CS-GY 9223 D: Lecture 6
Online and Stochastic Gradient Decent

NYU Tandon School of Engineering, Prof. Christopher Musco

- If you don't have a project partner by the end of today, please email me.
- Take home midterm **week of October 26th**.
  - 2 hours, self-proctored. Design for 1.25 hours.
  - Can take anytime during that week.
  - Administered either via email or another option.
  - Solutions can be hand-written and scanned.
  - I will post some review questions.
- Need volunteers to present at **10/26 reading group** (in 2 weeks). Sign-up sheet on course webpage.

**First Order Optimization:** Given a function $f$ $:\mathbb{B}^d \to \mathbb{R}$ and a constraint set $\mathcal{S}$, assume we have:

- **Function oracle**: Evaluate $f(\mathbf{x})$ for any $\mathbf{x}$.
- **Gradient oracle**: Evaluate $\nabla f(\mathbf{x})$ for any $\mathbf{x}$.
- **Projection oracle**: Evaluate $P_{\mathcal{S}}(\mathbf{x})$ for any $\mathbf{x}$.

   **Goal:** Find $\hat{\mathbf{x}} \in \mathcal{S}$ such that $f(\hat{\mathbf{x}}) \leq \min_{\mathbf{x} \in \mathcal{S}} f(\mathbf{x}) + \epsilon.$

Projected gradient descent:

- Select starting point $x^{(0)}$, learning rate $\eta$.
- For $i = 0, \ldots, T$:
    - $z = x^{(i)} - \eta \nabla f(x^{(i)})$
    - $x^{(i+1)} = P_{\mathcal{S}}(z)$
- Return $\hat{x} = \arg \min_i f(x^{(i)})$.

Conditions for convergence:

- **Convexity:** $f$ is a convex function, $\mathcal{S}$ is a convex set.
- **Bounded initial distant:**

$$\|\underline{x^{(0)}} - x^*\|_2 \leq \boxed{R}$$

- **Bounded gradients (Lipschitz function):**

$$\|\nabla f(x)\|_2 \leq G \text{ for all } x \in \mathcal{S}.$$

**Theorem:** Projected Gradient Descent returns $\hat{x}$ with $f(\hat{x}) \leq \min_{x \in \mathcal{S}} f(x) + \epsilon$ after

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

iterations.
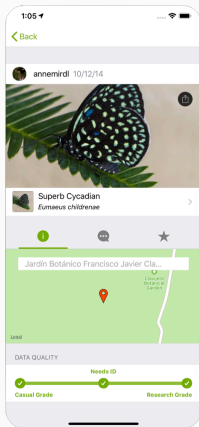
## ONLINE AND STOCHASTIC GRADIENT DESCENT

Today:

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

Many machine learning problems are solved in an online setting with constantly changing data.

- Spam filters are incrementally updated and adapt as they see more examples of spam over time.
- Image classification systems learn from mistakes over time (often based on user feedback).
- Content recommendation systems adapt to user behavior and clicks (which may not be a good thing...)
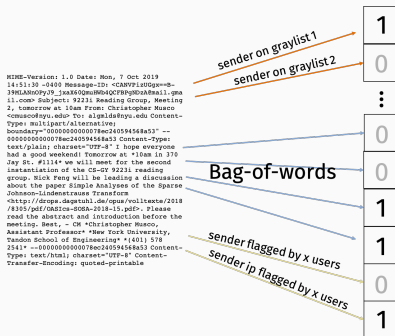
Plant identification via iNaturalist app.

(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.

ML based email spam/scam filtering.



Markers for spam change overtime, so model might change.

## ML based email spam/scam filtering.



Markers for spam change overtime, so model might change.

