

---

# AN INTRODUCTION TO OPTIMIZATION

SOLUTIONS MANUAL

---

Fourth Edition

Edwin K. P. Chong and Stanislaw H. Żak



A JOHN WILEY & SONS, INC., PUBLICATION

## 1. Methods of Proof and Some Notation

### 1.1

---

A	B	not A	not B	$A \Rightarrow B$	$(\text{not } B) \Rightarrow (\text{not } A)$
F	F	T	T	T	T
F	T	T	F	T	T
T	F	F	T	F	F
T	T	F	F	T	T

### 1.2

---

A	B	not A	not B	$A \Rightarrow B$	not (A and (not B))
F	F	T	T	T	T
F	T	T	F	T	T
T	F	F	T	F	F
T	T	F	F	T	T

### 1.3

---

A	B	not (A and B)	not A	not B	(not A) or (not B)
F	F	T	T	T	T
F	T	T	T	F	T
T	F	T	F	T	T
T	T	F	F	F	F

### 1.4

---

A	B	A and B	A and (not B)	(A and B) or (A and (not B))
F	F	F	F	F
F	T	F	F	F
T	F	F	T	T
T	T	T	F	T

### 1.5

---

The cards that you should turn over are 3 and  $A$ . The remaining cards are irrelevant to ascertaining the truth or falsity of the rule. The card with  $S$  is irrelevant because  $S$  is not a vowel. The card with 8 is not relevant because the rule does not say that if a card has an even number on one side, then it has a vowel on the other side.

Turning over the  $A$  card directly verifies the rule, while turning over the 3 card verifies the contraposition.

## 2. Vector Spaces and Matrices

### 2.1

---

We show this by contradiction. Suppose  $n < m$ . Then, the number of columns of  $\mathbf{A}$  is  $n$ . Since  $\text{rank } \mathbf{A}$  is the maximum number of linearly independent columns of  $\mathbf{A}$ , then  $\text{rank } \mathbf{A}$  cannot be greater than  $n < m$ , which contradicts the assumption that  $\text{rank } \mathbf{A} = m$ .

### 2.2

---

$\Rightarrow$ : Since there exists a solution, then by Theorem 2.1,  $\text{rank } \mathbf{A} = \text{rank}[\mathbf{A}; \mathbf{b}]$ . So, it remains to prove that  $\text{rank } \mathbf{A} = n$ . For this, suppose that  $\text{rank } \mathbf{A} < n$  (note that it is impossible for  $\text{rank } \mathbf{A} > n$  since  $\mathbf{A}$  has only  $n$  columns). Hence, there exists  $\mathbf{y} \in \mathbb{R}^n$ ,  $\mathbf{y} \neq \mathbf{0}$ , such that  $\mathbf{A}\mathbf{y} = \mathbf{0}$  (this is because the columns of

$\mathbf{A}$  are linearly dependent, and  $\mathbf{A}\mathbf{y}$  is a linear combination of the columns of  $\mathbf{A}$ ). Let  $\mathbf{x}$  be a solution to  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . Then clearly  $\mathbf{x} + \mathbf{y} \neq \mathbf{x}$  is also a solution. This contradicts the uniqueness of the solution. Hence,  $\text{rank } \mathbf{A} = n$ .

$\Leftarrow$ : By Theorem 2.1, a solution exists. It remains to prove that it is unique. For this, let  $\mathbf{x}$  and  $\mathbf{y}$  be solutions, i.e.,  $\mathbf{A}\mathbf{x} = \mathbf{b}$  and  $\mathbf{A}\mathbf{y} = \mathbf{b}$ . Subtracting, we get  $\mathbf{A}(\mathbf{x} - \mathbf{y}) = \mathbf{0}$ . Since  $\text{rank } \mathbf{A} = n$  and  $\mathbf{A}$  has  $n$  columns, then  $\mathbf{x} - \mathbf{y} = \mathbf{0}$  and hence  $\mathbf{x} = \mathbf{y}$ , which shows that the solution is unique.

### 2.3

Consider the vectors  $\bar{\mathbf{a}}_i = [1, \mathbf{a}_i^\top]^\top \in \mathbb{R}^{n+1}$ ,  $i = 1, \dots, k$ . Since  $k \geq n + 2$ , then the vectors  $\bar{\mathbf{a}}_1, \dots, \bar{\mathbf{a}}_k$  must be linearly independent in  $\mathbb{R}^{n+1}$ . Hence, there exist  $\alpha_1, \dots, \alpha_k$ , not all zero, such that

$$\sum_{i=1}^k \alpha_i \mathbf{a}_i = \mathbf{0}.$$

The first component of the above vector equation is  $\sum_{i=1}^k \alpha_i = 0$ , while the last  $n$  components have the form  $\sum_{i=1}^k \alpha_i \mathbf{a}_i = \mathbf{0}$ , completing the proof.

### 2.4

a. We first postmultiply  $\mathbf{M}$  by the matrix

$$\begin{bmatrix} \mathbf{I}_k & \mathbf{O} \\ -\mathbf{M}_{m-k,k} & \mathbf{I}_{m-k} \end{bmatrix}$$

to obtain

$$\begin{bmatrix} \mathbf{M}_{m-k,k} & \mathbf{I}_{m-k} \\ \mathbf{M}_{k,k} & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{I}_k & \mathbf{O} \\ -\mathbf{M}_{m-k,k} & \mathbf{I}_{m-k} \end{bmatrix} = \begin{bmatrix} \mathbf{O} & \mathbf{I}_{m-k} \\ \mathbf{M}_{k,k} & \mathbf{O} \end{bmatrix}.$$

Note that the determinant of the postmultiplying matrix is 1. Next we postmultiply the resulting product by

$$\begin{bmatrix} \mathbf{O} & \mathbf{I}_k \\ \mathbf{I}_{m-k} & \mathbf{O} \end{bmatrix}$$

to obtain

$$\begin{bmatrix} \mathbf{O} & \mathbf{I}_{m-k} \\ \mathbf{M}_{k,k} & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{O} & \mathbf{I}_k \\ \mathbf{I}_{m-k} & \mathbf{O} \end{bmatrix} = \begin{bmatrix} \mathbf{I}_k & \mathbf{O} \\ \mathbf{O} & \mathbf{M}_{k,k} \end{bmatrix}.$$

Notice that

$$\det \mathbf{M} = \det \left( \begin{bmatrix} \mathbf{I}_k & \mathbf{O} \\ \mathbf{O} & \mathbf{M}_{k,k} \end{bmatrix} \right) \det \left( \begin{bmatrix} \mathbf{O} & \mathbf{I}_k \\ \mathbf{I}_{m-k} & \mathbf{O} \end{bmatrix} \right),$$

where

$$\det \left( \begin{bmatrix} \mathbf{O} & \mathbf{I}_k \\ \mathbf{I}_{m-k} & \mathbf{O} \end{bmatrix} \right) = \pm 1.$$

The above easily follows from the fact that the determinant changes its sign if we interchange columns, as discussed in Section 2.2. Moreover,

$$\det \left( \begin{bmatrix} \mathbf{I}_k & \mathbf{O} \\ \mathbf{O} & \mathbf{M}_{k,k} \end{bmatrix} \right) = \det(\mathbf{I}_k) \det(\mathbf{M}_{k,k}) = \det(\mathbf{M}_{k,k}).$$

Hence,

$$\det \mathbf{M} = \pm \det \mathbf{M}_{k,k}.$$

b. We can see this on the following examples. We assume, without loss of generality that  $\mathbf{M}_{m-k,k} = \mathbf{O}$  and let  $\mathbf{M}_{k,k} = 2$ . Thus  $k = 1$ . First consider the case when  $m = 2$ . Then we have

$$\mathbf{M} = \begin{bmatrix} \mathbf{O} & \mathbf{I}_{m-k} \\ \mathbf{M}_{k,k} & \mathbf{O} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 2 & 0 \end{bmatrix}.$$

Thus,

$$\det \mathbf{M} = -2 = \det (-\mathbf{M}_{k,k}).$$

Next consider the case when  $m = 3$ . Then

$$\det \begin{bmatrix} \mathbf{O} & \mathbf{I}_{m-k} \\ \mathbf{M}_{k,k} & \mathbf{O} \end{bmatrix} = \det \begin{bmatrix} 0 & \vdots & 1 & 0 \\ 0 & \vdots & 0 & 1 \\ \dots & \dots & \dots & \dots \\ 2 & \vdots & 0 & 0 \end{bmatrix} = 2 \neq \det (-\mathbf{M}_{k,k}).$$

Therefore, in general,

$$\det \mathbf{M} \neq \det (-\mathbf{M}_{k,k})$$

However, when  $k = m/2$ , that is, when all sub-matrices are square and of the same dimension, then it is true that

$$\det \mathbf{M} = \det (-\mathbf{M}_{k,k}).$$

See [121].

