

6.006

Lecture 20

11/29/16

Topic

Continuous Optimization I

Today

Continuous Optimization

Gradient Descent Method

- Convexity & Optimality

Sofar in class:

Combinatorial optimization:

- optimize over large finite (discrete) collection of objects
e.g. APSP-paths
sorting-orderings of lists
- brute force is possible, just not efficient

Today:

Continuous optimization:

- optimize over continuous (infinite) set of solutions
- usually: some subset of \mathbb{R}^n

Isn't everything on a computer finite?

still useful:

- finite, but huge (solve PDE's)
- scientific computing
- solve systems of linear equations
- robotics, control
- machine learning - deep learning
- graph algorithms
- ...

A basic, very general, problem:

Unconstrained Minimization

Given real-valued function

OBJECTIVE
FUNCTION $\Rightarrow f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$\text{find } \min_{x \in \mathbb{R}^n} f(x) \quad \Leftarrow \quad \begin{array}{l} x \text{ is vector in } \mathbb{R}^n \\ x = (x_1, \dots, x_n) \end{array}$$

Why do we call it "unconstrained"?

we let x be any vector in \mathbb{R}^n

OPTIMAL
SOLUTION $x^* \equiv \operatorname{argmin}_{x \in \mathbb{R}^n} f(x)$

Comments:

- This is pretty general!

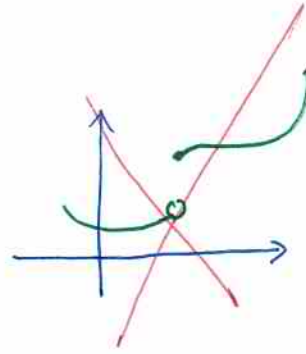
e.g. • to solve $\max_{x \in \mathbb{R}^n} f(x)$, just solve $\min_{x \in \mathbb{R}^n} -f(x)$

- can encode any continuous optimization task by making f sufficiently complex

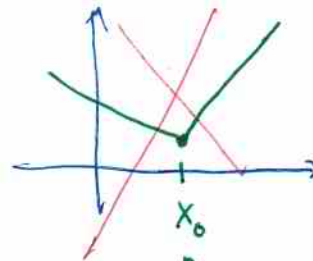
e.g. if want constrained optimization, make $f(x) = \infty$ when x doesn't satisfy constraints

Assumptions on f :

- f is continuous



- f is (infinitely) differentiable



not differentiable at x_0

assumptions can be relaxed

How to Solve? luckily we have a powerful approach coming to the rescue ...

GRADIENT DESCENT

Key Idea: Greedy local search

- start with some initial point x^0
- in each iteration, move a bit in direction that reduces value of f the most
- guarantees $f(x^{i+1}) \leq f(x^i)$

← "local"

← "greedy"

Warmup:

Gradient descent in 1-D ($n=1$):

Algorithm idea:

$i \leftarrow 0$

$x^0 \leftarrow$ arbitrary point

While $f'(x^i) \neq 0$ do

if $f'(x^i) > 0$: move left

if $f'(x^i) < 0$: move right

$i \leftarrow i+1$

if $f'(x^i) \approx 0$ then x is local extremum

local
optimality
condition

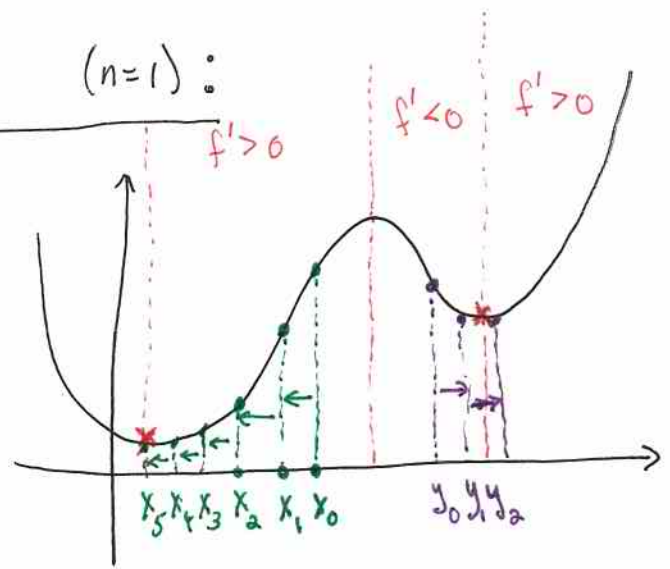
no more
local
improvement
possible

minimum (😊)

or

maximum (😞)

↑
only happens
if unlucky



Important note:

different starting points x^0 give different local maxima

Questions:

• how far to move?