Choose some model $M_{\mathbf{x}}$ parameterized by parameters $\mathbf{x}$ and some loss function $\ell$. At time steps $1, \ldots, T$, receive data vectors $\underline{\mathbf{a}^{(1)}, \ldots, \mathbf{a}^{(T)}}$. $\in \mathbb{B}^d$

$i \in 1, \ldots T$
$\rightarrow$ not randomly

- At each time step, we pick ("play") a parameter vector $\mathbf{x}^{(i)}$.
- Make prediction $\tilde{y}^{(i)} = M_{\mathbf{x}^{(i)}}(\mathbf{a}_i)$. $\in \mathbb{R}$
- Then told true value or label $\underline{y^{(i)}}$. $\underline{\ell\left(x^{(i)}, a^{(i)}, y^{(i)}\right)}$
- Goal is to minimize cumulative loss:

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

For example, for a regression problem we might use the $\ell_2$ loss:

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

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

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

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

$$f_1(x) = \left| q^{(1)^T} \underline{x} - \underline{y^{(1)}} \right|^2$$

We make no assumptions that $f_1, \ldots, f_T$ are related to each other at all!

$$f_2(x) = \left| q^{(2)^T} x - y^{(2)} \right|^2$$

Online Gradient descent: *learning rate*

- Choose $\underline{\mathbf{x}^{(1)}}$ and $\underline{\eta = \text{\it lll}}$
- For $i = 1, \ldots, T$:
  - Play $\mathbf{x}^{(i)}$.
  - Observe $f_i$ and incur cost $f_i(\underline{\mathbf{x}^{(i)}})$.
  $\underline{\mathbf{x}^{(i+1)}} = \underline{\mathbf{x}^{(i)}} - \boxed{\eta \nabla f_i(\mathbf{x}^{(i)})}$

$$f_{i+1}\left(x^{(i+1)}\right)$$

If $f_1, \ldots, f_T = f$ are all the same, this looks a lot like regular gradient descent. We update parameters using the gradient $\nabla f$ at each step.

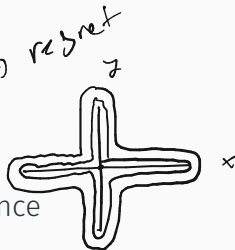If $f_1, \ldots, f_T$ are very different it might seem like nonsense right now...

$$f(x) - f(y) \leq \nabla f(x)^T (x - \partial)$$

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

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

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

regret

Here $\epsilon$ is called the **regret** of our solution sequence $x^{(1)}, \ldots, x^{(T)}$.

$\min (\partial - x) = 1$

This guarantee might seem a bit unfair. Why?

14

Regret compares to the best <u>fixed</u> solution in hindsight.

Possible
but weaker

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

It's very possible that $\sum_{i=1}^{T} f_i(\mathbf{x}^{(i)}) < \left[ \min_{\mathbf{x}} \sum_{i=1}^{T} f_i(\mathbf{x}) \right]$. Could we hope for something strong?

Exercise: Argue that the following is impossible to achieve:

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

$f_1, \ldots, f_T = f$

$$f_1(x) \quad \ldots \quad f_T(x)$$

$$f_i(x) = |x - h| \quad \text{where} \quad h = \begin{cases} 0 & \text{w/ prob } \frac{1}{2} \\ 1 & \text{w/ prob } \frac{1}{2} \end{cases}$$

$$\sum_{i=1}^{T} f_i(x^{(i)}) \leq \sum_{i=1}^{T} \min_x f_i(x^{(i)}) + \varepsilon$$

$$\underbrace{\phantom{\sum_{i=1}^{T} f_i(x^{(i)})}} \qquad \underbrace{\phantom{\sum_{i=1}^{T} \min_x f_i(x^{(i)})}}$$

$$\approx T/2 \qquad\qquad 0$$

$$h_1, \ldots h_T = 0, 0, 1, 0, 1, 1, 0$$

$$h_1, \ldots, h_T = 0, 1, 0, 1, 0, 1 \ldots -$$

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

Beautiful balance:

- Either $f_1, \ldots, f_T$ are similar, so an method like Online Gradient Descent will effectively minimize $\sum_{i=1}^{T} f_i(\mathbf{x}^{(i)})$.

- 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.

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

**Assume:**

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

**Online Gradient descent:**

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

$$x^{(1)} \ldots x^{(T)} \qquad x^*$$

$$\sum_{i=1}^{T} f_i(x^{(i)}) = \text{objective for online optimization}$$

$$= \arg\min_x f_i(x)$$

18

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

**Theorem (OGD Regret Bound)**

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

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

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

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.

