CS-GY 6763: Lecture 6 Gradient Descent and Projected Gradient Descent

NYU Tandon School of Engineering, Prof. Christopher Musco

CONTINUOUS OPTIMIZATION

Given function
$$f : \mathbb{R}^d \to \mathbb{R}^{\mathbf{0}}$$
. Find $\hat{\mathbf{x}}$ such that:
 $f(\hat{\mathbf{x}}) \leq \min_{\mathbf{x}} f(\mathbf{x}) + \epsilon$.

Let $\mathbf{a}_{-}^{(1)}, \ldots, \mathbf{a}_{-}^{(n)} \in \mathbb{R}^{d}$ be a collection of data points and $y^{(1)}, \ldots, y^{(n)}$ be a collection of target values.

- Model: $M_{\mathbf{X}}(\mathbf{a}) = \mathbf{x}^T \mathbf{a}$. **x** contains the regression coefficients.

• Function to minimize: $f(\mathbf{x}) = |z - y|^2$. • Function to minimize: $f(\mathbf{x}) = \sum_{i=1}^{n} |\mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)}|^2$

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

where **A** is a matrix with $\mathbf{a}^{(i)}$ as its *i*th row and **y** is a vector with $y^{(i)}$ as its i^{th} entry.

The choice of algorithm to minimize $f(\mathbf{x})$ will depend on:

- The form of $f(\mathbf{x})$ (is it linear, is it quadratic, does it have finite sum structure, etc.)
- If there are any additional constraints imposed on **x**. E.g. $\|\mathbf{x}\|_2 \leq c$.

Gradient descent: A greedy algorithm for minimizing functions of multiple variables that often works amazingly well.



For i = 1, ..., d, let x_i be the i^{th} entry of **x**. Let $e^{(i)}$ be the i^{th} standard basis vector.

Partial derivative: $\int \frac{\partial f}{\partial x_i}(\mathbf{x}) = \lim_{t \to 0} \frac{f(\mathbf{x} + t\mathbf{e}^{(i)}) - f(\mathbf{x})}{t}$ Directional derivative: Χ, $\underline{Dof(\mathbf{x})} = \lim_{t \to 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t}$ X ~ $\bigcup_{p(i)} f(\kappa)$

CALCULUS REVIEW



Directional derivative:



FIRST ORDER OPTIMIZATION

Given a function *f* to minimize, assume we have:

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

We view the implementation of these oracles as black-boxes, but they can often require a fair bit of computation.

7 unit

Linear least-squares regression:

 $\mathcal{O}(nd)$

- Given $\mathbf{a}^{(1)}, \dots \mathbf{a}^{(n)} \in \mathbb{R}^d$, $y^{(1)}, \dots y^{(n)} \in \mathbb{R}$.
- Want to minimize:

$$f(\mathbf{x}) = \sum_{i=1}^{n} \left(\mathbf{x}^{\mathsf{T}} \mathbf{a}^{(i)} - \mathbf{y}^{(i)} \right)^2 = \| \mathbf{A} \mathbf{x} - \mathbf{y} \|_2^2.$$

EXAMPLE GRADIENT EVALUATION

ar Linear least-squares regression • Given $\mathbf{a}^{(1)}, \ldots \mathbf{a}^{(n)} \in \mathbb{R}^d$, $y^{(1)}, \ldots y^{(n)} \in \mathbb{R}^d$. Want to minimize: $f(\mathbf{x}) = \sum_{i=1}^{n} \left(\mathbf{x}^{T} \mathbf{a}^{(i)} - \mathbf{y}^{(i)} \right)^{2} = \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_{2}^{2}.$ NXJ $\frac{\partial f}{\partial x_i} = \sum_{i=1}^{n} 2\left(\mathbf{x}^T \mathbf{a}^{(i)} - \mathbf{y}^{(i)}\right) \left(a_j^{(i)}\right) = 2\alpha^{(i)T} (\mathbf{A}\mathbf{x} - \mathbf{y})$ where $\alpha^{(j)}$ is the *j*th column of **A**. (dxy)(yx1)=dx1 $\nabla f(\mathbf{x}) = 2\mathbf{A}^{\mathsf{T}} (\mathbf{A}\mathbf{x} - \mathbf{y})^{\mathsf{T}}$ Ohzd) What is the time complexity of a gradient oracle for $\nabla f(\mathbf{x})$? Mud) O(ud) 2.0 (nd) 10

Greedy approach: Given a starting point \mathbf{x} , make a small adjustment that decreases $f(\mathbf{x})$. In particular, $\mathbf{x} \leftarrow \mathbf{x} + \eta \mathbf{v}$.