• What does " \approx " mean?

(5)

A closer look:

Taylor Expansion for $x' = x + \epsilon$

$$f(x') \approx f(x + \epsilon) = f(x) + f'(x) \cdot \epsilon + \underbrace{\frac{\epsilon^2}{2!} f''(x) + \frac{\epsilon^3}{3!} f'''(x) + \dots}_{\text{error of pretending } f \text{ is linear}}$$

use this to
approximate f
"locally" as
a linear fctn

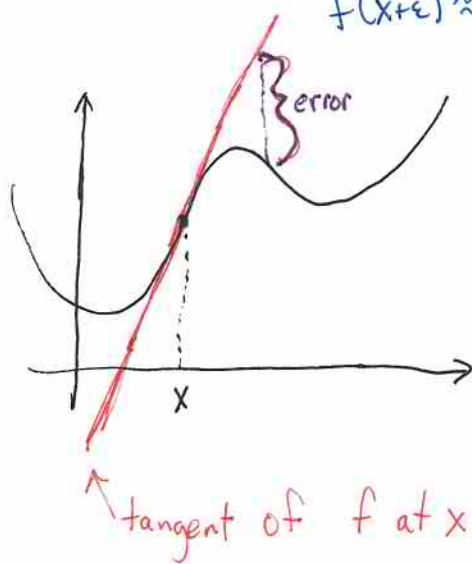
error of pretending
 f is linear
 \uparrow
assuming ϵ small
 $\approx \epsilon^2 f''(x)$



will use small ϵ
+ assume

$$f(x + \epsilon) \approx f(x) + f'(x) \cdot \epsilon$$

so treat f as a
linear function !!



What should the step size be?

- needs to be small enough such that approximation is still "good"
- needs to be large enough such that we make some progress?

Step size/direction

- move in direction of $-f'(x)$
- move by amount $-\eta f'(x)$

stepsize / learning rate
 η depends on f

why should step size depend on f ? later we will see that this works really well for analysis!

Claim good choice of η is $\eta = \frac{1}{f''(x)}$

why?

good thing to know otherwise you aren't necessarily getting better

\Rightarrow make sure the improvement is bigger than the error!

improvement:

$$\begin{aligned} f(x^{i+1}) &= f(x^i - \eta f'(x^i)) \\ &\approx f(x^i) - \underbrace{(\eta f'(x^i)) f'(x^i)}_{\varepsilon} + \text{Error term} \\ &\approx f(x^i) - \eta f'(x^i)^2 \end{aligned}$$

improvement / progress

error: bound by twice next term in Taylor expansion

$$\cancel{2} \cdot \frac{\varepsilon^2}{2!} f''(x) = \underbrace{\eta^2 f'(x)^2 f''(x)}_{\text{approx error}}$$

more exact calculations $\Rightarrow \eta \approx \frac{1}{2f''(x)}$

So progress in each step:
 $\approx \eta (f'(x))^2 \approx \frac{f'(x)^2}{f''(x)}$

So need

$$\cancel{\eta} f'(x)^2 \approx \cancel{\eta^2} f'(x)^2 f''(x) \Rightarrow \eta \approx 1/f''(x)$$

Commercial break:

Learning (Deep and Shallow...)

Given a bunch of points
 $(x_1, y_1) (x_2, y_2) \dots$

e.g. $x_i = \text{picture}$ $y_i = 1 \text{ if it is a cat}$
 0 o.w.

$x_i = \text{vector}$ $y_i = A x$
 \nearrow
 unknown matrix

$x_i = \text{patient information}$ $y_i = \text{probability patient dies}$

Goal we want to "learn" a function, which given x_i computes a good prediction to y_i

In general, learning is not really possible
 can't learn a random function!