## 2.5

Let

$$\mathbf{M} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}$$

and suppose that each block is  $k \times k$ . John R. Sylvester [121] showed that if at least one of the blocks is equal to  $\mathbf{O}$  (zero matrix), then the desired formula holds. Indeed, if a row or column block is zero, then the determinant is equal to zero as follows from the determinant's properties discussed Section 2.2. That is, if  $\mathbf{A} = \mathbf{B} = \mathbf{O}$ , or  $\mathbf{A} = \mathbf{C} = \mathbf{O}$ , and so on, then obviously  $\det \mathbf{M} = 0$ . This includes the case when any three or all four block matrices are zero matrices.

If  $\mathbf{B} = \mathbf{O}$  or  $\mathbf{C} = \mathbf{O}$  then

$$\det \mathbf{M} = \det \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \det (\mathbf{AD}).$$

The only case left to analyze is when  $\mathbf{A} = \mathbf{O}$  or  $\mathbf{D} = \mathbf{O}$ . We will show that in either case,

$$\det \mathbf{M} = \det (-\mathbf{BC}).$$

Without loss of generality suppose that  $\mathbf{D} = \mathbf{O}$ . Following arguments of John R. Sylvester [121], we premultiply  $\mathbf{M}$  by the product of three matrices whose determinants are unity:

$$\begin{bmatrix} \mathbf{I}_k & -\mathbf{I}_k \\ \mathbf{O} & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{I}_k & \mathbf{O} \\ \mathbf{I}_k & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{I}_k & -\mathbf{I}_k \\ \mathbf{O} & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{O} \end{bmatrix} = \begin{bmatrix} -\mathbf{C} & \mathbf{O} \\ \mathbf{A} & \mathbf{B} \end{bmatrix}.$$

Hence,

$$\begin{aligned} \det \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{O} \end{bmatrix} &= \det \begin{bmatrix} -\mathbf{C} & \mathbf{O} \\ \mathbf{A} & \mathbf{B} \end{bmatrix} \\ &= \det (-\mathbf{C}) \det \mathbf{B} \\ &= \det (-\mathbf{I}_k) \det \mathbf{C} \det \mathbf{B}. \end{aligned}$$

Thus we have

$$\det \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{O} \end{bmatrix} = \det (-\mathbf{BC}) = \det (-\mathbf{CB}).$$

## 2.6

---

We represent the given system of equations in the form  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , where

$$\mathbf{A} = \begin{bmatrix} 1 & 1 & 2 & 1 \\ 1 & -2 & 0 & -1 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}, \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}.$$

Using elementary row operations yields

$$\mathbf{A} = \begin{bmatrix} 1 & 1 & 2 & 1 \\ 1 & -2 & 0 & -1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 1 & 2 & 1 \\ 0 & -3 & -2 & -2 \end{bmatrix}, \quad \text{and}$$
$$[\mathbf{A}, \mathbf{b}] = \begin{bmatrix} 1 & 1 & 2 & 1 & 1 \\ 1 & -2 & 0 & -1 & -2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 1 & 2 & 1 & 1 \\ 0 & -3 & -2 & -2 & -3 \end{bmatrix},$$

from which  $\text{rank } \mathbf{A} = 2$  and  $\text{rank}[\mathbf{A}, \mathbf{b}] = 2$ . Therefore, by Theorem 2.1, the system has a solution.

We next represent the system of equations as

$$\begin{bmatrix} 1 & 1 \\ 1 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 - 2x_3 - x_4 \\ -2 + x_4 \end{bmatrix}$$

Assigning arbitrary values to  $x_3$  and  $x_4$  ( $x_3 = d_3$ ,  $x_4 = d_4$ ), we get

$$\begin{aligned} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} &= \begin{bmatrix} 1 & 1 \\ 1 & -2 \end{bmatrix}^{-1} \begin{bmatrix} 1 - 2x_3 - x_4 \\ -2 + x_4 \end{bmatrix} \\ &= -\frac{1}{3} \begin{bmatrix} -2 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 - 2x_3 - x_4 \\ -2 + x_4 \end{bmatrix} \\ &= \begin{bmatrix} -\frac{4}{3}d_3 - \frac{1}{3}d_4 \\ 1 - \frac{2}{3}d_3 - \frac{2}{3}d_4 \end{bmatrix}. \end{aligned}$$

Therefore, a general solution is

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} -\frac{4}{3}d_3 - \frac{1}{3}d_4 \\ 1 - \frac{2}{3}d_3 - \frac{2}{3}d_4 \\ d_3 \\ d_4 \end{bmatrix} = \begin{bmatrix} -\frac{4}{3} \\ -\frac{2}{3} \\ 1 \\ 0 \end{bmatrix} d_3 + \begin{bmatrix} -\frac{1}{3} \\ -\frac{2}{3} \\ 0 \\ 1 \end{bmatrix} d_4 + \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix},$$

where  $d_3$  and  $d_4$  are arbitrary values.

## 2.7

---

1. Apply the definition of  $|-a|$ :

$$\begin{aligned} |-a| &= \begin{cases} -a & \text{if } -a > 0 \\ 0 & \text{if } -a = 0 \\ -(-a) & \text{if } -a < 0 \end{cases} \\ &= \begin{cases} -a & \text{if } a < 0 \\ 0 & \text{if } a = 0 \\ a & \text{if } a > 0 \end{cases} \\ &= |a|. \end{aligned}$$

2. If  $a \geq 0$ , then  $|a| = a$ . If  $a < 0$ , then  $|a| = -a > 0 > a$ . Hence  $|a| \geq a$ . On the other hand,  $|-a| \geq -a$  (by the above). Hence,  $a \geq -|-a| = -|a|$  (by property 1).

3. We have four cases to consider. First, if  $a, b \geq 0$ , then  $a + b \geq 0$ . Hence,  $|a + b| = a + b = |a| + |b|$ .  
 Second, if  $a, b \leq 0$ , then  $a + b \leq 0$ . Hence  $|a + b| = -(a + b) = -a - b = |a| + |b|$ .  
 Third, if  $a \geq 0$  and  $b \leq 0$ , then we have two further subcases:

1. If  $a + b \geq 0$ , then  $|a + b| = a + b \leq |a| + |b|$ .
2. If  $a + b \leq 0$ , then  $|a + b| = -a - b \leq |a| + |b|$ .

The fourth case,  $a \leq 0$  and  $b \geq 0$ , is identical to the third case, with  $a$  and  $b$  interchanged.

4. We first show  $|a - b| \leq |a| + |b|$ . We have

$$\begin{aligned} |a - b| &= |a + (-b)| \\ &\leq |a| + |-b| \quad \text{by property 3} \\ &= |a| + |b| \quad \text{by property 1.} \end{aligned}$$

To show  $||a| - |b|| \leq |a - b|$ , we note that  $|a| = |a - b + b| \leq |a - b| + |b|$ , which implies  $|a| - |b| \leq |a - b|$ . On the other hand, from the above we have  $|b| - |a| \leq |b - a| = |a - b|$  by property 1. Therefore,  $||a| - |b|| \leq |a - b|$ .

