

Lecture 9. Descent Method

9.1 Unconstrained optimization problems

We now study the general convex optimization problems. First, we consider the easiest case: no constraints. Namely, the optimization problem is

$$\min_{x \in \mathbb{R}^n} f(x)$$

where $f(x)$ is a convex function.

Recall that, the optimality condition for convex functions is

Theorem

Suppose $f : D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex function. Then x^* is a *global* minimum point of f iff

$$\forall y \in D, \quad \nabla f(x^*)^\top (y - x^*) \geq 0.$$

In particular, if $D = \mathbb{R}^n$, then x^* is a global minimum point iff $\nabla f(x^*) = \mathbf{0}$.

For convenience, we assume that $D = \mathbb{R}^n$, the objective function $f(x)$ is differentiable and has a finite minimum point x^* (and the minimum value f^*). For some simple cases, we can compute the minimum point by solving the equation $\nabla f(x^*) = 0$. However, in general we cannot expect that closed-form solutions always exist. So we introduce some algorithms to find optimal solutions.

9.2 Descent method

Analogously to the simplex method, we would like to move from a solution x to a better “neighbor” y . The convexity guarantees that

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x).$$

As we hope y is better, i.e., $f(y) < f(x)$, it requires that $\nabla f(x)^\top (y - x) < 0$.

Conversely, we know that if the directional derivative $\nabla f(x)^\top v < 0$, then there exists $\varepsilon > 0$ such that $f(x + \varepsilon v) < f(x)$. So $\nabla f(x)^\top v < 0$ is a reasonable requirement for the moving direction v .

This inspired the so-called *descent method*: start from a solution \mathbf{x}_0 and move to $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{v}_k$ iteratively, where t_k is the *step size* to be determined and \mathbf{v}_k is the moving direction satisfying $\nabla f(\mathbf{x}_k)^\top \mathbf{v}_k < 0$.

The first question is when we can stop? Of course the ideal stopping criterion is $\nabla f(\mathbf{x}_k) = \mathbf{0}$ for some k . If so, we know that \mathbf{x}_k is indeed a minimum point. However, in practice, we cannot expect this happens. So we usually use stopping criteria such as $\|\nabla f(\mathbf{x})\| < \delta$, $|f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k)| < \delta$, or 1000 iterations.

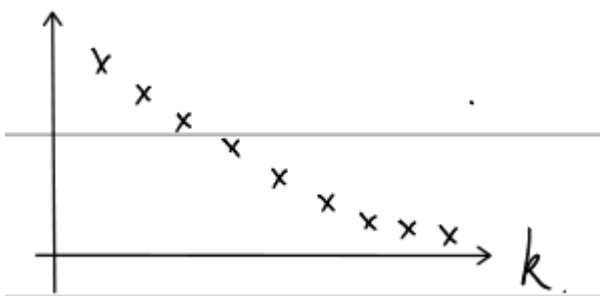
```

given a starting point  $\mathbf{x}_0$ 
repeat
    choose a proper step size  $t_k$ 
     $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + t_k \mathbf{v}_k$  where  $\nabla f(\mathbf{x}_k)^\top \mathbf{v}_k < 0$ 
     $k \leftarrow k + 1$ 
until  $\|\nabla f(\mathbf{x}_k)\| \leq \delta$  for some sufficiently small  $\delta$ 

```

The next question is, does this algorithm converge to an optimal solution? In fact, we claim that if we assume that t_k, \mathbf{v}_k only depend on \mathbf{x}_k , and the choice of t_k satisfies $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$ for every $\mathbf{x}_k \notin \arg \min f(\mathbf{x})$ (note that the optimal solution may not be unique), then the value of objective functions $\{f(\mathbf{x}_k)\}$ generated by the descent method converge to the minimum value f^* . (However, $\{\mathbf{x}_k\}$ may not converge and we will give an example later.)

We assume $f(\mathbf{x})$ has a finite minimum value f^* , and $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$. So $f(\mathbf{x}_k)$ has a limit as k goes to infinity. Now we would like to show that the limit is f^* .



