

Optimization Intro

MTA 3007 Notebook

Youthy WANG

- Terminology:

1) Euclidean inner product: $\langle x, y \rangle = x^T y$

2) Euclidean norm of a vector / point $\|x\|_2 = \|x\| = \sqrt{x^T x}$

3) $g(x) = \begin{bmatrix} g_1(x) \\ g_2(x) \\ \vdots \\ g_m(x) \end{bmatrix}$, then $g(x) \leq 0$ means $g_i(x) \leq 0, \forall i=1, 2, \dots, m$

4) Feasible point: a decision satisfying all constraints

Feasible set / region Ω — a set of all feasible pts

5) Optimal solution: a feasible point / decision variable that attains an objective value, as good as other feasible pts.
 may be more than 1

Optimal value: the objective value of any optimal solution.
 only 1

6) Formal definition of extrema:

local minimizer / maximizer: $x^* \in \Omega$ & $\exists \varepsilon > 0$ s.t. $f(x) \geq f(x^*)$ / $f(x) \leq f(x^*)$
 for all $x \in \Omega \cap B_\varepsilon(x^*)$ ($B_\varepsilon(y) := \{x \in \mathbb{R}^n \mid \|x-y\| \leq \varepsilon\}$)

strict local \sim (change \geq, \leq to $>, <$)

global minimizer / maximizer: $x^* \in \Omega$ & \Rightarrow It holds that $f(x) \geq f(x^*)$ / $f(x) \leq f(x^*)$,
 for all $x \in \Omega$
 global minimizer / maximizer = global solution = optimal solution

- Facts & Properties

Optimization & Statistics are two foundations of machine learning.

Modeling — transfer primitive questions into optimization problems.

The feasible set can be empty (called infeasible / degenerated)

The optimal value may be unbounded

Classification

Infeasible

Feasible; finite optimal value; not attainable

Feasible; finite optimal value, attainable

Feasible; unbounded optimal value.

Classifications of Optimization Problems:

{ Unconstrained v.s. Constrained

Linear (LP) v.s. Non-linear (NLP)

Integer/Discrete (IP) v.s. Continuous

• Real example analyses:

? Definition (math def of optimization problem):

$\min/\max f(x)$

$\underset{x}{\text{subject to}} \quad x \in S$

$(\min_x f(x) = -\max_x -f(x))$
same

$\min/\max f(x)$

$\underset{x}{\text{subject to}} \quad g_i(x) \leq 0, \forall i=1,2,\dots,m$

$h_j(x) = 0, \forall j=1,2,\dots,p$

(Unconstrained: $S = \mathbb{R}^n$ or $m=p=0$)

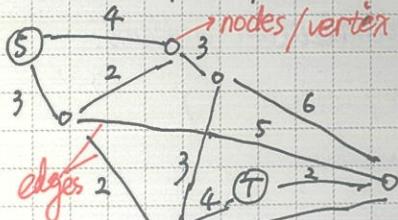
{ Decision ~ decision variables: unknown & needs choosing

Objective ~ objective functions

Constraints ~ constraint functions: equalities & inequalities

3 main factors

eg. 1 (Shortest path problem)



Find shortest path from S to T

$V = \text{set of nodes}$

$E \subseteq V \times V = \text{set of existing edges}$

w_{ij} expresses the distance from i to j

Choice: $x_{ij} = \begin{cases} 1, & (i,j) \text{ edge is chosen} \\ 0, & \text{otherwise} \end{cases}$

optimization model:

$$\underset{\{x_{ij}\}}{\text{minimize}} \sum_{(i,j) \in E} w_{ij} x_{ij}$$

$$\text{subject to } x_{ij} \in \{0, 1\}, \forall (i, j) \in E$$

constraints for a path $\leftarrow \begin{cases} \sum_j x_{sj} = 1 = \sum_j x_{jt} & (\text{for the pinnacle}) \\ \sum_j x_{ij} = \sum_j x_{ji}, \forall i \neq s, t \end{cases}$

eg. 2, for the same diagram in 1, find the smallest set of vertices that touch every edge.

$$\text{optimization model: } x_i = \begin{cases} 1, & \text{node } i \text{ is chosen} \\ 0, & \text{otherwise} \end{cases}$$

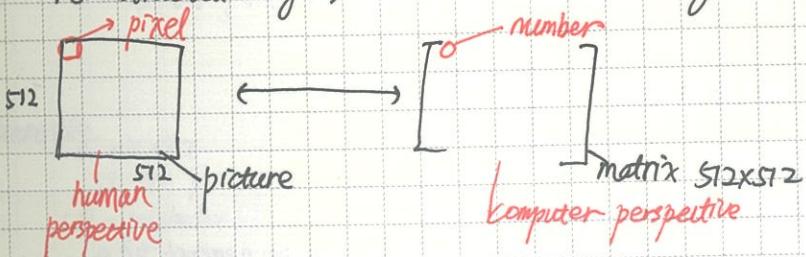
$$\underset{\{x_i\}}{\text{minimize}} \sum_i x_i \rightarrow (\text{size of set})$$

$$\text{subject to } x_i + x_j \geq 1, \text{ for } (i, j) \in E \hookrightarrow \text{the requirement of question}$$

$$x_i \in \{0, 1\}, \forall i \in V$$

eg. 3 (removing scratches in a picture)

Decision - the whole image; Objective - remove scratches and obtain a "nicer" reconstructed image; Constraints - not change undamaged ones.



Lecture 2 Application & Linear Programming

- Machine Learning:

To extract important patterns, understand data & use it to make predictions or decisions.

Classification & regression are two fundamental tasks in machine learning

Support Vector Machine (SVM) is one of the most important tools for classification.

eg. (for classification)

Learning to predict if one applicant should be approved a credit card.

	age	gender	salary	citizenship	years in job		approve
app 1	2.5	1	10	3	1	$\rightarrow x_i$	+1
app 2	2.8	0	8	6	5	(a row(i-th) containing a person's data)	+1
app 3	1.6	0	0	2	0		-1
:	:	:	:	:	:		:
app 5	3	1	8	2	1		-1

data matrix — representing by \mathbf{x}

(label indicates supervision)

Label — y

For all training samples = $\overset{\uparrow}{\text{row}} \text{label} : \overset{\uparrow}{\mathbf{x}_i}$
 (draw as (\mathbf{x}_i, y_i))

Idea: learn a function $f: \mathbb{R}^n \rightarrow \{-1, 1\}$
base on training samples

① Linear classification: $f(\mathbf{x}) := \mathbf{x}^T \mathbf{w} + b$ s.t. $f(\mathbf{x}_i) y_i \geq 1, \forall i$
transfer minimize \mathbf{w}, b $\circlearrowleft 0$ \rightarrow feasibility problem
 subject to $y_i (\mathbf{x}_i^T \mathbf{w} + b) \geq 1, \forall i$

To make the line be the best, (not naive)

SVM → ① find support vectors (above & below) =
transfer (draw parallel line/plane, passing nearest pt.)

② let two line/planes be $\mathbf{x}^T \mathbf{w} + b = 1$ & $\mathbf{x}^T \mathbf{w} + b = -1$
 (by normalize & choosing \mathbf{x}, b)

③ maximize $\frac{2}{\|\mathbf{w}\|}$ (distance)

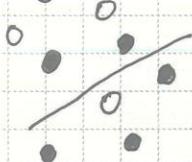
better case
transfer maximize $\frac{2}{\|\mathbf{w}\|}$ margin
 subject to $y_i (\mathbf{x}_i^T \mathbf{w} + b) \geq 1, \forall i$ $\overset{\uparrow}{\text{subject to}}$ $y_i (\mathbf{x}_i^T \mathbf{w} + b) \geq 1, \forall i$
non-linear, constrained, continuous

(use $\|\mathbf{w}\|^2$ because $\mathbf{w} \mapsto \|\mathbf{w}\|$ is not differentiable at 0)

maximum-margin hyperplane

2) Non-linear-separable data

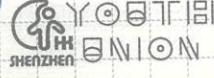
Define Hinge-loss $\max\{0, 1 - y_i (\mathbf{x}_i^T \mathbf{w} + b)\}$



- For points on the true side, 0
- For points on the wrong side, the farther, the larger.

Final Edition

$$\underset{w, b}{\text{minimize}} \quad \frac{\lambda}{2} \|w\|^2 + \sum_{i=1}^m \max\{0, 1 - y_i(x_i^T w + b)\}$$



w, b arbitrary, λ chosen to balance the margin & misclassification (weighted sum)

a special case — when $\lambda = 0$, $t_i = \max\{0, 1 - y_i(x_i^T w + b)\}$
(linear programming)

transfer $\underset{w, b, t}{\text{minimize}} \sum_i t_i$

$$\text{subject to } t_i = (1 - y_i(x_i^T w + b))^+, \forall i$$

can be written as $t_i \geq (1 - y_i(x_i^T w + b))^+$

$$\Leftrightarrow t_i \geq 0, t_i \geq (1 - y_i(x_i^T w + b)), \forall i$$

• Linear Programming & Standard Form

LP / linear optimization problem: where the objective function and all constraint functions are linear in the decision variables (independent from others)

1) Compact Way of expression:

$$\underset{x}{\text{minimize}} \quad c^T x \quad (\text{subject to } A_1 x \geq b, A_2 x \leq d, A_3 x = e)$$

$$\text{subject to } A_1 x \geq b$$

$$(A_1 = \begin{bmatrix} a_1^T \\ a_2^T \\ \vdots \\ a_m^T \end{bmatrix}, b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix})$$

$$A_2 x \leq d$$

$$\text{e.g. } (i.e. -a_i^T x \geq b_i)$$

$$A_3 x = e$$

(use matrices to express)

$$x_i \geq 0 \quad \forall i \in N_1$$

$$x_i \leq 0 \quad \forall i \in N_2$$

$$x_i \text{ free} \quad \forall i \in N_3$$

2) Standard Form: (of LP)

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad c^T x$$

(A — $m \times n$ matrix ($m < n$))

$$\text{subject to } Ax = b$$

$$(A \in \mathbb{R}^{m \times n}, m < n)$$

$$x \geq 0$$

infinite solutions

Transfer techs =

$$\textcircled{1} \max \rightarrow \min \quad \max c^T x \Leftrightarrow -\min -c^T x$$

$$\textcircled{2} Ax \leq b \text{ or } Ax \geq b \quad Ax + s = b, s \geq 0 \quad \& \quad Ax - s = b, s \geq 0$$

(s — slack variables)

$$\textcircled{3} x_i \leq 0, \text{ define } y_i = -x_i, y_i \geq 0 \quad x_i \leq 0 \quad (\text{every place } x_i \leftarrow -y_i)$$

$$\textcircled{4} x_i \text{ free}, x_i = x_i^+ - x_i^- \quad (\text{use two variables to substitute one})$$

$$x_i^+, x_i^- \geq 0$$

Examples.

eg1. $\max_{x_1, x_2} x_1 + 2x_2$ compact way
 subject to $x_1 \leq 100$ Standard form x_1, x_2, s_1, s_2, s_3
 $x_2 \leq 200$
 $x_1 + x_2 \leq 150$
 $x_1, x_2 \geq 0$

★ Standard form

$$\left[\begin{array}{ccc|c} 1 & 0 & 1 & 0 \\ 0 & 2 & 0 & 1 \\ 1 & 1 & 0 & 0 \end{array} \right] \quad \left[\begin{array}{c} x_1 \\ x_2 \\ s_1 \\ s_2 \\ s_3 \end{array} \right] = \left[\begin{array}{c} 100 \\ 200 \\ 150 \end{array} \right]$$

$x \geq 0 \rightarrow$ all variables!

subject to

$$\left[\begin{array}{cc} 1 & 0 \\ 0 & 2 \\ 1 & 1 \end{array} \right] \left[\begin{array}{c} x_1 \\ x_2 \end{array} \right] + \left[\begin{array}{c} s_1 \\ s_2 \\ s_3 \end{array} \right] = \left[\begin{array}{c} b_1 \\ b_2 \\ b_3 \end{array} \right]$$

$$(A \quad x + S = b)$$

eg2. $\min_{w, b, t} \sum_i t_i$
 subject to $y_i(x_i^T w + b) + t_i \geq 1$
 $t_i \geq 0$

standard form

$$\min_{w, b, t, b^+, b^-} \sum_i t_i$$

$$\text{subject to } y_i(x_i^T w^+ - x_i^T w^- + b^+ - b^-) + t_i - s_i = 1$$

$$w^+, w^-, b^+, b^- \geq 0$$

$$t_i, s_i \geq 0$$

- MATLAB Code:

CUX — similarly express to math
 ("CUX-begin" ... "CUX-end")

pay attention to $\left\{ \begin{array}{l} \text{"sum(sum(W.*X))"} \\ \text{"sum(X(1,:)) = 1"} \\ \text{"sum(X(:,1)} - sum(X(:,1)) = 0"} \end{array} \right.$
 (some high-level language)

Lecture 3 Linear Programming Modeling

- Modeling Transformation

ΔQ1: Air Traffic Control in OR field. (maximum problem)

Flights land in the order $1, 2, \dots, n$, time interval $[a_i, b_i] \forall i=1, 2, \dots, n$

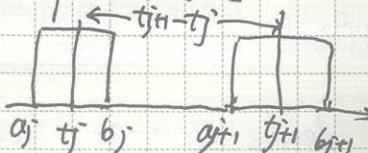
Objective — maximize the minimum separation time

$$\max_{t_j} \min_{t_i} \{t_{j+1} - t_j\}$$

$$\text{subject to } a_j \leq t_j \leq b_j \quad \forall j$$

$$t_j \leq t_{j+1} \quad \forall j \leq n-1$$

$$\text{Trick: } \Delta := \min_{j=1,2,\dots,n-1} \{t_{j+1} - t_j\}$$



LP: maximize Δ

Δ, t

subject to $t_j \geq -t_j \geq \Delta \quad \forall j = 1, 2, \dots, n-1$

 $b_j \geq t_j \geq a_j \quad \forall j$
 $t_j \geq -t_j \geq 0 \quad \forall j \leq n-1$

Generalize: maximin problem & minimax problem

maximize $\min_{i=1,2,\dots,n} \{c_i^T x + d_i\}$

subject to $Ax = b$

$x \geq 0$
no general way to do

with $\sum |x_i|$ problem (Q_2)

minimize $\max_{i=1,2,\dots,n} \{c_i^T x + d_i\}$

subject to $Ax = b$

$x \geq 0$

\rightarrow has LP forms

(piece-wise linear, convex form)

$y = \max_{i=1,2,\dots,n} \{c_i^T x + d_i\}$

$c_i^T x + d_i \geq y$

""

$\min_{x,t} \sum_{i=1}^n t_i$

$t_i \geq x_i$

$Ax = b$

ΔQ_2 : Minimize Absolute Values

$\min_x \sum_{i=1}^n |x_i|$

subject to $Ax = b$

(basis pursuit in signal processing)

(for $\max \sum |c_i^T x + d_i| \rightarrow$ non-convexity, non-LP!)

Generalize: Lemma — a modeling tool

minimize $\sum_{i=1}^n f_i(x)$

subject to $x \in S_2$

minimize $\sum_{i=1}^n t_i$

subject to $t_i \geq f_i(x) \quad \forall i$

(same optimal value)

ΔQ_3 : Linear Fractional Programming (LFP)

minimize $\frac{c^T x + d}{e^T x + f}$

assume $e^T x + f \neq 0$ for all $Ax \leq b$

subject to $Ax \leq b$

Trick: define $y = \frac{x}{e^T x + f}$ $z = \frac{1}{e^T x + f}$

whole proof process:

LP: minimize $c^T y + dz$

subject to $Ay - bz \leq 0$

$e^T y + fz = 1$

$z \geq 0$

sustituation, (y, z) is feasible with $z \geq 0$,

$x = \frac{y}{z}$ is feasible, (vice versa)

optimal value of LFP \geq ~ of origin ~ of origin

\Rightarrow The must be equivalent \geq not!

Note: when $z=0$, $x_0 + ty$ all feasible, (x_0 is feasible)

- Graphically Solving LP

Definitions:

1) **Polyhedron**: a set which can be written as $\{x \in \mathbb{R}^n \mid Ax \geq b\}$
 with $A \in \mathbb{R}^{m \times n}$ ($m \times n$ matrix), $b \in \mathbb{R}^m$

Corollary: $Ax = b, x \geq 0$ is a polyhedron

(Standard form)
 reason $\left\{ \begin{array}{l} Ax = b \\ x \geq 0 \end{array} \right. \Leftrightarrow \left\{ \begin{array}{l} Ax \geq b \\ -Ax \geq -b \\ x \geq 0 \end{array} \right. \Leftrightarrow \left[\begin{array}{c|c} A & b \\ -A & -b \\ I & 0 \end{array} \right] x \geq \left[\begin{array}{c} 0 \\ 0 \\ b \end{array} \right] \rightarrow \text{polyhedron}$

2) **Convex Set**: a set $S \subseteq \mathbb{R}^n$, for any $x, y \in S, \exists \lambda \in [0, 1]$
 (Affine Set) $\lambda x + (1-\lambda)y \in S$

3) **Convex Combination**: $\forall x_1, x_2, \dots, x_n \& \lambda_1, \lambda_2, \dots, \lambda_n \geq 0$ satisfying
 $\sum_{i=1}^n \lambda_i = 1, \sum_{i=1}^n \lambda_i x_i$ — convex combination
 of x_1, x_2, \dots, x_n

4) **Extreme point**: $x \in$ polyhedron P , s.t. we cannot find two elements
 $y, z \in P$ with $y, z \neq x$ & a scalar $\lambda \in [0, 1]$, satisfying
 (vertex/corner of a polyhedron) $x = \lambda y + (1-\lambda)z$

Finding extreme point in LP

minimize $C^T x$ $x \in \mathbb{R}^n$. A — $m \times n$ matrix with $m < n$
 subject to $Ax = b$ $b \in \mathbb{R}^m$
 $x \geq 0,$

* generally assume that A has full row rank m .
 (Otherwise, either can be redundant or empty)

Basic Solution:

x is called a basic solution of P if & only if:

(1) $Ax = b$; (2) \exists indices $B(1), B(2), \dots, B(m)$ s.t.

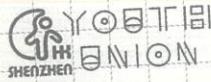
$A_B = [A_{B(1)} \ A_{B(2)} \ \dots \ A_{B(m)}]$ has all independent columns & $x_r = 0$,

$\forall i = B(1), B(2), \dots, B(m)$. (* note: A_i — i th column of A)

$$x_B = A_B^{-1} b \quad (\text{invertible})$$

$B = \{B(1), B(2), \dots, B(m)\}$ — basic indices

A_B — basic submatrix of A



* No more than m non-zeros could one have in basic solution (m constraints)
at most $\binom{n}{m} = \frac{n!}{m!(n-m)!} < +\infty$ basic solutions

Basic Feasible Solutions (BFS):

A basic solution x s.t. $x \geq 0$ can serve as a BFS.

• **Theorem:** For the standard LP polyhedron, $P = \{x \in \mathbb{R}^n \mid Ax = b, x \geq 0\}$
 x is an extreme point of $P \iff x$ is a BFS of P

Fundamental LP Theorem:

consider a LP in standard form, assume A — full row rank

(1) If the feasible set is non-empty, \exists (there exists) a BFS.

(2) If there's an optimal solution, \exists an optimal solution also BFS.

(Just look among BFSs to find optimal solution if existing)

If LP with m constraints has an optimal solution, \exists (must) an optimal solution with no more than m entries. (positive)

Tutorial — more modeling examples:

L_p -norm of $x = \|x\|_p = (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{\frac{1}{p}}$

exe 1: minimize $\|x\|_1$,
subject to $\|Ax - b\|_\infty \leq 1$ (ask for LP form)

$$\|x\|_\infty = \lim_{p \rightarrow \infty} (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{\frac{1}{p}} \in \lim_{p \rightarrow \infty} \max_{i=1,2,\dots,n} \{|x_i|\}^{p \cdot \frac{1}{p}} = \lim_{p \rightarrow \infty} (p \max_{i=1,2,\dots,n} \{|x_i|\})^{p \cdot \frac{1}{p}}$$

$$= \lim_{p \rightarrow \infty} \dots \text{(by squeezing theorem)} = \max_{i=1,2,\dots,n} \{|x_i|\}$$

$$\text{idea: } s_i = |x_i|$$

$$\overbrace{t = \max_{i=1,2,\dots,n} \{s_i\}}$$

minimize t^s

$s_i, s \in \mathbb{R}^+$

subject to $-s \leq x \leq s$

(substitution & transform L_p -norm)

$$\Rightarrow [\begin{smallmatrix} s \\ -s \end{smallmatrix}] \leq Ax - b \leq [\begin{smallmatrix} s \\ s \end{smallmatrix}]$$

$$\min_{x^+, x^-} t^s$$

$$\text{subject to } [\begin{smallmatrix} s \\ -s \end{smallmatrix}] \leq Ax^+ - Ax^- - b \leq [\begin{smallmatrix} s \\ s \end{smallmatrix}]$$

(understand, give $|x_i|$, either x_i^+ or x_i^- should be zero to minimize t^s)

\hookrightarrow find implicit x

$$\min_{x^+, x^-} C^T x$$

$$\text{subject to } Ax = b$$

Solutions: use the background knowledge $C = A^T \lambda + C'$ in the nullspace of A
in rowspace of A (given)

(i) if $b \notin C(A)$, not feasible / infeasible

(ii) if $b \in C(A)$ with $C' = 0$, then

$$C^T x = A^T A x = A^T b \text{ is a specific #.}$$

(iii) if $b \in C(A)$ with $\lambda = 0$, then $C = C'$, note that for $x = x_0 - t\zeta$

feasible for all t , $C^T x = C^T x_0 - t C^T \zeta$ can be $-\infty$

(iv) $\lambda \neq 0$, same with (iii) (unbounded)

$$\text{optimal value } p^* = \begin{cases} \infty, & b \notin C(A) \\ A^T b, & b \in C(A), C = A^T \lambda \\ -\infty, & \text{elsewhere} \end{cases}$$

exe 3. How to express projection of a vector P on plane S .

Projection \Rightarrow minimize the distance from P to S

$$\text{minimize } \|P - X\|^2$$

$$\begin{matrix} x \\ \text{subject to } X \in S \end{matrix}$$

exe 4. (Knapsack Problem)

Totally N numbers of goods, with weight $1, 2, \dots, N$ & value $1, 2, \dots, N$
 (w_i) (v_i) ,

The package can pack at most weight W , how to get more value.

(use $c_i \in \{0, 1\}$ for choice & not choosing)

$$\text{maximize } \sum_{i=1}^N c_i v_i$$

$$\text{subject to } \sum_{i=1}^N c_i w_i \leq W$$

$$c_i \in \{0, 1\}, \forall i = 1, 2, \dots, N$$

Lecture 4 The Simplex Method

- IDEA: proceeding from one BFS to a neighbouring one
 to continuously improve (increase/reduce) the value → standard LP form

• Neighboring / Adjacent Basic Solutions.

Two basic solutions are neighboring / adjacent if they differ by exactly one basic index. (one basic index differs \Leftrightarrow one non-basic index differs)

Efficient way to check "all" neighbors — (avard $m(n-m)$ choice of every time taking inverse)

1) First, rewrite $Ax=b$ as $[A_B \ A_N] \begin{bmatrix} x_B \\ x_N \end{bmatrix} = b$ (re-arrange these columns)

2) Finding a random number in non-basic set N , j , wanting to increase x_j to get the neighbors. Move $x_j \rightarrow x_j + \theta d_j$ (\tilde{x}), $\theta \geq 0$

d satisfies: 1) $d_j=1$; 2) $d_{j'}=0, \forall j' \in$ Basic "feasible" set

3) $A\tilde{x}=b$, therefore, $Ad=0 \Rightarrow \underbrace{A_B d_B}_0 + \underbrace{A_N d_N}_0 = 0$

d_B (specific) = $\tilde{A}_B^{-1} \tilde{A}_j$, $d = \begin{bmatrix} d_B \\ d_N \end{bmatrix}$, where d_N has $d_j=1$, all others zero. \uparrow full rank \tilde{A}_j

d is called — j th basic direction

3) Making \tilde{x} feasible, $\tilde{x} = x + \theta d \geq 0, \theta \geq 0$

(non-basic variable i , $x_i=0$. basic variable j , small θ satisfying $x+ \theta d \geq 0$)

\uparrow these contain just the case $x_B \geq 0$ strictly
when $x_B > 0$ (some entries=0). degeneracy

• Reduced Costs.

