CS-GY 6763: Lecture 5 Gradient Descent and Projected Gradient Descent

NYU Tandon School of Engineering, Prof. Christopher Musco

PROJECT

- Choose your partner and email me by end of this week (deadline was originally today).
- Sign-up to present or lead discussion for 1 reading group slot. We need presenters for next week!

Have some function $f : \mathbb{R}^d \to \mathbb{R}$. Want to find \mathbf{x}^* such that:

 $f(\mathbf{x}^*) = \min_{\mathbf{x}} f(\mathbf{x}).$

Or at least $\underline{\hat{x}}$ which is close to a minimum. E.g. $\underline{f(\hat{x})} \le \min_{x} \overline{f(x)} + \epsilon$

Often we have some additional constraints:

$$\cdot \mathbf{X} > \mathbf{0}.$$

$$\cdot \|\mathbf{x}\|_2 \leq R, \|\mathbf{x}\|_1 \leq R.$$

•
$$\mathbf{a}^T \mathbf{x} > c$$
.

CONTINUOUS OPTIMIZATION



Dimension d = 2:



Continuouos optimization is the foundation of modern machine learning.

Supervised learning: Want to learn a model that maps inputs

- numerical data vectors
- images, video
- text documents

to predictions

- numerical value (probability stock price increases)
- label (is the image a cat? does the image contain a car?)
- decision (turn car left, rotate robotic arm)

Let $\underline{M}_{\mathbf{x}}$ be a model with parameters $\mathbf{x} = \{\underline{x_1, \dots, x_k}\}$, which takes as input a data vector **a** and outputs a prediction.

Example:

$$M_{\mathbf{a}}(\mathbf{a}) = \operatorname{sign}(\mathbf{a}^{\mathsf{T}} \mathbf{x})$$

MACHINE LEARNING MODEL

Example:



 $x \in \mathbb{R}^{(\text{\# of connections})}$ is the parameter vector containing all the network weights.

Classic approach in <u>supervised learning</u>: Find a model that works well on data that you already have the answer for (labels, values, classes, etc.).

- Model M_x parameterized by a vector of numbers x.
- Dataset $\mathbf{a}^{(1)}, \ldots, \mathbf{a}^{(n)}$ with outputs $\underline{y}^{(1)}, \ldots, \underline{y}^{(n)}$.

Want to find $\hat{\mathbf{x}}$ so that $M_{\hat{\mathbf{x}}}(\mathbf{a}^{(i)}) \approx \mathbf{y}^{(i)}$ for $i \in 1, ..., n$. How do we turn this into a function minimization problem? **Loss function** $L(M_x(\mathbf{a}), y)$: Some measure of distance between prediction $M_x(\mathbf{a})$ and target output y. Increases if they are further apart.

- Squared (ℓ_2) loss: $|M_x(\mathbf{a}) \underline{y}|^2$
- Absolute deviation (ℓ_1) loss: $|M_x(\mathbf{a}) y|$
- Hinge loss: $1 y \cdot M_x(a)$
- Cross-entropy loss (log loss).
- Etc.

EMPIRICAL RISK MINIMIZATION

Empirical risk minimization:

$$\underbrace{f(\mathbf{x})}_{i=1} = \sum_{i=1}^{n} L\left(M_{\mathbf{x}}(\mathbf{a}^{(i)}), \underline{y}^{(i)}\right)$$

Solve the optimization problem $\min_{\mathbf{x}} f(\mathbf{x})$.

EXAMPLE: LINEAR REGRESSION



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.

ALGORITHMS FOR CONTINUOUS OPTIMIZATION

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. $\|x\|_2 \le c.$ Shure (e)ed

What are some example algorithms for continuous Auouting optimization?

Singlex. Live Programming

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



(and sometimes we can prove it works)

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

Partial derivative:

$$\frac{\partial f}{\partial x_i}(\mathbf{x}) = \lim_{t \to 0} \frac{f(\mathbf{x} + t\mathbf{e}^{(i)}) - f(\mathbf{x})}{t}$$

Directional derivative:

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

Gradient:

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

Directional derivative:

$$D_{\mathbf{v}}f(\mathbf{x}) = \lim_{t \to 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} = \left[\underbrace{\nabla f(\mathbf{x})}_{\mathbf{v}}^{\mathsf{T}} \mathbf{v} \right].$$

$$\Upsilon = e.$$

۱

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

- · Function oracle: Evaluate f(x) for any x. Z unit cost
- 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.

EXAMPLE GRADIENT EVALUATION



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

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

Leading question: When η is small, what's an approximation for $f(\mathbf{x} + \eta \mathbf{v}) - f(\mathbf{x})$?