What property do I want in **v**?

P Step 5.20

Leading question: When η is small, what's an approximation for $f(\mathbf{x} + \eta \mathbf{v}) - f(\mathbf{x})$? $V = - \nabla f(\mathbf{x})$ $f(\mathbf{x} + \eta \mathbf{v}) - f(\mathbf{x}) \approx . D_{\mathbf{v}} f(\mathbf{x}) \cdot \mathbf{M}$ $= \nabla f(\mathbf{x})^T \mathbf{v} \cdot \mathbf{M}$ $= - \nabla f(\mathbf{x})^T \nabla f(\mathbf{x}) \cdot \mathbf{M} = - |\nabla f(\mathbf{x})||_{\eta_1}^{\eta_1}$

DIRECTIONAL DERIVATIVES

$$D_{\mathbf{v}}f(\mathbf{x}) = \lim_{t \to 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} = \nabla f(\mathbf{x})^T \mathbf{v}.$$

So:

$$f(\mathbf{x} + \eta \mathbf{v}) - f(\mathbf{x}) \approx$$

How should we choose v so that $f(x + \eta v) < f(x)$?

Prototype algorithm:

• Choose starting point
$$\mathbf{x}^{(0)}$$
. $\Rightarrow \mathbf{0}$

Μ

• For
$$i = 0, ..., T$$
:

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

• Return $\mathbf{x}^{(T)}$.

 η is a step-size parameter, which is often adapted on the go. For now, assume it is fixed ahead of time.

GRADIENT DESCENT INTUITION



2 dimensional example:



For a convex function $f(\mathbf{x})$: For sufficiently small η and a sufficiently large number of iterations T, gradient descent will converge to a near global minimum:

$$f(\mathbf{X}^{(T)}) \leq f(\mathbf{X}^*) + \epsilon.$$

Examples: least squares regression, logistic regression, kernel regression, SVMs. $\nabla f(x) = \vec{o}$

For a non-convex function $f(\mathbf{x})$: For sufficiently small η and a sufficiently large number of iterations *T*, gradient descent will converge to a near stationary point:

$$\|\nabla f(\mathbf{x}^{(T)})\|_2 \leq \epsilon.$$

Examples: neural networks, matrix completion problems (

CONVEX VS. NON-CONVEX



One issue with non-convex functions is that they can have **local minima**. Even when they don't, convergence analysis requires different assumptions than convex functions.

We care about <u>how fast</u> gradient descent and related methods converge, not just that they do converge.

- Bounding iteration complexity requires placing some assumptions on *f*(**x**).
- Stronger assumptions lead to better bounds on the convergence.

Understanding these assumptions can help us design faster variants of gradient descent (there are many!).

Today, we will start with **convex** functions.

CONVEXITY



GRADIENT DESCENT



It is easy but not obvious how to prove the equivalence between these definitions. A short proof can be found in Karthik Sridharan's lecture notes here:

http://www.cs.cornell.edu/courses/cs6783/2018fa/lec16supplement.pdf

Assume:

- f is convex.
- Lipschitz function: for all \mathbf{x} , $\|\nabla f(\mathbf{x})\|_2$
- Starting radius: $\|\mathbf{x}^* \mathbf{x}^{(0)}\|_2 \leq \mathbf{x}^*$

Gradient descent:

- Choose number of steps T.
- Starting point $\mathbf{x}^{(0)}$. E.g. $\mathbf{x}^{(0)} = \vec{0}$.
- $\begin{array}{c} \underbrace{\eta}_{i} = \underbrace{\binom{R}{G\sqrt{T}}}_{For \ i = 0, \dots, T:} \\ \cdot \mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} \eta \nabla f(\mathbf{x}^{(i)}) \end{array}$
- Return $\hat{\mathbf{x}} = \arg \min_{\mathbf{x}^{(i)}} f(\mathbf{x}^{(i)}).$



Claim (GD Convergence Bound)

If we run GD for $T \ge \frac{R^2 G^2}{\epsilon^2}$ iterations with step-size $\eta = \frac{R}{G\sqrt{T}}$, then $f(\hat{\mathbf{x}}) \le f(\mathbf{x}^*) + \epsilon$.