Let $c = \lim_{k \rightarrow \infty} f(\mathbf{x}_k)$. Intuitively, if $c > f^*$, as we hope $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$ as long as $f(\mathbf{x}_k) \neq f^*$, we can argue that $f(\mathbf{x}_{k+1})$ still decreases too fast even if $f(\mathbf{x}_k)$ is sufficiently close to c .

Rigorously, let $S = \{\mathbf{x} \mid c \leq f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$. Then S is a *compact set*, if we assume $f(\infty) = \infty$ for convenience (otherwise S may not be necessarily bounded). Let $g(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function defined by

$$g(\mathbf{y}) = f(\mathbf{y} + t_{\mathbf{y}}\mathbf{v}_{\mathbf{y}}) - f(\mathbf{y}),$$

where $t_{\mathbf{y}}, \mathbf{v}_{\mathbf{y}}$ are the step size and the direction we choose if $\mathbf{x}_k = \mathbf{y}$. That is, $g(\mathbf{y})$ measures the difference between $f(\mathbf{x}_{k+1})$ and $f(\mathbf{x}_k)$ if we set $\mathbf{x}_k = \mathbf{y}$.

By our assumption $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$ as long as $f(\mathbf{x}_k) \neq f^*$, and noting that $\mathbf{x} \in S$ if $f(\mathbf{x}) \geq c > f^*$, we conclude that $g(\mathbf{x}) > 0$ for all $\mathbf{x} \in S$. Applying the extreme value theorem, there exists

$$\delta = \min_{\mathbf{x} \in S} g(\mathbf{x}) > 0,$$

which implies that $f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \delta$ for every $\mathbf{x}_k \in S$. This contradicts our assumption that there exists $\{\mathbf{x}_k\}$ such that $f(\mathbf{x}_k) \downarrow c$, and thus completes the proof.

Tip

In fact, it is not necessary to define g as the difference between the function values. Analogously to the *amortised analysis* for some data structures, we may define g to measure the difference between some *potential function*. So this argument above is a simplified result of the *Lyapunov's global stability theorem in discrete time*.

Suppose $d_{k+1} = \rho(d_k)$ where $\rho : \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuous function and $\rho(0) = 0$. If there exists a continuous (Lyapunov) function $\ell : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

1. $\ell(0) = 0, \ell(x) > 0$ for all $x \neq 0$, (*positivity*)
2. $\ell(x) \rightarrow \infty$ as $\|x\| \rightarrow \infty$, (*radical unboundedness*)
3. $\ell(\rho(x)) < \ell(x)$ for all $x \neq 0$. (*strict decrease*)

Then for all $d_0 \in \mathbb{R}^n$, we have $d_k \rightarrow 0$ as $k \rightarrow \infty$.

For our setting, just select an optimal solution \mathbf{x}^* , and set $d_k = \mathbf{x}_k - \mathbf{x}^*$, $\rho(d_k) = d_k + t_k \mathbf{v}_k$ and $\ell(d_k) = f(d_k + \mathbf{x}^*) - f^*$.

9.3 Gradient descent

We now consider a specific descent method, the *gradient descent*, where we select $\mathbf{v}_k = -\nabla f(\mathbf{x}_k)$. Then trivially $\nabla f(\mathbf{x}_k)^\top \mathbf{v}_k < 0$.

