#### PhD Public Defense

#### On the Complexity of Optimization: Curved Spaces and Benign Landscapes

Candidate: Chris Criscitiello

*Advisor*: Nicolas Boumal



OPTIM, Chair of Continuous Optimization

#### Outline

What is optimization?

Optimization algorithm?

Steepest descent

Best possible algorithm?

#### Outline

What is optimization?

Optimization algorithm?

Steepest descent

Best possible algorithm?

#### Based on:

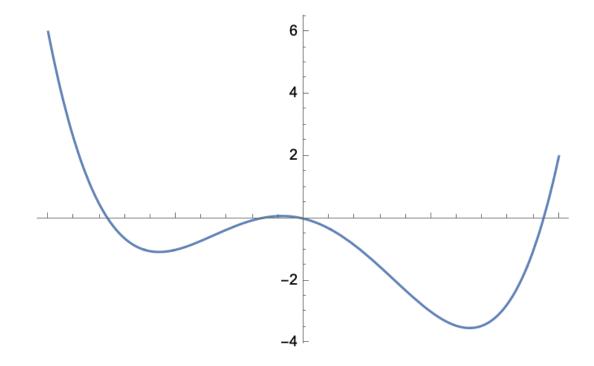
- "Negative curvature obstructs acceleration for strongly geodesically convex optimization" C & **Boumal** COLT'22
- "Curvature and Complexity: Better lower bounds for geodesically convex optimization" C & **Boumal** COLT'23
- "Synchronization on circles and spheres with nonlinear interactions" C, Rebjock, McRae, Boumal under review
- "The sensor network localization problem has benign landscape under mild rank relaxation" C, **Rebjock, McRae, Boumal** under review

Find minimum of function f over  $M = \{\text{set of possibilities}\}$   $\min_{x \in M} f(x)$ 