5. We have four cases. First, if  $a, b \geq 0$ , we have  $ab \geq 0$  and hence  $|ab| = ab = |a||b|$ . Second, if  $a, b \leq 0$ , we have  $ab \geq 0$  and hence  $|ab| = ab = (-a)(-b) = |a||b|$ . Third, if  $a \leq 0, b \geq 0$ , we have  $ab \leq 0$  and hence  $|ab| = -ab = a(-b) = |a||b|$ . The fourth case,  $a \geq 0$  and  $b \leq 0$ , is identical to the third case, with  $a$  and  $b$  interchanged.

6. We have

$$\begin{aligned} |a + b| &\leq |a| + |b| \quad \text{by property 3} \\ &\leq c + d. \end{aligned}$$

7.  $\Rightarrow$ : By property 2,  $-a \leq |a|$  and  $a \leq |a|$ . Therefore,  $|a| < b$  implies  $-a \leq |a| < b$  and  $a \leq |a| < b$ .

$\Leftarrow$ : If  $a \geq 0$ , then  $|a| = a < b$ . If  $a < 0$ , then  $|a| = -a < b$ .

For the case when “ $<$ ” is replaced by “ $\leq$ ”, we simply repeat the above proof with “ $<$ ” replaced by “ $\leq$ ”.

8. This is simply the negation of property 7 (apply DeMorgan’s Law).

## 2.8

Observe that we can represent  $\langle \mathbf{x}, \mathbf{y} \rangle_2$  as

$$\langle \mathbf{x}, \mathbf{y} \rangle_2 = \mathbf{x}^\top \begin{bmatrix} 2 & 3 \\ 3 & 5 \end{bmatrix} \mathbf{y} = (\mathbf{Q}\mathbf{x})^\top (\mathbf{Q}\mathbf{y}) = \mathbf{x}^\top \mathbf{Q}^2 \mathbf{y},$$

where

$$\mathbf{Q} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}.$$

Note that the matrix  $\mathbf{Q} = \mathbf{Q}^\top$  is nonsingular.

1. Now,  $\langle \mathbf{x}, \mathbf{x} \rangle_2 = (\mathbf{Q}\mathbf{x})^\top (\mathbf{Q}\mathbf{x}) = \|\mathbf{Q}\mathbf{x}\|^2 \geq 0$ , and

$$\begin{aligned} \langle \mathbf{x}, \mathbf{x} \rangle_2 = 0 &\Leftrightarrow \|\mathbf{Q}\mathbf{x}\|^2 = 0 \\ &\Leftrightarrow \mathbf{Q}\mathbf{x} = \mathbf{0} \\ &\Leftrightarrow \mathbf{x} = \mathbf{0} \end{aligned}$$

since  $\mathbf{Q}$  is nonsingular.

2.  $\langle \mathbf{x}, \mathbf{y} \rangle_2 = (\mathbf{Q}\mathbf{x})^\top (\mathbf{Q}\mathbf{y}) = (\mathbf{Q}\mathbf{y})^\top (\mathbf{Q}\mathbf{x}) = \langle \mathbf{y}, \mathbf{x} \rangle_2$ .

3. We have

$$\begin{aligned} \langle \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle_2 &= (\mathbf{x} + \mathbf{y})^\top \mathbf{Q}^2 \mathbf{z} \\ &= \mathbf{x}^\top \mathbf{Q}^2 \mathbf{z} + \mathbf{y}^\top \mathbf{Q}^2 \mathbf{z} \\ &= \langle \mathbf{x}, \mathbf{z} \rangle_2 + \langle \mathbf{y}, \mathbf{z} \rangle_2. \end{aligned}$$

$$4. \langle r\mathbf{x}, \mathbf{y} \rangle_2 = (r\mathbf{x})^\top \mathbf{Q}^2 \mathbf{y} = r\mathbf{x}^\top \mathbf{Q}^2 \mathbf{y} = r\langle \mathbf{x}, \mathbf{y} \rangle_2.$$

### 2.9

We have  $\|\mathbf{x}\| = \|(\mathbf{x} - \mathbf{y}) + \mathbf{y}\| \leq \|\mathbf{x} - \mathbf{y}\| + \|\mathbf{y}\|$  by the Triangle Inequality. Hence,  $\|\mathbf{x}\| - \|\mathbf{y}\| \leq \|\mathbf{x} - \mathbf{y}\|$ . On the other hand, from the above we have  $\|\mathbf{y}\| - \|\mathbf{x}\| \leq \|\mathbf{y} - \mathbf{x}\| = \|\mathbf{x} - \mathbf{y}\|$ . Combining the two inequalities, we obtain  $|\|\mathbf{x}\| - \|\mathbf{y}\|| \leq \|\mathbf{x} - \mathbf{y}\|$ .

### 2.10

Let  $\epsilon > 0$  be given. Set  $\delta = \epsilon$ . Hence, if  $\|\mathbf{x} - \mathbf{y}\| < \delta$ , then by Exercise 2.9,  $|\|\mathbf{x}\| - \|\mathbf{y}\|| \leq \|\mathbf{x} - \mathbf{y}\| < \delta = \epsilon$ .

## 3. Transformations

### 3.1

Let  $\mathbf{v}$  be the vector such that  $\mathbf{x}$  are the coordinates of  $\mathbf{v}$  with respect to  $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$ , and  $\mathbf{x}'$  are the coordinates of  $\mathbf{v}$  with respect to  $\{\mathbf{e}'_1, \mathbf{e}'_2, \dots, \mathbf{e}'_n\}$ . Then,

$$\mathbf{v} = x_1 \mathbf{e}_1 + \dots + x_n \mathbf{e}_n = [\mathbf{e}_1, \dots, \mathbf{e}_n] \mathbf{x},$$

and

$$\mathbf{v} = x'_1 \mathbf{e}'_1 + \dots + x'_n \mathbf{e}'_n = [\mathbf{e}'_1, \dots, \mathbf{e}'_n] \mathbf{x}'.$$

Hence,

$$[\mathbf{e}_1, \dots, \mathbf{e}_n] \mathbf{x} = [\mathbf{e}'_1, \dots, \mathbf{e}'_n] \mathbf{x}'$$

which implies

$$\mathbf{x}' = [\mathbf{e}'_1, \dots, \mathbf{e}'_n]^{-1} [\mathbf{e}_1, \dots, \mathbf{e}_n] \mathbf{x} = \mathbf{T} \mathbf{x}.$$

### 3.2

a. We have

$$[\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3] = [\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3] \begin{bmatrix} 1 & 2 & 4 \\ 3 & -1 & 5 \\ -4 & 5 & 3 \end{bmatrix}.$$

Therefore,

$$\mathbf{T} = [\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3]^{-1} [\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3] = \begin{bmatrix} 1 & 2 & 4 \\ 3 & -1 & 5 \\ -4 & 5 & 3 \end{bmatrix}^{-1} = \frac{1}{42} \begin{bmatrix} 28 & -14 & -14 \\ 29 & -19 & -7 \\ -11 & 13 & 7 \end{bmatrix}.$$

b. We have

$$[\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3] = [\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3] \begin{bmatrix} 1 & 2 & 3 \\ 1 & -1 & 0 \\ 3 & 4 & 5 \end{bmatrix}.$$

Therefore,

$$\mathbf{T} = \begin{bmatrix} 1 & 2 & 3 \\ 1 & -1 & 0 \\ 3 & 4 & 5 \end{bmatrix}.$$

### 3.3

We have

$$[\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3] = [\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3] \begin{bmatrix} 2 & 2 & 3 \\ 1 & -1 & 0 \\ -1 & 2 & 1 \end{bmatrix}.$$

Therefore, the transformation matrix from  $\{\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3\}$  to  $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}$  is

$$\mathbf{T} = \begin{bmatrix} 2 & 2 & 3 \\ 1 & -1 & 0 \\ -1 & 2 & 1 \end{bmatrix},$$