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

DIRECTIONAL DERIVATIVES

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

$$f(\mathbf{x} + \eta\mathbf{v}) - f(\mathbf{x}) = \int \nabla f(\mathbf{x}) \frac{f(\mathbf{x})}{t} \nabla f(\mathbf{x}) \frac{f(\mathbf{x})}{t} \nabla f(\mathbf{x})^{\mathsf{T}} \nabla f(\mathbf{x})^{\mathsf{T}} \nabla f(\mathbf{x})^{\mathsf{T}}}$$
How should we choose v so that $f(\mathbf{x} + \eta\mathbf{v}) < f(\mathbf{x})$?

$$\int (-\eta \nabla f(\mathbf{x})) \frac{f(\mathbf{x})}{t} \nabla f(\mathbf{x}) \frac{f(\mathbf{x})}{t} \nabla f(\mathbf{x}) \frac{f(\mathbf{x})}{t} \nabla f(\mathbf{x})}$$

Prototype algorithm:

- Choose starting point $\mathbf{x}^{(0)}$.
- For $i = 0, \dots, \overline{T}$

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

• 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.

1 dimensional example:



0

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**:

 $f(\mathbf{x}^{(T)}) < f(\mathbf{x}^*) + \epsilon$.

Examples: least squares regression, logistic regression, kernel regression, SVMs.

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, mixture models.

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** only.

CONVEXITY

Definition (Convex)

A function *f* is convex iff for any $\mathbf{x}, \mathbf{y}, \lambda \in [0, 1]$:

$$(1-\lambda)\cdot f(\mathbf{x}) + \lambda \cdot f(\mathbf{y}) \geq f((1-\lambda)\cdot \mathbf{x} + \lambda \cdot \mathbf{y})$$



GRADIENT DESCENT

Definition (Convex)

A function *f* is convex if and only if for any **x**, **y**:

Equivalently:





 $\chi^{(1)}$

Assume:

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

Gradient descent:

- Choose number of steps T.
- Starting point $\mathbf{x}^{(0)}$. E.g. $\mathbf{x}^{(0)} = \vec{0}$.
- $\cdot \eta = \frac{R}{G_{2}/T}$
- For i = 0, ..., T:

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

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

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





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 $\min_i f(\mathbf{x}^{(i)}) \leq \frac{1}{T} \sum_{i=0}^{T-1} f(\mathbf{x}^{(i)})$, so this is what we return:

$$f(\hat{\mathbf{x}}) = \min_{i \in 1, \dots, T} f(\mathbf{x}^{(i)}) \le f(\mathbf{x}^*) + \epsilon.$$

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

 $\min_{\mathbf{x}\in\mathcal{S}}f(\mathbf{x})$

where S is a **convex set**.



Which of these is convex?

CONSTRAINED CONVEX OPTIMIZATION



Definition (Convex set)

A set S is convex if for any $\mathbf{x}, \mathbf{y} \in S, \lambda \in [0, 1]$:



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 $\mathbf{x}^{(0)} \in S$, there is no guarantee that $\mathbf{x}^{(0)} - \eta \nabla f(\mathbf{x}^{(0)})$ 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 \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} .

$$P_{\mathcal{S}}(\mathbf{x}) = \arg\min_{\mathbf{y}\in\mathcal{S}} \|\mathbf{x} - \mathbf{y}\|_{2}$$
PROJECTION ORACLES



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
$$i = 0, \dots, T$$
:

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

$$\overbrace{\chi}^{\bullet} = \alpha \mathcal{S}_{\mathcal{M}} \circ \eta \quad f(x)$$

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

$$\chi^{4} - \chi^{(i+i)} |_{\gamma}^{2} \in ||\chi^{4} - \chi|_{\gamma}^{2}$$

Claim (PGD Convergence Bound)

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

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



Claim (PGD Convergence Bound)

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



Same telescoping sum argument:

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

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):

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

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

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

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).



JX

SMOOTHNESS

Definition (β -smoothness)

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

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

After some calculus (see Lem. 3.4 in **Bubeck's book**), this implies: $\mathcal{O} \leq [f(\mathbf{y}) - f(\mathbf{x})] - \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) \leq \frac{\beta}{2} ||\mathbf{x} - \mathbf{y}||_2^2$ $\{f_0\}_{\substack{i \leq i \\ i \leq i \\ j \leq$ Recall from definition of convexity that:

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

So now we have an upper and lower bound.

