

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Fall 2020.

Lecture 24

- Problem Set 4 is due tomorrow at 8pm.
- Optional Problem Set 5 will be released tomorrow, due 11/30.
- Exam will span December 3-4. Any two hour period.
- Exam review guide, practice problems, logistical details have been posted under the schedule tab on the course page.
- I am holding an optional SRTI (course reviews) for this class and would really appreciate your feedback (closes Dec 6).
- <http://owl.umass.edu/partners/courseEvalSurvey/uma/>.

Last Class:



Lipschitz.

- Analysis of gradient descent for optimizing convex functions.
- Introduction to convex sets and projection functions.
- (The same) analysis of projected gradient descent for optimizing under convex functions under (convex) constraints. $\theta \in S$

This Class:

- Online learning, regret, and online gradient descent.
- Application to stochastic gradient descent.

$$2x + 3y$$

$$\mathbb{R}^d \rightarrow \mathbb{R}$$

Consider the function $f(\vec{\theta}) = \vec{x}^T \vec{\theta}$ for

$x = [1, -1, -2]$. Give the minimum value of G such that $f(\vec{\theta})$ is G -Lipschitz

$$\| \nabla f(\vec{\theta}) \|_2 \leq G \quad \forall \theta$$

$$\nabla f(\vec{\theta}) = \vec{x}$$

$$= \begin{bmatrix} \frac{\partial f(\theta)}{\partial \theta(1)} \\ \frac{\partial f(\theta)}{\partial \theta(2)} \\ \vdots \end{bmatrix} = \begin{bmatrix} x(1) \\ x(2) \\ \vdots \\ x(d) \end{bmatrix}$$

$$f(\theta) = \vec{x}^T \theta = \sum_{i=1}^d \underline{x(i)} \cdot \theta(i)$$

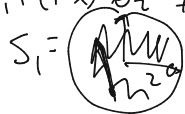
$$\frac{\partial f(\theta)}{\partial \theta(i)} = x(i)$$

$$\Rightarrow \underline{\| \nabla f(\theta) \|_2} = \| \vec{x} \| = \underline{\underline{\sqrt{6}}}$$

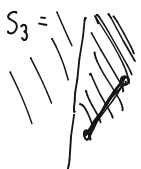
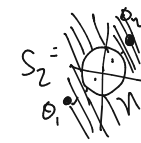
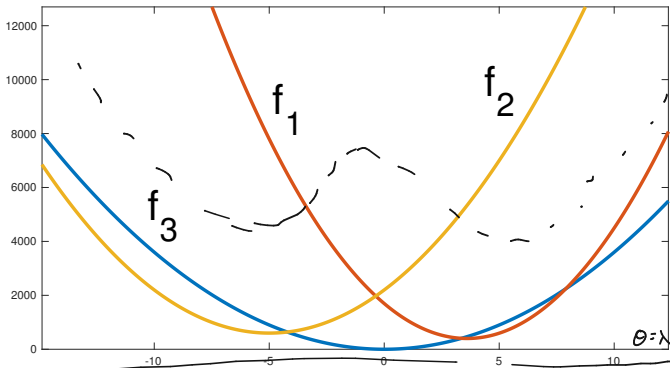
TEST OF INTUITION

$\theta_1, \theta_2 \in S_2, \lambda \in [0, 1]$ s.t. $\lambda\theta_1 + (1-\lambda)\theta_2 \neq S$

What does $f_1(\theta) + f_2(\theta) + f_3(\theta)$ look like?



$\lambda \in [0, 1]$



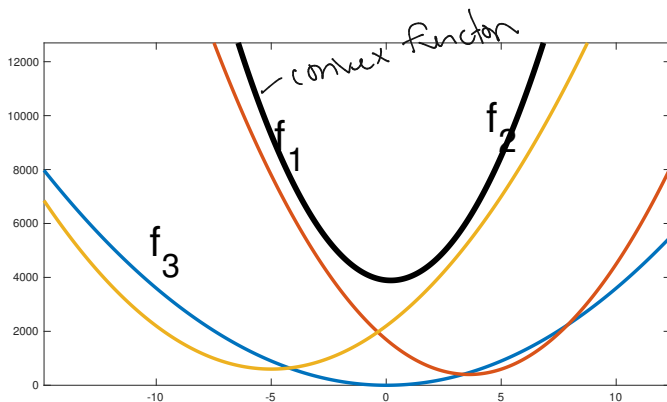
$\theta_1 \in S_3, \theta_2 \in S_3$
 $\theta_1(10) \geq 2, \theta_2(10) \geq 2$
 $\theta = \lambda\theta_1 + (1-\lambda)\theta_2$
 $\theta(10) \geq 2$

Q2: $\sqrt{S_1}: \{x : \|x\|_2 \leq 2\}$ $\sqrt{S_3}: \{x : x(10) \geq 2\}$

$\times S_2: \{x : \|x\|_2 \geq 2\}$
 $\theta_1 = [10, 0], \theta_2 = [-10, 0], \lambda = 1/2 \quad \lambda\theta_1 + (1-\lambda)\theta_2 = 0$

TEST OF INTUITION

What does $f_1(\theta) + f_2(\theta) + f_3(\theta)$ look like?

