

# COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

---

Cameron Musco

University of Massachusetts Amherst. Fall 2021.

Lecture 24

- Problem Set 5 is posted. It is due 12/13. It is optional and can be used to replace your lowest problem set grade.
- Quiz due Monday, 8pm. Reminder that lowest quiz grade is dropped.
- The final will be on 12/16 from 10:30am-12:30pm. In the class.
- Final review sheet is posted under the 'Schedule Tab'. I may continue to add to this and we plan to post a practice exam(s).

Several extra office hours will be held before the final. Times TBD.

## Last Class:

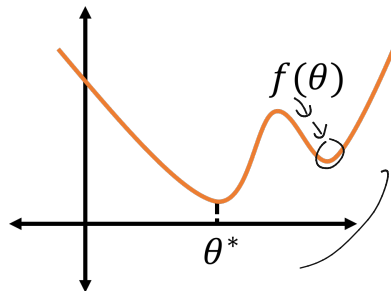
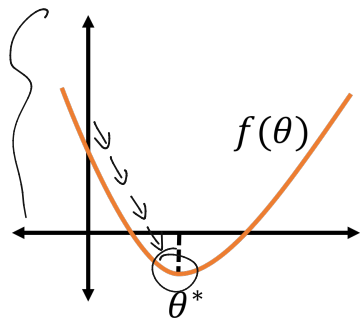
- Multivariable calculus review and gradient computation.
- Introduction to gradient descent. Motivation as a greedy algorithm.

## This Class:

- Conditions under which we will analyze gradient descent: convexity and Lipschitzness.
- Analysis of gradient descent for Lipschitz, convex functions.
- Extension to projected gradient descent for **constrained optimization**.

# WHEN DOES GRADIENT DESCENT WORK?

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

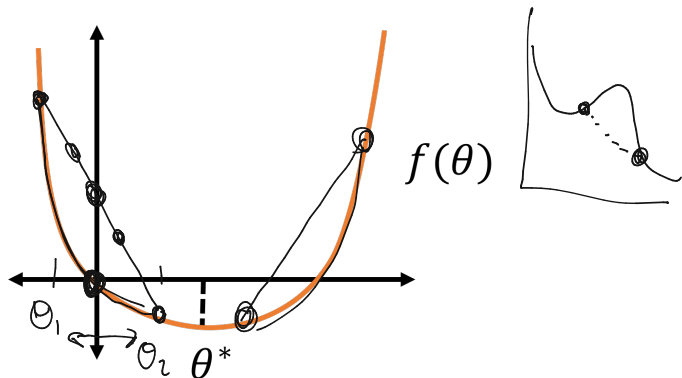


$$f'(\theta_i)$$

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

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

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

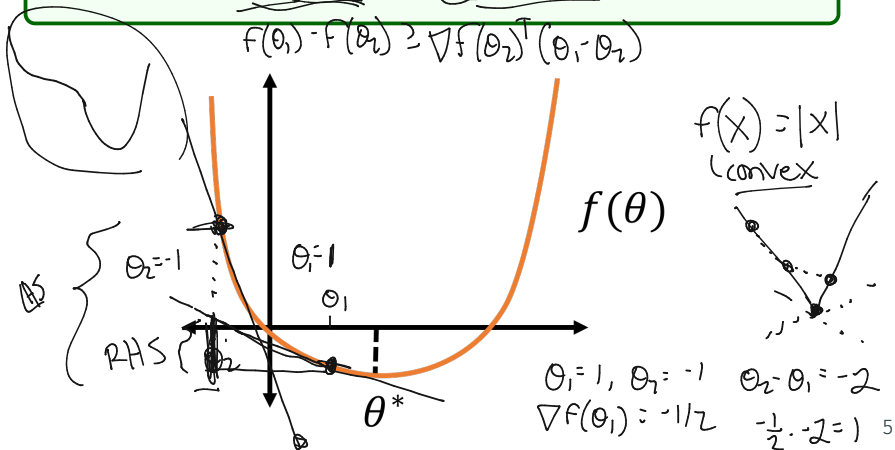


# CONVEXITY

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

$$\underbrace{f(\vec{\theta}_2) - f(\vec{\theta}_1)}_{\text{LHS}} \geq \underbrace{\nabla f(\vec{\theta}_1)^T (\vec{\theta}_2 - \vec{\theta}_1)}_{\text{RHS}}$$

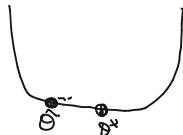
$$f(\theta_1) - f(\theta_2) \geq \nabla f(\theta_2)^T (\theta_1 - \theta_2)$$



## CONDITIONS FOR GRADIENT DESCENT CONVERGENCE

**Convex Functions:** After sufficient iterations, if the step size  $\eta$  is chosen appropriately, gradient descent will converge to an **approximate minimizer**  $\hat{\theta}$  with:

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



Examples: least squares regression, logistic regression, sparse regression (lasso), regularized regression, SVMs,...

# CONDITIONS FOR GRADIENT DESCENT CONVERGENCE

**Convex Functions:** After sufficient iterations, if the step size  $\eta$  is chosen appropriately, gradient descent will converge to an **approximate minimizer**  $\hat{\theta}$  with:

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

Examples: least squares regression, logistic regression, sparse regression (lasso), regularized regression, SVMs,...