Proof is made tricky by the fact that $f(\mathbf{x}^{(i)})$ does not improve monotonically. We can "overshoot" the minimum.

We will prove that the <u>average</u> solution value is low after $T = \frac{R^2G^2}{\epsilon^2}$ iterations. I.e. that:

$$\frac{1}{T}\sum_{i=0}^{T-1}\left[\underline{f(\mathbf{x}^{(i)})}-\underline{f(\mathbf{x}^*)}\right] \leq \epsilon.$$

Of course the best solution found, \hat{x} is only better than the average. $f(\mathcal{R}) - f(x^*) \leq \varepsilon$

Claim (GD Convergence Bound)

If we run GD for $T \ge \frac{R^2 G^2}{\epsilon^2}$ iterations with step-size $\eta = \frac{R}{G\sqrt{T}}$, then $f(\hat{\mathbf{x}}) \le f(\mathbf{x}^*) + \epsilon$.

Claim 1: For all i = 0, ..., T,

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

Claim 1(a): For all *i* = 0, ..., *T*,

$$\nabla f(\mathbf{x}^{(i)})^{\mathsf{T}}(\underline{\mathbf{x}}^{(i)} - \underline{\mathbf{x}}^{*}) \leq \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{*}\|_{2}^{2} - \|\mathbf{x}^{(i+1)} - \mathbf{x}^{*}\|_{2}^{2}}{2\eta} + \frac{\eta G^{2}}{2}$$

Claim 1 follows from Claim 1(a) by definition of convexity.

$$f(\mathbf{x}^{(i)}) - f(\mathbf{x}^{*}) \leq \nabla f(\mathbf{x}^{(i)})^{\mathsf{T}} \left(\mathbf{x}^{(i)} - \mathbf{x}^{*} \right)$$

$\|Q - b\|_{2}^{2} = \|o\|_{1}^{2} + \|b\|_{2}^{2} - 2QTb$

Claim (GD Convergence Bound)

If we run GD for $T \ge \frac{R^2G^2}{\epsilon^2}$ iterations with step size $\eta = \frac{R}{G\sqrt{T}}$, then $f(\hat{\mathbf{x}}) \le f(\mathbf{x}^*) + \epsilon$.

Claim 1(a): For all
$$i = 0, ..., T$$
,

$$\nabla f(\mathbf{x}^{(i)})^{T}(\mathbf{x}^{(i)} - \mathbf{x}^{*}) \leq \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{*}\|_{2}^{2} - (\|\mathbf{x}^{(i+1)} - \mathbf{x}^{*}\|_{2}^{2})}{2\eta} + \frac{\eta G^{2}}{2}$$

$$\| \underbrace{\mathbf{x}^{(i)}}_{2m} - \underbrace{\mathbf{x}^{*}}_{2m} \|_{2}^{2} = \| \mathbf{x}^{(i)} - \mathbf{x}^{*}\|_{2}^{2} + m^{2} \| \mathbf{x}^{f}(\mathbf{x}^{(i)}) \|_{2}^{2} - 2m \nabla f(\mathbf{x}^{(i)})^{T}(\mathbf{x}^{(i)} - \mathbf{x}^{*})$$

$$\leq \| \mathbf{x}^{(i)} - \mathbf{x}^{*}\|_{2}^{2} + \frac{m^{2} G^{2}}{2m} - 2m \nabla f(\mathbf{x}^{(i)})^{T}(\mathbf{x}^{(i)} - \mathbf{x}^{*})$$



27



Claim (GD Convergence Bound) If $T \ge \frac{R^2G^2}{\epsilon^2}$ and $\eta = \frac{R}{G\sqrt{T}}$, then $f(\hat{\mathbf{x}}) \le f(\mathbf{x}^*) + \epsilon$.