There is an advantage to choose $-\nabla f(x_k)$ since it is the direction of *steepest descent*, namely, the value of f decreases most rapidly: For any *unit* length vector v , the directional derivative $\nabla f(x)^\top v$ satisfies

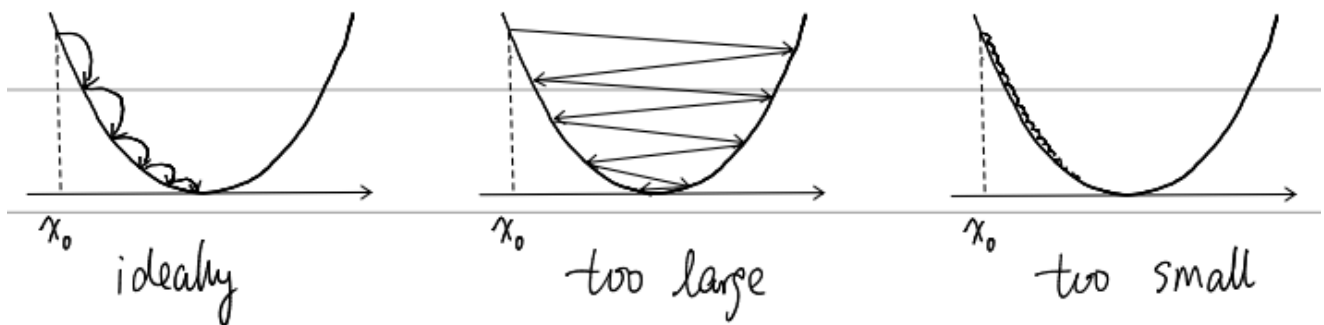
$$-\|v\| \cdot \|\nabla f(x)\| \leq \nabla f(x)^\top v \leq \|v\| \cdot \|\nabla f(x)\|$$

by the Cauchy-Schwarz inequality, and the equality holds iff $v = \pm \nabla f(x) / \|\nabla f(x)\|$.

Applying this choice of directions, we obtain the *gradient descent method*:

given a starting point x_0
 repeat
 choose a proper step size t_k
 $x_{k+1} \leftarrow x_k - t_k \nabla f(x_k)$
 $k \leftarrow k + 1$
 until $\|\nabla f(x_k)\| \leq \delta$ for some sufficiently small δ

We now consider how to choose the step size t_k . Intuitively, the choice of step size can effect the converge rate of the algorithm.



Let's start from an easy example: $f(x) = ax^2$ where $a > 0$. Since we hope $f(x_{k+1}) < f(x_k)$, it requires that $|x_{k+1}| < |x_k|$, which is equivalent to

$$|(1 - 2at_k)x_k| < |x_k|.$$

So $t_k < 1/a$ suffices.

Next, consider the multivariate function $f(\mathbf{x}) = \mathbf{x}^\top Q \mathbf{x}$ where $Q \succeq 0$. Now $\mathbf{x}_{k+1} = \mathbf{x}_k - 2t_k Q \mathbf{x}_k$. So

$$f(\mathbf{x}_{k+1}) = \mathbf{x}_k^\top Q \mathbf{x}_k + 4t_k^2 (Q \mathbf{x}_k)^\top Q (Q \mathbf{x}_k) - 4t_k (Q \mathbf{x}_k)^\top (Q \mathbf{x}_k).$$

It is sufficient to find a value of t_k such that for all $\mathbf{v} \in \mathbb{R}^n$, $t_k \mathbf{v}^\top Q \mathbf{v} < \mathbf{v}^\top \mathbf{v}$. We need the following lemma.

Lemma (Rayleigh quotient)

Let $Q \succeq 0$ be a positive semi-definite matrix, and λ_{\min} and λ_{\max} be its minimum and maximum eigenvalues, respectively. Then for all $x \in \mathbb{R}^n$, we have

$$\lambda_{\min} \|x\|_2^2 \leq x^T Q x \leq \lambda_{\max} \|x\|_2^2$$

Proof

Since $Q \in \mathbb{R}^{n \times n}$ is symmetric, consider its eigen-decomposition $Q = U \Lambda U^T$, where $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_n\}$ is the diagonal matrix consisting of Q 's eigenvalues, and $U = (u_1, \dots, u_n)$ consists of corresponding unit-length eigenvectors. It is easy to see that $U U^T = I$.

