

COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2023.

Lecture 23

- Problem Set 5 is posted. It can be turned in up to 12/11 (next Monday) at 11:59pm It is **optional** – the core problems can replace the lowest of your previous four core problem grades.
- The final will be on 12/14 in this room, 10:30am-12:30pm.
- See Piazza for additional final review office ours schedule.
- See website/Moodle for final prep material.

Summary

Last Class:

- Multivariable calculus review and gradient computation.
- Introduction to gradient descent. Motivation as a greedy algorithm.
- Convex functions

This Class:

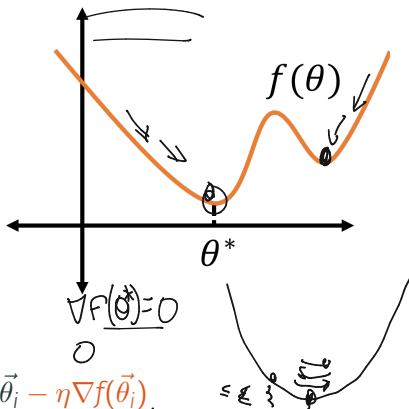
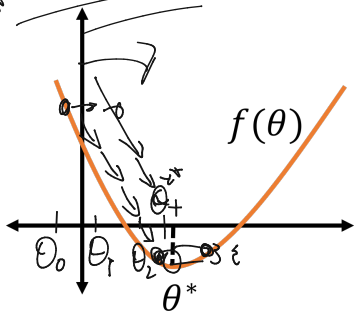
- Analysis of gradient descent for Lipschitz, convex functions.
- Extension to projected gradient descent for constrained optimization.

When Does Gradient Descent Work?

Convex functions

$$\theta \in \mathbb{R} \quad \nabla f(\theta) \in \mathbb{R}$$

$\epsilon \downarrow$ $\eta \downarrow$ $\text{steps} \uparrow$



Gradient Descent Update: $\vec{\theta}_{i+1} = \vec{\theta}_i - \eta \nabla f(\vec{\theta}_i)$

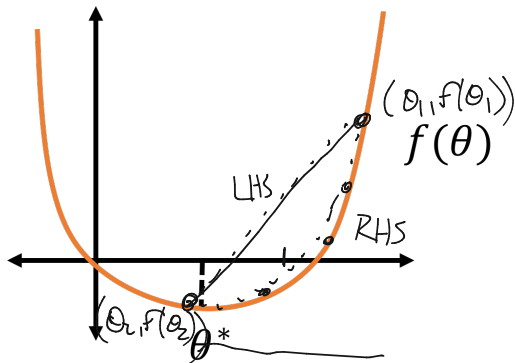
$$f(\theta_t) - f(\theta_*) \leq \epsilon$$

Convexity

Definition – Convex Function: A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathbb{R}^d$ and $\lambda \in [0, 1]$:

$$(1 - \lambda) \cdot f(\vec{\theta}_1) + \lambda \cdot f(\vec{\theta}_2) \geq f\left(\underbrace{(1 - \lambda) \cdot \vec{\theta}_1 + \lambda \cdot \vec{\theta}_2}\right)$$

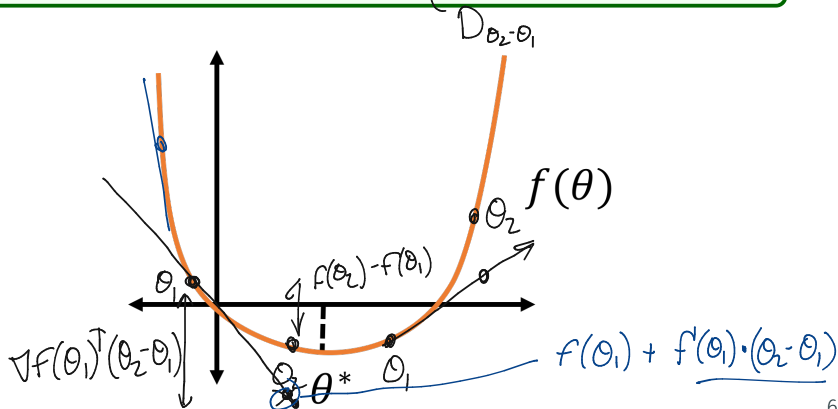
x^2
CX



Convexity

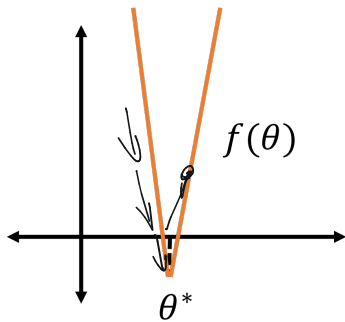
Corollary – Convex Function: A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathbb{R}^d$ and $\lambda \in [0, 1]$:

$$f(\vec{\theta}_2) - f(\vec{\theta}_1) \geq \underbrace{\nabla f(\vec{\theta}_1)^T}_{\substack{f'(\theta_1) \\ |x \times c \\ |dx|}} (\vec{\theta}_2 - \vec{\theta}_1)$$



Lipschitz Functions

$$\theta \in \mathbb{R} \quad \nabla f(\theta) \in \mathbb{R}$$



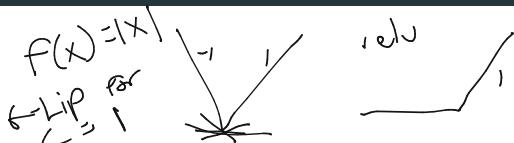
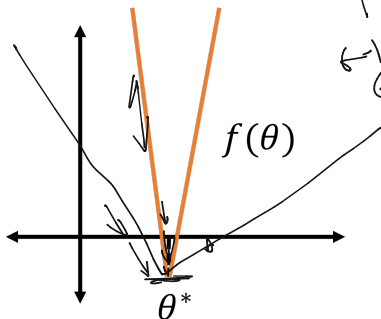
Gradient Descent Update:

$$\vec{\theta}_{i+1} = \vec{\theta}_i - \underline{\eta \nabla f(\vec{\theta}_i)}$$

bigger step size
= larger error

Lipschitz Functions

$$\theta \in \mathbb{R} \quad \nabla f(\theta) \in \mathbb{R}$$



Gradient Descent Update:

$$\vec{\theta}_{i+1} = \vec{\theta}_i - \eta \nabla f(\vec{\theta}_i)$$

is $f(x) = x^2 \leftarrow \text{Lip.}$
 $f'(x) = 2x$ for some $G < \infty$

Need to assume that the function is **Lipschitz** (size of gradient is bounded): There is some G s.t.:

$$\forall \vec{\theta} : \quad \|\vec{\nabla} f(\vec{\theta})\|_2 \leq G \Leftrightarrow \forall \vec{\theta}_1, \vec{\theta}_2 : \quad |f(\vec{\theta}_1) - f(\vec{\theta}_2)| \leq G \cdot \|\vec{\theta}_1 - \vec{\theta}_2\|_2$$

$$|f'(\theta)| \leq G$$

Well-Behaved Functions

Definition – Convex Function: A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathbb{R}^d$ and $\lambda \in [0, 1]$:

$$(1 - \lambda) \cdot f(\vec{\theta}_1) + \lambda \cdot f(\vec{\theta}_2) \geq f\left((1 - \lambda) \cdot \vec{\theta}_1 + \lambda \cdot \vec{\theta}_2\right)$$

Corollary – Convex Function: A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathbb{R}^d$ and $\lambda \in [0, 1]$:

$$f(\vec{\theta}_2) - f(\vec{\theta}_1) \geq \nabla f(\vec{\theta}_1)^T (\vec{\theta}_2 - \vec{\theta}_1)$$

Definition – Lipschitz Function: A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is G -Lipschitz if $\|\vec{\nabla} f(\vec{\theta})\|_2 \leq G$ for all $\vec{\theta}$.

GD Analysis – Convex Functions

Assume that:

- f is convex.
- f is G -Lipschitz.
- $\|\vec{\theta}_1 - \vec{\theta}_*\|_2 \leq R$ where $\vec{\theta}_1$ is the initialization point.

$$\theta_1 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad R = \|\theta_*\|$$

Gradient Descent

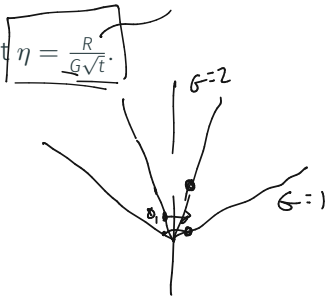
- Choose some initialization $\vec{\theta}_1$ and set $\eta = \frac{R}{G\sqrt{t}}$.