Now, consider a linear transformation  $L : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ , and let  $\mathbf{A}$  be its representation with respect to  $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}$ , and  $\mathbf{B}$  its representation with respect to  $\{\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3\}$ . Let  $\mathbf{y} = \mathbf{A}\mathbf{x}$  and  $\mathbf{y}' = \mathbf{B}\mathbf{x}'$ . Then,

$$\mathbf{y}' = \mathbf{T}\mathbf{y} = \mathbf{T}(\mathbf{A}\mathbf{x}) = \mathbf{T}\mathbf{A}(\mathbf{T}^{-1}\mathbf{x}') = (\mathbf{T}\mathbf{A}\mathbf{T}^{-1})\mathbf{x}'.$$

Hence, the representation of the linear transformation with respect to  $\{\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3\}$  is

$$\mathbf{B} = \mathbf{T}\mathbf{A}\mathbf{T}^{-1} = \begin{bmatrix} 3 & -10 & -8 \\ -1 & 8 & 4 \\ 2 & -13 & -7 \end{bmatrix}.$$

### 3.4

We have

$$[\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3, \mathbf{e}'_4] = [\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3, \mathbf{e}_4] \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

Therefore, the transformation matrix from  $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3, \mathbf{e}_4\}$  to  $\{\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3, \mathbf{e}'_4\}$  is

$$\mathbf{T} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix}^{-1} = \begin{bmatrix} 1 & -1 & 0 & 0 \\ 0 & 1 & -1 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

Now, consider a linear transformation  $L : \mathbb{R}^4 \rightarrow \mathbb{R}^4$ , and let  $\mathbf{A}$  be its representation with respect to  $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3, \mathbf{e}_4\}$ , and  $\mathbf{B}$  its representation with respect to  $\{\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3, \mathbf{e}'_4\}$ . Let  $\mathbf{y} = \mathbf{A}\mathbf{x}$  and  $\mathbf{y}' = \mathbf{B}\mathbf{x}'$ . Then,

$$\mathbf{y}' = \mathbf{T}\mathbf{y} = \mathbf{T}(\mathbf{A}\mathbf{x}) = \mathbf{T}\mathbf{A}(\mathbf{T}^{-1}\mathbf{x}') = (\mathbf{T}\mathbf{A}\mathbf{T}^{-1})\mathbf{x}'.$$

Therefore,

$$\mathbf{B} = \mathbf{T}\mathbf{A}\mathbf{T}^{-1} = \begin{bmatrix} 5 & 3 & 4 & 3 \\ -3 & -2 & -1 & -2 \\ -1 & 0 & -1 & -2 \\ 1 & 1 & 1 & 4 \end{bmatrix}.$$

### 3.5

Let  $\{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4\}$  be a set of linearly independent eigenvectors of  $\mathbf{A}$  corresponding to the eigenvalues  $\lambda_1, \lambda_2, \lambda_3$ , and  $\lambda_4$ . Let  $\mathbf{T} = [\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4]$ . Then,

$$\begin{aligned} \mathbf{A}\mathbf{T} &= \mathbf{A}[\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4] = [\mathbf{A}\mathbf{v}_1, \mathbf{A}\mathbf{v}_2, \mathbf{A}\mathbf{v}_3, \mathbf{A}\mathbf{v}_4] \\ &= [\lambda_1\mathbf{v}_1, \lambda_2\mathbf{v}_2, \lambda_3\mathbf{v}_3, \lambda_4\mathbf{v}_4] = [\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4] \begin{bmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ 0 & 0 & \lambda_3 & 0 \\ 0 & 0 & 0 & \lambda_4 \end{bmatrix}. \end{aligned}$$

Hence,

$$\mathbf{A}\mathbf{T} = \mathbf{T} \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix},$$



or

$$\mathbf{T}^{-1}\mathbf{A}\mathbf{T} = \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix}.$$

Therefore, the linear transformation has a diagonal matrix form with respect to the basis formed by a linearly independent set of eigenvectors.

Because

$$\det(\mathbf{A}) = (\lambda - 2)(\lambda - 3)(\lambda - 1)(\lambda + 1),$$

the eigenvalues are  $\lambda_1 = 2$ ,  $\lambda_2 = 3$ ,  $\lambda_3 = 1$ , and  $\lambda_4 = -1$ .

From  $\mathbf{A}\mathbf{v}_i = \lambda_i\mathbf{v}_i$ , where  $\mathbf{v}_i \neq 0$  ( $i = 1, 2, 3$ ), the corresponding eigenvectors are

$$\mathbf{v}_1 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}, \quad \mathbf{v}_3 = \begin{bmatrix} 0 \\ 2 \\ -9 \\ 1 \end{bmatrix}, \text{ and } \quad \mathbf{v}_4 = \begin{bmatrix} 24 \\ -12 \\ 1 \\ 9 \end{bmatrix}.$$

Therefore, the basis we are interested in is

$$\{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3\} = \left\{ \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 2 \\ -9 \\ 1 \end{bmatrix}, \begin{bmatrix} 24 \\ -12 \\ 1 \\ 9 \end{bmatrix} \right\}.$$

### 3.6

Suppose  $\mathbf{v}_1, \dots, \mathbf{v}_n$  are eigenvectors of  $\mathbf{A}$  corresponding to  $\lambda_1, \dots, \lambda_n$ , respectively. Then, for each  $i = 1, \dots, n$ , we have

$$(\mathbf{I}_n - \mathbf{A})\mathbf{v}_i = \mathbf{v}_i - \mathbf{A}\mathbf{v}_i = \mathbf{v}_i - \lambda_i\mathbf{v}_i = (1 - \lambda_i)\mathbf{v}_i$$

which shows that  $1 - \lambda_1, \dots, 1 - \lambda_n$  are the eigenvalues of  $\mathbf{I}_n - \mathbf{A}$ .

Alternatively, we may write the characteristic polynomial of  $\mathbf{I}_n - \mathbf{A}$  as

$$\pi_{\mathbf{I}_n - \mathbf{A}}(1 - \lambda) = \det((1 - \lambda)\mathbf{I}_n - (\mathbf{I}_n - \mathbf{A})) = \det(-[\lambda\mathbf{I}_n - \mathbf{A}]) = (-1)^n \pi_{\mathbf{A}}(\lambda),$$

which shows the desired result.

### 3.7

Let  $\mathbf{x}, \mathbf{y} \in \mathcal{V}^\perp$ , and  $\alpha, \beta \in \mathbb{R}$ . To show that  $\mathcal{V}^\perp$  is a subspace, we need to show that  $\alpha\mathbf{x} + \beta\mathbf{y} \in \mathcal{V}^\perp$ . For this, let  $\mathbf{v}$  be any vector in  $\mathcal{V}$ . Then,

$$\mathbf{v}^\top(\alpha\mathbf{x} + \beta\mathbf{y}) = \alpha\mathbf{v}^\top\mathbf{x} + \beta\mathbf{v}^\top\mathbf{y} = 0,$$

since  $\mathbf{v}^\top\mathbf{x} = \mathbf{v}^\top\mathbf{y} = 0$  by definition.

### 3.8

The null space of  $\mathbf{A}$  is  $\mathcal{N}(\mathbf{A}) = \{\mathbf{x} \in \mathbb{R}^3 : \mathbf{A}\mathbf{x} = \mathbf{0}\}$ . Using elementary row operations and back-substitution, we can solve the system of equations:

$$\begin{bmatrix} 4 & -2 & 0 \\ 2 & 1 & -1 \\ 2 & -3 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 4 & -2 & 0 \\ 0 & 2 & -1 \\ 0 & -2 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 4 & -2 & 0 \\ 0 & 2 & -1 \\ 0 & 0 & 0 \end{bmatrix} \Rightarrow \begin{array}{rcl} 4x_1 - 2x_2 & = & 0 \\ 2x_2 - x_3 & = & 0 \end{array}$$

$$\Rightarrow \quad x_2 = \frac{1}{2}x_3, \quad x_1 = \frac{1}{2}x_2 = \frac{1}{4}x_3 \quad \Rightarrow \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} \frac{1}{4} \\ \frac{1}{2} \\ 1 \end{bmatrix} x_3.$$

Therefore,

$$\mathcal{N}(\mathbf{A}) = \left\{ \begin{bmatrix} 1 \\ 2 \\ 4 \end{bmatrix} c : c \in \mathbb{R} \right\}.$$

### 3.9

Let  $\mathbf{x}, \mathbf{y} \in \mathcal{R}(\mathbf{A})$ , and  $\alpha, \beta \in \mathbb{R}$ . Then, there exists  $\mathbf{v}, \mathbf{u}$  such that  $\mathbf{x} = \mathbf{A}\mathbf{v}$  and  $\mathbf{y} = \mathbf{A}\mathbf{u}$ . Thus,

$$\alpha\mathbf{x} + \beta\mathbf{y} = \alpha\mathbf{A}\mathbf{v} + \beta\mathbf{A}\mathbf{u} = \mathbf{A}(\alpha\mathbf{v} + \beta\mathbf{u}).$$

Hence,  $\alpha\mathbf{x} + \beta\mathbf{y} \in \mathcal{R}(\mathbf{A})$ , which shows that  $\mathcal{R}(\mathbf{A})$  is a subspace.