1) rewrite $C^T = [C_B^T \ C_N^T]$, then $C^T d$ (same as above)

2) $C^T d = [C_B^T \ C_N^T] \begin{bmatrix} d_B \\ d_N \end{bmatrix} = C_B^T d_B + C_N^T d_N = \boxed{c_j - c_B^T \tilde{A}_B^{-1} \tilde{A}_j =: \bar{c}_j}$ reduced cost

\Rightarrow positive \bar{c}_j — increase; do not want to go

negative \bar{c}_j — decrease; go in that direction (indicators of where to go)

Extension: j in B ($j=B(i)$), $\bar{c}_{B(i)} = c_{B(i)} - \underbrace{c_B^T \tilde{A}_B^{-1} \tilde{A}_B(i)}_{A_i} = 0$

Stopping Point / Criterion:

Consider a basic solution x associated with the basic $B(1), B(2), \dots, B(m)$,

\bar{c} — the corresponding vector of reduced costs. If $c \geq 0$, then x must be optimal.

- Choosing a stepsize (θ) (Change of Basis)

Want θ to be as large as possible — reduce more

$$\star \theta^* = \max \{ \theta \geq 0 \mid x + \theta d \geq 0 \} \rightarrow \text{can get } \left\{ \begin{array}{l} d \geq 0, \theta = \infty, \text{ unbounded} \\ \theta^* = \min_{\{i: d_i < 0\}} -\frac{x_i}{d_i} \rightarrow \text{some } d_i < 0, \text{ this equation} \end{array} \right.$$

New basis $\tilde{x} = x + \theta^* d$ (someone $B(l)$ out & j in)
 $x_j > 0, x_{B(l)} = 0$

- The whole Process of Simplex Method

Step 1. Reduced Costs $C \rightarrow$ can we continue? which way to choose?

Step 2. j th basic directions \rightarrow given j , how to go? (d, θ^*)

Step 3. Move on to $x + \theta^* d$ & then repeat.

↑ make the whole process complete

- Additions: degenerate cases & choosing strategy

1) Theorem — properties of \tilde{x} : Let x be a non-degenerate BFS ($x \geq 0$) with basic indices B , and let $y = x + \theta^* d$ be generated by simplex iteration. Then, y is a basic feasible solution. (A_B is full rank) (with new indices)

2) Theorem — convergence: assume that the feasible set is non-empty & every BFS is non-degenerate. Then, the simplex method terminates after a finite # iterations, with the following options,

1) Stop with a BFS as optimal solution ; 2) $d \geq 0, C^T d < 0$, optimal value $-\infty$

3) Treat with degeneracy =

Still change basic index from i to j , with $\theta^* = 0$ to next iteration

Although objective value does not change, basis $B \rightarrow$ basis \hat{B} , do changes

$\Rightarrow \min_{\{i: d_i < 0\}} -\frac{x_i}{d_i}$ may change \Rightarrow finite times then to non-degeneracy.

4) Pivot Rules (avoid cycle in degeneracy, efficiently)

* choosing j to go → direction, enter

(Most common) - Smallest Index Rule (smallest j s.t. $\bar{c}_j < 0$)

Most negative Rule (smallest $\bar{c}_j < 0$)

Most Decrement Rule (smallest $\theta^* \bar{c}_j < 0$)

* Choosing θ^* (one index to leave) → leave

Smallest Index Rule (smallest j s.t. $\theta^* = -\frac{x_i}{a_{ij}}$, when $-\frac{x_i}{a_{ij}}$ is one of the minimum)
at degenerate cases, $\geq k$ planes intersect in (at a point) k dimension.

Theorem - Bland's Rule:

Use smallest rules (index) for choosing both the entering basis & the leaving basis,

No cycle will occur! — stop within finite times

Lecture 5 The Simplex Tableau

- Proof of property 1 in lecture 4:

$$m = \text{rank}(I) = \text{rank}(A_B^{-1} A_B) \leq \min \{ \text{rank}(A_B^{-1}), \text{rank}(A_B) \} \leq m$$

$$\Rightarrow \text{rank}(A_B^{-1}) = m. \quad \checkmark \text{rank}$$

$$\text{Because } \text{rank}(A_B^{-1} A_B) = \text{rank} \left(\sum_{i \in B \setminus \{B \cap C\}} A_B^{-1} A_i + A_B^{-1} A_j \right) = m$$

$$\text{Then } m = \text{rank}(A_B^{-1} A_B) \leq \min \{ \text{rank}(A_B^{-1}), \text{rank}(A_B) \} \leq m$$

$$\Rightarrow \text{rank}(A_B^{-1}) = m$$

- Finding Initial BFS:

1) Directly finding way: (adding slack variables (totally m amounts))

(Because the matrix which slack variables face is I) \Rightarrow BFS choosing all variables

2) Two-phase Method (use simplex method for two times) $\begin{cases} \text{more common} \\ \text{one} \end{cases}$

① minimize $\mathbf{1}^T \mathbf{y}$

x, y

subject to $Ax + y = b$ ← auxiliary problem

$x, y \geq 0$

② Standard LP form

change the row

Explanation: (with the loss of generality, assume $b \neq 0$, otherwise $-A^T x = -b$)

Initial BFS for auxiliary problem: $y = b, x = 0$ (written as $[A \quad I_m] \begin{bmatrix} x \\ y \end{bmatrix} = b$)

$\underbrace{\quad}_{\text{BFS, } m \text{ lines}}$

\widehat{A}

\widehat{x}

Apply simplex method to find the optimal solution

★ Theorem: Feasibility

The original problem is feasible if & only if the optimal value of the auxiliary problem is 0.

Optimal value $> 0 \Rightarrow$ cannot find x s.t. $Ax=b$ ($Ax < b$)

Optimal value = 0 \Rightarrow $Ax^*=b$, $x^* \geq 0$ & BFS x^* has new x_B^*

For degenerate case (x^* has less than m positive entries),
pick other columns to supplement, get an $m \times m$ matrix with rank = m

3) The Big-M Method

(transform to another auxiliary)

$$\underset{x,y}{\text{minimize}} \quad C^T x + M \cdot \mathbf{1}^T y$$

$$\text{subject to } Ax + y = b$$

$$x, y \geq 0$$

Explanation: BFS (initial) = ($y=b \geq 0$)

M should be large enough to avoid $y > 0$ case (infeasible)
 $y=0$ (feasible), optimal value = $C^T x$ (unbounded)

• Simplex Tableau

(avoid A_B^{-1} , simply solving by hand)

$C^T - C_B^T A_B^{-1} A$	$-C_B^T A_B^{-1} b$	→ current "optimal value"
--------------------------	---------------------	---------------------------

$A_B^{-1} A$	$A_B^{-1} \cdot b$	→ current BFS
--------------	--------------------	---------------

all $-d_s' \leftarrow (d = -A_B^{-1} A)$

$$A_B^{-1} b = A_B^{-1} A x = A_B^{-1} [A_B \quad A_N] \begin{bmatrix} x_B \\ x_N \end{bmatrix} = I x_B = x_B$$

The simplex

Om^T	$C_B^T - C_A^T A_B^{-1} A_N$	$-C_B^T X_B$	canonical form
I_m	$A_B^{-1} A_N$	X_B	

★ How to use simplex tableau:

(reduced costs - incoming basic index) - θ^* & outgoing $B(l)$ - new basis

pivoting — transform between two canonical forms

e.g. $B \mid -1 \quad -2 \quad 0 \quad 0 \quad 0 \quad 1 \quad 0$ $\xrightarrow{\text{choose } 1 = A_B^{-1} A_N}$ outgoing column/index
 $\xrightarrow{\text{pivot row}} \text{production plan in lecture 1}$

pivot $\xrightarrow{\text{pivot}} \begin{matrix} 3 \\ 4 \\ 5 \end{matrix} \mid \boxed{1} \quad 0 \quad 1 \quad 0 \quad 0 \quad 100$ pivot row
 $\xrightarrow{\text{pivot column}} \begin{matrix} 1 \\ 2 \\ 0 \\ 1 \end{matrix} \mid 0 \quad 1 \quad 0 \quad 200$ pivot column
 $\xrightarrow{A_B^{-1} A_N} \begin{matrix} 1 \\ 1 \\ 1 \end{matrix} \mid 0 \quad 0 \quad 1 \quad 150$

$\Rightarrow X_B = \begin{bmatrix} 0 \\ 100 \\ 200 \\ 150 \end{bmatrix}$ (Bland's Rule)

choosing θ^* (Minimal Ratio Test / MRT)

$\star \theta^* = \min_i \left\{ \frac{b_i}{A_{ij}} \mid A_{ij} > 0 \right\}$ (original $\theta^* = \min_{\text{all } i \in B} -\frac{x_i}{d_i}$)

★ Unbounded — $A_{ij} \leq 0$ for all i in one j , ($d \geq 0$)

Reason: $Q A_B^{-1} = A_B^{-1}$

$\xrightarrow{*Q \Leftrightarrow \text{do elimination}}$ $\xrightarrow{Q_{ii}=0 \text{ & all others } 1, 0 \dots}$

Steps — Elimination: $\star \text{pivot} = 1$, other entries = 0 in the column.
 (updating the tableau) \xrightarrow{I}

$B \mid 0 \quad -2 \quad 1 \quad 0 \quad 0 \quad 100$ first step \xrightarrow{I}
 $\xrightarrow{1 \quad 1 \quad 0 \quad 1 \quad 0 \quad 100} \xrightarrow{2 \quad 0 \quad 2 \quad 0 \quad 1 \quad 200}$
 $\xrightarrow{5 \quad 0 \quad 1 \quad -1 \quad 0 \quad 50} \xrightarrow{4 \quad 0 \quad 1 \quad -1 \quad 0 \quad 150}$ second step

$\Rightarrow \dots \text{ until } C \geq 0 \text{ and } d \geq 0$

Lecture 6 Duality Theory(I)

- Complexity Theory

Complexity, big-Oh Notation, how to compute

(omitted, refer to 2020-21 Term 2 CSC1001 Notes)

Polynomial-Time Algorithm = (some)

- 1) Gaussian Elimination $O(n^3)$
- 2) Fast method for matrix inversion $O(n^{2.373})$
- 3) Naive method of sorting $O(n^2)$

Non-polynomial-Time Algorithm = (some)

- 1) Enumeration method for traveling salesman problem $O(n!)$
- 2) Dynamic Programming for TSP $\downarrow O(2^n n^2)$

Given list of cities & distances between each two, asking for shortest route, visiting every city exactly once & return to the primitive one.

★ Definition:

P — a problem, for which \exists a polynomial-time algorithm. (solvable)

$NP\text{-hard problems}$ — a class of problems with currently no polynomial-time algorithm to solve

$\Rightarrow LP$ is a polynomial-time solvable problem.

Ellipsoid Method "Worst" \leftarrow contradict
 (theoretically good case) \leftarrow (practically bad to manipulate)
 Simplex Method
 (theoretically bad)
 (practically good to use)

Both good (practice & theory) = interior-point Method

• Duality For Linear Programming

(transform)

$$\begin{array}{l} \text{minimize}_{\mathbf{x}} \mathbf{c}^T \mathbf{x} \\ \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{b} \\ \quad \mathbf{x} \geq 0 \end{array}$$

primal problem

$$\begin{array}{l} \text{minimize}_{\mathbf{x}} \mathbf{c}^T \mathbf{x} \\ \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{b} \\ \quad \mathbf{x} \geq 0 \end{array}$$

$$\begin{array}{l} \text{minimize}_{\mathbf{x}} \mathbf{c}^T \mathbf{x} + \max_{\mathbf{y} \geq 0} \mathbf{y}^T (\mathbf{A}\mathbf{x} - \mathbf{b}) \\ \text{subject to } \mathbf{x} \geq 0 \end{array}$$

for some \mathbf{y}

$\left\{ \begin{array}{l} \mathbf{A}\mathbf{x} \neq \mathbf{b}, \mathbf{0} \\ \mathbf{Ax} = \mathbf{b}, \mathbf{0} \end{array} \right. \text{ (give up)}$

(reason)

$$\begin{array}{l} \left\{ \begin{array}{l} \min_{\mathbf{x}} \max_{\mathbf{y}} \mathbf{c}^T \mathbf{x} + \mathbf{y}^T (\mathbf{A}\mathbf{x} - \mathbf{b}) \\ \text{subject to } \mathbf{x} \geq 0 \end{array} \right. \end{array}$$

(reason later)

$$\begin{array}{l} \left\{ \begin{array}{l} \max_{\mathbf{y}} \mathbf{b}^T \mathbf{y} + \min_{\mathbf{x} \geq 0} \mathbf{x}^T (\mathbf{c} - \mathbf{A}^T \mathbf{y}) \end{array} \right. \end{array}$$

$\Leftrightarrow \text{maximize}_y b^T y$ $A^T y > c, -\infty \text{ (give up)}$
 subject to $A^T y \leq c$ $A^T y \leq c, 0 \leq x \geq 0$
 (reason) $\text{Lagrange (saddle point)}$
 dual problem
 y - dual variables . P (primal), D (dual), have the same buyer optimal value. x^*, p^* - optimal value

★ Dual Problem of all kinds of LP form (can be proved use tricks shown above)

$$\begin{array}{|l}
 \hline
 P \\
 \min_x c^T x \\
 \text{subject to } a_i^T x \geq b_i, i \in M_1 \\
 a_i^T x \leq b_i, i \in M_2 \\
 a_i^T x = b_i, i \in M_3 \\
 x_j \geq 0, j \in N_1 \\
 x_j \leq 0, j \in N_2 \\
 x_j \text{ free}, j \in N_3 \\
 (a_i^T \Rightarrow \text{row } i \text{ of } A) \\
 \hline
 \end{array}
 \quad
 \begin{array}{|l}
 \hline
 D \\
 \max_y b^T y \\
 \text{subject to } y_i \geq 0, i \in M_1 \\
 y_i \leq 0, i \in M_2 \\
 y_i \text{ free}, i \in M_3 \\
 A_j^T y \leq c_j, j \in N_1 \\
 A_j^T y \geq c_j, j \in N_2 \\
 A_j^T y = c_j, j \in N_3 \\
 (A_j^T \Rightarrow \text{column } j \text{ of } A) \\
 \hline
 \end{array}$$

Simple form	
P	$\min_x c^T x$
constraint	$\geq b_i$
-TS	$\leq b_i$
	$= b_i$
	≥ 0
	≤ 0
	free
	$\geq c_j$
	$\leq c_j$
	$= c_j$

(eg. Soft SVM)
-margin

$$\text{minimize}_{w, b, t} \sum_{i=1}^m t_i$$

$$\text{Subject to } y_i(x_i^T w + b) + t_i \geq 1, \forall i = 1, 2, \dots, m$$

$$t_i \geq 0$$

$$\begin{aligned}
 A &= \begin{bmatrix} y_1 x_1^T & y_1 \\ \vdots & \vdots \\ y_m x_n^T & y_m \end{bmatrix} \stackrel{\text{block I}}{\rightarrow} \begin{bmatrix} w \\ b \\ t \end{bmatrix} \stackrel{\text{block II}}{\rightarrow} y \\
 &\stackrel{\text{block III}}{\rightarrow} I_m \\
 &\text{block I} \rightarrow \text{diag}(y) X^T, \text{ with } X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \\
 &= [\text{diag}(y) X^T \ y \ I_m] \begin{bmatrix} w \\ b \\ t \end{bmatrix}
 \end{aligned}$$

$$\begin{array}{|l}
 \hline
 \text{Dual problem} \\
 (\text{of } P) \quad \text{maximize}_{u} \sum_{i=1}^m u_i \\
 \text{subject to } X \text{diag}(y) \cdot u = 0 \\
 y^T u = 0 \\
 0 \leq u \leq 1 \\
 \hline
 \end{array}
 \quad
 \left. \begin{array}{l} \\
 \end{array} \right\} \text{separately written way}$$

• Weak Duality

★ Theorem (duality under transformations) :

Transform an LP to an equivalent one will get equivalent two problems.

★ Theorem (doubly duality) :

Doubly transform the primal problem (i.e. transform the dual to its dual)
will get the identity (primal one).

★ Weak Duality theorem :

If x, y are feasible to corresponding problem respectively,
 then we have $b^T y \leq c^T x$

(any dual feasible solution will give a lower bound on primal optimal value.
 any primal feasible solution will give an upper bound on dual optimal value.)

$$\text{proof: } b^T y = (Ax)^T y = x^T (A^T y) \leq x^T c = c^T x$$

corollary I: primal problem $\begin{cases} \text{unbounded} \\ \text{infeasible} \end{cases} \rightleftharpoons \begin{cases} \text{infeasible} \\ \text{unbounded} \end{cases} \leftarrow \text{dual problem}$

corollary II: x, y - feasible pts (let be) of the primal & dual problem
 respectively. If $c^T x = b^T y$, then x, y must be optimal solutions,
 correspondingly.

From corollary II $\begin{cases} Ax = b, x \geq 0 \\ \text{sufficient optimality conditions } A^T y \leq c \\ c^T x = b^T y \end{cases} \Rightarrow \text{optimal solution / value}$
 (Unify finding linear opt & linear feasible pt.)
 same hard theoretically

Lecture 7 Duality Theory (II)

- Strong Duality (Theory)

DOES NOT ALWAYS HOLD!!!

It holds when duality gap=0.

In the case of LP,
 it always works

If the linear program has an optimal solution, so does its dual & the optimal values of the two are equal.

i.e. the optimal conditions are "exactly when" conditions

$$\begin{cases} Ax^* = b, x^* \geq 0 \\ A^T y^* \leq c \\ b^T y^* = c^T x^* \end{cases} \Leftrightarrow x^*, y^* \text{ optimal solution}$$

exactly when

Solving LPs

\Leftrightarrow solving linear system

★ (optimization)

\Leftrightarrow solving systems

Proof (under simplex method):

suppose x^* is an optimal solution & BFS. (Fundamental LP)

construct $y^T = C_B^{-1} A_B^{-1} b$ (namely, $y = (A_B^{-1})^T C_B$)

$$\begin{cases} b^T y = y^T b = C_B^T A_B^{-1} b = C_B^T A_B^{-1} A_B x_B = C_B^T x_B = C^T x^* & (\text{optimal } y \text{ because of} \\ \text{the weak duality.}) \\ A^T y \leq c \text{ because of } y^T A - C^T = C_B^T A_B^{-1} A - C^T = -\bar{c} \leq 0 \rightarrow \text{feasible} \\ y^T = C_B^T A_B^{-1} b \rightarrow \text{dual optimal solution} \end{cases}$$

* All (nearly) LP alg. (simplex, interior pt., ellipsoid) have corresponding y^T / y (primal & dual can be solved simultaneously)

Possible relations:

	P	Finite opt.	Unbounded	Infeasible
P	✓	—	—	
Unbounded	—	—	✓	
Infeasible	—	✓	✓	

- 1) If both primal & dual are feasible, they have bounded opt. values.
- 2) Their opt. values same by strong duality.

• Complementarity Conditions

Let x, y be feasible pts of the primal & dual problem, respectively.

Then, x & y are optimal solutions exactly when

$$x_i \cdot (c_i - A_i^T y) = 0, \forall i=1,2,\dots,n \quad (\text{or} \Leftrightarrow \begin{cases} 1) x_i > 0, c_i = A_i^T y \\ 2) x_i = 0, c_i > A_i^T y \\ 3) \text{both zeros>equals} \end{cases})$$

proof: by strong duality, optimal solutions $x \neq y$

$$\Leftrightarrow 0 = b^T y + c^T x \Leftrightarrow -(Ax)^T y + c^T x = 0 \Leftrightarrow -x^T A^T y + x^T c = 0$$

$$\Leftrightarrow x^T (c - A^T y) = 0 \Leftrightarrow \begin{cases} x_i > 0, c_i \geq A_i^T y \\ \sum_{i=1}^m x_i (c_i - A_i^T y) = 0 \end{cases} \Leftrightarrow \text{all } x_i \cdot (c_i - A_i^T y) = 0$$

Complementarity conditions common (as shown above)
 $\sim \text{slackness} \sim = A^T y + s = c, s \geq 0 \quad \& \quad x_i \cdot s_i = 0, \forall i$
 $\sim \text{general} \sim = y_i(a_i^T - b_i) = 0, \forall i \quad \& \quad x_j(A^T y - c_j) = 0, \forall j$
 ↳ "the primal-dual transformation form".

Then, we can get an equivalent set of optimal conditions

$$\begin{cases} Ax = b, x \geq 0 \text{ (or } x \text{ is feasible)} \\ A^T y \leq c, (or y \text{ is feasible}) \\ x_i(A^T y - c_i) = 0, \forall i, \text{ (or the complementarity conditions hold)} \end{cases}$$

can be used to solve some problems

$$(\text{by the way, } A_B^T y = c_B \implies y = (A_B^T)^{-1} c_B, x_N(A^T y - c_N) = 0 \text{ (because } x_N = 0))$$

• Duality via Simplex Tableau

If the simplex tableau is got by adding m slack variables, namely

$$\begin{array}{c|cc|c} C^T & 0_m & 0 & \\ \hline A & I_m & b & \\ \hline \text{original } x_B & & & \\ & & & & & \end{array}, \text{ then, after achieving optimal form } \rightarrow y^* \\ \begin{array}{c|cc|c} C^T - C_B^T A_B^{-1} A & -C_B^T A_B^{-1} & -C_B^T A_B^{-1} b & \\ \hline A_B^{-1} A & A_B^{-1} & A_B^{-1} b & \\ \hline \end{array} \\ (\text{correspond } [A^T; I] \\ C^T - C_B^T A_B^{-1} A, 0^T - C_B^T A_B^{-1} I_m) \\ (\text{by elimination})$$

• Modeling Examples – application with dual

Exp 1 (production planning)

	Steel	Iron	Copper	Profit
alloy I	1	0	1	\$1
alloy II	0	2	1	\$2
Resources	100	200	150	

$$\text{maximize } x_1 + 2x_2$$

$$x_1, x_2$$

$$\text{subject to } x_1 \leq 100$$

$$2x_2 \leq 200$$

$$x_1 + x_2 \leq 150$$

$$x_1, x_2 \geq 0$$

Transform to dual problem:

$$\text{minimize } 100P_1 + 200P_2 + 150P_3$$

$$P_1, P_2, P_3$$

$$\text{subject to } P_1 + P_2 \geq 1$$

$$2P_2 + P_3 \geq 2$$

$$P_1, P_2, P_3 \geq 0$$

Physical meaning at price p , selling resources.

From perspective of buyers, minimize total cost \rightarrow objective function
the selling price \geq profit \rightarrow constraints
otherwise the company (selling) will stop.

(Meaningful question)

Exp 2 (Multi-Firm Alliance)

Suppose firm $1, 2, 3, \dots, m$, making same set of products
alliance — they put all resources in a pool & use them jointly.

$$\Rightarrow \begin{array}{ll} \text{maximize}_x & V^i = C^T x \\ \text{subject to} & Ax \leq b_i \\ & x \geq 0 \end{array} \quad \text{for every company without alliance}$$

\hookrightarrow consumption matrix

$$\Rightarrow \begin{array}{ll} \text{maximize}_x & V^S = C^T x \\ \text{subject to} & Ax \leq \sum_{i \in S} b_i \\ & x \geq 0 \end{array} \quad \text{for every subset of company.} \quad (S \subseteq \{1, 2, \dots, m\})$$

Then the original LP becomes:

$$\begin{array}{ll} \text{maximize}_{x, z} & V^* = C^T x \\ \text{subject to} & Ax \leq \sum_{i=1}^m b_i \\ & x \geq 0 \\ (\text{allocation}) & \sum_{i=1}^m z_i = V^* \\ (\text{constraints}) & \sum_{i \in S} z_i \geq V^S \quad \forall S \subseteq \{1, 2, \dots, m\}. \text{ totally } 2^m \text{ inequalities.} \\ & \text{otherwise they will leave} \end{array}$$

Dual. minimize $\left(\sum_{i=1}^m b_i \right)^T y$
subject to $Ay \geq C$
 $y \geq 0$

set $z_i = b_i^T y^*$ then $\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix}$ is the core of grand alliance

Check: ① By strong duality $\sum_{i=1}^m z_i = V^*$
② By weak duality $\sum_{i \in S} z_i \geq V^S$

Lecture 8 Sensitivity Analysis

- Modeling examples (cont'd)

Exp 3 (alternative systems)

CQ: how to verify $A^T y \leq c$ does NOT have solutions?
(prove non-existence)

Answer: If $\bar{P} = \{x | Ax=0, x \geq 0, c^T x < 0\}$ is NOT empty,
then there must be no solution to $A^T y \leq c$.

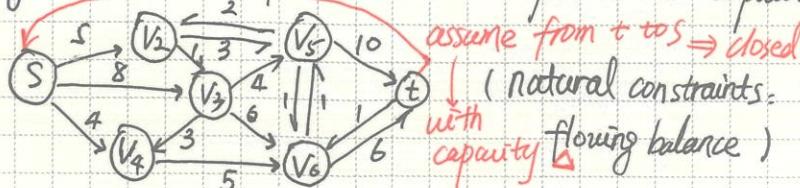
(Reason: duality theory)

(primal) $\min_{x \geq 0} c^T x$ subject to $Ax=0$ (dual, max $0^T y$)
 subject to $Ay \leq c$ or weak duality
 $0^T y \leq c^T x$ (any x)
 infeasible unbounded
 ↙ unbounded ↛ infeasible
 (If x satisfied, any $ax, a \geq 0$ satisfied,
 then $c^T x$ unbounded.)

