

*Believe Quiz  
Median*

## CS-GY 9223 I: Lecture 7

Smooth functions

Preconditioning, acceleration, coordinate

descent, etc.

Stagnant functions

---

$$\frac{df}{dx} = \lim_{t \rightarrow 0} \frac{f(x) - f(x+t)}{t}$$

NYU Tandon School of Engineering, Prof. Christopher Musco

$$2A^T(Ax - b) = f(x)$$

$$\lim_{z \rightarrow 0} \left[ 2A^T A x - A^T b \right] - \left[ 2A^T A (x+z) \right] - A^T b = \underline{\underline{2A^T Az}}$$

## LOGISTICS

- Self-proctored, 2-hour midterm to be taken anytime next week.
- No Collaboration allowed at all. Or outside resources. Just use your own notes and material from the class.
- Sample problems are available on course website. We can review during office hours tomorrow or next week.
- You should have received an invite to Gradescope. Hopefully tonight/tomorrow I can upload a "practice test" to make sure their system works.
- Demos + Calculator + Coding is allowed  
+ Wof from Alpha

## GRADIENT DESCENT

Conditions:

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

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

- Bounded gradients (Lipschitz function):

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

### Theorem

*GD Convergence Bound] (Projected) Gradient Descent returns  $\hat{\mathbf{x}}$  with  $f(\hat{\mathbf{x}}) \leq \min_{\mathbf{x} \in \mathcal{S}} f(\mathbf{x}) + \epsilon$  after*

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

## ONLINE GRADIENT DESCENT

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

Conditions:

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

Theorem (OGD Regret Bound)

After  $T$  steps,  $\left[ \sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \right] - \boxed{\left[ \sum_{i=1}^T f_i(\mathbf{x}^*) \right]} \leq RG\sqrt{T}$ . I.e. the average regret  $\frac{1}{T} \left[ \sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \right]$  is  $\leq \epsilon$  after:

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

# STOCHASTIC GRADIENT DESCENT

Conditions:

- Finite sum structure:  $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$ , with  $f_1, \dots, f_n$  all convex.
- Lipschitz functions: for all  $\mathbf{x}, j$ ,  $\|\nabla f_j(\mathbf{x})\|_2 \leq \frac{G'}{n}$ .
- Starting radius:  $\|\mathbf{x}^* - \mathbf{x}^{(1)}\|_2 \leq R$ .



## Theorem (SGD Regret Bound)

Stochastic Gradient Descent returns  $\hat{\mathbf{x}}$  with

$\mathbb{E}[f(\hat{\mathbf{x}})] \leq \min_{\mathbf{x} \in S} f(\mathbf{x}) + \epsilon$  after

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

We always have that  $G' > G$  but iterations are typically cheaper by a factor of  $n$ .

## BEYOND THE BASIC BOUNDS

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

### Goals:

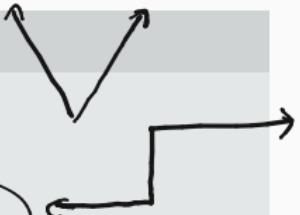
- Improve  $\epsilon$  dependence below  $1/\epsilon^2$ .
  - Ideally  $1/\epsilon$  or  $\log(1/\epsilon)$ .
- Reduce or eliminate dependence on  $G$  and  $R$ .
- Further take advantage of structure in the data (e.g. repetition in features in addition to data points).

## SMOOTHNESS

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

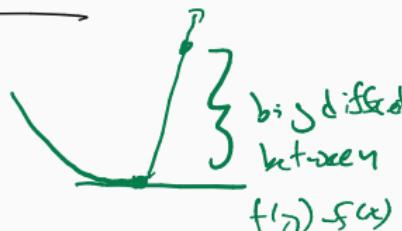
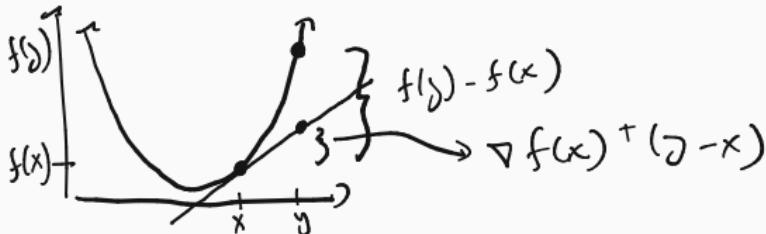
A function  $f$  is  $\beta$  smooth if, for all  $x, y$