Final step:

$$\frac{1}{T}\sum_{i=0}^{T-1} \left[f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \right] \le \epsilon$$
$$\left[\frac{1}{T}\sum_{i=0}^{T-1} f(\mathbf{x}^{(i)}) \right] - f(\mathbf{x}^*) \le \epsilon$$

We always have that $f(\hat{\mathbf{x}}) = \min_i f(\mathbf{x}^{(i)}) \le \frac{1}{T} \sum_{i=0}^{T-1} f(\mathbf{x}^{(i)})$, which gives the final bound:

$$f(\hat{\mathbf{X}}) \leq f(\mathbf{X}^*) + \epsilon.$$

Typical goal: Solve a <u>convex minimization problem</u> with additional <u>convex constraints</u>.



Which of these is convex?

CONSTRAINED CONVEX OPTIMIZATION



Definition (Convex set) A set S is convex if for any $\underline{x}, \underline{y} \in S$, $\underline{\lambda} \in [0, 1]$: $(1 - \lambda)\mathbf{x} + \underline{\lambda}\mathbf{y} \in S$.

CONSTRAINED CONVEX OPTIMIZATION

Examples:

- Norm constraint: minimize $||Ax b||_2$ subject to $||x||_2 \le \lambda$. Used e.g. for regularization, finding a sparse solution, etc.
- Positivity constraint: minimize $f(\mathbf{x})$ subject to $\mathbf{x} \ge 0$. Used e.g. in finding an optimal allocation for a portfolio into different assets.
- Linear constraint: minimize Subject to <u>Ax ≤ b</u>. Linear program used in training support vector machines, industrial optimization, subroutine in integer programming, etc.

PROBLEM WITH GRADIENT DESCENT

Gradient descent:

• For i = 0, ..., T:

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

• Return $\hat{\mathbf{x}} = \arg\min_i f(\mathbf{x}^{(i)})$.



Even if we start with $\underline{x}^{(0)} \in S$, there is no guarantee that $\underline{x}^{(0)} - \eta \nabla f(\underline{x}^{(0)})$ will remain in our set.

Extremely simple modification: Force $\mathbf{x}^{(i)}$ to be in S by **projecting** onto the set.

Given a function f to minimize and a convex constraint set 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 $\mathcal{P}_{\mathcal{S}}(\mathbf{x})$ for any \mathbf{x} . $\mathcal{P}_{\mathcal{S}}(\mathbf{x}) = \arg\min_{\mathbf{y} \in \mathcal{S}} \frac{\|\mathbf{x} - \mathbf{y}\|_2}{\|\mathbf{y} \in \mathcal{S}}$

PROJECTION ORACLES



• How would you implement P_S for $S = \{y : y = Qz\}$.

3



Given function $f(\mathbf{x})$ and set S, such that $\|\nabla f(\mathbf{x})\|_2 \leq G$ for all $\mathbf{x} \in S$ and starting point $\mathbf{x}^{(0)}$ with $\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2 \leq R$.

Projected gradient descent:

- Select starting point $\mathbf{x}^{(0)}$, $\eta = \frac{R}{G\sqrt{T}}$.
- For i = 0, ..., T:

$$\cdot \underbrace{\mathbf{z} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})}_{\mathbf{x}^{(i+1)} = P_{\mathcal{S}}(\mathbf{z}) }$$

• Return $\hat{\mathbf{x}} = \operatorname{arg\,min}_i f(\mathbf{x}^{(i)}).$

Claim (PGD Convergence Bound)

If f, S are convex and $T \ge \frac{R^2 G^2}{\epsilon^2}$, then $f(\hat{\mathbf{x}}) \le f(\mathbf{x}^*) \underbrace{f(\mathbf{x}^*)}_{\epsilon > \epsilon}$

Analysis is almost identical to standard gradient descent! We just need one additional claim:

Claim (Contraction Property of Convex Projection) If S is convex, then for any $y \in S$,

 $\|\mathbf{y} - P_{\mathcal{S}}(\mathbf{x})\|_{2} \le \|\mathbf{y} - \mathbf{x}\|_{2}.$





Claim (PGD Convergence Bound)

If f, S are convex and $T \ge \frac{R^2 G^2}{\epsilon^2}$, then $f(\hat{\mathbf{x}}) \le f(\mathbf{x}^*) + \epsilon$.

