

Gradient Descent and Convergence Analysis

1 Descent method for 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.

Suppose we know a solution x , which is not optimal. Can we find a better solution 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 direction $v = y - x$.

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 ideal stopping criterion is $\nabla f(\mathbf{x}_k) = \mathbf{0}$ for some k . If so, we know that 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  $\delta$ 

```

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 $-\nabla f(\mathbf{x}_k)$ 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*:

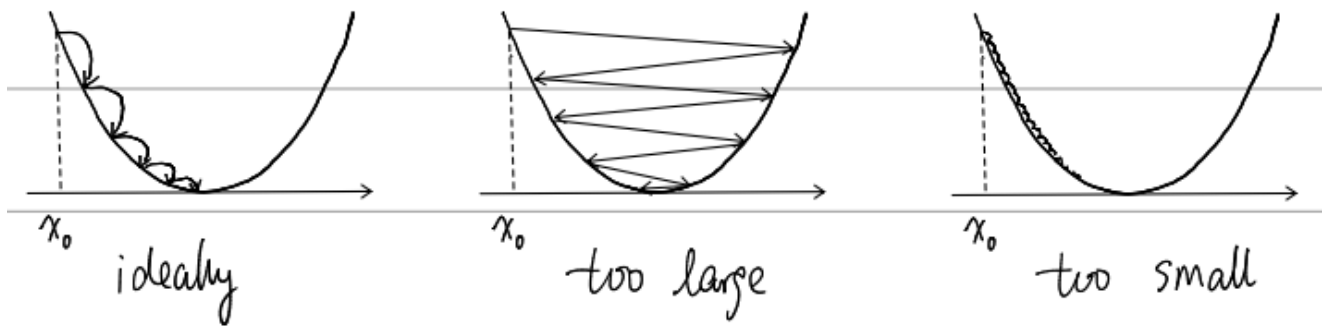
```

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

```

The key problem is how to choose an appropriate step size t_k to guarantee that the $\{x_k\}$ converges to the optimal solution. Intuitively, the choice of step size can effect

the correctness and the converge rate of the algorithm.



The easiest way is to choose a constant. 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

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

Proof

Since $Q \in \mathbb{R}^{n \times n}$ is symmetric, consider its eigen-decomposition $Q = U \Lambda U^\top$, 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 easy to see that $U U^\top = I$.

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

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

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

$$\|\mathbf{y}\|^2 = \mathbf{y}^\top \mathbf{y} = \mathbf{x}^\top U^\top U \mathbf{x} = \mathbf{x}^\top \mathbf{x} = \|\mathbf{x}\|^2,$$

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?

2. Smoothness

