CS-GY 6763: Lecture 10 Dimension Dependent Optimization, Linear Programming

NYU Tandon School of Engineering, Prof. Christopher Musco

First Order Optimization: Given a convex function f and a convex set S,

Goal: Find $\hat{\mathbf{x}} \in S$ such that $f(\hat{\mathbf{x}}) \leq (\min_{\mathbf{x} \in S} f(\mathbf{x}) + \epsilon)$.

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} .

Gradient descent requires $O\left(\frac{R^2G^2}{\epsilon^2}\right)$ calls to each oracle to solve the problem.

We were only able to improve the ϵ dependence by making stronger assumptions on f (strong convexity, smoothness).

typically convex

Alternatively, we can get much better bounds if we are willing to depend on the problem dimension. I.e. on d if $f(\mathbf{x})$ is a function mapping d-dimensional vectors to scalars.

We already know how to do this for a few special functions:



3

Let $f(\mathbf{x})$ be bounded between ([-B, B]) on $S_{\mathbf{x}}$ **Theorem (Dimension Dependent Convex Optimization)** There is an algorithm (the Center-of-Gravity Method) which finds $\hat{\mathbf{x}}$ satisfying $f(\hat{\mathbf{x}}) \leq \min_{\mathbf{x} \in S} f(\mathbf{x}) + \epsilon$ using $O(d \log(B/\epsilon))$ calls to a function and gradient oracle for convex f.

Caveat: Assumes we have some representation of S, not just a projection oracle. We will discuss this more later.

Note: For an unconstrained problem with known starting radius *R*, can take *S* to be the ball of radius *R* around $\mathbf{x}^{(0)}$. If $\|\nabla f(\mathbf{x})\|_2 \leq G$, we always have B = O(RG).

Natural "cutting plane" method. Developed simultaneous on opposite sides of iron curtain.



Not used in practice (we will discuss why) but the basic idea underlies many popular algorithms, including the famous

CENTER OF GRAVITY METHOD

A few basic ingredients:



1. The center-of-gravity of a convex set \mathcal{S} is defined as:

$$c = \underbrace{\int_{x \in \mathcal{S}} x \, dx}_{\operatorname{vol}(\mathcal{S})} = \frac{\int_{x \in \mathcal{S}} x \, dx}{\int_{x \in \mathcal{S}} 1 \, dx}$$

2. For two convex sets \mathcal{A} and \mathcal{B} , $\mathcal{A} \cap \mathcal{B}$ is convex. Proof by picture:



CENTER OF GRAVITY METHOD



Natural "cutting plane" method.

- $\cdot \ \mathcal{S}_1 = \mathcal{S}$
- For t = 1, ..., T:
 - $\mathbf{c}_t = \text{center of gravity of } \mathcal{S}_t.$
 - Compute $\nabla f(\mathbf{c}_t)$.
 - $\mathcal{H} = \{ \mathbf{x} | \langle \nabla f(\mathbf{c}_t), \mathbf{x} \mathbf{c}_t \rangle \leq 0 \}.$
 - $\cdot \ \mathcal{S}_{t+1} = \mathcal{S}_t \cap H$
- Return $\hat{\mathbf{x}} = \arg\min_t f(\mathbf{c}_t)$



Natural "cutting plane" method.

- $\cdot \ \mathcal{S}_1 = \mathcal{S}$
- For t = 1, ..., T:
 - $\mathbf{c}_t = \text{center of gravity of } \mathcal{S}_t$.
 - Compute $\nabla f(\mathbf{c}_t)$.
 - $\mathcal{H} = \{ \mathbf{x} | \langle \nabla f(\mathbf{c}_t), \mathbf{x} \mathbf{c}_t \rangle \leq 0 \}.$
 - $\cdot \ \mathcal{S}_{t+1} = \mathcal{S}_t \cap H$
- Return $\hat{\mathbf{x}} = \arg\min_t f(\mathbf{c}_t)$