Claim 1: For all $i = \underline{0}, ..., T$, $f(\mathbf{x}^{(i)}) - f(\mathbf{x}^{*}) \leq \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{*}\|_{2}^{2} - \|\mathbf{z} - \mathbf{x}^{*}\|_{2}^{2}}{2\eta} + \frac{\eta G^{2}}{2}$ $\leq \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{*}\|_{2}^{2} - \|\mathbf{x}^{(i+1)} - \mathbf{x}^{*}\|_{2}^{2}}{2\eta} + \frac{\eta G^{2}}{2}$ $\chi^{\text{GHV}} = \mathcal{P}_{S}(\mathbf{z}) \qquad \chi^{\text{T}} \text{ is } \eta \leq S$

Same telescoping sum argument:

$$\left[\frac{1}{T}\sum_{i=0}^{T-1}f(\mathbf{x}^{(i)})\right]-f(\mathbf{x}^*)\leq \frac{R^2}{2T\eta}+\frac{\eta G^2}{2}.$$

Conditions:

- **Convexity:** f is a convex function, S is a convex set.
- · Bounded initial distant:

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

• Bounded gradients (Lipschitz function):

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

Theorem (GD Convergence Bound)

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



Can our convergence bound be tightened for certain functions? Can it guide us towards faster algorithms?

Goals:

- Improve ϵ dependence below $1/\epsilon^2$. • Ideally $1/\epsilon$ or $\log(1/\epsilon)$.
- Reduce or eliminate dependence on *G* and *R*.
- **Next class:** Take advantage of additional problem structure (e.g. repetition in features and data points in ML problems).

SMOOTHNESS





A function f is β smooth if, for all \mathbf{x}, \mathbf{y}

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

For a scalar valued function f, equivalent to $f''(x) \leq \beta$. After

some calculus (see Lem. 3.4 in Bubeck's book), this implies:

$$\left(\underbrace{[f(\mathbf{y}) - f(\mathbf{x})] - \nabla f(\mathbf{x})^{\mathsf{T}}(\mathbf{y} - \mathbf{x})}_{\mathbb{Z}} \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|_{2}^{2} \right)$$

Recall from convexity that $f(\mathbf{y}) - f(\mathbf{x}) \ge \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x})$.



CONVERGENCE GUARANTEE

Theorem (GD convergence for β -smooth functions.) Let f be a β smooth convex function and assume we have $\|\mathbf{x}^* - \mathbf{x}^{(\mathbf{0})}\|_2 \leq R$. If we run GD for T steps, we have:

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

Corollary: If
$$\underline{T} = O\left(\frac{\beta R^2}{\epsilon}\right)$$
 we have $f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \in \epsilon$.
Compare this to $T = O\left(\frac{G^2 R^2}{\epsilon^2}\right)$ without a smoothness assumption.

GUARANTEED PROGRESS

M~-

Why do you think gradient descent might be faster when a function is β -smooth? Think about scalar case, in which case smoothness means $f''(x) \leq \beta$.



Previously learning rate/step size η depended on G. Now choose it based on β : M= $f(x^{(+*)}) \leq f(x^{(+)})$ $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - \int \frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})$ 1. $[f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)})] - \nabla f(\mathbf{x}^{(t)})^{\mathsf{T}} (\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}) \leq \frac{\beta}{2} || (\mathbf{x}^{(t)} - \mathbf{x}^{(t+1)}) ||_{2}^{2}.$ 2. $[f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)})] + \frac{1}{\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \frac{1}{62} \|\mathbf{\xi} \nabla f(\mathbf{x}^{(t)})\|_2^2.$ 3. $f(\mathbf{x}^{(t)}) - f(\mathbf{x}^{(t+1)}) \neq$ $\nabla f(x^{(+1)})^{1}(-\frac{1}{6})\nabla f(x^{(+)})$ 45

CONVERGENCE GUARANTEE

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

Again getting this result from the previous page is not hard, but also not obvious/direct. A concise proof can be found in Robert Gower's notes.



Where did we use convexity in this proof?

