

CS-GY 6763: Lecture 9

Finish Second Order Conditions, Online and Stochastic Gradient Descent

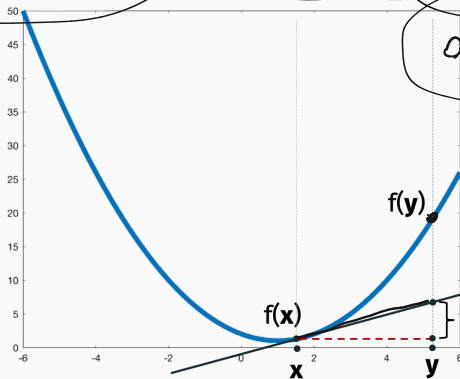
NYU Tandon School of Engineering, Prof. Christopher Musco

SECOND ORDER CONDITIONS

A function is α -strongly convex and β -smooth if for all x, y :

$$\frac{\alpha}{2} \|y - x\|_2^2 \leq [f(y) - f(x)] - \nabla f(x)^\top (y - x) \leq \frac{\beta}{2} \|y - x\|_2^2$$

$$\alpha \leq f''(x) \leq \beta$$



ALTERNATIVE DEFINITION OF SMOOTHNESS

Definition (β -smoothness)

A function f is β smooth if and only if, for all $\underline{x}, \underline{y}$

$$\|\nabla f(\underline{x}) - \nabla f(\underline{y})\|_2 \leq \beta \|\underline{x} - \underline{y}\|_2 \quad \rangle$$

I.e., the gradient function is a β -Lipschitz function.

We won't use this definition directly, but it's good to know.

Easy to prove equivalency to previous definition (see Lem. 3.4 in Bubeck's book).

$$y = x^\star$$

$$\|\nabla f(x) - \nabla f(x^\star)\|_2 \leq \beta \|x - x^\star\|_2$$

$$\|\nabla f(x)\|_2 \leq \beta \|x - x^\star\|_2$$

IMPROVING GRADIENT DESCENT

$$f(x^2) \leq f(x^*) + \epsilon$$

Having either an upper and lower bound on the second derivative helps convergence. Having both helps a lot.

Number of iterations for ϵ error:

	G-Lipschitz	β -smooth
R bounded start	$O\left(\frac{G^2 R^2}{\epsilon^2}\right)$	$O\left(\frac{\beta R^2}{\epsilon}\right)$
α -strong convex	$O\left(\frac{G^2}{\alpha \epsilon}\right)$	$O\left(\frac{\beta}{\alpha} \log(1/\epsilon)\right)$

Gradient descent for β -smooth functions:

- Select starting point $\mathbf{x}^{(0)}$, $\eta = 1/\beta$.
- For $i = 0, \dots, T$:
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$
- Return $\hat{\mathbf{x}} = \arg \min_i f(\mathbf{x}^{(i)})$.

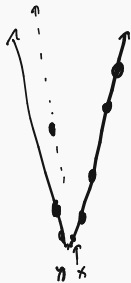
$$\|\nabla f(\mathbf{x})\|_r \leq G$$

$$\beta$$

$$\frac{G}{\sqrt{T}}$$

$$-\eta \nabla f(\mathbf{x}^{(i)})$$

Why do you think gradient descent might be faster when a function is β -smooth?



GUARANTEED PROGRESS

Previously learning rate/step size η depended on G . Now choose it based on β :

$$\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - \frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})$$

Progress per step of gradient descent:

$$1. [f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)})] - \underbrace{\nabla f(\mathbf{x}^{(t)})^T (\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)})}_{-1/\beta \nabla f(\mathbf{x}^{(t)})} \leq \frac{\beta}{2} \|\underbrace{\mathbf{x}^{(t)} - \mathbf{x}^{(t+1)}}_{+1/\beta \nabla f(\mathbf{x}^{(t)})^T \nabla f(\mathbf{x}^{(t)})}\|_2^2.$$

$$2. \underbrace{[f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)})]} + \underbrace{\frac{1}{\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2}_{-1/\beta \nabla f(\mathbf{x}^{(t)})^T \nabla f(\mathbf{x}^{(t)})} \leq \frac{\beta}{2} \underbrace{\|\frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})\|_2^2}_{= \frac{1}{2\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2}.$$

$$\left(3. f(\mathbf{x}^{(t)}) - f(\mathbf{x}^{(t+1)}) \geq \frac{1}{2\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2. \right) \quad f(\mathbf{x}^{(t+1)}) < f(\mathbf{x}^{(t)})$$

Where did we use convexity in this proof?

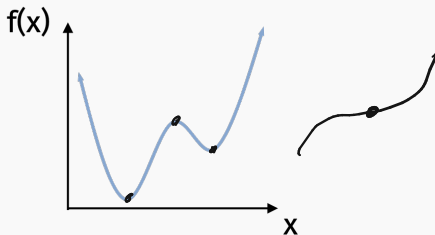
CONVERGENCE TO STATIONARY POINT

Theorem (Convergence to Near-Stationary Point)

For any β -smooth differentiable function f (convex or not), if we run GD for T steps, we can find a point \hat{x} such that:

$$\| \nabla f(\hat{x}) \|_2^2 \leq \frac{2\beta}{T} \left(\underbrace{f(x^{(0)})}_{\substack{\uparrow \\ T = \frac{2\beta}{\epsilon} f(x^{(0)} - f(x^*)}}) - \underbrace{f(x^*)}_{\epsilon} \right) \leq \epsilon$$

local/global minima - local/global maxima - saddle points



TELESCOPING SUM PROOF

We have that $\frac{1}{2\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq f(\mathbf{x}^{(t)}) - f(\mathbf{x}^{(t+1)})$. So:

$$\left(\sum_{t=0}^{T-1} \frac{1}{2\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \underbrace{f(\mathbf{x}^{(0)}) - f(\mathbf{x}^{(t)})}_{f(\mathbf{x}^{(0)}) - f(\mathbf{x}^{(1)}) + f(\mathbf{x}^{(1)}) - f(\mathbf{x}^{(2)}) + f(\mathbf{x}^{(2)}) - \dots}$$

