

# CS-GY 6763: Lecture 6

## Gradient Descent and Projected Gradient Descent

---

NYU Tandon School of Engineering, Prof. Christopher Musco

Given function  $f: \mathbb{R}^d \rightarrow \mathbb{R}^n$ . Find  $\hat{\mathbf{x}}$  such that:

$$f(\hat{\mathbf{x}}) \leq \min_{\mathbf{x}} f(\mathbf{x}) + \epsilon.$$

## EXAMPLE: LEAST SQUARES REGRESSION

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

- Model:  $M_{\mathbf{x}}(\mathbf{a}) = \mathbf{x}^T \mathbf{a}$ .  $\mathbf{x}$  contains the regression coefficients.
- Loss function:  $L(z, y) = |z - y|^2$ .
- Function to minimize:  $f(\mathbf{x}) = \sum_{i=1}^n |\mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)}|^2$

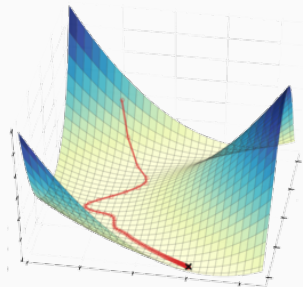
$$f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{y}\|_2^2$$

where  $\mathbf{A}$  is a matrix with  $\mathbf{a}^{(i)}$  as its  $i^{\text{th}}$  row and  $\mathbf{y}$  is a vector with  $y^{(i)}$  as its  $i^{\text{th}}$  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  $\mathbf{x}$ . E.g.  $\|\mathbf{x}\|_2 \leq c$ .

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



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

Partial derivative:

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

Directional derivative:

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

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:

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

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

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



Linear least-squares regression:

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

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

What is the time complexity to implement a function oracle for  $f(\mathbf{x})$ ?

## Linear least-squares regression:

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

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

$$\frac{\partial f}{\partial x_j} = \sum_{i=1}^n 2 \left( \mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)} \right) \cdot a_j^{(i)} = 2\boldsymbol{\alpha}^{(j)T} (\mathbf{Ax} - \mathbf{y})$$

where  $\boldsymbol{\alpha}^{(j)}$  is the  $j^{\text{th}}$  column of  $\mathbf{A}$ .

$$\nabla f(\mathbf{x}) = 2\mathbf{A}^T (\mathbf{Ax} - \mathbf{y})$$

What is the time complexity of a gradient oracle for  $\nabla f(\mathbf{x})$ ?

**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  $\mathbf{v}$ ?

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

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

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

So:

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

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

## Prototype algorithm:

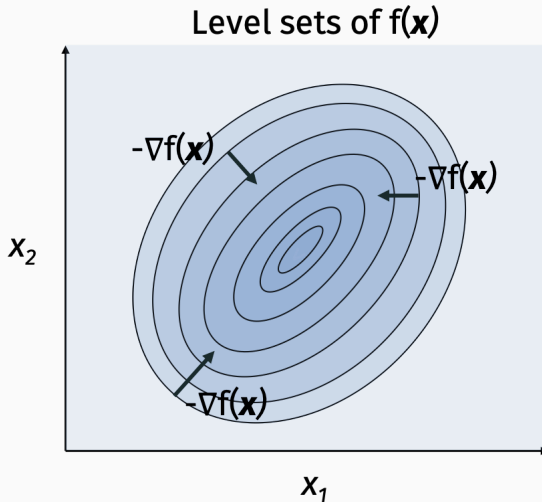
- Choose starting point  $\mathbf{x}^{(0)}$ .
- For  $i = 0, \dots, T$ :
  - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$
- Return  $\mathbf{x}^{(T)}$ .

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

1 dimensional example:

# GRADIENT DESCENT INTUITION

2 dimensional example:



## KEY RESULTS

**For a convex function  $f(\mathbf{x})$ :** For sufficiently small  $\eta$  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.