Progress per step of gradient descent:

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

STATIONARY POINTS

Definition (Stationary point)

For a differentiable function *f*, a <u>stationary point</u> is any **x** with:

$$\nabla f(\mathbf{x}) = \mathbf{0}$$

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

Theorem (Convergence to Stationary Point)

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

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

Corollary: If $T \ge \frac{2\beta}{\epsilon}$, then $\|\nabla f(\hat{\mathbf{x}})\|_2^2 \le \epsilon (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*))$. **f(x)**

Theorem (Convergence to 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:

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

We have that $f(\mathbf{x}^{(t)}) - f(\mathbf{x}^{(t+1)}) \ge \frac{1}{2\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2$ So: $\sum_{t=0}^{T-1} \frac{1}{2\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \le f(\mathbf{x}^{(0)}) - f(\mathbf{x}^{(t)})$ $(\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \le (\frac{2\beta}{T} (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*)))$ $\min_t \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \le \frac{2\beta}{T} (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*))$ If GD can find a stationary point, are there algorithms which find a stationary point faster using preconditioning, acceleration, stochastic methods, etc.? What if my function only has global minima and saddle points? Randomized methods (SGD, perturbed gradient methods, etc.) can provably "escape" saddle points.

Example: $\min_{x} \frac{-x^T A^T A x}{x^T x}$

- **Global minimum**: Top eigenvector of **A**^T**A** (i.e., top principal component of **A**).
- Saddle points: All other eigenvectors of A.

Useful for lots of other matrix factorization problems beyond vanilla PCA.

I said it was a bit tricky to prove that $f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \leq \frac{2\beta R^2}{T}$ for convex functions. But we just easily proved that $\|\nabla f(\hat{\mathbf{x}})\|_2^2$ is small. Why doesn't this show we are close to the minimum?

STRONG CONVEXITY

Definition (α -strongly convex)

A convex function f is α -strongly convex if, for all \mathbf{x}, \mathbf{y}

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

 α is a parameter that will depend on our function. For a twice-differentiable scalar function *f*, equivalent to $f''(x) \ge \alpha$.

When f is convex, we always have that $f''(x) \ge 0$, so larger values of α correspond to a "stronger" condition.

Gradient descent for strongly convex functions:

- Choose number of steps T.
- For i = 1, ..., 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 α -strongly convex functions.) Let f be an α -strongly convex function and assume we have that, for all \mathbf{x} , $\|\nabla f(\mathbf{x})\|_2 \leq \mathbf{G}$. If we run GD for T steps (with adaptive step sizes) we have:

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

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

We could also have that f is both β -smooth and α -strongly convex.

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

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

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:

$$\|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2^2 \le \frac{2}{\alpha} \left[f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*) \right]$$
$$\|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2 \ge \frac{2}{\beta} \left[f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \right]$$

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:

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \le \frac{\beta}{\alpha} e^{-T\frac{\alpha}{\beta}} \cdot \left[f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*) \right]$$

Corollary: If $T = O\left(\frac{\beta}{\alpha}\log(\beta/\alpha\epsilon)\right) = O(\kappa\log(\kappa/\epsilon))$ we have: $f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \le \epsilon \left[f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*)\right]$

Alternative 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 $1/\epsilon^2$!

Convexity:

$$0 \leq [f(\mathbf{y}) - f(\mathbf{x})] - \nabla f(\mathbf{x})^{\mathsf{T}}(\mathbf{y} - \mathbf{x})$$

 α -strong-convexity and β -smoothness:

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

Number of iterations for ϵ error:

	G-Lipschitz	eta-smooth
<i>R</i> bounded start	$O\left(\frac{G^2R^2}{\epsilon^2}\right)$	$O\left(\frac{\beta R^2}{\epsilon}\right)$
$\alpha\text{-}strong\ convex$	$O\left(\frac{G^2}{\alpha\epsilon}\right)$	$O\left(\frac{\beta}{\alpha}\log(1/\epsilon)\right)$

THE HESSIAN

Let *f* be a twice differentiable function from $\mathbb{R}^d \to \mathbb{R}$. Let the Hessian $H = \nabla^2 f(\mathbf{x})$ contain all of its second derivatives at a point \mathbf{x} . So $H \in \mathbb{R}^{d \times d}$. We have:

$$\mathbf{H}_{j,k} = \left[\nabla^2 f(\mathbf{x})\right]_{j,k} = \frac{\partial^2 f}{\partial x_j x_k}.$$