Exp 4 (Maximum Flow Problem)

- > Given a directed, weighted graph $G=(V, E)$ & a pair of nodes s, t .
- > There is an edge capacity on each edge c_{ij} .

Q: largest flow sent from s to t subject to the capacity



A: maximize $\Delta \rightarrow$ flow from t to s / total flows leaving s

subject to $\sum_{j: (j, i) \in E} x_{ji} - \sum_{j: (i, j) \in E} x_{ij} = 0 \quad \forall i \neq s, t, i \in V$
 (balance)

$$\sum_{j: (j, s) \in E} x_{js} + \Delta - \sum_{j: (s, j) \in E} x_{sj} = 0$$

$$\sum_{j: (j, t) \in E} x_{jt} - (\sum_{j: (t, j) \in E} x_{tj} + \Delta) = 0$$

$$0 \leq x_{ij} \leq c_{ij} \quad \forall (i, j) \in E$$

Take the dual problem as follows,

derive

$$\sum_{(i,j) \in E} (-w_{ij}x_{ij} + z_{ij}(x_{ij} - c_{ij})) + \sum_{k \in V \setminus \{i\}} y_k (\sum_{(i,k) \in E} x_{ik} - \sum_{(k,j) \in E} x_{kj})$$

GO YOUTUBE
SHENZHEN UNION

$$\Leftrightarrow \sum_{(i,j) \in E} M_{ij}(z, w, y) \leq \sum_{(j) \in E} z_{ij} c_{ij}$$

want this to be the primal (optimal) objective value

i.e. $\sum M_{si} = z_{si} - w_{si} + x_i = 1$

want $\begin{cases} M_{ij} = z_{ij} - w_{ij} + x_j - x_i = 0 \\ M_{jt} = z_{jt} - w_{jt} - x_j = 0 \end{cases}$

$$\Rightarrow \underset{\substack{z \in R^{|E|}, w \in R^{|V|}, (i,j) \in E}}{\text{minimize}} \sum z_{ij} c_{ij}$$

subject to

$$z_{ij} \geq x_i - x_j \quad (\text{with } w_{ij} \geq 0)$$

$$x_s = 1 \quad \& \quad x_t = 0 \quad / \quad x_s - x_t = 1$$

$$z_{ij} \geq 0$$

suppose $x_i \in \{0, 1\}$

$$\begin{cases} \text{if } x_i > x_j, z_{ij} = x_i - x_j = 1 \\ \text{if } x_i \leq x_j, z_{ij} = 0 \end{cases}$$

$$\begin{array}{c} \text{or } \\ \text{a cut } S \end{array} \begin{array}{c} 0 \\ 0/1 \\ 1 \end{array} \begin{array}{c} \text{all } x_i = 1 \\ S' \text{ --- all } x_j = 0 \end{array}$$

" $\sum_{i \in S, j \in S'} c_{ij}$ " min-cut problem

minimize the weights of a cut

("the tightest bottle neck of the network")

• Sensitivity Analysis

CQ: when inputs change (i.e. A/b/c changes in standard form),
how does the optimal value change?

* Local sensitivity (V-optimal value, with some caps fixed)

{ If the dual has a unique optimal solution y^* , then $\nabla V(b) = y^*$ (gradient)

If the primal has a unique optimal solution x^* , then $\nabla V(c) = x^*$

holds for Standard form

a way to judge uniqueness:

"In simplex tableau, $\bar{c}_i = 0, \bar{t}_i \in B; \bar{c}_j > 0, \bar{t}_j \in N$ " (Ax $\leq b$)

Local sensitivity with changes:

Under the conditions above, with SMALL change in b/c
the change of objective value $\Delta V = \Delta b y^* \text{ or } \Delta c^T x^*$

Note: under inactive constraints, ΔV can be zero (because of $q_i^* = 0$ / $x_i^* = 0$)

★ Global Sensitivity:

CQ: Under what scale of change can the optimal basis hold?

1) Change b : $b \rightarrow b + \Delta b$

- (i) Reduced costs: $\bar{c} = c^T - C_B^T A_B^{-1} A$ independent of b
- (ii) Feasible? — $\tilde{x}_B = A_B^{-1}(b + \Delta b) = x_B^* + A_B^{-1}\Delta b$ $\begin{cases} \geq 0, & B - \text{still opt. basis} \\ \text{otherwise, restart.} & V(B) = V^* + \Delta b^T y^* \quad (y = A_B^{-1} c_B) \end{cases}$

Real operation: write Δb as λe_i , with e_i having the (some cases) i^{th} entry = 1.

Check whether $\underline{x}_B^* + \lambda A_B^{-1} e_i \geq 0$

2) Change c : $c \rightarrow c + \Delta c$

- (i) Feasible? — $\tilde{x}_B = A_B^{-1} b$ independent of c .

$$\text{(ii) reduced costs: } \bar{c} = c^T - C_B^T A_B^{-1} A = \begin{cases} 0, & \text{for } B \\ C_N^T - C_B^T A_B^{-1} A_N, & \text{for } N \end{cases}$$

$\Rightarrow \begin{cases} \Delta c = \lambda e_j \text{ (same as above)} \\ j \in B: \bar{c} = C_B^T - C_N^T A_B^{-1} A_N - \lambda e_j^T A_B^{-1} A_N = r_N^T - \lambda e_j^T A_B^{-1} A_N \end{cases}$

$$\begin{cases} \geq 0, & \text{still opt. basis} \\ \text{otherwise, restart} & \end{cases}$$

$$\begin{cases} j \in N: \bar{c} = C_B^T + \lambda e_j^T - C_N^T A_B^{-1} A_N = r_N^T + \lambda e_j^T \geq 0, & \text{still opt. basis} \\ \text{otherwise, restart} & \end{cases}$$

(Note: sign of λ will change for "max" prob.)

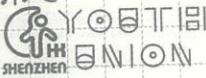
What if the change outside the range?

Change c — reduced cost renew

Change b — first transfer to the dual — reduced cost for the dual renew

3) Change A =

$\left\{ \begin{array}{l} \text{change in } A_j, j \in N, \text{ then compute } \bar{G} \geq 0, \text{ optimal solution still} \\ \text{otherwise, update for} \\ \text{the } j^{\text{th}} \text{ column} \\ \text{change in } A_j, j \in B, \text{ then } \bar{x}_B = \bar{A}_B^{-1} \bar{b} \text{ changes, no simple way.} \end{array} \right.$



- Some additions:

1) proof of "extreme pts \Leftrightarrow BFS":

BFS \Rightarrow extreme points (argue by contradiction & definition)

extreme pts \Rightarrow BFS (argue by contradiction & construction)

\hookrightarrow If not BFS, then at least more than n entries > 0 (otherwise, BFS)

By the definition of linear (in)dependence, $\exists \alpha \neq 0$ s.t. $\alpha_1 A_{B(1)} + \dots + \alpha_K A_{B(K)} = 0$

Construct two feasible pt. $\bar{x} = \bar{x}_B^* - \varepsilon \alpha$ & $\bar{x} = \bar{x}_B^* + \varepsilon \alpha$ (with small enough $\varepsilon > 0$)
(basis part)

then $\bar{x}^* = \frac{1}{2}\bar{x}^- + \frac{1}{2}\bar{x}^+$, contradictory with extreme points.

2) proof of "the fundamental LP theorem":

(a kind of induction method)

(i) a solution/feasible \rightarrow at least a BFS: every time reduce at least one entry to zero.

with linear (in)dependence, $\exists \alpha \neq 0$, s.t. $\sum_{i=1}^k \alpha_i A_{B(i)} = 0$ with at least 1 $\alpha_i > 0$

construct $\tilde{x}_B = \bar{x}_B - \varepsilon \alpha$, with $\varepsilon = \min \{ \alpha_i > 0 \mid \frac{\bar{x}_{B(i)}}{\alpha_i} \}$ until it's a BFS.

(ii) an optimal solution \rightarrow an optimal solution also BFS: same as part (i)

the same α ($* G^T \alpha_B$ must be zero, otherwise, either $\bar{x}_B + \varepsilon \alpha$ or $\bar{x}_B - \varepsilon \alpha$ with $\varepsilon > 0$
($\alpha_B = \alpha$) is better.) Repeat the step until it's a BFS.

3) small techs when using simplex tableau

whether in two-phase method or slacks cases,

\bar{A}_B^\top can be gotten by

$\bar{A}_B^\top Q = M$ in simplex tableau

\rightarrow diagonal matrix which can be chosen

$$\bar{A}_B^\top = M D^\top$$

from original A

Lecture 9 The Interior Point Method

- The uniform expression of LP algorithms

Optimal conditions for LP

- 1) Primal feasibility
- 2) Dual feasibility
- 3) Complementarity conditions

\Rightarrow Simplex Method
 Maintain 1 & 3), find 2) when optimal
 $(x_i \leftarrow \underbrace{c_i - A^T y_j}_{(A^T)^T c_i}, y = \underbrace{(A^T)^T c_i}_{(A^T)^T c_i})$

Dual-Simplex Method

Maintain 2) & 3), find 1) when optimal

The Interior-point Method: Maintain 1) & 2), find 3) when optimal

(characteristics: start from the interior point of the feasible set)

- IDEA & Properties of the Interior Point Method

High-level idea: find x, y

$$\begin{cases} Ax = b, x \geq 0 \\ A^T y + s = c, s \geq 0 \\ x_i s_i \leq \mu_i, \forall i \end{cases}$$

$\mu \geq 0$ — complementarity gap
 keep decreasing μ — a solution for LP

Complexity $\approx O(n^{3.5})$

Several variants (primal-dual for this class)

* Theorem (Quality of Solutions)

The interior point method will always find the opt. solution with the maximum possible # non-zeros. (high-rank solutions)

Software choice:

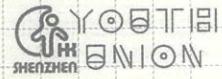
- MATLAB — linprog to choose method
- CVX — interior pt. method
- Excel — simplex method

- Tutorial — a theorem of alternative

(Gordan's Theorem) (a) $Ax \geq 0$ & (b) $A^T y = 0, y \geq 0$

exactly one is feasible. ($a \neq b \Leftrightarrow a \geq b$ but $a \neq b$)

$Ax \geq 0$ feasible $\Leftrightarrow Ax \geq 1$ feasible (*Note: $Ax \geq 0 \nRightarrow Ax \geq 1$)
 (when $Ax = S$, $A \frac{x}{S} = 1$)



$$(P) \min c^T x \quad | \quad (D) \max u^T y \\ \text{subject to } Ax \geq 1 \quad | \quad \text{subject to } A^T y = 0$$

$(y \not\geq 0) \rightarrow$ change this constraint
 or transform to $y \geq 0$

use this for the duality

(P) feasible \Rightarrow (D) feasible
 $Ax \geq 0 \Rightarrow$ with $u^T y^* = 0 \Rightarrow y^* = 0 \Rightarrow A^T y = 0, y \geq 0$ infeasible
 feasible

$Ax \geq 0$ infeasible \Rightarrow (P) infeasible by strong/weak duality
 (D) unbounded $\Rightarrow y \geq 0, A^T y = 0$ feasible
 But (D) feasible ($y = 0$)

(Same techs can be used for "Farka's Lemma": (a) $Ax = b, x \geq 0$; (b) $A^T y \geq 0, b^T y < 0$)
 exactly one holds

- Techniques - (more exercise)

① Use duality theory to transform a robust LP to a common one.

$$\text{maximize } c^T x$$

$$\text{subject to } a^T x \leq b + u^T x, \forall u \in U \cup V \leq w$$

"tech" $x \geq 0$

for $u^T x$, we have $a^T x \leq b + \underbrace{\min_{U \cup V} u^T x = x^T u}_{\text{duality theory}}$

Answer

$$\begin{aligned} &\Rightarrow \text{maximize } c^T x \\ &\text{subject to } a^T x \leq b + w^T y \\ &\quad A^T y = x \\ &\quad x, y \geq 0 \end{aligned}$$

$$\begin{aligned} &\text{maximize } w^T y \\ &\text{subject to } A^T y = x \\ &\quad y \leq 0 \end{aligned}$$

parts of the question

② (Duality in Chebyshev approximation)

Consider minimizing $\|Ax - b\|_\infty$ over all $x \in \mathbb{R}^n$, let v be the value of optimal cost.

(a) let p be any vector satisfying $\sum_{i=1}^m |p_i| \leq 1$ & $p^T A = 0^T$. Show that $p^T b \leq v$

(b) Show that the optimal cost in this problem is v .

$$\text{maximize } p^T b$$

$$\text{subject to } p^T A = 0^T$$

$$\sum_{i=1}^m |p_i| \leq 1$$

Sol. (a) $-v \leq Ax^* - b \leq v$, therefore $p^T(Ax^* - b) \leq v \sum_{i=1}^m |p_i|$

& $p^T(Ax^* - b) \geq -v \sum_{i=1}^m |p_i|$ (x^* \Rightarrow the optimal solution)

$$\therefore p^T b \leq v \sum_{i=1}^m |p_i| \leq v$$

(b) * By observation, we can suppose that the original LP & the LP in (b.) are primal-dual.

(check the assumption/hypothesis = (a) - weak duality
 min/max relationship strong duality ")

proof: the primal minimize $\|Ax-b\|_0$

$$\Leftrightarrow \underset{x,t}{\text{minimize}} \quad \underset{\text{see tutorial-1}}{\text{minimize}} [0^T, t] \begin{bmatrix} x \\ t \end{bmatrix}$$

subject to $Ax-b \geq t \mathbb{1}$ transform subject to $[A \mathbb{1}] \begin{bmatrix} x \\ t \end{bmatrix} \geq b \rightarrow \geq$
 $Ax-b \leq t \mathbb{1}$ $[A-\mathbb{1}] \begin{bmatrix} x \\ t \end{bmatrix} \leq b \rightarrow \leq$
 $t \geq 0$ $x \text{ free}$ $t \geq 0 \rightarrow \leq$

Take the dual

$$\Rightarrow \underset{\text{subject to}}{\text{maximize}} [b^T, b^T] \begin{bmatrix} y \\ z \end{bmatrix}$$

\star substitution \downarrow abs(p) in MATLAB

 $\begin{bmatrix} A^T & A^T \end{bmatrix} \begin{bmatrix} y \\ z \end{bmatrix} = 0$
 $y = \frac{1}{2}(1p + p) \geq 0$
 $[1^T - 1^T] \begin{bmatrix} y \\ z \end{bmatrix} \leq 1$
 $z = \frac{1}{2}(-1p + p) \leq 0$
 $y \geq 0, z \leq 0$
 $\Rightarrow \underset{\sum_i^m (p_i) \leq 1}{\text{maximize}} p^T b$

subject to $A^T p = 0$ (i.e. $p^T A = 0^T$)

Lecture 10 NLP - Intro

• Background Knowledge:

Gradient: $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ (x is an $n \times 1$ column)

$$\nabla f(x) = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right) - \text{an } n \times 1 \text{ column vector}$$

First-order Taylor Expansion:

$$f(x+td) = f(x) + t \nabla f(x)^T d + o(t)$$

$(x, d \in \mathbb{R}^n, t \in \mathbb{R})$ $\underset{t \rightarrow 0}{\text{when}}$ (compared with when $x \in \mathbb{R}^n$)

Hessian Matrix of f (when f is twice partial differentiable):

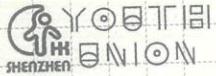
$$\nabla^2 f(x) = \left(\frac{\partial^2 f}{\partial x_i \partial x_j} \right)_{ij} \quad (\text{an } n \times n \text{ matrix with } \frac{\partial^2 f}{\partial x_i \partial x_j})$$

Second-order Taylor Expansion:

$$f(x+td) = f(x) + t \nabla f(x)^T d + \frac{1}{2} t^2 d^T \nabla^2 f(x) d + o(t^2)$$

$(x, d \in \mathbb{R}^n, t \in \mathbb{R})$ $(t \rightarrow 0)$

• Optimality Conditions for Unconstrained NLP



★ First-Order Necessary Conditions (FONC)

If x^* is a local minimizer of the unconstrained problem

$$\min_{x \in \mathbb{R}^n} f(x), \text{ then we must have } \nabla f(x^*) = 0 \quad \text{critical point / stationary point}$$

Proof. Consider the First-order Taylor Expansion

$$f(x+td) = f(x) + t \nabla f(x)^T d + o(t), \quad t \rightarrow 0$$

Suppose $\nabla f(x) = 0$, then let $d = -\nabla f(x)$

$$f(x+td) = f(x) - t(\|\nabla f(x)\|^2 - \frac{o(t)}{t}), \quad t \rightarrow 0$$

$$\text{so } \lim_{t \rightarrow 0} f(x+td) = \lim_{t \rightarrow 0} f(x) - t(\|\nabla f(x)\|^2 - \frac{o(t)}{t}) \leq f(x) \text{ contradictory!}$$

★ Example: Least-Square approximation (how to get $A^T A \hat{x} = A^T b$ in linear algebra)

Given m points $\{(x_{i1}, x_{i2}, \dots, x_{in}, y_i) \mid i=1, 2, \dots, m\}$

($m > n$) want $\beta \in \mathbb{R}^n$ such that $y \approx X\beta$ (with Least-square)

$$\begin{aligned} \text{minimize}_{\beta} \quad & \sum_{i=1}^m (y_i - \sum_{j=1}^n \beta_j x_{ij})^2 = \|y - X\beta\|^2 \\ & = y^T y - 2y^T X\beta + X\beta^T X = \beta^T X^T X \end{aligned}$$

Then, by FONC, we have $\nabla f(\beta) = -2X^T y + 2X^T X\beta = 0$

(if $f(x) = C^T x$, $\nabla f(x) = C$; if $f(x) = x^T M x$, $\nabla f(x) = 2Mx = 2M^T x$)

$\Rightarrow X^T X \beta = X^T y$ (least-square fitting with line) (candidates of local minimizer)

★ Second-Order Necessary Conditions (SONC)

If x^* is a local minimizer of f , then it holds that:

(1) $\nabla f(x^*) = 0$; (2) For all $d \in \mathbb{R}^n$: $d^T \nabla^2 f(x^*) d \geq 0$ (or $\nabla^2 f(x^*)$ is PSD)

{ Semidefiniteness: We call a (symmetric) matrix A positive/negative semidefinite (PSD/NSD)
iff $\forall x$, we have $x^T A x \geq 0 \leq 0$)

Proof: Consider the Second-order Taylor Expansion

$$f(x+td) = f(x) + t \nabla f(x)^T d + \frac{1}{2} t^2 d^T \nabla^2 f(x) d + o(t^2) \quad t \rightarrow 0$$

$$\stackrel{\text{FONC}}{=} f(x) + \frac{1}{2} t^2 d^T \nabla^2 f(x) d + o(t^2)$$

$$= f(x) + \underline{\frac{1}{2} d^T \nabla^2 f(x) d} + o\left(\frac{t^2}{t^2}\right) \quad t \rightarrow 0$$

If $d^T \nabla^2 f(x)d \geq 0$ does not hold, (i.e., $d^T \nabla^2 f(x)d < 0$,

we have $f(x+td) = f(x) + \frac{1}{2}t^2 \left(d^T \nabla^2 f(x)d + \underbrace{\frac{\partial t^2}{\partial}}_{\geq 0} \right) \leq f(x)$
contradictory!

Additions.

★ Positive Semidefinite Matrices (PSD)

1) If a matrix A is not symmetric, use $\frac{1}{2}(A+A^T)$ to define.

(Because $X^T A X = \frac{1}{2}X^T A X + \frac{1}{2}X^T A^T X = \frac{1}{2}X^T (A+A^T)X$
 $\Downarrow X^T (X^T A^T)^T = X^T A X$)

2) A symmetric matrix is PSD \iff all the eigenvalues are non-negative.

3) For any matrix A , $A^T A$ is PSD. (Because $X^T A^T A X = \|Ax\|^2 \geq 0$)

(If $f(x) = X^T M X$ with symmetric M , then $\nabla^2 f(x) = 2M$.)

★ Second-order Sufficient Conditions (SOSC)

If x^* satisfies: (with twice partial differential function f)

(1) $\nabla f(x^*) = 0$; (2) For all $d \in \mathbb{R}^n \setminus \{0\}$: $d^T \nabla^2 f(x^*)d \geq 0$ (or $\nabla^2 f(x^*)$ is PD)

(Definite matrices:

We call a symmetric matrix A positive/negative definite (PD/ND)

$\iff \forall x \neq 0: X^T A X \geq 0 / < 0$

Same as above. A symmetric matrix is PD. \Leftrightarrow all its eigenvalues positive

Proof. Lemma (Rayleigh Quotient)

Let $A \in \mathbb{R}^{n \times n}$ & symmetric, then

$$\lambda_{\min}(A) \|x\|^2 \leq X^T A X \leq \lambda_{\max}(A) \|x\|^2, \quad \forall x \in \mathbb{R}^n$$

the smallest eigenvalue of A the largest eigenvalue of A

By Taylor's Expansion

$$f(x^*+h) = f(x^*) + \frac{1}{2}h^T \nabla^2 f(x^*)h + o(\|h\|^2) \quad h \rightarrow 0$$

$$\Rightarrow f(x^*+h) \geq f(x^*) + \frac{1}{2}\underbrace{\|h\|^2}_{h^T \nabla^2 f(x^*)h} + o(\|h\|^2)$$

$h^T \nabla^2 f(x^*)h \geq \lambda_{\min}(A) \|h\|^2$

$$\Rightarrow f(x^* + h) \geq f(x^*) + \|h\|^2 \left(\frac{\mu}{2} + \frac{o(\|h\|^2)}{\|h\|^2} \right), \quad \|h\| \rightarrow 0$$

so $f(x^* + h) > f(x^*)$ (Because $\lim_{\|h\| \rightarrow 0} \frac{o(\|h\|^2)}{\|h\|^2} = 0$)

* For maximization:

$$FONC \rightarrow \nabla f(x^*) = 0$$

$$SONC \rightarrow \text{NSD}$$

If $\nabla^2 f(x^*)$ is indefinite, with $\nabla f(x^*) = 0$,

$$SOSC \rightarrow \begin{cases} ND & \text{all EVs} \leq 0 \\ & \text{all EVs} < 0 \end{cases} \quad x^* \text{ is a saddle point.}$$

Lecture 11 Optimality Conditions for Constrained Problems

- Existence of Solutions

* (Weierstrass Theorem) (Extreme value)

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuous function & let $S \subseteq \mathbb{R}^n$ be bounded, closed & non-empty set. Then, f attains a global maximum and minimum on S .

Def: (1) closed: For every convergent sequence x_k with $x_k \in S, \forall k$
 then $\lim_{k \rightarrow \infty} x_k \in S$ (condition: $\lim_{k \rightarrow \infty} x_k = x$ (convergent))

(2) bounded: There is a $B > 0$ with $\|x\| \leq B, \forall x \in S$ (no infinity ones)

(3) Compact (a closed & bounded set)

* Coercivity (for unconstrained problems):

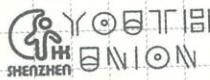
A continuous function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is coercive if

$$\lim_{\|x\| \rightarrow \infty} f(x) = +\infty \quad (\text{e.g. } x^2, x^4, |x| - \text{coercive})$$

(Trick - finding suitable lower bound for $x, e^x, x^3 - \text{non-coercive}$
 sufficiently large x)

* Theorem (coercivity & existence of solutions)

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuous & coercive function. Then, $\forall \alpha > 0$,
 the level set $L_{\leq \alpha} := \{x \in \mathbb{R}^n \mid f(x) \leq \alpha\}$ is compact & f has at least
 one global minimizer.



• Optimality Conditions for Constrained problems

Def.

(1) Feasible Directions

(1) **Feasible Direction**
 Given $x \in S_2$, we call d a feasible direction at x ($d \in \mathbb{R}^n$), if $\exists \bar{\epsilon} > 0$ s.t.

(e.g. (i) $\Omega = \{x \mid Ax = b\}$, then all d's: $\{d \mid Ad = 0\}$)
 (ii) $\Omega = \{x \mid Ax \geq b\}$, then all d's: $\{d \mid ad \geq 0 \text{ s.t. } a^T x = b\}$)
 (inactive $a^T x \geq b$, sufficiently small ϵ can satisfy)

(2) Descent Direction

Let f be continuously differentiable. Then, d is called a descent direction at x iff $\nabla f(x)^T d < 0$.