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

GUARANTEED PROGRESS

Previously learning rate/step size η depended on <u>G</u>. 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:

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

$$= \frac{1}{\mathcal{B}} \nabla f(\mathbf{x}^{(t)})$$

$$\left[f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)}) \right] + \frac{1}{\beta} \|\nabla f(\mathbf{x}^{(t)})\|_{2}^{2} \le \frac{4\beta}{\mathcal{B}} \|\mathcal{F} \nabla f(\mathbf{x}^{(t)})\|_{2}^{2}$$

$$= \frac{1}{\mathcal{B}} \|\nabla f(\mathbf{x}^{(t)})\|_{2}^{2}$$

45

CONVERGENCE GUARANTEE

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

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

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

STRONG CONVEXITY

Definition (α -strongly convex) A convex function f is α -strongly convex if, for all x, y $[f(\mathbf{y}) - f(\mathbf{x})] - \nabla f(\mathbf{x})^{\mathsf{T}} (\mathbf{y} - \mathbf{x}) \geq \frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$ S B TT α is a parameter that will depend on our function. አ For a twice-differentiable scalar valued function f_{i} equivalent F"(X) < B to $f''(x) > \alpha$.

Gradient descent for strongly convex functions:

- Choose number of steps T.
- For i = 1, ..., T: • $\eta = \underbrace{\begin{pmatrix} 2 \\ \underline{\alpha \cdot (i+1)} \end{pmatrix}}_{\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$

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

CONVERGENCE GUARANTEE

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

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

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

$$q \leq \int^{d} (x) \leq B$$

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

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:

$$\frac{2}{6}\left[f(x^{(T)})-f(x^{0})\right] \leq \|\mathbf{x}^{(T)}-\mathbf{x}^{*}\|_{2}^{2} \leq e^{-(T-1)}\left\|\mathbf{x}^{(1)}-\mathbf{x}^{*}\|_{2}^{2}\right] \leq e^{-T}\left\|\mathbf{x}^{(T)}-\mathbf{x}^{*}\|_{2}^{2} \leq e^{-T}\left\|\mathbf{x}^{(T)}-\mathbf{x}^{*}\|_{2}^{2}\right]$$

 $\kappa = \frac{\beta}{\alpha}$ is called the "condition number" of *f*. Is it better if κ is large or small?

SMOOTH AND STRONGLY CONVEX

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

$$\sum_{n=1}^{\infty} \|\mathbf{x} - \mathbf{y}\|_2^2 \leq \mathbf{y} + \mathbf{y} - [f(\mathbf{x}) - f(\mathbf{y})] \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2,$$

we have:

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

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:

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

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

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \le \epsilon \left[f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*)\right]$$
Alternative Corollary: If $T = O\left(\frac{\beta}{\alpha}\log(R\beta/\epsilon)\right)$ we have:

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

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



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 **x**. So $\mathbf{H} \in \mathbb{R}^{d \times d}$. We have: $\mathbf{H} = \mathcal{P} \mathbf{A}^{\dagger} \mathbf{A}$ $\mathbf{H}_{i,j} = \left[\nabla^2 f(\mathbf{x})\right]_{i,j} = \frac{\partial^2 f}{\partial x_i x_i}.$ $\mathbf{q}'_i = i \mathcal{H}_{i,j} \quad \text{of } \mathcal{A}$ **Example:** Let $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$. Recall that $\nabla f(\mathbf{x}) = 2\mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{b})$. $\frac{1}{2} \frac{\partial f}{\partial x_{i}} = \left(2\sigma_{i}^{T} (Ax - b) \right)$ $= \left(\frac{1}{2} \sigma_{i}^{T} (Ax - b) \right)$ $= \left(\frac{1}{2} \sigma_{i}^{T} (Ax - b) - 2\sigma_{i}^{T} (Ax - b) \right)$ $= \left(\frac{1}{2} \sigma_{i}^{T} A = \frac{1}{2} \sigma_{i}^{T} A$ $a_1 | a_2 | ...$ φ'j = 2 α;τα;

h

Α

55

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 $\underline{\mathbf{y}} \in \mathbb{R}^d$, $\mathbf{y}^T \mathbf{H} \mathbf{y} \ge 0$.

This is a natural notion of "positivity" for symmetric matrices. To denote that **H** is PSD we will typically use "<u>Loewner orde</u>r" 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 partial ordering 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: $f(\mathbf{x}) = \|\underline{\mathbf{A}\mathbf{x} - \mathbf{b}}\|_2^2$, $\mathbf{H} = \nabla^2 f(\mathbf{x}) = \underline{2\mathbf{A}^T \mathbf{A}}$ for all \mathbf{x} . Is \mathbf{H} PSD?

2= AZ

If *f* is β -smooth and α -strongly convex then at any point **x**, $\mathbf{H} = \nabla^2 f(\mathbf{x})$ satisfies:

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

where $I_{d \times d}$ 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**,

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

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 $d_{1}^{2} = \frac{1}{3}, d_{2}^{2} = 2$.

Any symmetric matrix **H** has an <u>orthogonal</u>, real valued eigendecomposition.