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



After some calculus (see Lem 3.4 in Rübeck's book), this implies:

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



For a scalar valued function  $f$ , equivalent to  $\underline{f''(x)} \leq \beta$ .  $\nabla f(x)^T (y - x)$

## SMOOTHNESS

Recall from definition of convexity that:

$$\underbrace{f(y) - f(x)}_{\text{upper bound}} \geq \underbrace{\nabla f(x)^T(y - x)}_{\text{lower bound}}$$

So now we have an upper and lower bound.

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

## GUARANTEED PROGRESS

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

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

Progress per step of gradient descent:

*definition of smooth*

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

$\downarrow$

-  $\frac{1}{\beta} \nabla f(\underline{x^{(t)}})$

$\gamma_\beta \nabla f(\underline{x^{(t)}})$

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

$\frac{1}{2\beta}$

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

## CONVERGENCE GUARANTEE

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:

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

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

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

## STRONG CONVEXITY

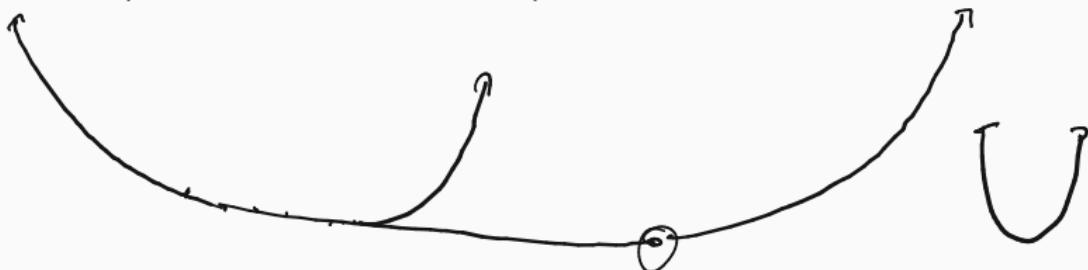
Definition ( $\alpha$ -strongly convex)

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

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

$\geq 0$

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



For a twice-differentiable scalar valued function  $f$ , equivalent to  $f''(x) \geq \alpha$ .

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

Gradient descent for strongly convex functions:

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

## CONVERGENCE GUARANTEE

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

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

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

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

$$\overbrace{A^T}^{g}$$

$$A^S \checkmark$$

$$\overbrace{f(x)}^{g}$$

$$(A^T A) \xrightarrow{\text{nd}} \underbrace{A}_{d^2}$$

$$\text{nd} \cdot 3$$

$$\frac{nd^2 \log(d)}{1}$$

$$A \cdot (A \cdot (A \cdot b))$$

$$\overbrace{\text{nd} g}^{g}$$

$$\overbrace{R}^{\text{R}} \quad \overbrace{A^2}^{\text{A}^2} \quad \overbrace{A^4}^{\text{A}^4}$$

## CONVERGENCE GUARANTEE

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

$$x - y \leq -(\gamma - \alpha)$$

$$\begin{aligned} \frac{\alpha}{2} \|x - y\|_2^2 &\leq \nabla f(x)^T (x - y) - [f(x) - f(y)] \leq \frac{\beta}{2} \|x - y\|_2^2. \\ &= \left\{ f(y) - f(x) \right\} - \nabla f(x)^T (y - x) \end{aligned}$$

## CONVERGENCE GUARANTEE

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

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:

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

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

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

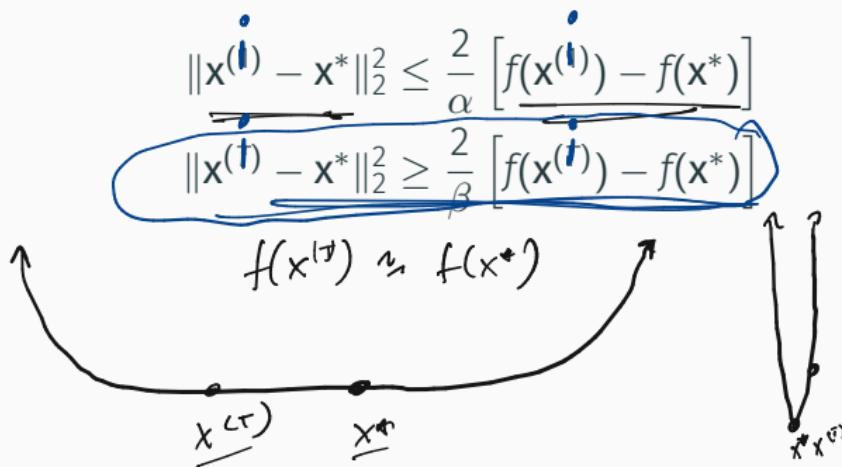
$$T = \Theta\left(\frac{\beta}{\alpha} \log \frac{1}{\epsilon}\right) \quad e^{-\frac{1}{2} \log \frac{1}{\epsilon}} = e^{\log \frac{1}{\epsilon}} = \frac{1}{\epsilon}$$

## SMOOTH AND STRONGLY CONVEX

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

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

we have:



## CONVERGENCE GUARANTEE

**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:

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

$\epsilon \approx \frac{\beta \beta^2}{2} \|\mathbf{x}^{(1)} - \mathbf{x}^*\|^2$

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

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \epsilon [f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*)]$$

$\epsilon \approx \kappa$

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

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \epsilon \approx n \log\left(\frac{R\beta^{2/2}}{\alpha}\right)$$

17

## THE LINEAR ALGEBRA OF CONDITIONING

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

$$H_{i,j} = [\nabla^2 f(x)]_{i,j} \xrightarrow{\text{symmetric matrix}} = \frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_i};$$

For vector  $x, y$ :

$$\underline{\nabla f(x) - \nabla f(y)} \approx \underbrace{[\nabla^2 f(x)]}_{H^{d \times d}} \underbrace{(x - y)}_{\mathbb{R}^d}.$$

# THE LINEAR ALGEBRA OF CONDITIONING

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

$$\frac{\partial}{\partial x_i} \left( \frac{\partial f}{\partial x_i} \right)$$

$$H_{i,j} = [\nabla^2 f(x)]_{i,j} = \frac{\partial^2 f}{\partial x_i \partial x_j}.$$

$$\begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}$$

**Example:** Let  $f(x) = \|Ax - b\|_2^2$ . Recall that  $\nabla f(x) = 2A^T(Ax - b)$ .

$$A \quad \begin{matrix} a_1 & a_2 & \dots & a_d \\ \cdot & \cdot & \cdot & \cdot \end{matrix} \quad x$$

$$\begin{aligned} \nabla^2 f(x) &= 2A^T A \\ [\nabla f(x)]_i &= \frac{\partial f}{\partial x_i} = 2a_i^T (Ax - b) \\ \frac{\partial^2 f}{\partial x_i \partial x_j} &= \lim_{t \rightarrow 0} \frac{2a_i^T (Ax - b) - 2a_i^T (A(x + te_j) - b)}{t} \\ &\stackrel{?}{=} 2a_i^T A e_j \end{aligned}$$

**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 positive semidefinite (PSD) for any vector  $\mathbf{y} \in \mathbb{R}^d$ ,  $\mathbf{y}^T \mathbf{H} \mathbf{y} \geq 0$ .

$$\begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix} \quad \mathbf{y}^T \mathbf{I}_d \mathbf{y}$$

$$= \mathbf{y}^T \mathbf{y} = \|\mathbf{y}\|_2^2$$

This is a natural notion of “positivity” for symmetric matrices. To denote that  $\mathbf{H}$  is PSD we will typically use “Loewner order” notation (\succeq in LaTeX):

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

We write  $\mathbf{B} \succeq \mathbf{A}$  or equivalently  $\mathbf{A} \not\succeq \mathbf{B}$  to denote that  $(\mathbf{B} - \mathbf{A})$  is positive semidefinite. This gives a partial ordering on matrices.

$$\mathbf{B} = \begin{bmatrix} 1 & 0 \\ 0 & -2 \end{bmatrix}$$

$$\mathbf{A} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$

**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 positive semidefinite (PSD) for any vector  $\mathbf{y} \in \mathbb{R}^d$ ,  $\mathbf{y}^T \mathbf{H} \mathbf{y} \geq 0$ .