(If d is a descent direction at x , then $\exists \bar{r} > 0$, s.t. $f(x + rd) < f(x)$, $\forall 0 < r < \bar{r}$)
 (by Taylor Expansion)

★ FONC for Constrained problem

Let x^* be a local minimum of $\min_{x \in S} f(x)$, then for any feasible direction d at x^* , we must have $Df(x^*)d \geq 0$.

(Proof: $f(x^* + td) = f(x^*) + t \nabla f(x^*)^T d + o(t)$, $t \rightarrow 0$
 $\quad \quad \quad (t \in \mathbb{E})$ assume $\nabla f(x^*)^T d < 0$, then $f(x^* + td) < f(x^*)$ contradict)

With (1) & (2), FONC can be written as $S_d(x^*) \cap S_0(x^*) = \emptyset$

• General Optimal Conditions

(General Non-linear Optimization Problem) (Implicit open constraints
 minimize $f(x)$ can be omitted in Lagrangian
 $x \in R^n$ function)
 subject to $g_i(x) \leq 0, \forall i=1,2,\dots,m$

The feasible set $\Omega = \{x \in \mathbb{R}^n | g(x) \leq 0, h(x) = 0\}$

Def. Active Set at a point $x \in \Omega$, the set $A(x) = \{i | g_i(x) = 0\}$ denotes the set of active constraints.

◎ YOUNION
SHENZHEN

Inactive set at a point $x \in \Omega$, the set $I(x) = \{i | g_i(x) < 0\}$ denotes the set of inactive constraints.

* FONC for linearly constrained problem $\begin{cases} \min_x f(x) \\ \text{subject to } Ax \geq b \end{cases}$

If x^* is a local minimum of the problem,

the \exists some $y \in \mathbb{R}^m$ with $\begin{bmatrix} y \geq 0; \nabla f(x^*) - A^T y = 0; y_i(a_i^T x^* - b_i) = 0 \forall i \end{bmatrix}$
gradient of Lagrangian function

Proof.

$$\left(\begin{array}{l} \min_d \nabla f(x^*)^T d \\ \text{subject to } a_i^T d \geq 0, \forall i \in A(x^*) \end{array} \right) \geq 0 \Leftrightarrow \text{FONC: } a_i^T d \geq 0 \quad \forall i \in A(x^*)$$

$$\uparrow \quad \left(\begin{array}{l} \min_d \nabla f(x^*)^T d \\ \text{subject to } C^T d \geq 0 \quad \text{with } C = \begin{bmatrix} -c_1^T \\ -c_2^T \\ \vdots \\ -c_n^T \end{bmatrix}, \quad c_i^T = \begin{cases} a_i^T, & i \in A(x^*) \\ 0, & i \in I(x^*) \end{cases} \end{array} \right) \geq 0$$

Take the dual problem (D) $\max_t \mathbf{0}^T t = 0$
Subject to $C^T t = \nabla f(x^*)$

By weak duality $t \geq 0$

$(P) \geq 0 \Leftrightarrow (D)$ feasible, finite opt. 0

$\Leftrightarrow t \geq 0, A^T t = \nabla f(x^*), t_i(a_i^T x^* - b_i) = 0 \quad \forall i$
complementarity, let $t_i = 0, \forall i \in I(x^*)$

* Special case, change $Ax \geq b$ to $Ax = b$: (if x^* is a local minimum)

$\exists y \in \mathbb{R}^n$ with $A^T y = \nabla f(x^*)$



Consider the non-linear program

$$\underset{x \in \Omega^n}{\text{minimize}} \quad f(x)$$

subject to $g_i(x) \leq 0, \forall i = 1, 2, \dots, m$

(for $h(x) = 0 \Leftrightarrow \begin{cases} h(x) \leq 0 \\ -h(x) \leq 0 \end{cases}$ substitute)

(f.g continuous differentiable)

Lemma - No Feasible Descent for Inequality constraints

Let x^* be a local minimum of above, there does not exist a vector

$d \in \mathbb{R}^n$ s.t. $\nabla f(x^*)^T d < 0 \quad \& \quad \nabla g_i(x^*)^T d < 0 \quad \forall i \in I(x^*)$

Why $\nabla g_i^T(x^*)d < 0$ implies feasibility? ($i \in A(x^*)$)

Reason = Taylor's expansion $g_i(x^*+td) = g_i(x^*) + t \nabla g_i^T(x^*)d + o(t) \xrightarrow{t \rightarrow 0}$
 $\nabla g_i^T(x^*)d < 0$ indicates that $g_i(x^*+td) < g_i(x^*) = 0$

Proof: assume $\exists d$ s.t. $\nabla f(x^*)d < 0$, $\nabla g_i^T(x^*)d < 0$, $\forall i \in A(x^*)$

$$\Rightarrow f(x^*+td) < f(x^*) \xrightarrow{t \rightarrow 0} \& g_i(x^*+td) < g_i(x^*) = 0, \quad (t < \bar{t})$$

which means that x^*+td is a local minimum (\exists small enough \bar{t} , $f(x^*+\bar{t}d) < 0$)
contradictory!! ($f(x^*+td) < f(x^*)$, $t \in (0, \min\{\bar{t}, \bar{t}\})$)

Lecture 12 KKT Conditions

• Fritz-John Conditions (FJ Conditions) — another form of FONC

Let x^* be a local minimum of $\left(\begin{array}{c} \text{min}_x f(x) \\ \text{subject to } g_i(x) \leq 0, \forall i=1, \dots, m \end{array} \right)$, then $\exists \lambda_0, \lambda_1, \dots$

$\lambda_m \geq 0$, and are not all zeros s.t.

$$\text{gradient of Lagrangian function} \quad \lambda_0 \nabla f(x^*) + \sum_{i=1}^m \lambda_i \nabla g_i(x^*) = 0$$

$$\lambda_i g_i(x^*) = 0 \quad \forall i=1, 2, \dots, m$$

complementarity conditions

Simple Proof: By the lemma above,

$$\nexists d: \nabla f(x^*)d < 0, \quad \nabla g_i^T(x^*)d < 0, \quad \forall i \in A(x^*)$$

which means $\left(\begin{array}{c} \text{minimize } t \\ \text{subject to } \nabla f(x^*)d \leq t, \quad \nabla g_i^T(x^*)d \leq t, \quad \forall i \in A(x^*) \end{array} \right) \geq 0$

Because t can be only chosen to be ≥ 0 to promise feasible of d .

$$\Leftrightarrow \left(\begin{array}{c} \text{minimize } t \\ \text{subject to } \nabla f(x^*)d \leq t, \quad (Cd \leq t) \end{array} \right) \geq 0$$

$$\text{with } C = \begin{bmatrix} -\nabla f^T \\ \vdots \\ -\nabla g_i^T \end{bmatrix}, \quad C^T = \begin{cases} \nabla g_i^T, & i \in A(x^*) \\ 0, & i \in I(x^*) \end{cases}$$

Take the dual

$$\text{maximize } 0 \quad \text{subject to } \lambda_0 \leq 0, \quad \lambda \leq 0$$

$$\begin{aligned} \lambda_0 \nabla f(x^*) + C^T \lambda &= 0 \\ -\lambda_0 - \lambda &= 1 \end{aligned}$$

By weak duality, the dual must be feasible.
use $-\lambda_0, \lambda$ to substitute for λ_0, λ , respectively.

$$\lambda_0 \geq 0, \lambda \geq 0 \iff \lambda_0, \lambda_1, \dots, \lambda_m \geq 0$$

$$\left\{ \begin{array}{l} \lambda_0 \nabla f(x^*) + C^\top \lambda = 0 \\ \lambda_0 \nabla f(x^*) + \sum_{i=1}^m \lambda_i g_i(x^*) = 0 \end{array} \right. \quad \left\{ \begin{array}{l} \lambda_0 \nabla f(x^*) + \sum_{i=1}^m \lambda_i g_i(x^*) = 0 \\ \lambda_i g_i(x^*) = 0, \forall i=1,2,\dots,m \end{array} \right. \text{ (make } \lambda_i = 0, \forall i > 2(x^*))$$

$$\lambda_0 + \lambda \cdot \lambda = 1 \iff \text{cannot be all zeros}$$

Defects of FJ conditions: λ_0 can be zero — may have too many pts not local minimum (when $\lambda_0=0$, no contacts with $\nabla f(x)/f(x)$).

↓ make it more precise

- Karush-Kuhn-Tucker Conditions (KKT conditions)

— another form of FONC

Let x^* be a local minimum of $\min_x f(x)$
subject to $g_i(x) \leq 0, \forall i=1, \dots, m$

& suppose vectors $\{\nabla g_i(x^*) | i \in I(x^*)\}$ are linearly independent.

Then, $\exists \lambda_1, \dots, \lambda_m \geq 0$ s.t. $\nabla f(x^*) + \sum_{i=1}^m \lambda_i \nabla g_i(x^*) = 0$
 $\lambda_i g_i(x^*) = 0, \forall i=1, \dots, m$

Proof. Use FJ conditions, $\lambda_0 \neq 0$ because
of linear independence, then divide each side by λ_0 .

- KKT Conditions (General Setting)

Background: (1) general optimization problem form:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x)$$

$$\text{subject to } g_i(x) \leq 0, \forall i=1, 2, \dots, m$$

$$h_j(x) = 0, \forall j=1, \dots, p$$

(2) Lagrangian function

$$L(x, \lambda, \mu) := f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{j=1}^p \mu_j h_j(x)$$

(λ_i, μ_j — can be seen as
dual variables)

$$(f(x) + \lambda^\top g(x) + \mu^\top h(x))$$

★ General KKT Conditions. If x^* is a local minimizer of a constraint qualification then $\exists \lambda \& u$ s.t.

(1) Main Condition

$$\nabla L(x, \lambda, u) = \nabla f(x^*) + \sum_{i=1}^m \lambda_i \nabla g_i(x^*) + \sum_{j=1}^p u_j \nabla h_j(x^*) = 0$$

(con x^*)

(2) Complementarity Conditions

$$\lambda_i g_i(x^*) = 0, \forall i=1, \dots, m ; \underbrace{\mu_j h_j(x^*) = 0, \forall j=1, \dots, p}_{\text{often omitted}}$$

(3) Primal Feasibility

$$g_i(x^*) \leq 0, \forall i=1, \dots, m ; \underbrace{h_j(x^*) = 0, \forall j=1, \dots, p}_{\text{often omitted}}$$

(4) Dual Feasibility

$$\lambda_i \geq 0, \forall i=1, \dots, m ; \underbrace{\mu_j \text{ free}, \forall j=1, \dots, p}_{\text{often omitted}}$$

→ Constraint Qualifications:

1) Linearly Independent Constraint Qualification:

$$(LICQ) \quad \left\{ \nabla g_i(x^*) \mid i \in I(x^*) \right\} \cup \left\{ \nabla h_j(x^*) \mid j = 1, \dots, p \right\}$$

are linearly independent

2) Others: ACQ, GCQ, MFCQ, PLICQ, Slater's condition.

→ Comments

(i) KKT is a kind of FONC, for general constrained optimization problems.

(ii) KKT point / stationary point = a point satisfies KKT conditions.

★ Second-Order Conditions (KKT form) for constrained problem

1) SONC Let x^* be a regular point & local min, then

(1) The KKT conditions hold.

(2) $d^T \nabla_x^2 L(x^*, \lambda, u) d \geq 0, \forall d \in C(x^*)$
(PSD on the critical cone)

Comments:

(i) Hessian of the Lagrangian function

$$\nabla_{xx}^2 \mathcal{L}(x^*, \lambda, \mu) = \nabla^2 f(x) + \sum_{i=1}^m \lambda_i \nabla^2 g_i(x) + \sum_{j=1}^p \mu_j \nabla^2 h_j(x)$$

(ii) Critical cone ($C(x^*)$):

$$C(x) := \{d \in \mathbb{R}^n \mid \nabla f(x)d = 0, \nabla g_i^\top(x)d \leq 0, \forall i \in A(x), \nabla h_j^\top(x)d = 0, \forall j\}$$

2) **SOSC**: Let x^* be a KKT point with multiplier λ, μ

(i.e.) (i) KKT conditions work.

$$(2) d^\top \nabla_{xx}^2 \mathcal{L}(x^*, \lambda, \mu) d \geq 0, \forall d \in C(x^*) \setminus \{0\}$$

(PD on the critical cone)

Then, x^* is a strict local minimizer

Lecture 13 Convexity

- Definition (convex & concave)

(1) Convex Function

A convex function f on a convex set Ω satisfies:

for every $x_1, x_2 \in \Omega$, and any $\lambda \in [0, 1]$

$$f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$$

(2) Concave Function

A concave function f on a convex set Ω satisfies:

$$(-f \text{ is convex}) \text{ or } f(\lambda x_1 + (1-\lambda)x_2) \geq \lambda f(x_1) + (1-\lambda)f(x_2)$$

- Test: ① Convexity via Hessian / Second-order condition

Let f be twice differentiable, then f is convex on Ω iff its Hessian Matrix is PSD (i.e.) $d^\top \nabla^2 f(x) d \geq 0, \forall d \in \mathbb{R}^n, x \in \Omega$

(PD - strictly convex)

- ② First-order condition

If $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable, then f is convex on its convex domf iff

$$f(y) \geq f(x) + \nabla f(x)^\top (y-x), \quad \forall x, y \in \text{dom} f \quad (> - \text{strictly convex})$$

Proof: ② suppose $f(x)$ is convex,

$$f(\theta y + (1-\theta)x) \leq \theta f(y) + (1-\theta)f(x), \quad \forall \theta \in [0,1]$$

$$\Leftrightarrow f(x + \theta(y-x)) - f(x) \leq \theta(f(y) - f(x))$$

$$\Rightarrow \lim_{\theta \rightarrow 0} \frac{f(x + \theta(y-x)) - f(x)}{\theta} \leq f(y) - f(x)$$

$$\nabla f(x)^T (y-x) \text{ definition}$$

Suppose first-order condition satisfied,

$$\text{Let } z = \theta x + (1-\theta)y, \text{ then } f(x) \geq f(z) + \nabla f(z)^T (x-z)$$

$$\forall \theta \in [0,1]$$

$$f(y) \geq f(z) + \nabla f(z)^T (y-z)$$

$$\Rightarrow \theta f(x) + (1-\theta)f(y) \geq f(z) \text{ (sum)}$$

① Suppose $f(x)$ is convex, by Taylor expansion

$$f(x+td) = f(x) + t \nabla f(x)^T d + \frac{1}{2} t^2 d^T \nabla^2 f(x) d + O(t^3) \quad t \rightarrow 0$$

$$\text{By First-order condition } f(x+td) \geq f(x) + \nabla f(x)^T (td)$$

$$\Rightarrow \frac{1}{2} t^2 \left(d^T \nabla^2 f(x) d + \frac{O(t^3)}{t^2} \right) \geq 0 \Rightarrow d^T \nabla^2 f(x) d \geq 0 \quad (\nabla^2 f(x) \succeq 0)$$

Suppose $\nabla^2 f(x) \succeq 0$, then for any $x, y \in \text{dom } f$

$$f(y) = f(x) + \nabla f(x)^T (y-x) + \underbrace{\frac{1}{2} (y-x)^T \nabla^2 f(x) (y-x)}_{\frac{1}{2} (y-x)^T \nabla^2 f(x) (y-x)} + O(\|y-x\|^3)$$

(Taylor's theorem)

$$\Rightarrow f(y) \geq f(x) + \nabla f(x)^T (y-x), \text{ by First-order condition}$$

z between $x, y \quad (z = \theta x + (1-\theta)y)$

- Operations that preserve convexity

1) **Nonnegative weighted sums**: If $a_1, \dots, a_m \geq 0$, f_1, \dots, f_m are convex functions

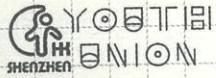
the $\sum_{i=1}^m a_i f_i$ is still convex.

2) **Composition with an affine mapping**: Suppose $f: \mathbb{R}^m \rightarrow \mathbb{R}$ is convex,

$A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$ given, define $g: \mathbb{R}^n \rightarrow \mathbb{R}$, $g(x) = f(Ax+b)$
is convex.

3) Positive maximum & supremum: (piece-wise - convex)

If f_1, \dots, f_m are convex functions, then $f(x) = \max\{f_1(x), \dots, f_m(x)\}$ is still convex. (Can be extended to infinitely/uncountably many)



• Convexity & Optimality

(1) Local & Global minimizer:

Let $f: \Omega \rightarrow \mathbb{R}$ be a convex function & $\Omega \subset \mathbb{R}^n$ be a convex set. Then any local minimizer of problem $\left(\begin{array}{c} \min_x f(x) \\ \text{subject to } x \in \Omega \end{array} \right)$ is also a global minimizer.

Proof: suppose x is not a global minimizer. $\exists y$ s.t. $f_0(y) < f_0(x)$, $\|y-x\|_2 > R$ (because $f_0(x) = \inf\{f_0(z) | z \text{ feasible}, \|z-x\|_2 \leq R\}$)

Construct $z = (1-\theta)x + \theta y$, $\theta = \frac{R}{2\|y-x\|_2}$, then $z = x + \theta(y-x)$

$\|z-x\|_2 = \theta \|y-x\|_2 = \frac{R}{2} < R$, by convexity of f_0

$f_0(z) \leq (1-\theta)f_0(x) + \theta f_0(y) < f_0(x)$ contradict!

(2) Stationary & Global minimizer:

Let f be the convex & suppose $\Omega := \{x | g_i(x) \leq 0, h_j(x) = 0\}$ is a convex set

Then the KKT conditions for $\left(\begin{array}{c} \min_x f(x) \\ \text{subject to } x \in \Omega \end{array} \right)$ are sufficient for global optimality

Proof: Firstly, consider the reduced case.

(Let $S \subset \mathbb{R}^n$ be an open convex set. Suppose that $f: S \rightarrow \mathbb{R}$ is a convex on S & continuously differentiable at \bar{x} , then \bar{x} is a global min iff $\nabla f(\bar{x}) = 0$)

\Rightarrow proof of this proposition: (i) sufficient part, \bar{x} is a global min $\Rightarrow \bar{x}$ is a local min

(ii) necessary part, by convexity & first-order conditions $\Rightarrow \bar{x}$ satisfies FONC $\Rightarrow \nabla f(\bar{x}) = 0$

$$f(y) \geq f(\bar{x}) + \nabla f(\bar{x})(y - \bar{x}), \forall y$$

Then, use the proposition

$$f(x) \geq f(\bar{x}) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{j=1}^p \mu_j h_j(x) \stackrel{\text{complementarity}}{=} \min_x \{f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{j=1}^p \mu_j h_j(x)\}, \quad (L(x, \bar{x}, \bar{\mu}))$$

$$\stackrel{\substack{\text{local/global} \\ \text{relation}}}{\leq} \min_{\substack{g_i(x) \leq 0 \\ h_j(x) = 0}} \{f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{j=1}^p \mu_j h_j(x)\} \leq \min_{x \in \Omega} f(x)$$

(3) A Lemma — how to recognize some convex sets?

Let f be a convex/concave function, then for any c , the level set $L \leq c = \{x | f(x) \leq c\} / L \geq c = \{x | f(x) \geq c\}$ is a convex set.

(Tricks: transformations of monotone/variable substitutions)

Lecture 14 Algorithms for Unconstrained Problems

- General ideas of algorithms:

The nature of optimization algorithms is **iterative procedures**.

{ starting from some point x_0 , generate a sequence $\{x_k\}$

{ terminates when attain / satisfactory, every step $f(x_{k+1}) < f(x_k)$

hopefully, $\{x_k\}$ converges to a local min x^*

Def.(convergence). Let $\{x_k\}$ be a sequence of real vectors. Then $\{x_k\}$ converges to x^* iff $\forall \varepsilon > 0, \exists$ a pos. integer K s.t. $\|x_k - x^*\| < \varepsilon, \forall k \geq K$.

- Algorithm for single variable problems:

★ **Bisection Method (binary search)** — $g(x) = f'(x)$, then, need x^* s.t $g(x^*) = 0$

Start from x_l, x_r s.t. $g(x_l) < 0 \& g(x_r) > 0$

(1) Define $x_m = \frac{1}{2}(x_l + x_r)$

(2) If $g(x_m) = 0$, output x_m

(3) Otherwise { $g(x_m) > 0$, let $x_r = x_m$
 $g(x_m) < 0$ let $x_l = x_m$

(4) If $|x_r - x_l| < \varepsilon$, stop & output $\frac{1}{2}(x_l + x_r)$, otherwise iterate.

} usage - find
 approx. stationary / critical
 points.

★ **Golden Section Method** — do not need to use $f'(x)$

(unimodal — f has only one critical point)

Start from $[x_l, x_r]$ (suppose min in $[x_l, x_r]$). $\phi \in (0, 0.5)$

(1) Set $x_l' = \phi x_r + (1-\phi)x_l, x_r' = (1-\phi)x_r + \phi x_l$ ($x_r' > x_l'$)

(2) If $f(x_l') < f(x_r')$, which means $[x_l, x_r] \rightarrow [x_l, x_r'] \Rightarrow$ let $x_r = x_r'$

| If $f(x_l') > f(x_r')$, which means $[x_l, x_r] \rightarrow [x_l, x_r'] \Rightarrow$ let $x_l = x_l'$

(3) If $|x_r - x_l| < \varepsilon$, output $\frac{1}{2}(x_l + x_r)$, otherwise iterate.

Q - why golden section? / what is ϕ ?

Suppose we update $x_r = x_r'$, we want $x_e' = \text{new } x_r'$ for less calculation.

$$\rightarrow \phi x_r + (1-\phi) x_e = \phi x_e + (1-\phi)((1-\phi)x_r + \phi x_e)$$

$$\left\{ \begin{array}{l} (1-\phi)^2 = \phi \\ \phi(2-\phi) = 1-\phi \end{array} \right. \Rightarrow \frac{\phi^2 - 3\phi + 1}{= 0} \Rightarrow \boxed{\phi = \frac{3-\sqrt{5}}{2} \quad (\phi < 0.5)}$$

$$\Rightarrow 1-\phi = \frac{\sqrt{5}-1}{2} \approx 0.618 \text{ golden section!}$$

- High-dimensional problems - algorithms Descent Methods

(i) Gradient descent Method:

\Rightarrow IDEA: find search direction \rightarrow find step size & stopping criterion

$$\text{aim } x_{k+1} = x_k + \alpha_k d_k, f(x_{k+1}) < f(x_k) \text{ — main idea for all descent methods}$$

Descent Direction — $d \in \mathbb{R}^n$ s.t. $\nabla f(x)d < 0$ (to lower the function values) (gradient, steps)
Newton's

$$\rightarrow \text{to be definite, choose } -\nabla f(x) \text{ to be } d, \nabla f(x)(-\nabla f(x)) = \|\nabla f(x)\|^2 \leq 0$$

\Rightarrow (1) Initialization — select an initial point $x^0 \in \mathbb{R}^n$

(2) Pick a step size (constant / exact line search / backtracking) on the function

$$\phi(\alpha) = f(x_k - \alpha \nabla f(x_k)), \text{ get } \alpha_k \text{ tolerance}$$

(3) Set $x_{k+1} = x_k - \alpha_k \nabla f(x_k)$, when $\|\nabla f(x_{k+1})\| < \epsilon$ STOP. stopping criterion
otherwise iterate.

★ Choose step size α_k :

WAY I: $\alpha_k = \bar{\alpha}$ (commonly used for the sake of convenience)

WAY II: EXACT LINE SEARCH choose α_k s.t. $\alpha_k = \underset{\alpha > 0}{\operatorname{argmin}} f(x_k + \alpha d_k)$

(because $\underset{\alpha > 0}{\operatorname{argmin}} f(x_k + \alpha d_k)$ is a single variable function \rightarrow output argument which minimizes f)

\Rightarrow analytically get α

/ use golden section method

WAY III: BACKTRACKING / Armijo LINE SEARCH

Let $\alpha, r \in (0, 1)$, choose $\alpha_k \in \{1, \alpha, \alpha^2, \dots\}$ s.t. $f(x_k + \alpha_k d_k) - f(x_k) \leq r \alpha_k \nabla f(x_k)^\top d_k$
(as the largest element)

$\alpha_k \rightarrow$ finitely many step, starting from 1.

Armijo condition

"Feasibility": define $\hat{f}_k(\alpha) = f(x_k + \alpha d_k) - f(x_k) \Rightarrow \hat{f}'_k(\alpha) = \nabla f(x_k)^\top d_k$

Armijo condition $\Leftrightarrow \hat{f}_k(\alpha) \leq r \alpha \cdot \hat{f}'_k(0) \Leftrightarrow \frac{\hat{f}_k(\alpha)}{\alpha} \leq r \hat{f}'_k(0)$, by MVT, $\exists \delta \in (0, 1)$

* Convergence & properties of GD Method:

Global Convergence: Fixed critical points independent of the chosen initial point

Accumulation Point: a point x is an accumulation point of $\{x_k\}$ if for every $\epsilon > 0$, there are infinitely many numbers k with $x_k \in B_\epsilon(x)$

- Comments:
- 1) If x is an accumulation point of $\{x_k\}$, \exists a subsequence $\{x_{k_j}\}$ that converges to x .
 - 2) If $\{x_k\}$ converges to some $x^* \in \mathbb{R}^n$, it is the only unique accumulation point of $\{x_k\}$.
 - 3) A bounded sequence always possess at least an accumulation point.