- For $i = 1, \dots, t-1$

$$\vec{\theta}_{i+1} = \vec{\theta}_i - \eta \vec{\nabla} f(\vec{\theta}_i)$$

- Return $\hat{\theta} = \arg \min_{\vec{\theta}_1, \dots, \vec{\theta}_t} f(\vec{\theta}_i)$.

θ_+



Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$\underline{f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.}$$

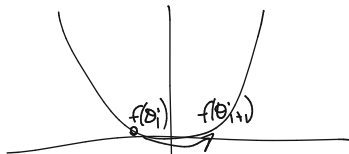
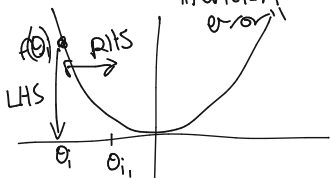
GD Analysis Proof

Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$\underline{f(\hat{\theta})} \leq f(\vec{\theta}_*) + \epsilon.$$

how much closer to θ^ I move in step i*

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \underbrace{\frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta}}_{\text{LHS}} + \underbrace{\left(\frac{\eta G^2}{2}\right)}_{\text{Small}}$ **Visually:**



Either $f(\theta_i)$ is small

OR I move much closer to θ^* .

GD Analysis Proof

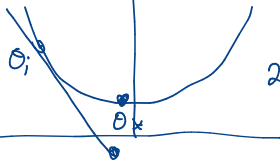
Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.$$

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \theta_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$. Formally:

$$\begin{aligned} \|\theta_{i+1} - \theta_*\|_2^2 &= \|\theta_i - \eta \nabla f(\theta_i) - \theta_*\|_2^2 \\ &= \|\theta_i - \theta_*\|_2^2 - 2\eta \nabla f(\theta_i)^\top (\theta_i - \theta_*) + \|\eta \nabla f(\theta_i)\|_2^2 \\ &\leq \|\theta_i - \theta_*\|_2^2 - 2\eta (f(\theta_i) - f(\theta_*)) + \eta^2 G^2 \end{aligned}$$

$$\begin{aligned} 2\eta \nabla f(\theta_i)^\top (\theta_i - \theta_*) &\leq \|\theta_i - \theta_*\|_2^L - \|\theta_{i+1} - \theta_*\|_2 + \eta^2 G^2 \\ \left[\nabla f(\theta_i)^\top (\theta_i - \theta_*) \right] &\leq \frac{\|\theta_i - \theta_*\|_2^L - \|\theta_{i+1} - \theta_*\|_2 + \eta G^2}{2} \end{aligned}$$



Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.$$

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

Step 1.1: $\underbrace{\vec{\nabla} f(\vec{\theta}_i)^T (\vec{\theta}_i - \vec{\theta}_*)}_{\leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}}$

Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.$$

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

Step 1.1: $\nabla f(\vec{\theta}_i)^T (\vec{\theta}_i - \vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2} \implies$ **Step 1 by convexity.**

$\nabla f(\theta_i)^T (\theta_i - \theta_*) \geq f(\theta_i) - f(\theta_*)$ by convexity

Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.$$

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$

GD Analysis Proof

Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of θ_* , outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\theta_*) + \epsilon.$$

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2} \implies$

Step 2: $\frac{1}{t} \sum_{i=1}^t f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{R^2}{2\eta \cdot t} + \frac{\eta G^2}{2}.$

$$\frac{1}{2\eta t} \sum_{i=1}^t \|\theta_i - \theta_*\|_2^2 - \|\theta_{t+1} - \theta_*\|_2^2 + \frac{\eta G^2}{2} \leq \frac{R^2}{2\eta t} + \frac{\eta G^2}{2} \leq \epsilon$$

$$\|\theta_1 - \theta_*\|_2^2 - \|\theta_2 - \theta_*\|_2^2 + \|\theta_2 - \theta_*\|_2^2 - \|\theta_3 - \theta_*\|_2^2 \dots - \|\theta_t - \theta_*\|_2^2$$

$$\|\theta_1 - \theta_*\|_2^2 - \|\theta_{t+1} - \theta_*\|_2^2 \leq \|\theta_1 - \theta_*\|_2^2 \leq R^2$$

GD Analysis Proof

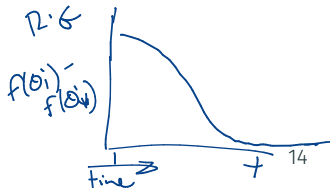
Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$ and starting point within radius R of θ_* , outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\theta_*) + \epsilon.$$