Let  $\mathbf{x}, \mathbf{y} \in \mathcal{N}(\mathbf{A})$ , and  $\alpha, \beta \in \mathbb{R}$ . Then,  $\mathbf{A}\mathbf{x} = \mathbf{0}$  and  $\mathbf{A}\mathbf{y} = \mathbf{0}$ . Thus,

$$\mathbf{A}(\alpha\mathbf{x} + \beta\mathbf{y}) = \alpha\mathbf{A}\mathbf{x} + \beta\mathbf{A}\mathbf{y} = \mathbf{0}.$$

Hence,  $\alpha\mathbf{x} + \beta\mathbf{y} \in \mathcal{N}(\mathbf{A})$ , which shows that  $\mathcal{N}(\mathbf{A})$  is a subspace.

### 3.10

Let  $\mathbf{v} \in \mathcal{R}(\mathbf{B})$ , i.e.,  $\mathbf{v} = \mathbf{B}\mathbf{x}$  for some  $\mathbf{x}$ . Consider the matrix  $[\mathbf{A} \ \mathbf{v}]$ . Then,  $\mathcal{N}(\mathbf{A}^\top) = \mathcal{N}([\mathbf{A} \ \mathbf{v}]^\top)$ , since if  $\mathbf{u} \in \mathcal{N}(\mathbf{A}^\top)$ , then  $\mathbf{u} \in \mathcal{N}(\mathbf{B}^\top)$  by assumption, and hence  $\mathbf{u}^\top \mathbf{v} = \mathbf{u}^\top \mathbf{B}\mathbf{x} = \mathbf{x}^\top \mathbf{B}^\top \mathbf{u} = \mathbf{0}$ . Now,

$$\dim \mathcal{R}(\mathbf{A}) + \dim \mathcal{N}(\mathbf{A}^\top) = m$$

and

$$\dim \mathcal{R}([\mathbf{A} \ \mathbf{v}]) + \dim \mathcal{N}([\mathbf{A} \ \mathbf{v}]^\top) = m.$$

Since  $\dim \mathcal{N}(\mathbf{A}^\top) = \dim \mathcal{N}([\mathbf{A} \ \mathbf{v}]^\top)$ , then we have  $\dim \mathcal{R}(\mathbf{A}) = \dim \mathcal{R}([\mathbf{A} \ \mathbf{v}])$ . Hence,  $\mathbf{v}$  is a linear combination of the columns of  $\mathbf{A}$ , i.e.,  $\mathbf{v} \in \mathcal{R}(\mathbf{A})$ , which completes the proof.

### 3.11

We first show  $\mathbf{V} \subset (\mathbf{V}^\perp)^\perp$ . Let  $\mathbf{v} \in \mathbf{V}$ , and  $\mathbf{u}$  any element of  $\mathbf{V}^\perp$ . Then  $\mathbf{u}^\top \mathbf{v} = \mathbf{v}^\top \mathbf{u} = 0$ . Therefore,  $\mathbf{v} \in (\mathbf{V}^\perp)^\perp$ .

We now show  $(\mathbf{V}^\perp)^\perp \subset \mathbf{V}$ . Let  $\{\mathbf{a}_1, \dots, \mathbf{a}_k\}$  be a basis for  $\mathbf{V}$ , and  $\{\mathbf{b}_1, \dots, \mathbf{b}_l\}$  a basis for  $(\mathbf{V}^\perp)^\perp$ . Define  $\mathbf{A} = [\mathbf{a}_1 \cdots \mathbf{a}_k]$  and  $\mathbf{B} = [\mathbf{b}_1 \cdots \mathbf{b}_l]$ , so that  $\mathbf{V} = \mathcal{R}(\mathbf{A})$  and  $(\mathbf{V}^\perp)^\perp = \mathcal{R}(\mathbf{B})$ . Hence, it remains to show that  $\mathcal{R}(\mathbf{B}) \subset \mathcal{R}(\mathbf{A})$ . Using the result of Exercise 3.10, it suffices to show that  $\mathcal{N}(\mathbf{A}^\top) \subset \mathcal{N}(\mathbf{B}^\top)$ . So let  $\mathbf{x} \in \mathcal{N}(\mathbf{A}^\top)$ , which implies that  $\mathbf{x} \in \mathcal{R}(\mathbf{A})^\perp = \mathbf{V}^\perp$ , since  $\mathcal{R}(\mathbf{A})^\perp = \mathcal{N}(\mathbf{A}^\top)$ . Hence, for all  $\mathbf{y}$ , we have  $(\mathbf{B}\mathbf{y})^\top \mathbf{x} = \mathbf{0} = \mathbf{y}^\top \mathbf{B}^\top \mathbf{x}$ , which implies that  $\mathbf{B}^\top \mathbf{x} = \mathbf{0}$ . Therefore,  $\mathbf{x} \in \mathcal{N}(\mathbf{B}^\top)$ , which completes the proof.

### 3.12

Let  $\mathbf{w} \in \mathcal{W}^\perp$ , and  $\mathbf{y}$  be any element of  $\mathcal{V}$ . Since  $\mathcal{V} \subset \mathcal{W}$ , then  $\mathbf{y} \in \mathcal{W}$ . Therefore, by definition of  $\mathbf{w}$ , we have  $\mathbf{w}^\top \mathbf{y} = 0$ . Therefore,  $\mathbf{w} \in \mathcal{V}^\perp$ .

### 3.13

Let  $r = \dim \mathcal{V}$ . Let  $\mathbf{v}_1, \dots, \mathbf{v}_r$  be a basis for  $\mathcal{V}$ , and  $\mathbf{V}$  the matrix whose  $i$ th column is  $\mathbf{v}_i$ . Then, clearly  $\mathcal{V} = \mathcal{R}(\mathbf{V})$ .

Let  $\mathbf{u}_1, \dots, \mathbf{u}_{n-r}$  be a basis for  $\mathcal{V}^\perp$ , and  $\mathbf{U}$  the matrix whose  $i$ th row is  $\mathbf{u}_i^\top$ . Then,  $\mathcal{V}^\perp = \mathcal{R}(\mathbf{U}^\top)$ , and  $\mathcal{V} = (\mathcal{V}^\perp)^\perp = \mathcal{R}(\mathbf{U}^\top)^\perp = \mathcal{N}(\mathbf{U})$  (by Exercise 3.11 and Theorem 3.4).