19

Theorem (OGD Regret Bound)

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

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

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

(Same proof as last class. Only uses convexity of $f_i$.)

**Theorem (OGD Regret Bound)**

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

*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 $i = 1, \ldots, T$,

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

**Telescoping Sum:**

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

$\geq 0$

$g(\eta) = \frac{R^2}{2\eta} + \frac{T\eta G^2}{2} \leq \frac{R^2}{2\eta} + \frac{T\eta G^2}{2} = RG\sqrt{T}$

$\frac{R^2}{2\eta} = \frac{T\eta G^2}{2}$

$\eta = \sqrt{\frac{R^2}{G^2 T}} = \frac{R}{G\sqrt{T}}$

21

## STOCHASTIC GRADIENT DESCENT (SGD)

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

$$f(x) = \|Ax - b\|_2^2$$
$$= \sum_{i=1}^{n} (\langle a_i, x \rangle - b_i)^2$$

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

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!

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

$$\nabla \left( \sum_{i=1}^{n} f_i(x) \right) = \left( \sum_{i=1}^{n} \nabla f_i(x) \right)$$

23

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

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

*uniformly random*

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

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

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

Trade slower convergence for cheaper iterations.

$$f(x) = \sum_{j=1}^{n} f_j(x) \qquad \mathbb{E}_{i=1,\ldots n}\left[f_i(x)\right] = \sum_{j=1}^{n} \frac{1}{n} f_j(x)$$
$$= \frac{1}{n} f(x)$$

24

## STOCHASTIC GRADIENT DESCENT

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

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

Computing $f(\mathbf{x})$ would take $n$ separate function queries.

**Stochastic Gradient descent:**    $f(x^{(i)})$

- Choose starting vector $\mathbf{x}^{(1)}$, learning rate $\eta$
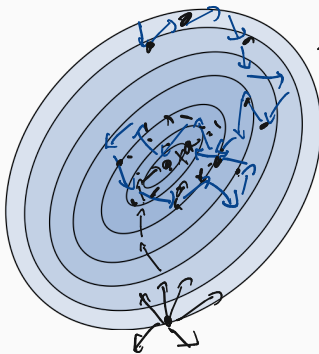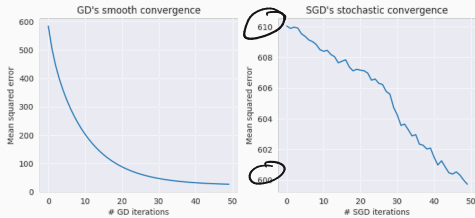- For $i = 1, \ldots, T$:
  - Pick random $j_i \in 1, \ldots, n$      $j_1, \ldots j_T \in 1, \ldots n$
  - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f_{j_i}(\mathbf{x}^{(i)})$
- Return $\hat{\mathbf{x}} = \frac{1}{T} \sum_{i=1}^{T} \mathbf{x}^{(i)}$

GD's smooth convergence

SGD's stochastic convergence

$\rightarrow$ level sets of $f(x)$

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

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

**Stochastic Gradient descent:**

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

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

$\min f(\hat{x})$

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

27

Via Markovs:

with prob 9/10

$$f(\hat{x}) - f(x^*) \leq 10\varepsilon$$

**Claim (SGD Convergence)**

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

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

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

random variable

**Claim 1:**

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

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

28

### Claim (SGD Convergence)

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

*where $\hat{\mathbf{x}} = \frac{1}{T}\sum_{i=1}^{T}\mathbf{x}^{(i)}$.*

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

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

$$f_{j_1}, \ldots, f_{j_T}$$

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

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

$$= \frac{RG'}{\sqrt{T}} < \epsilon$$

29

Number of iterations for error $\epsilon$:

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

**Always have** $G \leq G'$:

*triangle inequality*

$$\max_x \quad \|\nabla f(x)\|_2 \leq \|\nabla f_1(x)\|_2 + \ldots + \|\nabla f_n(x)\|_2 \leq n \cdot \frac{G'}{n} = G'.$$

$$\nabla f(x) = \nabla f_1(x) + \ldots + \nabla f_n(x)$$

So GD converges strictly faster than *SGD*.

