

SDSC6015 Stochastic Optimization for Machine Learning

Lu Yu

Department of Data Science, City University of Hong Kong

September 18, 2025

Gradient Descent and Subgradient Method

Convex Optimization Problems

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

where

- ▶ f is a **convex** function
- ▶ \mathbb{R}^d is convex
- ▶ x^* is the minimizer of function f :

$$x^* = \arg \min_{x \in \mathbb{R}^d} f(x)$$

Recap: Gradient Descent

Update rule for gradient descent:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_{k+1} \nabla f(\mathbf{x}_k)$$

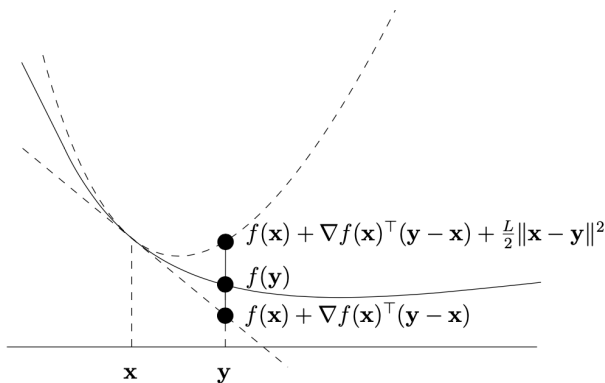
- ▶ \mathbf{x}_k : current point (parameters or variables).
- ▶ η_k : step size (learning rate), a positive scalar determining how far we move in the gradient direction.
- ▶ \mathbf{x}_{k+1} : next point after the update.

Recap: Smooth Functions

Definition

Let $f : \mathbf{dom}(f) \rightarrow \mathbb{R}$ be differentiable, $X \subseteq \mathbf{dom}(f)$, $L > 0$. f is called smooth (with parameter L) over X if

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



Subgradient

Recall: for convex and differentiable $f : \mathbb{R}^d \rightarrow \mathbb{R}$

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

Definition

A **subgradient** of a convex function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ at \mathbf{x} is any $g \in \mathbb{R}^d$ such that

$$f(\mathbf{y}) \geq f(\mathbf{x}) + g^\top (\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{y}.$$

Set of all subgradients of f is called the **subdifferential**:

$$\partial f(\mathbf{x}) = \{g \in \mathbb{R}^d : g \text{ is a subgradient of } f \text{ at } \mathbf{x}\}.$$

Subgradient Method

Now consider convex function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ convex, but **not necessarily differentiable**.

Subgradient method: like gradient descent, but replacing gradients with subgradients

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_{k+1} g_k$$

- ▶ \mathbf{x}_k : current point
- ▶ $g_k \in \nabla f(\mathbf{x}_k)$: any subgradient of f at \mathbf{x}_k
- ▶ $\eta_k > 0$: step size
- ▶ \mathbf{x}_{k+1} : next point after the update.

Caveat: Subgradient method is not necessarily a descent method!

e.g. $f(x) = |x|$ (non-smoothness causes oscillation)

Summary

f	Algorithm	Convergence	# Iterations
Convex L -Lipschitz	GD	$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq \frac{RL}{\sqrt{T}}$	$\frac{R^2 L^2}{\epsilon^2}$
Convex L -Smooth	GD	$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq \frac{R^2 L}{2T}$	$\frac{R^2 L}{2\epsilon}$
Convex L -Lipschitz	Subgrad	$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq \frac{LR}{\sqrt{T}}$	$\frac{R^2 L^2}{\epsilon^2}$

- ▶ Time horizon $T > 0$ is given
- ▶ $R := \|\mathbf{x}_0 - \mathbf{x}^*\|$
- ▶ $\mathbf{x}_{\text{best}}^{(T)} := \arg \min_{i=0,1,\dots,T} f(\mathbf{x}_i)$.