Find minimum of function f over  $M = \{\text{set of possibilities}\}$ 

$$\min_{x \in M} f(x)$$

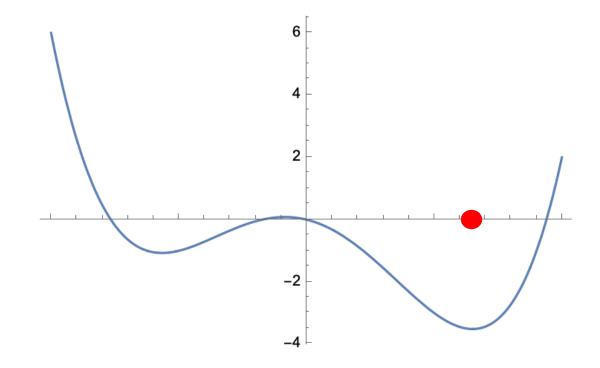
$$f(x) = x^4 - 3x^2 - x$$
  
 
$$M = \{\text{numbers}\}$$



Find minimum of function f over  $M = \{\text{set of possibilities}\}$ 

$$\min_{x \in M} f(x)$$

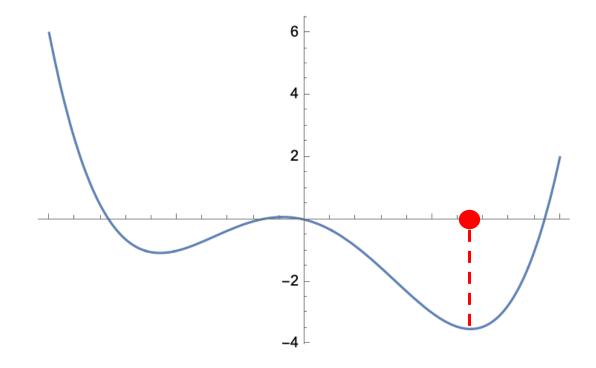
$$f(x) = x^4 - 3x^2 - x$$
  
 
$$M = \{\text{numbers}\}$$



Find minimum of function f over  $M = \{\text{set of possibilities}\}$ 

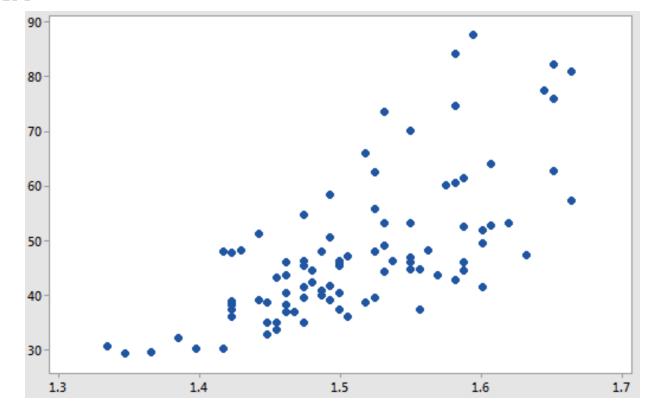
$$\min_{x \in M} f(x)$$

$$f(x) = x^4 - 3x^2 - x$$
  
 
$$M = \{\text{numbers}\}$$



Find minimum of function f over  $M = \{\text{set of possibilities}\}$ 

$$\min_{x \in M} f(x)$$



Find minimum of function f over  $M = \{\text{set of possibilities}\}$ 

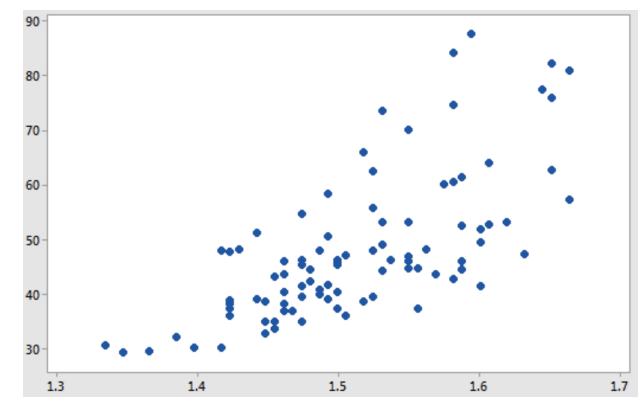
$$\min_{x \in M} f(x)$$

#### **Examples**:

$$M = \{all lines\}$$

f(line) = distance between points and line

$$f(\beta) = \|X\beta - y\|^2$$



Find minimum of function f over  $M = \{\text{set of possibilities}\}$ 

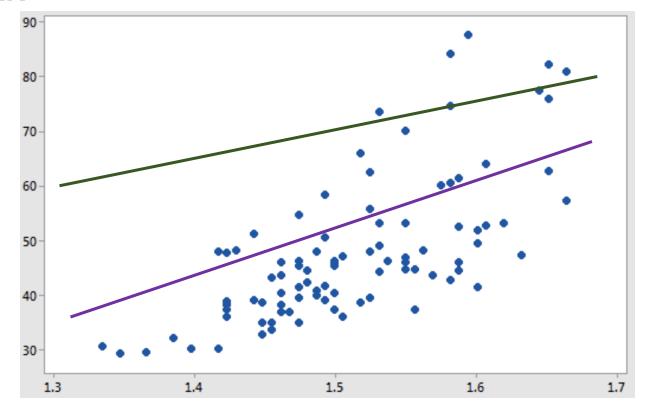
$$\min_{x \in M} f(x)$$

#### **Examples**:

 $M = \{all lines\}$ 

f(line) = distance between points and line

$$f(\beta) = \|X\beta - y\|^2$$



Find minimum of function f over  $M = \{\text{set of possibilities}\}$ 

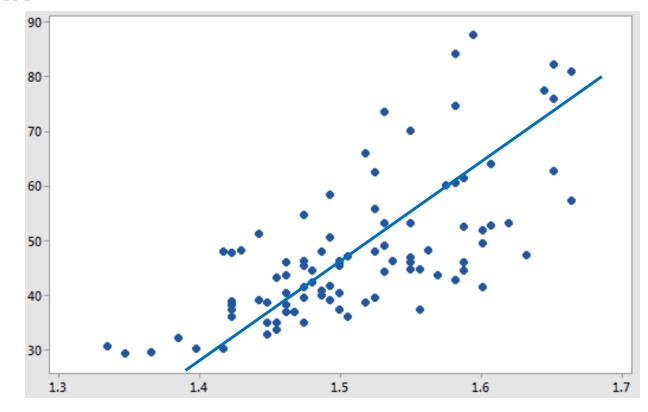
$$\min_{x \in M} f(x)$$

#### **Examples**:

$$M = \{all lines\}$$

f(line) = distance between points and line

$$f(\beta) = \|X\beta - y\|^2$$



Find minimum of function f over  $M = \{\text{set of possibilities}\}$ 

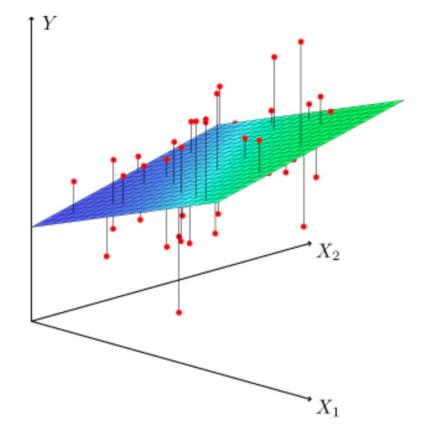
$$\min_{x \in M} f(x)$$

#### **Examples**:

Higher dimensional variants

Harder, can't visualize!

Easily encounter millions of features



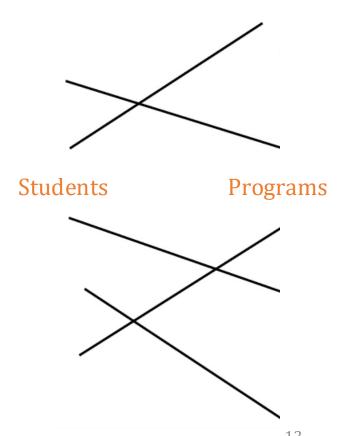
Find minimum of function f over  $M = \{\text{set of possibilities}\}$ 

$$\min_{x \in M} f(x)$$

#### **Examples**:

How to find *best* matching of medical students to residency programs?

 $M = \{all possible matchings\}$ 



Algorithm = method to solve an optim problem

Algorithm = method to solve an optim problem

Imagine super-complicated function *f* stored on a computer/server

Algorithm = method to solve an optim problem

Imagine super-complicated function f stored on a computer/server

What can you do to minimize f?

Algorithm = method to solve an optim problem

Imagine super-complicated function f stored on a computer/server

What can you do to minimize f?

• You can evaluate f on an input x, and get f(x)



Algorithm = method to solve an optim problem

Imagine super-complicated function f stored on a computer/server

What can you do to minimize *f*?

- You can evaluate f on an input x, and get f(x)
- You can evaluate f on an input x, and get "gradient"  $\nabla f(x)$

$$x \longrightarrow f(x), \nabla f(x)$$

Algorithm = method to solve an optim problem

Imagine super-complicated function f stored on a computer/server

What can you do to minimize f?

- You can evaluate f on an input x, and get f(x)
- You can evaluate f on an input x, and get "gradient"  $\nabla f(x)$

query 
$$x \longrightarrow f(x), \nabla f(x)$$

$$x_0 \longrightarrow f(x_0), \nabla f(x_0)$$

$$x_0 \longrightarrow f(x_0), \nabla f(x_0)$$

Algo 
$$\longrightarrow x_1$$



$$x_0 \longrightarrow f(x_0), \nabla f(x_0)$$

$$x_1 \longrightarrow f(x_1), \nabla f(x_1)$$

$$Algo \longrightarrow x_2 \longrightarrow f(x_2), \nabla f(x_2)$$

**Algorithm** = a method to choose queries  $x_0, x_1, x_2, ...$ 

$$x_0 \longrightarrow f(x_0), \nabla f(x_0)$$
 $x_1 \longrightarrow f(x_1), \nabla f(x_1)$ 
 $x_2 \longrightarrow f(x_2), \nabla f(x_2)$ 
 $x_T \text{ so that } f(x_t) \text{ is small}$ 