Intuitively, why does it make sense to search in $S_t \cap \mathcal{H}$ where:



Intuitively, why does it make sense to search in $\mathcal{S}_t \cap \mathcal{H}$ where:

$$\mathcal{H} = \{ \mathbf{x} | \langle \nabla f(\mathbf{c}_t), \mathbf{x} - \mathbf{c}_t \rangle \leq 0 \}?$$

$$\begin{aligned} & \{(\mathcal{G}) - f(c_{\star}) \not\equiv (\nabla f(c_{\star}), \mathcal{G} - \mathcal{G}) \\ & \text{By convexity,} \quad f(\mathbf{y}) - \mathcal{G}(\mathbf{x}) \not\equiv f'(\mathbf{x})(\mathcal{G} - \mathbf{x}) \\ & \underbrace{f(\mathbf{y}) \geq f(\mathbf{c}_{t}) + \langle \nabla f(\mathbf{c}_{t}), \mathbf{y} - \mathbf{c}_{t} \rangle}_{& If \ \mathbf{y} \notin \{\mathcal{S}_{t} \cap \mathcal{H}\} \text{ then } \overset{\mathcal{G}}{\mathcal{G}}_{t, 1} \rightarrow \mathbf{c}_{t} \rangle} \\ & \text{If } \mathbf{y} \notin \{\mathcal{S}_{t} \cap \mathcal{H}\} \text{ then } \overset{\mathcal{G}}{\mathcal{G}}_{t, 1} \rightarrow \mathbf{c}_{t} \rangle \\ & \underbrace{\langle \nabla f(\mathbf{c}_{t}), \mathbf{y} - \mathbf{c}_{t} \rangle}_{& f(\mathbf{c}_{t}), \mathbf{y} - \mathbf{c}_{t} \rangle} = \underbrace{\langle \mathbf{x} : \mathcal{L} \forall f(\mathbf{c}), \mathbf{x} - \mathbf{c}_{t} \rangle \leq \mathbf{0}^{2}}_{& Ievel \ sets \ of \ f(\mathbf{x})} \end{aligned}$$

CONVERGENCE THEOREM

Theorem (Center-of-Gravity Convergence) Let f be a convex function with values in [-B, B]. Let $\hat{\mathbf{x}}$ be the output of the center-of-gravity, method run for T iterations. Then: $f(\hat{\mathbf{x}}) - \underline{f(\mathbf{x}^*)} \le 2B\left(\left(1 - \frac{1}{e}\right)\right)^{T/d} \le 2Be^{-1}$ -T/31 = 10 + (6/2B) изч xeb(h,x^{o1}) If we set $T = 3d \log(2B/\epsilon)$, then $f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \leq \epsilon$. 0 T/3 d= log (26/2) T= 32 log (28/2)



KEY GEOMETRIC TOOL

Theorem (Grünbaum's Theorem)

For any convex set S with center-of-gravity \mathbf{c} , and any halfspace $\mathcal{Z} = \{\mathbf{x} | \langle \mathbf{a}, \mathbf{x} - \mathbf{c} \rangle \le 0\}$ then: $\underbrace{\operatorname{vol}(S \cap \mathcal{Z})}_{\operatorname{vol}(S)} \ge \frac{1}{e} \approx .368$

Want to argue that, at every step of the algorithm, we "cut off" a large portion of the convex set we are searching over.

Theorem (Grünbaum's Theorem)

For any convex set S with center-of-gravity c, and any halfspace $\mathcal{Z} = \{x | \langle a, x - c \rangle \le 0\}$ then:

$$\frac{\mathsf{vol}(\mathcal{S} \cap \mathcal{Z})}{\mathsf{vol}(\mathcal{S})} \ge \frac{1}{e} \approx .368$$

Let \mathcal{Z} be the compliment of \mathcal{H} from the algorithm. Then we cut off at least a 1/*e* fraction of the convex body on every iteration.