**Non-Convex Functions:** After sufficient iterations, gradient descent will converge to an **approximate stationary point**  $\hat{\theta}$  with:

$$\|\nabla f(\hat{\theta})\|_2 \leq \epsilon.$$





**Convex Functions:** After sufficient iterations, if the step size  $\eta$  is chosen appropriately, gradient descent will converge to an **approximate minimizer**  $\hat{\theta}$  with:

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

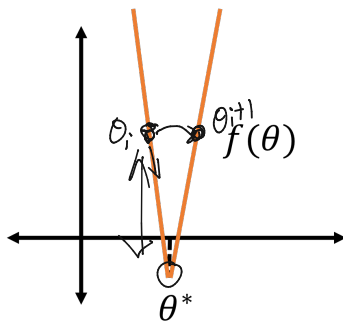
Examples: least squares regression, logistic regression, sparse regression (lasso), regularized regression, SVMs,...

**Non-Convex Functions:** After sufficient iterations, gradient descent will converge to an **approximate stationary point**  $\hat{\theta}$  with:

$$\|\nabla f(\hat{\theta})\|_2 \leq \epsilon.$$

Examples: neural networks, clustering, mixture models.

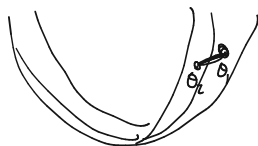
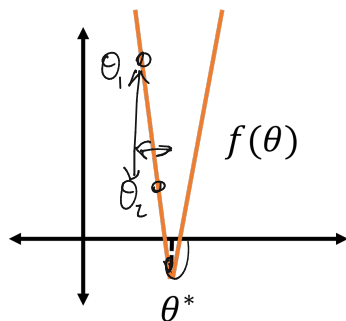
$$\theta \in \mathbb{R} \quad \nabla f(\theta) \in \mathbb{R}$$



Gradient Descent Update:

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

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



Gradient Descent Update:

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

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

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

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

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

**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:

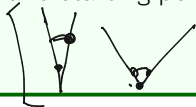
- $f$  is convex.
- $f$  is  $G$ -Lipschitz.
- $\|\vec{\theta}_1 - \vec{\theta}_*\|_2 \leq R$  where  $\vec{\theta}_1$  is the initialization point.

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:



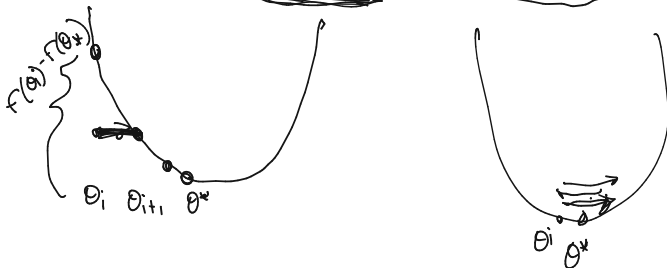
$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.$$

$$\|\theta_1 - \theta_*\| \leq R$$

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

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



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

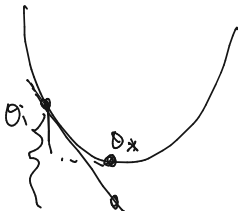
**Step 1.1:**  $\underline{\nabla} f(\vec{\theta}_i)^T (\vec{\theta}_i - \vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$

**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.



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

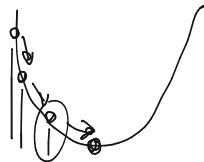
**Theorem – GD on Convex Lipschitz Functions:** For convex  $G$ -Lipschitz function  $f$ , GD run with  $t \geq \frac{R^2 G^2}{\epsilon^2}$  iterations,  $\eta = \frac{R}{G\sqrt{t}}$ , and starting point within radius  $R$  of  $\vec{\theta}_*$ , outputs  $\hat{\theta}$  satisfying:

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

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

~~$\frac{1}{t} \sum_{i=1}^t f(\vec{\theta}_i) - f(\vec{\theta}_*)$~~   
error at iteration  $i$



$$\underline{f(\hat{\theta}) - f(\theta_*)} \leq \frac{RL}{2\eta t} + \frac{\eta G^2}{2}$$

$$\underline{\hat{\theta}} = \underset{i=1, \dots, t}{\operatorname{argmin}} f(\theta_i)$$

$\hat{\theta}$  is our output  $\rightarrow$  approx minimizer.

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

$$\leq \frac{1}{t} \sum_{i=1}^t \left( \frac{\|\theta_i - \theta_*\|_2^2 - \|\theta_{i+1} - \theta_*\|_2^2}{2\eta} + \frac{\eta G^2}{2} \right)$$

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

$$\begin{aligned} & \|\theta_1 - \theta_*\|_2^2 - \|\theta_2 - \theta_*\|_2^2 + \|\theta_2 - \theta_*\|_2^2 - \|\theta_3 - \theta_*\|_2^2 + \dots \\ & = \|\theta_1 - \theta_*\|_2^2 - \|\theta_{t+1} - \theta_*\|_2^2 \leq \|\theta_1 - \theta_*\|_2^2 \leq R^2 \end{aligned}$$

not making it to opt

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

$$\leq \frac{R^2}{2\eta t} + \frac{\eta G^2}{2}$$

# 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**.

Previous slide:  $t = \frac{R^2 G^2}{\epsilon^2}$        $m = \frac{R}{G\sqrt{t}} = \frac{\epsilon}{G^2}$

$$\begin{aligned} f(\vec{\theta}) - f(\vec{\theta}^*) &\leq \frac{\epsilon^2}{2G^2 m} + \frac{m b^2}{2} \\ &\leq \frac{\epsilon^2}{2G^2 \frac{\epsilon}{G^2}} + \frac{\frac{\epsilon}{G^2} G^2}{2} = \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon \end{aligned}$$