24

gradient  $\nabla f(x)$  = direction of steepest descent

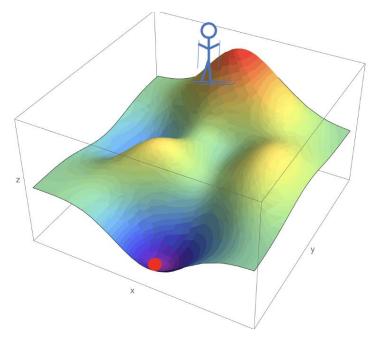
gradient  $\nabla f(x)$  = direction of steepest descent

Imagine you are skiing in a blizzard, or hiking in a dense forest

gradient  $\nabla f(x)$  = direction of steepest descent

Imagine you are skiing in a blizzard, or hiking in a dense forest

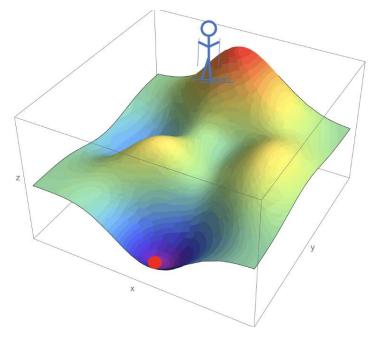
You want to minimize your elevation as quickly as possible



gradient  $\nabla f(x)$  = direction of steepest descent

Imagine you are skiing in a blizzard, or hiking in a dense forest

You want to minimize your elevation as quickly as possible



gradient  $\nabla f(x)$  = direction of steepest descent

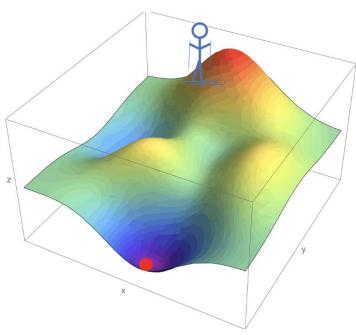
Imagine you are skiing in a blizzard, or hiking in a dense forest

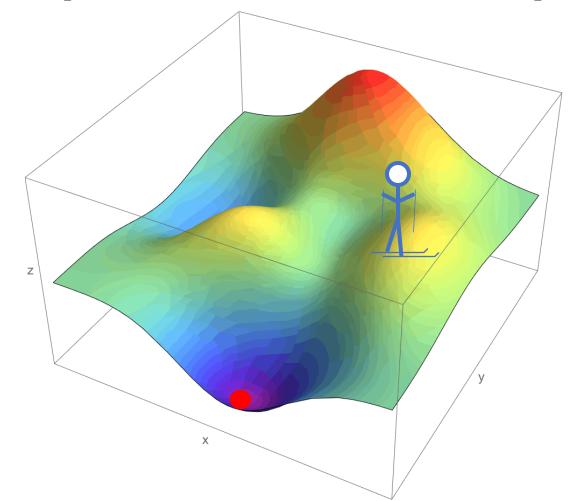
You want to minimize your elevation as quickly as possible

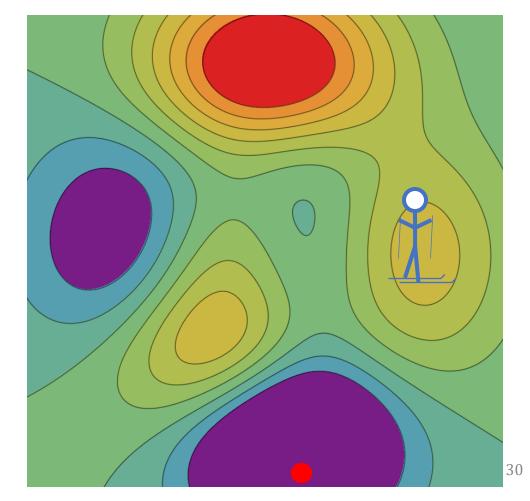
Simple idea: follow direction of steepest descent!