Corollary: After t steps,
$$\operatorname{vol}(\mathcal{S}_t) \leq \left(1 - \frac{1}{e}\right)^t \operatorname{vol}(\mathcal{S})$$
.
 $\operatorname{vol}(\mathcal{S}_t) \leq \left(1 - \frac{1}{e}\right) \operatorname{vo}(\mathcal{S}_{t-1})$

15

Let
$$\delta$$
 be any small error parameter. $x = \lambda + S^{\delta}$
Let $S^{\delta} = \{(1-\delta)x^{*} + \delta x \mid \text{for } x \in S\}$
 $S^{\delta} = (1 \cdot S) x^{*} + \delta S$
 $y_{0} \setminus (S^{\delta}) = S^{\delta} y_{0} \mid (S)$

Claim: Every point **y** in S^{δ} has good function value.









Theorem (Center-of-Gravity Convergence)

Let f be a convex function with values in [-B, B]. Let $\hat{\mathbf{x}}$ be the output of the center-of-gravity method run for T iterations. Then:

$$f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \leq 2B \left(1 - \frac{1}{e}\right)^{T/d} \leq 2Be^{-T/3d}.$$

If we set
$$T = O(d \log(B/\epsilon))$$
, then $f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \le \epsilon$.

In terms of <u>gradient-oracle</u> complexity, this is essentially optimal. So why isn't the algorithm used?

- In general computing the centroid is hard. #P-hard even when when S is an intersection of half-spaces (a polytope).
- Even if the problem isn't hard for your starting convex body $\mathcal{S},$ it likely will become hard for $\mathcal{S} \cap \mathcal{H}_1 \cap \mathcal{H}_2 \cap \mathcal{H}_3 \dots$
- So while the <u>oracle complexity</u> of dimension-dependent optimization was settled in the 60, basic questions remained regarding <u>computational complexity</u>.

We will see how to resolve this issue with an elegant cutting plane methods called the *(Ellipsoid Method)* that was introduced <u>by</u> Naum Shor in <u>1977</u>.





To talk about <u>runtime efficiency</u> we need to be more concrete about how our (convex) constraint set is even specified.

Seperation Oracle: For a convex set $\underline{\mathcal{K}} \subset \mathbb{R}^d$, a seperation oracle $\underline{S_{\mathcal{K}}}$ is a function that takes in points in $\underline{\mathbb{R}}^d$ and returns:



SEPARATION ORACLE

Example: How would you implement a separation oracle for a (0, x) 7, b polytope { $\mathbf{x} : \mathbf{Ax} \ge \mathbf{b}$ }. 7, 5 Suppor Ay 76. For some row Q: 64; , y) < b; but, Lai, X72 b: for all XES, 24

Instead of directly solving a constrained optimization problem, solve the <u>membership problem</u>. Given a separation oracle $S_{\mathcal{K}}$ for a convex set \mathcal{K} , determine if \mathcal{K} is empty, or otherwise return any point $\mathbf{x} \in \mathcal{K}$.



FROM MEMBERSHIP TO OPTIMIZATION

Original problem
$$\min_{x \in S} f(x)$$
. (Is fore some $X \in S$ such that $f(x) \leq C$.)
How to reduce to determining if a convex set K is empty or not?
Is funce any $X \in 2 \leq 3 \land 2 \leq 5 \leq 5 \leq 5$ (on we have)
constraint set S (on we have)
(on we

Original problem: $\min_{x \in S} f(x)$. How to reduce to determining if a convex set \mathcal{K} is empty or not?



Claim: Given any fixed value c_i can check if $f(\mathbf{x}^*) \le c$ and, if it is, find some \mathbf{x} with $f(\mathbf{x}) \le c$.

Approach: Solve membership problem on $\mathcal{K} = \mathcal{S} \cap \mathcal{C}$ where $\mathcal{C} = \{\mathbf{x} : f(\mathbf{x}) \leq c\}$. \mathcal{C} and \mathcal{S} are convex, so \mathcal{K} is as well.