For the least squares regression loss function:

$$f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2, \mathbf{H} = \nabla^2 f(\mathbf{x}) = 2\mathbf{A}^T \mathbf{A} \text{ for all } \mathbf{x}. \text{ Is } \mathbf{H} \text{ PSD?}$$

$$\begin{aligned} \mathbf{y}^T 2\mathbf{A}^T \mathbf{A} \mathbf{y} &= 2 \mathbf{y}^T \mathbf{A}^T \mathbf{A} \mathbf{y} = 2 (\mathbf{A} \mathbf{y})^T \mathbf{A} \mathbf{y} \\ &= 2 \|\mathbf{A} \mathbf{y}\|_2^2 \geq 0 \end{aligned}$$

## THE LINEAR ALGEBRA OF CONDITIONING

If  $f$  is  $\beta$ -smooth and  $\alpha$ -strongly convex then at any point  $x$ ,  
 $H = \nabla^2 f(x)$  satisfies:

$$\begin{aligned} & (\beta \cdot I - H) \succ 0 \\ & (H - \alpha I) \succ 0 \\ & \alpha I_{d \times d} \preceq H \preceq \beta I_{d \times d}, \\ & z^T H z - z^T \alpha I z \\ & = \underbrace{z^T H z}_{\geq 0} - \underbrace{\alpha \|z\|_2^2}_{\geq 0} \geq 0 \end{aligned}$$

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:

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

$$z^T \underline{\beta} \cdot I z - z^T H z \geq 0$$

$$\underline{\beta} \|z\|_2^2 - z^T H z \geq 0$$

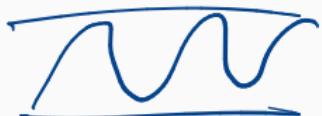
## SMOOTH AND STRONGLY CONVEX HESSIAN

$$S_{\text{smooth}}(x)$$

$$S_{\text{strongly convex}}(x) = -c \rightarrow (x)$$

$$S_{\text{strongly concave}}''(x) = -\underline{s''}(x)$$

$$\alpha I_{d \times d} \preceq H \preceq \beta I_{d \times d}$$



Equivalently for any  $z$ ,

$$\alpha \|z\|_2^2 \leq z^T H z \leq \beta \|z\|_2^2.$$

**Exercise:** Show that for  $f(x) = \|Ax - b\|_2^2$ ,

$$\underbrace{[f(\overset{\gamma}{y}) - f(\overset{x}{y})] - \nabla f(x)^T(y - x)}_{\text{Left side of the inequality}} = \underbrace{(x - y)^T [2A^T A](x - y)}_{\text{Right side of the inequality}}.$$

This would imply:

$$z = x - y$$

$$\frac{\alpha}{2} \|x - y\|_2^2 \leq \underbrace{[f(\overset{\gamma}{y}) - f(\overset{x}{y})] - \nabla f(x)^T(y - x)}_{\text{Left side of the inequality}} \leq \frac{\beta}{2} \|x - y\|_2^2$$

## SIMPLE EXAMPLE

Let  $f(x) = \|\underline{Dx - b}\|_2^2$  where  $D$  is a diagonal matrix. For now imagine we're in two dimensions:  $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ ,  $D = \begin{bmatrix} d_1 & 0 \\ 0 & d_2 \end{bmatrix}$ .