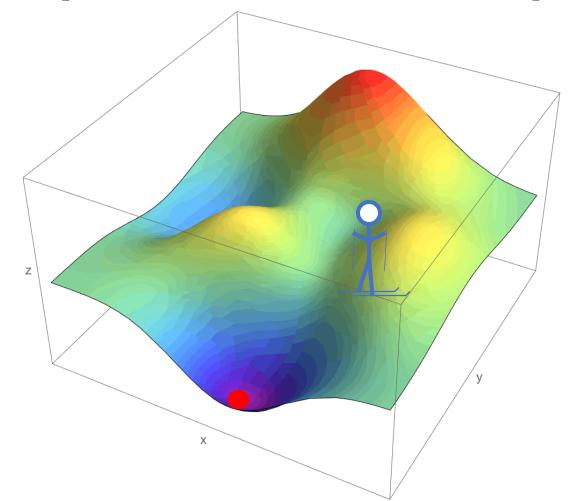
Algorithm = **Steepest Descent** 

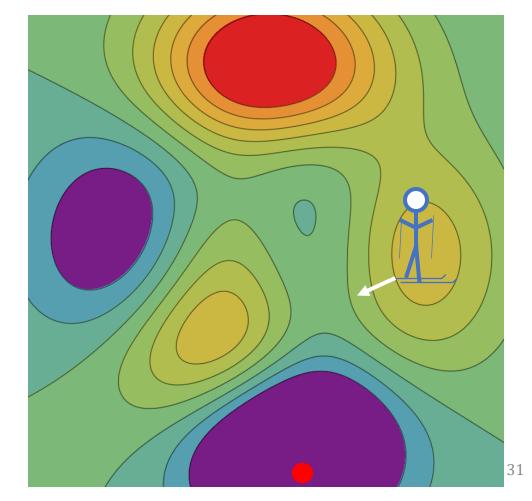
$$x_{k+1} = x_k - \eta \nabla f(x_k)$$

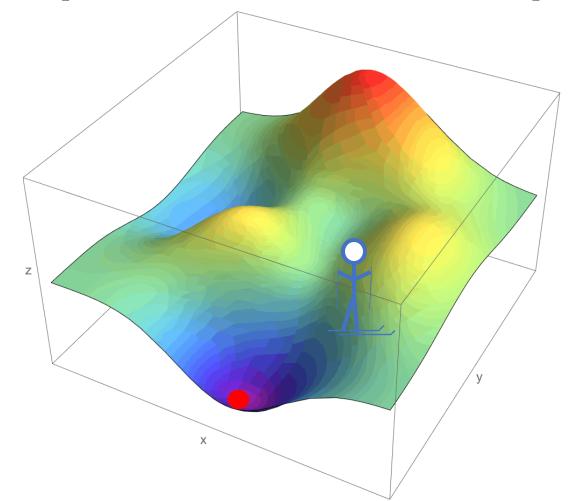


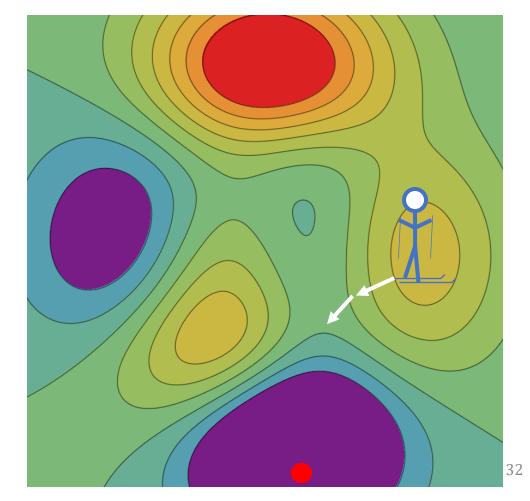


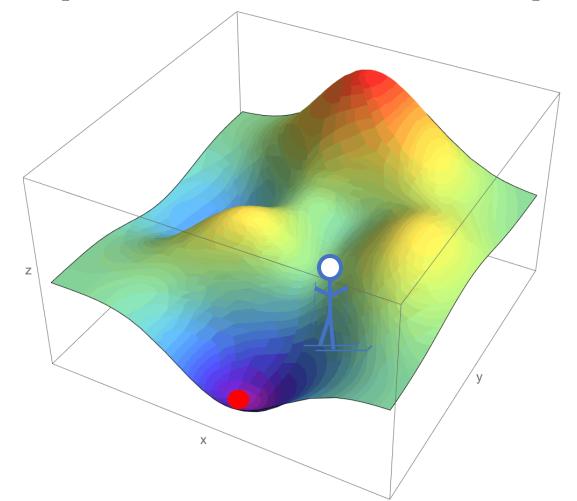


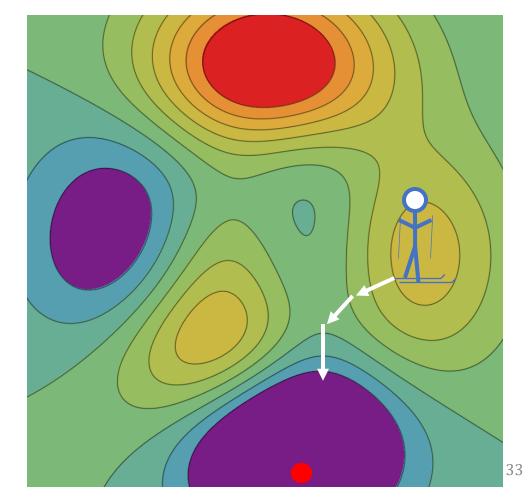


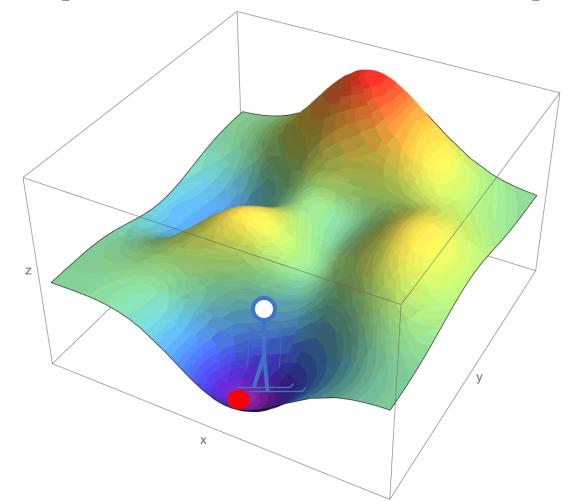


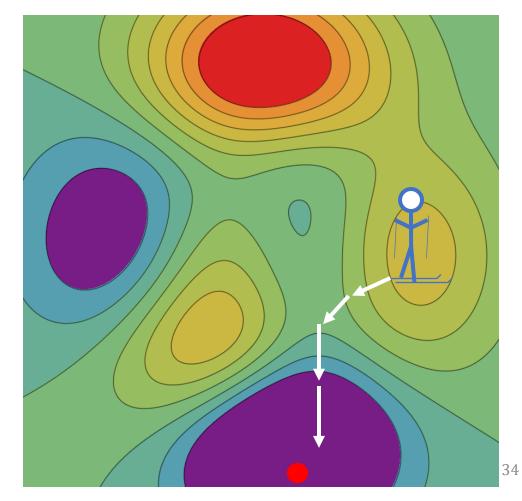


