

Lecture Notes for MAT3220

March 5, 2019

Abstract

This document is the typed lecture notes for MAT3220

1 Introduction to Linear Programming

1.1 Preliminaries

Standard Form We usually consider the standard linear programming (LP) model:

$$\begin{aligned} \max \quad & \sum_{j=1}^n c_j x_j \\ \text{s.t.} \quad & \sum_{j=1}^n a_{ij} x_j \leq b_i, \quad i = 1, \dots, m \\ & x_j \geq 0, j = 1, \dots, n \end{aligned} \tag{1.1}$$

Or more generally, the constraints with equalities:

$$\begin{aligned} \min \quad & \sum_{j=1}^n c_j x_j \\ \text{s.t.} \quad & \sum_{j=1}^n a_{ij} x_j \leq b_i, \quad i \in I \\ & \sum_{j=1}^n a_{ij} x_j = b_i, \quad i \in E \\ & x_j \geq 0, j = 1, \dots, n \end{aligned}$$

It's often convenient to write the LP (1.1) into the compact matrix form:

$$\begin{aligned} \max \quad & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{A} \mathbf{x} \leq \mathbf{b} \\ & \mathbf{x} \geq \mathbf{0} \end{aligned} \tag{1.2}$$

where $\mathbf{c} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$.

We also write \mathbf{A} as the column form:

$$\mathbf{A} = (\mathbf{a}_1, \quad \dots \quad \mathbf{a}_n)$$

where \mathbf{a}_i is the i -th column of \mathbf{A} . We also express the submatrix of \mathbf{A} , i.e., $\mathbf{A}_I \subset \mathbf{A}$ as:

$$\mathbf{A}_I := [a_i \mid i \in I],$$

where I is a subset of $\{1, 2, \dots, n\}$.

Dictionaries of an LP We can introduce slack variables to transform (1.1) into LP with equalities:

$$x_{n+i} := b_i - \sum_{j=1}^n a_{ij}x_j, \quad i = 1, \dots, m$$

Let $z = \sum_{j=1}^n c_j x_j$ be the objective function, and therefore we obtain a *dictionary* for the LP (1.1):

$$\left. \begin{aligned} x_{n+i} &= b_i - \sum_{j=1}^n a_{ij}x_j, \quad i = 1, \dots, m \\ z &= \sum_{j=1}^n c_j x_j \end{aligned} \right\} \quad \text{Dictionary} \quad (1.3)$$

Assume that $b_i \geq 0$ for $i = 1, \dots, m$. Therefore we obtain a *feasible solution* associated with the dictionary, say *dictionary solution*:

$$x_j = 0, \text{ for } j = 1, \dots, n \quad x_{n+i} = b_i \text{ for } i = 1, \dots, m$$

It's clear how to improve the current dictionary solution:

- If $c_j \leq 0, \forall j$, then we cannot possibly improve the dictionary solution
- If $c_j > 0$ for some $1 \leq j \leq n$, we increase the value for x_j from 0 into maximal value, while fixing $x_j = 0$ for $1 \leq k (\neq j) \leq n$. Keep implementing until $c_j \leq 0, \forall j$.

Example 1.1.1. Consider the optimization problem

$$\begin{aligned} \max \quad & 5x_1 + 4x_2 + 3x_3 \\ \text{s.t.} \quad & 2x_1 + 3x_2 + x_3 \leq 5 \\ & 4x_1 + x_2 + 2x_3 \leq 11 \\ & 3x_1 + 4x_2 + 2x_3 \leq 8 \\ & x_1 \geq 0, \quad x_2 \geq 0, \quad x_3 \geq 0 \end{aligned} \quad (1.4a)$$

We can find its dictionary:

$$\begin{aligned} x_4 &= 5 - 2x_1 - 3x_2 - x_3 \\ x_5 &= 11 - 4x_1 - x_2 - 2x_3 \\ x_6 &= 8 - 3x_1 - 4x_2 - 2x_3 \\ z &= 0 + 5x_1 + 4x_2 + 3x_3 \end{aligned} \quad (1.4b)$$

Since $c_1 > 0$, increasing value for x_1 suffices to consider the dictionary below instead:

$$\begin{aligned} x_1 &= \frac{5}{2} - \frac{3}{2}x_2 - \frac{1}{2}x_3 - \frac{1}{2}x_4 & x_1 &= \frac{5}{2} - \frac{3}{2}x_2 - \frac{1}{2}x_3 - \frac{1}{2}x_4 \\ x_5 &= 11 - 4x_1 - x_2 - 2x_3 & x_5 &= 1 + 5x_2 + 0x_3 + 2x_4 \\ x_6 &= 8 - 3x_1 - 4x_2 - 2x_3 & x_6 &= \frac{1}{2} + \frac{1}{2}x_2 - \frac{1}{2}x_3 + \frac{3}{2}x_4 \\ z &= 0 + 5x_1 + 4x_2 + 3x_3 & z &= \frac{25}{2} - \frac{7}{2}x_2 + \frac{1}{2}x_3 - \frac{5}{2}x_4 \end{aligned} \quad (1.4c)$$

Also, since $c_3 > 0$, increasing value for x_3 suffices to consider the dictionary below instead:

$$\begin{array}{ll}
x_1 = \frac{5}{2} - \frac{3}{2}x_2 - \frac{1}{2}x_3 - \frac{1}{2}x_4 & x_3 = 1 + x_2 + 3x_4 - 2x_6 \\
x_5 = 1 + 5x_2 + 0x_3 + 2x_4 & x_1 = 2 - x_2 - 2x_4 + x_6 \\
x_3 = 1 + x_2 + 3x_4 - 2x_6 & x_5 = 1 + 5x_2 + 2x_4 + 2x_6 \\
z = \frac{25}{2} - \frac{7}{2}x_2 + \frac{1}{2}x_3 - \frac{5}{2}x_4 & z = 13 - 3x_2 - x_4 - x_6
\end{array} \quad \Longleftrightarrow \quad (1.4d)$$

1.2 Simplex Method

Notations The general dictionary for the problem (1.1) can be expressed as:

$$\begin{aligned}
x_i &= \bar{b}_i - \sum_{j \in N} \bar{a}_{ij}x_j, \quad i \in B \\
z &= \zeta - \sum_{j \in N} \bar{c}_jx_j
\end{aligned} \quad (1.5)$$

where

1. the set B is called a *basis*, with $|B| = m$
2. the set N is called a *non-basis*, with $|N| = n - m$. Moreover, $B \cup N = \{1, \dots, n\}$.
3. the basis B is said to be *primal feasible* if $\bar{\mathbf{b}} \geq 0$, since in this case we can choose a primal feasible solution by setting non-basis variables to be zero and basis variables x_i to be \bar{b}_i .
4. the non-basis N is said to be *dual feasible* if $\bar{\mathbf{c}} \leq 0$, since in this case we can choose a dual feasible solution by setting non-basis variables to be \bar{c}_i and other variables to be 0.

One can verify that in (1.5),

$$\begin{aligned}
\bar{\mathbf{b}} &= \mathbf{A}_B^{-1}\mathbf{b} \\
\bar{\mathbf{c}}_N^T &= \mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{A}_B^{-1} \mathbf{A}_N \\
\bar{\mathbf{A}}_N &= \mathbf{A}_B^{-1} \mathbf{A}_N \\
\zeta &= \mathbf{c}_B^T \mathbf{A}_B^{-1} \mathbf{b}
\end{aligned}$$

One can verify that $(\mathbf{x}_B, \mathbf{x}_N) = (\mathbf{A}_B^{-1}\mathbf{b}, \mathbf{0})$ is a *basic solution*.

Algorithm 1 Framework for the one step of the Simplex Method

Input:

Primal feasible basic solution;

Output:

Improved feasible basic solution;

1: Find Entering Basis Variable j

- Search for $j \in N$ such that $\bar{c}_j > 0$
- If none exists then the current basic solution is optimal; otherwise choose one of such j .

2: Find Leaving Basis Variable i

- Search for $i \in B$ such that $\bar{a}_{ij} > 0$
- If none exists then the problem is unbounded; otherwise choose

$$i \in \arg \min \left\{ \frac{\bar{b}_i}{\bar{a}_{ij}} : \bar{a}_{ij} > 0, i \in B \right\}$$

3: Basis Update: $B \leftarrow B \cup \{j\} \setminus \{i\}$, and then form the corresponding basic solution.

Simplex Method Algorithm The assumption for the working of simplex method is that we are given a primal feasible basic solution, i.e., $\bar{b} \geq 0$. The framework for obtaining an improved solution is summarized in (1).

Remark 1.2.1. The *one-step* of the simplex method is also called a *pivot step*, i.e., choose one pivot variable entering the basis and one leaving the basis.

The objective value for a successful pivot is improved by $\frac{\bar{c}_j \bar{b}_i}{\bar{a}_{ij}}$. However, the simplex method may not necessarily increase the objective value at each pivot, e.g., the case $\bar{b}_i = 0$ could happen. In this case, the basic solution is said to be *degenerate*.

Since there are no more than $\binom{n}{m}$ (finite) possible bases, the simplex method will stop on two cases: (a) declaring the problem is unbounded; (b) finding a basic optimal solution.

Pivot Rules The *simplex method* specializes into a *simplex algorithm* if one specifies a *pivot rule* to determine which one variable to enter the basis and which one to leave, when there is a choice to make. Note that there exists some pivot rules that will make the problem face into cycling circumstance (see the example below), but here we list some examples of pivot rules that will be shown to definitely avoid cycling circumstance:

- Dantzig's pivot rule: *choose the largest positive coefficient to enter the basis.*
- The maximum improvement rule: *try all the combinations and pick the pivot pair with the largest improvement.*
- Bland's rule: *Among the candidates always pick the one with the smallest index.*

Example 1.2.1. This example shows that some pivot rules may let the problem face into cycling circumstance, i.e., the algorithm solves the problem in a loop and fails to go out:

$$\mathbf{x}_5 = -0.5x_1 + 5.5x_2 + 2.5x_3 - 9x_4$$

$$x_6 = -0.5x_1 + 1.5x_2 + 0.5x_3 - x_4$$

$$x_7 = 1 - x_1$$

$$z = \mathbf{10x_1} - 57x_2 - 9x_3 - 24x_4$$

Choosing x_1 to enter the basis and x_5 to leave gives:

$$x_1 = -2x_5 + 11x_2 + 5x_3 - 18x_4$$

$$\mathbf{x_6} = x_5 - 4x_2 - 2x_3 + 8x_4$$

$$x_7 = 1 + 2x_5 - 11x_2 - 5x_3 + 18x_4$$

$$z = -20x_5 + \mathbf{53x_2} + 41x_3 - 204x_4$$

Choosing x_2 to enter the basis and x_6 to leave gives:

$$\mathbf{x_1} = 0.75x_5 - 2.75x_6 - 0.5x_3 + 4x_4$$

$$x_2 = 0.25x_5 - 0.25x_6 - 0.5x_3 + 2x_4$$

$$x_7 = 1 - 0.75x_5 - 13.25x_6 + 0.5x_3 - 4x_4$$

$$z = -6.75x_5 - 13.25x_6 + \mathbf{14.5x_3} - 98x_4$$

Choosing x_3 to enter the basis and x_1 to leave gives:

$$x_3 = 1.5x_5 - 5.5x_6 - 2x_1 + 8x_4$$

$$\mathbf{x_2} = -0.5x_5 + 2.5x_6 + x_1 - 2x_4$$

$$x_7 = 1 - x_1$$

$$z = 15x_5 - 93x_6 - 29x_1 + \mathbf{18x_4}$$

Choosing x_4 to enter the basis and x_2 to leave gives:

$$\mathbf{x_3} = -0.5x_5 + 4.5x_6 + 2x_1 - 4x_2$$

$$x_4 = -0.25x_5 + 1.25x_6 + 0.5x_1 - 0.5x_2$$

$$x_7 = 1 - x_1$$

$$z = \mathbf{10.5x_5} - 70.5x_6 - 20x_1 - 9x_2$$

Choosing x_5 to enter the basis and x_3 to leave gives:

$$\begin{aligned}x_5 &= 9x_6 + 4x_1 - 8x_2 - 2x_3 \\ \mathbf{x}_4 &= -x_6 - 0.5x_1 + 1.5x_2 + 0.5x_3 \\ x_7 &= 1 - x_1 \\ z &= \mathbf{24x}_6 + 22x_1 - 93x_2 - 21x_3\end{aligned}$$

Choosing x_6 to enter the basis and x_4 to leave gives the *same dictionary as we started*:

$$\begin{aligned}\mathbf{x}_5 &= -0.5x_1 + 5.5x_2 + 2.5x_3 - 9x_4 \\ x_6 &= -0.5x_1 + 1.5x_2 + 0.5x_3 - x_4 \\ x_7 &= 1 - x_1 \\ z &= \mathbf{10x}_1 - 57x_2 - 9x_3 - 24x_4\end{aligned}$$

Theorem 1.2.1. Bland's pivot rule would avoid cycling.

Proof. We show this claim by contradiction. If Bland's pivot rule produces cycling, let's study one cycle. For a sequence of dictionaries that form a cycle, let's delete all the variables that neither leave nor enter the basis, then it will remain a cycle.

In all these dictionaries, all \bar{b}_i will be zero, since otherwise the objective value will be strictly increased.

Let's study the tableau of dictionaries. It's a matrix that stores all the coefficients of a dictionary:

\mathbf{I}_B	$\mathbf{A}_B^{-1}\mathbf{A}_N$	$\bar{\mathbf{b}}$
$\mathbf{0}^T$	$\bar{\mathbf{c}}_N^T$	$\mathbf{c}_B^T \mathbf{A}_B^{-1} \mathbf{b}$

Two vectors are of special interest. The last row of the tableau in left part can be written as

$$\bar{\mathbf{c}}^T = \mathbf{c} - \mathbf{c}_B^T \mathbf{A}_B^{-1} \mathbf{A}$$

For the chosen $j \in N$, the direction

$$\mathbf{d}_i^{(j)} = \begin{cases} -\bar{a}_{ij}, & i \in B \\ 0, & i \neq j \\ 1, & i = j \end{cases}$$

It's clear that $\bar{\mathbf{c}}^T \mathbf{d}^{(j)} = \bar{c}_j$.

Suppose that ℓ is the largest index of all variables that are involved in the cycle. Let (B, N) be the pivot where ℓ was about to enter the basis, $\mathbf{v} = \bar{\mathbf{c}}$ be the last row for that tableau at that point; let (B', N') be the pivot where ℓ was about to leave the basis, and k was to enter the basis at that point, $\mathbf{d}^{(k)}$ be the corresponding direction vector, \mathbf{u} the last row of that tableau.