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \frac{2\beta}{T} (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^{(T)})) \leq \frac{2\beta}{T} (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*))$$

$$\min_t \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \frac{2\beta}{T} (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*))$$

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\nabla f(\mathbf{x}^{(t)})\|_2^2$$

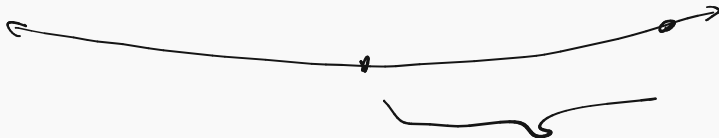
BACK TO CONVEX FUNCTIONS

For convex functions, we want to further prove that:

$$f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \leq \frac{2\beta R^2}{T}$$



Convex functions only have one stationary point: the global minimum \mathbf{x}^* . Ideally, we would argue that a near-stationary point is a near-minimizer. However, this isn't always the case!



CONVERGENCE GUARANTEE

Nevertheless, not too hard to obtain a proof from the progress condition. A concise version can be found on Page 15 in

(Garrigos and Gower's notes.)

$$T = O\left(\frac{L^2 R^2}{\epsilon^2}\right)$$

Theorem (GD convergence for β -smooth functions.)

Let f be a β smooth convex function and assume we have $\|\mathbf{x}^* - \mathbf{x}^{(0)}\|_2 \leq R$. If we run GD for T steps with $\eta = \frac{1}{\beta}$ we have:

$$\underline{f(\mathbf{x}^{(T)})} - \underline{f(\mathbf{x}^*)} \leq \frac{2\beta R^2}{T}$$

Corollary: If $T = O\left(\frac{\beta R^2}{\epsilon}\right)$ we have $\underline{f(\mathbf{x}^{(T)})} - \underline{f(\mathbf{x}^*)} \leq \epsilon$.

Note: This is not optimal! Can be improved to depend on $O(1/T^2)$ using a technique called acceleration.

CONVERGENCE GUARANTEE

What if f is both β -smooth and α -strongly convex? $f(x^{(T)}) \in f(x^*) - \epsilon$

Theorem (GD for β -smooth, α -strongly convex.) $f(x^{(T)}) - f(x^*)$

Let f be a β -smooth and α -strongly convex function. If we run GD for T steps (with step size $\eta = \frac{1}{\beta}$) we have:

$$\frac{1}{\beta} [f(x^{(T)}) - f(x^*)] \leq \underbrace{\|x^{(T)} - x^*\|_2^2}_{\leq R} \leq e^{-\frac{T\alpha}{\beta}} \underbrace{\|x^{(0)} - x^*\|_2^2}_{\leq R} \leq$$

$\kappa = \frac{\beta}{\alpha}$ is called the "condition number" of f .

Is it better if κ is large or small?

$$\beta \geq \alpha$$

$$\frac{\beta}{\alpha} \geq 1$$

$$e^{-T \frac{\alpha}{\beta}} = \epsilon$$

$$T \frac{\alpha}{\beta} = \log(1/\epsilon)$$

$$T = \frac{\beta}{\alpha} \log(1/\epsilon)$$

SMOOTH AND STRONGLY CONVEX

Converting to more familiar form: Using that fact the $\mathbf{x} = \mathbf{x}^*$
 $\nabla f(\mathbf{x}^*) = \mathbf{0}$ along with

$$\frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \leq [f(\mathbf{y}) - f(\mathbf{x})] - \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2,$$

we have: $\frac{\alpha}{2} \|\mathbf{x}^* - \mathbf{y}\|_2^2 \leq f(\mathbf{y}) - f(\mathbf{x}^*) \leq \frac{\beta}{2} \|\mathbf{x}^* - \mathbf{y}\|_2^2$

$$\frac{2}{\beta} [f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*)] \leq \|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2$$

$\mathbf{y} = \mathbf{x}^{(T)}$

We also assume

$$\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 \leq R^2.$$

CONVERGENCE GUARANTEE

Corollary (GD for β -smooth, α -strongly convex.)

Let f be a β -smooth and α -strongly convex function. If we run GD for T steps (with step size $\eta = \frac{1}{\beta}$) we have:

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \frac{\beta}{2} \underbrace{e^{-T \frac{\alpha}{\beta}} \cdot R^2}_{\log(1/\epsilon)}$$

Corollary: If $T = O\left(\frac{\beta}{\alpha} \log(R\beta/\epsilon)\right)$ we have:

$$\underline{f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*)} \leq \epsilon$$

$$T = O\left(\sqrt{\frac{\beta}{\alpha}} \cdot \log(R\beta/\epsilon)\right)$$

Only depend on $\log(1/\epsilon)$ instead of on $1/\epsilon$ or $1/\epsilon^2$!

Note: Can be further improved with acceleration!

After break we will prove the guarantee for the special case of:

$$f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2$$

Goal: Get some of the key ideas across, introduces important concepts like the Hessian, and show the connection between conditioning and linear algebra.

