

Intuitions for Optimization

pluskid

October 12, 2015

Contents

1	Projected Subgradient Descent for Convex Lipschitz Functions	1
1.1	Problem Setup	1
1.2	Algorithm and its Bounds	2
1.3	Intuitions and Analysis	2
1.3.1	Problem Assumptions	2
1.3.2	Convergence Analysis	2
1.4	Choosing Step Sizes	3
1.4.1	Constrained Optimization	4
2	Gradient Descent for Convex Smooth Functions	5
2.1	Sandwiching Smooth Convex Functions	5
2.2	Convergence Analysis and Tighter Sandwiching	7
3	Projected Gradient Descent for Convex Smooth Functions with Constraints	9

1 Projected Subgradient Descent for Convex Lipschitz Functions

1.1 Problem Setup

Problem 1. Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$, and $\mathcal{K} \subset \mathbb{R}^n$. Assume \mathcal{K} is compact and convex, included in a Euclidean ball of radius R :

$$\mathcal{K} \subset \mathcal{B}_2(0; R) \quad (1)$$

and f is convex and L -Lipschitz on \mathcal{K} :

$$|f(x) - f(y)| \leq L\|x - y\|_2, \quad x, y \in \mathcal{K} \quad (2)$$

Find a minimizer of f on \mathcal{K} :

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

Algorithm 1 Projected Subgradient Descent

```
randomly initialize  $x^0 \in \mathcal{K}$ 
for  $t \leftarrow 0, \dots, T-1$  do
     $y^{t+1} \leftarrow x^t - \eta_t g_t$ , where  $g_t \in \partial f(x^t)$ 
     $x^{t+1} \leftarrow \Pi_{\mathcal{K}}(y^{t+1})$ 
end for
```

1.2 Algorithm and its Bounds

Theorem 1. Running Algorithm 1 on Problem 1 for T iterations gives

$$\min_{0 \leq \tau \leq T} f(x^\tau) - f(x^*) \leq \frac{R^2 + G^2 \sum_{t=0}^{T-1} \eta_t^2}{\sum_{t=0}^{T-1} \eta_t} \quad (3)$$

In general, the optimal bound is around $O(GR/\sqrt{T})$ with stepsizes around $\eta_t \approx R/(G\sqrt{t})$. That means in order to get an approximate error of ε , we will need to run the algorithm for $O(1/\varepsilon^2)$ iterations.

1.3 Intuitions and Analysis

1.3.1 Problem Assumptions

We consider the unconstrained case first, i.e. $\mathcal{K} = \mathbb{R}^n$. Since f is convex know that x^* is a minimizer of f if and only if $0 \in \partial f(x^*)$. However, without making any extra assumptions, we have no idea of the behavior of f even in a small neighborhood of x^* . Consider for example $f(x) = C|x|$, with $C > 0$ a large constant. Assume we are currently very close to the optimal $x^* = 0$, say $x^t = \varepsilon, \varepsilon > 0$. f is differentiable at ε , so the only subgradient is $g_t = C$. Therefore,

$$x^{t+1} = x^t - \eta_t g_t = \varepsilon - \eta_t C < -\varepsilon, \quad \forall \eta_t > \frac{2\varepsilon}{C}$$

As we can see, unless the stepsize η_t is very tiny, we will overshoot, $f(x^{t+1}) > f(x^t)$. Moreover, if we use a constant stepsize $\eta_t = \eta$, then if $\eta > \varepsilon/C$, we will be jumping back and forth at ε and $\varepsilon - \eta C$ indefinitely.

In order to fix this, we need to make additional assumptions. We will see later in the case of smooth functions, the gradient changes continuously. So we know that at a local neighborhood of the optimal (gradient is 0), the gradient is also small. But here, we are working with non-differentiable functions, we will just assume f is L -Lipschitz.

Note f being L -Lipschitz in \mathcal{K} implies that $\|g\|_2 \leq L, \forall g \in \partial f(x), \forall x \in \mathcal{K}$. In order to avoid overshooting, we will have to move with tiny stepsizes.