* Theorem: Global Convergence

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be continuously differentiable, & let x_k be generated by gradient method for solving $\min_{x \in \mathbb{R}^n} f(x)$ (unconstrained) with exact / backtracking linear search

Then $f(x_k)$ is nonincreasing & every accumulation point of x_k is a stationary point of f .

Lipschitz Continuity: ∇f is Lipschitz Continuous over \mathbb{R}^n means

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|, \forall x, y \in \mathbb{R}^n$$

where $L > 0$ called Lipschitz constant. The class of functions with Lipschitz gradient with constant L is denoted by $C_L^{(0)}(\mathbb{R}^n) \xrightarrow{\text{continuously differentiable}} \text{Lipschitz continuous}$

(Comment: If f is twice differentiable, $f \in C_L^{(1)}(\mathbb{R}) \Leftrightarrow \|\nabla^2 f(x)\| \leq L, \forall x \in \mathbb{R}^n$)

$$\text{spectral norm} = \sqrt{\lambda_{\max}(H^T H)}$$

(largest singular value)

Linear Convergence: we say that x_k converges linearly

with rate $\gamma \in (0, 1)$ to $x^* \in \mathbb{R}^n$, if $\exists l > 0$

$$\text{s.t. } \|x_{k+l} - x^*\| \leq \gamma \|x_k - x^*\|, \forall k \geq l$$

why "linear"? \rightarrow rate/order of convergence \rightarrow { A sequence $\{x_n\}$ is said to have order of convergence $q \geq 1$ & rate of convergence μ

$$\text{if } \lim_{n \rightarrow \infty} \frac{|x_{n+1} - x^*|}{|x_n - x^*|^q} = \mu$$

* Theorem: Rates for convex problems

Let $f \in C_2^{++}$, & suppose $\exists \mu > 0$ s.t. $\mu \|d\|^2 \leq d^T \nabla f(x) d$ ($\forall d, \forall x$)
 Let x_k be generated by the gradient method (strong convexity) & let x^* be the solution of $\min_x f(x)$
 Then x_k converges linearly to x^* with rate $\eta = 1 - M^{-1}$

it follows $f(x_k) - f(x^*) \leq \eta^k (f(x_0) - f(x^*))$

$$\| \nabla f(x_k) \| \leq \sqrt{\frac{L}{\mu}} \eta^k \| \nabla f(x_0) \|, \quad \| x_k - x^* \| \leq \sqrt{\frac{L}{\mu}} \eta^k \| x_0 - x^* \|$$

(with $M = \begin{cases} \frac{2}{\alpha}(1 - \frac{L}{2\alpha}), & \text{constant step size} \\ \frac{1}{2L}, & \text{exact line search} \end{cases}$)

$\eta = \min\left\{1, \frac{2\alpha(1-\rho)}{L}\right\}$, backtracking ~

★ More properties: (GD Method)

When using exact line search, the directions between the consecutive steps are perpendicular. (i.e.) $(d^{k+1})^T d_k = 0$ (proof: FONC $\frac{\nabla f(x_k + \alpha d_k)^T d_k}{d^{k+1}^T} = 0$)

(ii) Newton's Method

Features: / converge much faster than the gradient descent method.
 required second-order informative
 more sensitive to the initial point

① in \mathbb{R} : Critical point x^* means $g(x^*) = 0$ (with $g'(x) = f'(x)$)

$$\Rightarrow x_{k+1} = x_k - \frac{g(x_k)}{g'(x_k)} = x_k - \frac{f(x_k)}{f'(x_k)}$$

a kind of descent method (Given $f''(x) = g''(x) \neq 0$ at each step)

★ Theorem: Convergence of Newton's Method (one dimension)

If g is twice cont. differentiable & x^* is the root of g where $g'(x^*) = 0$
 provided that $|x_0 - x^*|$ is sufficiently small, the sequence generated by Newton's method

will satisfy $|x_{k+1} - x^*| \leq C |x_k - x^*|^2$ (quadratic convergence)

with $C = \sup_x \frac{1}{2} \left| \frac{g''(x)}{g'(x)} \right|$ (converges much faster than gradient method)

Another interpretation of Newton's Method.

$f(x) \approx f(x_k) + f'(x_k)(x - x_k) + \frac{1}{2} f''(x_k)(x - x_k)^2$. (estimate by quadratic functions)

& get the argmin $\rightarrow x_{k+1} = x_k - \frac{f'(x_k)}{2f''(x_k)}$ (suppose $f''(x_k) > 0$)

② in \mathbb{R}^n : similar to \mathbb{R} $x_{k+1} = x_k - (\nabla^2 f(x_k))^{-1} \nabla f(x_k)$ search direction / Newton direction

When f is convex / $\nabla^2 f$ is pos. def., $-\nabla f(x) \nabla^2 f(x) \nabla f(x) < 0$ is the search descent direction

How to ensure Newton's Method's convergence with any initial point?

backtracking line search — use α_k as a step size.

⇒ Complete steps of Newton's Method

- (1) Initialization: select an initial point $x_0 \in \mathbb{R}^n$
- (2) Compute the Newton's direction d_k which satisfies $\nabla^2 f(x_k) d_k = -\nabla f(x_k)$
- (3) Choose a stepsize α_k by line search, calculating $x_{k+1} = x_k + \alpha_k d_k$
- (4) If $\|\nabla f(x_{k+1})\| \leq \varepsilon$, STOP.

★ Theorem: Convergence of Newton's Method (high/n dimension)

Let f be twice cont. diff. & x^* be the local minimizer of f . For some given $\varepsilon > 0$. assume that (i) f is strongly convex with constant m (i.e. $\nabla^2 f(x) \succeq mI$)

both or $\begin{cases} (i) \nabla^2 f \text{ is Lipschitz continuous for all } x, y \Leftrightarrow \nabla^2 f(x) \succeq mI \\ (ii) \exists B_\varepsilon(x^*) \text{ with constant } L, \text{ (i.e. } \|\nabla^2 f(x) - \nabla^2 f(y)\| \leq L \|x - y\| \text{ is pos. semidef.)} \end{cases}$

If $\{x_k\}$ is generated by Newton's Method, with assumptions above,

we have for $k = 0, 1, \dots$ $\|x_{k+1} - x^*\| \leq \frac{L}{2m} \|x_k - x^*\|^2$ (quadratic convergence)

& if $\|x^* - x^*\| \leq \frac{m}{L} \min\{1, \varepsilon\}$, we have $\|x_k - x^*\| \leq \frac{2m}{L} \left(\frac{1}{2}\right)^{2k}$

- Comparison between gradient method & Newton's method

Gradient Descent Method { less memory space, computational easy
converges slowly, more iterations } \rightarrow large-size problems

Newton's Method { converges fast, less iterations
more storage, computational cost } \rightarrow small-size problems

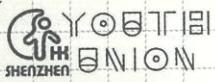
Advanced: algorithm in between — Quasi-Newton Methods

(approximate $\nabla^2 f(x)$)

algorithm for large-size ones — Stochastic gradient Methods

Lecture 15 Algorithms for Constrained Problems

(Intro)



- Ideas for the algorithm

Recall gradient descent method, what if $x_{k+1} = x_k + \alpha_k d_k$ outside the feasible region?
 \Rightarrow project it to the feasible region \rightarrow projected gradient method

- Euclidean / Orthogonal Projection

i) Def: $P_C(x_0) := \arg\min_x \{ \|x - x_0\| \mid x \in C \}$ (projection of x_0 onto a set (convex) C)

It can be seen as a minimization programming with unique minimizer x^*
 non-empty, closed & convex

$$\begin{aligned} &\min_x \frac{1}{2} \|x - x_0\|^2 \\ &\text{subject to } x \in C \end{aligned}$$

write $x^* = P_C(x_0)$

2) Special cases:

i) projection on a polyhedron $\{x \mid Ax \leq b\}$

\Rightarrow consider special cases: 1) $C = \{x \mid a_i^T x = b\}$, then $P_C(x_0) = x_0 + \frac{(b - a_i^T x_0) a_i}{a_i^T a_i}$
 (proof: KKT conditions)

2) $C = \{x \mid a_i^T x \leq b\}$, then $P_C(x_0) = \begin{cases} x_0 + \frac{(b - a_i^T x_0) a_i}{a_i^T a_i} a_i, & a_i^T x_0 \geq b \\ x_0, & a_i^T x_0 < b \end{cases}$
 (proof: KKT conditions)

3) $C = \{x \mid l \leq x \leq u\}$ (box constraints)

$P_C(x_0)_k = \max_{\text{in the } k^{\text{th}} \text{ item}} \{ \min\{x_0, u\}, l \}$

(ii) ball constraints

$C = \{x \in \mathbb{R}^n \mid \|x - m\| \leq r\}$ given $m, r \in \mathbb{R}^n, R$

$P_C(x_0) = \begin{cases} x_0, & \text{if } \|x_0 - m\| \leq r \\ m + \frac{r}{\|x_0 - m\|} (x_0 - m), & \text{if } \|x_0 - m\| > r \end{cases}$

(Proof: KKT conditions
 $(x - x_0) + \lambda(x - m) = 0$
 $\lambda(\|x - m\| - r) = 0$
 discuss $\lambda \neq 0, \lambda = 0 \Rightarrow$ get $\lambda \Rightarrow$ answer)

(iii) Matrix projection

(rank k matrices) $C = \{X \in \mathbb{R}^{m \times n} \mid \text{rank } X \leq k\}$ with $k \leq \min\{m, n\}$

For X_0 , by SVD, $X_0 = \sum_{i=1}^r \sigma_i u_i v_i^T$, with $r = \text{rank } X_0$

then $Y = \sum_{i=1}^{\min\{k, r\}} \sigma_i u_i v_i^T$ is a projection of X_0 on C .

• Optimization Problems with Convex Constraints

★ Theorem: FONC for Problems with convex constraints

Let f be cont. diff. on an open set that contains the convex, closed $S \subseteq \mathbb{R}^n$

Let $y^* \in S$ be a local minimizer of $(\min_{y \in S} f(y))$, then $\nabla f(y^*)^\top (y - y^*) \geq 0$, $\forall y \in S$

(Proof. Recall FONC for constrained problem - \forall feasible d $\nabla f(y^*)^\top d \geq 0$

& d satisfies $\exists \bar{t}$, $y^* + t\bar{d} \in S$, $t \in [0, \bar{t}]$ can always do that

use $y - y^*$ to substitute d . $y^* + t(y - y^*) = ty + (1-t)y^* \in S$, let $\bar{t} = 1$
($y \in S$)

\Rightarrow every feasible d corresponds to a $y \in S$ (one-to-one mapping)

★ Projection Theorem: Let S be a non-empty, closed & convex set

then a point y^* is the projection of x onto S , iff

$$(y^* - x)^\top (y - y^*) \geq 0, \forall y \in S$$

(Proof: $y^* = P_S(x) \iff \nabla f(y^*) = y^* - x \iff \nabla f(y^*)^\top (y - y^*) = 0, \forall y \in S$

$(f(y) = \frac{1}{2}\|y-x\|^2)$ $\begin{cases} y^* \text{ minimizer (projection pt)} \\ \text{(use convexity} \quad \Downarrow \\ \text{critical pt.} \Leftrightarrow \text{minimizer}) \end{cases}$ $\iff (y^* - x)^\top (y - y^*) \geq 0$)

COROLLARY:

(i) The mapping $P_S: \mathbb{R}^n \rightarrow \mathbb{R}^n$ is Lipschitz continuous with constant 1.

(Proof: $(P_S(x) - x)^\top (P_S(y) - P_S(x)) \geq 0$ (Because $P_S(y) \in S$)

$(P_S(y) - y)^\top (P_S(x) - P_S(y)) \geq 0$ (similarly)

Add up $\Rightarrow (P_S(x) - P_S(y))^\top (P_S(y) - P_S(x) + (x - y)) \geq 0$

$$\Rightarrow \|P_S(x) - P_S(y)\|^2 \leq (P_S(x) - P_S(y))^\top (x - y) \leq \|P_S(x) - P_S(y)\| \|x - y\|$$

$$\Rightarrow \|P_S(x) - P_S(y)\| \leq \|x - y\| \quad \begin{matrix} \downarrow \\ \text{Cauchy-Schwarz} \\ \text{inequality} \end{matrix}$$

(ii) The vector x^* is a critical point exactly when

$$x^* - P_S(x^* - \lambda \nabla f(x^*)) = 0, \forall \lambda > 0$$

$$\begin{aligned}
 & \text{Proof: } \nabla f(x^*)(y - x^*) \geq 0 \\
 & \quad y \in \mathbb{R} \Leftrightarrow \lambda \nabla f(x^*)(y - x^*) \geq 0 \\
 & \Leftrightarrow (x^* - (x^* - \lambda \nabla f(x^*))^\top (y - x^*) \geq 0 \Leftrightarrow x^* = P_{\mathbb{R}}(x^* - \lambda \nabla f(x^*))
 \end{aligned}$$

* The Projected Gradient Method:

1) Initialization: choose initial pt. $x_0 \in \mathbb{R}$ & $\alpha, \gamma \in (0, 1)$

2) Select $\lambda_k \geq 0$, computing $\nabla f(x_k)$ & new direction $d_k = P_{\mathbb{R}}(x_k - \lambda_k \nabla f(x_k)) - x_k$

3) If $\|d_k\| \leq \lambda_k \varepsilon$, STOP, output x_k

4) Choose a maximal step, using backtracking line search with α, γ

Armijo condition $f(x_k + \alpha_k d_k) - f(x_k) \leq \gamma \alpha_k \nabla f(x_k)^\top d_k$

5) Set $x_{k+1} = x_k + \alpha_k d_k$.

Comments: (i) $x_{k+1} = x_k + \alpha_k (P_{\mathbb{R}}(x_k - \lambda_k \nabla f(x_k)) - x_k) \in \mathbb{R}$ (in the line connected $x_k, P_{\mathbb{R}}(\cdot)$)

$$\begin{aligned}
 \text{(ii) descent direction} &= \nabla f(x_k)^\top (P_{\mathbb{R}}(x_k - \lambda_k \nabla f(x_k)) - x_k) = \nabla f(x_k)^\top (P_{\mathbb{R}}(x_k - \lambda_k \nabla f(x_k)) - P_{\mathbb{R}}(x_k)) \\
 &= \frac{1}{\lambda_k} (x_k - (x_k - \lambda_k \nabla f(x_k)))^\top (P_{\mathbb{R}}(x_k - \lambda_k \nabla f(x_k)) - P_{\mathbb{R}}(x_k)) \\
 &\stackrel{\text{(last page)}}{\leq} -\frac{1}{\lambda_k} \|P_{\mathbb{R}}(x_k - \lambda_k \nabla f(x_k)) - P_{\mathbb{R}}(x_k)\|^2 = -\frac{1}{\lambda_k} \|d_k\|^2 < 0
 \end{aligned}$$

(iii) Convergence: f is cont. diff., \mathbb{R} is non-empty, convex, promise descent

closed & stepsize λ_k bounded $0 < \lambda_k \leq \bar{\lambda} \leq \bar{\lambda}$

\Rightarrow every accumulation point is a critical point.
(of $\{x_k\}$)

Lecture 16 Integer Linear Programming (ILP)

- IDEA for solving integer programming.

1) General form: maximize $c^\top x$ (special case: binary integer programming)
 subject to $Ax = b$
 $x \geq 0$
 $x \in \mathbb{Z}^n$

maximize $c^\top x$
 subject to $Ax = b$
 $x \in \{0, 1\}^n$

2) Relaxation: LP relaxation:

Aim - to get an approximation opt. value & solutions

★ Strategy: For $x \in \mathbb{Z}^n$ — drop it directly

For $x_i \in \{0, 1\}$ — $0 \leq x_i \leq 1$

Directly rounding the optimal solution to integers may not yield a good solution!

⇒ However, the usage of LP relaxation — Bound!

- (i) For maximization IP, optimal value of relaxed LP provides an upper bound.
- (ii) ~ minimization ~, ... a lower bound

$$\text{Integrality gap} = |V^{IP} - V^{LP}|.$$

$$(0 \leq V^{IP} - V^{\text{rounding}} \leq V^{LP} - V^{\text{rounding}} \text{ (max case)} \text{ when RHS} = 0, V^{IP} = V^{LP} = V^{\text{rounding}})$$

Above is the ideal case, appearing in **total unimodularity** (TU) A with integer b

TU: a matrix A is **totally unimodular** if det of each sub-matrix is in $\{0, \pm 1\}$

A sufficient condition of TU: Let $A \in \mathbb{R}^{m \times n}$, then A is TU when

$$\text{eg. } \boxed{\begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix}}$$

(i) every column of A has no more than 2 non-zero entries

(ii) every entry in A belongs to $\{0, \pm 1\}$

(iii) The rows can be partitioned into 2 disjoint B & C

s.t. ~~two non-zero entries in the same column belong to every~~
~~with same sign different sets.~~
~~every two non-zero entries with different signs~~
~~in the same column belong to same set~~
~~the.~~

• Branch-and-Bound Method

1) IDEA: Branching — divide the feasible region into smaller ones & solve them.

 | Bounding — Use bounds (by relaxation) to reduce some branches.

2) Procedure (max case)

(i) Branching: Solving the LP relaxation. (Stop when getting integers)

* With optimal of LP x^* , if x_i^* (choose from small i) is fractional solution

 ⇒ branch them into $x_i \leq \lfloor x_i^* \rfloor$; $x_i \geq \lceil x_i^* \rceil$

 (eliminate $\lfloor x_i^* \rfloor \leq x_i \leq \lceil x_i^* \rceil$ — no influence) (2 branches)

* Each of 2 problems, same method to solve, opt. solution & value

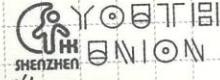
y_1^*, v_1^* ; y_2^*, v_2^* , compared to get opt. sol.

(DFS in practice)
 ↓ space, bound & simpleness

(ii) Bounding Procedure:

- { Upper bound - the original LP relaxation's opt. value
- Lower bound - every obj. value with integer feasible sol in every branch.

When the opt. value of some LP relaxation of IP branches \leq current lower bounds, (because of feasibility)
no need to consider these branches (discard)



3) Complexity: exponential, no polynomial-time algorithm in IP currently.

4) With softwares:
 { Gurobi can be used together with CVX for IPs
 { "intlinprog" function used in MATLAB for small size IPs.

(Another way - dynamic programming omitted)

Additional Topic — Duality Theory in General Convex cases

- The Lagrangian Dual Problems

Consider the Lagrangian problem in standard form

minimize $f_0(x)$

subject to $f_i(x) \leq 0, i=1, \dots, m$

$h_j(x) = 0, j=1, \dots, p$

① Define the Lagrangian associated with the problem $L: \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$

$$L(x, \lambda, \mu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^p \mu_j h_j(x)$$

with domain: $\text{dom } L = D \times \mathbb{R}^m \times \mathbb{R}^p$

dual variables / Lagrangian multipliers

② Define the Lagrange dual function $g: \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$ as the minimum value of Lagrangian over x , for $\lambda \in \mathbb{R}^m, \mu \in \mathbb{R}^p$

$$g(\lambda, \mu) = \inf_{x \in D} L(x, \lambda, \mu) = \inf_{x \in D} \left(f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^p \mu_j h_j(x) \right)$$

(when the Lagrangian is unbounded below in x , $g(\lambda, \mu) = -\infty$)

→ the dual function is concave even when Lagrangian

③ CQ: What's the best lower bound? (affine mapping & pointwise minimum)

given $\lambda \geq 0$?

Lagrange Dual Problem

maximize $g(\lambda, \mu)$

subject to $\lambda \geq 0$

eg.1 (LP Dual form derived with Lagrangian)

$$\text{minimize } C^T x$$

$$\text{subject to } Ax = b$$

$$x \geq 0$$

$$\text{Sol.} \Rightarrow \text{get the Lagrangian } L(x, \lambda, \mu) = C^T x + \lambda(-x) + \mu^T(Ax - b) \\ = -\mu^T b + (A^T \mu + c - \lambda)^T x$$

$$\text{The dual function } g(\lambda, \mu) = \begin{cases} -\mu^T b, & A^T \mu + c - \lambda = 0 \\ -\infty, & A^T \mu + c - \lambda \neq 0 \end{cases}$$

We get the dual form

$$\left(\begin{array}{l} \text{maximize } -\mu^T b \\ \text{subject to } A^T \mu + c - \lambda = 0 \\ \lambda \geq 0 \end{array} \right) \stackrel{g = -\mu^T b}{\Leftrightarrow} \left(\begin{array}{l} \text{maximize } b^T y \\ \text{subject to } c \geq A^T y \end{array} \right)$$

eg.2 (The two-way partitioning problem)

Consider the (non-convex) problem

$$\text{minimize } x^T W x$$

$$\text{subject to } x_i^2 = 1, \quad i=1, \dots, n \quad (W \in S^n)$$

Interpretation: two-way partitioning on a set $\{1, 2, \dots, n\}$ (with n elements / any set)

a feasible x corresponds to the partition $\{1, \dots, n\} = \{i \mid x_i = 1\} \cup \{j \mid x_j = -1\}$

The matrix coefficient W_{ij} — the cost of having elements i & j in the same partition,
 $-W_{ij}$ — different & objective fun — total cost.)

$$\text{Sol. Get the Lagrangian } L(x, \pi) = x^T W x + \sum_{i=1}^n \pi_i (x_i^2 - 1) \\ = x^T (W + \text{diag}(\pi)) x - \pi^T \pi$$

$$\text{The dual function } g(\pi) = \begin{cases} -\pi^T \pi, & W + \text{diag}(\pi) \succeq 0 \\ -\infty, & W + \text{diag}(\pi) \text{ not pos. semi-def.} \end{cases}$$

$$\Rightarrow \text{The dual problem} \left(\begin{array}{l} \text{maximize } -\pi^T \pi \\ \text{subject to } W + \text{diag}(\pi) \succeq 0 \end{array} \right) \text{ an SDP (semi-def program)}$$

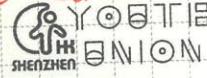
• Weak Duality & Strong Duality.

(1) Weak Duality = the optimal value of the Lagrange dual problem, denoted by d^* , is, by definition, the best lower bound on p^* obtained from the Lagrange dual function, then

$$d^* \leq p^* \quad \text{holds (even if non-convex)}$$

The optimal duality gap $\frac{p^* - d^*}{p^*}$ when d^*, p^* finite

(2) Strong duality: $d^* = p^*$. i.e. the optimal duality gap = zero
 NOT hold in general!!!



For convex problems with constraint qualification (CQ), it works.

★ **Slater's condition:** $\exists x \in \text{relint } D$ such that x is strictly feasible
 (i.e.) $f_i(x) < 0, h_j(x) = 0, \forall i, j$
 $(\text{relint } C = \{x \in C \mid B(x, r) \cap \text{aff } C \subseteq C, \text{ for some } r > 0\})$

e.g. (a non-convex quadratic problem with strong duality)

minimizing a nonconvex quadratic function over the unit ball / trust region problem.

$$\begin{aligned} & \text{minimize } x^T A x + 2b^T x \\ & \text{subject to } x^T x \leq 1, \quad \text{where } A \in \mathbb{R}^n \text{ & } A \neq 0, b \in \mathbb{R}^n \end{aligned}$$

Sol. the Lagrangian $L(x, \lambda) = x^T A x + 2b^T x + \lambda(x^T x - 1)$

$$= x^T (A + \lambda I) x + 2b^T x - \lambda \quad \text{pseudo-inverse}$$

the dual function is given by $g(\lambda) = \begin{cases} -b^T (A + \lambda I)^{-1} b - \lambda, & \lambda \in C(A + \lambda I), \\ -\infty, & \text{otherwise} \end{cases}$

⇒ the Lagrange dual problem is thus

$$\begin{aligned} & \text{maximize } -b^T (A + \lambda I)^{-1} b - \lambda \\ & \text{subject to } A + \lambda I \succeq 0, b \in C(A + \lambda I) \end{aligned}$$

diagonalization

$$\begin{aligned} & \text{maximize } -\frac{\sum_{i=1}^n \|q_i^T b\|^2}{\lambda + \lambda_i} - \lambda \\ & \text{subject to } \lambda \geq \max\{0, -\lambda_{\min}(A)\} \end{aligned}$$

(eigenvalues λ_i & eigenvectors q_i of A)

⇒ Conclusion: Strong duality holds for any opt problem with quadratic objective func & one quadratic inequality, provided slater's condition holds.