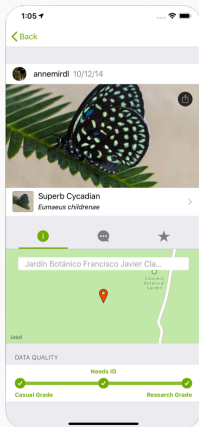
But first we will talk about online gradient descent and stochastic gradient descent.

- Basics of Online Learning + Optimization.
- Introduction to Regret Analysis.
- Application to analyzing Stochastic Gradient Descent.

Original motivation for online learning: Often need to train machine learning models on constantly updating/changing data. Do not want to restart from scratch.

Plant identification via iNaturalist app.

(California Academy of Science + National Geographic)



- When the app fails, image is classified via crowdsourcing (backed by huge network of amateurs and experts).
- Single model that is updated constantly, not retrained in batches.

EXAMPLE

Machine learning based email spam filtering.

```
MIME-Version: 1.0 Date: Mon, 7 Oct 2019
14:51:30 -0400 Message-ID: <CANV7LzUgXw==B-
39MLAN0Py19_jaaK60gmR04QCFHgyB0a@mail-spa
11.com> Subject: 92231 Reading Group, Meeting
2, tomorrow at 10am From: Christopher Musco
<cmusco@nyu.edu> To: algalids@nyu.edu Content-
Type: multipart/alternative;
boundary="0000000000078ec240594568a53" --
0000000000078ec240594568a53 Content-Type:
text/plain; charset="UTF-8" I hope everyone
had a good weekend! Tomorrow at *10am in 370
Jay St. #1114* we will meet for the second
instantiation of the CS-CV 92231 reading
group. Nick Feng will be leading a discussion
about the paper Simple Analyses of the Sparse
Johnson-Lindenstrauss Transform
<http://drops.dagstuhl.de/opus/volltexte/2018
/8305/pdf/OASiCS-SCOA-2018-15.pdf>. Please
read the abstract and introduction before the
meeting. Best, - CM *Christopher Musco,
Assistant Professor* *New York University,
Tandon School of Engineering* *4011 576
2541* --0000000000078ec240594568a53 Content-
Type: text/html; charset="UTF-8" Content-
Transfer-Encoding: quoted-printable
```

sender on graylist 1

sender on graylist 2

Bag-of-words

sender flagged by x users

sender ip flagged by x users

1
0
:
0
0
1
1
0
1

Markers for spam change overtime, so model might change.

EXAMPLE

Machine learning based email spam filtering.



New Report

Elon Musk's New Electricity Saving Invention Has Residents Saving Up to 90% Off Their Monthly Electric Bill. Electric Power Companies Are Demanding It Be Banned Immediately!

Do not pay your electric bill until you read this. As electricity prices continue to rise in , I realize that not everyone can afford solar panels, so we wanted to come up with a way that **EVERYONE** can save tons of money on their electric bill. Hurry up and learn this trick before the power companies get their way and it's gone." - Elon Musk

[READ MORE](#)

Markers for spam change overtime, so model might change.

ONLINE LEARNING FRAMEWORK

Choose some model M_x parameterized by parameters x and some loss function ℓ . At time steps $1, \dots, T$, receive data vectors $\underline{a}^{(1)}, \dots, \underline{a}^{(T)}$. $M_x(a)$

- At each time step, we pick ("play") a parameter vector $\underline{x}^{(i)}$.
- Make prediction $\underline{\tilde{y}}^{(i)} = M_{\underline{x}^{(i)}}(\underline{a}_i)$.
- Then told true value or label $\underline{y}^{(i)}$. Possibly use this information to choose a new $\underline{x}^{(i+1)}$.
- Goal is to minimize cumulative loss:

$$L = \sum_{i=1}^T \ell(\underline{x}^{(i)}, \underline{a}^{(i)}, y^{(i)}) \quad \downarrow \quad M_{x^{(i)}}(a^{(i)}) - y^{(i)}$$
$$L(x) = \sum_{i=1}^T \ell(x, a^{(i)}, y^{(i)})$$

For example, for a regression problem we might use the ℓ_2 loss:

$$\ell(\underline{x}^{(i)}, \underline{a}^{(i)}, y^{(i)}) = \left(\underline{M}_{\underline{x}^{(i)}}(\underline{a}_i) - y^{(i)} \right)^2.$$

$f(x)$

For classification, we could use logistic/cross-entropy loss.

Regret

Abstraction as optimization problem: Instead of a single objective function f , we have a unknown function $\underline{f}_1, \dots, \underline{f}_T : \mathbb{R}^d \rightarrow \mathbb{R}$ for each time step.

- For time step $i \in \underline{1}, \dots, \underline{T}$, select vector $\underline{\mathbf{x}}^{(i)}$.
- Observe \underline{f}_i and pay cost $\underline{f}_i(\underline{\mathbf{x}}^{(i)})$
- Goal is to minimize $\underline{\sum_{i=1}^T \underline{f}_i(\underline{\mathbf{x}}^{(i)})}$.

We make no assumptions that $\underline{f}_1, \dots, \underline{f}_T$ are related to each other at all!

REGRET BOUND

In offline optimization, we wanted to find \hat{x} satisfying $f(\hat{x}) \leq \min_x f(x) + \epsilon$. Ask for a similar thing here.

Objective: Choose $\underline{x}^{(1)}, \dots, \underline{x}^{(T)}$ so that:

$$\frac{1}{T} \sum_{i=1}^T f_i(x^{(i)}) \leq \underbrace{\left[\min_x \sum_{i=1}^T f_i(x) \right]}_{\text{regret}} + \frac{\epsilon}{T}$$

$\epsilon \leq \sqrt{T}$
 $\leq \log(T)$

Here ϵ is called the **regret** of our solution sequence $x^{(0)}, \dots, x^{(T)}$.