What are  $\alpha, \beta$  for this problem?

$$2 D^T D$$

$$\alpha \|z\|_2^2 \leq z^T H z \leq \beta \|z\|_2^2$$

$$\underline{\alpha} = \min(d_1^2, d_2^2)$$

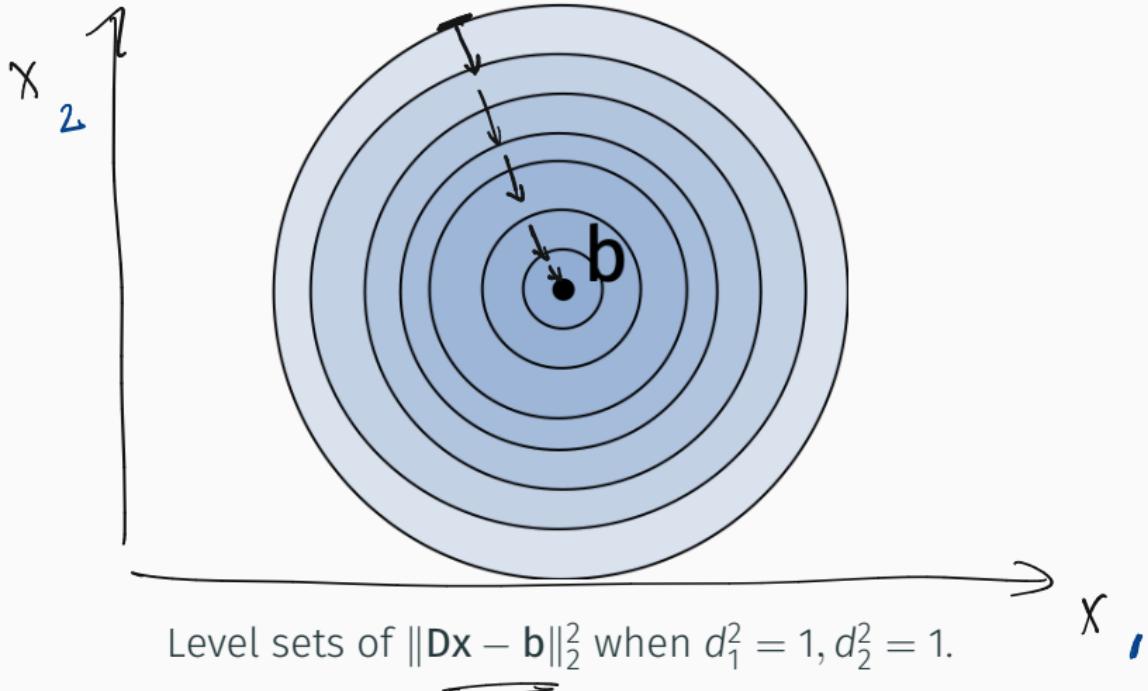
$$z = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\underline{\beta} = \max(d_1^2, d_2^2)$$

$$z = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad z^T (2 D^T D) z = 2 d_2^2 \quad \|z\|_2^2$$

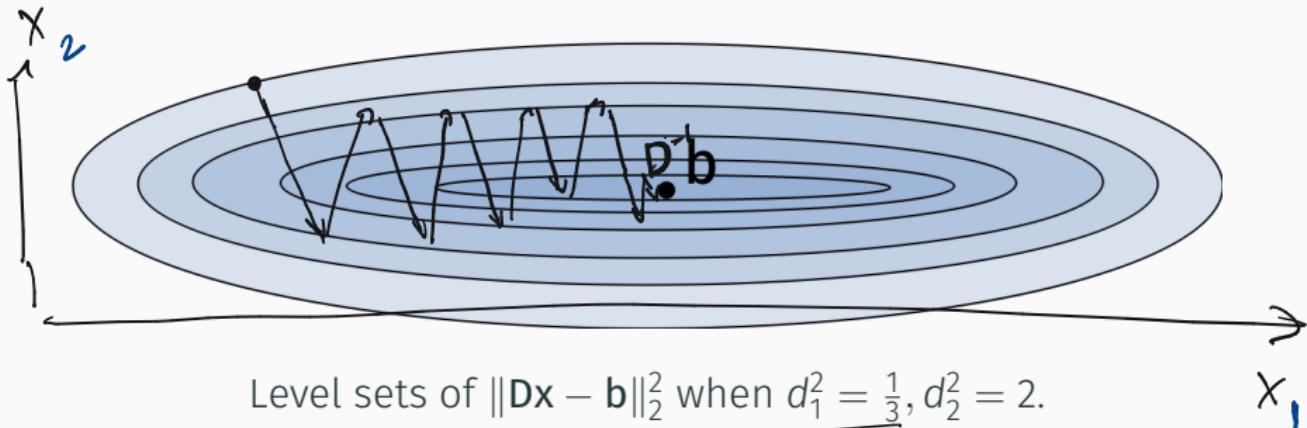
$$z = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad z^T (2 D^T D) z = 2 d_1^2 + d_2^2 \quad \|z\|_2^2$$