Assume $x = U y$ (i.e. $y = U^{-1} x = U^T x$). Then

$$x^T Q x = y^T U^T Q U y = y^T U^T U \Lambda U^T U y = y^T \Lambda y = \sum_{i=1}^n \lambda_i y_i^2.$$

So clearly we have $\lambda_{\min} \|y\|^2 \leq x^T Q x \leq \lambda_{\max} \|y\|^2$. Moreover, we have

$$\|y\|^2 = y^T y = x^T U^T U x = x^T x = \|x\|^2,$$

which completes the proof.

Note that in this proof we do not really need $Q \succeq 0$. This lemma holds for all symmetric Q . Applying this lemma, it gives that $t_k < 1/\lambda_{\max}$ suffices in the gradient descent method for quadratic functions.

However, for general cases, we cannot expect a universal condition for t_k . For example, consider the function $f(x) = |x|$. If we choose t_k to be a constant $t > 0$, no matter what value t is, the algorithm does not work as long as $|x_k| < t$.

Question

Under which assumptions can we choose a constant as the step size?

9.4 L -smooth functions

We would like to avoid functions similar to $|x|$, where $\nabla f(x)$ changes too drastically near x^* .

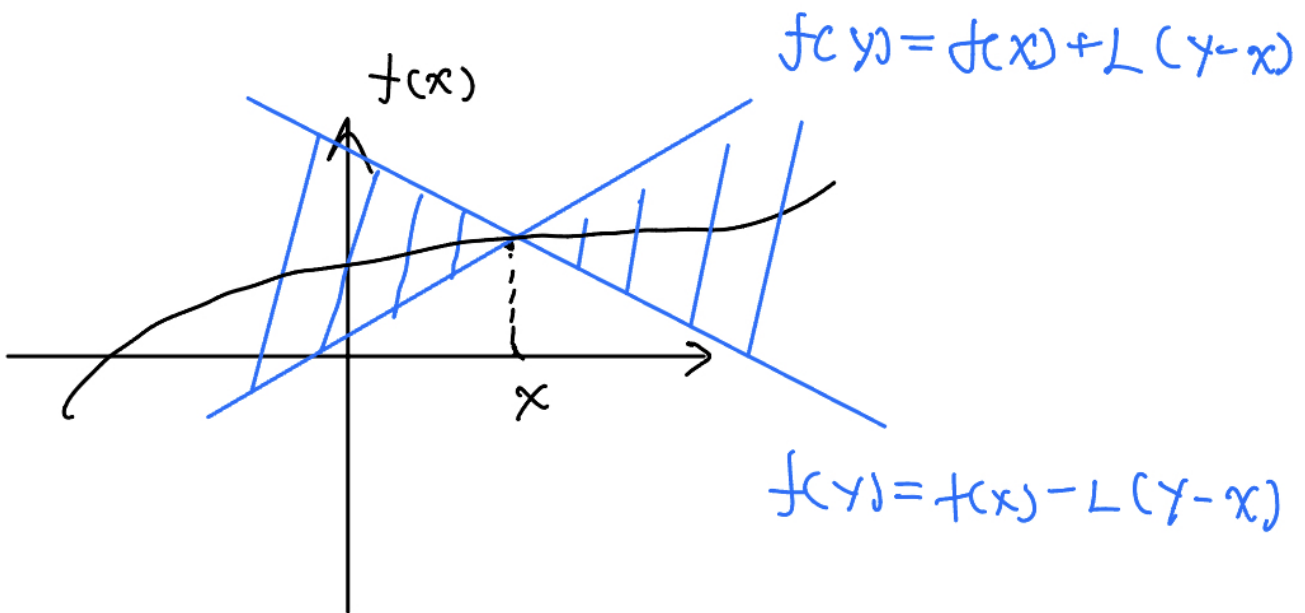
Definition (Lipschitz continuity)

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is L -Lipschitz, if for all $x, y \in \text{dom } f$,

$$\|f(x) - f(y)\| \leq L\|x - y\|.$$

We usually use L^2 -norm, unless otherwise specified.