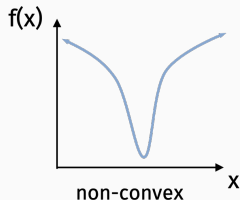
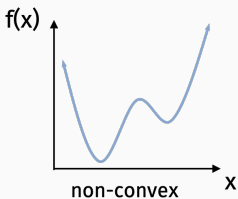
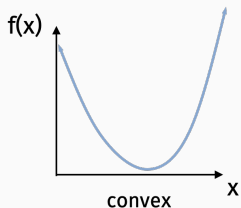
**For a non-convex function  $f(\mathbf{x})$ :** For sufficiently small  $\eta$  and a sufficiently large number of iterations  $T$ , gradient descent will converge to a **near stationary point**:

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

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.

## APPROACH FOR THIS UNIT

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

- Bounding iteration complexity requires placing some assumptions on  $f(\mathbf{x})$ .
- Stronger assumptions lead to better bounds on the convergence.

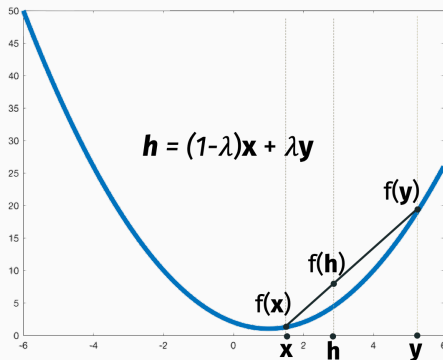
Understanding these assumptions can help us design faster variants of gradient descent (there are many!).

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

## Definition (Convex)

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

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



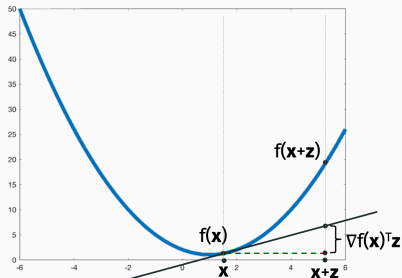
## Definition (Convex)

A function  $f$  is convex if and only if for any  $\mathbf{x}, \mathbf{y}$ :

$$f(\mathbf{x} + \mathbf{z}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{z}$$

Equivalently:

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



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/lec16-supplement.pdf>

Assume:

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

Gradient descent:

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

### Claim (GD Convergence Bound)

*If we run GD for  $T \geq \frac{R^2 G^2}{\epsilon^2}$  iterations then  $f(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*) + \epsilon$ .*

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

### Claim (GD Convergence Bound)

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

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 average solution value is low after  $T = \frac{R^2 G^2}{\epsilon^2}$  iterations. I.e. that:

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

Of course the best solution found,  $\hat{\mathbf{x}}$  is only better than the average.



**Claim (GD Convergence Bound)**

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

**Claim 1:** For all  $i = 0, \dots, T$ ,

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

**Claim 1(a):** For all  $i = 0, \dots, T$ ,

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

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

## Claim (GD Convergence Bound)

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

**Claim 1(a):** For all  $i = 0, \dots, T$ ,

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

## Claim (GD Convergence Bound)

If  $T \geq \frac{R^2 G^2}{\epsilon^2}$  and  $\eta = \frac{R}{G\sqrt{T}}$ , then  $f(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*) + \epsilon$ .

**Claim 1:** For all  $i = 0, \dots, T$ ,

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

Telescoping sum:

$$\begin{aligned} \sum_{i=0}^{T-1} [f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*)] &\leq \frac{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2} \\ &+ \frac{\|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(2)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2} \\ &+ \frac{\|\mathbf{x}^{(2)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(3)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2} \\ &\vdots \\ &+ \frac{\|\mathbf{x}^{(T-1)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2} \end{aligned}$$

## Claim (GD Convergence Bound)

If  $T \geq \frac{R^2 G^2}{\epsilon^2}$  and  $\eta = \frac{R}{G\sqrt{T}}$ , then  $f(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*) + \epsilon$ .