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

Definition (Lipschitz continuity)

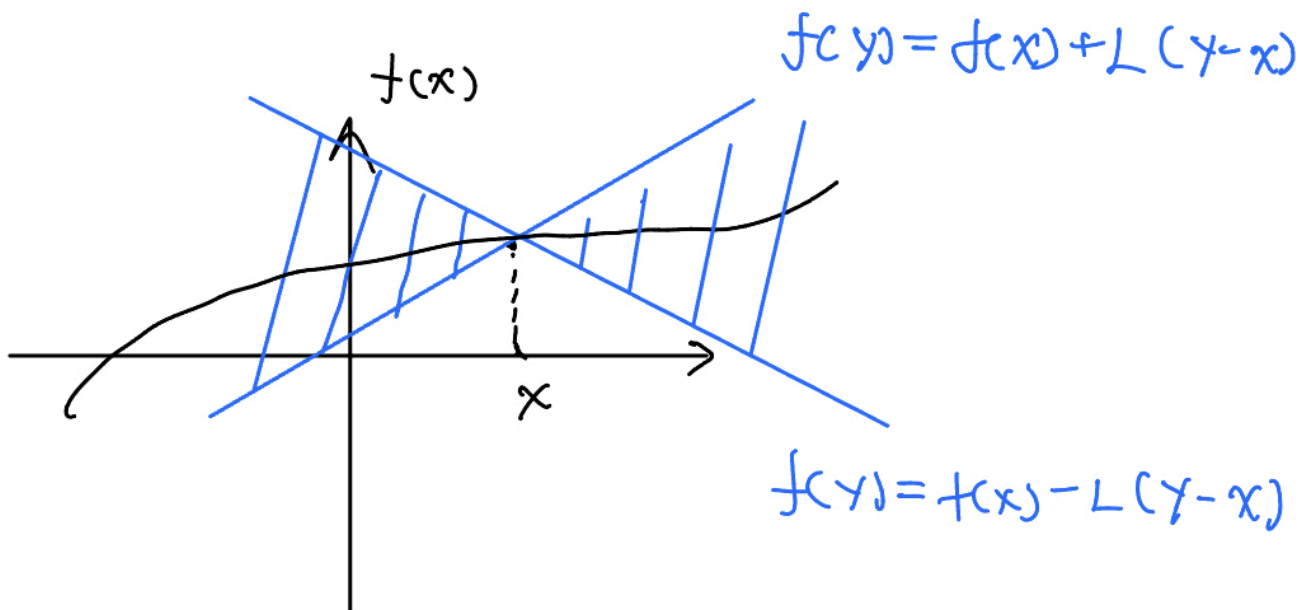
A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is L -Lipschitz, if for all $x, y \in \text{dom } f$,

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

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

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

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

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

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

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

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 $-L\mathbf{I}_n \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}_n$ for all $\mathbf{x} \in \mathbb{R}^n$, where \mathbf{I}_n is the $n \times n$ identity matrix. Namely, for all $\mathbf{x} \in \mathbb{R}^n$, $|\lambda_i(\nabla^2 f(\mathbf{x}))| \leq L$, 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 $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. Let $\varphi : [0, 1] \rightarrow \mathbb{R}$ be a function defined by

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

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

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

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.

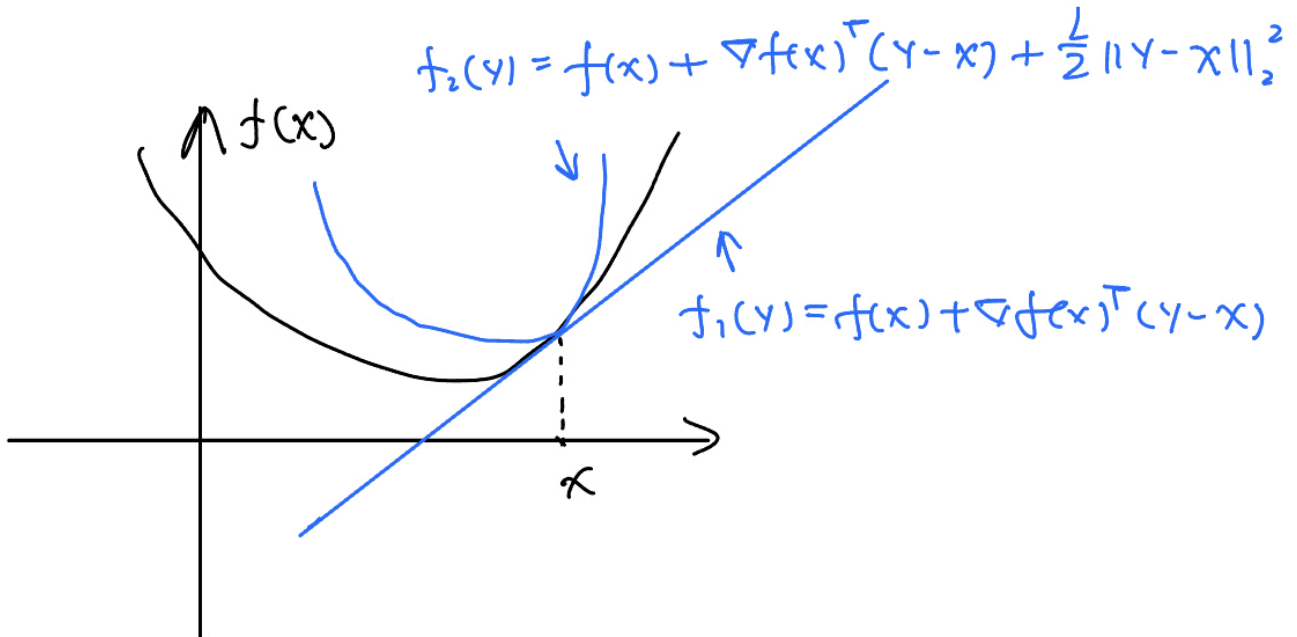
Lemma

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

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

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

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

3. Convergence analysis with smoothness

We first introduce a continuous version of the gradient descent instead, which is easier to analyze.