FROM MEMBERSHIP TO OPTIMIZATION



FROM MEMBERSHIP TO OPTIMIZATION



Claim: Given any fixed value c, can check if $f(\mathbf{x}^*) \leq c$ and, if it is, find some \mathbf{x} with $f(\mathbf{x}) \leq c$.

Final algorithm: Assuming *f* is positive, just run exponential/binary search to find $\tilde{c} \leq f(\mathbf{x}^*) + \epsilon!$



ELLIPSOID METHOD SKETCH

 $f(x) \leq f(x^*) + \epsilon$ $B(c_R,R)$ 01 К B(c"r) K Application to original problem: Lots of details to consider. Assume for simplicity we known $f(x^*)$ and that have no constraint set. Goal is to solve membership problem on $\tilde{C} = \{\mathbf{x} : f(\mathbf{x}) \le f(\mathbf{x}^*) + \epsilon\}$. For a convex function f such that $\|\nabla f(\mathbf{x})\|_2 \le G$, it can be checked that \tilde{C} contains a ball of radius ϵ/G .

n = measure of Iterative method similar to center-of-gravity: 1. Check if center \mathbf{c}_R of $B(\mathbf{c}_R, R)$ is in \mathcal{K} . 2. If it is, we are done. 3. If not, cut search space in half, using separating phyperplane. S (CA) $B(c_R,R)$.K o B(dr,r) 32

Key insight: Before moving on, approximate new search region by something that we can easily compute the centroid of. Specifically an ellipse!



Produce a sequence of ellipses that <u>always contain</u> \mathcal{K} and decrease in volume: $B(\mathbf{c}_R, R) = E_1, E_2, \dots$ Once we get to an ellipse with volume $\leq B(\mathbf{c}_r, r)$, we know that \mathcal{K} must be empty.

ELLIPSE

An ellipse is a convex set of the form: $\{\underline{x} : \|A(\underline{x} - \underline{c})\|_2^2 \le \alpha\}$ for some constant α' and matrix A. The center-of-mass is c.

 $\{x: ||I(x-c)|| < \alpha\} \quad \{x: ||D(x-c)|| < \alpha\} \quad \{x: ||A(x-c)|| < \alpha\}$ **_**C **_**C Q= (+ A + A)-1 $\|X-C\|^{1} \leq Q$ $\|A(\mathbf{x} \cdot \mathbf{c})\|_{\nu}^{\nu} = (\mathbf{x} \cdot \mathbf{c})^{\dagger} A^{\dagger} A(\mathbf{x} \cdot \mathbf{c}) \leq 1$ Often re-parameterized to say that the ellipse is all **x** with $\{x : (x - c)^T Q^{-1} (x - c) \le 1\}$ RELXA

R There is a closed form solution for the equation of the smallest ellipse containing a given half-ellipse. I.e. let \mathbf{E}_i have Ci = center E; parameters $\mathbf{Q}_i, \mathbf{c}_i$ and consider the half-ellipse: $\left(\underline{\mathsf{E}}_i \cap \{ \mathsf{x} : \mathsf{a}_i^T \mathsf{x} \leq \mathsf{a}_i^T \mathsf{c}_i \}. \right)$ Then E_{i+1} is the ellipse with parameters: $\frac{d^2}{d^2-1}\left(\underline{\mathbf{Q}}_i - \frac{2}{d+1}\mathbf{h}\mathbf{h}^T\right) \underbrace{\mathbf{c}_{i+1}}_{\underline{\mathbf{C}}_i} = \underline{\mathbf{c}}_i - \frac{1}{n+1}$ where $\mathbf{h} = \sqrt{a} \mathbf{Q}_i \mathbf{a}_i$ $\cdot \mathbf{a}_i$.

Computing the update takes $O(d^2)$ time.

GEOMETRIC OBSERVATION





Not as good as the $(1 - \frac{1}{e})$ constant-factor volume reduction we got from center-of-gravity, but still very good!