### 3.14

a. Let  $\mathbf{x} \in \mathcal{V}$ . Then,  $\mathbf{x} = \mathbf{P}\mathbf{x} + (\mathbf{I} - \mathbf{P})\mathbf{x}$ . Note that  $\mathbf{P}\mathbf{x} \in \mathcal{V}$ , and  $(\mathbf{I} - \mathbf{P})\mathbf{x} \in \mathcal{V}^\perp$ . Therefore,  $\mathbf{x} = \mathbf{P}\mathbf{x} + (\mathbf{I} - \mathbf{P})\mathbf{x}$  is an orthogonal decomposition of  $\mathbf{x}$  with respect to  $\mathcal{V}$ . However,  $\mathbf{x} = \mathbf{x} + \mathbf{0}$  is also an orthogonal decomposition of  $\mathbf{x}$  with respect to  $\mathcal{V}$ . Since the orthogonal decomposition is unique, we must have  $\mathbf{x} = \mathbf{P}\mathbf{x}$ .

b. Suppose  $\mathbf{P}$  is an orthogonal projector onto  $\mathcal{V}$ . Clearly,  $\mathcal{R}(\mathbf{P}) \subset \mathcal{V}$  by definition. However, from part a,  $\mathbf{x} = \mathbf{P}\mathbf{x}$  for all  $\mathbf{x} \in \mathcal{V}$ , and hence  $\mathcal{V} \subset \mathcal{R}(\mathbf{P})$ . Therefore,  $\mathcal{R}(\mathbf{P}) = \mathcal{V}$ .

### 3.15

To answer the question, we have to represent the quadratic form with a symmetric matrix as

$$\mathbf{x}^\top \left( \frac{1}{2} \begin{bmatrix} 1 & -8 \\ 1 & 1 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -8 & 1 \end{bmatrix} \right) \mathbf{x} = \mathbf{x}^\top \begin{bmatrix} 1 & -7/2 \\ -7/2 & 1 \end{bmatrix} \mathbf{x}.$$

The leading principal minors are  $\Delta_1 = 1$  and  $\Delta_2 = -45/4$ . Therefore, the quadratic form is indefinite.

### 3.16

The leading principal minors are  $\Delta_1 = 2$ ,  $\Delta_2 = 0$ ,  $\Delta_3 = 0$ , which are all nonnegative. However, the eigenvalues of  $\mathbf{A}$  are  $0, -1.4641, 5.4641$  (for example, use Matlab to quickly check this). This implies that the matrix  $\mathbf{A}$  is indefinite (by Theorem 3.7). An alternative way to show that  $\mathbf{A}$  is not positive semidefinite is to find a vector  $\mathbf{x}$  such that  $\mathbf{x}^\top \mathbf{A} \mathbf{x} < 0$ . So, let  $\mathbf{x}$  be an eigenvector of  $\mathbf{A}$  corresponding to its negative eigenvalue  $\lambda = -1.4641$ . Then,  $\mathbf{x}^\top \mathbf{A} \mathbf{x} = \mathbf{x}^\top (\lambda \mathbf{x}) = \lambda \mathbf{x}^\top \mathbf{x} = \lambda \|\mathbf{x}\|^2 < 0$ . For this example, we can take  $\mathbf{x} = [0.3251, 0.3251, -0.8881]^\top$ , for which we can verify that  $\mathbf{x}^\top \mathbf{A} \mathbf{x} = -1.4643$ .

### 3.17

a. The matrix  $\mathbf{Q}$  is indefinite, since  $\Delta_2 = -1$  and  $\Delta_3 = 2$ .

b. Let  $\mathbf{x} \in \mathcal{M}$ . Then,  $x_2 + x_3 = -x_1$ ,  $x_1 + x_3 = -x_2$ , and  $x_1 + x_2 = -x_3$ . Therefore,

$$\mathbf{x}^\top \mathbf{Q} \mathbf{x} = x_1(x_2 + x_3) + x_2(x_1 + x_3) + x_3(x_1 + x_2) = -(x_1^2 + x_2^2 + x_3^2).$$

This implies that the matrix  $\mathbf{Q}$  is negative definite on the subspace  $\mathcal{M}$ .

### 3.18

a. We have

$$f(x_1, x_2, x_3) = x_2^2 = [x_1, x_2, x_3] \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}.$$

Then,

$$\mathbf{Q} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

and the eigenvalues of  $\mathbf{Q}$  are  $\lambda_1 = 0$ ,  $\lambda_2 = 1$ , and  $\lambda_3 = 0$ . Therefore, the quadratic form is positive semidefinite.

b. We have

$$f(x_1, x_2, x_3) = x_1^2 + 2x_2^2 - x_1x_3 = [x_1, x_2, x_3] \begin{bmatrix} 1 & 0 & -\frac{1}{2} \\ 0 & 2 & 0 \\ -\frac{1}{2} & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}.$$

Then,

$$\mathbf{Q} = \begin{bmatrix} 1 & 0 & -\frac{1}{2} \\ 0 & 2 & 0 \\ -\frac{1}{2} & 0 & 0 \end{bmatrix}$$

and the eigenvalues of  $\mathbf{Q}$  are  $\lambda_1 = 2$ ,  $\lambda_2 = (1 - \sqrt{2})/2$ , and  $\lambda_3 = (1 + \sqrt{2})/2$ . Therefore, the quadratic form is indefinite.

c. We have

$$f(x_1, x_2, x_3) = x_1^2 + x_3^2 + 2x_1x_2 + 2x_1x_3 + 2x_2x_3 = [x_1, x_2, x_3] \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}.$$

Then,

$$\mathbf{Q} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

and the eigenvalues of  $\mathbf{Q}$  are  $\lambda_1 = 0$ ,  $\lambda_2 = 1 - \sqrt{3}$ , and  $\lambda_3 = 1 + \sqrt{3}$ . Therefore, the quadratic form is indefinite.

### 3.19

We have

$$\begin{aligned} f(x_1, x_2, x_3) &= 4x_1^2 + x_2^2 + 9x_3^2 - 4x_1x_2 - 6x_2x_3 + 12x_1x_3 \\ &= [x_1, x_2, x_3] \begin{bmatrix} 4 & -2 & 6 \\ -2 & 1 & -3 \\ 6 & -3 & 9 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}. \end{aligned}$$

Let

$$\mathbf{Q} = \begin{bmatrix} 4 & -2 & 6 \\ -2 & 1 & -3 \\ 6 & -3 & 9 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = x_1\mathbf{e}_1 + x_2\mathbf{e}_2 + x_3\mathbf{e}_3,$$

where  $\mathbf{e}_1$ ,  $\mathbf{e}_2$ , and  $\mathbf{e}_3$  form the natural basis for  $\mathbb{R}^3$ .

Let  $\mathbf{v}_1$ ,  $\mathbf{v}_2$ , and  $\mathbf{v}_3$  be another basis for  $\mathbb{R}^3$ . Then, the vector  $\mathbf{x}$  is represented in the new basis as  $\tilde{\mathbf{x}}$ , where  $\mathbf{x} = [\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3]\tilde{\mathbf{x}} = \mathbf{V}\tilde{\mathbf{x}}$ .

Now,  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q} \mathbf{x} = (\mathbf{V}\tilde{\mathbf{x}})^\top \mathbf{Q} (\mathbf{V}\tilde{\mathbf{x}}) = \tilde{\mathbf{x}}^\top (\mathbf{V}^\top \mathbf{Q} \mathbf{V}) \tilde{\mathbf{x}} = \tilde{\mathbf{x}}^\top \tilde{\mathbf{Q}} \tilde{\mathbf{x}}$ , where

$$\tilde{\mathbf{Q}} = \begin{bmatrix} \tilde{q}_{11} & \tilde{q}_{12} & \tilde{q}_{13} \\ \tilde{q}_{21} & \tilde{q}_{22} & \tilde{q}_{23} \\ \tilde{q}_{31} & \tilde{q}_{32} & \tilde{q}_{33} \end{bmatrix}$$

and  $\tilde{q}_{ij} = \mathbf{v}_i \mathbf{Q} \mathbf{v}_j$  for  $i, j = 1, 2, 3$ .