We typically ϵ to be growing sublinearly in T .

$$x^* = \underset{x}{\operatorname{argmin}} \sum_{i=1}^T f_i(x) \quad x^* = \text{"best solution in hindsight"} \\ = \text{"best fixed solution"}. \quad 21$$

Regret compares to the best fixed solution in hindsight.

$$\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \leq \left[\min_{\mathbf{x}} \sum_{i=1}^T f_i(\mathbf{x}) \right] + \epsilon.$$

It is very possible that $\left[\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \right] < \left[\min_{\mathbf{x}} \sum_{i=1}^T f_i(\mathbf{x}) \right]$. Could we hope for something stronger?

Exercise: Argue that the following is impossible to achieve:

$$\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \leq \left[\sum_{i=1}^T \min_{\mathbf{x}} f_i(\mathbf{x}) \right] + \epsilon. \quad \rightarrow O(T)$$

\mathbf{x}_i $\sum_{i=1}^T f_i(\mathbf{x}_i)$

HARD EXAMPLE FOR ONLINE OPTIMIZATION

Convex functions:

1/2

$$f_1(x) = |x - h_1|$$

\vdots

$$f_T(x) = |x - h_T|$$

If $h_1 = 0$

choose $x_1^* = 0$

If $h_2 = 1$

choose $x_2^* = 1$

where h_1, \dots, h_T are i.i.d. uniform $\{0, 1\}$.

$$\sum_{i=1}^T \min_x f_i(x) = 0$$

$$\mathbb{E} \left(\sum_{i=1}^T f_i(x^{(i)}) \right) \geq T/2$$

$$\begin{aligned} \mathbb{E} f_i(x^{(i)}) &= \frac{1}{2} |0 - x^{(i)}| + \frac{1}{2} |1 - x^{(i)}| \\ &= \frac{1}{2} |x^{(i)}| + \frac{1}{2} |1 - x^{(i)}| \\ &\geq \frac{1}{2} \end{aligned}$$

$$\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \leq \left[\min_{\mathbf{x}} \sum_{i=1}^T f_i(\mathbf{x}) \right] + \epsilon.$$

Beautiful balance:

- Either f_1, \dots, f_T are similar or changing slowly, so we can learn/predict f_i from earlier functions.
- Or f_1, \dots, f_T are very different, in which case $\min_{\mathbf{x}} \sum_{i=1}^T f_i(\mathbf{x})$ is large, so regret bound is easy to achieve.
- Or we live somewhere in the middle.

Follow-the-leader algorithm:

$$x^{(i)}$$

- Choose $x^{(0)}$.

$$f_1, \dots, f_{i-1}$$

- For $i = 1, \dots, T$:

- Let $x^{(i)} = \arg \min_x \sum_{j=1}^{i-1} f_j(x)$.

$$\arg \min_x \sum_{j=1}^{i-1} \alpha^j f_j(x)$$

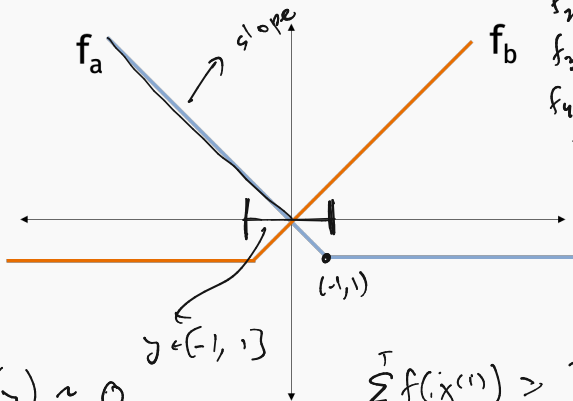
- Play $x^{(i)}$.

- Observe f_i and incur cost $f_i(x^{(i)})$.

Simple and intuitive, but there are two issues with this approach. One is computational, one is related to the accuracy.

FOLLOW-THE-LEADER

Hard case with regret $O(T)$



$f_1 = f_a$
 $f_2 = f_b$
 $f_3 = f_b$
 $f_4 = f_a$
 \vdots
 f_b
 \vdots
 f_b
 \vdots

$$\sum_{i=1}^T f_i(y) \approx 0$$

$$\sum_{i=1}^T f_i(x^{(i)}) > T/2$$

<https://www.desmos.com/calculator/3t8bfowo3j>

ONLINE GRADIENT DESCENT

Online Gradient descent:

- Choose $\mathbf{x}^{(1)}$ and η .
- For $i = \underline{1}, \dots, \underline{T}$:
 - Play $\mathbf{x}^{(i)}$. !
 - Observe f_i and incur cost $f_i(\mathbf{x}^{(i)})$.
 - $\mathbf{x}^{(i+1)}$ = $\mathbf{x}^{(i)}$ - $\eta \nabla f_i(\mathbf{x}^{(i)})$

return 3:35

If $f_1, \dots, f_T = f$ are all the same, this is the same as regular gradient descent. We update parameters using the gradient ∇f at each step.

ONLINE GRADIENT DESCENT (OGD)

$\mathbf{x}^* = \arg \min_{\mathbf{x}} \sum_{i=1}^T f_i(\mathbf{x})$ (the offline optimum)

Assume:

"best hindsight solution" $\|\nabla f(\mathbf{x})\|_2 \leq G$

- f_1, \dots, f_T are all convex.)
- Each is G-Lipschitz: for all \mathbf{x}, i , $\|\nabla f_i(\mathbf{x})\|_2 \leq G$.
- Starting radius: $\|\mathbf{x}^* - \mathbf{x}^{(1)}\|_2 \leq R$.