Telescoping sum:

$$\sum_{i=0}^{T-1} [f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*)] \leq \frac{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{T\eta G^2}{2}$$
$$\frac{1}{T} \sum_{i=0}^{T-1} [f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*)] \leq \frac{R^2}{2T\eta} + \frac{\eta G^2}{2}$$

## Claim (GD Convergence Bound)

If  $T \geq \frac{R^2 G^2}{\epsilon^2}$  and  $\eta = \frac{R}{G\sqrt{T}}$ , then  $f(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*) + \epsilon$ .

Final step:

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

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

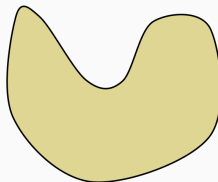
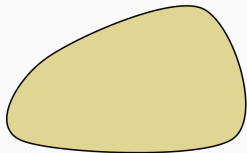
$$f(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*) + \epsilon.$$

## CONSTRAINED CONVEX OPTIMIZATION

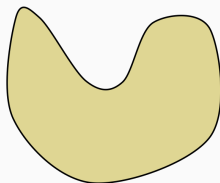
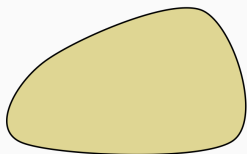
Typical goal: Solve a convex minimization problem with additional convex constraints.

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

where  $\mathcal{S}$  is a **convex set**.



Which of these is convex?



## Definition (Convex set)

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

$$(1 - \lambda)\mathbf{x} + \lambda\mathbf{y} \in \mathcal{S}.$$

### Examples:

- **Norm constraint:** minimize  $\|\mathbf{Ax} - \mathbf{b}\|_2$  subject to  $\|\mathbf{x}\|_2 \leq \lambda$ . Used e.g. for regularization, finding a sparse solution, etc.
- **Positivity constraint:** minimize  $f(\mathbf{x})$  subject to  $\mathbf{x} \geq 0$ . Used e.g. in finding an optimal allocation for a portfolio into different assets.
- **Linear constraint:** minimize  $\mathbf{c}^T \mathbf{x}$  subject to  $\mathbf{Ax} \leq \mathbf{b}$ . Linear program used in training support vector machines, industrial optimization, subroutine in integer programming, etc.



### Gradient descent:

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

Even if we start with  $\mathbf{x}^{(0)} \in \mathcal{S}$ , 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  $\mathcal{S}$  by **projecting** onto the set.

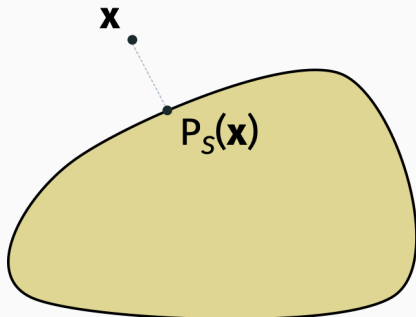
Given a function  $f$  to minimize and a convex constraint set  $\mathcal{S}$ , assume we have:

- **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

- How would you implement  $P_S$  for  $S = \{\mathbf{y} : \|\mathbf{y}\|_2 \leq 1\}$ .
- How would you implement  $P_S$  for  $S = \{\mathbf{y} : \mathbf{y} = \mathbf{Qz}\}$ .



## PROJECTED GRADIENT DESCENT

Given function  $f(\mathbf{x})$  and set  $\mathcal{S}$ , such that  $\|\nabla f(\mathbf{x})\|_2 \leq G$  for all  $\mathbf{x} \in \mathcal{S}$  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  $i = 0, \dots, T$ :
  - $\mathbf{z} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$
  - $\mathbf{x}^{(i+1)} = P_{\mathcal{S}}(\mathbf{z})$
- Return  $\hat{\mathbf{x}} = \arg \min_i f(\mathbf{x}^{(i)})$ .

