{0, \dots, T\}$. Then,

$$\min_{t \in \{0, \dots, T\}} \|\nabla f(x_t)\|_2 \leq \sqrt{\frac{2L}{T+1}} (f(x_0) - f(x^*))^{\frac{1}{2}} = O\left(\frac{1}{\sqrt{T}}\right).$$

Proof: We know that

$$(T+1) \cdot \min_t \|\nabla f(x_t)\|_2^2 \leq \sum_{t=0}^T \|\nabla f(x_t)\|_2^2.$$

Then,

$$\begin{aligned} \frac{T+1}{2L} \cdot \min_t \|\nabla f(x_t)\|_2^2 &\leq \frac{1}{2L} \sum_{t=0}^T \|\nabla f(x_t)\|_2^2 \leq f(x_0) - f(x^*) \Rightarrow \\ \min_t \|\nabla f(x_t)\|_2^2 &\leq \frac{2L}{T+1} \cdot (f(x_0) - f(x^*)) \\ \min_t \|\nabla f(x_t)\|_2 &\leq \sqrt{\frac{2L}{T+1}} \cdot (f(x_0) - f(x^*))^{\frac{1}{2}} \\ &= O\left(\frac{1}{\sqrt{T}}\right). \end{aligned}$$

□

This is called *sublinear convergence rate*. Just to provide a perspective of what it means, focus on Figure 30. Given other convergence rates known, this is a rather pessimistic result. However, have in mind that we made no assumptions other than differentiability of f , and f being a L -smooth function. We will see that making more assumptions helps improving the convergence radically.

Footnote on convergence rates. As a footnote, there are two notations for convergence rate, one using an error level ε that our stopping criterion is based on, and the other using the number of iterations T . For the moment, we know that gradient descent has a convergence rate, with respect to the norm of the gradients, $O\left(\frac{1}{\sqrt{T}}\right)$. Pick a small ε , and assume we require $\min_t \|\nabla f(x_t)\|_2 \leq \varepsilon$. This translates into:

$$\begin{aligned} \sqrt{\frac{2L}{T+1}} \cdot (f(x_0) - f(x^*))^{\frac{1}{2}} &\leq \varepsilon \Rightarrow \\ T+1 &\geq \frac{2L}{\varepsilon^2} \cdot (f(x_0) - f(x^*)) \Rightarrow \\ T &\geq \left\lceil \frac{2L}{\varepsilon^2} \cdot (f(x_0) - f(x^*)) - 1 \right\rceil \end{aligned}$$