It's clear that

- \mathbf{v} is everywhere non-positive except for one position $v_\ell > 0$
- $\mathbf{d}^{(k)}$ is everywhere non-negative except for one position $d_\ell^{(k)} < 0$

Moreover, $\mathbf{v} - \mathbf{u} \in \mathcal{R}(\mathbf{A}^T)$ and $\mathbf{d}^{(k)} \in \mathcal{N}(\mathbf{A})$, which implies

$$0 = (\mathbf{v} - \mathbf{u})^T \mathbf{d}^{(k)} = \mathbf{v}^T \mathbf{d}^{(k)} - \mathbf{u}_k < 0,$$

which is a contradiction. \square

Lemma 1.2.2. Given the condition that the LP (1.2) has one basic feasible solution, then the LP (1.2) with perturbations, i.e.,

$$\begin{aligned} \max \quad & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{A}\mathbf{x} \leq \mathbf{b} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_m \end{pmatrix} \\ & \mathbf{x} \geq \mathbf{0} \end{aligned} \tag{1.6}$$

will face no degeneracy for $\forall \boldsymbol{\varepsilon} \in (\mathbf{0}, \boldsymbol{\varepsilon}_1)$ for some $\varepsilon_1 > 0$.

Proof. For any basis B , the feasible solution for LP (1.6) is $\mathbf{A}_B^{-1}(\bar{\mathbf{b}} + \boldsymbol{\varepsilon})$. Suppose its i -th component is zero, i.e., $0 + 0\varepsilon_1 + \dots + 0\varepsilon_m$.

However, its i -th component is $\mathbf{e}_i^T \mathbf{A}_B^{-1}(\bar{\mathbf{b}} + \boldsymbol{\varepsilon})$, which implies $\mathbf{e}_i^T \mathbf{A}_B^{-1} = \mathbf{0}$, which is a contradiction. \square

Question: what's the conclusion for page 19 in slides 1?

Two-Phase Simplex Method Given a dictionary

$$x_{n+i} = b_i - \sum_{j=1}^n a_{ij}x_j, \quad i = 1, \dots, m,$$

with some $b_i < 0$, the question is how to choose an initial basic feasible solution? The *two-phase simplex method* proceeds as follows:

1. Introduce a new variable x_0

$$x_{n+i} = b_i - \sum_{j=1}^n a_{ij}x_j + x_0, \quad i = 1, \dots, m,$$

and an objective $-x_0$ to maximize

2. Suppose that $b_i < 0$ is the smallest value. Perform a pivot on x_0 and thus x_{n+i} will turn the dictionary into a feasible one.

3. This *non-cycling pivots* will lead to either (a) an optimal basis where x_0 is within the basis, and we conclude this problem is infeasible; (b) or we have x_0 out of the basis, and we just delete x_0 and plug back the original objective, and go from there.

Example 1.2.2. Given the dictionary

$$\begin{aligned}x_4 &= 4 - 2x_1 + x_2 - 2x_3 \\x_5 &= -5 - 2x_1 + 3x_2 - x_3 \\x_6 &= -1 + x_1 - x_2 + 2x_3\end{aligned}$$

We first add the new variable x_0 and an objective $-x_0$:

$$\begin{aligned}x_4 &= 4 - 2x_1 + x_2 - 2x_3 + x_0 \\x_5 &= -5 - 2x_1 + 3x_2 - x_3 + x_0 \\x_6 &= -1 + x_1 - x_2 + 2x_3 + x_0 \\z &= -x_0\end{aligned}$$

Choosing x_0 entering the basis and x_5 leaving the basis, we obtain:

$$\begin{aligned}x_4 &= 9 - x_2 + x_3 + x_5 \\x_0 &= 5 + 2x_1 - 3x_2 - x_3 + x_5 \\x_6 &= 4 + 3x_1 - 4x_2 + 3x_3 + x_5 \\w &= -5 - 2x_1 + 3x_2 + x_3 - x_5\end{aligned}$$

and our feasible solution is $(x_1, x_2, x_3, x_4, x_5, x_6) = (0, 0, 0, 9, 0, 4)$.

1.3 Duality Results

Theorem 1.3.1. A linear programming problem can only be (i) *feasible*; or (ii) *infeasible*. In case (i), then *there exists a basic feasible solution*, and further with two possibilities: (i.a) *an optimal solution exists, in that case a basic optimal solution exists* (i.b) *the problem is unbounded*.

Duality problem is the best possible upper bounding problem Consider the primal problem

$$(P) \quad \begin{aligned} \max \quad & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & \mathbf{x} \geq 0 \end{aligned}$$

Take any $\mathbf{y} \geq 0$ such that $\mathbf{y}^T \mathbf{A} \geq \mathbf{c}^T$, and thus $\mathbf{y}^T \mathbf{b}$ becomes an upper bound for the optimal value. Therefore the best possible upper bounding problem becomes:

$$(D) \quad \begin{array}{ll} \max & \mathbf{b}^T \mathbf{y} \\ \text{s.t.} & \mathbf{A}^T \mathbf{y} \geq \mathbf{c} \\ & \mathbf{y} \geq 0 \end{array}$$

which is known as the *dual* problem.

The proceed above can be summarized as the weak duality theorem:

Theorem 1.3.2 (Weak Duality). Let \mathbf{x}, \mathbf{y} be the primal feasible, and dual feasible solution to (P) and (D), respectively, then we always have $\mathbf{b}^T \mathbf{y} \geq \mathbf{c}^T \mathbf{x}$.

Theorem 1.3.3 (Strong Duality). If (P) has an optimal solution, then (D) has an optimal solution. Moreover, the optimal values coincide.

Proof. Let B be an optimal basis for (P), then we have

$$\mathbf{A}_B^{-1} \mathbf{b} \geq 0, \quad \begin{bmatrix} \mathbf{c}^T & \mathbf{0}_m^T \end{bmatrix} - \mathbf{c}_B^T \mathbf{A}_B^{-1} \begin{bmatrix} \mathbf{A} & \mathbf{I} \end{bmatrix} \leq 0$$

Therefore we construct the dual feasible solution $\mathbf{y} := \mathbf{A}_B^{-1} \mathbf{c}_B$, which implies $\mathbf{b}^T \mathbf{y} = \mathbf{c}_B^T \mathbf{A}_B^{-1} \mathbf{b}$. Therefore $\mathbf{b}^T \mathbf{y}$ should be the optimal solution for (D). \square

Complentarity Slackness

Theorem 1.3.4 (Complentarity Condition). Consider the primal and dual problem

$$(P) \quad \begin{array}{ll} \max & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} & \mathbf{A} \mathbf{x} + \mathbf{s} = \mathbf{b} \\ & \mathbf{x} \geq 0, \mathbf{s} \geq 0 \end{array}, \quad (D) \quad \begin{array}{ll} \min & \mathbf{b}^T \mathbf{y} \\ \text{s.t.} & \mathbf{A}^T \mathbf{y} - \mathbf{w} = \mathbf{c} \\ & \mathbf{y} \geq 0, \mathbf{w} \geq 0 \end{array}$$

If (P) has an optimal solution (\mathbf{x}, \mathbf{s}) and (D) has an optimal solution (\mathbf{y}, \mathbf{w}) , then

$$\begin{aligned} \mathbf{s} \circ \mathbf{y} &= \mathbf{0} \\ \mathbf{w} \circ \mathbf{x} &= \mathbf{0} \end{aligned}$$

Remark 1.3.1. 1. If (P) is feasible and unbounded, then (D) must be infeasible.

2. The dual of the dual problem is the primal problem

3. There is possibility that both (P) and (D) are infeasible. Consider the self-dual problem for example:

$$\begin{array}{ll} \max & x_1 - x_2 \\ \text{s.t.} & \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \leq \begin{pmatrix} -1 \\ 1 \end{pmatrix} \\ & x_1 \geq 0, x_2 \geq 0 \end{array}$$

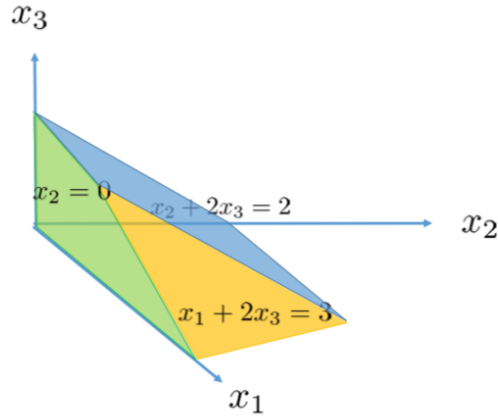


Figure 2.1: Illustration for polyhedral geometry

4. Therefore, the relationship for primal and dual problems can be summarized in the table below:

	Feasible	Unbounded	Infeasible
Feasible	Y	N	N
Unbounded	N	N	Y
Infeasible	N	Y	Y

2 Geometry and Duality for Linear Programming

2.1 The polyhedral geometry

The constraint of a LP forms a polyhedron. One example for a LP with ternary variables is shown in the Fig (2.1)

Let's introduce some terminologies formally:

Definition 2.1.1. • *Hyperplane* is the set $\{\mathbf{x} \mid \mathbf{a}^T \mathbf{x} = \mathbf{b}\}$

• *Half-space* is the set $\{\mathbf{x} \mid \mathbf{a}^T \mathbf{x} \leq \mathbf{b}\}$

• The *polyhedron* P is the intersection of *finite* number of half-spaces:

$$P = \left\{ \mathbf{x} \mid \mathbf{a}_i^T \mathbf{x} \leq \mathbf{b}_i, i = 1, \dots, m \right\}$$

• The *dimension* of a *polyhedron* is defined as the *lowest* dimension *affine space* containing P

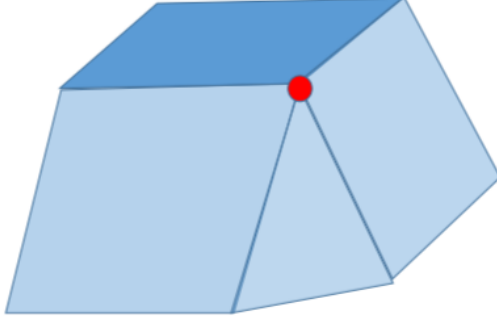


Figure 2.2: Over-determination results in Degeneracy

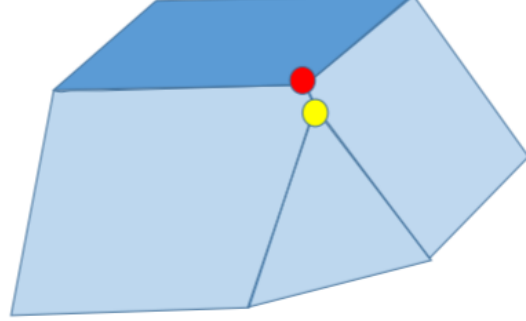


Figure 2.3: Perturbation diminishes over-determination

- The *face* of a polyhedron is defined as

$$\{\mathbf{x} \mid \mathbf{a}^T \mathbf{x} = \mathbf{b}\} \cap P,$$

where $P \subseteq \{\mathbf{x} \mid \mathbf{a}^T \mathbf{x} \leq \mathbf{b}\}$.

- Note that the face of a polyhedron is also a polyhedron. Therefore we define *facet* is the face of P that is one dimensional lower than that of P ; the *vertex* of P is the face of P that has dimension 0.

Remark 2.1.1. In space \mathbb{R}^n , normally n hyperplanes intersect at one point

If P is full dimensional (i.e., with dimension n), then a vertex of P is an intersection of n facets.

However, sometimes there is a case that more than n hyperplanes intersect at one point, say a vertex, which creates *degeneracy* (show in Fig (2.4)). In such case, adding regularization, i.e., perturbation diminishes over-determination.

Definition 2.1.2. Given a full dimensional P , we say two distinct vertices of P are *adjacent* if they are in the same $n - 1$ hyperplane.

Remark 2.1.2. Every update for simplex pivots move from a vertice to one of its adjacent position.

Definition 2.1.3. A *polyhedral cone* is defined as an intersection of a finite number of half-spaces, i.e., $K = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} \geq \mathbf{0}\}$, where $\mathbf{A} \in \mathbb{R}^{m \times n}$. In geometry shown in Fig (2.6), a polyhedral cone is a cone where its boundaries are polyhedral

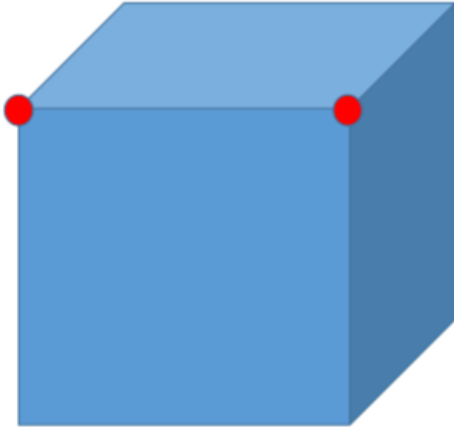


Figure 2.4: Illustration for *adjacent* vertices

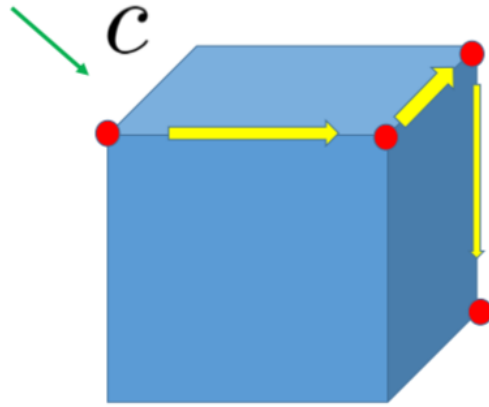


Figure 2.5: Path for simplex pivots

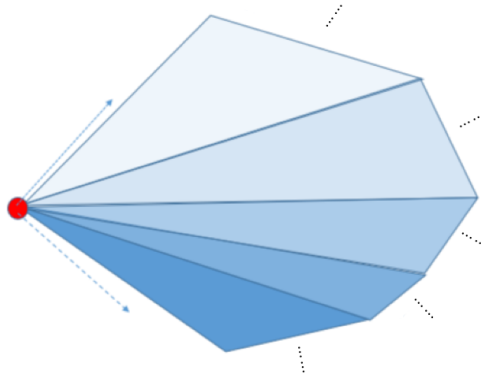


Figure 2.6: Illustration for a polyhedral cone

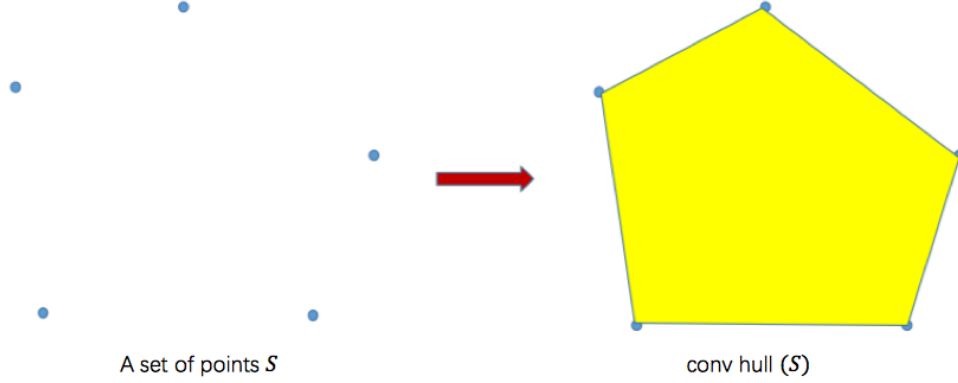


Figure 2.7: Illustration for a convex hull

Definition 2.1.4. Given a set of points S , we can define its convex hull as

$$\text{conv hull}(S) := \left\{ \sum_{i=1}^m \lambda_i \mathbf{y}_i \mid \forall \mathbf{y}_i \in S, \lambda_i \geq 0, \sum_{i=1}^m \lambda_i = 1 \right\}$$

In other words, a convex hull of S is a set, in which each point is a convex combination of points in S .