## GEOMETRIC VIEW



$$= \|x - b\|_2^2 \quad \alpha = \beta \quad K = \frac{\beta}{\alpha} = 1$$

## GEOMETRIC VIEW

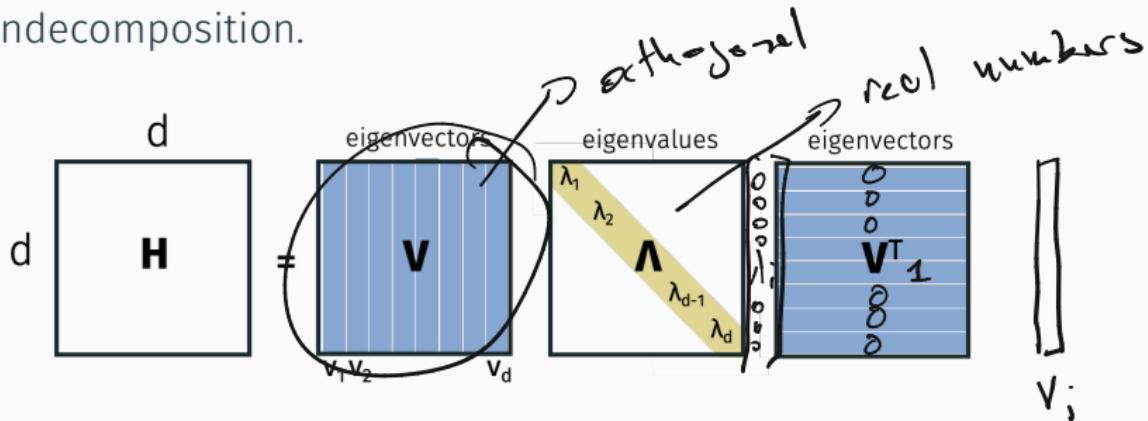


Level sets of  $\|Dx - b\|_2^2$  when  $d_1^2 = \frac{1}{3}$ ,  $d_2^2 = 2$ .

$$\alpha = \frac{1}{3}, \quad \beta = 2, \quad h = \frac{2}{\sqrt{3}} = \frac{2}{3}\sqrt{3} = 6$$

## EIGENDECOMPOSITION VIEW

Any symmetric matrix  $\mathbf{H}$  has an orthogonal, real valued eigendecomposition.



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

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

That's what makes  $\mathbf{v}_1, \dots, \mathbf{v}_d$  eigenvectors.

## EIGENDECOMPOSITION VIEW

Recall  $VV^T = V^TV = I$ .

$$\sqrt{H} = \begin{pmatrix} f_1 & f_2 & \dots & f_d \end{pmatrix}$$

$$H = V \underbrace{\sqrt{\Lambda}}_{A^T} \cdot \underbrace{\sqrt{\Lambda}}_{A} V^T$$