- Saddle-point interpretation

- (1) max-min characterization

consider the more general case $\begin{cases} \text{minimize } f_0(x) \\ \text{subject to } f_i(x) \leq 0 \quad \forall i = 1, \dots, n \end{cases}$

We know $p^* = \inf_{x \in D} \sup_{\lambda \geq 0} L(x, \lambda)$ (because $\sup_{\lambda \geq 0} L(x, \lambda) = \begin{cases} f_0(x), & f_i(x) \leq 0, i = 1, \dots, n \\ +\infty, & \text{otherwise} \end{cases}$)

★ $d^* = \sup_{\lambda \geq 0} \inf_{x \in D} L(x, \lambda)$ (by the definition)

weak duality says $\sup_{\lambda \geq 0} \inf_x L(x, \lambda) \leq \inf_x \sup_{\lambda \geq 0} L(x, \lambda)$

(In reality, it works for every function $f(x, \lambda)$,

because $f(\tilde{x}, \lambda) \geq \inf_x f(x, \lambda)$, for any \tilde{x} , given λ

$$\Rightarrow \sup_{\lambda \geq 0} f(\tilde{x}, \lambda) \geq \sup_{\lambda \geq 0} \inf_x f(x, \lambda), \text{ then let } \tilde{\pi} = \arg \min_x (\sup_{\lambda \geq 0} f(x, \lambda))$$

$$\Rightarrow \left[\inf_{\lambda \geq 0} \sup_x f(x, \lambda) \geq \sup_{\lambda \geq 0} \inf_x f(x, \lambda) \right] \text{ max-min inequality}$$

strong duality says

$$\sup_{\lambda \geq 0} \inf_x L(x, \lambda) = \inf_x \sup_{\lambda \geq 0} L(x, \lambda)$$

(2) Saddle-point interpretation

we refer to a pair $(\hat{w}, \hat{z}) \in W \times Z$ as a saddle point of f

if $f(\hat{w}, \hat{z}) \leq f(w, \hat{z}) \leq f(\hat{w}, z)$, $\forall w \in W, z \in Z$

$$\Rightarrow \left[\begin{array}{l} f(\hat{w}, \hat{z}) = \inf_{w \in W} f(w, \hat{z}) \\ \sup_{z \in Z} f(\hat{w}, z) \end{array} \right] \xrightarrow{\text{implies}} \text{strong duality holds!}$$

$\Rightarrow (x^*, \lambda^*)$ are primal & dual optimal pts with strong duality exactly when (x^*, λ^*) is a saddle-point of the Lagrangian $L(x, \lambda)$

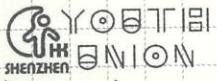
Explanation:

$$\sup_{z \in Z} f(\hat{w}, z) \geq \inf_{w \in W} \sup_{z \in Z} f(w, z) \quad \checkmark \leftarrow \text{Weak Duality}$$

$$\sup_{z \in Z} f(\hat{w}, z) = \inf_{w \in W} f(w, \hat{z}) \leq \sup_{w \in W} \inf_{z \in Z} f(w, z)$$

$\Rightarrow \sup \inf = \inf \sup$, when applied in x & λ , we can get that strong duality holds for $L(x, \lambda)$, $\lambda \geq 0$ & $x \in R$

Distributed Optimization & Statistical Learning via ADMM



(2021-22 Term 1 additional topic)

Part 1 precursors

- The Conjugate Function

(1) Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$, define the function $f^*: \mathbb{R}^n \rightarrow \mathbb{R}$ as the conjugate of f .
where $f(y)^* := \sup_{x \in \text{dom } f} (y^T x - f(x))$.

$$\text{dom } f^* = \text{all } y \in \mathbb{R}^n \text{ s.t. } \sup_{x \in \text{dom } f} (y^T x - f(x)) < \infty$$

For general matrix
 $\langle X, Y \rangle = \text{tr}(X^T Y)$
For square
 $= \text{tr}(XY)$

(2) Examples:

(Log-determinant) consider $f(X) = \log \det X$ on S_+^n (pos. def)

The conjugate of f - defined as $f^*(Y) = \sup_{X \succ 0} (\text{tr}(YX) + \log \det X)$

dom $f^* = -S_+^n$ because if $Y \leq 0$, then

\exists an eigenvector U of Y , with $\|U\|=1$, and corresponding eigenvalue $\lambda \geq 0$

Take $X = I + tUU^T \succ 0$

$$\text{tr}(YX) + \log \det X = \text{tr}(Y) + t\lambda + \log \det(I + tUU^T) = \text{tr}(Y) + t\lambda + \log(1+t) \rightarrow \infty \text{ as } t \rightarrow \infty.$$

If $Y \prec 0$, find the maximizing X by gradient

$$\nabla_X (\text{tr}(YX) + \log \det X) = Y + X^T = 0, \text{ yield that } Y = -X^T.$$

(3) Basic Properties

Fenchel's inequality

$$f(x) + f^*(y) \geq y^T x$$

$$(f(x) + f^*(y) = f(x) + \sup_{x \in \text{dom } f} (y^T x - f(x)) \geq y^T x)$$

"=" holds when $\sup_{x \in \text{dom } f} (y^T x - f(x)) = \max_{x \in \text{dom } f} (y^T x - f(x))$

Legendre Transform: if f is convex, differentiable with $\text{dom } f = \mathbb{R}^n$
then any maximizer x^* satisfies $y = \nabla f(x^*)$

Let $z \in \mathbb{R}^n$, $y = \nabla f(z)$,

$$\Rightarrow f^*(y) = z^T \nabla f(z) - f(z) \quad (\text{in this case})$$

• Dual Decomposition

(1) Dual Ascent (like gradient descent)

consider the equality-constrained convex problem

$$\begin{array}{ll} \text{minimize } f(x) \\ \text{subject to } Ax = b \end{array} \quad (*)$$

with $x \in \mathbb{R}^n$, $A \in \mathbb{R}^{m \times n}$, $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is convex.

We get the Lagrangian $L(x, y) = f(x) + y^T(Ax - b)$

$$\text{Dual Func. } g(y) = \inf_x f(x) + y^T(Ax - b) = -f^*(-A^T y) - b^T y$$

\Rightarrow Dual Problem $\max_{y \in \mathbb{R}^m} g(y)$

* Assume strong duality holds $\Rightarrow x^* = \arg \min_x L(x, y^*)$
 (provided only one minimizer of $L(x, y)$)

Iterations of dual ascent method

$$x^{k+1} := \arg \min_x L(x, y^k)$$

$$y^{k+1} := y^k + \alpha^k (Ax^{k+1} - b) \quad \text{evaluate } \nabla g(y) \text{ as }$$

where $\alpha^k > 0$ is the step size, $k/k+1$ - iteration counter

with appropriate choice, $g(y^{k+1}) > g(y^k)$

then y^k converges to optimal sol. when $k \rightarrow \infty$.

(2) Dual Decomposition

(The major benefit of dual ascent)

(i) A Decentralized Algorithm.

Suppose, for instance, objective f is separable, which means

$$f(x) = \sum_{i=1}^n f_i(x_i), \quad \text{where } x = (x_1, \dots, x_n) \text{ & } x_i \in \mathbb{R}^{n_i}$$

are subvectors of x .

Partitioning the matrix A as $[A_1 \cdots A_n]$

$$\text{so that } Ax = \sum_{i=1}^n A_i x_i$$

The Lagrangian can be written as

$$L(x, y) = \sum_{i=1}^N L_i(x_i, y) = \sum_{i=1}^N (f_i(x_i) + y^T A_i x_i - \frac{1}{N} y^T b)$$

(ii) The Iterations of dual decomposition

$$\begin{aligned} x_i^{k+1} &:= \underset{x_i}{\operatorname{argmin}} L_i(x_i, y^k) && \text{broadcast} \\ y^{k+1} &:= y^k + \alpha^k (A x^{k+1} - b) && \text{gather} \end{aligned}$$

> 2 steps

- Augmented Lagrangian & the Method of Multipliers

(1) Aim: Bring robustness to the dual ascent, yielding convergence without assumptions like strict convexity / fitness of f .

(2) The augmented Lagrangian for (x) is

$$L_p(x, y) = f(x) + y^T (Ax - b) + \frac{\rho}{2} \|Ax - b\|_2^2 \quad (\text{Namely } L_0(x, y))$$

(subject to $Ax = b$)

$Ax - b = 0 \quad + \frac{\rho}{2} \|Ax - b\|_2^2$
for feasible x

(3) The Iterations

$$\begin{aligned} x^{k+1} &:= \underset{x}{\operatorname{argmin}} L_p(x, y^k) \\ y^{k+1} &:= y^k + \rho(A x^{k+1} - b) \end{aligned}$$

hard to separate because of $\|Ax - b\|_2^2$
method of multipliers

Part II Alternating Direction Method of Multipliers (ADMM)

- Algorithm

(1) Aim: Blend the decomposability of dual ascent with the superior convergence properties of the method of multipliers.

Solve the problem in the form

$$\begin{aligned} &\underset{x, z}{\operatorname{minimize}} f(x) + g(z) \quad (x \in \mathbb{R}^n, z \in \mathbb{R}^m) \\ &\text{Subject to } Ax + Bz = c \quad (A \in \mathbb{R}^{pn}, B \in \mathbb{R}^{pm}) \end{aligned}$$

(Assume f, g are both convex)

The augmented Lagrangian $L_p(x, z, y) = f(x) + g(z) + y^T (Ax + Bz - c) + \frac{\rho}{2} \|Ax + Bz - c\|_2^2$

(2) The Iteration steps:
(of ADMM)

$$x^{k+1} := \arg \min_x L_p(x, z^k, y^k)$$

$$z^{k+1} := \arg \min_z L_p(x^{k+1}, z, y^k)$$

$$y^{k+1} := y^k + p(Ax^{k+1} + Bz^{k+1} - c)$$

updated in
an alternating fashion

Scaled Form

Let residual $r = Ax + Bz - c \rightarrow$ primal residual

$$y^T r + \frac{P}{2} \|r\|_2^2 = \frac{P}{2} \|r + \frac{1}{P} y\|_2^2 - \frac{1}{2P} \|y\|_2^2, \text{ Let } u = \frac{1}{P} y$$

$$\text{then } y^T r + \frac{P}{2} \|r\|_2^2 = \frac{P}{2} \|r + u\|_2^2 - \frac{P}{2} \|u\|_2^2.$$

$$x^{k+1} := \arg \min_x (f(x) + \frac{P}{2} \|Ax + Bz^k - c + u^k\|_2^2)$$

$$z^{k+1} := \arg \min_z (g(z) + \frac{P}{2} \|Ax^{k+1} + Bz - c + u^k\|_2^2)$$

$$u^{k+1} := u^k + r^{k+1} (\text{i.e., } Ax^{k+1} + Bz^{k+1} - c)$$

• Optimality Conditions & Stopping Criterion

(i) Necessary & Sufficient optimality conditions

$$\left. \begin{array}{l} \text{primal feasibility: } Ax^* + Bz^* - c = 0 \\ \text{dual feasibility: } 0 \in \partial f(x^*) + A^T y^* \end{array} \right\} \quad (1)$$

$$0 \in \partial f(x^*) + B^T y^* \quad (2)$$

$$0 \in \partial g(z^*) + A^T y^* \quad (3)$$

(Addition. "∂" denotes the subdifferential operator)

(Defn) $\partial f(x) := \{y \mid f(y) \geq f(x) + g_x^T(y - x)\}$ $f(x)$ does NOT need to be differentiable at x .

(ii) Because z^{k+1} minimize $L_p(x^{k+1}, z, y^k) \Rightarrow 0 \in \partial f(z^{k+1}) + B^T y^{k+1}$

which means (3) holds for all $(x^{k+1}, z^{k+1}, y^{k+1})$, $\forall k \in \mathbb{N}$

Similarly $s^{k+1} = P A^T B (z^{k+1} - z^k) \in \partial f(x^{k+1}) + A^T y^{k+1}$

[dual residual] ($s^{k+1} \rightarrow 0$ & $r^{k+1} \rightarrow 0$ as ADMM iterates)

(iii) Stopping Criteria

when r^k, s^k are both small, the objective suboptimality must be small

A reasonable termination criterion: $\|r^k\|_2 \leq \epsilon_{\text{pri}}^{\text{stop}}$; $\|s^k\|_2 \leq \epsilon_{\text{dual}}^{\text{stop}}$

Chosen way: (absolute & relative criterion)

© Y@BTG
SHENZHEN UNION

ϵ^{pri} tolerance for the primal conditions

$$= \sqrt{p} \epsilon^{\text{abs}} + \epsilon^{\text{rel}} \max \{ \|Ax^k\|_2, \|Bz^k\|_2, \|C\|_2 \}$$

ϵ^{dual} (tolerance for the dual conditions)

$$= \sqrt{n} \epsilon^{\text{abs}} + \epsilon^{\text{rel}} \|A^T y^k\|_2$$

ϵ^{abs} = absolute tolerance

ϵ^{rel} = relative tolerance

p, n : l_2 norms are in R^p, R^n (respectively)

• Proof for convergence

2 assumptions: (1) $f: R^n \rightarrow R \cup \{-\infty\}$ & $g: R^m \rightarrow R \cup \{+\infty\}$ are closed, proper & convex.

(closed: epigraph $\text{epi } f = \{(x, t) \in R^n \times R \mid f(x) \leq t\}$ is closed.)

or sublevel sets $\{x \in \text{dom } f \mid f(x) \leq \alpha\}$ are closed $\forall \alpha \in R$)

(proper convex function: epigraph ($\text{epi } f = \{(x, t) \in R^n \times R \mid f(x) \leq t\}$) is non-empty & contains no vertical lines)

(i.e., $f(x) < +\infty$ for at least one x
& $f(x) > -\infty$ for all x)

criterion of proper: $f(x)$ is proper iff $\text{dom } f$ is convex, non-empty & restriction of f to C is finite.

(2) Unaugmented Lagrangian $L_o = f(x) + g(z) + y^T(Ax + Bz - c)$ has a saddle point (i.e., $\exists (x^*, z^*, y^*)$, not necessarily unique).

$$\text{s.t. } L_o(x^*, z^*, y) \leq L_o(x^*, z^*, y^*) \leq L_o(x, z, y^*) \quad \forall x, y, z$$

Need to show

$$p^* - p^{k+1} \leq y^{*T} r^{k+1} \quad (\text{A.1}) \quad (\text{If } r^{k+1} \rightarrow 0, p^{k+1} \rightarrow p^*)$$

$$p^{k+1} - p^* \leq - (y^{k+1})^T r^{k+1} - p(B(z^{k+1} - z^k))^T (-r^{k+1} + B(z^{k+1} - z^k)) \quad (\text{A.2})$$

$$V^k \leq V^* - \beta \|r^k\|_2^2 - \beta \|B(z^k - z^*)\|_2^2 \quad (\Rightarrow r^k \rightarrow 0, B(z^k - z^*) \rightarrow 0, \Rightarrow s^k \rightarrow 0)$$

$$\text{where } V^k := \frac{1}{2} \|y^k - y^*\|_2^2 + \beta \|B(z^k - z^*)\|_2^2$$

is a Lyapunov Function for the algorithm.

(PROOF.)

$$\text{Since, } L_o(x^*, z^*, y^*) \leq L_o(x^{k+1}, z^k, y^*)$$

saddle point property $\Rightarrow p^* \leq p^{k+1} + y^{*T} r^{k+1}$ (which is A.1)

Since x^{k+1} minimize $L_o(x, z^k, y^k)$, $0 \in \partial_x L_o(x^{k+1}, z^k, y^k)$ (closed, proper, convex \Rightarrow subdifferential)

$$\partial f(x^{k+1}) + A^T y^k + \partial f'(Ax^{k+1} + Bz^k - c)$$

$$y^{k+1} = y^k + \rho r^{k+1} \rightarrow$$

SHENZHEN UNION $0 \in \partial f(x^{k+1}) + A^T(y^{k+1} - \rho B(z^{k+1} - z^k))$
 $\therefore z^{k+1}$ minimizes $f(x) + (y^{k+1} - \rho B(z^{k+1} - z^k))^T A x$

Similarly, z^k minimizes $g(z) + (y^{k+1})^T B z$

$$\begin{cases} f(x^{k+1}) + (y^{k+1} - \rho B(z^{k+1} - z^k))^T A x^{k+1} \leq f(x^*) + (y^{k+1} - \rho B(z^{k+1} - z^k))^T A x^* \\ g(z^{k+1}) + (y^{k+1})^T B z^{k+1} \leq g(z^*) + (y^{k+1})^T B z^* \end{cases}$$

Add the two together, $\rho r^{k+1} + (y^{k+1})^T r^{k+1} - (\rho B(z^{k+1} - z^k))^T A x^{k+1} \leq \rho^* - (\rho B(z^{k+1} - z^k))^T A x^*$
 (which is A.2)

Adding A.1 & A.2, regrouping terms (& multiply by 2)

$$2(y^{k+1} - y^*)^T r^{k+1} - 2\rho(B(z^{k+1} - z^k))^T r^{k+1} + 2\rho(B(z^{k+1} - z^k))^T (B(z^{k+1} - z^*)) \leq 0$$

Substitute $y^{k+1} = y^k + \rho r^{k+1}$, $r^{k+1} = \frac{1}{\rho}(y^{k+1} - y^k)$, rewrite the 1st term

$$2(y^{k+1} - y^*)^T r^{k+1} = \frac{1}{\rho} \|y^{k+1} - y^*\|_2^2 + \rho \|r^{k+1}\|_2^2 - \frac{1}{\rho} \|y^k - y^*\|_2^2$$

$$\begin{aligned} & -2\rho(B(z^{k+1} - z^k))^T r^{k+1} + 2\rho(B(z^{k+1} - z^k))^T B(z^{k+1} - z^*) + \rho \|r^{k+1}\|_2^2 \quad (z^{k+1} - z^* = z^{k+1} - z^k + z^k - z^*) \\ & = \rho \|r^{k+1} - B(z^{k+1} - z^*)\|_2^2 + \rho \|B(z^{k+1} - z^k)\|_2^2 + 2\rho(B(z^{k+1} - z^k))^T (B(z^k - z^*)) \end{aligned}$$

Substitute $z^{k+1} - z^k = (z^{k+1} - z^*) - (z^k - z^*)$ in the last 2 terms

$$\Rightarrow V^k - V^{k+1} \geq \rho \|r^{k+1} - B(z^{k+1} - z^k)\|_2^2 \geq \rho \|r^{k+1}\|_2^2 + \rho \|B(z^{k+1} - z^k)\|_2^2$$

(Because z^{k+1} minimizes $g(z) + (y^{k+1})^T B z$, z^k minimizes $g(z) + (y^k)^T B z$)

$$\Rightarrow (y^{k+1} - y^k)^T B(z^{k+1} - z^k) = \rho(r^{k+1})^T B(z^{k+1} - z^k) \leq 0.$$

Part III Consensus & Distributed Model Fitting

- Background Knowledge

① soft thresholding consider $f(x) = \lambda \|x\|_1$, with $\lambda > 0$ & $A = I$ (in ADMM)

In this case x_i -update is $x_i^+ := \arg \min_{x_i} (\lambda |x_i| + \frac{\rho}{2} (x_i - v_i)^2)$

The solution is $x_i^+ = S_K(v_i)$, where $S_K(a) = \begin{cases} a - x, & a > K \\ 0, & |a| \leq K \\ a + K, & a < -K \end{cases}$

(Sof. sub-gradient method.)

$$(i.e., S_K(a) = (1 - \frac{K}{|a|})_+ a, \forall a \neq 0)$$

② proximity operator (of f with penalty ρ) (denoted as $\text{prox}_{f, \rho}(u)$)

$$\text{prox}_{f, \rho}(u) := \arg \min_x (f(x) + \frac{\rho}{2} \|x - u\|_2^2)$$

• Global variable consensus optimization

① general case: $\underset{\text{rewrite}}{\text{minimize } f(x) = \sum_{i=1}^n f_i(x_i)}, x \in \mathbb{R}^n, f_i: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$

SHENZHEN UNION

are convex \downarrow can add some constraints

$$\left(\begin{array}{l} \text{minimize } f(x) = \sum_{i=1}^n f_i(x_i) \\ \text{subject to } x_i - z = 0, i=1, \dots, N \end{array} \right) \xrightarrow[\text{Form}]{} \left\{ \begin{array}{l} x_i^{k+1} = \arg \min_{x_i} (f_i(x_i) + y_i^T(x_i - \bar{x}^k)) \\ y_i^{k+1} = y_i^k + \rho(x_i^{k+1} - \bar{x}^k) \end{array} \right.$$

② with regularization

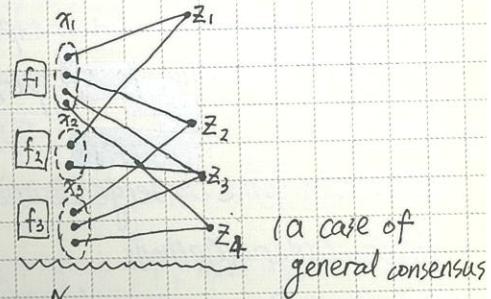
$$\left(\begin{array}{l} \text{minimize } f(x) + g(z) = \sum_{i=1}^n f_i(x_i) + g(z) \\ \text{subject to } x_i - z = 0, i=1, \dots, N \end{array} \right)$$

$$\xrightarrow[\text{Form}]{} \left\{ \begin{array}{l} x_i^{k+1} = \arg \min_{x_i} (f_i(x_i) + \frac{\rho}{2} \|x_i - z^k + u_i^k\|_2^2) \\ z^{k+1} = \arg \min_z (g(z) + \frac{\rho}{2} \|z - \bar{x}^k - u^k\|_2^2) \\ u_i^{k+1} = u_i^k + x_i^{k+1} - z^{k+1} \end{array} \right.$$

• General Form consensus optimization

- ① Each component of each local variable corresponds to some global variable component z_g .

(i.e., \exists a mapping $g = G(i,j)$ s.t. $(x_i)_j = \sum g(i,j)$)



Let $\bar{z}_i \in \mathbb{R}^{nc}$, with $(\bar{z}_i)_j = \sum g(i,j) \Rightarrow \left(\begin{array}{l} \text{minimize } \sum_{i=1}^N f_i(x_i) \\ \text{subject to } x_i - \bar{z}_i = 0, i=1, \dots, N \end{array} \right)$

$$\xrightarrow[\text{Form}]{} \left\{ \begin{array}{l} x_i^{k+1} = \arg \min_{x_i} (f_i(x_i) + y_i^T x_i + \frac{\rho}{2} \|x_i - \bar{z}_i^k\|_2^2) \\ z_g^{k+1} = \frac{\sum g(i,j) \cdot (\bar{z}_i^{k+1})_j + \rho(y_g^k)}{\sum g(i,j) = g \cdot 1} \rightarrow \text{avg} = 0 = \bar{z}_g^k \sum g(i,j) = g(\bar{z}_i^{k+1})_j \\ y_i^{k+1} = y_i^k + \rho(x_i^{k+1} - \bar{z}_i^{k+1}) \end{array} \right.$$

② Sharing:

$$\underset{\substack{x_i \in \mathbb{R}^n \\ i=1, \dots, N}}{\text{minimize}} \sum_{i=1}^N f_i(x_i) + g\left(\sum_{i=1}^N x_i\right) \xrightarrow[\text{Form}]{\text{ADMM}} \left\{ \begin{array}{l} x_i^{k+1} = \arg \min_{x_i} (f_i(x_i) + \frac{\rho}{2} \|x_i - x^k + \bar{x}^k - \bar{z}^k + u_i^k\|_2^2) \\ \bar{z}^{k+1} = \arg \min_{\bar{z}} (g(N\bar{z}) + \frac{\rho}{2} \|\bar{z} - u^k - \bar{x}^k\|_2^2) \\ u_i^{k+1} = u_i^k + \bar{x}^k - \bar{z}^{k+1} \end{array} \right.$$

• Distribute Model Fitting

1) General Convex model fitting problem

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} l(Ax - b) + R(x) \quad \begin{array}{l} \text{regularization} \\ \text{function} \end{array}$$

Assume additive $l = \sum_{i=1}^m l_i(a^T x - b)$, loss function $Tikhonov \sim = \lambda \|x\|_2^2$
common regularizations / lasso penalty = $\lambda \|x\|_1$

2) Examples: (supervised learning)

① Regression: with the form $b_i = a_i^T \omega + v_i$

(a_i - i^{th} feature vector, v_i - measurement noises (indep. with log-concave densities p_i)), then $l_i(\omega) = -\log p_i(\omega)$.

| $\tau = 0$, MLE of ω under noise model p

| $\tau > 0$, negative log prior density of ω , MAP estimation

② Classification: $q_i \in \{-1, 1\}$ & $q_i = \text{sign}(p_i^T \omega + v)$

(Find weight vector $\omega \in \mathbb{R}^{n+1}$ & offset $v \in \mathbb{R}$) Then, $l_i = l_i(q_i | p_i^T \omega + v)$

This can be modeled as $\hat{m} \sum_{i=1}^m l_i(q_i | p_i^T \omega + v) + \tau w^T \omega$ penalized erm

| l_i - hinge loss ($1 - \mu_i + \lambda_2$ penalty) \Rightarrow SVM

(empirical risk minimization)

| l_i - logistic loss $\log(1 + \exp(-\mu_i))$ \Rightarrow logistic regression

| l_i - exponential loss $\exp(-\mu_i)$ \Rightarrow boosting

margin

Additions (some convergence analyses & introductions)

• First-Order Methods

1) Gradient Descent Method: (convergence analyses) "df"

Assume that f is convex, differentiable & Lipschitz continuous with const $L > 0$.