GEOMETRIC OBSERVATION



After O(d) iterations, we reduce the volume by a constant. In total require $O(d^2 \log(R/r))$ iterations to solve the problem.

Complexity for solving $\min_{\mathbf{x}\in\mathcal{S}} f(\mathbf{x})$ is roughly $\tilde{O}(d^4 \log(R/\epsilon))$, hiding logarithmic factors.

Linear programs (LPs) are one of the most basic convex constrained, convex optimization problems:

Let $\mathbf{c} \in \mathbb{R}^d$, $\mathbf{b} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{n \times d}$ be fixed vectors that define the problem, and let \mathbf{x} be our variable parameter.

 $\min f(\mathbf{x}) = \mathbf{c}^{\mathsf{T}} \mathbf{x}$
subject to $\mathbf{A}\mathbf{x} \ge \mathbf{b}$.

Think about $Ax \ge b$ as a union of half-space constraints:

$$\{\mathbf{x} : \mathbf{a}_1^T \mathbf{x} \ge b_1\}$$
$$\{\mathbf{x} : \mathbf{a}_2^T \mathbf{x} \ge b_2\}$$
$$\vdots$$
$$\{\mathbf{x} : \mathbf{a}_n^T \mathbf{x} \ge b_n\}$$



LINEAR PROGRAMMING APPLICATIONS

- Classic optimization applications: industrial resource optimization problems were killer app in the 70s.
- Robust regression: $min_x \|Ax b\|_1$.
- L1 constrained regression: $\min_{\mathbf{x}} \|\mathbf{x}\|_1$ subject to $A\mathbf{x} = \mathbf{b}$. Lots of applications in sparse recovery/compressed sensing.
- $\cdot \ \text{Solve min}_x \, \|Ax b\|_\infty.$
- Polynomial time algorithms for Markov Decision Processes.
- Many combinatorial optimization problems can be solved via <u>LP relaxations</u>.

Theorem (Khachiyan, 1979)

Assume n = d. The ellipsoid method solves any linear program with L-bit integer valued constraints exactly in (n^4) time.

A Soviet Discovery Rocks World of Mathematics

By MALCOLM W. BROWNE

A surprise discovery by an obscure Soviet mathematician has rocked the world of mathematics and computer analysis, and experts have begun exploring its practical applications.

Mathematicians describe the discoverv by L.G. Khachian as a method by which computers can find guaranteed solutions to a class of very difficult problems that have hitherto been tackled on a kind of hit-or-miss basis.

Apart from its profound theoretical interest, the discovery may be applicable sometimes involves so many steps that it

in weather prediction, complicated indus- could take billions of years to compute. trial processes, petroleum refining, the scheduling of workers at large factories, secret codes and many other things.

"I have been deluged with calls from virtually every department of government for an interpretation of the significance of this," a leading expert on computer methods, Dr. George B. Dantzig of Stanford University, said in an interview.

The solution of mathematical problems by computer must be broken down into a series of steps. One class of problem

The Russian discovery offers a way by which the number of steps in a solution can be dramatically reduced. It also offers the mathematician a way of learning quickly whether a problem has a solution or not, without having to complete the entire immense computation that may be required.

According to the American journal Sci-

Continued on Page A20, Column 3

ONLY \$10.00 A MONTH 24 Hr. Phone Answering Service, Totally New Concept" Increable!! 279-3870-ADVT.

Front page of New York Times, November 9, 1979.

Theorem (Karmarkar, 1984)

Assume n = d. The <u>linterior point method</u> solves any linear program with L-bit integer valued constraints in $O(n_{3.5}^{3.5}L)$ time.



Front page of New York Times, November 19, 1984.

Lecture notes are posted on the website (optional reading).



Projected Gradient Descent Optimization Path

Lecture notes are posted on the website (optional reading).



Ideal Interior Point Optimization Path

Both results had a huge impact on the theory of optimization, although at the time neither the ellipsoid method or interior point method were faster than a heuristic known at the Simplex Method.