2.2 Fundamental Theorems for Linear Programming

Theorem 2.2.1 (Caratheodory's Theorem). Let $P \subseteq \mathbb{R}^n$ be a set of points. For each $\mathbf{x} \in \text{conv}(P)$, there exists a set $P' \subseteq P$ of cardinality at most $n + 1$, such that $\mathbf{x} \in \text{conv}(P')$.

Definition 2.2.1. The point \mathbf{x} in a convex set S is an *extremal point* if

1. $\mathbf{x} \in S$
2. \mathbf{x} is a convex combination of $\mathbf{y}, \mathbf{z} \in S$ implies $\mathbf{y} = \mathbf{z} = \mathbf{x}$.

Sometimes it is difficult to define the extremal point, if the convex set covers infinite area. Therefore, we define the points in the boundary of vertices instead:

Definition 2.2.2. A ray d in a convex cone \mathcal{K} is an *extremal ray* if

1. $d \in \mathcal{K}$

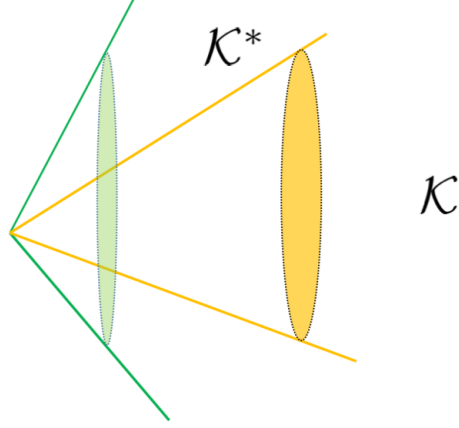


Figure 2.8: Illustration for a convex cone and its dual cone

2. If d can be written as non-negative summation of two rays $\xi_1, \xi_2 \in \mathcal{K}$, then $d = \xi_1 = \xi_2$

Definition 2.2.3. 1. A *polytope* is a convex hull of a *finite* number of points

2. A *polytope* is also a *polyhedron*

3. In general, a polyhedron can be written as

$$\left\{ \sum_{i=1}^m \lambda_i p_i + \sum_{j=1}^{\ell} \mu_j d_j \mid \lambda_i \geq 0, \sum_{i=1}^m \lambda_i = 1, \mu_j \geq 0 \right\}$$

In other words,

$$H = P + C,$$

where H, P, C denote collections of polyhedrons, polytopes, polyhedral cones, respectively.

Definition 2.2.4. If \mathcal{K} is a convex cone, then its *dual cone* is defined as

$$\mathcal{K}^* = \{ \mathbf{s} \mid \langle \mathbf{s}, \mathbf{x} \rangle \geq 0, \forall \mathbf{x} \in \mathcal{K} \}.$$

The Fig (2.8) represents one illustration for convex cone and its dual cone, but note that the dual cone may not necessarily larger than the convex cone itself. (Counter-example: Lorentz cone)

A polyhedral cone can be represented in two ways:

$$\{ \mathbf{A}\mathbf{x} \mid \mathbf{x} \geq 0 \} \quad \text{and} \quad \{ \mathbf{x} \mid \mathbf{B}^T \mathbf{x} \geq 0 \}$$

Here we consider two cones $\{ \mathbf{A}\mathbf{x} \mid \mathbf{x} \geq 0 \}$ and $\{ \mathbf{y} \mid \mathbf{A}^T \mathbf{y} \geq 0 \}$:

Theorem 2.2.2. Two cones $\mathcal{K}_1 = \{\mathbf{Ax} \mid \mathbf{x} \geq 0\}$ and $\mathcal{K}_2 = \{\mathbf{y} \mid \mathbf{A}^T \mathbf{y} \geq 0\}$ are actually *duality pairs*.

Proof. It's obvious that $\{\mathbf{Ax} \mid \mathbf{x} \geq 0\} \subseteq \{\mathbf{y} \mid \mathbf{A}^T \mathbf{y} \geq 0\}^*$.

Given that $\mathbf{b} \notin \{\mathbf{Ax} \mid \mathbf{x} \geq 0\}$, there exists \mathbf{y} such that $\mathbf{A}^T \mathbf{y} \geq 0$ and $\langle \mathbf{b}, \mathbf{y} \rangle < 0$ (by Farkas Lemma), i.e., \mathbf{b} cannot be in $\{\mathbf{y} \mid \mathbf{A}^T \mathbf{y} \geq 0\}^*$. \square

Theorem 2.2.3 (Farkas Lemma). If $\mathbf{b} \notin \{\mathbf{Ax} \mid \mathbf{x} \geq 0\}$, then there exists \mathbf{y} such that $\mathbf{A}^T \mathbf{y} \geq 0$ and $\langle \mathbf{b}, \mathbf{y} \rangle < 0$.

An equivalent form of Farkas Lemma is as follows (known as the *theorem of alternatives*):

Either $\mathbf{Ax} = \mathbf{b}, \mathbf{x} \geq 0$ has a solution, or $\mathbf{A}^T \mathbf{y} \geq 0, \langle \mathbf{b}, \mathbf{y} \rangle = -1$ has a solution, but not neither, nor both.

Proof. Consider the primal LP

$$(P) \quad \begin{array}{ll} \max & -s \\ \text{such that} & \mathbf{Ax} + s\mathbf{b} = \mathbf{b}, \\ & \mathbf{x} \geq 0, s \geq 0 \end{array}$$

The dual LP is

$$(D) \quad \begin{array}{ll} \min & \langle \mathbf{b}, \mathbf{y} \rangle \\ \text{such that} & \mathbf{A}^T \mathbf{y} \geq 0 \\ & \langle \mathbf{b}, \mathbf{y} \rangle \geq -1 \end{array}$$

Since the primal problem is feasible and bounded, and so an optimal solution exists. By strong duality theorem, $\{\mathbf{x} \mid \mathbf{Ax} = \mathbf{b}, \mathbf{x} \geq 0\} = \emptyset$ iff the (D) admits negative optimal value. \square

Remark 2.2.1. The Farkas Lemma essentially claims that given a polyhedron $\mathcal{K}_1 = \{\mathbf{x} \mid \mathbf{Ax} = \mathbf{b}, \mathbf{x} \geq 0\}$ and a point \mathbf{y} outside \mathcal{K}_1 , we can always find an affine that separates \mathcal{K}_1 and \mathbf{y} (see Fig (2.9))

Remark 2.2.2. Consider instead the case $\{\mathbf{x} \mid \mathbf{Ax} = \mathbf{b}\} = \emptyset$. We can imply $\mathbf{b} \notin \mathcal{R}(\mathbf{A})$, i.e., there exists \mathbf{y} such that $\mathbf{A}^T \mathbf{y} = \mathbf{0}$ and $\langle \mathbf{b}, \mathbf{y} \rangle \neq 0$.

Actually, we can generalize this separation into general polyhedron.

Theorem 2.2.4 (Separation Theorem). Let $H \subseteq \mathbb{R}^n$ be a general *polyhedron*. Suppose that $\mathbb{R}^n \ni \mathbf{p} \notin H$, then there exists an affine $f(\mathbf{x}) = \mathbf{c}^T \mathbf{x} + \mathbf{d}$ satisfying

$$f(\mathbf{p}) < 0, \quad f(\mathbf{x}) = \mathbf{c}^T \mathbf{x} + \mathbf{d} > 0, \forall \mathbf{x} \in H$$

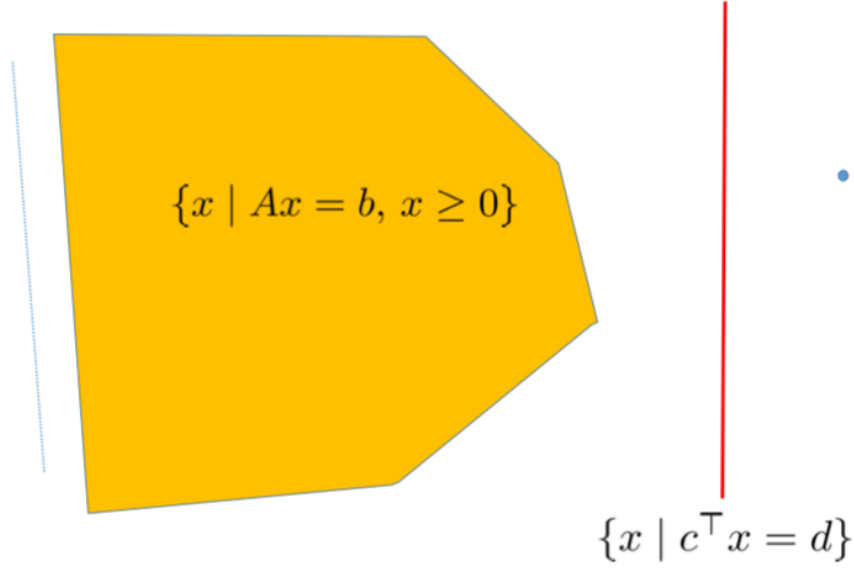


Figure 2.9: Illustration for the separation of \mathcal{K}_1 and y

Proof. Consider a general polyhedron H with vertices $\{p_i \mid i = 1, \dots, m\}$ and extreme rays $\{d_j \mid j = 1, \dots, \ell\}$. Since $p \notin H$, the system (2.1) does not have a solution:

$$\begin{aligned}
 p &= \sum_{i=1}^m \lambda_i p_i + \sum_{j=1}^{\ell} \mu_j d_j \\
 1 &= \sum_{i=1}^m \lambda_i \\
 0 &\leq \lambda_i, \quad i = 1, \dots, m \\
 0 &\leq \mu_j, \quad j = 1, \dots, \ell
 \end{aligned} \tag{2.1}$$

By Farkas Lemma, there exists $y \in \mathbb{R}^n$ and $s \in \mathbb{R}$ such that

$$\begin{aligned}
 s + p_i^T y &\geq 0, \quad i = 1, \dots, m \\
 d_j^T y &\geq 0, \quad j = 1, \dots, \ell \\
 s + p^T y &< 0
 \end{aligned}$$

□

Remark 2.2.3. It's also easy to derive the Farkas Lemma from the Separation Theorem:

Suppose $\mathbf{b} \notin \{\mathbf{Ax} \mid \mathbf{x} \geq 0\}$, then there is a separating hyperplane $\mathbf{c}^T \mathbf{z} = d$, such that

$$\mathbf{c}^T \mathbf{b} < d, \quad \mathbf{c}^T \mathbf{z} > d, \forall \mathbf{z} = \mathbf{Ax}, \text{ with } \mathbf{x} \geq 0$$

which implies that $\mathbf{A}^T \mathbf{c} \geq 0, d < 0, \mathbf{c}^T \mathbf{b} < d$.

There are different ways to show the same result. Our previous proof of Theorem (2.2.3) is based on the *finiteness of the simplex method*, i.e., the duality of LP. Terence Tao applied a different and more direct method to show the Farkas Lemma, let's study it in detail.

First we aim to show the Farkas Lemma in the dual way:

Theorem 2.2.5 (Rephrased Form of Farkas Lemma). Let $P_i(x) = \sum_{j=1}^n p_{ij}x_j - r_i, i = 1, \dots, m$. If the system $P_i(x) \geq 0$ for $i = 1, \dots, m$ does not have a solution, then there exists $y_i \geq 0$ such that $\sum_{i=1}^m y_i P_i(x) = -1$.

Proof. We can show this result by induction on n :

- When $n = 1$, we can scale $P_i(x) \geq 0$ possibly into three cases:

$$x_1 - a_i \geq 0, \quad i \in I_+; \quad -x_1 + b_j \geq 0, \quad j \in I_-; \quad c_k \geq 0, \quad k \in I_0$$

If there exists $k \in I_0$ with $c_k < 0$, we let $y_k = -1/c_k$ and $y_\ell = 0, \forall \ell \neq k$; otherwise there exists $i \in I_+$ and $j \in I_-$ with $a_i > b_j$, we let $y_i = y_j = 1/(a_i - b_j)$ and $y_\ell = 0, \forall \ell \neq i, j$. Therefore, we obtain $\sum_{i=1}^m y_i P_i(x) = -1$.

- Now suppose that this theorem is shown for dimension no more than n . Consider the case where we have $n + 1$ variables: $\bar{\mathbf{x}} = (\mathbf{x}, x_{n+1}), \mathbf{x} \in \mathbb{R}^n$. We can scale the original inequalities over $P_i(\bar{\mathbf{x}})$ according to the coefficients of x_{n+1} to obtain the following equivalent system of inequalities:

$$\begin{cases} P_i(\bar{\mathbf{x}}) := x_{n+1} - Q_i(\mathbf{x}) \geq 0, & i \in I_+ \\ P_j(\bar{\mathbf{x}}) := -x_{n+1} - Q_j(\mathbf{x}) \geq 0, & j \in I_- \\ P_k(\bar{\mathbf{x}}) := Q_k(\mathbf{x}) \geq 0, & k \in I_0 \end{cases} \quad (2.2)$$

Note that the system

$$\begin{cases} -Q_i(x) + Q_j(x) \geq 0, & i \in I_+, j \in I_- \\ Q_k(x) \geq 0, & k \in I_0 \end{cases} \quad (2.3)$$

has a solution implies that the original system would have a solution too. Due to the hypothesis of the lemma, (2.3) does not have a solution.