Online Gradient descent:

- Choose $\mathbf{x}^{(1)}$ and $\eta = \frac{R}{G\sqrt{T}}$.
- For $i = 1, \dots, T$:
 - Play $\mathbf{x}^{(i)}$.
 - Observe f_i and incur cost $f_i(\mathbf{x}^{(i)})$.
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f_i(\mathbf{x}^{(i)})$

Theorem (OGD Regret Bound)

If the conditions of the previous slide hold, then after T steps,

$$\epsilon = \underbrace{\frac{1}{T} \left[\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \right]}_{\text{Average regret}} - \underbrace{\frac{1}{T} \left[\sum_{i=1}^T f_i(\mathbf{x}^*) \right]}_{\text{Minimum possible regret}} \leq RG\sqrt{T}.$$

Average regret overtime is bounded by $\frac{\epsilon}{T} \leq \frac{RG}{\sqrt{T}}$.

Goes $\rightarrow 0$ as $T \rightarrow \infty$.

All this with (no assumptions on how f_1, \dots, f_T) relate to each other! They could have even been chosen **adversarially** – e.g. with f_i depending on our choice of \mathbf{x}_i and all previous choices.

$$f_i \quad \mathbf{x}^{(i-1)}$$

Theorem (OGD Regret Bound)

If the conditions of the previous slide hold, then after T steps,
 $\epsilon = \left[\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \right] - \left[\sum_{i=1}^T f_i(\mathbf{x}^*) \right] \leq RG\sqrt{T}.$

Claim 1: For all $i = 1, \dots, T$,

$$\left(f_i(\mathbf{x}^{(i)}) - \underline{f_i(\mathbf{x}^*)} \leq \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2} \right)$$

(Same proof for standard GD. Only uses convexity of f_i .)

$$\nearrow \mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f_i(\mathbf{x}^{(i)})$$

Theorem (OGD Regret Bound)

\bar{x}

After T steps, $\epsilon = \left[\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \right] - \left[\sum_{i=1}^T f_i(\mathbf{x}^*) \right] \leq RG\sqrt{T}$.

Claim 1: For all $i = 1, \dots, T$,

$$\underline{f_i(\mathbf{x}^{(i)})} - \underline{f_i(\mathbf{x}^*)} \leq \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$$

Telescoping Sum:

$$\begin{aligned} \sum_{i=1}^T \left[\underline{f_i(\mathbf{x}^{(i)})} - \underline{f_i(\mathbf{x}^*)} \right] &\leq \frac{\|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{T\eta G^2}{2} \\ &\leq \frac{R^2}{2\eta} + \frac{T\eta G^2}{2} = \frac{RG\sqrt{T}}{2} + \frac{RG\sqrt{T}}{2} \\ &= RG\sqrt{T} \end{aligned}$$

$\mu = \frac{R}{G\sqrt{T}}$

That's it!

STOCHASTIC GRADIENT DESCENT (SGD)

Efficient offline optimization method for functions f with (finite sum structure):

$$L(x) = \sum_{i=1}^n \ell(a^{(i)}, y^{(i)}, x)$$

of training examples

$$f(x) = \sum_{i=1}^n f_i(x).$$

Goal is to find \hat{x} such that $f(\hat{x}) \leq f(x^*) + \epsilon$.

- The most widely use optimization algorithm in modern machine learning.
- Easily analyzed as a special case of online gradient descent!

STOCHASTIC GRADIENT DESCENT

Recall the machine learning setup. In empirical risk minimization, we can typically write:

$$f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$$

$$\begin{pmatrix} \mathbf{a}^{(1)} \\ \vdots \\ \mathbf{a}^{(n)} \end{pmatrix} \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{pmatrix}$$

where f_i is the loss function for a particular data example $(\mathbf{a}^{(i)}, y^{(i)})$.

Example: least squares linear regression.

$$= \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$$

$$O(nd)$$

$$2\mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{y})$$

$$f(\mathbf{x}) = \sum_{i=1}^n (\mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)})^2$$

f_i

Note that by linearity, $\nabla f(\mathbf{x}) = \sum_{i=1}^n \nabla f_i(\mathbf{x})$.

$$2(\mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)}) \cdot \mathbf{a}^{(i)} \quad O(d)$$

STOCHASTIC GRADIENT DESCENT

Main idea: Use random approximate gradient in place of actual gradient. $\nabla f(x)$

Pick random $j \in 1, \dots, n$ and update x using $\nabla f_j(x)$.

$$\mathbb{E} [\nabla f_j(x)] = \frac{1}{n} \nabla f(x).$$

$$\sum_{j=1}^n \frac{1}{n} \nabla f_j(x) = \frac{1}{n} \sum_{j=1}^n \nabla f_j(x) = \frac{1}{n} \nabla \left[\sum_{j=1}^n f_j(x) \right] = \frac{1}{n} \nabla f(x)$$

$\nabla f_j(x)$ is an unbiased estimate for the true gradient $\nabla f(x)$, but can typically be computed in a $1/n$ fraction of the time!

Trade slower convergence for cheaper iterations.

γn
(cost)

STOCHASTIC GRADIENT DESCENT

Stochastic first-order oracle for $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$.