$$G \leq G'$$

But for a fair comparison:

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

$$\|x + y\|_2 \leq \|x\|_2 + \|y\|_2$$

30

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

$$G << G'$$

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

$$\nabla f_1(x) = \nabla f_2(x) = \ldots = \nabla f_n(x)$$

$$\|\nabla f(x)\|_2 = n \cdot \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(x) = G'$$

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?

$$\mathbb{E}\left[\|\nabla f_i(\mathbf{x})\|_2^2\right] = \sum_{j=1}^{d} z_j^2 \quad \text{where} \quad z \sim \mathcal{N}(0,1)$$
$$= \boxed{d}$$

$$\mathbb{E}\left[\|\nabla f(\mathbf{x})\|_2^2\right] = \mathbb{E}\left[\|\sum_{i=1}^{n} \nabla f_i(\mathbf{x})\|_2^2\right] = \sum_{j=1}^{d} s_j^2 \quad \text{where} \quad s_j \sim \mathcal{N}(0,n)$$
$$= \boxed{dn}$$

Each entry of $\sum_{i=1}^{n} \nabla f_i(x)$ is the sum of $n$ guassians, so distributed as $\mathcal{N}(0,n)$.

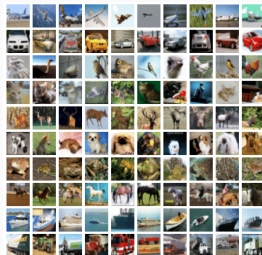So, $G' \approx n \cdot \sqrt{d}$ and $G \propto \sqrt{dn}$ . $\longrightarrow$ roughly same complexity.

SGD. takes $O\left(\frac{R^2 n^2 d}{\epsilon^2}\right)$ iterations, GD takes $O\left(\frac{R^2 nd}{\epsilon^2}\right)$.

32

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

Can our convergence bounds be tightened for certain functions? Can they guide us towards faster algorithms?

### Goals:

- Improve $\epsilon$ dependence below $1/\epsilon^2$.
    - Ideally $1/\epsilon$ or $\log(1/\epsilon)$.
- Reduce or eliminate dependence on *G* and *R*.
- Further take advantage of structure in the data (e.g. repetition in features in addition to data points).

### Definition ($\beta$-smoothness)

A function $f$ is $\beta$ smooth if, for all x, y

$$\|\nabla f(\mathsf{x}) - \nabla f(\mathsf{y})\|_2 \leq \beta \|\mathsf{x} - \mathsf{y}\|_2$$

After some calculus (see Lemma 3.4 in **Bubeck's book**), this implies:

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

Recall from definition of convexity that:

$$f(\mathbf{x}) - f(\mathbf{y}) \leq \nabla f(\mathbf{x})^T(\mathbf{x} - \mathbf{y})$$

So now we have an upper and lower bound.

$$0 \leq \nabla f(\mathbf{x})^T(\mathbf{x} - \mathbf{y}) - [f(\mathbf{x}) - f(\mathbf{y})] \leq \frac{\beta}{2}\|\mathbf{x} - \mathbf{y}\|_2^2$$

Previously learning rate/step size $\eta$ depended on $G$. Now choose it based on $\beta$:

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

Progress per step of gradient descent:

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

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

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

**Theorem (GD convergence for $\beta$-smooth functions.)**

*Let $f$ be a $\beta$ 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:*

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \frac{2\beta R^2}{T - 1}$$

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

Definition ($\alpha$-strongly convex)

A convex function $f$ is $\alpha$-strongly convex if, for all x, y

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

$\alpha$ is a parameter that will depend on our function.

**Completing the picture:** If $f$ is $\alpha$ strongly convex and $\beta$ smooth,

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

Gradient descent for strongly convex functions:

- Choose number of steps $T$.
- For $i = 1, \ldots, T$:
    - $\eta = \frac{2}{\alpha \cdot (i+1)}$
    - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$
- Return $\hat{\mathbf{x}} = \arg\min_{\mathbf{x}^{(i)}} f(\mathbf{x}^{(i)})$.

