CS-GY 6763: Lecture 9 Finish Second Order Conditions, Online and Stochastic Gradient Descent

NYU TandonSchool of Engineering, Prof. Christopher Musco

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

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



Definition (β -smoothness) A function f is β smooth if and only if, for all \mathbf{x}, \mathbf{y} $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \le \beta \|\mathbf{x} - \mathbf{y}\|_2$

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**). Having <u>either</u> an upper and lower bound on the second derivative helps convergence. Having both helps a lot.

Number of iterations for ϵ error:



Gradient descent for β -smooth functions:

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

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

Previously learning rate/step size η depended on G. Now choose it based on β : $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - \frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})$

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

2. $[f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)})] + \frac{1}{\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \le \frac{\beta}{2} \|\frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})\|_2^2$.

3. $f(\mathbf{x}^{(t)}) - f(\mathbf{x}^{(t+1)}) \ge \frac{1}{2\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2$.

Where did we use convexity in this proof?

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 $\hat{\mathbf{x}}$ such that:

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

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



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

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

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

$$\min_{t} \|\nabla f(\mathbf{x}^{(t)})\|_{2}^{2} \leq \frac{2\beta}{T} \left(f(\mathbf{x}^{(0)}) - f(\mathbf{x}^{*}) \right)$$

For convex functions, we want to further prove that:

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

Convex functions only have one stationary point: the global minimum **x***. Ideally, we would argue that a <u>near-stationary</u> point is a <u>near-minimizer</u>. However, this isn't always the case!

Nevertheless, not to hard to obtain a proof from the progress condition. A concise version can be found on Page 15 in Garrigos and Gower's notes.

Theorem (GD convergence for β -smooth functions.)

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

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

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

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

What if f is both β -smooth and α -strongly convex?

Theorem (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:

$$\|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2 \le e^{-T\frac{lpha}{eta}} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2$$

 $\kappa = \frac{\beta}{\alpha}$ is called the "condition number" of *f*. Is it better if κ is large or small? Converting to more familiar form: Using that fact the $\nabla f(x^*) = 0$ along with

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

we have:

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

We also assume

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

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:

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \le \frac{\beta}{2}e^{-T\frac{\alpha}{\beta}} \cdot R^2$$

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

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

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:

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

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 <u>online gradient descent</u> and <u>stochastic gradient descent</u>.

ONLINE AND STOCHASTIC GRADIENT DESCENT

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

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.

Machine learning based email spam filtering.



Markers for spam change overtime, so model might change.

EXAMPLE

Machine learning based email spam filtering.

Re:SAFTY CORONA VIRUS AWARENESS WHO				
wo	World Health Organization			
	World Health Organization			
0	Dear Sir,			
(Go through the attached document on safety measures regarding the spreading of corona virus.			
(Click on the button below to download			
	Safety measures			
1	Symptoms common symptoms include fever,coughcshortness of breath and breathing difficulties.			
ę	Regards,			
5	Dr. Stella Chungong Specialist wuhan-virus-advisory			

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 <u>EVERVONE</u> 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.

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

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

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

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

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

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

Abstraction as optimization problem: Instead of a single objective function f, we have a <u>unknown</u> function $f_1, \ldots, f_T : \mathbb{R}^d \to \mathbb{R}$ for each time step.

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

We make <u>no assumptions</u> that f_1, \ldots, f_T are related to each other at all!

REGRET BOUND

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

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

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

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

We typically ϵ to be growing <u>sublinearly</u> in *T*.

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:

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

Convex functions:

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

$$\vdots$$

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

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

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

Beautiful balance:

- Either f_1, \ldots, f_T are similar or changing slowly, so we can learn/predict f_i from earlier functions.
- Or f_1, \ldots, 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:

- Choose $\mathbf{x}^{(0)}$.
- For i = 1, ..., T:
 - Let $\mathbf{x}^{(i)} = \arg \min_{\mathbf{x}} \sum_{j=1}^{i-1} f_j(\mathbf{x})$.
 - Play $\mathbf{x}^{(i)}$.
 - Observe f_i and incur cost $f_i(\mathbf{x}^{(i)})$.

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

Hard case with regret O(T):



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

Online Gradient descent:

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

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

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

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

Assume:

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

Online Gradient descent:

- Choose $\mathbf{x}^{(1)}$ and $\eta = \frac{R}{G\sqrt{T}}$.
- For i = 1, ..., 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, $\epsilon = \left[\sum_{i=1}^{T} f_i(\mathbf{x}^{(i)})\right] - \left[\sum_{i=1}^{T} f_i(\mathbf{x}^*)\right] \le RG\sqrt{T}.$

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, \ldots, 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.

Theorem (OGD Regret Bound)

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

Claim 1: For all i = 1, ..., T,

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

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

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

Claim 1: For all
$$i = 1, ..., T$$
,

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

Telescoping Sum:

$$\sum_{i=1}^{T} \left[f_i(\mathbf{x}^{(i)}) - f_i(\mathbf{x}^*) \right] \le \frac{\|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{T\eta G^2}{2}$$
$$\le \frac{R^2}{2\eta} + \frac{T\eta G^2}{2}$$

That's it!