By the induction hypothesis, there exists $y_{ij}, y_k \geq 0$ such that

$$\sum_{i \in I_+, j \in I_-} y_{ij} (-Q_i(x) + Q_j(x)) + \sum_{k \in I_0} y_k Q_k(x) = -1$$

Therefore, we imply

$$\begin{aligned} \sum_{i \in I_+} \left(\sum_{j \in I_-} y_{ij} \right) (x_{n+1} - Q_i(x)) + \sum_{j \in I_-} \left(\sum_{i \in I_+} y_{ij} \right) (-x_{n+1} + Q_j(x)) \\ + \sum_{k \in I_0} y_k Q_k(x) = -1 \end{aligned}$$

The Farkas Lemma is shown by this induction argument. \square

Then we are ready to show the following form of separation theorem:

Theorem 2.2.6 (Repharsed Form of Separation Theorem). If polytopes P_1, P_2 do not intersect, then there is an affine $f(\mathbf{x}) = \langle \mathbf{y}, \mathbf{x} \rangle + y_0$ such that

$$f(\mathbf{x}) \geq 1, \forall \mathbf{x} \in P_1 \quad f(\mathbf{x}) \leq -1, \forall \mathbf{x} \in P_2$$

Proof. Suppose the extreme points of P_1, P_2 are $\{\mathbf{p}_1, \dots, \mathbf{p}_s\}$ and $\{\mathbf{q}_1, \dots, \mathbf{q}_t\}$, respectively. Therefore, it suffices to show that

$$\langle \mathbf{y}, \mathbf{p}_i \rangle + y_0 \geq 1, \quad i = 1, \dots, s; \quad \langle \mathbf{y}, \mathbf{q}_j \rangle + y_0 \leq -1, \quad j = 1, \dots, t \quad (2.4)$$

Suppose on the contrary that (2.4) does not have a solution (\mathbf{y}, y_0) . Then the Farkas Lemma asserts that there exists $u_i \geq 0, i = 1, \dots, s$ and $v_j \geq 0, j = 1, \dots, t$ such that

$$\sum_{i=1}^s u_i (\langle \mathbf{y}, \mathbf{p}_i \rangle + y_0 - 1) + \sum_{j=1}^t v_j (-\langle \mathbf{y}, \mathbf{q}_j \rangle - y_0 - 1) = -1$$

which implies

$$\begin{aligned} \sum_{i=1}^s u_i \mathbf{p}_i - \sum_{j=1}^t v_j \mathbf{q}_j &= 0 \\ \sum_{i=1}^s u_i - \sum_{j=1}^t v_j &= 0 \\ \sum_{i=1}^s u_i + \sum_{j=1}^t v_j &= 1 \end{aligned}$$

Therefore, $\sum_{i=1}^s u_i = \sum_{j=1}^t v_j = 1/2$, and thus

$$P_1 \ni 2 \sum_{i=1}^s u_i \mathbf{p}_i = 2 \sum_{j=1}^t v_j \mathbf{q}_j \in P_2$$

which contradicts to the assumption that $P_1 \cap P_2 = \emptyset$. \square

The same argument applies if *polytopes* is replaced by *polyhedron*.

2.3 More Theorems of Alternatives: the case for polyhedrons

Some more refined forms of the theorems of alternatives exist. Here we list some examples.

Notations Let \mathbf{x}, \mathbf{y} be two vectors, we denote $\mathbf{x} \gneq \mathbf{y}$ to be “ $\mathbf{x} \geq \mathbf{y}$ and $\mathbf{x} \neq \mathbf{y}$ ”, and similarly for $\mathbf{x} \lneq \mathbf{y}$.

Theorem 2.3.1 (Gordan). Either $\mathbf{Ax} > \mathbf{0}$ has a solution, or $\mathbf{A}^T \mathbf{y} = \mathbf{0}, \mathbf{y} \gneq \mathbf{0}$ has a solution.

Theorem 2.3.2 (Stiemke). Either $\mathbf{Ax} \gneq \mathbf{0}$ has a solution, or $\mathbf{A}^T \mathbf{y} = \mathbf{0}, \mathbf{y} > \mathbf{0}$ has a solution.

Theorem 2.3.3 (Gale). Assuming $\mathbf{Ax} \leq \mathbf{b}$ is feasible, then either $\mathbf{Ax} \lneq \mathbf{b}$ has a solution, or $\mathbf{A}^T \mathbf{y} = \mathbf{0}, \mathbf{b}^T \mathbf{y} = 0, \mathbf{y} > \mathbf{0}$ has a solution.

Theorem 2.3.4 (Tucker). Suppose that $\mathbf{A} \neq \mathbf{0}$. Either $\mathbf{Ax} \gneq \mathbf{0}, \mathbf{Bx} \geq \mathbf{0}, \mathbf{Cx} = \mathbf{0}$ has a solution, or $\mathbf{A}^T \mathbf{u} + \mathbf{B}^T \mathbf{v} + \mathbf{C}^T \mathbf{w} = \mathbf{0}, \mathbf{u} > \mathbf{0}, \mathbf{v} \geq \mathbf{0}$ has a solution.

Theorem 2.3.5 (Motzkin). Suppose that $\mathbf{A} \neq \mathbf{0}$. Either

$$\mathbf{Ax} > \mathbf{0}, \mathbf{Bx} \geq \mathbf{0}, \mathbf{Cx} = \mathbf{0}$$

has a solution, or

$$\mathbf{A}^T \mathbf{u} + \mathbf{B}^T \mathbf{v} + \mathbf{C}^T \mathbf{w} = \mathbf{0}, \mathbf{u} \gneq \mathbf{0}, \mathbf{v} \geq \mathbf{0}$$

has a solution.

3 Computational Complexity for Linear Programming

3.1 A case study

Suppose we aim to find a shortest path on a directed network, from the source node s to the sink node t . The network is supposed to be $\mathcal{N} = (V, E; w)$, where $V = \{v_1, \dots, v_n\}$ is the set of nodes, E is the set of edges, and w is the vectors of weights, i.e., w_{ij} denotes the weight on the edge (v_i, v_j) .

We need to input this instance into the computer. One way is to use the *node-arc* incidence matrix:

$$a_{ij} = \begin{cases} -1, & \text{if } v_i \text{ is the head of arc } e_j \\ 0, & \text{if } v_i \text{ is not related to arc } e_j \\ +1, & \text{if } v_i \text{ is the tail of arc } e_j \end{cases}$$

Input size of a problem

1. Suppose $|V| = n$ and $|E| = m$, then inputting the matrix A requires to touch the keyboard $mn + 2m$ times.
2. To tell the vector w into computer, we need w to be *integer-valued* (since otherwise the problem will be much more difficult. question: do w need to be positive?). Typing w_{ij} requires at most $\lfloor \log_2(|w_{ij}| + 1) \rfloor + 1$ bits. Therefore typing the vector w requires totally

$$\sum_{(v_i, v_j) \in E} \lfloor \log_2(|w_{ij}| + 1) \rfloor + |E| \text{ bits}$$

3. Therefore, the input-length (size) of this problem is

$$L = mn + 3m + \sum_{(v_i, v_j) \in E} \lfloor \log_2(|w_{ij}| + 1) \rfloor.$$

The Running Time Complexity It's logical that the total amount of perations required by an algorithm to solve this problem is *dependent* on L .

Consider applying the *Dijkstra* algorithm to find the shortest path from the source node s to the sink node t . Define a relevant working set S that contains all the nodes with *known shortest distance* to t . Initially $S = \{t\}$. At each iteration, by *dynamic programming principle*, one node is identified to be in S after at most $\mathcal{O}(|V|)$ comparisions. Therefore, the overall computational complexity is $\mathcal{O}(|V|^2)$.

3.2 Computational Complexity

Polynomial Time Algorithm If an algorithm solves a combinatorial problem with input size L with total number of operations no more than a polynomial of L , then it is called a *polynomial-time* algorithm.

The Dijkstra algorithm requires no more than $\mathcal{O}(L^2)$ (question: $\mathcal{O}(L)$ or $\mathcal{O}(L^2)$?) operations to terminate, and therefore it is a polynomial-time algorithm.

However, there are many combinatorial problems for which no polynomial time algorithms are known, e.g., the longest path problem, the Hamitonian circle problem, the maximum cut, and many more.

3.2.1 P & NP

Definition 3.2.1 (NP problem). If a combinatorial decision problem with input size L is such that:

if the answer to the problem is yes, then a yes-certificate exists and the length (size) of which is polynomial in L .

then we call the problem is NP .

Definition 3.2.2 (co-NP problem). If a combinatorial decision problem with input size L is such that:

if the answer to the problem is no, then a no-certificate exists and the length (size) of which is polynomial in L .

then we call the problem is $co-NP$.

Definition 3.2.3 (P problem). A *polynomially solvable* problem is called to be P .

It's clear that $P \subseteq NP$ and $P \subseteq co-NP$.

Definition 3.2.4 (NP-Complete problem). The problem in NP such that any other problem in NP can be reduced to it is called to be NP -Complete.

It's *strongly believed* (*haven't shown*) that NP-Complete problems are not in P.

3.2.2 Dimensions and Parameters

A decision problem typically involves *dimension* and the *parameters* of the problem, though they are mixed in the definition of input-length L .

- In the case study, $n = |V|$ and $m = |E|$ are known to be the problem dimensions, and the weight w is known as the problem parameter.
- For the linear programming with integer-valued parameters

$$\begin{array}{ll} \max & \sum_{j=1}^n c_j x_j \\ \text{such that} & \sum_{j=1}^n a_{ij} x_j \leq b_i, \quad i = 1, \dots, m \\ & x_j \geq 0, \quad j = 1, \dots, n \end{array}$$

the problem dimension would be m and n , and the input-length is

$$L = mn + m + n + \sum_{j=1}^n \lfloor \log_2(|c_j| + 1) \rfloor + \sum_{i=1}^m \lfloor \log_2(|b_i| + 1) \rfloor + \sum_{i=1}^m \sum_{j=1}^n \lfloor \log_2(|a_{ij}| + 1) \rfloor.$$

3.2.3 Other terminologies

Definition 3.2.5 (Weak & Strong Polynomial Time Algorithms). If a *polynomial-time* algorithm requires number of operations to be *polynomial* in dimensions, then the algorithm is called the *strongly polynomial*, otherwise it is only *weakly polynomial*.

Remark 3.2.1. Note that the input size of a problem defines the polynomial time, and the dimension specializes the strong & weak polynomial time. The number of operations is different from the number of *bit* operations. The Dijkstra algorithm is strongly polynomial.

Definition 3.2.6 (Weak & Strong NP-Completeness). If an *NP-complete* is such that *even when it is restricted to be the case where all the parameters are constants*, it still remains to be *NP-complete*, then it is called *strongly NP-Complete*, otherwise it is called *weakly NP-complete*.

Example: the 2-partition problem is weakly NP-complete; the Hamiltonian circle problem is strongly NP-complete.

Definition 3.2.7. If a decision version of the problem is *NP-complete* or *co-NP-complete*, then the problem is called *NP-Hard*.

For example, the travelling salesman (TRS) problem is NP-Hard.

3.3 Complexity of Linear Programming

3.3.1 LP is both NP and Co-NP

The decision version of LP is: does there exists \mathbf{x} satisfying

$$\mathbf{c}^T \mathbf{x} \geq \mathbf{v}, \quad \mathbf{A} \mathbf{x} \leq \mathbf{b}, \quad \mathbf{x} \geq 0?$$

1. If the answer is yes, then a certificate is such \mathbf{x} , whose size can be bounded by a polynomial of the input-length.
2. If the answer is no, then the above system is infeasible. By Farkas Lemma, there exists $y_0 \geq 0$ and $\mathbf{y} \geq 0$ such that

$$-y_0 \mathbf{c}^T + \mathbf{y}^T \mathbf{A} \geq 0, \quad -y_0 \mathbf{v} + \mathbf{y}^T \mathbf{b} < 0.$$

In this case, the size of (y_0, \mathbf{y}) can be bounded by a polynomial of the input-length.

3.3.2 Is LP in P?

In the complexity theory, we believe that $P \subsetneq NP$, i.e., $P \neq NP$ -Complete and $NP \cap Co-NP = P$. Therefore, a natural question arises: Is linear programming in P? The answer is no, but the simplex method is *not* a polynomial time algorithm. Let's study an example first.

The Klee-Minty Example (n=3) Consider solving a LP using the Largest Coefficient Rule:

$$\begin{aligned}
& \max && 100x_1 + 10x_2 + x_3 \\
& \text{such that} && x_1 \leq 1 \\
& && 20x_1 + x_2 \leq 100 \\
& && 200x_1 + 20x_2 + x_3 \leq 10000 \\
& && x_1 \geq 0, x_2 \geq 0, x_3 \geq 0
\end{aligned} \tag{3.1a}$$

We can find its dictionary:

$$\begin{aligned}
\mathbf{x}_4 &= 1 - x_1 \\
x_5 &= 100 - 20x_1 - x_2 \\
x_6 &= 10000 - 200x_1 - 20x_2 - x_3 \\
z &= \mathbf{100x}_1 + 10x_2 + x_3
\end{aligned} \tag{3.1b}$$

If we choose the variable with largest coefficient to enter the basis, i.e., x_1 , we obtain:

$$\begin{aligned}
x_1 &= 1 - x_4 \\
\mathbf{x}_5 &= 80 + 2 - x_4 - x_2 \\
x_6 &= 9800 + 200x_4 - 20x_2 - x_3 \\
z &= 100 - 100x_4 + \mathbf{10x}_2 + x_3
\end{aligned} \tag{3.1c}$$

If we choose the variable with largest coefficient to enter the basis, i.e., x_2 , we obtain:

$$\begin{aligned}
\mathbf{x}_1 &= 1 - x_4 \\
x_2 &= 80 + 20x_4 - x_5 \\
x_6 &= 8200 - 200x_4 + 20x_5 - x_3 \\
z &= 900 + \mathbf{100x}_4 - 10x_5 + x_3
\end{aligned} \tag{3.1d}$$

If we choose the variable with largest coefficient to enter the basis, i.e., x_4 , we obtain:

$$\begin{aligned}
x_4 &= 1 - x_1 \\
x_2 &= 100 - 20x_1 - x_5 \\
\mathbf{x}_6 &= 8000 + 200x_1 + 20x_5 - x_3 \\
z &= 1000 - 100x_1 - 10x_5 + \mathbf{x}_3
\end{aligned} \tag{3.1e}$$

If we choose the variable with largest coefficient to enter the basis, i.e., x_3 , we obtain:

$$\begin{aligned}
\mathbf{x}_4 &= 1 - x_1 \\
x_2 &= 100 - 20x_1 - x_5 \\
x_3 &= 8000 + 200x_1 + 20x_5 - x_6 \\
z &= 9000 + \mathbf{100x}_1 + 10x_5 - x_6
\end{aligned} \tag{3.1f}$$

If we choose the variable with largest coefficient to enter the basis, i.e., x_1 , we obtain:

$$\begin{aligned} x_1 &= 1 - x_4 \\ \mathbf{x}_2 &= 80 + 20x_4 - x_5 \\ x_3 &= 8200 - 200x_4 + 20x_5 - x_6 \\ z &= 9100 - 100x_4 + \mathbf{10x}_5 - x_6 \end{aligned} \tag{3.1g}$$

If we choose the variable with largest coefficient to enter the basis, i.e., x_5 , we obtain:

$$\begin{aligned} \mathbf{x}_1 &= 1 - x_4 \\ x_5 &= 80 + 20x_4 - x_2 \\ x_3 &= 9800 + 200x_4 - 20x_2 - x_6 \\ z &= 9900 + \mathbf{100x}_4 - 10x_2 - x_6 \end{aligned} \tag{3.1h}$$

If we choose the variable with largest coefficient to enter the basis, i.e., x_4 , we obtain:

$$\begin{aligned} x_4 &= 1 - x_1 \\ x_5 &= 100 - 20x_1 - x_2 \\ x_3 &= 10000 - 200x_1 - 20x_2 - x_6 \\ z &= 10000 - 100x_1 - 10x_2 - x_6 \end{aligned} \tag{3.1i}$$

After total 7 simplex pivot steps, we get the optimal solution. This process is mazy, while it only contains 3 variables.