Definition (Gradient flow)

A *gradient flow* is a curve $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^n$ following the direction of steepest descent of a function. Given a smooth convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and a point $\hat{\mathbf{x}} \in \mathbb{R}^n$, the *gradient flow* of f with initial point $\hat{\mathbf{x}}$ is the solution to the following differential equation

$$\frac{d}{dt}\mathbf{x}_t = -\nabla f(\mathbf{x}_t), \quad \mathbf{x}_0 = \hat{\mathbf{x}}.$$

Here we use the notation $\mathbf{x}_t = \mathbf{x}(t)$ for convenience.

Applying the chain rule, $f(\mathbf{x}_t)$ is decreasing since

$$\frac{d}{dt}f(\mathbf{x}_t) = \left\langle \nabla f(\mathbf{x}_t), \frac{d}{dt}\mathbf{x}_t \right\rangle = -\langle \nabla f(\mathbf{x}_t), \nabla f(\mathbf{x}_t) \rangle \leq 0.$$

Now we can take the derivative

$$\begin{aligned} \frac{d}{dt}\|\mathbf{x}_t - \mathbf{x}^*\|^2 &= 2 \left\langle \mathbf{x}_t - \mathbf{x}^*, \frac{d}{dt}\mathbf{x}_t \right\rangle \\ &= -2 \langle \mathbf{x}_t - \mathbf{x}^*, \nabla f(\mathbf{x}_t) \rangle \\ &= 2 \langle \mathbf{x}^* - \mathbf{x}_t, \nabla f(\mathbf{x}_t) \rangle \\ &\leq 2(f(\mathbf{x}^*) - f(\mathbf{x}_t)) \end{aligned}$$

by convexity. Then integrating both sides, we obtain that

$$\|\mathbf{x}_T - \mathbf{x}^*\|^2 - \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \leq 2Tf(\mathbf{x}^*) - 2 \int_0^T f(\mathbf{x}_t) dt \leq 2T(f(\mathbf{x}^*) - f(\mathbf{x}_T)),$$

which further gives that

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

We compare the gradient descent with the gradient flow. Assume the gradient descent iterates with a fixed step size η and an initial point $\hat{\mathbf{x}}$, i.e.,

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \nabla f(\mathbf{x}_k), \quad \mathbf{x}_0 = \hat{\mathbf{x}}.$$

For the gradient descent,

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \int_{k\eta}^{(k+1)\eta} \nabla f(\mathbf{x}_k) dt.$$

For the gradient flow,

$$\mathbf{x}_{(k+1)\eta} = \mathbf{x}_{k\eta} - \int_{k\eta}^{(k+1)\eta} \nabla f(\mathbf{x}_t) dt.$$

Intuitively we know that, if ∇f does not change too fast, the gradient descent approximates the gradient flow.

