

# COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

---

Cameron Musco

University of Massachusetts Amherst. Spring 2020.

Lecture 23

## SUMMARY

### Last Class:

- Multivariable calculus review and gradient computation.
- Introduction to gradient descent. Motivation as a greedy algorithm.
- Conditions under which we will analyze gradient descent: convexity and Lipschitzness.

### This Class:

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

**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  $\lambda \in [0, 1]$ :

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

**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  $\lambda \in [0, 1]$ :

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

**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  $\vec{\theta}$ .

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.

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

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

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

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

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

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

**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.

**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  $\mathcal{S}$ .

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

**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:

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

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