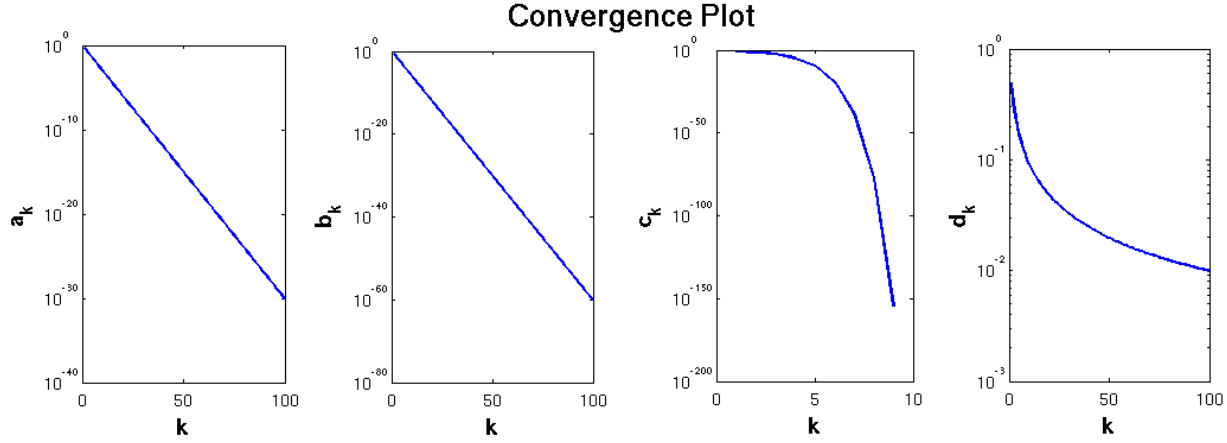


Fig. 12. Borrowed from Wikipedia. Illustration of different convergence rates. Note that y-axis is in logarithmic scale for all the plots, while the x-axis has a linear scale. The y-axis denotes a metric that dictates to the optimum point; think for example $\|x_k - x^*\|_2$ (*We use k as an iteration subscript here*). The x-axis represent the iteration count k . The first two plots represent *linear* convergence rates: it is called linear as a convention to match the linear curve in the *logarithmic* y-axis scale. While the second plot depicts a more preferable behavior, in the big-Oh notation, the two plots are equivalent. For an error level ε , linear convergence rate implies $O(\log \frac{1}{\varepsilon})$. The third plot depicts a *quadratic* convergence rate. For an error level ε , linear convergence rate implies $O(\log \log \frac{1}{\varepsilon})$. Finally, the fourth plot represents the *sublinear* convergence rate; much slower than the linear rate. Some typical rates are: $O(1/\varepsilon^2)$, $O(1/\varepsilon)$, $O(1/\sqrt{\varepsilon})$.

Usually, in order for our convergence rates to make sense, we pick a small value for ε : e.g., let $\varepsilon = 10^{-3}$. Our result dictates that in order to get a solution with $\min_t \|\nabla f(x_t)\|_2 \leq 10^{-3}$, we will need approximately $O(1/\varepsilon^2) = O(10^6)$ iterations (hiding all other constants). This is the meaning of a sublinear convergence rate: in order to get ε accuracy in some sense, we require $1/\varepsilon^2$ iterations. In this course, we will discuss how to achieve better sublinear rates, or even better rates than linear.

Example: Logistic regression. We already discussed the case of linear regression, where the objective $f(x) = \frac{1}{2} \|Ax - b\|_2^2$ has Lipschitz continuous gradients, with constant $L := \|A^\top A\|_2$. Here, we consider another famous—and less straightforward—objective: that of logistic regression. We know that logistic regression is based on the following premise for binary classification:

Given a sample feature vector $\alpha_i \in \mathbb{R}^p$ and a binary class $y_i \in \{\pm 1\}$, define the conditional probability of y_i given α_i as:

$$\mathbb{P}[y_i | \alpha_i, x^*] \propto \frac{1}{1 + \exp(-y_i \alpha_i^\top x^*)}.$$

The above generative assumption leads to the following objective:

$$\min_{x \in \mathbb{R}^p} \left\{ f(x) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i \alpha_i^\top x)) \right\}.$$

Following the same recipe with linear regression, one can compute the gradient and Hessian as

$$\begin{aligned} \nabla f(x) &= \frac{1}{n} \sum_{i=1}^n \nabla [\log(1 + \exp(-y_i \alpha_i^\top x))] \\ &= \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + \exp(-y_i \alpha_i^\top x)} \cdot \nabla_x [\exp(-y_i \alpha_i^\top x)] \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\exp(-y_i \alpha_i^\top x)}{1 + \exp(-y_i \alpha_i^\top x)} \cdot \nabla_x [-y_i \alpha_i^\top x] \\ &= \frac{1}{n} \sum_{i=1}^n \frac{-y_i}{1 + \exp(y_i \alpha_i^\top x)} \alpha_i^\top \end{aligned}$$

and

$$\begin{aligned} \nabla^2 f(x) &= \frac{1}{n} \sum_{i=1}^n \frac{y_i}{(1 + \exp(y_i \alpha_i^\top x))^2} \cdot \nabla [1 + \exp(y_i \alpha_i^\top x)] \cdot \alpha_i^\top \\ &= \frac{1}{n} \sum_{i=1}^n \frac{y_i^2}{(1 + \exp(y_i \alpha_i^\top x))^2} \cdot \exp(y_i \alpha_i^\top x) \cdot \alpha_i \alpha_i^\top \\ &= \frac{1}{n} \sum_{i=1}^n \underbrace{\frac{1}{(1 + \exp(y_i \alpha_i^\top x))^2}}_{\text{scalar}} \cdot \underbrace{\exp(y_i \alpha_i^\top x) \cdot \alpha_i \alpha_i^\top}_{\in \mathbb{R}^{p \times p}} \end{aligned}$$