Step 2: $\frac{1}{t} \sum_{i=1}^t f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{R^2}{2\eta t} + \frac{\eta G^2}{2} \leq \epsilon$

$$\frac{R^2}{2\eta t} + \frac{\eta G^2}{2} = \frac{R^2}{2 \cdot \frac{R}{G\sqrt{t}} \cdot t} + \frac{\frac{R}{G\sqrt{t}} \cdot G^2}{2} = \frac{R^2 G}{2\sqrt{t}} + \frac{R G}{2\sqrt{t}} = \frac{R G}{\sqrt{t}} = \frac{R G}{\frac{R G}{\epsilon}} = \epsilon$$

$$f(x) = (x-a)^2$$
$$f(\hat{\theta}) = \min_{i=1..t} f(\theta_i) = \frac{1}{t} \sum f(\theta_i)$$



Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.$$

Step 2: $\frac{1}{t} \sum_{i=1}^t f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{R^2}{2\eta \cdot t} + \frac{\eta G^2}{2}.$

Constrained Convex Optimization

Often want to perform convex optimization with convex constraints.

$$\vec{\theta}^* = \arg \min_{\vec{\theta} \in \mathcal{S}} f(\vec{\theta}),$$

where \mathcal{S} is a convex set.

Constrained Convex Optimization

Often want to perform **convex optimization with convex constraints**.

$$\vec{\theta}^* = \arg \min_{\vec{\theta} \in \mathcal{S}} f(\vec{\theta}),$$

where \mathcal{S} is a **convex set**.

Definition – Convex Set: A set $\mathcal{S} \subseteq \mathbb{R}^d$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathcal{S}$ and $\lambda \in [0, 1]$:

$$(1 - \lambda)\vec{\theta}_1 + \lambda \cdot \vec{\theta}_2 \in \mathcal{S}$$

Constrained Convex Optimization

Often want to perform **convex optimization with convex constraints**.

$$\vec{\theta}^* = \arg \min_{\vec{\theta} \in \mathcal{S}} f(\vec{\theta}),$$

where \mathcal{S} is a **convex set**.

Definition – Convex Set: A set $\mathcal{S} \subseteq \mathbb{R}^d$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathcal{S}$ and $\lambda \in [0, 1]$:

$$(1 - \lambda)\vec{\theta}_1 + \lambda \cdot \vec{\theta}_2 \in \mathcal{S}$$

E.g. $\mathcal{S} = \{\vec{\theta} \in \mathbb{R}^d : \|\vec{\theta}\|_2 \leq 1\}$.

Projected Gradient Descent

For any convex set let $P_{\mathcal{S}}(\cdot)$ denote the projection function onto \mathcal{S} .

- $P_{\mathcal{S}}(\vec{y}) = \arg \min_{\vec{\theta} \in \mathcal{S}} \|\vec{\theta} - \vec{y}\|_2.$

Projected Gradient Descent

For any convex set let $P_{\mathcal{S}}(\cdot)$ denote the projection function onto \mathcal{S} .

- $P_{\mathcal{S}}(\vec{y}) = \arg \min_{\vec{\theta} \in \mathcal{S}} \|\vec{\theta} - \vec{y}\|_2$.
- For $\mathcal{S} = \{\vec{\theta} \in \mathbb{R}^d : \|\vec{\theta}\|_2 \leq 1\}$ what is $P_{\mathcal{S}}(\vec{y})$?

Projected Gradient Descent

For any convex set let $P_{\mathcal{S}}(\cdot)$ denote the projection function onto \mathcal{S} .

- $P_{\mathcal{S}}(\vec{y}) = \arg \min_{\vec{\theta} \in \mathcal{S}} \|\vec{\theta} - \vec{y}\|_2$.
- For $\mathcal{S} = \{\vec{\theta} \in \mathbb{R}^d : \|\vec{\theta}\|_2 \leq 1\}$ what is $P_{\mathcal{S}}(\vec{y})$?
- For \mathcal{S} being a k dimensional subspace of \mathbb{R}^d , what is $P_{\mathcal{S}}(\vec{y})$?

Projected Gradient Descent

For any convex set let $P_{\mathcal{S}}(\cdot)$ denote the projection function onto \mathcal{S} .