Theorem

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex and L -smooth function. Choose $\eta \leq 1/L$, and let the gradient descent iterate with a fixed step size η . Then it holds that

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

Proof

Analogously to the gradient flow, we calculate $\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 - \|\mathbf{x}_k - \mathbf{x}^*\|^2$. Note that

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &= \langle \mathbf{x}_{k+1} - \mathbf{x}^*, \mathbf{x}_{k+1} - \mathbf{x}^* \rangle \\ &= \langle \mathbf{x}_k - \mathbf{x}^* - \eta \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* - \eta \nabla f(\mathbf{x}_k) \rangle \\ &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\eta \langle \mathbf{x}_k - \mathbf{x}^*, \nabla f(\mathbf{x}_k) \rangle + \eta^2 \|\nabla f(\mathbf{x}_k)\|^2. \end{aligned}$$

Similarly, it suffices to bound $-2\eta \langle \mathbf{x}_k - \mathbf{x}^*, \nabla f(\mathbf{x}_k) \rangle + \eta^2 \|\nabla f(\mathbf{x}_k)\|^2$ by $f(\mathbf{x}_k) - f(\mathbf{x}^*)$. Since $\eta \leq 1/L$, we have

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

by the descent lemma. In addition, the convexity of f gives that

$$-2\eta \langle \mathbf{x}_k - \mathbf{x}^*, \nabla f(\mathbf{x}_k) \rangle = 2\eta \langle \mathbf{x}^* - \mathbf{x}_k, \nabla f(\mathbf{x}_k) \rangle \leq 2\eta (f(\mathbf{x}^*) - f(\mathbf{x}_k)).$$

Thus, we obtain that

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 - \|\mathbf{x}_k - \mathbf{x}^*\|^2 &= -2\eta \langle \mathbf{x}_k - \mathbf{x}^*, \nabla f(\mathbf{x}_k) \rangle + \eta^2 \|\nabla f(\mathbf{x}_k)\|^2 \\ &\leq 2\eta (f(\mathbf{x}^*) - f(\mathbf{x}_k)) + 2\eta (f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})) \\ &= 2\eta (f(\mathbf{x}^*) - f(\mathbf{x}_{k+1})). \end{aligned}$$

Summing over both sides from 0 to $T - 1$, it implies that

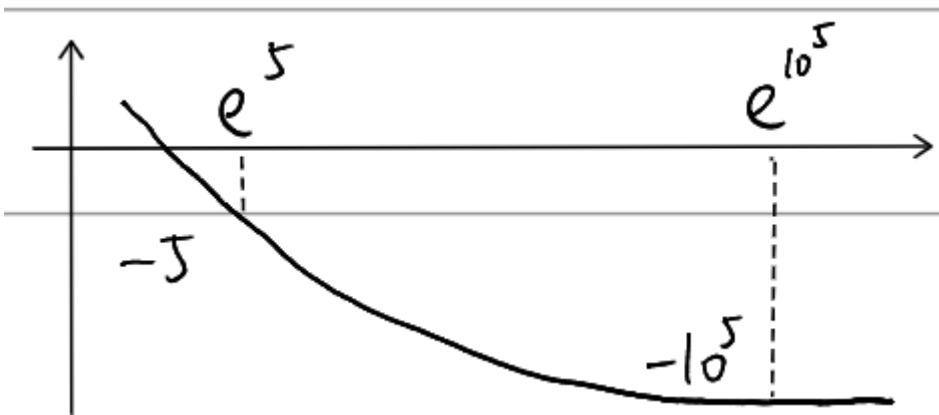
$$\|\mathbf{x}_T - \mathbf{x}^*\|^2 - \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \leq \sum_{k=1}^T 2\eta (f(\mathbf{x}^*) - f(\mathbf{x}_k)) \leq 2T\eta (f(\mathbf{x}^*) - f(\mathbf{x}_T)),$$

which is equivalent to

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

If we hope $|f(\mathbf{x}_T) - f^*| < \varepsilon$, we need to run the gradient descent $T = \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2\varepsilon\eta} = O(1/\varepsilon)$ steps. If the initial point \mathbf{x}_0 is far from \mathbf{x}^* , and ε is sufficiently small, the gradient descent is slow. Unfortunately, consider the following function

$$f(x) = \begin{cases} -\log x, & x < e^{10^5} \\ -10^5, & x \geq e^{10^5} \end{cases}.$$



This function is convex and 1-smooth. Hence $x_{k+1} = x_k + \eta/x_k \leq x_k + 1/x_k$ and the convergence rate of the gradient descent will be very small if $x_0 = 1$.