We will find a basis  $\{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3\}$  such that  $\tilde{q}_{ij} = 0$  for  $i \neq j$ , and is of the form

$$\begin{aligned} \mathbf{v}_1 &= \alpha_{11}\mathbf{e}_1 \\ \mathbf{v}_2 &= \alpha_{21}\mathbf{e}_1 + \alpha_{22}\mathbf{e}_2 \\ \mathbf{v}_3 &= \alpha_{31}\mathbf{e}_1 + \alpha_{32}\mathbf{e}_2 + \alpha_{33}\mathbf{e}_3 \end{aligned}$$

Because

$$\tilde{q}_{ij} = \mathbf{v}_i \mathbf{Q} \mathbf{v}_j = \mathbf{v}_i \mathbf{Q} (\alpha_{j1}\mathbf{e}_1 + \dots + \alpha_{jj}\mathbf{e}_j) = \alpha_{j1}(\mathbf{v}_i \mathbf{Q} \mathbf{e}_1) + \dots + \alpha_{jj}(\mathbf{v}_i \mathbf{Q} \mathbf{e}_j),$$

we deduce that if  $\mathbf{v}_i \mathbf{Q} \mathbf{e}_j = 0$  for  $j < i$ , then  $\mathbf{v}_i \mathbf{Q} \mathbf{v}_j = 0$ . In this case,

$$\tilde{q}_{ii} = \mathbf{v}_i \mathbf{Q} \mathbf{v}_i = \mathbf{v}_i \mathbf{Q} (\alpha_{i1}\mathbf{e}_1 + \dots + \alpha_{ii}\mathbf{e}_i) = \alpha_{i1}(\mathbf{v}_i \mathbf{Q} \mathbf{e}_1) + \dots + \alpha_{ii}(\mathbf{v}_i \mathbf{Q} \mathbf{e}_i) = \alpha_{ii}(\mathbf{v}_i \mathbf{Q} \mathbf{e}_i).$$

Our task therefore is to find  $\mathbf{v}_i$  ( $i = 1, 2, 3$ ) such that

$$\begin{aligned} \mathbf{v}_i \mathbf{Q} \mathbf{e}_j &= 0, & j < i \\ \mathbf{v}_i \mathbf{Q} \mathbf{e}_i &= 1, \end{aligned}$$

and, in this case, we get

$$\tilde{\mathbf{Q}} = \begin{bmatrix} \alpha_{11} & 0 & 0 \\ 0 & \alpha_{22} & 0 \\ 0 & 0 & \alpha_{33} \end{bmatrix}.$$

▪ Case of  $i = 1$ .

From  $\mathbf{v}_1^\top \mathbf{Q} \mathbf{e}_1 = 1$ ,

$$(\alpha_{11}\mathbf{e}_1)^\top \mathbf{Q} \mathbf{e}_1 = \alpha_{11}(\mathbf{e}_1^\top \mathbf{Q} \mathbf{e}_1) = \alpha_{11}q_{11} = 1.$$

Therefore,

$$\alpha_{11} = \frac{1}{q_{11}} = \frac{1}{\Delta_1} = \frac{1}{4} \quad \Rightarrow \quad \mathbf{v}_1 = \alpha_{11}\mathbf{e}_1 = \begin{bmatrix} \frac{1}{4} \\ 0 \\ 0 \\ 0 \end{bmatrix}.$$

- Case of  $i = 2$ .

From  $\mathbf{v}_2^\top \mathbf{Q} \mathbf{e}_1 = 0$ ,

$$(\alpha_{21} \mathbf{e}_1 + \alpha_{22} \mathbf{e}_2)^\top \mathbf{Q} \mathbf{e}_1 = \alpha_{21} (\mathbf{e}_1^\top \mathbf{Q} \mathbf{e}_1) + \alpha_{22} (\mathbf{e}_2^\top \mathbf{Q} \mathbf{e}_1) = \alpha_{21} q_{11} + \alpha_{22} q_{21} = 0.$$

From  $\mathbf{v}_2^\top \mathbf{Q} \mathbf{e}_2 = 1$ ,

$$(\alpha_{21} \mathbf{e}_1 + \alpha_{22} \mathbf{e}_2)^\top \mathbf{Q} \mathbf{e}_2 = \alpha_{21} (\mathbf{e}_1^\top \mathbf{Q} \mathbf{e}_2) + \alpha_{22} (\mathbf{e}_2^\top \mathbf{Q} \mathbf{e}_2) = \alpha_{21} q_{12} + \alpha_{22} q_{22} = 1.$$

Therefore,

$$\begin{bmatrix} q_{11} & q_{21} \\ q_{12} & q_{22} \end{bmatrix} \begin{bmatrix} \alpha_{21} \\ \alpha_{22} \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}.$$

But, since  $\Delta_2 = 0$ , this system of equations is inconsistent. Hence, in this problem  $\mathbf{v}_2^\top \mathbf{Q} \mathbf{e}_2 = 0$  should be satisfied instead of  $\mathbf{v}_2^\top \mathbf{Q} \mathbf{e}_2 = 1$  so that the system can have a solution. In this case, the diagonal matrix becomes

$$\tilde{\mathbf{Q}} = \begin{bmatrix} \alpha_{11} & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & \alpha_{33} \end{bmatrix},$$

and the system of equations become

$$\begin{bmatrix} q_{11} & q_{21} \\ q_{12} & q_{22} \end{bmatrix} \begin{bmatrix} \alpha_{21} \\ \alpha_{22} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow \begin{bmatrix} \alpha_{21} \\ \alpha_{22} \end{bmatrix} = \begin{bmatrix} \frac{1}{2} \\ 1 \end{bmatrix} \alpha_{22},$$

where  $\alpha_{22}$  is an arbitrary real number. Thus,

$$\mathbf{v}_2 = \alpha_{21} \mathbf{e}_1 + \alpha_{22} \mathbf{e}_2 = \begin{bmatrix} \frac{1}{2} \\ 1 \\ 0 \end{bmatrix} a,$$

where  $a$  is an arbitrary real number.

- Case of  $i = 3$ .

Since in this case  $\Delta_3 = \det(\mathbf{Q}) = 0$ , we will have to apply the same reasoning of the previous case and use the condition  $\mathbf{v}_3^\top \mathbf{Q} \mathbf{e}_3 = 0$  instead of  $\mathbf{v}_3^\top \mathbf{Q} \mathbf{e}_3 = 1$ . In this way the diagonal matrix becomes

$$\tilde{\mathbf{Q}} = \begin{bmatrix} \alpha_{11} & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

Thus, from  $\mathbf{v}_3^\top \mathbf{Q} \mathbf{e}_1 = 0$ ,  $\mathbf{v}_3^\top \mathbf{Q} \mathbf{e}_2 = 0$  and  $\mathbf{v}_3^\top \mathbf{Q} \mathbf{e}_3 = 0$ ,

$$\begin{aligned} \begin{bmatrix} q_{11} & q_{21} & q_{31} \\ q_{12} & q_{22} & q_{32} \\ q_{13} & q_{23} & q_{33} \end{bmatrix} \begin{bmatrix} \alpha_{31} \\ \alpha_{32} \\ \alpha_{33} \end{bmatrix} &= \mathbf{Q}^\top \begin{bmatrix} \alpha_{31} \\ \alpha_{32} \\ \alpha_{33} \end{bmatrix} = \mathbf{Q} \begin{bmatrix} \alpha_{31} \\ \alpha_{32} \\ \alpha_{33} \end{bmatrix} \\ &= \begin{bmatrix} 4 & -2 & 6 \\ -2 & 1 & -3 \\ 6 & -3 & 9 \end{bmatrix} \begin{bmatrix} \alpha_{31} \\ \alpha_{32} \\ \alpha_{33} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}. \end{aligned}$$

Therefore,

$$\begin{bmatrix} \alpha_{31} \\ \alpha_{32} \\ \alpha_{33} \end{bmatrix} = \begin{bmatrix} \alpha_{31} \\ 2\alpha_{31} + 3\alpha_{33} \\ \alpha_{33} \end{bmatrix},$$

where  $\alpha_{31}$  and  $\alpha_{33}$  are arbitrary real numbers. Thus,

$$\mathbf{v}_3 = \alpha_{31}\mathbf{e}_1 + \alpha_{32}\mathbf{e}_2 + \alpha_{33}\mathbf{e}_3 = \begin{bmatrix} b \\ 2b + 3c \\ c \end{bmatrix},$$

where  $b$  and  $c$  are arbitrary real numbers.