**Theorem (GD convergence for $\alpha$-strongly convex functions.)**

*Let f be an $\alpha$-strongly convex function and assume we have that, for all x, $\|\nabla f(\mathbf{x})\|_2 \leq G$. If we run GD for T steps (with adaptive step sizes) we have:*

$$f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \leq \frac{2G^2}{\alpha(T-1)}$$

**Corollary**: If $T = O\left(\frac{G^2}{\alpha\epsilon}\right)$ we have $f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \leq \epsilon$

What if $f$ is both $\beta$-smooth and $\alpha$-strongly convex?

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

### Theorem (GD for $\beta$-smooth, $\alpha$-strongly convex.)

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

$$\|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2 \leq e^{-(T-1)\frac{\alpha}{\beta}} \|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2^2$$

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

Is it better if $\kappa$ is large or small?

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

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

we have:

$$\|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2^2 \leq \frac{2}{\alpha}\left[f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*)\right]$$

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

### Corollary (GD for $\beta$-smooth, $\alpha$-strongly convex.)

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

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \frac{\beta}{\alpha} e^{-(T-1)\frac{\alpha}{\beta}} \cdot \left[ f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*) \right]$$

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

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \epsilon \left[ f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*) \right]$$
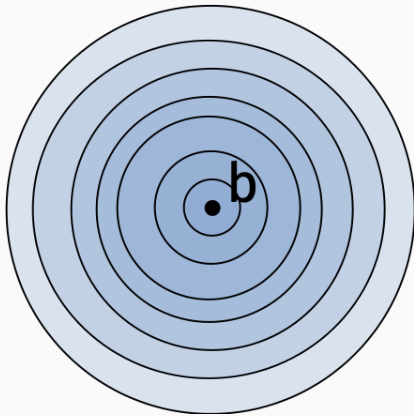
## UNDERSTANDING CONDITIONING

Let $f(\mathbf{x}) = \|\mathbf{D}\mathbf{x} - \mathbf{b}\|_2^2$ where $\mathbf{D}$ is a diagaonl matrix. For now imagine we're in two dimensions: $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$, $\mathbf{D} = \begin{bmatrix} d_1 & 0 \\ 0 & d_2 \end{bmatrix}$.
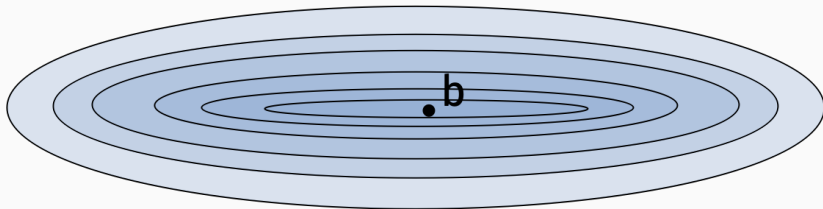
### What are $\alpha, \beta$ for this problem?

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

**Exercise:** Show that:

$$\nabla f(\mathbf{x})^T(\mathbf{x} - \mathbf{y}) - [f(\mathbf{x}) - f(\mathbf{y})] = (\mathbf{x} - \mathbf{y})^T \mathbf{D}^2(\mathbf{x} - \mathbf{y})$$
$$= \|\mathbf{D}(\mathbf{x} - \mathbf{y})\|_2^2$$

Level sets of $\|\mathbf{D}\mathbf{x} - \mathbf{b}\|_2^2$ when $d_1 = 1, d_2 = 1$.

Level sets of $\|\mathbf{D}\mathbf{x} - \mathbf{b}\|_2^2$ when $d_1 = \frac{1}{3}, d_2 = 2$.

## UNDERSTANDING CONDITIONING

Steps to convergence $\approx O\left(\kappa \log(1/\epsilon)\right) = O\left(\frac{\max(\mathbf{D}^2)}{\min(\mathbf{D}^2)} \log(1/\epsilon)\right)$.

For general regression problems $\|\mathbf{Ax} - \mathbf{b}\|_2^2$,

$$\beta = \lambda_{max}(\mathbf{A}^T\mathbf{A})$$
$$\alpha = \lambda_{min}(\mathbf{A}^T\mathbf{A})$$