An L -Lipschitz function is continuous, but may not be differentiable. Intuitively, for a Lipschitz continuous function, there exists a double cone (white) whose origin can be moved along the graph so that the whole graph always stays outside the double cone.



Example

- $f(x) = kx$ where $x \in \mathbb{R}$ is $|k|$ -Lipschitz.
- $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ where $\mathbf{x} \in \mathbb{R}^n$ is $\|\mathbf{w}\|$ -Lipschitz
- $f(\mathbf{x}) = Q\mathbf{x}$ where $\mathbf{x} \in \mathbb{R}^n$ is $\lambda_{\max}(Q^\top Q)^{1/2}$ -Lipschitz, since

$$\begin{aligned} \|f(\mathbf{x}) - f(\mathbf{y})\| &= \|Q(\mathbf{x} - \mathbf{y})\| = ((\mathbf{x} - \mathbf{y})^\top Q^\top Q(\mathbf{x} - \mathbf{y}))^{1/2} \\ &\leq \lambda_{\max}(Q^\top Q)^{1/2} \|\mathbf{x} - \mathbf{y}\| \end{aligned}$$

by the bound for the Rayleigh quotient. In particular, if Q is symmetric,

$$\lambda_{\max}(Q^\top Q)^{1/2} = \max\{|\lambda_{\min}(Q)|, |\lambda_{\max}(Q)|\}.$$

Recall that we hope $\nabla f(x)$ does not change rapidly. So we define the following notion of "smoothness".

Definition (Smoothness)

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is L -smooth if ∇f is L -Lipschitz, i.e., for all x, y ,

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|.$$

Example

$f(x) = x^T Q x$ with $Q \succeq 0$ is $2\lambda_{\max}(Q)$ -smooth ($\nabla f(x) = 2Qx$).

We use the notation $A \succeq B$ if $A - B \succeq 0$. Then we have the following equivalent definitions.

Lemma

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a twice differentiable function. Then f is L -smooth iff $-L\mathbf{I}_n \preceq \nabla^2 f(x) \preceq L\mathbf{I}_n$ for all $x \in \mathbb{R}^n$, where \mathbf{I}_n is the $n \times n$ identity matrix. Namely, for all $x \in \mathbb{R}^n$, $|\lambda_i(\nabla^2 f(x))| \leq L$, where $\lambda_1, \dots, \lambda_n$ are n eigenvalues.

Note that if $f : \mathbb{R} \rightarrow \mathbb{R}$, we can easily prove the " \Leftarrow " direction since the mean value theorem gives that $f'(x) - f'(y) = f''(z)(x - y)$ for some z . However, there is no such theorem for vector-valued functions.

Proof \checkmark

- " \Leftarrow " direction. We would like to restrict the vector-valued function ∇f to a line. Fix any $x, y \in \mathbb{R}^n$. Let $\varphi : [0, 1] \rightarrow \mathbb{R}$ be a function defined by

$$\varphi(t) = \langle \nabla f(y) - \nabla f(x), \nabla f(x + t(y - x)) \rangle.$$

Then, $\varphi(1) = \langle \nabla f(y), \nabla f(y) - \nabla f(x) \rangle$ and

$\varphi(0) = \langle \nabla f(x), \nabla f(y) - \nabla f(x) \rangle$. By the mean value theorem, there exists $t \in [0, 1]$ such that $\varphi(1) - \varphi(0) = \varphi'(t)$. Note that

$$\begin{aligned} \varphi'(t) &= \langle \nabla f(y) - \nabla f(x), \nabla^2 f(x + t(y - x))(y - x) \rangle \\ &\leq \|\nabla f(y) - \nabla f(x)\| \cdot \|\nabla^2 f(x + t(y - x))(y - x)\| \end{aligned}$$

by the Cauchy-Schwarz inequality. It implies that

$$\begin{aligned}\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|^2 &= \varphi(1) - \varphi(0) \\ &\leq \|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\| \cdot \|\nabla^2 f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))(\mathbf{y} - \mathbf{x})\|,\end{aligned}$$

which further gives that

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\| \leq \|\nabla^2 f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))(\mathbf{y} - \mathbf{x})\| \leq L\|\mathbf{y} - \mathbf{x}\|.$$

The last inequality follows from the third example of Lipschitz functions.