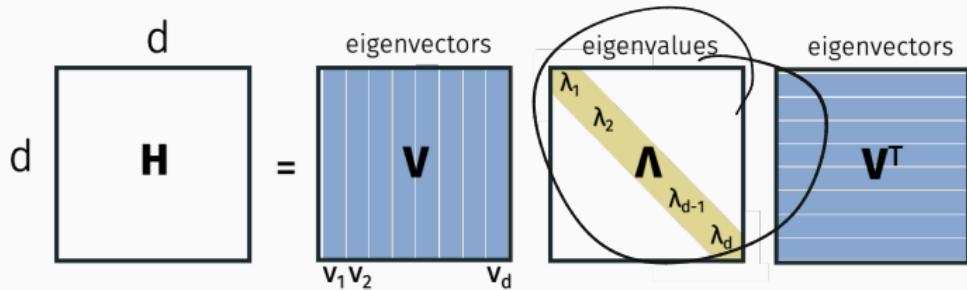
$$d \begin{matrix} H \\ d \end{matrix} = \begin{matrix} \text{eigenvectors} \\ V \\ v_1 v_2 \dots v_d \end{matrix} \begin{matrix} \text{eigenvalues} \\ \Lambda \\ \lambda_1 \lambda_2 \dots \lambda_d \end{matrix} \begin{matrix} \text{eigenvectors} \\ V^T \\ V^T \end{matrix} \begin{matrix} \gamma^T H \gamma \\ = \|A\gamma\|_2^2 \end{matrix}$$

Claim:  $H \Leftrightarrow \lambda_1, \dots, \lambda_d \geq 0$ .   
 ↪  $\text{PSD}$

$$\underline{v_2^T H v_2} = v_2 \cdot \lambda_2 v_2 = \lambda_2 \|v_2\|_2^2$$

## EIGENDECOMPOSITION VIEW

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



Claim:  $\alpha I \preceq H \preceq \beta I \Leftrightarrow \underbrace{\alpha \leq \lambda_1, \dots, \lambda_d \leq \beta}$ .

$$H \succcurlyeq \alpha I \quad H - \alpha I \succcurlyeq 0$$

$$= V \Lambda V^T - \alpha V V^T = V \underbrace{(I - \alpha I)}_{\text{symmetric}} V^T$$

## EIGENDECOMPOSITION VIEW

Recall  $\mathbf{V}\mathbf{V}^T = \mathbf{V}^T\mathbf{V} = \mathbf{I}$ .

$$\begin{matrix} d \\ d \end{matrix} \mathbf{H} = \begin{matrix} \text{eigenvectors} \\ \mathbf{V} \end{matrix} \begin{matrix} \text{eigenvalues} \\ \Lambda \end{matrix} \begin{matrix} \text{eigenvectors} \\ \mathbf{V}^T \end{matrix}$$

The diagram illustrates the eigen decomposition of a  $d \times d$  matrix  $\mathbf{H}$ . The matrix  $\mathbf{H}$  is shown as a large square divided into four quadrants. The top-left quadrant is labeled  $\mathbf{H}$ , and the bottom-right quadrant contains the eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_d$  on a yellow diagonal. The left column of the matrix is labeled  $v_1, v_2, \dots, v_d$  at the bottom, and the right column is labeled  $v_1, v_2, \dots, v_d$  at the top. The middle column is labeled  $v$  at both the top and bottom. The word "eigenvectors" is written above the first and last columns, and the symbol  $\Lambda$  is placed between the two columns of eigenvalues.

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

$$\mathbf{z}^T \mathbf{H} \mathbf{z} \leq \lambda_{\max}(\mathbf{H}) \cdot \|\mathbf{z}\|^2$$

$$\mathbf{z}^T \mathbf{H} \mathbf{z} \geq \lambda_{\min}(\mathbf{H}) \cdot \|\mathbf{z}\|^2$$

## EIGENDECOMPOSITION VIEW

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

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

$$\mathbf{H} = 2\mathbf{A}^\top \mathbf{A}$$

## POLYNOMIAL VIEW POINT

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:

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

$$\frac{\beta}{2} = \lambda_{\max}(A^T A)$$

$$\frac{1}{2\lambda_{\max}(A^T A)} = \frac{1}{\beta}$$

$$\underline{\beta} = \lambda_{\max}(2A^T A)$$

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

$$\frac{\beta}{2} = \lambda_{\max}(A^T A)$$

$$\left( I - \frac{1}{\lambda_{\max}} A^T A \right)$$

$$VV^T - \frac{1}{\lambda_{\max}} V \Lambda V^T =$$

$$V \left( I - \frac{1}{\lambda_{\max}} \Lambda \right) V^T$$

## ALTERNATIVE VIEW OF GRADIENT DESCENT

$\lambda_{\max} = \lambda_{\max}(A^T A)$   
Richardson Iteration view:

$$\omega$$

$$\frac{\beta}{\alpha} = \lambda_{\max}$$

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

$$(x^{(T+1)} - x^*) = \underbrace{x^{(T)} - \frac{1}{2\lambda_{\max}} \cdot \nabla f(x^{(T)})}_{2A^T A x^{(T)} - 2A^T b} - x^*$$

$$= \underbrace{x^{(T)} - \frac{1}{\lambda_{\max}} A^T A x^{(T)}}_{A^T b} - \frac{1}{\lambda_{\max}} A^T b - x^*$$

$$\underline{A^T b} = \underline{A^T A x^*} \quad \nabla(A^T A x^* - A^T b) = 0$$

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

$$1 - \frac{1}{\lambda_{\max}}$$

$$1 - \frac{\lambda_1}{\lambda_{\max}} = 0$$

$$1 - \frac{\lambda_{\min}}{\lambda_{\max}}$$

## UNROLLED GRADIENT DESCENT

$$(x^{(T+1)} - x^*) = \underbrace{\left( I - \frac{1}{\lambda_{\max}} A^T A \right)}_{\text{Matrix}} \underbrace{(x^{(1)} - x^*)}_{\text{Initial vector}} + \text{number of iterations}$$

$\kappa = \frac{\lambda_{\max}}{\lambda_{\min}}$

What is the maximum eigenvalue of the symmetric matrix

$$\left( I - \frac{1}{\lambda_{\max}} A^T A \right) \stackrel{?}{=} \left( I - \frac{\lambda_i}{\lambda_{\max}} \right)^+ \quad \left( I - \frac{\lambda_{\min}}{\lambda_{\max}} \right)^+ \quad \underbrace{\left( \left( 1 - \frac{1}{\kappa} \right) \right)^+}_e$$

$$(V \Lambda V^T)(V \Lambda V^T) \dots (V \Lambda V^T) \underbrace{\left( \left( 1 - \frac{1}{\kappa} \right) \right)^+}_e$$

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

$$\left( \frac{1}{\kappa} \right)^{T/\kappa} = e^{-T/\kappa}$$

## IMPROVING GRADIENT DESCENT

$\checkmark \quad (\mathbb{I} - \frac{1}{\lambda_{\max}} A^T A) \rightsquigarrow (x^{(t+)} - x^*)$        $t = \# \text{ of iterations}$

$w \cdot w \cdot w \cdot w \dots w$        $z^{2+}$

We now have a really good understanding of gradient descent.

Number of iterations for  $\epsilon$  error:

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

How do we use this understanding to design faster algorithms?

$$z = (x^{(1)} - x^*) \quad r^{(t+)} = \underline{w^+ z} \quad r^{(t+)} = x^{(t+)} - x^*$$

$$\|r^{(t+)}\|_2^2 = \underline{z^T w^{2+} z} \leq \lambda_{\max}(w^{2+}) \|z\|_2^2 \\ = \lambda_{\max}(w^{2+}) \|x^{(1)} - x^*\|_2^2$$

## IMPORTANCE SAMPLING

Often it doesn't make sense to sample  $i$  uniformly at random:

$$A = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 & 0 \\ 0 & 0 & -0.5 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & -2 & 0 & 0 & 0 \end{bmatrix} \quad b = \begin{bmatrix} 10 \\ 42 \\ -11 \\ -51 \\ 34 \\ -22 \end{bmatrix}$$

Select indices  $i$  proportional to  $\|a_i\|_2^2$ :

$$\Pr[\text{select index } i \text{ to update}] = \frac{\|a_i\|_2^2}{\sum_{j=1}^d \|a_j\|_2^2} = \frac{\|a_i\|_2^2}{\|A\|_2^2}$$