1.3.2 Convergence Analysis

By the property of subgradient, we have

$$0 \leq f(x^t) - f(x^*) \leq g_t^\top (x^t - x^*) = -g_t^\top (x^* - x^t) \quad (4)$$

where the left hand side is due to the optimality of x^* . This inequality indicates that the vector $x^* - x^t$ is non-negatively correlated with $-g_t$, the direction of a subgradient descent.

Unlike in the case of differentiable functions, in which we can guarantee that the function values decreases when moving along the direction of negative gradient with a small enough step size; here we do not know much about the function values. But as we can see from Figure 1, when the angle between $-g_t$ and $x^* - x^t$ is greater than or equal to $\pi/2$, we will move away from x^* with any positive step size. However, if $f(x^t) - f(x^*) > 0$, i.e. we are not already at the optimal, the angle is strictly less than $\pi/2$, so if we move with a small enough step size, we will get closer to x^* . Algebraically,

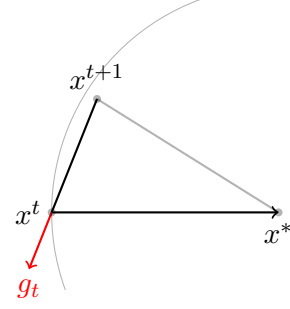


Figure 1: Demonstration of subgradient descent.

$$\begin{aligned} \|x^{t+1} - x^*\|^2 &= \|x^{t+1} - x^t + x^t - x^*\|^2 = \eta_t^2 \|g_t\|^2 + \|x^t - x^*\|^2 - 2\eta_t g_t^\top (x^t - x^*) \\ &\leq \eta_t^2 \|g_t\|^2 + \|x^t - x^*\|^2 - 2\eta_t (f(x^t) - f(x^*)) \end{aligned} \quad (5)$$

Since $\|g_t\| \leq L$ by our assumption, the **red term** decays quadratically, while the **blue term** only decays linearly. So when η_t is small enough, we will have $\|x^{t+1} - x^*\|^2 \leq \|x^t - x^*\|^2$. Furthermore, the progress we make by moving towards x^* is characterized by $f(x^t) - f(x^*)$. So if we are still far away from the optimal, we will be making quite a lot progress in each step. On the other hand, when $f(x^t) - f(x^*)$ is small, our progress might be small, but at that point we are already close to the optimal function value $f(x^*)$.

Actually, to get a bound on the algorithm, we can just sum up the previous inequality for all $t = 0, \dots, T-1$,

$$\begin{aligned} 0 \leq \|x^T - x^*\|^2 &\leq \|x^0 - x^*\|^2 + \sum_{t=0}^{T-1} \eta_t^2 \|g_t\|^2 - 2 \sum_{t=0}^{T-1} \eta_t (f(x^t) - f(x^*)) \\ &\leq R^2 + G^2 \sum_{t=0}^{T-1} \eta_t^2 - 2 \left(\min_{0 \leq \tau \leq T} f(x^\tau) - f(x^*) \right) \sum_{t=0}^{T-1} \eta_t \end{aligned}$$

It then implies

$$\min_{0 \leq \tau \leq T} f(x^\tau) - f(x^*) \leq \frac{R^2 + G^2 \sum_{t=0}^{T-1} \eta_t^2}{\sum_{t=0}^{T-1} \eta_t} \quad (6)$$

which proved Theorem 1 for the case of unconstrained optimization ($\mathcal{K} = \mathbb{R}^n$).

1.4 Choosing Step Sizes

Note (3) holds for any choices of step sizes η_t (though some of them will give completely trivial bounds). So we could actually optimize the right hand side to get an “optimal” bound. To make the problem easier, we choose a fixed stepsize $\eta_t = \eta$ for $t = 0, \dots, T-1$. So the right hand side becomes

$$\frac{R^2 + G^2 T \eta^2}{T \eta} = \frac{R^2}{T \eta} + G^2 \eta \geq \frac{2GR}{\sqrt{T}} \quad (7)$$

where the inequality holds with equality when

$$\eta = \frac{R}{G\sqrt{T}} \quad (8)$$

Note the choice of step size depends on several factors:

G : As we described in Section 1.3.1, large G will force us to be careful and move with small step sizes. Our intuition is consistent here.