Question

Under which assumptions the gradient descent converges rapidly?

4. Strong convexity

Recall that, if we run the gradient descent for a quadratic function $f(x) = ax^2$ where $a \in \mathbb{R}_{>0}$, it gives that $x_{k+1} = (1 - 2a\eta)x_k$ and thus $f(x_k) = a(1 - 2a\eta)^{2k}x_0^2$. Clearly $f(x_k)$ converges to the optimal value 0 at an exponential rate.

We now introduce the following definition, which requires the function is a bit "better" than some quadratic function.

Definition (Strong convexity)

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is *strongly convex* with $\mu > 0$ if $f(x) - \frac{\mu}{2}\|x\|^2$ is convex.

There are some other forms of quadratic functions. Why don't we choose other functions such as $x^\top Qx$ or $\|x - y_0\|^2$ for some given y_0 ? In fact, these functions mentioned can just achieve a similar effect to $\frac{\mu}{2}\|x\|^2$. For example, $x^\top Qx$ is almost equivalent to $\lambda_{\max}(Q)x^\top x = \lambda_{\max}(Q)\|x\|^2$. In addition,

$$\|x - y_0\|^2 = \|x\|^2 + \underbrace{\|y_0\|^2}_{\text{constant}} - \underbrace{2\langle x, y_0 \rangle}_{\text{don't affect convexity}}.$$

Hence, all quadratic functions achieve similar effects to $\|x\|^2$.

Recall that, a function is convex iff its hessian matrix is positive semidefinite. The hessian matrix of $f(x) - \frac{\mu}{2}\|x\|^2$ is

$$\nabla^2 f(x) - \frac{\mu}{2}\nabla^2(\|x\|^2) = \nabla^2 f(x) - \frac{\mu}{2}\nabla^2(x^\top x) = \nabla^2 f(x) - \mu I.$$

Lemma

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a twice continuously differentiable function. Then f is μ -strongly convex iff $\nabla^2 f(x) \succeq \mu I_n$. Namely, for all $x \in \mathbb{R}^n$, $\lambda_{\min}(\nabla^2 f(x)) \geq \mu$.

We also have the following lemma similar to the first order condition for convexity and smoothness.

Lemma

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

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

Proof

Let $g(x) = f(x) - \frac{\mu}{2}\|x\|^2$. By the first order condition for convexity, $g(x)$ is convex iff for all $x, y \in \mathbb{R}^n$,

$$g(y) \geq g(x) + \langle \nabla g(x), y - x \rangle.$$

Note that $\nabla g(\mathbf{x}) = \nabla f(\mathbf{x}) - \mu\mathbf{x}$. So it gives that $g(\mathbf{x})$ is convex iff for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

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

The last inequality is equivalent to

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

Rearranging it, we obtain that

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

As a corollary, above lemma implies that $f(\mathbf{y}) > f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$ for any $\mathbf{x} \neq \mathbf{y}$. Hence, f is strictly convex.

Example

1. An affine functions $f(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$ can not be strongly convex since it is not strictly convex.
2. $f(x) = -\log x$ can not be strongly convex since $f'(x) = -\frac{1}{x}$ and $f''(x) = \frac{1}{x^2}$ and we can not find out such μ when $x \rightarrow 0$.
3. $f(x) = ax^2, a > 0$ is $2a$ -strongly convex.
4. $f(x) = x^4$ is not strongly convex since $f'(x) = 4x^3, f''(x) = 12x^2$ and we can not find such $\mu > 0$.