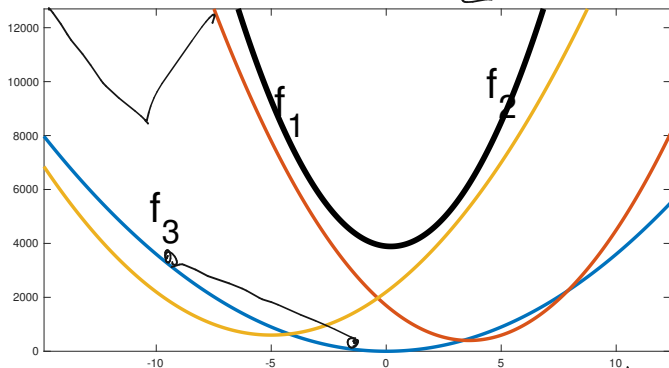


TEST OF INTUITION

$$\lambda \in [0, 1] \quad \theta_1, \theta_2$$

What does $f_1(\theta) + f_2(\theta) + f_3(\theta)$ look like?

$$[f(\lambda\theta_1, (1-\lambda)\theta_2) \leq \lambda f(\theta_1) + (1-\lambda)f(\theta_2)]$$



✓ prove directly from definition of convexity

A sum of convex functions is always convex (good exercise).

In reality many learning problems are online.

- Websites optimize ads or recommendations to show users, given continuous feedback from these users.
- Spam filters are incrementally updated and adapt as they see more examples of spam over time.
- Face recognition systems, other classification systems, learn from mistakes over time.

In reality many learning problems are online.

- Websites optimize ads or recommendations to show users, given continuous feedback from these users.
- Spam filters are incrementally updated and adapt as they see more examples of spam over time.
- Face recognition systems, other classification systems, learn from mistakes over time.

Want to minimize some global loss $L(\vec{\theta}, \mathbf{X}) = \sum_{i=1}^n \ell(\vec{\theta}, \vec{x}_i)$, when data points are presented in an online fashion $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$ (similar to streaming algorithms)

In reality many learning problems are online.

- Websites optimize ads or recommendations to show users, given continuous feedback from these users.
- Spam filters are incrementally updated and adapt as they see more examples of spam over time.
- Face recognition systems, other classification systems, learn from mistakes over time.

Want to minimize some global loss $L(\vec{\theta}, \mathbf{X}) = \sum_{i=1}^n \ell(\vec{\theta}, \vec{x}_i)$, when data points are presented in an online fashion $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$ (similar to streaming algorithms)

Stochastic gradient descent is a special case: when data points are considered a random order for computational reasons.

Online Optimization: In place of a single function f , we see a different objective function at each step:

$$f_1, f_2, \dots, f_t : \mathbb{R}^d \rightarrow \mathbb{R}$$

Online Optimization: In place of a single function f , we see a different objective function at each step:

$$f_1, f_2, \dots, f_t : \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\theta^{(1)} \quad \theta^{(2)} \quad \dots \quad \theta^{(t)}$$

- At each step, first pick (play) a parameter vector $\vec{\theta}^{(i)}$.
- Then are told f_i and incur cost $f_i(\vec{\theta}^{(i)})$.
- **Goal:** Minimize total cost $\sum_{i=1}^t f_i(\vec{\theta}^{(i)})$.

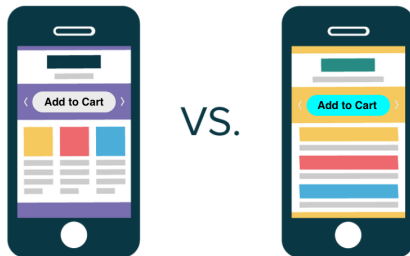
Online Optimization: In place of a single function f , we see a different objective function at each step:

$$f_1, f_2, \dots, f_t : \mathbb{R}^d \rightarrow \mathbb{R}$$

- At each step, first pick (play) a parameter vector $\vec{\theta}^{(i)}$.
- Then are told f_i and incur cost $f_i(\vec{\theta}^{(i)})$.
- **Goal:** Minimize total cost $\sum_{i=1}^t f_i(\vec{\theta}^{(i)})$.