- $P_{\mathcal{S}}(\vec{y}) = \arg \min_{\vec{\theta} \in \mathcal{S}} \|\vec{\theta} - \vec{y}\|_2$.
- For $\mathcal{S} = \{\vec{\theta} \in \mathbb{R}^d : \|\vec{\theta}\|_2 \leq 1\}$ what is $P_{\mathcal{S}}(\vec{y})$?
- For \mathcal{S} being a k dimensional subspace of \mathbb{R}^d , what is $P_{\mathcal{S}}(\vec{y})$?

Projected Gradient Descent

- Choose some initialization $\vec{\theta}_1$ and set $\eta = \frac{R}{G\sqrt{t}}$.
- For $i = 1, \dots, t - 1$
 - $\vec{\theta}_{i+1}^{(out)} = \vec{\theta}_i - \eta \cdot \vec{\nabla} f(\vec{\theta}_i)$
 - $\vec{\theta}_{i+1} = P_{\mathcal{S}}(\vec{\theta}_{i+1}^{(out)})$.
- Return $\hat{\theta} = \arg \min_{\vec{\theta}_i} f(\vec{\theta}_i)$.

Convex Projections

Projected gradient descent can be analyzed identically to gradient descent!

Convex Projections

Projected gradient descent can be analyzed identically to gradient descent!

Theorem – Projection to a convex set: For any convex set $\mathcal{S} \subseteq \mathbb{R}^d$, $\vec{y} \in \mathbb{R}^d$, and $\vec{\theta} \in \mathcal{S}$,

$$\|P_{\mathcal{S}}(\vec{y}) - \vec{\theta}\|_2 \leq \|\vec{y} - \vec{\theta}\|_2.$$

Projected Gradient Descent Analysis

Theorem – Projected GD: For convex G -Lipschitz function f , and convex set \mathcal{S} , Projected GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon = \min_{\vec{\theta} \in \mathcal{S}} f(\vec{\theta}) + \epsilon$$

Projected Gradient Descent Analysis

Theorem – Projected GD: For convex G -Lipschitz function f , and convex set \mathcal{S} , Projected GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon = \min_{\vec{\theta} \in \mathcal{S}} f(\vec{\theta}) + \epsilon$$

Recall: $\vec{\theta}_{i+1}^{(out)} = \vec{\theta}_i - \eta \cdot \vec{\nabla} f(\vec{\theta}_i)$ and $\vec{\theta}_{i+1} = P_{\mathcal{S}}(\vec{\theta}_{i+1}^{(out)})$.

Projected Gradient Descent Analysis

Theorem – Projected GD: For convex G -Lipschitz function f , and convex set \mathcal{S} , Projected GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon = \min_{\vec{\theta} \in \mathcal{S}} f(\vec{\theta}) + \epsilon$$

Recall: $\vec{\theta}_{i+1}^{(out)} = \vec{\theta}_i - \eta \cdot \vec{\nabla} f(\vec{\theta}_i)$ and $\vec{\theta}_{i+1} = P_{\mathcal{S}}(\vec{\theta}_{i+1}^{(out)})$.

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \theta_*\|_2^2 - \|\vec{\theta}_{i+1}^{(out)} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

Projected Gradient Descent Analysis

Theorem – Projected GD: For convex G -Lipschitz function f , and convex set \mathcal{S} , Projected GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon = \min_{\vec{\theta} \in \mathcal{S}} f(\vec{\theta}) + \epsilon$$

Recall: $\vec{\theta}_{i+1}^{(out)} = \vec{\theta}_i - \eta \cdot \vec{\nabla} f(\vec{\theta}_i)$ and $\vec{\theta}_{i+1} = P_{\mathcal{S}}(\vec{\theta}_{i+1}^{(out)})$.

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1}^{(out)} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

Step 1.a: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

Projected Gradient Descent Analysis

Theorem – Projected GD: For convex G -Lipschitz function f , and convex set \mathcal{S} , Projected GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon = \min_{\vec{\theta} \in \mathcal{S}} f(\vec{\theta}) + \epsilon$$

Recall: $\vec{\theta}_{i+1}^{(out)} = \vec{\theta}_i - \eta \cdot \vec{\nabla} f(\vec{\theta}_i)$ and $\vec{\theta}_{i+1} = P_{\mathcal{S}}(\vec{\theta}_{i+1}^{(out)})$.

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1}^{(out)} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

Step 1.a: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

Step 2: $\frac{1}{t} \sum_{i=1}^t f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{R^2}{2\eta \cdot t} + \frac{\eta G^2}{2} \implies$ Theorem.