Then, we get, with fixed step size $t \leq \frac{1}{L}$, then $f(x^k) - f^* \leq \frac{\|x^0 - x^*\|_2^2}{2tK}$

(proof.) By 1st-order condition of convexity,

$$f^* = f(x^*) \geq f(x^t) + \nabla f(x^t)^T (x^* - x^t) \quad \textcircled{1}$$

Because of $C_L^{(1)}$, we get $\nabla^2 f \leq L I$ (if $f \in C^2$), then it's easy to get

$$\text{that } f(x^{t+1}) \leq f(x^t) + \nabla f(x^t)^T (x^{t+1} - x^t) + \frac{L}{2} \|x^{t+1} - x^t\|_2^2 = f(x^t) - t(1 - \frac{L}{2}) \|\nabla f(x^t)\|_2^2 \quad \textcircled{2}$$

$$\text{combine } \textcircled{1} \text{ \& } \textcircled{2}, \text{ we have } f(x^{t+1}) - f^* \leq \nabla f(x^t)^T (x^t - x^*) - t(1 - \frac{L}{2}) \|\nabla f(x^t)\|_2^2$$

$$\Rightarrow f(x^{t+1}) - f^* \leq \frac{1}{t} \|x^t - x^*\|_2^2 - \frac{1}{t} (x^{t+1} - x^*)^T (x^t - x^*) - \frac{1}{t} (1 - \frac{L}{2}) \|x^t - x^{t+1}\|_2^2$$

$$\therefore f(x^{t+1}) - f^* \leq \frac{1}{2t} (\|x^t - x^*\|_2^2 - \|x^{t+1} - x^*\|_2^2) \quad \textcircled{3}$$

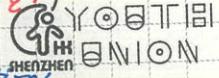
$$\text{Add } \textcircled{3} \text{ from } t=0 \text{ to } \infty, \text{ we get } \sum_{t=0}^{\infty} f(x^t) - f^* \leq \frac{1}{2t} \|x^0 - x^*\|_2^2$$

$\therefore f(x^t) - f(x^*) \rightarrow 0$ as $t \rightarrow \infty$, which means that it converges

Because $|f(x^t) - f(x^*)|$ decreases (according to $\textcircled{3}$), we get

$$K[f(x^t) - f(x^*)] \leq \sum_{i=1}^K f(x^i) - f^* \leq \sum_{i=1}^{\infty} f(x^i) - f^* \leq \frac{1}{2t} \|x^0 - x^*\|_2^2. \quad \square$$

We say gradient descent has convergence rate $O(1/k)$ (or $O(\epsilon)$)
 That is, it finds ϵ -suboptimal point ($f(x^k) - f^* \leq \epsilon$) in $O(1/\epsilon)$ iterations.



(For backtracking, with t replaced by γ/L , we can get similar results)

(ii) Assume [strong convexity] holds (with all other properties same as (i)), $\nabla^2 f(x) \succeq mI$
 Then, gradient descent with fixed $t \leq \frac{2}{m+L}$ or with backtracking line search
 satisfies $f(x^k) - f^* \leq \delta^k \frac{L}{2} \|x^0 - x^k\|_2^2$, where $\delta \in (0, 1)$. ($t > 0$)

(proof.) Similar to (i), we get $f(x^k) - f^* \leq f(x^0) - f^* - t(1 - \frac{tL}{2}) \|\nabla f(x^0)\|_2^2$

By Polyak-Lojasiewicz inequality, $\|\nabla f(x^0)\|_2^2 \geq 2m[f(x^0) - f^*]$

$$\Rightarrow f(x^k) - f^* \leq \left[1 - 2mt\left(1 - \frac{tL}{2}\right)\right] [f(x^0) - f^*] = \delta [f(x^0) - f^*]$$

$$\geq 1 - \frac{m}{L} > 0 \quad \leftarrow < 1 - 2mt \cdot \frac{m}{m+L} < 1, \forall t \in \mathbb{N}$$

$$\therefore f(x^k) - f^* \leq \delta^k [f(x^0) - f(x^k)] \leq \delta^k \frac{L}{2} \|x^0 - x^k\|_2^2. \square$$

(For backtracking, let $\delta = 1 - \min\{2m\alpha, \frac{2m\alpha m}{L}\}$, for exact, $\delta = 1 - \frac{m}{L}$)

We say rate under strong convexity is $O(\delta^k)$ (or $O(\log(1/\epsilon))$), that is,
 it finds ϵ -suboptimal point in $O(\log(1/\epsilon))$ iterations (Also called linear convergence)

★ First-order Method: iterative, update x^k in $x^0 + \text{span}\{\nabla f(x^0), \dots, \nabla f(x^{k-1})\}$

(i) (Thm-Nesterov) If $k \leq \frac{n-1}{2}$ (n -problem dimension), & any starting point x^0 , \exists a function in the problem class s.t. any 1st-order method satisfies $f(x^k) - f^* \geq \frac{3L}{32} \frac{\|x^0 - x^k\|_2^2}{(k+1)^2}$.

From this we can know that the best rate for 1st-order method is $O(1/\sqrt{\epsilon})$.

(over class of differentiable functions with $C_L^{1,1}$)

(ii) For non-convex condition, E.g. gradient-descent with fixed step size $t \leq 1/L$:

$$(\text{easily to show}) \min_{i=0, \dots, k} \|\nabla f(x^i)\|_2 \leq \sqrt{\frac{2(f(x^0) - f^*)}{t(k+1)}} \quad O(1/\epsilon^2)$$

2) Sub-gradient Descent Method (For Non-differentiable Functions):

(i) Operations: Every time let $x^k = x^{k-1} - t_k \cdot g^{k-1}$, where g^{k-1} is a sub-gradient at x^{k-1}

Because sub-gradient method is NOT necessarily a descent one, we take $[f(x_{\text{best}}) =$

$$\min_{i=0, \dots, K} f(x^i)]$$

(ii) Step-size choice = (rule) must be pre-specified, NOT adaptively computed.

Fixed step size $t_k = t$, all $k=1, 2, \dots$

Diminishing step-size $\sum_{k=1}^{\infty} t_k^2 < \infty$; $\sum_{k=1}^{\infty} t_k = \infty$ (e.g. harmonic series, $t_k = 1/k$)

iii) Convergence analysis:

Assume f is convex, & Lipschitz continuous with const $G > 0$ (i.e. $|f(x) - f(y)| \leq G\|x - y\|$)

Then, with fixed step size t , $\lim_{k \rightarrow \infty} f(x_{\text{best}}^k) - f^* \leq \frac{Gt^2}{2}$

diminishing step size, $\lim_{k \rightarrow \infty} f(x_{\text{best}}^k) = f^*$

$$(\text{proof.}) \|x^{k+1} - x^*\|_2^2 = \|x^k - x^* - t_k g^k\|_2^2 = \|x^k - x^*\|_2^2 - 2t_k g^k{}^T (x^k - x^*) + t_k^2 \|g^k\|_2^2$$

Because of convexity, $f^* - f(x^k) \geq g^k{}^T (x^* - x^k)$, thus $\|x^{k+1} - x^*\|_2^2 \leq \|x^k - x^*\|_2^2 - 2t_k [f(x^k) - f^*] + t_k^2 \|g^k\|_2^2$.

$$\therefore \|x^{k+1} - x^*\|_2^2 \leq \|x^0 - x^*\|_2^2 - 2 \sum_{l=0}^k t_l [f(x^l) - f^*] + \sum_{l=0}^{\infty} t_l^2 \|g^l\|_2^2, \forall k \in \mathbb{N} \quad (*)$$

Because $\|x^{k+1} - x^*\|_2^2 \geq 0$, $f(x^l) - f^* \geq f(x_{\text{best}}^k) - f^* \& \|g^l\|_2^2 \leq G^2$, let $\|x^0 - x^*\|_2 = R > 0$

$$\text{we get } 0 \leq R^2 - 2[f(x_{\text{best}}^k) - f^*] \sum_{l=0}^k t_l + G^2 \sum_{l=0}^k t_l^2. \quad \forall k \in \mathbb{N}$$

Fixed step size: $f(x_{\text{best}}^k) - f^* \leq \frac{R^2 + G^2 \sum_{l=0}^k t_l^2}{2 \sum_{l=0}^k t_l} = \frac{R^2}{2kt} + \frac{G^2 t}{2}, \forall k \in \mathbb{N}$

Diminishing step size: $\lim_{k \rightarrow \infty} f(x_{\text{best}}^k) - f^* \leq \frac{R^2 + G^2 \sum_{l=0}^{\infty} t_l^2}{2 \sum_{l=0}^{\infty} t_l} = 0.$

(Note.) Convergence rate: $\frac{R^2}{2kt} \leq \frac{\epsilon}{2} \& \frac{G^2 t}{2} \leq \frac{\epsilon}{2} \Rightarrow k = \frac{R^2}{t\epsilon} = \frac{R^2 G^2}{\epsilon^2}$.
 $O(\frac{1}{\epsilon^2})$ (compared with gradient method)

* Polyak Step-size $t_k = \frac{f(x^k) - f^*}{\|g^k\|_2^2}$ (when f^* is in hand)

(motivated by minimize RHS of $(*)$)

(iv) Non-smooth first-order Method: iterative, update x^k in

$$x^0 + \text{span}\{g^0, \dots, g^{k-1}\}$$

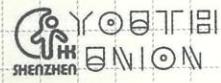
(Thm - Nesterov) $\forall k \leq n-1$ & any starting point x^0 , \exists a function in the problem

class s.t. any non-smooth 1st-order methods s.t. $f(x^k) - f^* \geq \frac{RG}{2(1 + \sqrt{k+1})}$

(Cannot do better than $O(\frac{1}{\epsilon^2})$)

3) Proximal Gradient Descent:

A method between gradient-descent & sub-gradient methods,
applied to special-structured functions.



(i) Suppose $f(x) = g(x) + h(x)$ (Decomposition)
 Convex convex
 differentiable (probably) Not-differentiable

We try to apply gradient-descent on "g", because $\tilde{x}^+ = x - t \nabla g(x) = \arg \min_{\mathbb{Z}} \{ g(x) + \nabla g(x)^T (\mathbb{Z} - x) + \frac{t}{2} \| \mathbb{Z} - x \|^2_2 \}$, approximate "g" by this quadratic function

(For the next step). Thus, we need $x^+ = \arg \min_{\mathbb{Z}} \{ \frac{1}{2t} \| \mathbb{Z} - (x - t \nabla g(x)) \|^2_2 + h(\mathbb{Z}) \}$
 $\stackrel{\triangle}{=} \text{prox}_{h, \frac{1}{t}}[x - t \nabla g(x)]$ (proximity operator, see in previous pages)

∴ Every step let $\underline{x^{k+1}} = \text{prox}_{h, \frac{1}{t}}(x^k - t \nabla g(x^k)) \stackrel{\triangle}{=} x^k - t \nabla g(x^k) G_{t \nabla g}(x^k)$,

where $G_{t \nabla g}(x^k) = \frac{1}{t} [x^k - \text{prox}_{h, \frac{1}{t}}(x^k - t \nabla g(x^k))]$.

(ii) Convergence Analysis:

(Thm) Assume $f = g + h$ with properties above, moreover, ∇g is Lipschitz continuous with constant $L > 0$ & $\text{prox}_{h, \frac{1}{t}}(x)$ can be evaluated. Then, with fixed stepsize $t \leq \frac{1}{L}$

we have $f(x^k) - f^* \leq \frac{\|x^0 - x^*\|^2_2}{2tk}$ (same as gradient descent!)

(proof.) Similarly to proofs for gradient descent. C_L'' for $g \Rightarrow g(x^{k+1}) - g(x^k)$

$\leq -t \nabla g(x^k)^T G_t(x^k) + \frac{Lt^2}{2} \|G_t(x^k)\|^2_2$. By defn of proximity operator, we get

$0 \in \partial h(x^{k+1}) + [x^{k+1} - x^k + t \nabla g(x^k)] \frac{1}{t} \Rightarrow G_t(x^k) - \nabla g(x^k) \in \partial h(x^{k+1})$, by convexity of h ,

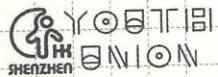
∴ $f(x^{k+1}) - f(x^k) = (g+h)(x^{k+1}) - (g+h)(x^k) \leq (\frac{Lt^2}{2} - t) \|G_t(x^k)\|^2_2 \leq -\frac{t}{2} \|G_t(x^k)\|^2_2$.
 $\stackrel{Lt \leq 1}{\leq}$ (Descent!)

In fact, we can find an even stronger inequality, i.e.,

$$\begin{aligned} f(x^{k+1}) - f(z) &= (g+h)(x^{k+1}) - (g+h)(z) \leq g(x^k) - t \nabla g(x^k)^T G_t(x^k) + \frac{Lt^2}{2} \|G_t(x^k)\|^2_2 - g(z) \\ &+ h(x^{k+1}) - h(z) \stackrel{\text{convex}}{\leq} \frac{Lt^2}{2} \|G_t(x^k)\|^2_2 - t \|G_t(x^k)\|^2_2 + G_t(x^k)^T (x^k - z). \end{aligned}$$

Take $z = x^*$, it's easy to get that $f(x^{k+1}) - f^* \leq G_t(x^k)^T (x^k - x^*) - \frac{t}{2} \|G_t(x^k)\|^2_2$

$\Rightarrow f(x^{k+1}) - f^* \leq \frac{1}{2t} (\|x^k - x^*\|^2_2 - \|x^{k+1} - x^*\|^2_2)$, similarly to gradient-descent



(iii) Acceleration: by accelerating proximal gradient descent, it's possible to achieve the optimal $O(1/\epsilon)$ convergence rate. (4 IDEAs from Nesterov)

One of the ideas: $x^{(0)} = x^{(1)} \in \mathbb{R}^n; \quad \left\{ \begin{array}{l} v^{k-1} = x^{k-1} + \frac{k-2}{k+1}(x^{k-1} - x^{k-2}) \\ x^{(k)} = \text{prox}_{h, \frac{1}{L}}(v^{k-1} - t_k \nabla g(v^{k-1})) \end{array} \right. \quad (*)$

(* v carries some momentum from previous iterations.)

Thm) With the same assumptions as prox-gradient, accelerated proximal gradient method with fixed step $t \leq 1/L$ satisfies $f(x^k) - f^* \leq \frac{2\|x^0 - x^*\|_2^2}{t(k+1)^2} (O(1/\epsilon))$

DDA 4300 Optimization (II) (PS & ML), CUHKSE
(2022-23, Spring)

• Math Formulations of Optimization

$$1) \min_{\vec{x}} f(\vec{x}) \quad (P)$$

subject to $\vec{x} \in X$

Classifications:

- unconstrained** $\rightarrow X = \mathbb{R}^n$ (whole set)
- linear / non-linear** \rightarrow both objective & constraints are affine / otherwise
- convex** $\rightarrow f$ - convex function & X - convex set
- conic linear** \rightarrow linear + \vec{x} in a cone
- (mixed) integer** \rightarrow some of variables are restricted to integer
- Stochastic** $\rightarrow f$ - expected function with random parameters
- fixed-point/ min-max** \rightarrow multiple agents with zero-sum objectives

2) Structured Optimization: Conic Linear Programms (CLP)

$$\begin{aligned} & \min_{\vec{x} \in K} \vec{c}^\top \vec{x} \\ & \text{subject to } A\vec{x} = \vec{b} \quad (A \in \mathbb{R}^{m \times n}, K \text{ is a cone}) \end{aligned} \quad \begin{array}{l} \vec{x} \in K \Rightarrow c\vec{x} \in K, \forall c \geq 0 \\ \text{(assume ext origin)} \end{array}$$

LP-cone

Linear Programming (LP): non-negative constant cone ($\vec{x} \geq 0$)

Second-order Cone Programming (SOCP): second-order cone ($x_1 \geq \|\vec{x}(2:n)\|$)

Semi-definite Cone Programming (SDP): semidefinite matrix cone s.psd.

(Dual Question is introduced) Different Cones give different solns.

e.g. $\begin{bmatrix} x_1 & x_2 \\ x_2 & x_3 \end{bmatrix} \succeq 0$

3) Facility Location Problem $\minimize \sum_j \|\vec{y} - \vec{z}_j\|_p$ $p=1$ Manhattan distance

$$(\Leftrightarrow \minize \sum_j s_j \quad \text{subject to } \vec{y} + \vec{x}_j = \vec{z}_j, \boxed{\|\vec{x}_j\|_p \leq s_j, \forall j} \quad \text{relaxation to be a cone})$$

p-order conic linear programming P-linear cone (POCP)

4) Sparse Linear Regression:

minimize $\|\vec{x}\|_1$, (LASSO)

$$\text{s.t. } A\vec{x} = \vec{b}$$

\Leftrightarrow
of feasible set
(or relaxation)

$$\min \sum_{i=1}^n y_i$$

s.t.

$$A\vec{x} = \vec{b}, -\vec{y} \leq \vec{x} \leq \vec{y}$$

SHENZHEN UNION

linear!

(proof: equivalence of those 3)

$$\min \sum_{i=1}^n x_{1,i} + x_{2,i}$$

s.t.

$$A(\vec{x}_1 - \vec{x}_2) = \vec{b}, \vec{x}_1 \geq 0, \vec{x}_2 \geq 0$$

* Better: minimize $\|\vec{x}\|_p$, ($0 < p < 1$) \rightarrow NOT a norm

$$\text{subject to } A\vec{x} = \vec{b}$$

(or add penalty) $\min \|A\vec{x} - \vec{b}\|^2 + \beta \|\vec{x}\|_p^p$
unconstrained, non-convex

(Note that previously minimize $\|\vec{x}\|_0 = \sum_j |x_j|$ if $x_j \neq 0$)

subject to $A\vec{x} = \vec{b}$ is the aim)

5) Data Classification: support vector machine & logistic regression

(i) A binary-classification SVM: (red data $\vec{a}_i \in \mathbb{R}^d$, $i=1, \dots, n_1$
blue data $\vec{b}_j \in \mathbb{R}^d$, $j=1, \dots, n_2$)

W.T.F (want to find) a hyper-plane, slope \vec{x} and intersect x_0
to separate 2 classes

\Rightarrow Subject to $\vec{a}_i^T \vec{x} + x_0 \geq 1, \forall i$ come from $\begin{cases} > 0 & \text{in optimization} \\ < 0 & \text{scaling} \end{cases} \geq 1$
 $\vec{b}_j^T \vec{x} + x_0 \leq -1, \forall j$ $\begin{cases} < 0 & \text{in optimization} \\ > 0 & \text{scaling} \end{cases} \leq -1$

If strict separation is impossible: minimize $\|\beta\|$ \rightarrow minimize error size

To make it strongly convex \rightarrow unique soln s.t. $\vec{a}_i^T \vec{x} + x_0 + \beta \geq 1, \forall i$

change the objective to $\beta + \mu \|\vec{x}\|^2$ $\vec{b}_j^T \vec{x} + x_0 - \beta \leq -1, \forall j$
 $\beta \geq 0$
(QP)

(ii) Ellipsoidal Separation: Find a form $\vec{y}^T \vec{X} \vec{y} + \vec{y}^T \vec{x} + x_0 = 0$ intersect

minimize $\text{tr}(X) + \|\vec{x}\|^2 \rightarrow$ scaling \rightarrow Hessian slope

s.t. $\vec{a}_i^T \vec{X} \vec{a}_i + \vec{a}_i^T \vec{x} + x_0 \geq 1, \forall i$ (can add β as loss, too)

$\vec{b}_j^T \vec{X} \vec{b}_j + \vec{b}_j^T \vec{x} + x_0 \leq -1, \forall j$

(SDP) $X \succeq 0$ (pos. semi-def)

(iii) Logistic Regression: use function $\frac{e^{\vec{y}^T \vec{x} + x_0}}{1 + e^{\vec{y}^T \vec{x} + x_0}}$ for separation

$$\frac{e^{\vec{a}_i^T \vec{x} + x_0}}{1 + e^{\vec{a}_i^T \vec{x} + x_0}} = \begin{cases} 1, & \text{if } \vec{a}_i \in C \\ 0, & \text{if } \vec{a}_i \notin C \end{cases}$$

Use MLE: to maximize $\prod_{\vec{a}_i \in C} \left(\frac{e^{\vec{a}_i^T \vec{x} + x_0}}{1 + e^{\vec{a}_i^T \vec{x} + x_0}} \right) \prod_{\vec{a}_i \notin C} \left(\frac{1}{1 + e^{\vec{a}_i^T \vec{x} + x_0}} \right)$
Likelihood functn

6) Portfolio Management

(i) minimize $\vec{x}^T V \vec{x}$ \rightarrow proportion on every asset
s.t. $\vec{1}^T \vec{x} \geq \mu$
 $\vec{e}^T \vec{x} = 1$ ($\vec{e} = \vec{I}$)
(Quadratic)

(ii) Robust cases: $\min \max_i \vec{x}^T V_i \vec{x}$

$$\text{s.t. } \min_i \vec{r}_i^T \vec{x} \geq \mu$$

Since V_i, \vec{r}_i can be estimated under various scenarios.

Rewrite the robust portfolio Management

minimize α

$$\Leftrightarrow \text{s.t. } \vec{r}^T \vec{x} \geq \mu, \forall i$$

$$\vec{x}^T V_i \vec{x} \leq \alpha^2, \forall i$$

$$\vec{e}^T \vec{x} = 1.$$

(QCQP)

quadratically constrained quadratic problem

(iii) Portfolio Selection

$$\text{minimize } \vec{x}^T V \vec{x}$$

$$\text{s.t. } \vec{r}^T \vec{x} \geq \mu$$

$$\vec{e}^T \vec{x} = 1$$

$$0 \leq \vec{x} \leq \vec{y}, \vec{e}^T \vec{y} \leq k, \vec{y} \in \{0,1\}^n$$

mixed integer quadratic programming

7) The transportation problem

$$\begin{matrix} d_1 \\ \vdots \\ d_n \end{matrix} \xleftarrow{\text{①}} \begin{matrix} s_1 \\ \vdots \\ s_m \end{matrix} \quad \text{minimize } \sum_{ij} c_{ij} x_{ij}$$

$$\text{s.t. } \sum_j x_{ij} = s_i, \forall i$$

$$\sum_i x_{ij} = d_j, \forall j$$

cost

amount of transport

(can add some null supply points $\Rightarrow m=n$ matching)

The minimal transportation

$$x_{ij} \geq 0$$

some kind of "distance"

cost is called Wasserstein Distance (WD) between 2 distributions \vec{s} and \vec{d}

(After normalization, can be interpreted as 2 probability distributions)

(ii) The [Wasserstein Barycenter Problem] is to find a distribution s.t. sum of its WD to every set of distributions is minimized.

to every set of distributions is minimized. different scenarios to estimate every with prob $\frac{1}{3}$

$$\text{(e.g. minimize } WD_{\vec{s}}(\vec{s}, \vec{d}_1) + WD_{\vec{s}}(\vec{s}, \vec{d}_m) + WD_{\vec{s}}(\vec{s}, \vec{d}_r))$$

"Hierarchy" \vec{s} s.t. $\sum_i s_i = S, s_i \geq 0, \forall i$

Stochastic optimization

"optimal of optimal" kind of "probability center"
a high-level decision! > change/reduce to simple case $\underbrace{\text{write WD}_{\vec{s}} \text{ as } \min}_{\dots} \text{ s.t. } \sim$)

Application = "consensus center of image" LP-formulation $\min WD$

(Image: distribution on pixels) > "statistical median"

(NOT simple algebraic mean!)