Thus, the subgradient method has convergence rate $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$
...compare this to $\mathcal{O}\left(\frac{1}{\epsilon}\right)$ rate of gradient descent

Faster Gradient Descent

Can we go even faster?

- ▶ So far: Error decreases with $1/\sqrt{T}$, or $1/T$...
- ▶ Could the error decrease exponentially in T , i.e.,

$$e^{-cT}, \quad \text{for some } c > 0,$$

rather than following a polynomial decay¹?

¹Polynomial decay refers to a rate of decrease that follows a power law, meaning a quantity $f(t)$ diminishes over time t as $f(t) = \mathcal{O}(t^{-\alpha})$ for some exponent $\alpha \geq 0$

Can we go even faster?

- Consider $f(x) = x^2$: step size $\eta = \frac{1}{2}$ (f is $L = 2$ - smooth)

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

Converge in one step!

- Same $f(x) = x^2$: step size $\eta = \frac{1}{4}$

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

so

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

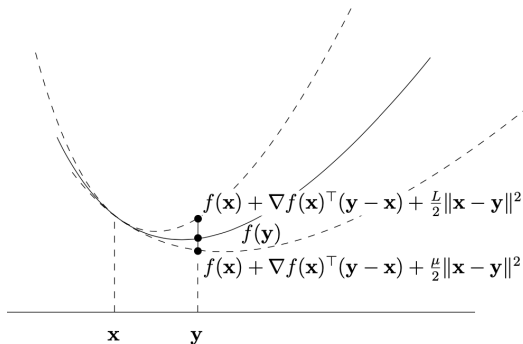
Exponential in t !

Strongly Convex Functions

Definition

Let $f : \mathbf{dom}(f) \rightarrow \mathbb{R}$ be a differentiable function, $\mathbf{dom}(f)$ is a convex set and a constant $\mu > 0$. f is called **strongly convex with parameter μ** if

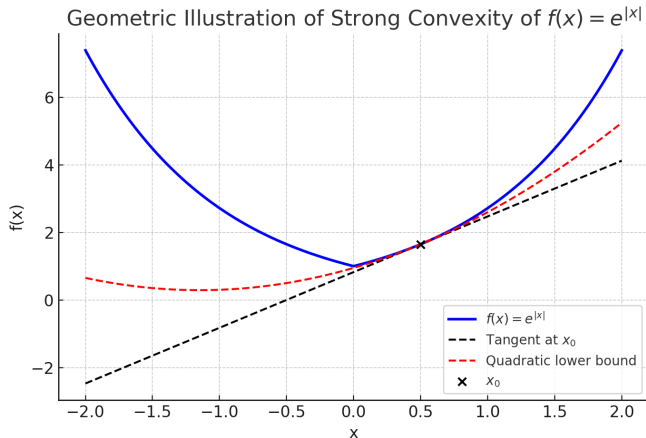
$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbf{dom}(f).$$



For any x , the graph of f lies above a tangent paraboloid at $(x, f(x))$.

Strongly Convex Functions

Example: $f(x) = e^{|x|}$ is strongly convex with parameter $\mu = 1$.



Strongly Convex Functions

Strong convexity enforces that the function grows at least a **quadratic rate** as you move away from the minimum.

This ensures:

- ▶ A unique minimizer – there is only one global minimum.

Lemma 1

If f is strongly convex, then f is strictly convex and has a unique global minimum.

- ▶ Faster optimization convergence – see later!

GD on Smooth and Strongly Convex Functions

► Aim to show $\lim_{t \rightarrow \infty} \mathbf{x}_t = \mathbf{x}^*$.