R : In our analysis, we only use R to bound $\|x^0 - x^*\|$. When R is large, we want to use large step size, otherwise we might never reach the optimal in the given time budget T . Generally when x^* is unknown, R can be bounded by the size of \mathcal{K} for the case of constrained optimization.

T : The inverse dependency on T can be interpreted as: when having a large time budget, we can be a little bit more careful and move slowly.

However, in general, the fact that the step size depends on the total number of iterations is strange. That means if I want to compute more iterations, I will have to start over again and use a different step size if I want to bound the performance with formula.

In general, we will prefer to use a decaying learning rate. Specifically, as long as $\sum_t \eta_t \rightarrow \infty$ and $\sum_t \eta_t^2$ is bounded or approaches infinity at a slower rate than $\sum_t \eta_t$, (3) will give a reasonable bound. For example, take $\eta_t = R/(G\sqrt{t+1})$, since

$$\begin{aligned} \sum_{t=0}^{T-1} \frac{1}{t+1} &\leq 1 + \int_1^T \frac{1}{x} dx = 1 + \log T \\ \sum_{t=0}^{T-1} \frac{1}{\sqrt{t+1}} &\geq \int_1^{T+1} \frac{1}{\sqrt{x}} dx = 2\sqrt{T+1} - 2 \end{aligned}$$

Plug-in to (3), we get

$$\min_{0 \leq \tau \leq T} f(x^\tau) - f(x^*) \leq GR \frac{1 + \log T}{2\sqrt{T+1} - 2} \lesssim \frac{GR \log T}{\sqrt{T}} \quad (9)$$

Comparing with the optimal bound we get with a fixed step size in (7), we lose a factor of $\log T$, but our step size does not depend on the total number of iterations any more.

1.4.1 Constrained Optimization

In the constrained case, we have an extra projection step. However, the same analysis naturally goes through because projection into convex set is a *contraction*. Specifically, we have the following lemma.

Lemma 1. Let $\mathcal{K} \subset \mathbb{R}^n$ be a closed convex set, let $x \in \mathcal{K}$ and $y \in \mathbb{R}^n$. Then

$$(x - \Pi_{\mathcal{K}}(y))^\top (y - \Pi_{\mathcal{K}}(y)) \leq 0 \quad (10)$$

which also implies

$$\|x - \Pi_{\mathcal{K}}(y)\|^2 + \|y - \Pi_{\mathcal{K}}(y)\|^2 \leq \|y - x\|^2 \quad (11)$$

This lemma could be proved with *supporting hyperplane theorem* of convex sets. (11) implies

$$\|x - \Pi_{\mathcal{K}}\|^2 \leq \|y - x\|^2$$

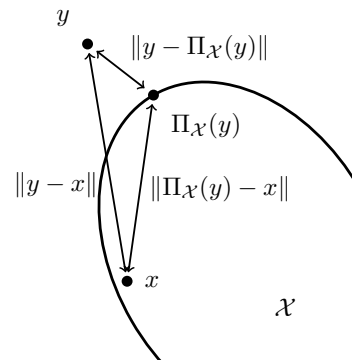


Figure 2: Illustration of convex projection. Figure source: Sébastien Bubeck, *Theory of Convex Optimization for Machine Learning*.

So if we go back to our analysis in Section 1.3.2, the only thing we need to modify is (5). Since $x^* \in \mathcal{K}$,

$$\|x^{t+1} - x^*\|^2 = \|\Pi_{\mathcal{K}}(y^{t+1}) - x^*\|^2 \leq \|y^{t+1} - x^*\|^2 = \eta_t^2 \|g_t\|^2 + \|x^t - x^*\|^2 - 2\eta_t g_t^\top (x^t - x^*)$$

and the rest of the analysis follows as before. So for constrained optimization, we get the same convergence rate as the unconstrained case.

2 Gradient Descent for Convex Smooth Functions

In this section, we look at smooth functions. Specifically, f is differentiable, and its gradient ∇f is Lipschitz continuous. By adding those assumptions, we know better about f than in the simple Lipschitz case. For example, we know two nearby points should have similar gradients. In particular, if x is close to the optimal x^* , since $\nabla f(x^*) = 0$, we know that $\nabla f(x)$ must be small (close to zero).