General Case Consider the general Klee-Minty example:

$$\begin{aligned} \max \quad & \sum_{j=1}^n 10^{n-j} x_j \\ \text{such that} \quad & 2 \sum_{j=1}^{i-1} 10^{i-j} x_j + x_i \leq 100^{i-1}, \quad i = 1, \dots, n \\ & x_j \geq 0, \quad j = 1, \dots, n \end{aligned}$$

Its input-length (in digits) should be:

$$\begin{aligned} & \sum_{j=1}^n (n - j - 1) + \sum_{i=1}^n (2(i - 1) + 1) + \sum_{i=1}^n \left(\sum_{j=1}^{i-1} (i - j - 1) + 1 \right) \\ &= \frac{n^3 + 12n^2 + 5n}{6} \\ &= \mathcal{O}(n^3) \end{aligned}$$

Let the slack variables be:

$$s_i := 100^{i-1} - 2 \sum_{j=1}^{i-1} 10^{i-j} x_j - x_i, \quad i = 1, \dots, n$$

One can show that in every feasible basis, either x_i or s_i must be in the basis, $i = 1, \dots, n$, which implies that there are 2^n feasible basis in total.

If applying the largest coefficient pivot rule, then after $2^{n-1} - 1$ iterations, the last row reads

$$z = 10 \left(100^{n-2} - \sum_{j=1}^{n-2} 10^{n-1-j} x_j - s_{n-1} \right) + x_n$$

After further 2^{n-1} iterations, the last row reads

$$z = 90 \cdot 100^{n-2} + 10 \left(\sum_{j=1}^{n-2} 10^{n-1-j} x_j + s_{n-1} \right) - s_n$$

After further 2^{n-1} iterations, the last row reads

$$z = 100^{n-1} - \sum_{j=1}^{n-1} 10^{n-j} x_j - s_n,$$

which corresponds to an optimal basic solution.

Remark 3.3.1. Since the total number of pivot steps for the Klee-Minty example is $\mathcal{O}(2^n - 1)$ for a problem whose input-length is $\mathcal{O}(n^3)$, it shows that in the worst case the simplex method with the largest coefficient pivot rule is exponential.

Most other conceivable simplex pivot rules all admit similar exponential examples. However, it remains open that whether we can find a special simplex pivot rule that is polynomial even in the worst case.

The following figure shows the Geometric process for the Klee-Minty example with $n = 3$, and here $(s_1, s_2, s_3) = (x_1, 100x_2, 10000x_3)$:

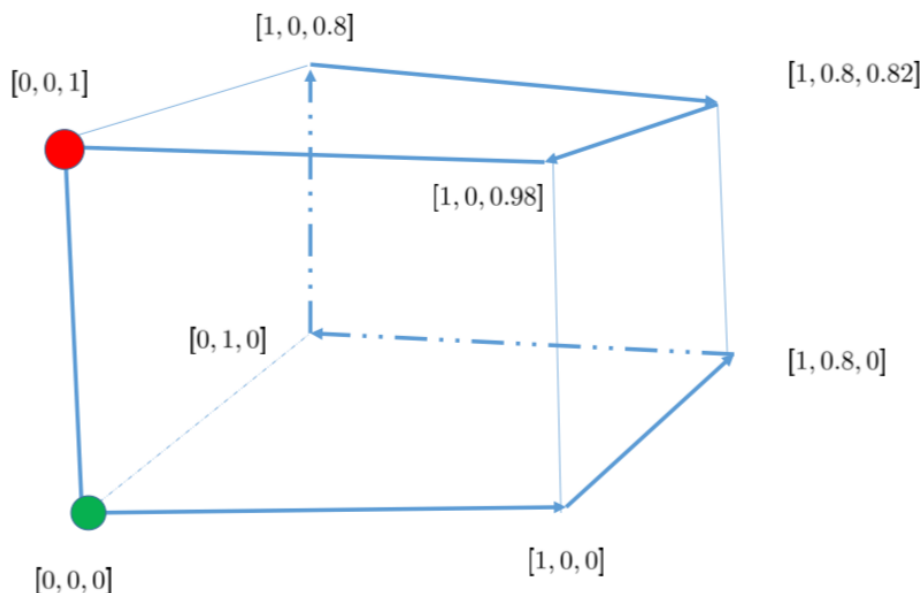


Figure 3.1: A Geometric Picture for the Klee-Minty example with $n = 3$

4 The KKT Condition for Nonlinear Programming

4.1 unconstrained Optimality Condition

Consider the unconstrained optimization

$$\min f(x) \tag{4.1}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is an m -th order *continuously differentiable* function. We aim to find the optimality condition for this problem.

Theorem 4.1.1 (First Order Necessary Condition). Given the condition that $f \in \mathcal{C}^1$, if x^* is a local minimum point, then $\nabla f(x^*) = 0$.

Theorem 4.1.2 (Second Order Necessary Condition). Given the condition that $f \in \mathcal{C}^2$, if x^* is a local minimum point, then $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*) \succeq 0$.

Theorem 4.1.3 (Second Order Sufficient Condition). Given the condition that $f \in \mathcal{C}^2$, if $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*) \succ 0$, then x^* is a local minimum point.

Theorem 4.1.4 (First Order Necessary and Sufficient Condition). If f is convex, then x^* is a local minimum point if and only if $0 \in \partial f(x^*)$.

Proof. All theorems above relies on the Taylor expansion and

$$f(x^* + \Delta x) \geq f(x^*), \forall \Delta x \iff \text{the point } x^* \text{ is minimum}$$

□

All the conditions above only involve the explicit quantities related to x . Now we turn into the constrained optimization case.

4.2 Constrained Optimality Condition

Consider the special constrained optimization

$$\begin{aligned} \min \quad & f(x) \\ \text{such that} \quad & x \in \mathcal{X} \end{aligned} \tag{4.2}$$

where $\mathcal{X} \subseteq \mathbb{R}^n$ is usually pre-assumed to be a closed convex set.

Definition 4.2.1 (Tangent Cone). Let $\hat{x} \in \mathcal{X}$. The *tangent cone* of \mathcal{X} at point \hat{x} is defined as

$$\mathcal{T}(\hat{x}) := \left\{ y \neq 0 \mid \exists x^k \in \mathcal{X} : x^k \rightarrow \hat{x} \ \& \ \frac{x^k - \hat{x}}{\|x^k - \hat{x}\|} \rightarrow \frac{y}{\|y\|} \right\} \cup \{0\}$$

It's clear that $\mathcal{T}(\hat{x})$ denotes all feasible directions from \hat{x} within \mathcal{X}

Therefore we obtain a necessary optimality condition for constrained optimization, which is beautiful in theory but not that useful in practice:

Theorem 4.2.1 (First Order Necessary Condition). If x^* is a local minimum point for the constraint (neither necessarily convex nor closed) set \mathcal{X} , then

$$\langle \nabla f(x^*), y \rangle \geq 0, \ \forall y \in \mathcal{T}(x^*),$$

i.e., $\nabla f(x^*) \in (\mathcal{T}(x^*))^*$.

Theorem 4.2.2. The above condition becomes necessary and sufficient if f is convex. In that case, we may weaken the condition into

$$\partial f(x^*) \cap (\mathcal{T}(x^*))^* \neq \emptyset.$$

Proof. The proof relies on the fact that f is convex iff

$$f(y) \geq f(x) + \langle d, (y - x) \rangle, \forall x, y \in \mathcal{X},$$

where $d \in \partial f(x^*)$. □

Theorem 4.2.3. Given further condition that \mathcal{X} is a convex set, if x^* is a local minimum point for the constraint set \mathcal{X} , then

$$\langle \nabla f(x^*), y - x^* \rangle \geq 0, \ \forall y \in \mathcal{X}.$$

This condition becomes sufficient if f is convex.

4.2.1 Characteration of Tangent Cone

1. Consider the polyhedral constraint set

$$\mathcal{X} = \{x \mid a_i^T x \leq b_i, \ i = 1, \dots, m\}$$

Let $\hat{x} \in \mathcal{X}$, and $I(\hat{x}) = \{i \mid a_i^T \hat{x} = b_i\}$. Then we have

$$\mathcal{T}(\hat{x}) = \{d \mid a_i^T d \leq 0, \ \forall i \in I(\hat{x})\}.$$

2. Consider another more general constraint set

$$\mathcal{X} = \left\{ x \mid \begin{array}{l} h_i(x) = 0, \ i = 1, \dots, m; \\ g_j(x) \leq 0, \ j = 1, \dots, r \end{array} \right\}.$$

For $\hat{x} \in \mathcal{X}$, define $I(\hat{x}) = \{j \mid g_j(\hat{x}) = 0\}$ likewise. Let's introduce a easily computable cone

$$\mathcal{C}(\hat{x}) = \left\{ d \mid \begin{array}{l} \langle \nabla h_i(\hat{x}), d \rangle = 0, \ i = 1, \dots, m; \\ \langle \nabla g_j(\hat{x}), d \rangle \leq 0, \ \forall j \in I(\hat{x}) \end{array} \right\}.$$

It's clear that $\mathcal{T}(\hat{x}) \subseteq \mathcal{C}(\hat{x})$. However, in general $\mathcal{T}(\hat{x}) \neq \mathcal{C}(\hat{x})$. (Consider $\mathcal{X} = \{(x_1, x_2) \mid x_1 \geq 0, x_2^2 \leq 0\}$).

A popular condition to ensure the equality is the *Linear Independence Constraint Qualification* (LICQ):

The vectors $\nabla h_i(x)$, $i = 1, \dots, m$; $\nabla g_j(x)$, $j \in I(x)$ are always linearly independent for $\forall x \in \mathcal{X}$.

Theorem 4.2.4. Under the LICQ, if the functions h_i, g_j are differentiable with Lipschitz continuous gradient, then $\mathcal{T}(\hat{x}) = \mathcal{C}(\hat{x})$.

Proof. Consider the case for which \mathcal{X} does not have equality constraints. Take $d \in \mathcal{C}(\hat{x})$. Applying LICQ and Theorem (2.3.1), there exists \hat{d} such that

$$\langle \nabla g_i(\hat{x}), \hat{d} \rangle < 0, \ i \in I(\hat{x}).$$

Consider $x(t) = \hat{x} + td + t^{1.5}\hat{d}$, we imply

$$g_i(x(t)) < 0, \ i \in I(\hat{x}), \ \text{for } 0 < t < \varepsilon,$$

where $\varepsilon > 0$ is sufficiently small. Therefore, $d \in \mathcal{T}(\hat{x})$. □

4.2.2 The KKT Condition

The optimization problem below admits another type of optimality condition:

$$\begin{aligned} \min \quad & f(x) \\ \text{such that} \quad & h_i(x) = 0, \quad i = 1, \dots, m; \\ & g_j(x) \leq 0, \quad j = 1, \dots, r \end{aligned} \tag{4.3}$$

Theorem 4.2.5 (Karush, Kuhn, and Tucker). Suppose that x^* is an optimal solution to the problem (4.3), then under some *regularity* condition (e.g., LICQ), there exists $\lambda \in \mathbb{R}^m$ and $\mu \in \mathbb{R}_+^r$ such that

$$\begin{cases} \nabla f(x^*) + \sum_{i=1}^m \lambda_i \nabla h_i(x^*) + \sum_{j=1}^r \mu_j \nabla g_j(x^*) = 0 \\ \mu_j g_j(x^*) = 0, \quad j = 1, \dots, r. \end{cases}$$

4.2.3 The Lagrangian Dual Problem

The KKT condition relates the explicit variable x with the implicit variable (λ, μ) . Now we take a close look at implicit variables, and they are called the *Lagrangian multipliers*.

Define the Lagrangian function

$$\mathcal{L}(x; \lambda, \mu) := f(x) + \sum_{i=1}^m \lambda_i h_i(x) + \sum_{j=1}^r \mu_j g_j(x), \quad \mu_j \geq 0.$$

Let $d(\lambda, \mu) := \min_x \mathcal{L}(x; \lambda, \mu)$, which implies the Lagrangian dual problem:

$$\begin{aligned} \max \quad & d(\lambda, \mu) \\ \text{such that} \quad & \lambda \in \mathbb{R}^m, \mu \in \mathbb{R}_+^r \end{aligned} \tag{4.4}$$

4.2.4 Duality for conic optimization

Consider the conic optimization problem

$$\begin{aligned} \min \quad & \langle \mathbf{c}, \mathbf{x} \rangle \\ \text{such that} \quad & \mathbf{A}\mathbf{x} = \mathbf{b} \\ & \mathbf{x} \in \mathcal{K} \end{aligned}$$

its Lagrangian dual problem is

$$\begin{aligned} \max \quad & \langle \mathbf{b}, \mathbf{y} \rangle \\ \text{such that} \quad & \mathbf{A}^T \mathbf{y} + \mathbf{s} = \mathbf{c} \\ & \mathbf{s} \in \mathcal{K}^* \end{aligned}$$

Note that linear programming is a special conic optimization.

The conic optimization problem is commonly solved by simultaneously optimizing the primal and the dual problem, due to the duality theorem:

1. Weak duality: any feasible primal solution value is an upper bound for any feasible dual solution value.
2. Strong duality: under some conditions, i.e., for some special conic optimization problems such as LP, the primal optimal value coincides with the dual optimal value.