Efficient <u>offline</u> optimization method for functions *f* with <u>finite</u> <u>sum structure</u>:

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

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

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

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

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

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

Example: least squares linear regression.

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

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

Main idea: Use random approximate gradient in place of actual gradient.

Pick <u>random</u> $j \in 1, ..., n$ and update **x** using $\nabla f_j(\mathbf{x})$.

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

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

Trade slower convergence for cheaper iterations.

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
- For i = 1, ..., T:
 - Pick random $j_i \in 1, \ldots, n$.
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} \eta \nabla f_{j_i}(\mathbf{x}^{(i)})$
- 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, \ldots, 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: $\|\mathbf{x}^* \mathbf{x}^{(1)}\|_2 \le R$.

Stochastic Gradient descent:

- Choose $\mathbf{x}^{(1)}$, steps *T*, step size $\eta = \frac{R}{G'\sqrt{T}}$.
- For i = 1, ..., T:
 - Pick random $j_i \in 1, \ldots, n$.

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

• 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}, \ldots, f_{j_T} .

Claim (SGD Convergence) After $T = \frac{R^2G'^2}{\epsilon^2}$ iterations: $\mathbb{E} [f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \le \epsilon.$

Will prove using:

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

JENSEN'S INEQUALITY

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

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

Claim (SGD Convergence) After $T = \frac{R^2 G'^2}{\epsilon^2}$ iterations: $\mathbb{E} [f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \le \epsilon.$

Claim 1:

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

Prove using Jensen's Inequality:

Claim (SGD Convergence) After $T = \frac{R^2 G'^2}{\epsilon^2}$ iterations: $\mathbb{E} [f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \le \epsilon.$

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

Claim (SGD Convergence) After $T = \frac{R^2 G'^2}{\epsilon^2}$ iterations: $\mathbb{E} [f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \le \epsilon.$

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

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

Claim (SGD Convergence)

After $T = \frac{R^2 G'^2}{\epsilon^2}$ iterations: $\mathbb{E} \left[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \right] \le \epsilon.$

$$\mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \leq \frac{1}{T} \sum_{i=1}^{T} \mathbb{E}\left[f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*)\right]$$
$$= \frac{1}{T} \sum_{i=1}^{T} n \mathbb{E}\left[f_{j_i}(\mathbf{x}^{(i)}) - f_{j_i}(\mathbf{x}^*)\right]$$
$$\leq \frac{n}{T} \cdot \mathbb{E}\left[\sum_{i=1}^{T} f_{j_i}(\mathbf{x}^{(i)}) - f_{j_i}(\mathbf{x}^{offline})\right]$$
$$\leq \frac{n}{T} \cdot \left(R \cdot \frac{G'}{n} \cdot \sqrt{T}\right) \qquad (by \ OGD \ guarantee.)$$

Number of iterations for error ϵ :

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

Always have $G \leq G'$:

$$\begin{split} \max_{\mathbf{x}} \|\nabla f(\mathbf{x})\|_{2} &\leq \max_{\mathbf{x}} \left(\|\nabla f_{1}(\mathbf{x})\|_{2} + \ldots + \|\nabla f_{n}(\mathbf{x})\|_{2} \right) \\ &\leq \max_{\mathbf{x}} \left(\|\nabla f_{1}(\mathbf{x})\|_{2} \right) + \ldots + \max_{\mathbf{x}} \left(\|\nabla f_{n}(\mathbf{x})\|_{2} \right) \\ &\leq n \cdot \frac{G'}{n} = G'. \end{split}$$

So GD converges strictly faster than SGD.

But for a fair comparison:

- SGD cost = (# of iterations) · O(1)
- GD cost = (# of iterations) · O(n)

We always have $G \le G'$. When it is <u>much smaller</u> 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'?

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?

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

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 $h(\mathbf{x}) : \mathbb{R}^d \to \mathbb{R}^d$ be an <u>invertible function</u>. 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 **P** be an invertible $d \times d$ matrix and let $g(\mathbf{x}) = f(\mathbf{Px})$.

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

Second Goal:

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

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

Third Goal: P should be easy to compute.

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

Example: Diagonal preconditioner.

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

P is often called a Jacobi preconditioner. Often works very well in practice!

DIAGONAL PRECONDITIONER

~	_
~	_

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

>> cond(A'*A)	<pre>>> P = sqrt(inv(diag(diag(A'*A)))); >> cond(P*A'*A*P)</pre>
ans =	ans =
8.4145e+07	10.3878

ADAPTIVE STEPSIZES

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 **P** is symmetric.

Gradient descent on g:

• For
$$t = 1, ..., T$$
,
• $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \mathbf{P} \left[\nabla f(\mathbf{P} \mathbf{x}^{(t)}) \right]$

• Return **Px**^(T)

Gradient descent on g:

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

When **P** is diagonal, this is just gradient descent with a <u>different step size for each parameter!</u>

ADAPTIVE STEPSIZES

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 <u>single random</u> <u>entry</u> of $\nabla f(\mathbf{x})$ on each iteration:

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

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