(iii) More applications: Combinatorial Auction

To buy a combination $(1, 1, \dots, 0, \dots, 0, 1, \dots)$ with some cost C_i ,

(say, b_i)

(win - give 1 and lose - give 0)

and quantity limit $q_i, \forall i$

The aim is to decide the fill ordered $\vec{x} = (x_i)$:

"option market"

(with quantity $0 \leq x_i \leq q_i$ being sold).

for designer

$$\text{Profit} = \sum_i c_i x_i - p_j \sum_i b_{ij} x_i \quad \begin{array}{l} \text{if No "p" in hand} \\ \text{risk-averse} \end{array}$$

SHENZHEN UNION

maximize the worst case

$$\rightarrow \text{maximize } \sum_i c_i x_i - \max_j (\sum_i b_{ij} x_i)$$

Formulation to LP: maximize $\vec{c}^T \vec{x} - \vec{y}$
 s.t. $B^T \vec{x} - \vec{y} \leq 0$ auxiliary variables
 $\vec{0} \leq \vec{x} \leq \vec{q}$

Dual to the LP: dual variable \vec{y} represents the consensus probability of winners!

customers decide cost to buy & also quantities
 and market developer decide how many to sell.
 kind of confidence of bet!

"Information/Prediction Markets"
 economics

8) Graph Realization & Sensor Network Location

Locate the place (coordinates) by using anchors & lights.

(RD)

location of "i" $\| \vec{x}_i - \vec{x}_j \|_2^2 = d_{ij}^2, \forall (i, j) \in N_A$ → realize of G , with distances between nodes, and also vertexes \vec{x}_i & \vec{x}_j

formulate anchor $\| \vec{x}_k - \vec{x}_j \|_2^2 = d_{kj}^2, \forall (k, j) \in N_A$
 $\min_{\vec{x}, \vec{v}_i} \sum_{(i, j) \in N_A} (\| \vec{x}_i - \vec{x}_j \|_2^2 - d_{ij}^2) + \sum_{(k, j) \in N_A} (\| \vec{x}_k - \vec{x}_j \|_2^2 - d_{kj}^2)$ non-linear, non-convex
 (do matrix relaxation! (SOCP or SDP))

9) Stochastic Optimizations & Learning

(i) Modeling: $\min_{\vec{x} \in X} \mathbb{E}_{\xi} [h(\vec{x}, \xi)]$ ~ some distributions

(ii) Learning with noises/distortions: (robustness)

(iii) Deep learning on NN (neural-network)

weight at layer l

Input vec $\vec{x} \equiv \vec{y}^0$
 Output vec (at layer l) $\vec{y}^l \rightarrow y_j^l = \max \{0, w_{0j}^l + \sum_i w_{ij}^l y_i^{l-1}\}, \forall j, l=1, \dots, L$

Aim: use massive \vec{x} & \vec{y}^L (last-layer) to train reasonable w_{ij}^l

& Net-work verification. $\min_{(\vec{x}, \vec{y}^L)} \|\vec{x} - \vec{y}^L\|^2 \rightarrow$ smallest distortion
 s.t. $\vec{y}^L(\vec{x}) \in$ a convex region, outside of $\vec{y}^L(\vec{x})$

Special case: linearly-constrained quadratic minimization

$$\min_{(\vec{x}, \vec{y}^L)} \|\vec{x} - \vec{y}^L\|^2 + \mu \sum_l \sum_j y_j^l (y_j^L - w_{0j}^l - \sum_i w_{ij}^l y_i^{l-1})$$

s.t. $\vec{y}^L \in$ a convex polyhedron, outside of $\vec{y}^L(\vec{x})$

$$y_j^L \geq w_{0j}^L + \sum_i w_{ij}^L y_i^{l-1}, y_j^L \geq 0, \forall j, l, \vec{y}^0 = \vec{x}$$

10) Reinforcement Learning: Markov Decision / Game Processes ^{sequential decision-making} outcome partly random & partly in control

(i) Defn: Given number of states, i , \rightarrow 2 people decide in turns, zero-sum game.

with A_i - collection of actions, every action $j \in A_i$, is assigned with a cost c_{ij} , and a probability distribution \vec{p}_j to transfer to all possible states at next time period.

A stationary policy $\vec{\pi} = (\pi_1, \dots, \pi_m)$, specifies action $\pi_i \in A_i$, that the decision maker will take, also leading to a cost-to-go value for every i .

(ii) MDP modeling = (discount factor $r \in [0, 1]$)

$$\underset{\pi}{\text{minimize}} \sum_{t=0}^{\infty} r^t \mathbb{E}[C^{T_t}(i_t, i_{t+1})] \quad \xrightarrow{\text{expected discounted sum}}$$

Stochastic game (Zero-sum): $\xrightarrow{\text{cost-to-go with } \pi_{i^*}}$

one partition to minimize, the other to maximize

(iii) Solve MDP: choose the optimal policy s.t. the cost-to-go vector \vec{y}^*

satisfying $y_i^* = \min_j \{c_j + r \vec{p}_j^T \vec{y}^*, \forall j \in A_i\}, \forall i$ ^{saddle-point optimization}
 optimal policy $\rightarrow (\pi_i^* = \arg \min_j \{c_j + r \vec{p}_j^T \vec{y}^*, \forall j \in A_i\}, \forall i)$

$$y_i^* = \min_j \{c_j + r \vec{p}_j^T \vec{y}^*, \forall j \in A_i\}, \forall i \in I^-$$

$$y_i^* = \max_j \{c_j + r \vec{p}_j^T \vec{y}^*, \forall j \in A_i\}, \forall i \in I^+$$

Reformulated as LP: $\underset{\vec{y}}{\text{maximize}} \vec{e}^T \vec{y} \rightarrow$ to make \vec{y} be saddle

$$\text{subject to } y_i - r \vec{p}_i^T \vec{y} \leq c_j, j \in A_i$$

$$\vdots$$

$$y_m - r \vec{p}_m^T \vec{y} \leq c_j, j \in A_m$$

[Thm] when $\vec{y}^T \vec{e}$ is maximized, there must be at least 1

inequality constraint in A_i which becomes equal, for every state i .

i.e. \vec{y}^* is a fixed-point soln.

• Math Foundations of Optimization

1) Cone: C is a cone if $x \in C$ implies $\alpha x \in C$. \forall scalar $\alpha \geq 0$

Dual of a cone $C^* := \{y | \langle x, y \rangle \geq 0, \forall x \in C\}$ ^{closed convex cone}

$(C^*)^*$ is the closure of convex hull of C

(e.g. Polyhedral cone $\{ \vec{x} | A\vec{x} \leq \vec{0} \}$

finite extreme points $\{ \vec{x} | A\vec{x} = \vec{0} \}$ conic combination

p-order cone $\{ (t, \vec{x}) | t \geq \| \vec{x} \|_p \}$ dual p'-order cone $\{ (t, \vec{x}) | t \geq \| \vec{x} \|_{p'} \}$

p, p' are dual exponent ($p + p' = 1$) \rightarrow proved by Hölder's inequality

Extreme points of convex set: points which cannot be expressed by convex combinations by other 2 points.

2) Convex functions: $f \circ \phi$ is convex, if ϕ is convex & f is convex, non-decreasing

(e.g.) $\{ \vec{x}(\vec{b}) := \underset{\vec{x}}{\text{minimize}} f(\vec{x}) \text{ s.t. } A\vec{x} = \vec{b}, \vec{x} \geq 0 \}$ is convex, where f is convex
 $(\log(1 + e^{\vec{a}^T \vec{x} + x_0}) \text{ is convex.})$ proof by defn $\vec{x}(\vec{b}) \triangleq f(\vec{x}(\vec{b}_1)) \leq f(\vec{x}(\vec{b}_1) + \vec{x}(\vec{b}_2))$
 $\vec{x}^{\text{composition}}$ unique soln. $\vec{x}(\vec{b}_1 + \vec{b}_2)$ \uparrow sub-additivity (inside f)

3) Carathéodory Theorem for generation of cone: & convexity of f

Given $A \in \mathbb{R}^{m \times n}$, let convex polyhedral cone $C = \{ A\vec{x} | \vec{x} \geq \vec{0} \}$.

For any $\vec{b} \in C$, $\vec{b} = \sum_{i=1}^d \vec{a}_{j_i} x_{j_i}$, $x_{j_i} \geq 0 \forall i$, $\vec{a}_{j_1}, \dots, \vec{a}_{j_d}$ are l.indep. from $\vec{a}_1, \dots, \vec{a}_n$

reduced/mean representation

$\{ \vec{x} | A\vec{x} = \vec{b}, \vec{x} \geq \vec{0} \}$, basic soln \vec{x}_B exists, due to above thm.

BFS must be an extreme or corner point, as a direct corollary from Carathéodory's.

Hyper-plane: $H \triangleq \{ \vec{x} | \vec{a}^T \vec{x} = b \}$

* Separating hyperplane theorem: Let C be a closed, convex set. \vec{b} is NOT in C , then

$\exists \vec{y} \text{ s.t. } \langle \vec{b}, \vec{y} \rangle > \sup_{\vec{x} \in C} \langle \vec{x}, \vec{y} \rangle$

($\Leftrightarrow \exists H \triangleq \{ \vec{z} | \langle \vec{y}, \vec{z} \rangle = s \}$, s.t. \vec{b} & C are separated)

(proof.) omitted. Note that consider dual by H .

$\{ \alpha \vec{b} | \alpha > 0 \}^* \notin C^*$, otherwise $(C^*)^*$ contains \vec{b} , contradictory with closedness & convexity.

4) Farka's lemma: proved from separating hyperplane theorem.

(Alternating system: $\{ \vec{x} | A\vec{x} = \vec{b}, \vec{x} \geq 0 \}$ & $\{ \vec{y} | A^T \vec{y} \leq 0, \vec{b}^T \vec{y} > 0 \}$, only 1 (exactly) is non-empty)

Farka's Lemma for general closed convex cones:

$\{ \vec{x} | A\vec{x} = \vec{b}, \vec{x} \in K \}$ & $\{ \vec{y} | -A^T \vec{y} \in K^*, \vec{b}^T \vec{y} = 0 \}$ cannot be alternating system,

generally. (Counter-e.g. (neither non-empty): consider $K = S_+^2$ (p.s.d, 2×2)

$A := (a_1, a_2)$, with $a_1 = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$. $a_2 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$

$A\vec{x} = \begin{bmatrix} \text{tr}(a_1 \vec{x}) \\ \text{tr}(a_2 \vec{x}) \end{bmatrix}$ (or $\begin{bmatrix} \langle a_1, \vec{x} \rangle \\ \langle a_2, \vec{x} \rangle \end{bmatrix}$) and $b = \begin{pmatrix} 0 \\ 2 \end{pmatrix}$

Note that, here $C := \{ A\vec{x} | \vec{x} \in K \}$ is NOT closed.)

* Let K be a closed convex cone, if $\exists \vec{y}$ s.t. $-\vec{y}^T A \in \text{int } K^*$,

then $C := \{Ax \mid x \in K\}$ is a closed convex cone.

Consequently, $\{Ax = \vec{b}, x \in K\} \& \{-\vec{y}^T A \in K^*, \vec{b}^T \vec{y} = 1 \text{ (or } > 0)\}$

are an alternative pair.

5) Dual of Conic LP:

(CL, P) minimize $\langle C, X \rangle$

subject to $\langle a_i, X \rangle = b_i, \forall i = 1, \dots, n$

$X \in K$ *closed, pointed*

(CL, D) maximize $\vec{b}^T \vec{y}$

subject to $\sum_{i=1}^m y_i a_i + s = c$

$s \in K^*$

(Theorem) The dual of the dual is the primal (in conic linear programming)

(i) Dual of transportation problem:

(P) minimize $\langle C, X \rangle = \text{tr}(C^T X)$

s.t. $\langle E_j, X \rangle = s_j$

$\langle E_i^T, X \rangle = d_i$

$X \in \text{non-negative constant cone}$

(D) maximize $\vec{u}^T \vec{s} + \vec{v}^T \vec{d}$

another way to change

s.t. $\sum_i u_i E_i + \sum_j v_j E_j^T \leq C$

$(\Leftrightarrow u_i + v_j \leq c_{ij}, \forall i, j)$ *no more than previous charge*

Interpretation: charge u_i when sending 1 unit from i , & ... v_j ... accepting 1 ... from j

(ii) The dual of MDP-LP.

(D) minimize $\sum_i \sum_{j_k \in A_i} c_{jk} x_{jk}$ *state action fluency / flux*
 subject to $\sum_i \sum_{j_k \in A_i} (r_{ik} - r p_{jk}) x_{jk} = \bar{r}$ *expected PV of # times that in i, taking action jk.*
 $(\text{flow in} - \text{flow out})$ *period 1 period 2 ...*
 $x_{jk} \geq 0, \forall i, j_k \in A_i$ *with price*

(iii) Nash-Equilibrium: (2-person zero-sum game)

e.g., P (payoff matrix) = $\begin{bmatrix} +3 & -1 & -4 \\ -3 & +1 & +4 \end{bmatrix}$ (2 people, column & row players)

random strategies: a vector of probabilities (with states)

* Nash Equilibrium: pure strategy — the expected payout (dominating) \geq payout for the other action;

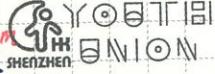
randomized strategy — the expected payout for each part \rightarrow other's choice of the action (in the strategy must be the same), & these \geq actions that are not part of the action.

"Column" Strategy:

$$\begin{aligned} & \underset{\vec{x}}{\text{maximize}} && V \\ & \text{s.t.} && \vec{v}\vec{e} \leq P\vec{x} \\ & && \vec{e}^T \vec{x} = 1 \end{aligned}$$

due to

Nash Equilibrium



"Row" strategy:

$$\begin{aligned} & \underset{\vec{y}}{\text{minimize}} && U \\ & && \vec{u}\vec{e} \geq P^T\vec{y} \\ & && \vec{e}^T \vec{y} = 1 \\ & && \vec{y} \geq 0 \end{aligned}$$

dual equal to each other

6) Duality Theorems for CLP.

(Weak Duality) $\langle C, x \rangle - \vec{b}^T \vec{y} = \langle x, s \rangle \geq 0$ (for CLP), \forall feasible x of (CLP), and (\vec{y}, s) of (CLD)

* Duality gap

(Strong Duality) Let $x^* \in F_p$ and $(\vec{y}^*, s^*) \in F_d$. Then, $\langle C, x^* \rangle = \vec{b}^T \vec{y}^*$, implies that x^* is optimal for (CLP) & (\vec{y}^*, s^*) ... for (CLD)

(proof.) Let x^* be a minimizer of (LP). Then, the following system:

$$\left\{ \begin{array}{l} \text{for LP} \quad A\vec{x}' - \vec{b}c = \vec{0} \\ \quad \quad \quad \vec{x}' \geq \vec{0} \\ \quad \quad \quad \vec{c}^T \vec{x}' - (\vec{c}^T \vec{x}^*)c = -1 < 0 \end{array} \right. \quad \begin{array}{l} \text{has no solution, for otherwise } \vec{x}'/c \text{ is feasible } (c > 0) \\ \quad \quad \quad "(\vec{x}', c)" \\ \quad \quad \quad \text{for (LP)} \quad \left(\begin{array}{l} \min \vec{c}^T \vec{x} \\ \text{s.t. } A\vec{x} = \vec{b} \\ \vec{x} \geq 0 \end{array} \right) \end{array}$$

$\vec{c}^T \vec{x}'/c < \vec{c}^T \vec{x}^*$, contradictory!

if $c = 0$, $\vec{x}' + \vec{x}^*$ is feasible for (LP), $\vec{c}^T(\vec{x}' + \vec{x}^*) = \vec{c}^T \vec{x}^* - 1 < \vec{c}^T \vec{x}^*$, contra!

∴ With Farka's lemma, $\exists \vec{y}^* \text{ s.t. } \vec{c} - A^T \vec{y}^* \geq 0, -\vec{c}^T \vec{x}^* + \vec{b}^T \vec{y}^* \geq 0$.

use weak-duality $\Rightarrow \exists$ such optimal \vec{y}^* .

For general CLP: \exists duality gap in some cases, i.e., strong duality does NOT always hold.

e.g. In general CLP, $C = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$, $a_1 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$, $a_2 = \begin{bmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 2 \end{bmatrix}$ & $\vec{b} = \begin{bmatrix} 0 \\ 2 \end{bmatrix}$

$K = S_+^3$ (SDP-cone)

(Zero-duality gap but NOT attainable)

e.g. $C = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$, $a_1 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$ & $b_1 = 2$. $K = S_+^2$. (The dual does NOT have interior)

condition

(proof.) (i) Let F_p be non-empty & has an interior, \underline{z}^* be its infimum.

for CLP Consider alternative system = $\vec{y}^T A + S = C$ \uparrow 0-duality gap but may not attainable

$$\left\{ \begin{array}{l} Ax - \vec{b}c = \vec{0} \\ \langle C, x \rangle - \underline{z}^* c < 0 \\ (x, c) \in K \times R_+ \end{array} \right. \quad \left\{ \begin{array}{l} \vec{y}^T A + S = C \\ -\vec{b}^T \vec{y} + K = -\underline{z}^* \\ (S, K) \in K^* \times R_+ \end{array} \right. \quad \begin{array}{l} \exists \vec{y}^*, S^* \& \\ \text{use weak duality, done} \end{array}$$

↑ infeasible

(ii) WTS: $\exists x \in F_p$ s.t. $\langle C, x \rangle = \underline{z}^*$ (attainable infimum) $\Leftrightarrow F_d$ feasible, $\text{int } F_d \neq \emptyset$

& \underline{z}^* is the sup of (CLD)

\uparrow 0-duality gap & attainable

↑ more condition

iii) Suppose $f_d = \emptyset$. f_p is feasible & $\text{int } f_p \neq \emptyset$. Then, $\exists \bar{x} \in \text{int } K, \bar{\epsilon} \geq 0$

$$\text{s.t. } A\bar{x} - b\bar{\epsilon} = 0, (\bar{x}, \bar{\epsilon}) \in \text{int}(K \times R_+).$$

$\forall \bar{z}^*$, consider again the alternative pair above. the latter becomes infeasible

\therefore The former has soln $(\bar{x}_{\bar{z}^*}, \bar{\epsilon}_{\bar{z}^*})$, $\forall \bar{z}^*$. If $\bar{\epsilon}_{\bar{z}^*} > 0$, $\langle C, \frac{\bar{x}_{\bar{z}^*}}{\bar{\epsilon}_{\bar{z}^*}} \rangle < \bar{z}^*$, if $\bar{\epsilon} = 0$
 $\bar{x} + \alpha \bar{x}_{\bar{z}^*}$ is feasible if \bar{x} is feasible, obj $\mapsto -\infty$ as $\alpha \mapsto \infty$ $\uparrow \bar{z}^*$
 (unbdd)

Summary: (i) Both (CLP) & (CLD) feasible \Rightarrow one has interior - 0-duality gap
 both have interior - attainable opt. soln
 (ii) One of (CLP), (CLD) infeasible,
 another is feasible & has interior \Rightarrow it is unbdd.

Possible relations of SDP-LP:

	Finite opt.	Unbdd	Infeasible
Finite opt.	✓	"new from LP"	✓
Unbdd	—	✓	✓
Infeasible	✓	✓	✓

Rules for taking dual of CLP:

$$(CLP) \min \sum_k \langle C_k, x_k \rangle \\ \text{s.t. } \sum_k A_k x_k = b \\ x_k \in K_k, \forall k$$

$$(CLD) \min b^T \bar{y} \\ \text{s.t. } \bar{y}_k^T A_k + s_k = C_k, \forall k \\ s_k \in K_k^*, \forall k$$

Optimality conditions: $\langle X, S \rangle = \text{tr}(X^T S) = 0$

(for SDP) $\left\{ \begin{array}{l} AX = b \\ -\bar{y}_k^T A - S = -C \\ X, S \succeq 0 \end{array} \right.$

$\Downarrow \langle X, S \rangle_{ij} = 0, \forall i, j$
 (under SDP cone, could be proved)

* Equivalence of Convex Optimization & CLP:

$$(CO) \min_{\vec{x}, \alpha} \alpha \\ \text{s.t. } C_i(\vec{x}) - \alpha \leq 0 \\ C_k(\vec{x}) \leq 0, \forall k$$

the same form

$$\min_{\vec{x}} \vec{c}^T \vec{x} \\ \text{s.t. } C_i(\vec{x}) \leq 0, \forall i$$

CLP form

$$\min_{\vec{z}, \vec{x}} \begin{bmatrix} 0 \\ \vec{z} \end{bmatrix}^T \begin{bmatrix} \vec{x} \\ \vec{z} \end{bmatrix} \\ \text{s.t. } \begin{bmatrix} 1 \\ \vec{z} \end{bmatrix}^T \begin{bmatrix} \vec{x} \\ \vec{z} \end{bmatrix} = 1 \\ (\vec{z}, \vec{x}) \in \bigcap_{k=1}^m K_k$$

closure cone

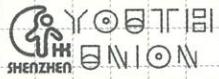
where $K_i := \{(\vec{z}, \vec{x}) \mid \vec{z} \geq 0 \text{ & } C_i(\vec{x}/\vec{z}) \leq 0\}$

How to construct the dual cone: (of $\cap K_i$)

i.e. $\{(K, S) | K \subset \bar{S}^T \bar{x} \geq 0, V(\bar{x}, \bar{s}) \in \{\bar{c}_i\}\}$, w.r.t.o.g. let $i = 1$.

$\Leftrightarrow \{(K, S) | K \subset \bar{S}^T \bar{x} \geq 0, \forall \bar{x} \text{ s.t. } c_i(\bar{x}) \leq 0, \forall i\}$

consider $\varphi(\bar{s}) := \inf \{ \bar{S}^T \bar{x} | c_i(\bar{x}) \leq 0, \forall i\}$, the dual cone $K^* = \{(K; \bar{s}) | K + \varphi(\bar{s}) \geq 0\}$



* Robust Optimization: application of duality choice of adversary

$$\text{Robust Min-max model} \quad \min_{\bar{x}} \max_{\bar{u}} \begin{cases} \bar{u}^T \bar{x} \\ (\bar{u} \geq 0, \bar{u} \leq \bar{e}) \end{cases} \quad \begin{matrix} (\bar{c} + \bar{C}\bar{u})^T \bar{x} \\ \text{s.t. } A\bar{x} = b; \bar{x} \geq 0 \end{matrix} \quad \text{choice of player}$$

(change adversary's decision to its dual)

$$\text{Robust Min-min model} \quad \min_{\bar{x}} \min_{\bar{y}} \begin{cases} \bar{c}^T \bar{x} + \bar{e}^T \bar{y} \\ \bar{y} \geq \bar{c}^T \bar{x}, \bar{y} \geq 0 \end{cases} \quad \begin{matrix} \text{put together, we get a pure} \\ \text{opt. question (min...)} \end{matrix}$$

Combinatorial Auction: the meaning of dual price information

$$(D\text{-CAP}) \quad \max_{\bar{p}, \bar{y}} \bar{q}^T \bar{y} \quad \begin{matrix} \text{decision / fraction} \\ \text{state price (consensus probability)} \end{matrix}$$

s.t. $\bar{B}(\bar{p} + \bar{y}) \geq \bar{c}$
 $\bar{e}^T \bar{p} = 1$
 $(\bar{p}, \bar{y}) \geq 0$

$$\star \text{Online L-P:} \quad \max_{\bar{x}} \sum_{t=1}^n \pi_t(x_t) \quad \begin{matrix} \text{decision variable} \\ \text{total amount} \end{matrix}$$

s.t. $\sum_{t=1}^n a_{it} x_t \leq b_i, \forall i$ only know (n, b) at the start.
 $0 \leq x_t \leq 1, \forall t$ Data of x_t is revealed sequentially

An algorithm A is c -competitive iff $E_0 \left[\sum_{t=1}^n \pi_t x_t(s, A) \right] \geq c \cdot \text{opt}(A, \pi), V(A, \pi)$

Strategy: $x_t = 0, 1 \leq t \leq \bar{t}$, solve $\max_{\bar{x}} \sum_{t=1}^{\bar{t}} \pi_t x_t$ off-line

$$\text{s.t. } \sum_{t=1}^{\bar{t}} a_{it} x_t \leq b_i, \forall i$$

\Rightarrow get dual optimal soln: \bar{p} Determine, future $x_t = \begin{cases} 0, & \text{if } \pi_t \leq \bar{p} \\ 1, & \text{otherwise} \end{cases}$

as long as $a_{it} x_t \leq b_i - \sum_{j=1}^{t-1} a_{ij} x_j, \forall i$, otherwise $x_t = 0$

7) Support-Size & Rank of CLP Solutions & Applications

(i) Let x^*, s^* be optimal solutions with 0-duality gap:

$$|\text{supp}(x^*)| + |\text{supp}(s^*)| \leq n \leftarrow \text{prob. size.}$$

[Strict Complementarity] If $(LP), (LD)$ are both feasible, \exists a pair of strict complementarity solutions $x^* \in F_p$ & $(y^*, s^*) \in F_d$ st. $\langle x^*, s^* \rangle = 0$

$$\& |\text{supp}(x^*)| + |\text{supp}(s^*)| = n$$