Step 1: Vanilla Analysis from the last lecture:

$$\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) = \frac{\eta}{2} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\eta} (\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2)$$

Step 2: Now use **stronger** lower bound on the left hand side, coming from **strong convexity**

$$\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \geq f(\mathbf{x}_t) - f(\mathbf{x}^*) + \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2$$

Step 3: Putting it together and rearranging gives

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

GD on Smooth and Strongly Convex Functions

$$\| \mathbf{x}_{t+1} - \mathbf{x}^* \|^2 \leq \underbrace{2\eta(f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \eta^2 \|\nabla f(\mathbf{x}_t)\|^2}_{\text{"noise"}} + (1 - \mu\eta) \|\mathbf{x}_t - \mathbf{x}^*\|^2$$

- ▶ Choose $\eta < \frac{1}{\mu}$
- ▶ Squared distance to \mathbf{x}^* goes down by a constant factor $(1 - \mu\eta)$, up to some “noise”

GD on Smooth and Strongly Convex Functions

Theorem 1

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable with a global minimum \mathbf{x}^* ; suppose that f is L -smooth and strongly convex with parameter μ . Choosing

$$\eta = \frac{1}{L},$$

gradient descent with arbitrary \mathbf{x}_0 satisfies

- Squared distances to \mathbf{x}^* are geometrically decreasing:

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

- The absolute error after T iterations is exponentially with T :

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad T \geq 0.$$

GD on Smooth and Strongly Convex Functions

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \underbrace{2\eta(f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \eta^2 \|\nabla f(\mathbf{x}_t)\|^2}_{\text{"noise"}} + (1 - \mu\eta) \|\mathbf{x}_t - \mathbf{x}^*\|^2$$

Proof of (i).

Bounding the noise Note that $\eta = 1/L$

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

Employing Lemma 3 (sufficient decrease) from the last lecture:

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

Hence, the noise is nonpositive

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

GD on Smooth and Strongly Convex Functions

Then, we get (i):

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

Proof of (ii). From (i):

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

Smoothness together with $\nabla f(\mathbf{x}^*) = \mathbf{0}$:

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

Putting it together:

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

GD on Smooth and Strongly Convex Functions

$$R^2 := \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$

$$T \geq \frac{L}{\mu} \ln \left(\frac{R^2 L}{2\varepsilon} \right) \Rightarrow \text{error} \leq \frac{L}{2} \left(1 - \frac{\mu}{L} \right)^T R^2 \leq \varepsilon$$

Conclusion: To reach absolute error at most ε , we only need $\mathcal{O}(\log \frac{1}{\varepsilon})$ iterations, e.g.

- ▶ $\frac{L}{\mu} \ln(50 \cdot R^2 L)$ iterations for error $\varepsilon = 0.01$
- ▶ ... as opposed to $50 \cdot R^2 L$ in the smooth case

In Practice:

What if we don't know the smoothness parameter L ?

Questions?

Summary

f	Algorithm	Convergence	# Iterations
Convex L -Lipschitz	GD	$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq \frac{RL}{\sqrt{T}}$	$\frac{R^2 L^2}{\varepsilon^2}$
Convex L -Smooth	GD	$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq \frac{R^2 L}{2T}$	$\frac{R^2 L}{2\varepsilon}$
μ -Strongly Convex L -Smooth	GD	$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq \frac{RL}{2}(1 - \frac{\mu}{L})^T$	$\frac{L}{\mu} \ln \left(\frac{R^2 L}{2\varepsilon} \right)$
Convex L -Lipschitz	Subgrad	$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq \frac{LR}{\sqrt{T}}$	$\frac{R^2 L^2}{\varepsilon^2}$

Faster Subgradient Method

Given a convex and L -Lipschitz function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ (f is not necessarily differentiable), the subgradient method with $R = \|\mathbf{x}_0 - \mathbf{x}^*\|$ satisfies

$$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq \frac{LR}{\sqrt{T}}$$

Question: Can we improve the convergence rate?

Optimality of First-order Methods

With all the convergence rates we have seen so far, a very natural question to ask is if these rates are the best possible or not. Surprisingly, the rate can indeed not be improved in general.