- " \implies " direction. Fix any $\mathbf{x}, \mathbf{v} \in \mathbb{R}^n$. Let $\psi : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$ be a function defined by

$$\psi(t) = \langle \nabla f(\mathbf{x} + t\mathbf{v}), \mathbf{v} \rangle.$$

Then, by the Cauchy-Schwarz inequality and the L -smoothness, we have

$$\begin{aligned}|\psi(t) - \psi(0)| &= |\langle \nabla f(\mathbf{x} + t\mathbf{v}) - \nabla f(\mathbf{x}), \mathbf{v} \rangle| \\ &\leq \|\nabla f(\mathbf{x} + t\mathbf{v}) - \nabla f(\mathbf{x})\| \cdot \|\mathbf{v}\| \\ &\leq tL\|\mathbf{v}\|^2,\end{aligned}$$

which further gives that $\left| \frac{\psi(t) - \psi(0)}{t} \right| \leq L\|\mathbf{v}\|^2$. Taking the limit $t \rightarrow 0$ on both sides, and applying the chain rule, we obtain that

$$|\mathbf{v}^\top \nabla^2 f(\mathbf{x}) \mathbf{v}| = |\psi'(0)| \leq L\|\mathbf{v}\|^2.$$

Thus, $-L\mathbf{I}_n \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}_n$.

An L -smooth functions may be not convex. If f is further convex, all absolute values are not necessary.

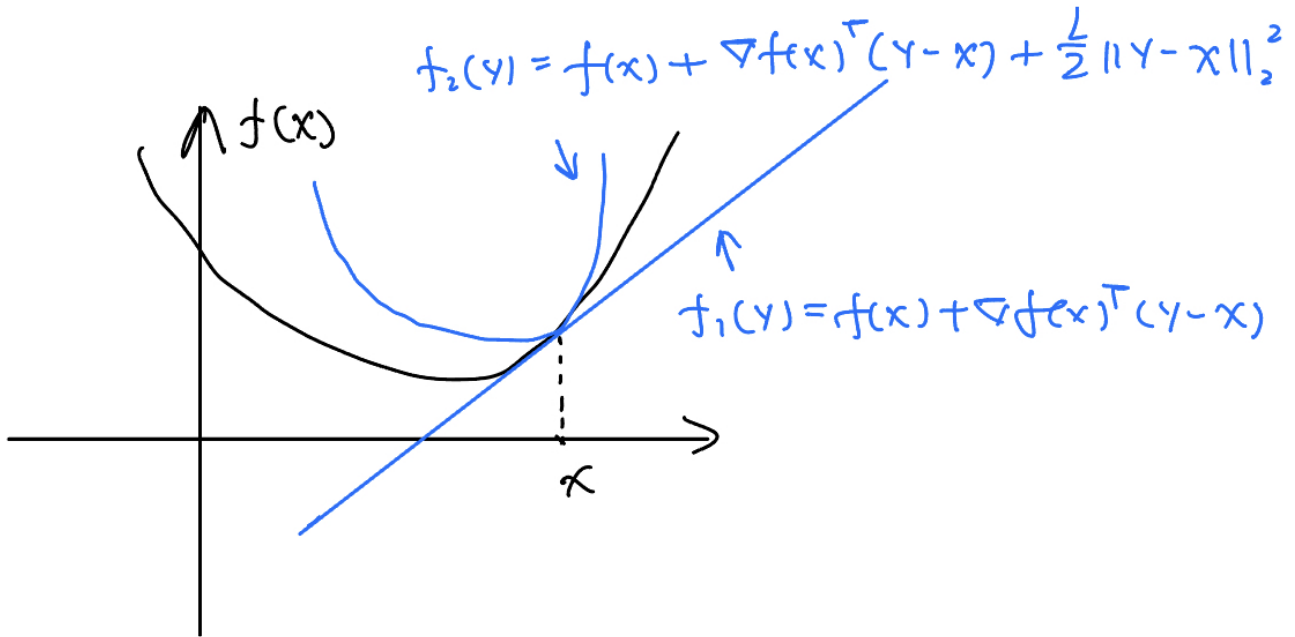
Lemma

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a differentiable function. Then f is L -smooth iff for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