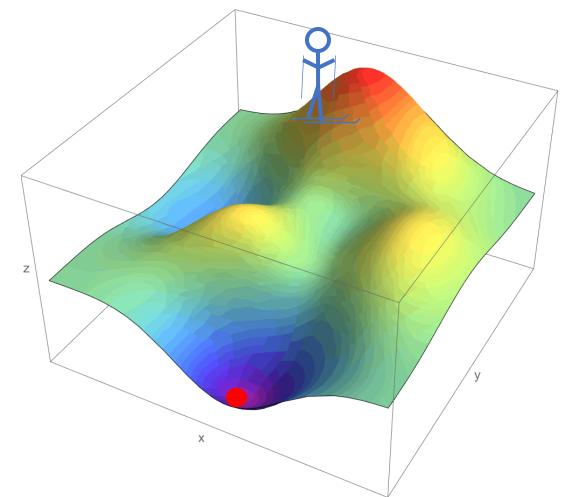


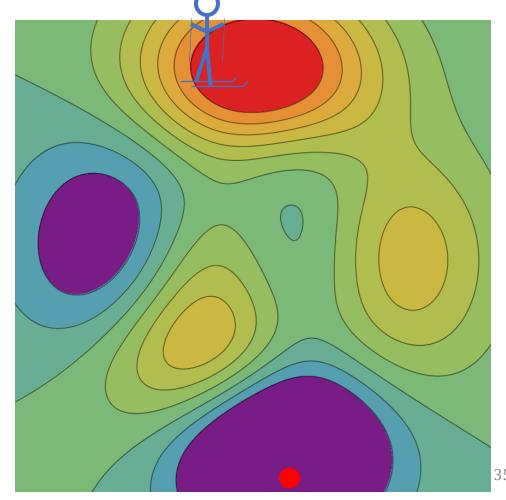


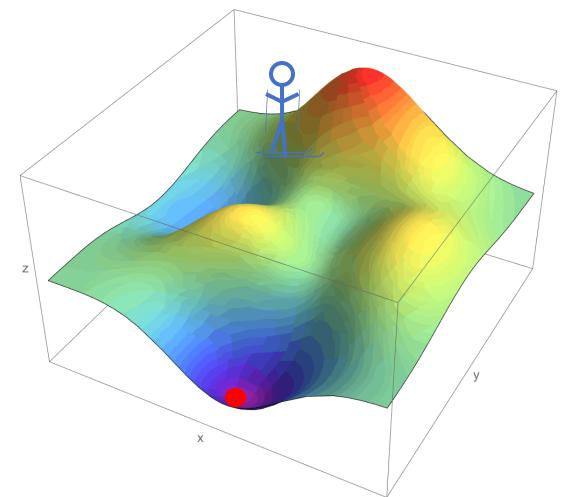


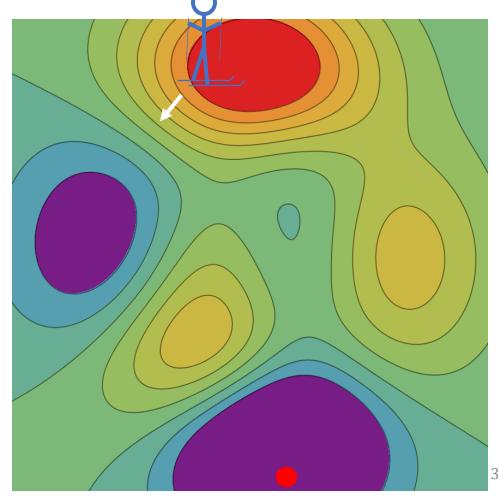




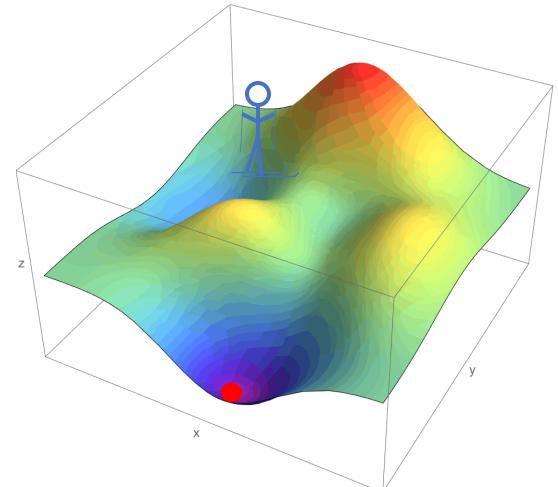


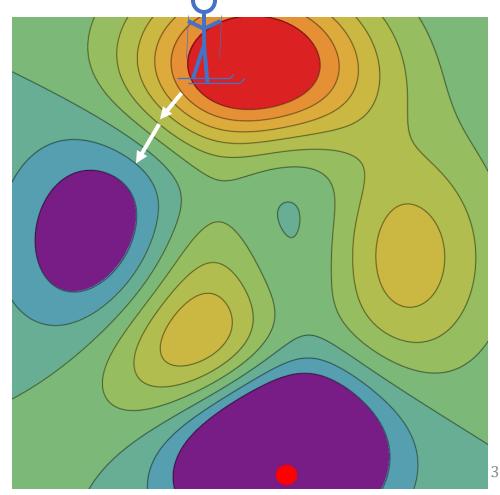




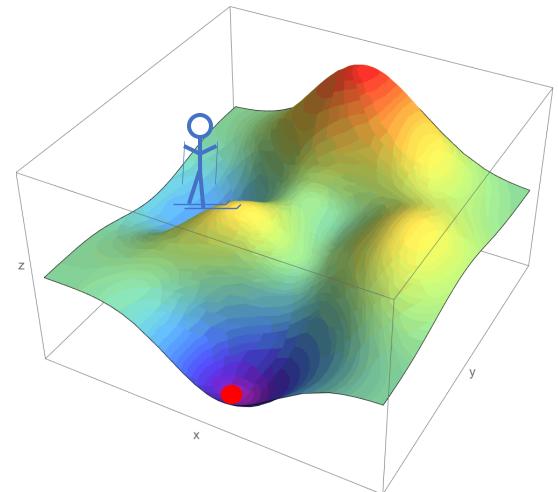


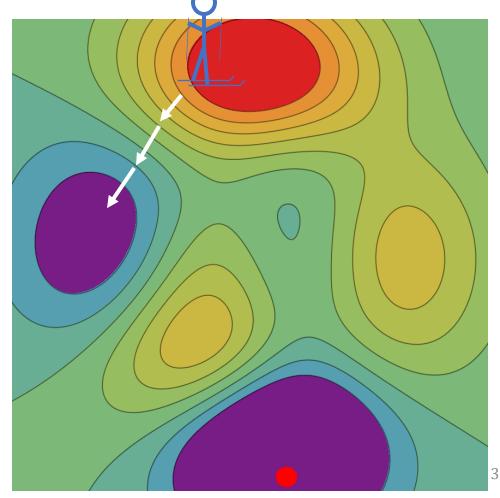
Simple idea: follow direction of steepest descent!