Our analysis will make no assumptions on how f_1, \dots, f_t are related to each other!

UI design via online optimization.



- Parameter vector $\vec{\theta}^{(i)}$: some encoding of the layout at step i .
- Functions f_1, \dots, f_t : $f_i(\vec{\theta}^{(i)}) = 1$ if user does not click 'add to cart' and $f_i(\vec{\theta}^{(i)}) = 0$ if they do click.
- Want to maximize number of purchases. I.e., minimize $\sum_{i=1}^t f_i(\vec{\theta}^{(i)})$.

$$f_i(\theta^i)$$

Home pricing tools.



$$\vec{x} = [\underbrace{\#baths}_+, \underbrace{\#beds}_-, \underbrace{\#floors}_+ \dots]$$

$$f_i(\theta) = (\langle \vec{x}_i, \theta \rangle - p_i)^2$$

linear model

$$\langle \vec{x}, \vec{\theta} \rangle$$



\$275,000

home i

- Parameter vector $\vec{\theta}^{(i)}$: coefficients of linear model at step i .
- Functions f_1, \dots, f_t : $f_i(\vec{\theta}^{(i)}) = (\langle \vec{x}_i, \vec{\theta}^{(i)} \rangle - \text{price}_i)^2$ revealed when home_i is listed or sold.
- Want to minimize total squared error $\sum_{i=1}^t f_i(\vec{\theta}^{(i)})$ (same as classic least squares regression).

In normal optimization, we seek $\hat{\theta}$ satisfying:

$$\underline{f(\hat{\theta})} \leq \underbrace{\min_{\vec{\theta}} f(\vec{\theta})} + \epsilon. \lesssim f(\theta^*) + \epsilon$$



In normal optimization, we seek $\hat{\theta}$ satisfying:

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

In online optimization we will ask for the same.

online cost

$$\sum_{i=1}^t f_i(\vec{\theta}^{(i)}) \leq \min_{\vec{\theta}} \sum_{i=1}^t f_i(\vec{\theta}) + \epsilon = \sum_{i=1}^t f_i(\vec{\theta}^{off}) + \epsilon$$

$\vec{\theta}^{off}$: optimal
offline solution

ϵ is called the **regret**.

In normal optimization, we seek $\hat{\theta}$ satisfying:

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

In online optimization we will ask for the same.

$$\sum_{i=1}^t f_i(\vec{\theta}^{(i)}) \leq \min_{\vec{\theta}} \sum_{i=1}^t f_i(\vec{\theta}) + \epsilon = \sum_{i=1}^t f_i(\vec{\theta}^{off}) + \epsilon$$

ϵ is called the **regret**.

- This error metric is a bit 'unfair'. **Why?**

- we haven't seen θ ahead of time

$$\left[\min_{\theta_1, \theta_2, \dots, \theta_t} \sum f_i(\theta^i) \right]$$

In normal optimization, we seek $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq \min_{\vec{\theta}} f(\vec{\theta}) + \epsilon. \quad \epsilon \geq 0$$

In online optimization we will ask for the same.

$$\sum_{i=1}^t f_i(\vec{\theta}^{(i)}) \leq \min_{\vec{\theta}} \sum_{i=1}^t f_i(\vec{\theta}) + \epsilon = \sum_{i=1}^t f_i(\vec{\theta}^{off}) + \epsilon$$

ϵ is called the **regret**.

- This error metric is a bit 'unfair'. **Why?**
- Comparing online solution to best fixed solution in hindsight. ϵ can be negative!

INTUITION CHECK

$$f_1(\theta) = \underline{|\theta - 1000|} \quad f_2(\theta) = \underline{|\theta + 1000|}, \dots$$

What if for $i = 1, \dots, t$, $f_i(\theta) = |\theta - 1000|$ or $f_i(\theta) = |\theta + 1000|$ in an alternating pattern?

How small can the regret ϵ be? $\sum_{i=1}^t f_i(\bar{\theta}^{(i)}) \leq \sum_{i=1}^t f_i(\bar{\theta}^{\text{off}}) + \epsilon \approx 0$

$$\theta = 0$$

$$\min_{\theta} \sum_{i=1}^t f_i(\theta) = \sum_{i=1}^t f_i(0) = \underline{\underline{1000t}}$$

$$\theta_{10} = -1000 \quad \theta_{11} = +1000 \quad \theta_{12} = -1000$$

$$f_{10}(\theta_{10}) = 0 \quad f_{11}(\theta_{11}) = 0 \quad \dots$$

INTUITION CHECK

$$f_1(\theta) = |\theta - 1000| \quad f_2(\theta) = |\theta + 1000| \quad f_3(\theta) = |\theta + 1000|$$

What if for $i = 1, \dots, t$, $f_i(\theta) = |\theta - 1000|$ or $f_i(\theta) = |\theta + 1000|$ in an alternating pattern?