The dual variables are some hidden, implicit, potential quantities that are implied by the primal variables at the optimality. The *optimality balance* is

$$\begin{cases} \mathbf{Ax} = \mathbf{b} \\ \mathbf{x} \in \mathcal{K} \\ \mathbf{A}^T \mathbf{y} + \mathbf{s} = \mathbf{c} \\ \mathbf{s} \in \mathcal{K}^* \\ \langle \mathbf{x}, \mathbf{s} \rangle = 0 \end{cases} \quad (4.5)$$

4.3 Projection Problem

Now we introduce a useful optimization problem, which can be applied to solve the system (4.5):

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{x} - \mathbf{v}\|^2 \\ \text{such that} \quad & \mathbf{x} \in \mathcal{X} \end{aligned} \quad (4.6)$$

This problem is called the *projection of \mathbf{v} onto the closed convex set \mathcal{X}*

4.3.1 Projection onto the nonnegative orthant

In particular, the problem is formulated as:

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{x} - \mathbf{v}\|^2 \\ \text{such that} \quad & \mathbf{x} \geq 0. \end{aligned} \quad (4.7)$$

Therefore, the KKT condition is

$$\begin{cases} \mathbf{x} - \mathbf{v} - \boldsymbol{\mu} = 0 \\ \mathbf{x} \geq 0, \boldsymbol{\mu} \geq 0, \langle \mathbf{x}, \boldsymbol{\mu} \rangle = 0 \end{cases}$$

In this case, the solution to the system above is

$$\mathbf{x} = [\mathbf{v}]_+, \quad \boldsymbol{\mu} = -[\mathbf{v}]_-$$

The corresponding dual problem is

$$\begin{aligned} \max \quad & -\frac{1}{2} \|\boldsymbol{\mu} + \mathbf{v}\|^2 + \frac{1}{2} \|\mathbf{v}\|^2 \\ \text{such that} \quad & \boldsymbol{\mu} \geq 0 \end{aligned}$$

Check First Order Optimality Condition Directly Another way is to apply Theorem (4.2.1) directly:

$$\mathcal{T}([v]_+) = \{d \mid d_i \geq 0, i \in J\}, \text{ where } J = \{i \mid v_i \leq 0\}.$$

Since $\nabla \big|_{x=[v]_+} (\frac{1}{2}\|x - v\|^2) = [v]_+ - v$, we have

$$\langle [v]_+ - v, d \rangle \geq 0, \forall d \in \mathcal{T}([v]_+).$$

Therefore, by Theorem (4.2.1), $[v]_+$ is optimal to the nonnegative cone projection problem.

4.3.2 Projection onto the affine linear space $Ax = b$

In particular, the problem is formulated as:

$$\begin{aligned} \min \quad & \frac{1}{2}\|x - v\|^2 \\ \text{such that} \quad & Ax = b. \end{aligned} \tag{4.8}$$

The solution for this problem is

$$\begin{aligned} x^* &= v + A^T(AA^T)^{-1}(b - Av) \\ &= (I - A^T(AA^T)^{-1}A)v + A^T(AA^T)^{-1}b \\ \lambda^* &= (AA^T)^{-1}(b - Av) \end{aligned}$$

The solution is implied from the first order optimality condition:

$$\begin{aligned} \langle \nabla f(x^*), d \rangle &= \langle x^* - v, d \rangle \\ &= (b - Av)^T(AA^T)^{-1}Ad \\ &= 0, \end{aligned}$$

for any feasible direction $d \in \mathcal{T}(x^*) = \{d \mid Ad = 0\}$.

4.3.3 General Projection Problem

Let \mathcal{X} be a closed convex set and suppose $v \notin \mathcal{X}$. Consider

$$\begin{aligned} \min \quad & \frac{1}{2}\|x - v\|^2 \\ \text{such that} \quad & x \in \mathcal{X} \end{aligned}$$

Suppose that x^* is the projection, i.e., optimal solution, then

$$\langle x^* - v, x - x^* \rangle \geq 0, \forall x \in \mathcal{X}.$$

Often, we denote the projection as $x^* = [v]_{\mathcal{X}}$.

4.4 The Augmented Lagrangian Dual

Now consider a more general optimization problem instead of projection:

$$\begin{aligned} \min \quad & f(x) \\ \text{such that} \quad & Ax = b \\ & x \in \mathcal{X} \end{aligned}$$

In projection problem we see that when the objective function involves the squared term, the optimization becomes easy because of the strong convexity of squared term. Therefore, we introduce an *augmented Lagrangian function*

$$\mathcal{L}_\gamma(x; \lambda) = f(x) + \lambda^\top (Ax - b) + \frac{\gamma}{2} \|Ax - b\|^2,$$

where $\gamma > 0$ denotes the penalty parameter. Like (4.4), define $d_\gamma(\lambda) = \min_{x \in \mathcal{X}} \mathcal{L}_\gamma(x; \lambda)$, and consider the dual problem

$$\begin{aligned} \max \quad & d_\gamma(\lambda) \\ \text{such that} \quad & \lambda \in \mathbb{R}^m \end{aligned} \tag{4.9}$$

question: what's this dual problem relates to the primal problem?

Proposition 4.4.1 (Convexity). $d_\gamma(\lambda)$ is concave in λ . Therefore the dual problem (4.9) is always a convex optimization problem.

Proof. The dual function $h(y) = \min_x a(x) + y^\top b(x)$ is always concave in y . (question: do we need to specify that \mathcal{X} is a convex set?) \square

Proposition 4.4.2 (Smooth). $d_\gamma(\lambda)$ is differentiable in λ .

Proof. Note that $x_y \in \arg \min a(x) + y^\top b(x)$ implies $b(x_y) \in \partial h(y)$, since

$$h(y') \leq h(y) + (y' - y)^\top b(x_y), \quad \forall y'.$$

Therefore, for $x_\lambda \in \arg \min f(x) + \lambda^\top (Ax - b) + \frac{\gamma}{2} \|Ax - b\|^2$, we have

$$Ax_\lambda - b \in \partial d_\gamma(\lambda).$$

To show the differentiability suffices to show the uniqueness of the sub-gradient. Suppose there exists two solutions x'_λ and x''_λ . By optimality condition,

$$\begin{cases} \langle \nabla f(x'_\lambda) + A^\top \lambda + \gamma A^\top (Ax'_\lambda - b), x''_\lambda - x'_\lambda \rangle \geq 0, \\ \langle \nabla f(x''_\lambda) + A^\top \lambda + \gamma A^\top (Ax''_\lambda - b), x'_\lambda - x''_\lambda \rangle \geq 0 \end{cases}$$

Adding them up and noting that $\langle \nabla f(x') - \nabla f(x''), x' - x'' \rangle \geq 0$, we have

$$Ax'_\lambda - b = Ax''_\lambda - b \implies \nabla d_\gamma(\lambda) = Ax_\lambda - b.$$

\square

Proposition 4.4.3 (Continuous differentiable). $\nabla d_\gamma(\lambda)$ is Lipschitz continuous with a Lipschitz constant $1/\gamma$.

Proof. Let λ' and λ'' be two vectors, similar as previous proof, we have:

$$\begin{cases} \langle \nabla f(x_{\lambda'}) + A^T \lambda' + \gamma A^T (Ax_{\lambda'} - b), x_{\lambda''} - x_{\lambda'} \rangle \geq 0, \\ \langle \nabla f(x_{\lambda''}) + A^T \lambda'' + \gamma A^T (Ax_{\lambda''} - b), x_{\lambda'} - x_{\lambda''} \rangle \geq 0 \end{cases}$$

Adding them up, we have

$$\gamma \|A(x_{\lambda''} - x_{\lambda'})\|^2 \leq (\lambda' - \lambda'')^T A(x_{\lambda''} - x_{\lambda'})$$

By Cauchy-Schwarz inequality,

$$\|\nabla d_\gamma(\lambda') - \nabla d_\gamma(\lambda'')\| \leq \frac{1}{\gamma} \|\lambda' - \lambda''\|.$$

□

5 Basic Algorithms for Nonlinear Programming

5.1 Gradient Algorithms

5.1.1 Preliminaries: convergence analysis

Consider an iterative algorithm for solving the optimization problem $\min f(x)$, producing iterates $\{x^0, x^1, \dots\}$.

1. The possible error measurements are as follows. The stopping criteria depends on these error measurements.

- $e(x^k) := \|x^k - x^*\|$;
- $e(x^k) = f(x^k) - f(x^*)$;

where x^* denotes the underlying optimal solution.

2. We say the algorithm converges if $\lim_{k \rightarrow \infty} e(x^k) = 0$
3. There are different types of convergence rate:

- (a) R-linear convergence: there exists $a \in (0, 1)$ such that $e(x^k) \leq Ca^k$;
- (b) Q-linear convergence: there exists $a \in (0, 1)$ such that $\frac{e(x^{k+1})}{e(x^k)} \leq a$;
- (c) Sub-linear convergence: $e(x^k) \leq C/k^p$ for some $p > 0$.

question: when say about convergence rate, do we need to specify which error measurements we use?

5.1.2 The (Sub)gradient algorithm for Unconstrained Optimization

Consider an unconstrained optimization problem $\min f(x)$, where f may not necessarily be smooth. Let $\{t_k > 0 \mid k = 0, 1, \dots\}$ be a sequence of step-sizes. Let's study the simplest first order optimization algorithm.

Algorithm 2 The (Sub)gradient Algorithm

Input: Initial guess $x^0 \in \mathcal{X}$

Output: Optimal solution \hat{x}

For $k = 0, 1, \dots$, **do**

- Take $d^k \in \partial f(x^k)$;
- $x^{k+1} \leftarrow x^k - t_k d^k$

end for.

Worst Case Bounds Consider a *convex optimization model* where f is a completely unknown function. The first order type algorithm essentially produces a sequence of iterates $\{x^k \mid k = 0, 1, 2, \dots\}$ in such a way that x^k is in the *affine space* spanned by

$$x^0, g(x^0), \dots, g(x^{k-1}), \text{ where } g(\cdot) = \partial f(\cdot).$$

- Suppose f is *Lipschitz continuous* and no other information is known, we can construct an example such that

$$\min_{x \in \text{Span}\{x^0, g(x^0), \dots, g(x^{k-1})\}} f(x) - f(x^*) \geq \mathcal{O}\left(\frac{1}{\sqrt{k}}\right), \forall k = 1, 2, \dots, \lfloor \frac{n}{2} \rfloor$$

Therefore, the first order type algorithm can never reach the convergence rate faster than $\mathcal{O}(\frac{1}{\sqrt{k}})$.

- Additionally, if we know f is *differentiable* and ∇f is *Lipschitz continuous*, then we can construct an example such that

$$\min_{x \in \text{Span}\{x^0, g(x^0), \dots, g(x^{k-1})\}} f(x) - f(x^*) \geq \mathcal{O}\left(\frac{1}{k^2}\right), \forall k = 1, 2, \dots, \lfloor \frac{n}{2} \rfloor$$

Therefore, the first order type algorithm can never reach the convergence rate faster than $\mathcal{O}(\frac{1}{k^2})$ for optimizing this class of function.

5.1.3 Gradient Algorithm with Exact Line-Search

First we discuss the optimization with a uniform convex function. This assumption is by default unless specifically mentioned. A nice Q-linear convergence result is obtained:

Theorem 5.1.1. Suppose there exists $0 < m \leq M$ such that $0 \succ mI \succeq \nabla^2 f(x) \succeq MI$ (i.e., f is *uniformly convex*), and an exact line search is performed per iteration:

$$t_k := \arg \min_t f(x^k - t \nabla f(x^k)),$$

then

$$f(x^{k+1}) - f(x^*) \leq \left(1 - \frac{m}{M}\right) [f(x^k) - f(x^*)] \quad (5.1)$$

Proof. • **(Uniform Convexity implies Strongly Convexity)** For $\forall \mathbf{x}_1, \mathbf{x}_2 \in \text{dom}(f)$, by mean-value theorem,

$$f(\mathbf{x}_2) = f(\mathbf{x}_1) + \langle \nabla f(\mathbf{x}_1), \mathbf{x}_2 - \mathbf{x}_1 \rangle + \frac{1}{2}(\mathbf{x}_2 - \mathbf{x}_1)^T \nabla^2 f(\boldsymbol{\xi})(\mathbf{x}_2 - \mathbf{x}_1),$$

where $\boldsymbol{\xi}$ is some number between \mathbf{x}_2 and \mathbf{x}_1 . Applying the uniform convexity of f , we derive the strongly convexity property:

$$\frac{m}{2} \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \leq f(\mathbf{x}_2) - f(\mathbf{x}_1) - \langle \nabla f(\mathbf{x}_1), \mathbf{x}_2 - \mathbf{x}_1 \rangle \leq \frac{M}{2} \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \quad (5.2)$$

- **(Applying Strongly Convexity Property)** On the one hand, by setting $\mathbf{x}_1 = \mathbf{x}^*$ and $\mathbf{x}_2 = \mathbf{x}$ in (5.2), we obtain:

$$\frac{m}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2 \leq f(\mathbf{x}) - f(\mathbf{x}^*) \leq \frac{M}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2 \quad (5.3)$$

On the other hand, by setting $\mathbf{x}_1 = \mathbf{x}$ and $\mathbf{x}_2 = \mathbf{x}^*$ in (5.2), we obtain

$$\begin{aligned} \frac{m}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2 &\leq f(\mathbf{x}^*) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle \\ &\leq f(\mathbf{x}^*) - f(\mathbf{x}) + \|\nabla f(\mathbf{x})\| \cdot \|\mathbf{x}^* - \mathbf{x}\| \\ &\leq -\frac{m}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2 + \|\nabla f(\mathbf{x})\| \cdot \|\mathbf{x}^* - \mathbf{x}\| \end{aligned}$$

which implies $m\|\mathbf{x} - \mathbf{x}^*\| \leq \|\nabla f(\mathbf{x})\|$. Similarly, we get

$$m\|\mathbf{x} - \mathbf{x}^*\| \leq \|\nabla f(\mathbf{x})\| \leq M\|\mathbf{x} - \mathbf{x}^*\| \quad (5.4)$$

- **(Upper Bounding left and right side of (5.1))** Moreover, we upper bounding the left side of (5.1) by setting $\mathbf{x}_2 = \mathbf{x}^{k+1}$ and $\mathbf{x}_1 = \mathbf{x}^k$ in (5.2):

$$\begin{aligned} f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k) &\leq \langle \nabla f(\mathbf{x}^k), \mathbf{x}^{k+1} - \mathbf{x}^k \rangle + \frac{M}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \\ &\leq -\frac{1}{2M} \|\nabla f(\mathbf{x}^k)\|^2 \end{aligned} \quad (5.5)$$

where the second inequality is active when $\mathbf{x}^{k+1} = \mathbf{x}^k - \frac{1}{M}\nabla f(\mathbf{x}^k)$. On the other hand, by setting $\mathbf{x}_2 = \mathbf{x}^*$ and $\mathbf{x}_1 = \mathbf{x}^k$ in (5.2), we obtain

$$\begin{aligned} f(\mathbf{x}^k) - f(\mathbf{x}^*) &\leq \langle \nabla f(\mathbf{x}^k), \mathbf{x}^k - \mathbf{x}^* \rangle - \frac{m}{2} \|\mathbf{x}^k - \mathbf{x}^*\|_2^2 \\ &\leq \|\nabla f(\mathbf{x}^k)\| \|\mathbf{x}^k - \mathbf{x}^*\| - \frac{m}{2} \|\mathbf{x}^k - \mathbf{x}^*\|_2^2 \\ &\leq \frac{1}{2m} \|\nabla f(\mathbf{x}^k)\|^2 \end{aligned} \tag{5.6}$$

Therefore, substituting (5.6) into (5.5), we obtain

$$f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k) \leq -\frac{m}{M} [f(\mathbf{x}^k) - f(\mathbf{x}^*)]$$

Or equivalently,

$$f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*) \leq \left(1 - \frac{m}{M}\right) [f(\mathbf{x}^k) - f(\mathbf{x}^*)]$$

□

question: this proof also holds for $t_k = \frac{1}{M}$. Thus what is the intuition behind the line search.

question: is uniformly convex and strongly convex talking about the same thing?