Observe that, for $\beta \in \mathbb{R}$,

$$\frac{1}{(1 + \exp(\beta))^2} \cdot \exp(\beta) = \frac{1}{1 + \exp(\beta)} \cdot \frac{\exp(\beta)}{1 + \exp(\beta)} = \frac{1}{1 + \exp(\beta)} \cdot \frac{1}{1 + \exp(-\beta)}$$

Define $h(\beta) = \frac{1}{1 + \exp(-\beta)}$, and observe that h maps to $(0, 1)$. Also observe that $h(-\beta) = 1 - h(\beta)$. Then, one can check (*Check it!*) that $h(\beta) \cdot h(-\beta) \leq \frac{1}{4}$.

Going back to our Hessian derivations:

$$\begin{aligned} \nabla^2 f(x) &= \frac{1}{n} \sum_{i=1}^n h(y_i \alpha_i^\top x) \cdot h(-y_i \alpha_i^\top x) \cdot \alpha_i \alpha_i^\top \\ &\leq \frac{1}{n} \sum_{i=1}^n \frac{\alpha_i \alpha_i^\top}{4} = \frac{1}{4n} \cdot A^\top A. \end{aligned}$$

where A accumulates all α_i 's as rows. Thus

$$\|\nabla^2 f(x)\|_2 \leq \frac{1}{4n} \|A^\top A\|_2 := L.$$

Example: $f(x) = x^2 + 3 \sin^2(x)$. This is a less practical example, but it is an example that does not satisfy some nice properties that linear regression and logistic regression satisfy (*Spoiler alert: they are both convex objectives*). The objective looks like:

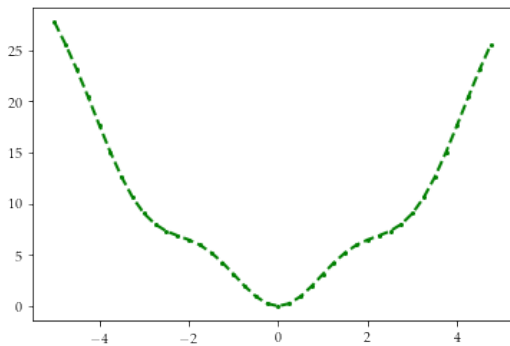


Fig. 13. $f(x) = x^2 + 3 \sin^2(x)$

Let us compute the first and second derivatives of this function:

$$f'(x) = 2x + 6 \sin(x) \cdot \cos(x)$$

and

$$f''(x) = 2 + 6 \cos^2(x) - 6 \sin^2(x)$$

Plotting the Hessian function, we obtain:

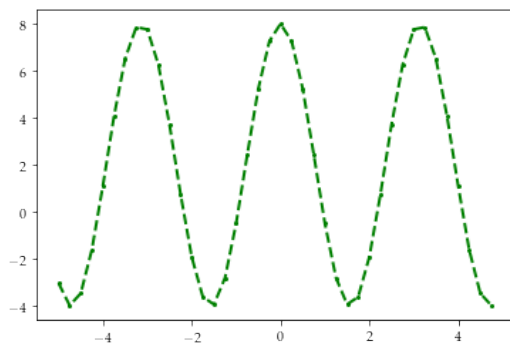


Fig. 14. $f''(x) = 2 + 6 \cos^2(x) - 6 \sin^2(x)$

By inspection (and based on the periodicity of the Hessian function), we can bound:

$$|f''(x)| \leq 8 := L.$$