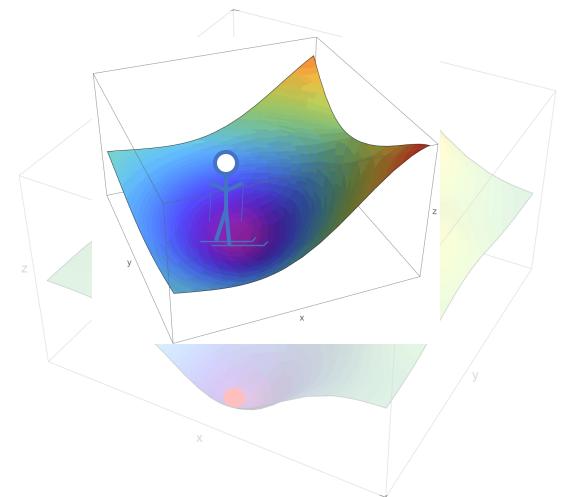


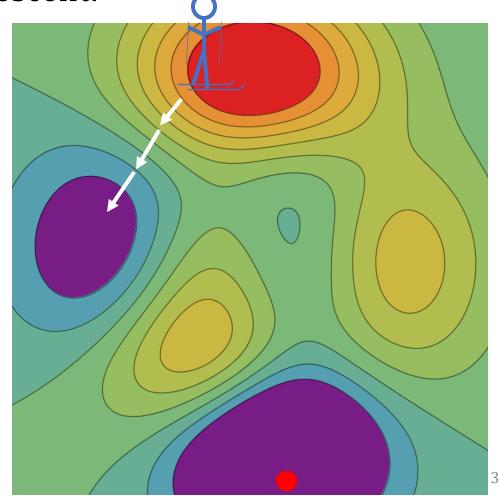
Simple idea: follow direction of steepest descent!

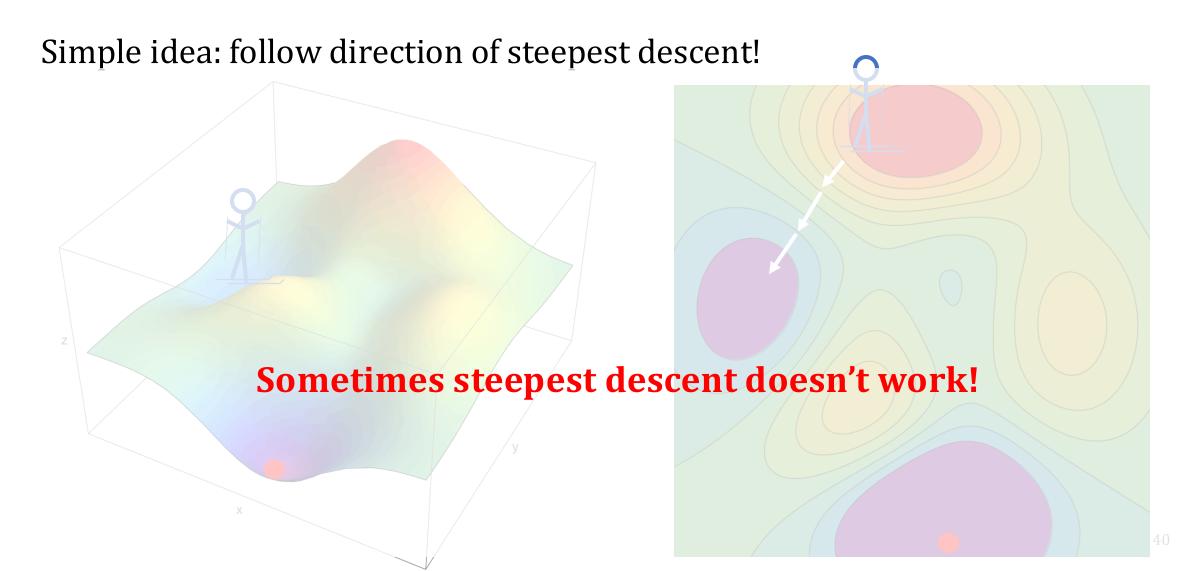




Simple idea: follow direction of steepest descent!







**Q**: For which optim problems does steepest descent work?

**Q**: For which optim problems does steepest descent work?

→ see Part II of my PhD thesis!

**Q**: For which optim problems does steepest descent work?

→ see Part II of my PhD thesis!

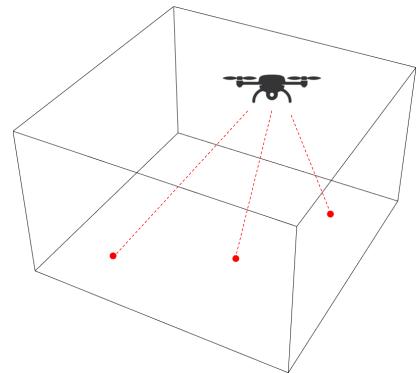
Simple application from thesis: Locating a drone from distance sensors

**Q**: For which optim problems does steepest descent work?

→ see Part II of my PhD thesis!

Simple application from thesis: Locating a drone from distance sensors

**Given**: distances between **drone** and **sensors** 



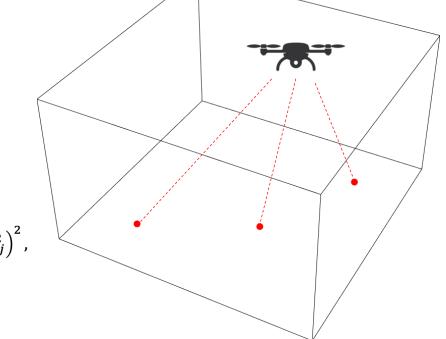
**Q**: For which optim problems does steepest descent work?

→ see Part II of my PhD thesis!

Simple application from thesis: Locating a drone from distance sensors

Given: distances between drone and sensors

**Goal**: Find location of drone, i.e., x, y, z-coordinates



$$\min \sum_{ij \in E} (\|z_i - z_j\|^2 - d_{ij}^2)^2,$$
over  $z_1, z_2, ..., z_n \in \mathbb{R}^k$ 

**Q**: For which optim problems does steepest descent work?

→ see Part II of my PhD thesis!