5.1.4 Gradient Algorithm with Diminishing Step Sizes

Consider a pre-scribed diminishing step size $\{\alpha_k\} \rightarrow 0$ but satisfies the infinite travel condition $\sum_{k=1}^{\infty} \alpha_k = \infty$.

In this case, for sufficiently large k , we have $\alpha_k \leq \frac{1}{M}$ and similar to the idea in (5.5),

$$f(\mathbf{x}^{k+1}) \leq f(\mathbf{x}^k) - \frac{\alpha_k}{2} \|\nabla f(\mathbf{x}^k)\|^2$$

which implies that $\nabla f(\mathbf{x}^k)$ cannot be bounded away from 0 whenever $f(\mathbf{x}^k)$ is finitely lower bounded. In other words, if a finite minimum exists for $f(\mathbf{x}^k)$, then the iterates satisfy $\lim_{k \rightarrow \infty} \inf \|\nabla f(\mathbf{x}^k)\| = 0$.

We can further show the whole sequence $f(\mathbf{x}^k)$ converges:

Proof. w.l.o.g., assume the inequality below holds for $k = 1, 2, \dots$, i.e., $\alpha_k \leq \frac{1}{M}$:

$$f(\mathbf{x}^{k+1}) \leq f(\mathbf{x}^k) - \frac{\alpha_k}{2} \|\nabla f(\mathbf{x}^k)\|^2$$

Therefore, for any $k = 1, 2, \dots$,

$$\begin{aligned} f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*) &\leq f(\mathbf{x}^k) - f(\mathbf{x}^*) - \frac{\alpha_k}{2} \|\nabla f(\mathbf{x}^k)\|^2 \\ &\leq (1 - m\alpha_k) [f(\mathbf{x}^k) - f(\mathbf{x}^*)] \end{aligned}$$

where the second inequality is by applying (5.6). It follows that

$$f(\mathbf{x}^n) - f(\mathbf{x}^*) \leq [f(\mathbf{x}^1) - f(\mathbf{x}^*)] \prod_{k=1}^n (1 - m\alpha_k) \rightarrow 0,$$

i.e., $\lim_{n \rightarrow \infty} f(\mathbf{x}^n) = f(\mathbf{x}^*)$.

There is another way to show the convergence of $\{\nabla f(\mathbf{x}^k)\}$:

$$\|\nabla f(\mathbf{x}^n)\|^2 \leq 2M[f(\mathbf{x}^n) - f(\mathbf{x}^*)] \implies \lim_{n \rightarrow \infty} \nabla f(\mathbf{x}^k) = \mathbf{0}.$$

□

We summarize the results above as a theorem for the convergence of the gradient algorithm with diminishing step sizes:

Theorem 5.1.2. Suppose there exists $0 < m \leq M$ such that $0 \succ mI \succeq \nabla^2 f(x) \succeq MI$ (i.e., f is *uniformly convex*), and the diminishing step size is performed per iteration:

$$\alpha_k \rightarrow 0, \text{ but } \sum_{k=1}^{\infty} \alpha_k = \infty,$$

then either $f(\mathbf{x}^k) \rightarrow -\infty$ or else $\{f(\mathbf{x}^k)\}$ converges to a finite value and $\nabla f(\mathbf{x}^k) \rightarrow \mathbf{0}$.

5.1.5 Gradient Algorithm with Armijo's Rule

Consider a general iterative descent algorithm $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k$. The Armijo's rule for choosing step sizes is as follows:

Let $\gamma \in (0, 1)$ (question: $1/2$?). Start with $s > 0$ and continue with $\beta s, \beta^2 s, \dots$, until $\beta^\ell s$ falls within the set of α with the condition

$$f(\mathbf{x}^k) - f(\mathbf{x}^k + \alpha \mathbf{d}^k) \geq -\gamma \alpha \cdot \nabla^T f(\mathbf{x}^k) \mathbf{d}^k$$

In this case we have $\alpha_k = s\beta^\ell$ and

$$f(\mathbf{x}^k) \geq f(\mathbf{x}^k + \alpha_k \mathbf{d}^k) - \gamma \alpha_k \nabla^T f(\mathbf{x}^k) \mathbf{d}^k \quad (5.7a)$$

$$f(\mathbf{x}^k) < f(\mathbf{x}^k + \alpha_k / \beta \mathbf{d}^k) - \gamma \alpha_k / \beta \cdot \nabla^T f(\mathbf{x}^k) \mathbf{d}^k \quad (5.7b)$$

We can analysis the convergence result for gradient algorithm, i.e., $\mathbf{d}^k = -\nabla f(\mathbf{x}^k)$:

- From the (5.7b) and the Taylor expansion on $f(\mathbf{x}^k + \alpha_k \mathbf{d}^k)$ we obtain:

$$f(\mathbf{x}^k) + \gamma \alpha_k / \beta \cdot \nabla^T f(\mathbf{x}^k) \mathbf{d}^k < f(\mathbf{x}^k) + \alpha_k / \beta \nabla^T f(\mathbf{x}^k) \mathbf{d}^k + \frac{M}{2} (\alpha_k / \beta)^2 \|\mathbf{d}^k\|^2$$

Or equivalently, $\alpha_k > \frac{2\beta(1-\gamma)}{M}$

- Combining the (5.7a), (5.6) and the bound on α_k , we obtain

$$f(\mathbf{x}^k + \alpha_k \mathbf{d}^k) \leq f(\mathbf{x}^k) - 4\beta\gamma(1 - \gamma)\frac{m}{M}[f(\mathbf{x}^k) - f(\mathbf{x}^*)]$$

Therefore, we get the Q-linear convergence for Armijo's rule:

$$f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*) \leq \left(1 - 4\beta\gamma(1 - \gamma)\frac{m}{M}\right)[f(\mathbf{x}^k) - f(\mathbf{x}^*)]$$

5.1.6 The Gradient Algorithm for non-strongly convex case

The estimations of convergence so far are based on the assumption that $m > 0$. Now we discuss the case where $m = 0$. The function is convex but not necessarily strongly convex.

Assume that the set of optimal solutions is a bounded set, and that there is a bounded *level set*. If still apply the exact line search, the iterates will be bounded. Note that the inequalities below still hold:

$$\begin{aligned} f(\mathbf{x} + \alpha \mathbf{d}) &\leq f(\mathbf{x}) - \frac{1}{2M}\|\nabla f(\mathbf{x})\|^2 \\ f(\mathbf{x}) - f(\mathbf{x}^*) &\leq \|\nabla f(\mathbf{x})\| \cdot \|\mathbf{x} - \mathbf{x}^*\| \end{aligned}$$

Assume that $\|\mathbf{x}^k - \mathbf{x}^*\| \leq C$, and let $e(\mathbf{x}^k) = f(\mathbf{x}^k) - f(\mathbf{x}^*)$, Using the inequalities above, it's easy to show that

$$e(\mathbf{x}^{k+1}) \leq e(\mathbf{x}^k) - c[e(\mathbf{x}^k)]^2, \text{ where } c = \frac{1}{2MC^2}.$$

which follows that

$$\begin{aligned} \frac{1}{e(\mathbf{x}^{k+1})} &\geq \frac{1}{e(\mathbf{x}^k)} + \frac{c}{1 - c \cdot e(\mathbf{x}^k)} \\ &\geq \frac{1}{e(\mathbf{x}^k)} + c \\ &\geq \dots \\ &\geq \frac{1}{e(\mathbf{x}^1)} + k \cdot c \end{aligned}$$

Therefore, we obtain the *sublinear rate of convergence*:

$$e(\mathbf{x}^{k+1}) \leq \frac{e(\mathbf{x}^1)}{1 + k(c \cdot e(\mathbf{x}^1))}$$

5.1.7 Linear Convergence without Second Order Differentiability

Acutally, the assumptions on the existence of $\nabla^2 f$ is unnecessary in Theorem (5.1.1). We can weaken the condition by the inequality below to obtain the same linear convergence result:

$$\sigma\|\mathbf{x} - \mathbf{y}\|^2 \leq \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \leq L\|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y}, \quad (5.8)$$

where $0 < \sigma \leq L < \infty$.

Remark 5.1.1. • The condition (5.8) can be implied by uniform convexity.

- The interpretation of (5.8) is that, restricting f to any line segment between \mathbf{x} and \mathbf{y} , the function $h(t) := f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))$ satisfies

$$0 \leq \frac{h'(t) - h'(s)}{t - s} \leq L, \quad \forall 0 \leq s < t \leq 1,$$

i.e., the slope of ∇f is bounded.

- The condition (5.8) implies the strong convexity, which can be shown by applying the directional derivative and (5.8):

$$\frac{\sigma}{2} \|\mathbf{y} - \mathbf{x}\|^2 \leq f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \leq \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2$$

Therefore, we can use the same logic to show the following inequalities:

$$\begin{aligned} f(\mathbf{x} - \alpha \nabla f(\mathbf{x})) - f(\mathbf{x}) &\leq -\frac{1}{2L} \|\nabla f(\mathbf{x})\|^2 \\ \sigma \|\mathbf{x} - \mathbf{x}^*\|^2 &\leq \|\nabla f(\mathbf{x})\| \|\mathbf{x} - \mathbf{x}^*\| \\ f(\mathbf{x}^*) &\geq f(\mathbf{x}) - \frac{1}{2\sigma} \|\nabla f(\mathbf{x})\|^2 \end{aligned}$$

and therefore,

$$f(\mathbf{x} - \alpha \nabla f(\mathbf{x})) - f(\mathbf{x}^*) \leq \left(1 - \frac{\sigma}{L}\right) [f(\mathbf{x}) - f(\mathbf{x}^*)].$$

5.2 The Pure Newton's Method

Now we discuss a particularly important method in optimization: Newton's method.

Motivation This method is a *linearization scheme* for solving a nonlinear equation.

- For scalar form of nonlinear equation $g(x) = 0$, we apply Taylor's expansion on the root \hat{x} :

$$g(\hat{x}) = g(x) + g'(x)(\hat{x} - x) + o(|\hat{x} - x|)$$

Ignoring the high order part we get an approximation, i.e., iterative formula

$$\bar{x} = x - \frac{g(x)}{g'(x)}.$$

- Consider a n -dimensional equation $g_{1:n}(x_1, \dots, x_n) = 0$, we have a similar solution

$$g(\hat{\mathbf{x}}) = g(\mathbf{x}) + J(g(\mathbf{x})) \cdot (\hat{\mathbf{x}} - \mathbf{x}) + o(\|\hat{\mathbf{x}} - \mathbf{x}\|)$$

where $J(g(\mathbf{x}))$ denotes the Jacobian matrix of g :

$$\mathbb{R}^{n \times n} \ni J(g(\mathbf{x})) := \left[\frac{\partial g_i(\mathbf{x})}{\partial x_j} \right]$$

Therefore, the unconstrained optimization problem suffices to solve a nonlinear equation $\nabla f(\mathbf{x}) = 0$, and the iterative formula is

$$\bar{\mathbf{x}} = \mathbf{x} - [\nabla^2 f(\mathbf{x})]^{-1} \nabla f(\mathbf{x}), \quad (\text{Newton's Method})$$

- Remark 5.2.1.**
1. Newton's direction may not necessarily exist;
 2. It is a descent direction for strongly convex functions;
 3. However, the function may not necessarily decrease even for strongly convex function.
 4. It minimizes a strongly convex *quadratic* function in just *one* step.
 5. The pure form of Newton's method can be modified by taking another step length.

5.2.1 Local Convergence Analysis

We analysis the convergence rate for Newton's method under the convexity and continuity conditions first:

Assumption: The function f is *convex, twice continuously differentiable*, and that $\nabla^2 f(\mathbf{x}^*)$ is non-singular for local minimum \mathbf{x}^* .

A key inequality for the analysis is

$$\nabla f(\mathbf{y}) = \nabla f(\mathbf{x}) + \int_0^1 \nabla^2 f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) \cdot (\mathbf{y} - \mathbf{x}) dt$$

Suppose that \mathbf{x}^k is close to \mathbf{x}^* enough, then $\nabla^2 f(\mathbf{x}^k)$ is non-singular as well due to the continuity of determinant function. It follows that

$$\begin{aligned} \mathbf{x}^{k+1} - \mathbf{x}^* &= \mathbf{x}^k - \mathbf{x}^* - [\nabla^2 f(\mathbf{x}^k)]^{-1} \nabla f(\mathbf{x}^k) \\ &= [\nabla^2 f(\mathbf{x}^k)]^{-1} [\nabla^2 f(\mathbf{x}^k)(\mathbf{x}^k - \mathbf{x}^*) - \nabla f(\mathbf{x}^k)] \\ &= [\nabla^2 f(\mathbf{x}^k)]^{-1} \left[\nabla^2 f(\mathbf{x}^k)(\mathbf{x}^k - \mathbf{x}^*) - \int_0^1 \nabla^2 f(\mathbf{x}^* + t(\mathbf{x}^k - \mathbf{x}^*))(\mathbf{x}^k - \mathbf{x}^*) dt \right] \\ &= [\nabla^2 f(\mathbf{x}^k)]^{-1} \left\{ \int_0^1 [\nabla^2 f(\mathbf{x}^k) - \nabla^2 f(\mathbf{x}^* + t(\mathbf{x}^k - \mathbf{x}^*))](\mathbf{x}^k - \mathbf{x}^*) dt \right\} \end{aligned}$$

Therefore,

$$\|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq \|\mathbf{x}^k - \mathbf{x}^*\| \cdot \|[\nabla^2 f(\mathbf{x}^k)]^{-1}\| \cdot \int_0^1 \|\nabla^2 f(\mathbf{x}^k) - \nabla^2 f(\mathbf{x}^* + t(\mathbf{x}^k - \mathbf{x}^*))\| dt$$

Since \mathbf{x}^k is close to \mathbf{x}^* , $\|[\nabla^2 f(\mathbf{x}^k)]^{-1}\|$ is bounded. Since $\nabla^2 f(\mathbf{x})$ is continuous, the integration term goes to zero as $\|\mathbf{x}^k - \mathbf{x}^*\| \rightarrow 0$. Thus we imply $\|\mathbf{x}^{k+1} - \mathbf{x}^*\| = o(\|\mathbf{x}^k - \mathbf{x}^*\|)$, ensuring a superlinear convergence.