How small can the regret ϵ be? $\underbrace{\sum_{i=1}^t f_i(\vec{\theta}^{(i)})}_{\approx 1000t} \leq \underbrace{\sum_{i=1}^t f_i(\vec{\theta}^{\text{off}})}_{\underline{\underline{1000t}}} + \epsilon$

What if for $i = 1, \dots, t$, $f_i(\theta) = |\theta - 1000|$ or $f_i(\theta) = |\theta + 1000|$ in **no particular pattern**? How can any online learning algorithm hope to achieve small regret?

$$\begin{array}{c} \theta_i \\ \parallel \\ 0 \end{array} \quad f_i$$

$$\epsilon \geq 0$$

but ϵ small.

Assume that:

→ f_1, \dots, f_t are all convex.

→ Each f_i is G -Lipschitz (i.e., $\|\vec{\nabla} f_i(\vec{\theta})\|_2 \leq G$ for all $\vec{\theta}$.)

• $\|\underline{\theta}^{(1)} - \underline{\theta}^{off}\|_2 \leq R$ where $\theta^{(1)}$ is the first vector chosen.

$$\Theta' = \mathbf{0}$$

Assume that:

$$|\Theta| \leq 1000 \quad \text{+ steps} \quad \frac{1000 \pm 0}{2} \cdot t = 1000t \quad \text{SVRG}$$

$$\theta^i = 0 \quad f_i(\theta^i) = 1000$$

$$\theta^i : \pm 1000 \quad \mathbb{E} f_i(\theta^i) = 1000$$

- f_1, \dots, f_t are all convex.
- Each f_i is G -Lipschitz (i.e., $\|\vec{\nabla} f_i(\vec{\theta})\|_2 \leq G$ for all $\vec{\theta}$.)
- $\|\vec{\theta}^{(1)} - \vec{\theta}^{\text{off}}\|_2 \leq R$ where $\theta^{(1)}$ is the first vector chosen.

Online Gradient Descent

- Pick some initial $\vec{\theta}^{(1)}$.
 - Set step size $\eta = \frac{R}{G\sqrt{t}}$.
 - For $i = 1, \dots, t$
- Play $\vec{\theta}^{(i)}$ and incur cost $f_i(\vec{\theta}^{(i)})$.
- $\vec{\theta}^{(i+1)} = \vec{\theta}^{(i)} - \eta \cdot \vec{\nabla} f_i(\vec{\theta}^{(i)})$

Theorem – OGD on Convex Lipschitz Functions: For convex G -Lipschitz f_1, \dots, f_t , OGD initialized with starting point $\theta^{(1)}$ within radius R of θ^{off} , using step size $\eta = \frac{R}{G\sqrt{t}}$, has regret bounded by:

$$\frac{1}{t} \left[\underbrace{\sum_{i=1}^t f_i(\theta^{(i)})}_{\text{actual cost}} - \underbrace{\sum_{i=1}^t f_i(\theta^{off})}_{\text{optimal cost}} \right] \leq \frac{RG\sqrt{t}}{t} = \frac{RG}{\sqrt{t}}$$

Theorem – OGD on Convex Lipschitz Functions: For convex G -Lipschitz f_1, \dots, f_t , OGD initialized with starting point $\theta^{(1)}$ within radius R of θ^{off} , using step size $\eta = \frac{R}{G\sqrt{t}}$, has regret bounded by:

$$\frac{1}{t} \left[\sum_{i=1}^t f_i(\theta^{(i)}) - \sum_{i=1}^t f_i(\theta^{off}) \right] \leq \frac{RG\sqrt{t}}{t} = \frac{RG}{\sqrt{t}}$$

Upper bound on **average regret** goes to 0 and $t \rightarrow \infty$.

at step i set $\theta^{(i)}$ to the minimum of $\sum_{j=1}^i f_j(\theta)$

Theorem – OGD on Convex Lipschitz Functions: For convex G -Lipschitz f_1, \dots, f_t , OGD initialized with starting point $\theta^{(1)}$ within radius R of θ^{off} , using step size $\eta = \frac{R}{G\sqrt{t}}$, has regret bounded by:

$$\left[\sum_{i=1}^t f_i(\theta^{(i)}) - \sum_{i=1}^t f_i(\theta^{off}) \right] \leq RG\sqrt{t}$$

Upper bound on **average regret** goes to 0 and $t \rightarrow \infty$. No assumptions on f_1, \dots, f_t !