**Claim (PGD Convergence Bound)**

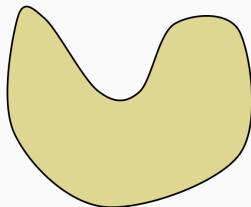
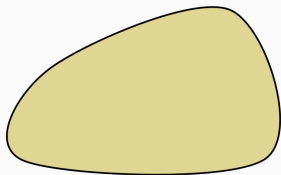
*If  $f, \mathcal{S}$  are convex and  $T \geq \frac{R^2 G^2}{\epsilon^2}$ , then  $f(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*) + \epsilon$ .*

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

### Claim (Contraction Property of Convex Projection)

If  $\mathcal{S}$  is convex, then for any  $\mathbf{y} \in \mathcal{S}$ ,

$$\|\mathbf{y} - P_{\mathcal{S}}(\mathbf{x})\|_2 \leq \|\mathbf{y} - \mathbf{x}\|_2.$$



## Claim (PGD Convergence Bound)

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

**Claim 1:** For all  $i = 0, \dots, T$ ,

$$\begin{aligned} f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) &\leq \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{z} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2} \\ &\leq \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2} \end{aligned}$$

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,  $\mathcal{S}$  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

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

Can our convergence bound be tightened for certain functions? Can it 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  $R$ .
- **Next class:** Take advantage of additional problem structure (e.g. repetition in features and data points in ML problems).



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

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

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

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

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

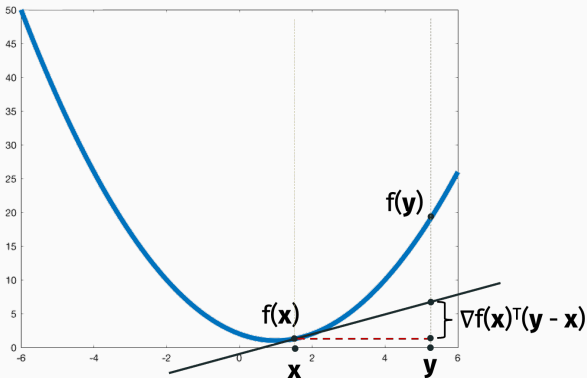
$$[f(\mathbf{y}) - f(\mathbf{x})] - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$$

# SMOOTHNESS

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

So now we have an upper and lower bound.

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



### 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, 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  $f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \epsilon$ .

Compare this to  $T = O\left(\frac{G^2 R^2}{\epsilon^2}\right)$  without a smoothness assumption.

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

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:

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

### 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:

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

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

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:

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

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

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

### Definition (Stationary point)

For a differentiable function  $f$ , a stationary point is any  $\mathbf{x}$  with:

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

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



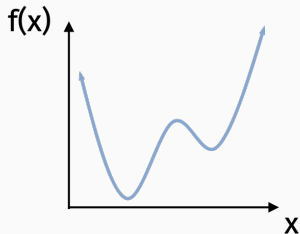
## CONVERGENCE TO STATIONARY POINT

### Theorem (Convergence to Stationary Point)

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

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

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



### Theorem (Convergence to Stationary Point)

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

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

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

$$\begin{aligned} \sum_{t=0}^{T-1} \frac{1}{2\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 &\leq f(\mathbf{x}^{(0)}) - f(\mathbf{x}^{(T)}) \\ \frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 &\leq \frac{2\beta}{T} (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*)) \\ \min_t \|\nabla f(\mathbf{x}^{(t)})\|_2^2 &\leq \frac{2\beta}{T} (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*)) \end{aligned}$$

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?

### Definition ( $\alpha$ -strongly convex)

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

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

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

When  $f$  is convex, we always have that  $f''(x) \geq 0$ , so larger values of  $\alpha$  correspond to a “stronger” condition.