Here V is square and orthogonal, so $V^T V = V V^T = I$. And for each v_i , we have:

 $\mathbf{H}\mathbf{v}_i = \lambda_i \mathbf{v}_i.$

By definition, that's what makes $\mathbf{v}_1, \ldots, \mathbf{v}_d$ eigenvectors.

Recall $VV^{T} = V^{T}V = I$.



Claim: H is PSD $\Leftrightarrow \lambda_1, ..., \lambda_d \ge 0$.

Recall $VV^{T} = V^{T}V = I$.



Claim: $\alpha I \preceq H \preceq \beta I \Leftrightarrow \alpha \leq \lambda_1, ..., \lambda_d \leq \beta$.

Recall $VV^{T} = V^{T}V = I$.



In other words, if we let $\lambda_{max}(H)$ and $\lambda_{min}(H)$ be the smallest and largest eigenvalues of H, then for all z we have:

$$\begin{split} \mathbf{z}^{\mathsf{T}}\mathbf{H}\mathbf{z} &\leq \lambda_{\mathsf{max}}(\mathbf{H}) \cdot \|\mathbf{z}\|^2 \\ \mathbf{z}^{\mathsf{T}}\mathbf{H}\mathbf{z} &\geq \lambda_{\mathsf{min}}(\mathbf{H}) \cdot \|\mathbf{z}\|^2 \end{split}$$

If the maximum eigenvalue of $\mathbf{H} = \nabla^2 f(\mathbf{x}) = \beta$ and the minimum eigenvalue of $\mathbf{H} = \nabla^2 f(\mathbf{x}) = \alpha$ then $f(\mathbf{x})$ is β -smooth and α -strongly convex.

 $\lambda_{\max}(\mathsf{H}) = \beta$ $\lambda_{\min}(\mathsf{H}) = \alpha$

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 $\eta = \frac{2}{\beta}$) we have:

$$\|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2 \le e^{-T/\kappa} \|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2$$

Goal: Prove for $f(x) = ||Ax - b||_2^2$.

Let $\lambda_{\max} = \lambda_{\max}(\mathbf{A}^T \mathbf{A})$. Gradient descent update is:

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \frac{1}{2\lambda_{\max}} 2\mathbf{A}^{\mathsf{T}} (\mathbf{A}\mathbf{x}^{(t)} - \mathbf{b})$$

Richardson Iteration view:

$$(\mathbf{x}^{(t+1)} - \mathbf{x}^*) = \left(\mathbf{I} - \frac{1}{\lambda_{\max}} \mathbf{A}^T \mathbf{A}\right) (\mathbf{x}^{(t)} - \mathbf{x}^*)$$

What is the maximum eigenvalue of the symmetric matrix $\left(\mathbf{I} - \frac{1}{\lambda_{\max}} \mathbf{A}^T \mathbf{A}\right)$ in terms of the eigenvalues $\lambda_{\max} = \lambda_1 \ge \ldots \ge \lambda_d = \lambda_{\min}$ of $\mathbf{A}^T \mathbf{A}$?

UNROLLED GRADIENT DESCENT

$$(\mathbf{x}^{(T+1)} - \mathbf{x}^*) = \left(\mathbf{I} - \frac{1}{\lambda_{\max}} \mathbf{A}^T \mathbf{A}\right)^T (\mathbf{x}^{(1)} - \mathbf{x}^*)$$

What is the maximum eigenvalue of the symmetric matrix $\left(I - \frac{1}{\lambda_{max}} \mathbf{A}^T \mathbf{A}\right)^T$?

So we have
$$\|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2 \leq 1$$

We now have a pretty good understanding of gradient descent. Number of iterations for ϵ error:

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

ACCELERATION
ACCELERATED GRADIENT DESCENT

Nesterov's accelerated gradient descent:

•
$$x^{(1)} = y^{(1)} = z^{(1)}$$

For
$$t = 1, ..., T$$

• $\mathbf{y}^{(t+1)} = \mathbf{x}^{(t)} - \frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})$
• $\mathbf{x}^{(t+1)} = \left(1 + \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}\right) \mathbf{y}^{(t+1)} + \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \left(\mathbf{y}^{(t+1)} - \mathbf{y}^{(t)}\right)$

Theorem (AGD for β **-smooth,** α **-strongly convex.)** Let f be a β -smooth and α -strongly convex function. If we run AGD for T steps we have:

$$f(\mathbf{x}^{(t)}) - f(\mathbf{x}^*) \le \kappa e^{-(t-1)\sqrt{\kappa}} \left[f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*) \right]$$

Corollary: If $T = O(\sqrt{\kappa} \log(\kappa/\epsilon))$ achieve error ϵ .

INTUITION BEHIND ACCELERATION



Level sets of $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$.

Other terms for similar ideas:

- Momentum
- Heavy-ball methods

What if we look back beyond two iterates?