5. $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q}\mathbf{x}$ where $\mathbf{Q} \succ 0$ is strongly convex. Because $\nabla^2 f(\mathbf{x}) = 2\mathbf{Q}$, f is $2\lambda_{\min}(\mathbf{Q})$ -strongly convex.

Recall the property of monotone gradient for convex functions. We have a similar corollary.

Corollary

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

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \mu \|\mathbf{x} - \mathbf{y}\|^2.$$

5. Convergence analysis with strong convexity

We now establish the convergence of gradient descent with strong convexity. First consider the gradient flow again. By strong convexity, we can bound the derivative as follows

$$\begin{aligned}\frac{d}{dt} \|\mathbf{x}_t - \mathbf{x}^*\|^2 &= 2 \left\langle \mathbf{x}_t - \mathbf{x}^*, \frac{d}{dt} \mathbf{x}_t \right\rangle \\ &= -2 \left\langle \mathbf{x}_t - \mathbf{x}^*, \nabla f(\mathbf{x}_t) \right\rangle \\ &= -2 \left\langle \mathbf{x}_t - \mathbf{x}^*, \nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}^*) \right\rangle \\ &\leq -2\mu \|\mathbf{x}_t - \mathbf{x}^*\|^2.\end{aligned}$$

For a time-continuous non-negative process $u_t = u(t)$, if $\frac{d}{dt} u_t = -\alpha u_t$, then we have $u_t = u_0 \exp(-\alpha t)$. The same result holds if we replace the equality by an inequality.

Theorem (Gronwall's lemma)

For a time-continuous non-negative process $u_t = u(t)$, if $\frac{d}{dt} u_t \leq -\alpha_t u_t$, then we have

$$u_T \leq u_0 \exp\left(-\int_0^T \alpha_t dt\right).$$

Applying the Gronwall's lemma, we conclude $\|\mathbf{x}_T - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \exp(-2\mu T)$ immediately, which gives an exponential decay rate. Intuitively, as the discretization version of the gradient flow, the gradient descent for strongly convex functions should also follow the exponential decay.

Theorem

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is L -smooth and μ -strongly convex function. Choose $\eta \leq 1/L$, and let the gradient descent iterate with a fixed step size η . Then it holds that

$$\|\mathbf{x}_T - \mathbf{x}^*\|^2 \leq (1 - \mu\eta)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

Proof

By strong convexity,

$$\begin{aligned}
\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &= \langle (\mathbf{x}_k - \mathbf{x}^*) - \eta \nabla f(\mathbf{x}_k), (\mathbf{x}_k - \mathbf{x}^*) - \eta \nabla f(\mathbf{x}_k) \rangle \\
&= \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \eta^2 \|\nabla f(\mathbf{x}_k)\|^2 - 2\eta \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle \\
&\leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \eta^2 \|\nabla f(\mathbf{x}_k)\|^2 - 2\eta \left(f(\mathbf{x}_k) - f(\mathbf{x}^*) + \frac{\mu}{2} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \right) \\
&= (1 - \mu\eta) \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \eta^2 \|\nabla f(\mathbf{x}_k)\|^2 - 2\eta (f(\mathbf{x}_k) - f(\mathbf{x}^*)) \\
&\leq (1 - \mu\eta) \|\mathbf{x}_k - \mathbf{x}^*\|^2 + 2\eta (f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})) - 2\eta (f(\mathbf{x}_k) - f(\mathbf{x}^*)) \\
&= (1 - \mu\eta) \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\eta (f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*)) \\
&\leq (1 - \mu\eta) \|\mathbf{x}_k - \mathbf{x}^*\|^2
\end{aligned}$$

where the second inequality is due to the *descent lemma*

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

The function value also has an exponential decay. Since f is L -smooth, we have

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

which gives the following corollary.

Corollary

$$f(\mathbf{x}_T) - f^* \leq \frac{L}{2} (1 - \mu\eta)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$