Definition 1 (β -smooth functions). A differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is β -smooth if its gradient ∇f is β -Lipschitz, that is

$$\|\nabla f(x) - \nabla f(y)\| \leq \beta \|x - y\|, \quad \forall x, y \in \mathbb{R}^d \quad (12)$$

Problem 2. Given a convex and β -smooth function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Find a minimizer of f .

Theorem 2. Solving Problem 2 using gradient descent with step size $\eta_t = 1/\beta$ for T iterations gives

$$f(x^T) - f(x^*) \leq \frac{2\beta \|x_1 - x^*\|^2}{T + 3} \quad (13)$$

That means to get an approximation error of ε , we will need to run the algorithm for $O(1/\varepsilon)$ iterations.

Note we get much faster convergence rate than Problem 1 by making more assumptions on f . The parameter β control the smoothness of f : smaller β means the gradient of f changes more slowly, and $\eta_t = 1/\beta$ means we could make aggressive movement in each iteration.

2.1 Sandwiching Smooth Convex Functions

In the case of convex Lipschitz function, we use the property of subgradient to lower bound f . $\forall x, y \in \mathbb{R}^n$ and $\forall g \in \partial f(x)$

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

this leads to (4). And in the proof, we use this to lower bound $f(x^*)$, making sure that $f(x^t) - f(x^*)$ is controlled, i.e. we are not too far away from the optimal.

In the scenario of β -smooth function, we can get the same lower bound, because $\nabla f(x)$ is always the unique subgradient at any point x . Moreover, the β -smoothness gives us an upper bound:

$$f(y) \leq f(x) + \nabla f(x)^\top (y - x) + \frac{\beta}{2} \|y - x\|^2$$

and this could be used to get a lower bound on the decrement $f(x^t) - f(x^{t+1})$ at each iteration. See Fig. 3 for an illustration. We state the conclusions below formally.

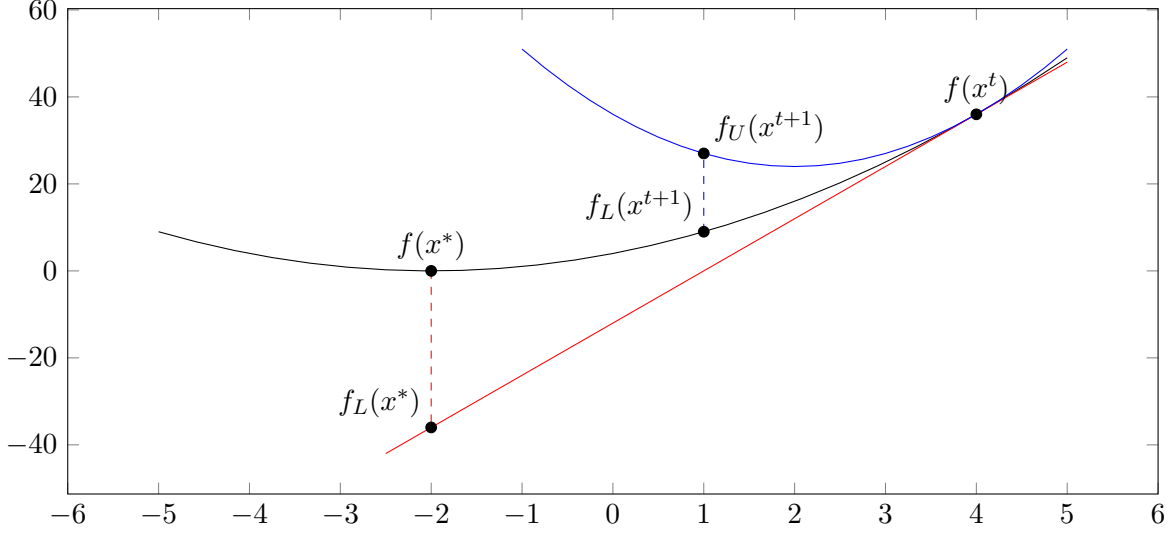


Figure 3: Illustration of a convex β -smooth function $f(x)$, with its lower bound $f^L(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle$ and upper bound $f^U(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + 0.5\beta\|x - x_0\|^2$. The lower bound makes sure $f(x^t)$ is not too far away from $f(x^*)$, while the upper bound makes sure some progress $f(x^t) - f(x^{t+1})$ are made in each iteration.