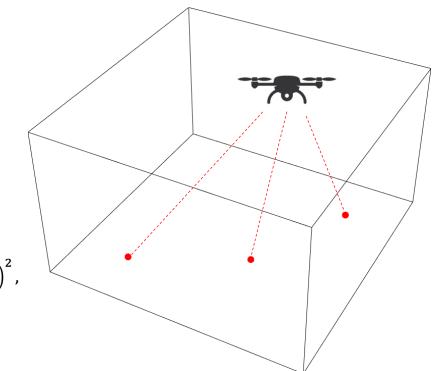
Simple application from thesis: Locating a drone from distance sensors

Given: distances between drone and sensors

**Goal**: Find location of drone, i.e., x, y, z-coordinates

Thm: Steepest descent always works! (under assumptions)

$$\min \sum_{ij \in E} (\|z_i - z_j\|^2 - d_{ij}^2)^2,$$
over  $z_1, z_2, ..., z_n \in \mathbb{R}^k$ 



**Q**: For which optim problems does steepest descent work?

→ see Part II of my PhD thesis!

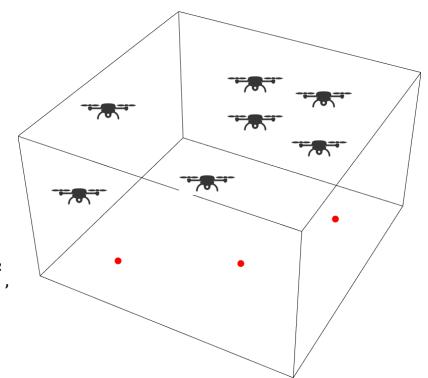
Simple application from thesis: Locating a drone from distance sensors

Given: distances between drone and sensors

**Goal**: Find location of drone, i.e., x, y, z-coordinates

Thm: Steepest descent always works! (under assumptions)

$$\min \sum_{ij \in E} (\|z_i - z_j\|^2 - d_{ij}^2)^2,$$
over  $z_1, z_2, ..., z_n \in \mathbb{R}^k$ 



Other algorithms than steepest descent?

$$x_0 \longrightarrow f(x_0), \nabla f(x_0)$$
Algo  $\longrightarrow x_1 \longrightarrow f(x_1), \nabla f(x_1)$ 

$$\longrightarrow f(x_2), \nabla f(x_2)$$
Algo  $\longrightarrow x_T$  so that  $f(x_t)$  is small

Other algorithms than steepest descent?

Yes! SGD, Nesterov acceleration, Trust regions, ...



$$x_0 \longrightarrow f(x_0), \nabla f(x_0)$$
Algo  $\longrightarrow x_1 \longrightarrow f(x_1), \nabla f(x_1)$ 

$$\longrightarrow f(x_2), \nabla f(x_2)$$

Algo  $\longrightarrow x_T$  so that  $f(x_t)$  is small

Other algorithms than steepest descent?

Yes! SGD, Nesterov acceleration, Trust regions, ...



Which algo is **best**? [Smallest time T to finish. Eg, assume 1 query per second]

Other algorithms than steepest descent?

Yes! SGD, Nesterov acceleration, Trust regions, ...



Which algo is best? [Smallest time T to finish. Eg, assume 1 query per second]

• Space of algos is huge (infinite), can't test all of them on your problem!

Other algorithms than steepest descent?

Yes! SGD, Nesterov acceleration, Trust regions, ...



Which algo is **best**? [Smallest time T to finish. Eg, assume 1 query per second]

- Space of algos is huge (infinite), can't test all of them on your problem!
- That's where math enters *prove* there is no better algorithm

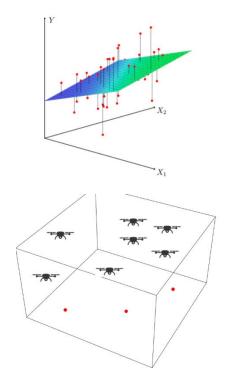
```
Consider a class of optimization problems \{f_1, f_2, ...\} \min_{x \in M} f_i(x)
```

Consider a **class** of optimization problems  $\{f_1, f_2, ...\}$ 

$$\min_{x \in M} f_i(x)$$

**Examples**: Linear regression with different sets of data,

Drone localization with different distance measurements



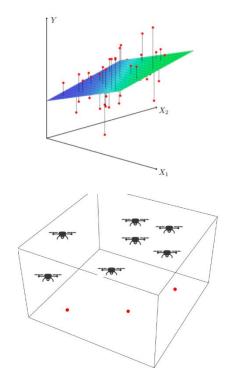
Consider a **class** of optimization problems  $\{f_1, f_2, ...\}$ 

$$\min_{x \in M} f_i(x)$$

**Examples**: Linear regression with different sets of data,

Drone localization with different distance measurements

Amanda, Chris, Joe each have an algo

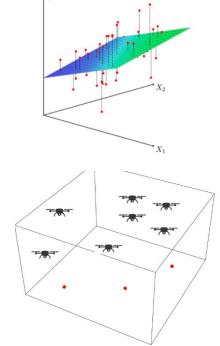


Consider a **class** of optimization problems  $\{f_1, f_2, ...\}$ 

$$\min_{x \in M} f_i(x)$$

**Examples**: Linear regression with different sets of data,

Drone localization with different distance measurements



Amanda, Chris, Joe each have an algo

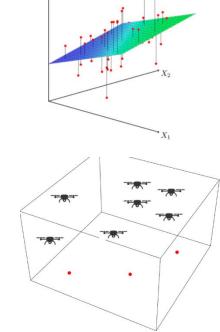
Amanda: Mine minimizes every problem in the class in at most 200 seconds.

Consider a **class** of optimization problems  $\{f_1, f_2, ...\}$ 

$$\min_{x \in M} f_i(x)$$

**Examples**: Linear regression with different sets of data,

Drone localization with different distance measurements



Amanda, Chris, Joe each have an algo

Amanda: Mine minimizes every problem in the class in at most 200 seconds.

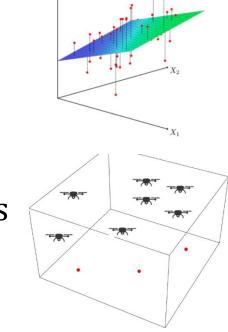
Chris: " 100 seconds.