Extra Assumption: The term $\nabla^2 f(\mathbf{x})$ is *Lipschitz continuous*: there exists $L_2 > 0$ such that

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\| \leq L_2 \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y}.$$

This extra assumption will ensure a *quadratic convergence rate*:

$$\begin{aligned} \|\mathbf{x}^{k+1} - \mathbf{x}^*\| &\leq \|[\nabla^2 f(\mathbf{x}^k)]^{-1}\| \cdot \|\mathbf{x}^k - \mathbf{x}^*\| \cdot \int_0^1 \|\nabla^2 f(\mathbf{x}^k) - \nabla^2 f(\mathbf{x}^* + t(\mathbf{x}^k - \mathbf{x}^*))\| dt \\ &\leq \frac{L_2}{2} \|[\nabla^2 f(\mathbf{x}^k)]^{-1}\| \cdot \|\mathbf{x}^k - \mathbf{x}^*\|^2. \end{aligned}$$

Further Assumption: Based on the previous two assumptions, we assume that f is *strongly convex*.

In this case, it is easy to show that

$$\|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq \frac{L_2}{2m} \|\mathbf{x}^k - \mathbf{x}^*\|^2.$$

This inequality introduces a *region of attraction*, i.e., as soon as \mathbf{x}^k falls into the neighborhood of \mathbf{x}^* with radius $2m/L_2$, the iterates will be trapped in the neighborhood and converge to \mathbf{x}^* *quadratically*.

Remark 5.2.2. The pure form of Newton's method, however, has several drawbacks:

1. It in general does not guarantees global convergence if no additional assumption is given. Fortunartely, if f is strongly convex, then the Newton's method with *line-search* (e.g., with Armijo's step-length rule) will be globally convergent with a globally linear convergence rate.
2. If f is not *strictly* convex, then $\nabla^2 f$ may be singular. Even worse, if f is not convex, then the Newton's direction may not be a *descent direction*. In next section we will discuss how to handle such a situation.

5.3 Practical Implementation of Newton's method

5.3.1 Cholesky Factorization

First let's introduce a technique in optimization algorithms that can reduce computational complexity: the *Cholesky factorization*.

Consider the case where $\nabla^2 f(\mathbf{x}^k) \succ 0$, and the Newton's direction can be found by solving the linear system

$$\nabla^2 f(\mathbf{x}^k) \mathbf{d} = -\nabla f(\mathbf{x}^k).$$

Directly computing the inverse of $\nabla^2 f(\mathbf{x}^k)$ is computationally expensive, which motivates us to apply the *Cholesky factorization* as follows:

1. First apply the Cholesky factorization to get $\nabla^2 f(\mathbf{x}^k) = \mathbf{L}_k \mathbf{L}_k^T$, where \mathbf{L}_k is a lower triangular matrix, resulting in the following Newton's equation

$$\mathbf{L}_k \mathbf{L}_k^T \mathbf{d} = -\nabla f(\mathbf{x}^k).$$

2. Firstly solve the lower triangular system below by forward substitution:

$$\mathbf{L}_k \mathbf{y} = -\nabla f(\mathbf{x}^k)$$

The complexity for this process is $\mathcal{O}(n^2)$.

3. Then solve the triangular system below by backward substitution:

$$\mathbf{L}_k^T \mathbf{d} = \mathbf{y}_k$$

Again, this step takes complexity $\mathcal{O}(n^2)$.

The basic Cholesky factorization algorithm is as follows:

Remark 5.3.1. If \mathbf{A} is not positive semidefinite, then at a certain stage we will encounter a j such that

$$a_{jj} - \sum_{k=1}^{j-1} l_{jk}^2 < 0.$$

In that case, the Cholesky decomposition cannot proceed. Note that the Cholesky decomposition takes about $\mathcal{O}(n^3)$ operations.

5.3.2 Modified Newton's method

In case the Hessian matrix is not positive definite, the following remedies can be applied:

If there occurs $a_{jj} - \sum_{k=1}^{j-1} l_{jk}^2 < 0$ for a certain j , then we simply increase a_{jj} so that the quantity becomes positive again. (question: increase how much?)

Algorithm 3 Basic Cholesky factorization Algorithm

Input: A positive definite $n \times n$ matrix \mathbf{A}

Output: Lower triangular matrix \mathbf{L} such that $\mathbf{A} = \mathbf{L}\mathbf{L}^T$

For $j = 1 : n$, **do**

• **For** $i = j + 1 : n$, **do**

$$- l_{ij} = \left(a_{ij} - \sum_{k=1}^{j-1} l_{jk} l_{ik} \right) / l_{jj}$$

end for.

$$l_{jj} = \left(a_{jj} - \sum_{k=1}^{j-1} l_{jk}^2 \right)^{1/2}.$$

end for.

This remedy has the same effect of changing $\nabla^2 f(\mathbf{x}^k)$ into $\nabla^2 f(\mathbf{x}^k) + \Delta^k \succ 0$, where Δ^k is non-negative (question: positive or non-negative?) diagonal, which suffices to solve the regularized equation

$$(\nabla^2 f(\mathbf{x}^k) + \Delta^k) \mathbf{d} = -\nabla f(\mathbf{x}^k).$$

Moreover, we may use the direction with Armijo's line search technique to guarantee the global convergence. (how to show?)

5.3.3 The Trust Region Approach

Another way to handle the case that $\nabla^2 f(\mathbf{x}^k)$ is indefinite is to use the *trust region approach*. It is the complement of the line search approach.

The direction \mathbf{d}^k for each iteration suffices to consider the trust region subproblem

$$\begin{aligned} \min \quad & \langle \nabla f(\mathbf{x}^k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{x}^k) \mathbf{d} \\ \text{such that} \quad & \|\mathbf{d}\| \leq \delta \end{aligned}$$

where $\delta > 0$ is called the trust region radius.

Remark 5.3.2. It can be shown that when δ is sufficiently small, $f(\mathbf{x}^k + \mathbf{d}^k) < f(\mathbf{x}^k)$, i.e., \mathbf{d}^k is the descent direction. This trust region subproblem can be efficiently solved (question: which method? curious about it)

5.3.4 Implementation of Least Squares Problem

Consider solving the *nonlinear least square problem* (NLSP)

$$\min \quad f(x) = \frac{1}{2} \sum_{i=1}^m f_i^2(x)$$

Firstly note that

$$\begin{aligned}\nabla f(x) &= \sum_{i=1}^m f_i(x) \nabla f_i(x) \\ \nabla^2 f(x) &= \sum_{i=1}^m [\nabla f_i(x) \nabla^T f_i(x) + f_i(x) \nabla^2 f_i(x)]\end{aligned}$$

The so-called Gauss-Newton method is a *quasi-Newton's method*, specialized to this NLSP:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \left(\sum_{i=1}^m \nabla f_i(\mathbf{x}^k) \nabla^T f_i(\mathbf{x}^k) \right)^{-1} \left(\sum_{i=1}^m f_i(\mathbf{x}^k) \nabla f_i(\mathbf{x}^k) \right)$$

Remark 5.3.3. It works well when f_i 's are not *too linear*, or when at the optimality, f_i 's are close to zero.

A variation of the Gauss-Newton's method operates as follows:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \left(\sum_{i=1}^m \nabla f_i(\mathbf{x}^k) \nabla^T f_i(\mathbf{x}^k) + \lambda_k \mathbf{I} \right)^{-1} \left(\sum_{i=1}^m f_i(\mathbf{x}^k) \nabla f_i(\mathbf{x}^k) \right)$$

which is called the *Levenberg-Marquardt method*.

Note that if consider solving the equation $f_{1:n}(\mathbf{x}_{1:n}) = \mathbf{0}$, the Gauss-Newton direction is just the Newton direction itself.

6 Primal-Dual Interior Point Methods (PDIPM)

6.1 PDIPM for Linear Programming

Consider a linear programming problem

$$\begin{aligned}\min \quad & \mathbf{c}^T \mathbf{x} \\ \text{such that} \quad & \mathbf{A} \mathbf{x} = \mathbf{b} \quad (P) \\ & \mathbf{x} \geq 0\end{aligned} \tag{6.1}$$

with its dual problem

$$\begin{aligned}\max \quad & \mathbf{b}^T \mathbf{y} \\ \text{such that} \quad & \mathbf{A}^T \mathbf{y} + \mathbf{s} = \mathbf{c} \quad (D) \\ & \mathbf{s} \geq 0\end{aligned} \tag{6.2}$$

We assume that the *primal-dual Slater condition* (what is it?) holds. Consider solving the barriered problem of (P) for $\mu > 0$:

$$\begin{aligned}\min \quad & \mathbf{c}^T \mathbf{x} - \mu \sum_{i=1}^n \ln x_i \\ \text{such that} \quad & \mathbf{A} \mathbf{x} = \mathbf{b} \quad (P_\mu)\end{aligned} \tag{6.3}$$

6.1.1 Duality Gap

Let $\mathbf{x}(\mu)$ be the optimal solution for (P_μ) . By KKT condition, there exists $\mathbf{y}(\mu) \in \mathbb{R}^m$ such that

$$\mathbf{c} - \mu \mathbf{x}(\mu)^{-1} - \mathbf{A}^T \mathbf{y}(\mu) = \mathbf{0}.$$

Define the slack variable $\mathbf{s}(\mu) := \mu \mathbf{x}(\mu)^{-1} - \mathbf{A}^T \mathbf{y}(\mu)$. It's clear that:

1. $\mathbf{x}(\mu)$ is primal-feasible to (P_μ) ; while $(\mathbf{y}(\mu), \mathbf{s}(\mu))$ is dual-feasible to (D_μ)
2. The duality gap between $\mathbf{x}(\mu)$ and $(\mathbf{y}(\mu), \mathbf{s}(\mu))$ is

$$\mathbf{x}(\mu)^T \mathbf{s}(\mu) = n\mu.$$

6.1.2 Convergence of Barrier Problem P_μ

The set $\{\mathbf{x}(\mu) \mid \mu > 0\}$ is known as the *primal analytic central path*; and the set $\{(\mathbf{y}(\mu), \mathbf{s}(\mu)) \mid \mu > 0\}$ is known as the *dual analytic central path*.

Proposition 6.1.1. The set $\{\mathbf{x}(\mu) \mid 0 < \mu \leq 1\}$ and $\{(\mathbf{y}(\mu), \mathbf{s}(\mu)) \mid 0 < \mu \leq 1\}$ are bounded.

Proof. Let $\tilde{\mathbf{x}}$ be an interior for (P) and $(\tilde{\mathbf{y}}, \tilde{\mathbf{s}})$ be an interior for (D) . Then $\tilde{\mathbf{x}} - \mathbf{x}(\mu) \in \mathcal{N}(\mathbf{A})$ and $\tilde{\mathbf{s}} - \mathbf{s}(\mu) \in \text{Range}(\mathbf{A}^T)$, which implies

$$0 = (\tilde{\mathbf{x}} - \mathbf{x}(\mu))^T (\tilde{\mathbf{s}} - \mathbf{s}(\mu)) = \tilde{\mathbf{x}}^T \tilde{\mathbf{s}} - \tilde{\mathbf{x}}^T \mathbf{s}(\mu) - \tilde{\mathbf{s}}^T \mathbf{x}(\mu) + n\mu.$$

Therefore, $\{\mathbf{x}(\mu) \mid 0 < \mu \leq 1\}$ and $\{(\mathbf{y}(\mu), \mathbf{s}(\mu)) \mid 0 < \mu \leq 1\}$ must be bounded. (question: how to specify?) \square

Proposition 6.1.2. The set $\{(\mathbf{x}(\mu), \mathbf{y}(\mu), \mathbf{s}(\mu)) \mid 0 < \mu \leq 1\}$ converges to $(\hat{\mathbf{x}}, \hat{\mathbf{y}}, \hat{\mathbf{s}})$ as $\mu \rightarrow 0$. Moreover, the active sets $B := \{i \mid \hat{x}_i > 0\}$ and $N = \{j \mid \hat{s}_j > 0\}$ form a partition of $\{1, \dots, n\}$.

Proof. • Due to the boundness of $\{(\mathbf{x}(\mu), \mathbf{y}(\mu), \mathbf{s}(\mu)) \mid 0 < \mu \leq 1\}$, we imply there exists a subsequence μ_k such that

$$\lim_{k \rightarrow \infty} (\mathbf{x}(\mu), \mathbf{y}(\mu), \mathbf{s}(\mu)) = (\hat{\mathbf{x}}, \hat{\mathbf{y}}, \hat{\mathbf{s}})$$

It's clear that $\hat{\mathbf{x}}$ is primal optimal and $(\hat{\mathbf{y}}, \hat{\mathbf{s}})$ is dual optimal. By complementarity condition, $B \cap N = \emptyset$.

- At the same time, consider the equality again that

$$0 = (\hat{\mathbf{x}} - \mathbf{x}(\mu_k))^T (\hat{\mathbf{s}} - \mathbf{s}(\mu_k)) = - \sum_{i \in B} \hat{x}_i s_i(\mu_k) - \sum_{j \in N} \hat{s}_j x_j(\mu_k) + n\mu_k.$$

Or equivalently,

$$n = \sum_{i \in B} \frac{\hat{x}_i}{x_i(\mu_k)} + \sum_{j \in N} \frac{\hat{s}_j}{s_j(\mu_k)}.$$

Taking $k \rightarrow \infty$, we imply $|B| + |N| = n$, i.e., B and N form a partition of $\{1, \dots, n\}$.

- Next we show that $\hat{\mathbf{x}}$ is unique. We shall show that $\hat{\mathbf{x}}_B$ is the optimal solution to

$$\begin{aligned} \max \quad & \sum_{i \in B} \ln x_i \\ \text{such that} \quad & \mathbf{A}_B \mathbf{x}_B = \mathbf{b} \end{aligned} \quad (O)$$

and then the uniqueness of $\hat{\mathbf{x}}_B$ is proved due to the strict concavity of the objective function. (question: why the uniqueness of $\hat{\mathbf{x}}_B$ implies the uniqueness of $\hat{\mathbf{x}}$?)

The KKT condition for (O) gives

$$\mathbf{x}_B^{-1} \in \text{Range}(\mathbf{A}_B^T) \quad \mathbf{A}_B \mathbf{x}_B = \mathbf{b}$$

We may verify that $\hat{\mathbf{x}}_B$ satisfies the condition above.

□

Remark 6.1.1. The particular optimal solution \mathbf{x} , denoted by $\mathbf{x}(0)$ is called the *analytic center of the optimal face*. If taking $\mu \rightarrow \infty$, then $\mathbf{x}(\mu)$ converges to the optimal solution of

$$\begin{aligned} \max \quad & \sum_{i=1}^n \ln x_i \\ \text{such that} \quad & \mathbf{A}_B \mathbf{x}_B = \mathbf{b} \end{aligned}$$

which is known as the *analytic center of the feasible region*.

Note that $\mathbf{x}(0) + \mathbf{s}(0) > 0$, i.e., they are *strictly complementary*.

7 The Distribution and Installation

7.1 Pre-requisites

You will need a working LaTeX installation. We recomend using pdflatex to process the files. You will also need biber.exe installed. This is distributed as part of the latest versions of LiveTex and MikTex. If you have problems, please let us know.

7.2 The Distribution

The distribution contains 2 folders: `nowfnt` and `nowfnttexmf`.