- **Function Query:** For any chosen j, \mathbf{x} , return $f_j(\mathbf{x})$
- **Gradient Query:** For any chosen j, \mathbf{x} , return $\nabla f_j(\mathbf{x})$

Stochastic Gradient descent:

- Choose starting vector $\mathbf{x}^{(1)}$, step size η

$$j_i \sim 1, \dots, n$$

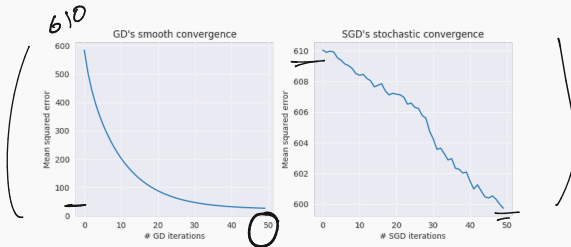
- For $i = 1, \dots, T$:

- Pick random $j_i \in 1, \dots, n$.

- $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f_{j_i}(\mathbf{x}^{(i)})$

- Return $\hat{\mathbf{x}} = \frac{1}{T} \sum_{i=1}^T \mathbf{x}^{(i)}$

VISUALIZING SGD



STOCHASTIC GRADIENT DESCENT

Assume:

- Finite sum structure: $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$, with f_1, \dots, f_n all convex.
- Lipschitz functions: for all \mathbf{x}, j , $\|\nabla f_j(\mathbf{x})\|_2 \leq \frac{G'}{n}$.
 - What does this imply about Lipschitz constant of f ?

- Starting radius: $\|\mathbf{x}^* - \mathbf{x}^{(1)}\|_2 \leq R$.

$$\begin{aligned}\|\nabla f(\mathbf{x})\|_2 &= \left\| \sum_{j=1}^n \nabla f_j(\mathbf{x}) \right\|_2 \\ &\leq \sum_{j=1}^n \|\nabla f_j(\mathbf{x})\|_2 \\ &\leq n \cdot \frac{G'}{n} = G'\end{aligned}$$

Stochastic Gradient descent:

- Choose $\mathbf{x}^{(1)}$, steps T , step size $\eta = \frac{R}{G'\sqrt{T}}$.
- For $i = 1, \dots, T$:
 - Pick random $j_i \in 1, \dots, n$.
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f_{j_i}(\mathbf{x}^{(i)})$
- Return $\hat{\mathbf{x}} = \frac{1}{T} \sum_{i=1}^T \mathbf{x}^{(i)}$

Approach: View as online gradient descent run on function sequence f_{j_1}, \dots, f_{j_T} .

uniform $1, \dots, n$

Claim (SGD Convergence)

After $T = \frac{R^2 G^2}{\epsilon^2}$ iterations:

$$\mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \leq \epsilon.$$

Will prove using:

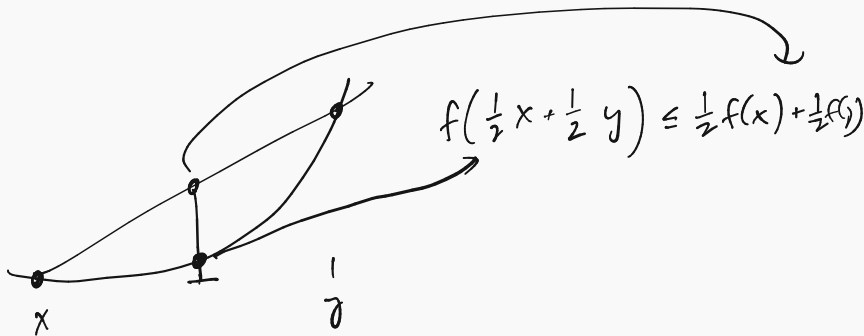
1. Black-box result for online gradient descent (already proven).
2. The fact that $\underline{n} \cdot \mathbb{E}[f_{j_i}(\mathbf{x}^{(i)})] = \underline{f(\mathbf{x}^{(i)})}$.
3. Jensen's inequality.

$$n \cdot \mathbb{E} f_{j_i}(y) = n \cdot \sum_{i=1}^n \frac{1}{n} f_i(y) = \sum_{i=1}^n f_i(y) = f(y)$$

JENSEN'S INEQUALITY

For a convex function f and points $\mathbf{x}^{(1)}$, ..., $\mathbf{x}^{(t)}$

$$f\left(\frac{1}{t} \cdot \mathbf{x}^{(1)} + \dots + \frac{1}{t} \cdot \mathbf{x}^{(t)}\right) \leq \frac{1}{t} \cdot \underline{f(\mathbf{x}^{(1)})} + \dots + \frac{1}{t} \cdot \underline{f(\mathbf{x}^{(t)})}$$



STOCHASTIC GRADIENT DESCENT ANALYSIS

Claim (SGD Convergence)

After $T = \frac{R^2 G^2}{\epsilon^2}$ iterations:

$$\hat{x} = \frac{1}{T} \sum_{i=1}^T x^{(i)}$$

$$\mathbb{E}[f(\hat{x}) - f(x^*)] \leq \epsilon.$$

Claim 1:

$$f(\hat{x}) - f(x^*) \leq \frac{1}{T} \sum_{i=1}^T [f(x^{(i)}) - f(x^*)]$$

Prove using Jensen's Inequality:

$$\begin{aligned} f\left(\frac{1}{T} \sum_{i=1}^T x^{(i)}\right) - \frac{1}{T} \sum_{i=1}^T f(x^*) &\leq \frac{1}{T} \sum_{i=1}^T f(x^{(i)}) - \frac{1}{T} \sum_{i=1}^T f(x^*) \\ &= \frac{1}{T} \sum_{i=1}^T [f(x^{(i)}) - f(x^*)] \end{aligned}$$

STOCHASTIC GRADIENT DESCENT ANALYSIS

Claim (SGD Convergence)