These days, improved interior point methods often outperform simplex.

Polynomial time linear programming algorithms have also had a huge impact of <u>combinatorial optimization</u>. They are often the work-horse behind approximation algorithms for NP-hard problems. Given a graph G with n nodes and edge set E. Each node is assigned a weight w_1, \ldots, w_n .



Goal: Select subset of nodes with minimum total weight that covers all edges.

Given a graph G with n nodes and edge set E. Each node is assigned a weight w_1, \ldots, w_n .

Formally: Denote if node *i* is selected by assigning variable x_i to 0 or 1. Let $\mathbf{x} = [x_1, \dots, x_n]$.

$$\min_{\mathbf{x}} \sum_{i=1}^{n} x_i w_i \quad \text{subject to} \quad x_i \in \{0, 1\} \text{ for all } i$$
$$x_i + x_j \ge 1 \text{ for all } (i, j) \in E$$

NP-hard to solve exactly. We will use convex optimization give a 2-approximation in polynomial time.

Function to minimize is linear (so convex) but constraint set is not convex. Why?

High level approach:

- <u>Relax</u> to a problem with convex constraints.
- <u>Round</u> optimal solution of convex problem back to original constraint set.



High level approach:

- \cdot <u>Relax</u> to a problem with convex constraints.
- <u>Round</u> optimal solution of convex problem back to original constraint set.



High level approach:

- <u>Relax</u> to a problem with convex constraints.
- <u>Round</u> optimal solution of convex problem back to original constraint set.

Let $\bar{S} \supseteq S$ be the relaxed constraint set. Let $\mathbf{x}^* = \arg \min_{\mathbf{x} \in S} f(\mathbf{x})$ and let $\bar{\mathbf{x}}^* = \arg \min_{\mathbf{x} \in \bar{S}} f(\mathbf{x})$. We always have that:

$$f(\bar{\mathsf{X}}^*) \leq f(\mathsf{X}^*).$$

So typically the goal is to round \bar{x}^* to ${\cal S}$ in such a way that we don't increase the function value too much.

RELAXING VERTEX COVER

Vertex Cover:



$$x_i + x_j \ge 1$$
 for all $(i, j) \in E$

Relaxed Vertex Cover:



The second problem is a linear program! It can be solved in poly(n) time!

Simple rounding procedure: If $\bar{x}_i^* \ge 1/2$, set $x_i = 1$, and set $x_i = 0$ otherwise.



Observation 1: All edges remain covered. I.e., the constraint $x_i + x_j \ge 1$ for all $(i, j) \in E$ is not violated.

Observation 2: Let **x** be the rounded version of $\bar{\mathbf{x}}^*$. We have $f(\mathbf{x}) \leq 2 \cdot f(\bar{\mathbf{x}})$, and thus $f(\mathbf{x}) \leq 2 \cdot f(\mathbf{x}^*)$.

Proof:

So, a polynomial time algorithm for solving LPs immediately yields a 2-approximation algorithm for the NP-hard problem of vertex cover.

- Proven that it is NP-hard to do better than a 1.36 approximation in [Dinur, Safra, 2002].
- Recently improved to $\sqrt{2} \approx$ 1.41 in [Khot, Minzer, Safra 2018], which proved the 2-to-2 games conjecture.
- Widely believed that doing better than 2
 e is NP-hard for any
 e > 0, and this is implied by Subhash Khot's Unique Games Conjecture.

There is a simpler greedy 2-approximation algorithm that doesn't use optimization at all. Try coming up with it on your own!

Next section of course: <u>Spectral methods</u> and <u>numerical linear</u> <u>algebra</u>.

Spectral methods generally refer to methods based on the "spectrum" of a matrix. I.e. on it's eigenvectors/eigenvalues and singular vectors/singular values. We will look at applications in:

- Low-rank approximation and dimensionality reduction.
- Data clustering and related problems.
- Constructing data embeddings (e.g. Word2Vec).