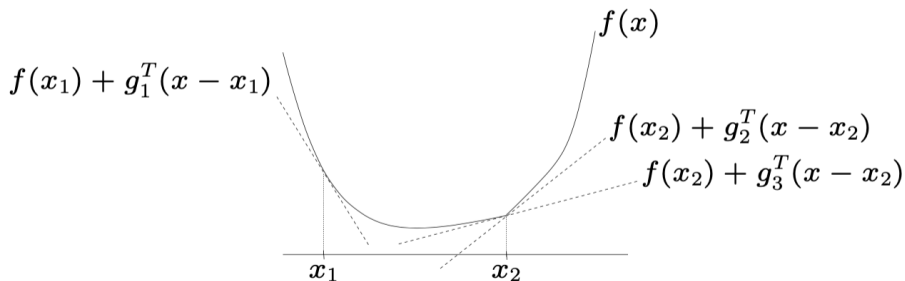
Theorem 2

For any $T \leq d - 1$ and starting point \mathbf{x}_0 , there is a function f in the problem class of **L -Lipschitz** functions over \mathbb{R}^d , such that any (sub)gradient method has an objective error at least

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \geq \frac{LR}{2(1 + \sqrt{T + 1})}$$

Smooth (non-differentiable) functions?

- ▶ They don't exist: A non-differentiable function cannot be smooth (e.g. $f(x) = |x|$).
- ▶ Can we still improve over $\mathcal{O}(1/\varepsilon^2)$ steps for Lipschitz functions?
- ▶ Yes, if we also require **strong convexity**.



Strongly Convex Functions

Straightforward generalization to the non-differentiable case:

Definition.

Let $f : \mathbf{dom}(f) \rightarrow \mathbb{R}$ be convex, and $\mu > 0$. Function f is called **strongly convex** with parameter μ if

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top(\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbf{dom}(f), \forall \mathbf{g} \in \partial f(\mathbf{x})$$

Strongly Convex Functions

Does **strong convexity alone** guarantee a fast convergence?

NO!

Reason: The subgradient method does not fully exploit strong convexity, since subgradients

- ▶ are not always well-aligned with the optimal descent direction
- ▶ lack curvature information (e.g. Hessian information)

Strongly Convex Functions

For fast convergence, we consider **additional** assumptions beyond strongly convexity.

- ▶ Smoothness? - Not an option in the non-differentiable case.
- ▶ Instead: assume that all subgradients \mathbf{g}_t during the algorithm are bounded in norm.
 - This is not always equivalent to Lipschitz (see notes)
 - Over \mathbb{R}^d , strong convexity and Lipschitz continuity contradict each other (see notes)

Subgradient Method on Strongly Convex Functions

Theorem 3

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be strongly convex with parameter $\mu > 0$ and let \mathbf{x}^* be the unique global minimum of f . With decreasing step size

$$\eta_t = \frac{2}{\mu(t+1)}, \quad t \geq 0,$$

subgradient descent yields

$$\underbrace{f\left(\frac{2}{T(T+1)} \sum_{t=1}^T t \cdot \mathbf{x}_t\right)}_{\text{convex combination of iterates}} - f(\mathbf{x}^*) \leq \frac{2B^2}{\mu(T+1)},$$

where $B = \max_{t=1}^T \|\mathbf{g}_t\|$.

Subgradient Method on Strongly Convex Functions

Proof.

Vanilla analysis: for any $\mathbf{g}_t \in \partial f(\mathbf{x}_t)$

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) = \frac{\eta_t}{2} \|\mathbf{g}_t\|^2 + \frac{1}{2\eta_t} (\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2).$$

Lower bound from **strong convexity**:

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) \geq f(\mathbf{x}_t) - f(\mathbf{x}^*) + \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2.$$

Putting it together with $\|\mathbf{g}_t\|^2 \leq B^2$:

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

Summing over $t = 1, \dots, T$: we used to have telescoping with $\eta_t = \eta, \mu = 0$ in the previous

Subgradient Method on Strongly Convex Functions

So far we have:

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

Plug in $\eta_t^{-1} = \mu(1+t)/2$ multiply with t on both sides:

$$\begin{aligned} t \cdot (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{B^2 t}{\mu(t+1)} + \frac{\mu}{4} \left(t(t-1) \|\mathbf{x}_t - \mathbf{x}^*\|^2 - t(t+1) \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right) \\ &\leq \frac{B^2}{\mu} + \frac{\mu}{4} \left(t(t-1) \|\mathbf{x}_t - \mathbf{x}^*\|^2 - t(t+1) \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right) \end{aligned}$$

Subgradient Method on Strongly Convex Functions

Now we get telescoping

$$\sum_{t=1}^T t \cdot (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{TB^2}{\mu} + \frac{\mu}{4} \left(0 - T(T+1) \|\mathbf{x}_{T+1} - \mathbf{x}^*\|^2 \right) \leq \frac{TB^2}{\mu}.$$

Since

$$\frac{2}{T(T+1)} \sum_{t=1}^T t = 1,$$

Jensen's inequality yields

$$f\left(\frac{2}{T(T+1)} \sum_{t=1}^T t \cdot \mathbf{x}_t\right) - f(\mathbf{x}^*) \leq \frac{2}{T(T+1)} \sum_{t=1}^T t \cdot (f(\mathbf{x}_t) - f(\mathbf{x}^*)).$$

Subgradient Method on Strongly Convex Functions

Putting together

$$f\left(\frac{2}{T(T+1)} \sum_{t=1}^T t \cdot \mathbf{x}_t\right) - f(\mathbf{x}^*) \leq \frac{2B^2}{\mu(T+1)}.$$

- ▶ **Weighted average** of iterates achieves the bound (later iterates have more weight)
- ▶ Bound is independent of initial distance $\|\mathbf{x}_0 - \mathbf{x}^*\|$?
 - Not really: B typically depends on $\|\mathbf{x}_0 - \mathbf{x}^*\|$
 - for example, $B = \mathcal{O}(\|\mathbf{x}_0 - \mathbf{x}^*\|)$ for quadratic functions

Subgradient Method on Strongly Convex Functions

$$T \geq \frac{2B^2}{\mu\varepsilon} \quad \Rightarrow \quad \text{error} \leq \frac{2B^2}{\mu T} \leq \varepsilon$$

Conclusion: To reach absolute error at most ε , we need $\mathcal{O}(\frac{1}{\varepsilon})$ iterations,

- ▶ Recall: we can only hope that B is small (can be checked while running the algorithm)
- ▶ What if we don't know the parameter μ of strong convexity?
 - Heuristic strategy: try some μ 's, pick best solution obtained
 - Choosing the step size without μ ...

Questions?

Summary of Convergence Rates for Gradient Descent and Subgradient Methods

f	Algorithm	Convergence	# Iterations
Convex L -Lipschitz	GD	$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq \frac{RL}{\sqrt{T}}$	$\frac{R^2 L^2}{\varepsilon^2}$
Convex L -Smooth	GD	$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq \frac{R^2 L}{2T}$	$\frac{R^2 L}{2\varepsilon}$
μ -Strongly Convex L -Smooth	GD	$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq \frac{RL}{2}(1 - \frac{\mu}{L})^T$	$\frac{L}{\mu} \ln \left(\frac{R^2 L}{2\varepsilon} \right)$
Convex L -Lipschitz	Subgrad	$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq \frac{LR}{\sqrt{T}}$	$\frac{R^2 L^2}{\varepsilon^2}$
μ -Strongly Convex $\ g\ \leq B$	Subgrad	$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq \frac{2B^2}{\mu(T+1)}$	$\frac{2B^2}{\mu\varepsilon}$

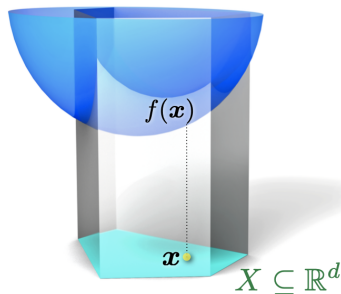
Here, $R = \|\mathbf{x}_0 - \mathbf{x}^*\|$