Consider a **class** of optimization problems  $\{f_1, f_2, ...\}$ 

$$\min_{x \in M} f_i(x)$$

**Examples**: Linear regression with different sets of data,

Drone localization with different distance measurements



Amanda, Chris, Joe each have an algo

Amanda: Mine minimizes every problem in the class in at most 200 seconds.

Chris: " 100 seconds.

Joe: " 50 seconds.

Amanda, Chris, Joe each have an algo

**Amanda**: Mine minimizes **every** problem in the **class** in at most 200 seconds.

Chris: " 100 seconds.

Joe: " 50 seconds.

Amanda, Chris, Joe each have an algo

**Amanda**: Mine minimizes **every** problem in the **class** in at most 200 seconds.

Chris: " 100 seconds.

Joe: " 50 seconds.

Amanda, Chris, Joe each have an algo

**Amanda**: Mine minimizes every problem in the class in at most 200 seconds.

Chris: " 100 seconds.

Joe: " 50 seconds.

Chris: I can prove that no other algo can do better than 100 seconds. Therefore,

My algo is the best (optimal)

Amanda, Chris, Joe each have an algo

**Amanda**: Mine minimizes **every** problem in the **class** in at most 200 seconds.

Chris: " 100 seconds.

Joe: " 50 seconds.

- My algo is the best (optimal)
- Amanda's is suboptimal

Amanda, Chris, Joe each have an algo

**Amanda**: Mine minimizes every problem in the class in at most 200 seconds.

Chris: " 100 seconds.

Joe: " 50 seconds.

- My algo is the best (optimal)
- Amanda's is suboptimal
- Joe is lying

Amanda, Chris, Joe each have an algo

**Amanda**: Mine minimizes every problem in the class in at most 200 seconds.

Chris: " 100 seconds.

Joe: " 50 seconds.

- My algo is the best (optimal)
- Amanda's is suboptimal
- Joe is lying

Consider a **class** of optimization problems  $\{f_1, f_2, ...\}$   $\min_{x \in M} f_i(x)$ 

Consider a **class** of optimization problems  $\{f_1, f_2, ...\}$   $\min_{x \in M} f_i(x)$ 

Chris: I can prove that no other algo can do better than 100 seconds.

What this means: For every algo A, I mathematically construct a problem f in the class so that running A on f takes at least 100 seconds to finish.

Consider a **class** of optimization problems  $\{f_1, f_2, ...\}$   $\min_{x \in M} f_i(x)$ 

Chris: I can prove that no other algo can do better than 100 seconds.

What this means: For every algo A, I mathematically construct a problem f in the class so that running A on f takes at least 100 seconds to finish.

→ Part I of my thesis, looks at a specific **class** (geodesically convex functions), and proves there is a fundamental limit for that class.

→ Part I of my thesis, looks at a specific **class** (geodesically convex functions), and proves there is a fundamental limit for that class.

→ Part I of my thesis, looks at a specific **class** (geodesically convex functions), and proves there is a fundamental limit for that class.

#### Some consequences:

• Steepest descent is optimal for that class (under some assumptions)

→ Part I of my thesis, looks at a specific **class** (geodesically convex functions), and proves there is a fundamental limit for that class.

#### Some consequences:

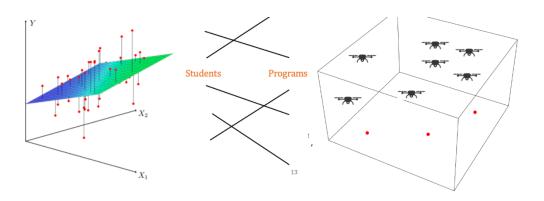
- Steepest descent is optimal for that class (under some assumptions)
- Surprising because expectation in community was that steepest descent is not optimal

Steepest Descent **Optimal!** 

Nesterov Accelerated Descent?
Not Possible!

## Some takeaways

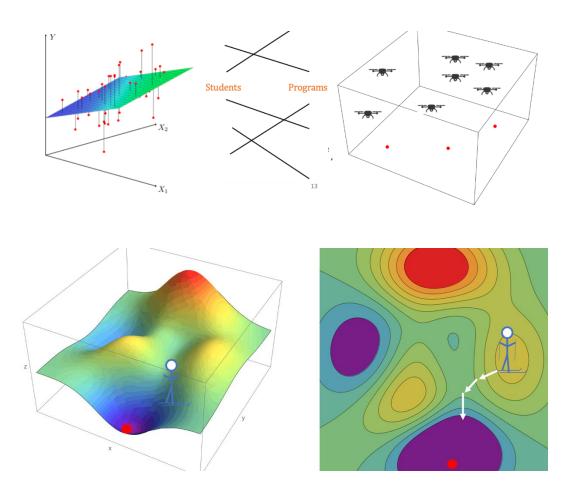
Optimization everywhere



# Some takeaways

Optimization everywhere

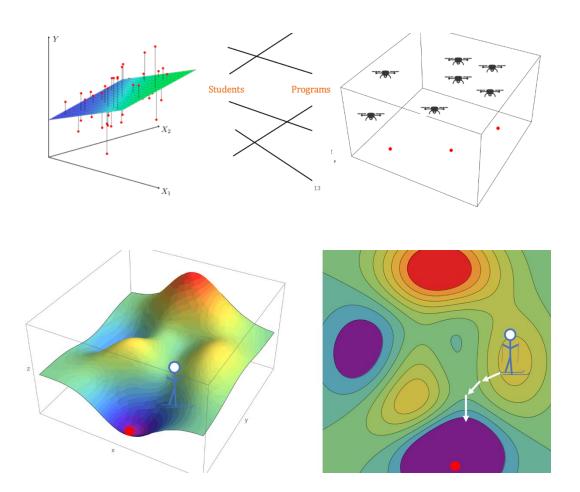
Steepest descent



## Some takeaways

Optimization everywhere

Steepest descent



Math in optim: best algorithm

Steepest Descent **Optimal!** 

Nesterov Accelerated Descent?

Not Possible!