For vector **x**, **v**:

$$\nabla f(\mathbf{x} + t\mathbf{v}) \approx \nabla f(\mathbf{x}) + t \left[\nabla^2 f(\mathbf{x})\right] \mathbf{v}.$$

THE HESSIAN

Let *f* be a twice differentiable function from $\mathbb{R}^d \to \mathbb{R}$. Let the Hessian $H = \nabla^2 f(\mathbf{x})$ contain all of its second derivatives at a point \mathbf{x} . So $H \in \mathbb{R}^{d \times d}$. We have:

$$\mathsf{H}_{j,k} = \left[\nabla^2 f(\mathsf{x})\right]_{j,k} = \frac{\partial^2 f}{\partial x_j x_k}.$$

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

$$\frac{\partial f}{\partial x_{j}} = \sum_{i=1}^{n} 2 (\mathbf{x}^{T} \mathbf{a}^{(i)} - y^{(i)}) \cdot a_{j}^{(i)}$$
$$\frac{\partial^{2} f}{\partial x_{k} \partial x_{j}} = \sum_{i=1}^{n} 2 a_{k}^{(i)} a_{j}^{(i)}$$
$$\mathbf{H} =$$

$$f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$
. Recall that $\nabla f(\mathbf{x}) = 2\mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{b})$.

Claim: If *f* is twice differentiable, then it is convex if and only if the matrix $\mathbf{H} = \nabla^2 f(\mathbf{x})$ is positive semidefinite for all \mathbf{x} .

Definition (Positive Semidefinite (PSD))

A square, symmetric matrix $\mathbf{H} \in \mathbb{R}^{d \times d}$ is <u>positive semidefinite</u> (PSD) for any vector $\mathbf{y} \in \mathbb{R}^d$, $\mathbf{y}^T \mathbf{H} \mathbf{y} \ge 0$.

This is a natural notion of "positivity" for symmetric matrices. To denote that **H** is PSD we will typically use "Loewner order" notation (**succeq** in LaTex):

$\mathbf{H} \succeq \mathbf{0}.$

We write $B \succeq A$ or equivalently $A \preceq B$ to denote that (B - A) is positive semidefinite. This gives a <u>partial ordering</u> on matrices.

Claim: If *f* is twice differentiable, then it is convex if and only if the matrix $\mathbf{H} = \nabla^2 f(\mathbf{x})$ is positive semidefinite for all \mathbf{x} .

Definition (Positive Semidefinite (PSD))

A square, symmetric matrix $\mathbf{H} \in \mathbb{R}^{d \times d}$ is <u>positive semidefinite</u> (PSD) for any vector $\mathbf{y} \in \mathbb{R}^d$, $\mathbf{y}^T \mathbf{H} \mathbf{y} \ge 0$.

For the least squares regression loss function: $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$, $\mathbf{H} = \nabla^2 f(\mathbf{x}) = 2\mathbf{A}^T \mathbf{A}$ for all \mathbf{x} . Is \mathbf{H} PSD? If *f* is β -smooth and α -strongly convex then at any point **x**, $\mathbf{H} = \nabla^2 f(\mathbf{x})$ satisfies:

 $\alpha \mathsf{I} \preceq \mathsf{H} \preceq \beta \mathsf{I},$

where I is a $d \times d$ identity matrix.

This is the natural matrix generalization of the statement for scalar valued functions:

 $\alpha \leq f''(\mathbf{x}) \leq \beta.$

$$\alpha \mathbf{I}_{d \times d} \preceq \mathbf{H} \preceq \beta \mathbf{I}_{d \times d}.$$

Equivalently for any **z**,

$$\alpha \|\mathbf{z}\|_2^2 \le \mathbf{z}^{\mathsf{T}} \mathbf{H} \mathbf{z} \le \beta \|\mathbf{z}\|_2^2.$$

Let $f(\mathbf{x}) = \|\mathbf{D}\mathbf{x} - \mathbf{b}\|_2^2$ where **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 α, β for this problem?

 $\alpha \|\mathbf{z}\|_2^2 \le \mathbf{z}^T \mathbf{H} \mathbf{z} \le \beta \|\mathbf{z}\|_2^2$

GEOMETRIC VIEW



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

GEOMETRIC VIEW



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