Projected Gradient Descent

Constrained Optimization

Constrained Optimization Problem

$$\begin{array}{ll}\text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in X\end{array}$$



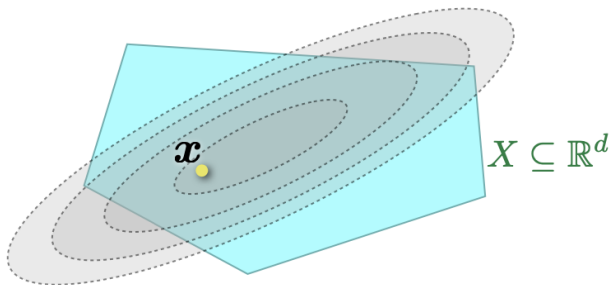
Constrained Optimization

Solving Constrained Optimization Problem

$$\begin{array}{ll}\text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in X\end{array}$$

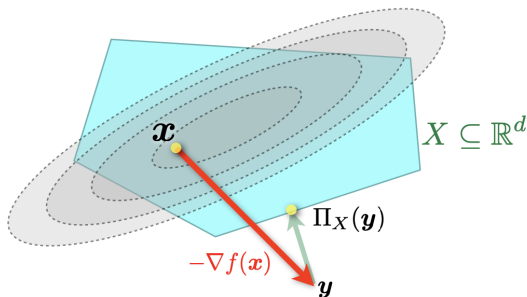
Solution:

- ▶ Projected Gradient Descent
- ▶ Transform it into an *unconstrained* problem



Projected Gradient Descent

Idea: project onto X after every step: $\Pi_X(\mathbf{y}) := \arg \min_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{y}\|$



Projected gradient descent: $\mathbf{x}_{t+1} = \Pi_X[\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)]$

Projected Gradient Descent

Projected gradient descent:

$$\mathbf{y}_{t+1} := \mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)$$

$$\mathbf{x}_{t+1} := \Pi_X(\mathbf{y}_{t+1}) = \arg \min_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{y}\|$$

for **stepsize** $\eta_t > 0$ and **timesteps** $t = 0, 1, \dots$

Properties of Projection

Fact 1.

Let $X \subseteq \mathbb{R}^d$ be closed and convex, $\mathbf{x} \in X, \mathbf{y} \in \mathbb{R}^d$. Then

- ▶ $(\mathbf{x} - \Pi_X(\mathbf{y}))^\top (\mathbf{y} - \Pi_X(\mathbf{y})) \leq 0$
- ▶ $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2$

Proof.

(i) $\Pi_X(\mathbf{y})$ is the minimizer of (differentiable) convex function $d_{\mathbf{y}} = \|\mathbf{x} - \mathbf{y}\|^2$ over X .

By first-order characterization of optimality (Lemma 4 from Lecture 2),

$$\begin{aligned} 0 &\leq \nabla d_{\mathbf{y}}(\Pi_X(\mathbf{y}))^\top (\mathbf{x} - \Pi_X(\mathbf{y})) \\ &= 2(\Pi_X(\mathbf{y}) - \mathbf{y})^\top (\mathbf{x} - \Pi_X(\mathbf{y})) \\ \Leftrightarrow 0 &\geq 2(\mathbf{y} - \Pi_X(\mathbf{y}))^\top (\mathbf{x} - \Pi_X(\mathbf{y})) \\ \Leftrightarrow 0 &\geq (\mathbf{x} - \Pi_X(\mathbf{y}))^\top (\mathbf{y} - \Pi_X(\mathbf{y})) \end{aligned}$$

Properties of Projection

Let $X \subseteq \mathbb{R}^d$ be closed and convex, $\mathbf{x} \in X, \mathbf{y} \in \mathbb{R}^d$. Then

- ▶ $(\mathbf{x} - \Pi_X(\mathbf{y}))^\top (\mathbf{y} - \Pi_X(\mathbf{y})) \leq 0$
- ▶ $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2$

Proof.