Lemma 2. Assume $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is β -smooth, then $\forall x, y \in \mathbb{R}^n$

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

Proof. By the fundamental theorem for line integrals,

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

Plugin $f(y) - f(x)$,

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

□

If we further know that f is convex, combining the lower bound from the first order condition of convexity, we get both lower bound and upper bound

$$0 \leq f(y) - f(x) - \langle \nabla f(x), y - x \rangle \leq \frac{\beta}{2} \|y - x\|^2 \quad (14)$$

See again Fig. 3 for an illustration.

2.2 Convergence Analysis and Tighter Sandwiching

To analyze the convergence, let $\Delta_t = f(x^t) - f(x^*)$ be the suboptimality gap at x^t . Let $y = x^*$, and $x = x^t$ in (14), and using the lower bound, we can upper bound Δ_t by

$$\Delta_t = f(x^t) - f(x^*) \leq -\langle \nabla f(x^t), x^* - x^t \rangle \leq \|\nabla f(x^t)\| \|x^* - x^t\| \leq R \|\nabla f(x^t)\| \quad (15)$$

where we define

$$R = \max_{1 \leq t \leq T} \|x^* - x^t\| \quad (16)$$

On the other hand, using the upper bound in (14) by letting $y = x^{t+1}$ and $x = x^t$, we could lower bound $\Delta_t - \Delta_{t+1}$ by

$$\begin{aligned} \Delta_t - \Delta_{t+1} &= f(x^t) - f(x^{t+1}) \geq -\langle \nabla f(x^t), x^{t+1} - x^t \rangle - \frac{\beta}{2} \|x^{t+1} - x^t\|^2 \\ &= \left(\eta_t - \frac{\beta \eta_t^2}{2} \right) \|\nabla f(x^t)\|^2 \end{aligned}$$

Naturally, we want to maximize the lower bound, so the step size is chosen to be $\eta_t = 1/\beta$. Combining with (15), we get

$$\Delta_t - \Delta_{t+1} \geq \frac{1}{2\beta} \|\nabla f(x^t)\|^2 \geq \frac{1}{2\beta R^2} \Delta_t^2 \quad (17)$$

Note the right hand side is non-negative, so $\Delta_t \geq \Delta_{t+1}$. To solve this recursion, divide both side by $\Delta_t \Delta_{t+1}$:

$$\frac{1}{\Delta_{t+1}} - \frac{1}{\Delta_t} \geq \frac{1}{2\beta R^2} \frac{\Delta_t}{\Delta_{t+1}} \geq \frac{1}{2\beta R^2} \quad (18)$$

Sum the recursion for $t = 2, \dots, T$, we get

$$\frac{1}{\Delta_T} \geq \frac{T-1}{2\beta R^2} + \frac{1}{\Delta_1} \geq \frac{T+3}{2\beta R^2}$$

where the last inequality is because Δ_1 can be controlled by the upper bound in (14). Let $x = x^*$ and $y = x^1$, notice $\nabla f(x^*) = 0$,

$$\Delta_1 = f(x^1) - f(x^*) \leq \frac{\beta}{2} \|x^1 - x^*\|^2 \leq \frac{\beta R^2}{2}$$

At this point, we almost proved Theorem 2, except that we have to control R . In the following, we will show that $\|x^t - x^*\|$ is actually decreasing at each iteration, and bound R by $\|x^1 - x^*\|$. Recall in the case of subgradient descent, we use the linear lower bound of the function by the subgradient to construct the inequality in (5). That inequality shows that when η_t is small enough, because the quadratic term decays faster, we get $\|x^{t+1} - x^*\|^2 \leq \|x^t - x^*\|^2$. However, in the case here, since $\eta_t = 1/\beta$, especially when β is small, we could be moving with very large step size. So the argument is no longer useful here.