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

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

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

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

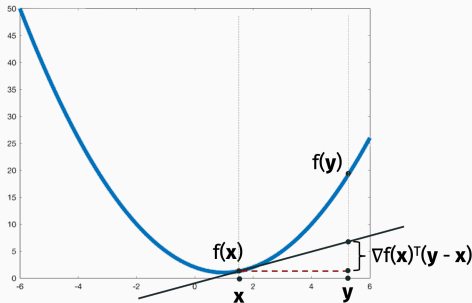
**Corollary:** If  $T = O\left(\frac{G^2}{\alpha\epsilon}\right)$  we have  $f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \leq \epsilon$



# CONVERGENCE GUARANTEE

We could also have that  $f$  is both  $\beta$ -smooth and  $\alpha$ -strongly convex.

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



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

**Theorem (GD for  $\beta$ -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:

$$\|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2 \leq e^{-T \frac{\alpha}{\beta}} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2$$

$\kappa = \frac{\beta}{\alpha}$  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

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

we have:

$$\begin{aligned} \|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2^2 &\leq \frac{2}{\alpha} [f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*)] \\ \|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2 &\geq \frac{2}{\beta} [f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*)] \end{aligned}$$

## 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:

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

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

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \epsilon [f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*)]$$

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

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

Only depend on  $\log(1/\epsilon)$  instead of on  $1/\epsilon$  or  $1/\epsilon^2$ !

Convexity:

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

$\alpha$ -strong-convexity and  $\beta$ -smoothness:

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

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	$O\left(\frac{G^2}{\alpha \epsilon}\right)$	$O\left(\frac{\beta}{\alpha} \log(1/\epsilon)\right)$

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

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

For vector  $\mathbf{x}, \mathbf{v}$ :

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

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

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

**Example:**  $f(\mathbf{x}) = \sum_{i=1}^n (\mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)})^2 = \|\mathbf{Ax} - \mathbf{y}\|_2^2$

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

$$\frac{\partial^2 f}{\partial x_k \partial x_j} = \sum_{i=1}^n 2 a_k^{(i)} a_j^{(i)}$$

$$\mathbf{H} =$$

## ALTERNATIVE DERIVATION

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



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

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

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

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

where  $\mathbf{I}$  is a  $d \times d$  identity matrix.

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

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

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

Equivalently for any  $\mathbf{z}$ ,

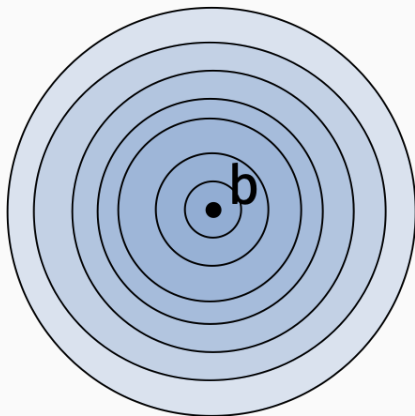
$$\alpha \|\mathbf{z}\|_2^2 \leq \mathbf{z}^T \mathbf{H} \mathbf{z} \leq \beta \|\mathbf{z}\|_2^2.$$

## SIMPLE EXAMPLE

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

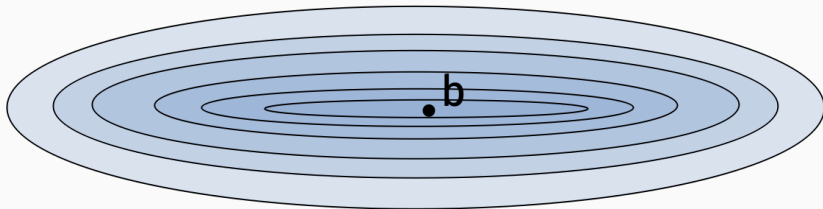
What are  $\alpha, \beta$  for this problem?

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



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  $d_1^2 = \frac{1}{3}, d_2^2 = 2$ .