(ii)

$$\mathbf{v} := (\mathbf{x} - \Pi_X(\mathbf{y})), \quad \mathbf{w} := (\mathbf{y} - \Pi_X(\mathbf{y})).$$

By (i),

$$\begin{aligned} 0 &\geq 2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2 \\ &= \|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 - \|\mathbf{x} - \mathbf{y}\|^2. \end{aligned}$$

Convergence Rate of Projected Gradient Descent

The same number of steps as a gradient over \mathbb{R}^d !

- ▶ Lipschitz convex functions over X : $\mathcal{O}(1/\varepsilon^2)$ steps
- ▶ Smooth convex functions over X : $\mathcal{O}(1/\varepsilon)$ steps
- ▶ Smooth and strongly convex functions over X : $\mathcal{O}(\log(1/\varepsilon))$ steps

We will adapt the previous proofs for gradient descent.

BUT:

- ▶ Each step involves a projection onto X
- ▶ may or may not be efficient. . .

Projected GD on Lipschitz Convex Functions

Assume that all gradients of f are bounded in norm over **closed and convex** X .

- ▶ Equivalent to f being Lipschitz over X
- ▶ Many interesting functions are Lipschitz over bounded sets X .

Theorem 4 (same as the unconstrained one, but more useful)

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and differentiable, $X \subseteq \mathbb{R}^d$ closed and convex, \mathbf{x}^* is the minimizer of f over X ; furthermore, suppose that $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$ with $\mathbf{x}_0 \in X$, and that $\|\nabla f(\mathbf{x})\| \leq B$ for all $\mathbf{x} \in X$.

Choosing the constant stepsize

$$\eta = \frac{R}{B\sqrt{T}},$$

projected gradient descent yields

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

Projected GD on Lipschitz Convex Functions

Proof.

- Replace \mathbf{x}_{t+1} in the vanilla analysis with \mathbf{y}_{t+1} (the unprojected gradient step):

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) = \frac{1}{2\eta} (\eta^2 \|\mathbf{g}_t\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}^*\|^2).$$

- Use Fact 1 (ii): $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2$
- With $\mathbf{x} = \mathbf{x}^*$, $\mathbf{y} = \mathbf{y}_{t+1}$ we have $\Pi_X(\mathbf{y}) = \mathbf{x}_{t+1}$, and hence

$$\|\mathbf{x}^* - \mathbf{x}_{t+1}\| \leq \|\mathbf{x}^* - \mathbf{y}_{t+1}\|$$

- We go back to the original vanilla analysis and continue from there as before:

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) \leq \frac{1}{2\eta} (\eta^2 \|\mathbf{g}_t\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2)$$

Smooth functions over X

Recall:

f is called smooth (with parameter L) over X if

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

Lemma 2

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable and smooth with parameter L over X .
Choosing stepsize

$$\eta = \frac{1}{L},$$

projected gradient descent with arbitrary $\mathbf{x}_0 \in X$ satisfies

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

Sufficient Decrease

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

Proof.





Use smoothness, $\mathbf{y}_{t+1} - \mathbf{x}_t = -\nabla f(\mathbf{x}_t)/L$, $2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2$:

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

Questions?

Assignment 1 will be released after this class
and is due on **October 1 at 11:59 PM.**

References

-  Stephen P Boyd, *Lecture notes for ee 264b,stanford university (2010-2011)*.
-  Sébastien Bubeck, *Convex optimization: Algorithms and complexity*, Foundations and Trends in Machine Learning **8** (2015), no. 3-4, 231–357.
-  Stephen P Boyd and Lieven Vandenberghe, *Convex optimization*, Cambridge University Press, 2004.
-  R. Tyrrell Rockafellar, *Convex analysis*, Princeton University Press.