To properly bound R , we will need to get a better lower bound of f than in (14). Actually, combining convexity and β -smoothness, the lower bound in (14) could be improved. Consider the extreme case when $f(x)$ is a linear function, then the lower bound is actually tight. In this case, we also have $\nabla f(x) = \nabla f(y)$. However, if $f(x)$ is not linear, $\nabla f(x) \neq \nabla f(y)$, we might observe a non-zero gap between $f(x)$ and its linear lower bound. It is also intuitive that the gap might be larger when the gradient $\nabla f(y)$ changed a lot from $\nabla f(x)$, so we are thinking about getting a better lower bound using the quantity $\|\nabla f(x) - \nabla f(y)\|$.

Lemma 3. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex and β -smooth function, then $\forall x, y \in \mathbb{R}^n$

$$\frac{1}{2\beta} \|\nabla f(x) - \nabla f(y)\|^2 \leq f(y) - f(x) - \langle \nabla f(x), y - x \rangle \leq \frac{\beta}{2} \|y - x\|^2 \quad (19)$$

Proof. In order to invite both $\nabla f(x)$ and $\nabla f(y)$ into play, we consider a third point $z \in \mathbb{R}^n$, and approximate $f(z)$ from below by $\nabla f(y)$ and from above by $\nabla f(x)$, respectively. Using (14)

$$\begin{aligned} f(z) - f(x) - \langle \nabla f(x), z - x \rangle &\geq 0 \\ f(z) - f(y) - \langle \nabla f(y), z - y \rangle &\leq \frac{\beta}{2} \|z - y\|^2 \end{aligned}$$

Multiply the first inequality by -1 and sum the two inequalities, we get

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

Re-write the inequality by moving the quantity we want to lower bound to the right,

$$\langle \nabla f(x), z - y \rangle - \langle \nabla f(y), z - y \rangle - \frac{\beta}{2} \|z - y\|^2 \leq f(y) - f(x) - \langle \nabla f(x), y - x \rangle$$

Inspecting the left hand side, if we let $z = y + \alpha(\nabla f(y) - \nabla f(x))$ for any $\alpha \in \mathbb{R}$, we get

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

Since the lower bound is a quadratic function in α , we can maximize the lower bound by taking $\alpha = 1/\beta$. And the conclusion follows. \square

With the improved lower bound of f in (19), we can now bound R by showing that

$$\begin{aligned} \|x^{t+1} - x^*\|^2 &= \|x^t - x^*\|^2 + \frac{1}{\beta^2} \|\nabla f(x^t)\|^2 - \frac{2}{\beta} \langle \nabla f(x^t), x^t - x^* \rangle \\ &\leq \|x^t - x^*\|^2 + \frac{1}{\beta^2} \|\nabla f(x^t)\|^2 - \frac{2}{\beta} \left(f(x^t) - f(x^*) + \frac{1}{2\beta} \|\nabla f(x^t) - \nabla f(x^*)\|^2 \right) \\ &\leq \|x^t - x^*\|^2 + \frac{1}{\beta^2} \|\nabla f(x^t)\|^2 - \frac{2}{\beta} \times \frac{1}{2\beta} \|\nabla f(x^t)\|^2 \\ &= \|x^t - x^*\|^2 \end{aligned}$$

Therefore, $R \leq \|x^1 - x^*\|$, which conclude the proof of Theorem 2. Actually, in the proof above, we only need a property of the gradient called *co-coercivity*.

Definition 2 (Co-coercive mapping). A mapping $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is co-coercive with parameter C if $\forall x, y \in \mathbb{R}^d$,

$$\langle F(x) - F(y), x - y \rangle \geq C \|F(x) - F(y)\|^2 \quad (20)$$

By switching the role of x and y in the lower bound of f in (19), it is easy to show the following property on the co-coercivity of $\nabla f(x)$.

Lemma 4 (Co-coercivity of gradient of convex smooth functions). Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex β -smooth function, then $\nabla f(x)$ is co-coercive mapping with parameter $1/\beta$. That is

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{1}{\beta} \|\nabla f(x) - \nabla f(y)\|^2 \quad (21)$$

We will re-visit some variants of this property in our later adventures of convergence analysis.

3 Projected Gradient Descent for Convex Smooth Functions with Constraints