After $T = \frac{R^2 G'^2}{\epsilon^2}$ iterations:

$$\mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \leq \epsilon.$$

$$\mathbb{E}[n \cdot f_{j_i}(\mathbf{x}^{(i)})] = \mathbb{E}[f(\mathbf{x}^{(i)})]$$

$$\rightarrow n \cdot f_{j_i}$$

$$\begin{aligned} \mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] &\leq \frac{1}{T} \sum_{i=1}^T \mathbb{E}[f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*)] \\ &= \frac{1}{T} \sum_{i=1}^T n \mathbb{E}[f_{j_i}(\mathbf{x}^{(i)}) - f_{j_i}(\mathbf{x}^*)] \end{aligned}$$

Claim (SGD Convergence)

$$f_{j_1} \dots f_{j_T}$$

After $T = \frac{R^2 G'^2}{\epsilon^2}$ iterations:

$$\mathbb{E} [f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \leq \epsilon.$$

$$\begin{aligned} \mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] &\leq \frac{1}{T} \sum_{i=1}^T \mathbb{E} [f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*)] \\ &= \frac{1}{T} \sum_{i=1}^T n \mathbb{E} [f_{j_i}(\mathbf{x}^{(i)}) - f_{j_i}(\mathbf{x}^*)] = \frac{n}{T} \mathbb{E} \sum_{i=1}^T (f_{j_i}(\mathbf{x}^{(i)}) - f_{j_i}(\mathbf{x}^*)) \\ &\leq \frac{n}{T} \cdot \mathbb{E} \left[\sum_{i=1}^T f_{j_i}(\mathbf{x}^{(i)}) - f_{j_i}(\mathbf{x}^{\text{offline}}) \right], \end{aligned}$$

where $\mathbf{x}^{\text{offline}} = \arg \min_{\mathbf{x}} \sum_{i=1}^T f_{j_i}(\mathbf{x})$.

$$\sum_{i=1}^T f_{j_i}(\mathbf{x}^*) \geq \sum_{i=1}^T f_{j_i}(\mathbf{x}^{\text{offline}})$$

STOCHASTIC GRADIENT DESCENT ANALYSIS

Claim (SGD Convergence)

After $T = \frac{R^2 G'^2}{\epsilon^2}$ iterations:

$$\mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \leq \epsilon.$$

$$\mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \leq \frac{1}{T} \sum_{i=1}^T \mathbb{E}[f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*)]$$

$$= \frac{1}{T} \sum_{i=1}^T n \mathbb{E}[f_{j_i}(\mathbf{x}^{(i)}) - f_{j_i}(\mathbf{x}^*)]$$

$$\leq \frac{n}{T} \cdot \mathbb{E} \left[\sum_{i=1}^T f_{j_i}(\mathbf{x}^{(i)}) - f_{j_i}(\mathbf{x}^{\text{offline}}) \right]$$

$$\leq \frac{n}{T} \cdot \left(R \cdot \frac{G'}{n} \cdot \sqrt{T} \right) \quad (\text{by OGD guarantee.})$$

$$= \frac{R G'}{\sqrt{T}}$$

$$T = \frac{R^2 G'^2}{\epsilon^2}$$

$$R G' \sqrt{T}$$

$$\frac{R \cdot G'}{n} \cdot \sqrt{T}$$

STOCHASTIC VS. FULL BATCH GRADIENT DESCENT

Number of iterations for error ϵ :

$$\|\nabla f(x)\|_2 \leq G$$

- Gradient Descent: $T = \frac{R^2 G^2}{\epsilon^2}$
- Stochastic Gradient Descent: $T = \frac{R^2 G'^2}{\epsilon^2}$

Always have $G \leq G'$:

$$\begin{aligned}
 G &= \max_x \|\nabla f(x)\|_2 \leq \max_x (\|\nabla f_1(x)\|_2 + \dots + \|\nabla f_n(x)\|_2) \\
 &\leq \max_x (\|\nabla f_1(x)\|_2) + \dots + \max_x (\|\nabla f_n(x)\|_2) \\
 &\leq n \cdot \frac{G'}{n} = \underline{\underline{G'}}
 \end{aligned}$$

So GD converges strictly faster than SGD.

But for a fair comparison:

$$f = f_1 + \dots + f_n$$

$$G = G'$$

- SGD cost = (# of iterations) $\cdot O(1)$
- GD cost = (# of iterations) $\cdot O(n)$

STOCHASTIC VS. FULL BATCH GRADIENT DESCENT

We always have $G \leq G'$. When it is much smaller then GD will perform better. When it is closer to this upper bound, SGD will perform better.

What is an extreme case where $G = G'$?

STOCHASTIC VS. FULL BATCH GRADIENT DESCENT

What if each gradient $\nabla f_i(\mathbf{x})$ looks like random vectors in \mathbb{R}^d ?

E.g. with $\mathcal{N}(0, 1)$ entries?

$\nabla f_i(\mathbf{x})$ is random Gaussian of dimension d .

$$\mathbb{E} [\|\nabla f_i(\mathbf{x})\|_2^2] = d \quad \|\nabla f_i(\mathbf{x})\|_1 \approx \sqrt{d}$$

$$\mathbb{E} [\|\nabla f(\mathbf{x})\|_2^2] = \mathbb{E} \left[\left\| \sum_{i=1}^n \nabla f_i(\mathbf{x}) \right\|_2^2 \right] = nd \quad \|\nabla f(\mathbf{x})\|_1 \approx \sqrt{nd}$$

$\sum_{i=1}^n \mathcal{N}(0, 1) = \mathcal{N}(0, n)$

$$\frac{G'}{n} \approx \sqrt{d}$$

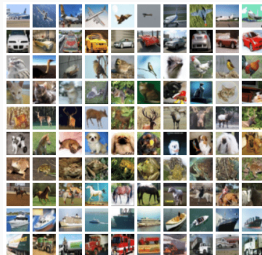
$$G' \approx n\sqrt{d}$$

$$G \approx \sqrt{nd}$$

$$G \leq \frac{1}{\sqrt{n}} G'$$

STOCHASTIC VS. FULL BATCH GRADIENT DESCENT

Takeaway: SGD performs better when there is more structure or repetition in the data set.