However, we could modify our goal. Here is an example:

Goal' find a linear function $\left\{ \begin{array}{l} \text{could be any} \\ \text{class of functions} \\ \text{e.g. deep learning} \\ \text{networks,} \\ \text{human brains,} \\ \text{simple network} \\ \text{of Boolean logic} \end{array} \right.$
 which best approximates (x_i, y_i)
 relationship under L_2 -error

more specifically!

for hypothesis linear function f_w

the loss of f_w on input is denoted $\text{Loss}(f_w(x_i), y_i)$

- could use $\text{loss}(f_w(x_i), y_i) = \|f_w(x_i) - y_i\|^2 \leftarrow \begin{array}{l} \text{"L}_2\text{-error"} \\ \text{"quadratic loss"} \\ \text{lots of other names} \end{array}$

$\left\{ \begin{array}{l} \text{could be any other} \\ \text{error} \\ \text{(L}_1, \text{information theoretic} \\ \text{loss functions)} \end{array} \right.$

Total loss of hypothesis on input:

$$\text{Total loss} = \sum_i \text{loss}(f_w(x_i), y_i)$$

But w is a set of parameters
(determining a linear function)

how do we find w that
minimizes total loss: for given $(x_1, y_1), (x_2, y_2), \dots$

GRADIENT DESCENT

Deep learning

f_w is not linear

but still, w is a bunch of parameters
to a network

back to work (now we motivated general case...)

9

General Case ($n > 1$)

Same philosophy: f is locally (multivariate) linear function

What is different? x is now a vector

"grad f " $\nabla f(x)$ is analogue of derivative = $\begin{bmatrix} \frac{\partial f}{\partial x_1}(x) \\ \frac{\partial f}{\partial x_2}(x) \\ \vdots \\ \frac{\partial f}{\partial x_n}(x) \end{bmatrix}$
it is also a vector!

$$\text{Still have: } f(x+\varepsilon) = f(x) + \underbrace{\nabla f(x)^T \varepsilon}_{\substack{\text{linear} \\ \text{approximation}}} + \underbrace{\frac{1}{2} \varepsilon^T \nabla^2 f(x) \varepsilon}_{\text{error}} + \dots$$

$\nwarrow \nearrow$
vector in \mathbb{R}^n

Notation:

$$\begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix}^T = [a_1, \dots, a_n]$$

"transpose"

Inner Product

$$a^T b = \sum_{i=1}^n a_i b_i = \langle a, b \rangle$$

Analogue of 2nd derivative:

"Hessian" how locally non linear the fctn is

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1}(x) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(x) & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(x) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(x) & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(x) & \dots & \dots & \frac{\partial^2 f}{\partial x_n \partial x_n}(x) \end{bmatrix}$$



Gradient Descent Algorithm

• start with x^0

• For $k = 0 \dots$ do

$$x^{k+1} = x^k - \eta \nabla f(x^k)$$

Coordinate-wise do this by

$$x_i^{k+1} = x_i^k - \eta \frac{\partial f}{\partial x_i}(x_i^k)$$

local optimality condition:

$$\nabla f(x) = 0$$

x is either
 local min or "saddle point"
 not local min

Example: Linear Function

$$f(x) = C^T x + b$$

↑↑↑ vectors!

Note: $\nabla f(x) = (c_1 \dots c_n) = C$
 $\nabla^2 f(x) = (0 \dots 0)$

↑↑ so no bound
 on step size η

what if we move in direction d ?

$$\begin{aligned} f(x+d) &= C^T(x+d) + b = C^T x + C^T d + b \\ &= f(x) + C^T d \\ &= f(x) + \|C\| \cdot \|d\| \cdot \cos \theta \end{aligned}$$

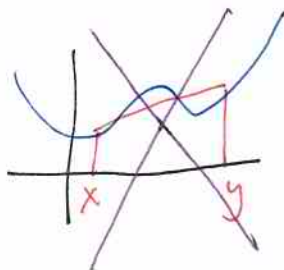
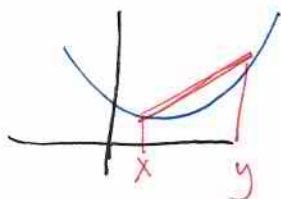


So, best choice of d : let $d = -C = -\nabla f(x)$ $\theta = 180^\circ$
 $\cos \theta = -1$

best bang for the buck: $f(x+d) = f(x) - \|C\|^2$

When does the local min equal the global min?

Important case: Convexity



f is "convex" iff

$$\forall \begin{array}{l} 0 \leq \lambda \leq 1 \\ x, y \in \mathbb{R}^n \end{array}$$

$$f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y)$$

every chord is above the fctn.



$$\forall x, y \quad f(x+y) \geq f(x) + \nabla f(x)^T (y-x)$$

Key fact convex \Rightarrow every local min is a global min