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

Recall that f is convex iff $f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \geq 0$, which shows that f is underestimated by an affine function. Now, if f is L -smooth, it is overestimated by

a quadratic function.



Proof \checkmark

- " \Leftarrow " direction. Fix $\mathbf{x} \in \mathbb{R}^n$. Define

$$g_1(\mathbf{y}) = f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2,$$

$$g_2(\mathbf{y}) = f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle - \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

Note that for all $\mathbf{y} \in \mathbb{R}^n$, $g_2(\mathbf{y}) \leq 0 \leq g_1(\mathbf{y})$, and $g_1(\mathbf{x}) = g_2(\mathbf{x}) = 0$. So \mathbf{x} is a local minimum point of g_1 , which gives that $\nabla^2 g_1(\mathbf{x}) \succeq 0$. Since $\nabla^2 g_1(\mathbf{y}) = \nabla^2 f(\mathbf{y}) + L\mathbf{I}_n$, we conclude that $\nabla^2 f(\mathbf{x}) \succeq -L\mathbf{I}_n$. Similarly, \mathbf{x} is a local maximum point of g_2 , and thus $\nabla^2 g_2(\mathbf{x}) = \nabla^2 f(\mathbf{x}) - L\mathbf{I}_n \preceq 0$.

- " \Rightarrow " direction. Fix $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. Let

$$h(\theta) = f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})).$$

It is clear that $h'(\theta) = \langle \nabla f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})), \mathbf{y} - \mathbf{x} \rangle$, and

$$f(\mathbf{y}) - f(\mathbf{x}) = h(1) - h(0) = \int_0^1 h'(\theta) d\theta.$$

Moreover, $\langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle = h'(0) = \int_0^1 h'(0) d\theta$. Therefore, it holds that

$$f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle = \int_0^1 h'(\theta) - h'(0) d\theta.$$

Note that

$$\begin{aligned} |h'(\theta) - h'(0)| &= |\nabla f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x}| \\ &\leq \|\nabla f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})\| \cdot \|\mathbf{y} - \mathbf{x}\| \\ &\leq \theta L \|\mathbf{y} - \mathbf{x}\|^2. \end{aligned}$$

We now have

$$\begin{aligned} |\langle f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| &\leq \int_0^1 |h'(\theta) - h'(0)| d\theta \\ &\leq \int_0^1 \theta L \|\mathbf{y} - \mathbf{x}\|^2 d\theta = \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2, \end{aligned}$$

which completes the proof.

Recall that, we hope to find the value of the step size t such that $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$. Now we assume that f is L -smooth. Then

$$\begin{aligned} f(\mathbf{x}_{k+1}) &= f(\mathbf{x}_k - t \cdot \nabla f(\mathbf{x}_k)) \\ &\leq f(\mathbf{x}_k) - \langle \nabla f(\mathbf{x}_k), t \cdot \nabla f(\mathbf{x}_k) \rangle + \frac{L}{2} \|t \cdot \nabla f(\mathbf{x}_k)\|^2 \\ &= f(\mathbf{x}_k) - t \left(1 - \frac{Lt}{2}\right) \|\nabla f(\mathbf{x}_k)\|^2 \\ &< f(\mathbf{x}_k) \end{aligned}$$

if we set $t < 2/L$. In particular, if we choose $t \leq 1/L$, it gives the following *descent lemma*.

Lemma (*Descent lemma*)

For an L -smooth differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ (not necessarily convex), and $t \leq 1/L$, we have

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