PRECONDITIONING

Main idea: Instead of minimizing $f(\mathbf{x})$, find another function $g(\mathbf{x})$ with the same minimum but which is better suited for first order optimization (e.g., is smoother, or has a smaller conditioner number).

Claim: Let $h(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be an invertible function. Let $g(\mathbf{x}) = f(h(\mathbf{x}))$. Then

$$\min_{\mathbf{x}} f(\mathbf{x}) = \min_{\mathbf{y}} g(\mathbf{y}) \quad \text{and} \quad \arg \min_{\mathbf{x}} f(\mathbf{x}) = h \left(\arg \min_{\mathbf{y}} g(\mathbf{y}) \right).$$

First Goal: We need $g(\mathbf{x})$ to still be convex.

Claim: Let \mathbf{P} be an invertible $d \times d$ matrix and let $g(\mathbf{x}) = f(\mathbf{P}\mathbf{x})$.

$g(\mathbf{x})$ is always convex.

Second Goal:

$g(\mathbf{x})$ should have better condition number κ than $f(\mathbf{x})$.

Example:

- $f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2$. $\kappa_f = \frac{\lambda_1(\mathbf{A}^T\mathbf{A})}{\lambda_d(\mathbf{A}^T\mathbf{A})}$.
- $g(\mathbf{x}) = \|\mathbf{APx} - \mathbf{b}\|_2^2$. $\kappa_g = \frac{\lambda_1(\mathbf{P}^T\mathbf{A}^T\mathbf{AP})}{\lambda_d(\mathbf{P}^T\mathbf{A}^T\mathbf{AP})}$.

Third Goal: \mathbf{P} should be easy to compute.

Many, many problem specific preconditioners are used in practice. Their design is usually a heuristic process.

Example: Diagonal preconditioner.

- Let $\mathbf{D} = \text{diag}(\mathbf{A}^T \mathbf{A})$
- Intuitively, we roughly have that $\mathbf{D} \approx \mathbf{A}^T \mathbf{A}$.
- Let $\mathbf{P} = \sqrt{\mathbf{D}^{-1}}$

\mathbf{P} is often called a **Jacobi preconditioner**. Often works very well in practice!

DIAGONAL PRECONDITIONER

A =

-734	1	33	9111	0
-31	-2	108	5946	-19
232	-1	101	3502	10
426	0	-65	12503	9
-373	0	26	9298	0
-236	-2	-94	2398	-1
2024	0	-132	-6904	-25
-2258	-1	92	-6516	6
2229	0	0	11921	-22
338	1	-5	-16118	-23

```
>> cond(A'*A)
```

ans =

8.4145e+07

```
>> P = sqrt(inv(diag(diag(A'*A))));
```

```
>> cond(P*A'*A*P)
```

ans =

10.3878

Another view: If $g(\mathbf{x}) = f(\mathbf{P}\mathbf{x})$ then $\nabla g(\mathbf{x}) = \mathbf{P}^T \nabla f(\mathbf{P}\mathbf{x})$.

$\nabla g(\mathbf{x}) = \mathbf{P} \nabla f(\mathbf{P}\mathbf{x})$ when \mathbf{P} is symmetric.

Gradient descent on g :

- For $t = 1, \dots, T$,
 - $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \mathbf{P} [\nabla f(\mathbf{P}\mathbf{x}^{(t)})]$
- Return $\mathbf{P}\mathbf{x}^{(T)}$

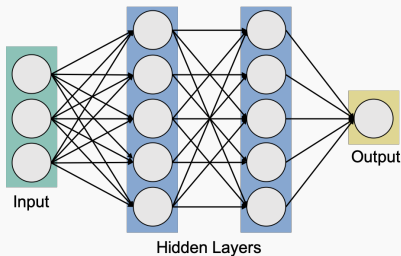
Gradient descent on g :

- For $t = 1, \dots, T$,
 - $\mathbf{y}^{(t+1)} = \mathbf{y}^{(t)} - \eta \mathbf{P}^2 [\nabla f(\mathbf{y}^{(t)})]$

When \mathbf{P} is diagonal, this is just gradient descent with a different step size for each parameter!

Algorithms based on this idea:

- AdaGrad
- RMSprop
- Adam optimizer



(Pretty much all of the most widely used optimization methods for training neural networks.)

COORDINATE DESCENT

Main idea: Trade slower convergence (more iterations) for cheaper iterations.

Stochastic Gradient Descent: When $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$, approximate $\nabla f(\mathbf{x})$ with $\nabla f_i(\mathbf{x})$ for randomly chosen i .

Main idea: Trade slower convergence (more iterations) for cheaper iterations.

Stochastic Coordinate Descent: Only compute a single random entry of $\nabla f(\mathbf{x})$ on each iteration:

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ \frac{\partial f}{\partial x_2}(\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_d}(\mathbf{x}) \end{bmatrix} \qquad \nabla_{if}(\mathbf{x}) = \begin{bmatrix} 0 \\ \frac{\partial f}{\partial x_i}(\mathbf{x}) \\ \vdots \\ 0 \end{bmatrix}$$

Update: $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} + \eta \nabla_{if}(\mathbf{x}^{(t)})$.