Finally,

$$\mathbf{V} = [\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3] = \begin{bmatrix} \frac{1}{4} & \frac{a}{2} & b \\ 0 & a & 2b + 3c \\ 0 & 0 & c \end{bmatrix},$$

where  $a$ ,  $b$ , and  $c$  are arbitrary real numbers.

### 3.20

We represent this quadratic form as  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q} \mathbf{x}$ , where

$$\mathbf{Q} = \begin{bmatrix} 1 & \xi & -1 \\ \xi & 1 & 2 \\ -1 & 2 & 5 \end{bmatrix}.$$

The leading principal minors of  $\mathbf{Q}$  are  $\Delta_1 = 1$ ,  $\Delta_2 = 1 - \xi^2$ ,  $\Delta_3 = -5\xi^2 - 4\xi$ . For the quadratic form to be positive definite, all the leading principal minors of  $\mathbf{Q}$  must be positive. This is the case if and only if  $\xi \in (-4/5, 0)$ .

### 3.21

The matrix  $\mathbf{Q} = \mathbf{Q}^\top > 0$  can be represented as  $\mathbf{Q} = \mathbf{Q}^{1/2} \mathbf{Q}^{1/2}$ , where  $\mathbf{Q}^{1/2} = (\mathbf{Q}^{1/2})^\top > 0$ .

1. Now,  $\langle \mathbf{x}, \mathbf{x} \rangle_Q = (\mathbf{Q}^{1/2} \mathbf{x})^\top (\mathbf{Q}^{1/2} \mathbf{x}) = \|\mathbf{Q}^{1/2} \mathbf{x}\|^2 \geq 0$ , and

$$\begin{aligned} \langle \mathbf{x}, \mathbf{x} \rangle_Q = 0 &\Leftrightarrow \|\mathbf{Q}^{1/2} \mathbf{x}\|^2 = 0 \\ &\Leftrightarrow \mathbf{Q}^{1/2} \mathbf{x} = \mathbf{0} \\ &\Leftrightarrow \mathbf{x} = \mathbf{0} \end{aligned}$$

since  $\mathbf{Q}^{1/2}$  is nonsingular.

2.  $\langle \mathbf{x}, \mathbf{y} \rangle_Q = \mathbf{x}^\top \mathbf{Q} \mathbf{y} = \mathbf{y}^\top \mathbf{Q}^\top \mathbf{x} = \mathbf{y}^\top \mathbf{Q} \mathbf{x} = \langle \mathbf{y}, \mathbf{x} \rangle_Q$ .

3. We have

$$\begin{aligned} \langle \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle_Q &= (\mathbf{x} + \mathbf{y})^\top \mathbf{Q} \mathbf{z} \\ &= \mathbf{x}^\top \mathbf{Q} \mathbf{z} + \mathbf{y}^\top \mathbf{Q} \mathbf{z} \\ &= \langle \mathbf{x}, \mathbf{z} \rangle_Q + \langle \mathbf{y}, \mathbf{z} \rangle_Q. \end{aligned}$$

4.  $\langle r\mathbf{x}, \mathbf{y} \rangle_Q = (r\mathbf{x})^\top \mathbf{Q} \mathbf{y} = r\mathbf{x}^\top \mathbf{Q} \mathbf{y} = r\langle \mathbf{x}, \mathbf{y} \rangle_Q$ .

### 3.22

We have

$$\|\mathbf{A}\|_\infty = \max\{\|\mathbf{A}\mathbf{x}\|_\infty : \|\mathbf{x}\|_\infty = 1\}.$$

We first show that  $\|\mathbf{A}\|_\infty \leq \max_i \sum_{k=1}^n |a_{ik}|$ . For this, note that for each  $\mathbf{x}$  such that  $\|\mathbf{x}\|_\infty = 1$ , we have

$$\begin{aligned} \|\mathbf{A}\mathbf{x}\|_\infty &= \max_i \left| \sum_{k=1}^n a_{ik} x_k \right| \\ &\leq \max_i \sum_{k=1}^n |a_{ik}| |x_k| \\ &\leq \max_i \sum_{k=1}^n |a_{ik}|, \end{aligned}$$

since  $|x_k| \leq \max_k |x_k| = \|\mathbf{x}\|_\infty = 1$ . Therefore,

$$\|\mathbf{A}\|_\infty \leq \max_i \sum_{k=1}^n |a_{ik}|.$$

To show that  $\|\mathbf{A}\|_\infty = \max_i \sum_{k=1}^n |a_{ik}|$ , it remains to find a  $\tilde{\mathbf{x}} \in \mathbb{R}^n$ ,  $\|\tilde{\mathbf{x}}\|_\infty = 1$ , such that  $\|\mathbf{A}\tilde{\mathbf{x}}\|_\infty = \max_i \sum_{k=1}^n |a_{ik}|$ . So, let  $j$  be such that

$$\sum_{k=1}^n |a_{jk}| = \max_i \sum_{k=1}^n |a_{ik}|.$$

Define  $\tilde{\mathbf{x}}$  by

$$\tilde{x}_k = \begin{cases} |a_{jk}|/a_{jk} & \text{if } a_{jk} \neq 0 \\ 1 & \text{otherwise} \end{cases}.$$

Clearly  $\|\tilde{\mathbf{x}}\|_\infty = 1$ . Furthermore, for  $i \neq j$ ,

$$\left| \sum_{k=1}^n a_{ik} \tilde{x}_k \right| \leq \sum_{k=1}^n |a_{ik}| \leq \max_i \sum_{k=1}^n |a_{ik}| = \sum_{k=1}^n |a_{jk}|$$

and

$$\left| \sum_{k=1}^n a_{jk} \tilde{x}_k \right| = \sum_{k=1}^n |a_{jk}|.$$

Therefore,

$$\|\mathbf{A}\tilde{\mathbf{x}}\|_\infty = \max_i \left| \sum_{k=1}^n a_{ik} \tilde{x}_k \right| = \sum_{k=1}^n |a_{jk}| = \max_i \sum_{k=1}^n |a_{ik}|.$$

### 3.23

We have

$$\|\mathbf{A}\|_1 = \max\{\|\mathbf{A}\mathbf{x}\|_1 : \|\mathbf{x}\|_1 = 1\}.$$

We first show that  $\|\mathbf{A}\|_1 \leq \max_k \sum_{i=1}^m |a_{ik}|$ . For this, note that for each  $\mathbf{x}$  such that  $\|\mathbf{x}\|_1 = 1$ , we have

$$\begin{aligned} \|\mathbf{A}\mathbf{x}\|_1 &= \sum_{i=1}^m \left| \sum_{k=1}^n a_{ik} x_k \right| \\ &\leq \sum_{i=1}^m \sum_{k=1}^n |a_{ik}| |x_k| \\ &\leq \sum_{k=1}^n |x_k| \sum_{i=1}^m |a_{ik}| \\ &\leq \left( \max_k \sum_{i=1}^m |a_{ik}| \right) \sum_{k=1}^n |x_k| \\ &\leq \max_k \sum_{i=1}^m |a_{ik}|, \end{aligned}$$

since  $\sum_{k=1}^n |x_k| = \|\mathbf{x}\|_1 = 1$ . Therefore,

$$\|\mathbf{A}\|_1 \leq \max_k \sum_{i=1}^m |a_{ik}|.$$

To show that  $\|\mathbf{A}\|_1 = \max_k \sum_{i=1}^m |a_{ik}|$ , it remains to find a  $\tilde{\mathbf{x}} \in \mathbb{R}^m$ ,  $\|\tilde{\mathbf{x}}\|_1 = 1$ , such that  $\|\mathbf{A}\tilde{\mathbf{x}}\|_1 = \max_k \sum_{i=1}^m |a_{ik}|$ . So, let  $j$  be such that

$$\sum_{i=1}^m |a_{ij}| = \max_k \sum_{i=1}^m |a_{ik}|.$$

Define  $\tilde{\mathbf{x}}$  by

$$\tilde{x}_k = \begin{cases} 1 & \text{if } k = j \\ 0 & \text{otherwise} \end{cases}.$$

Clearly  $\|\tilde{\mathbf{x}}\|_1 = 1$ . Furthermore,

$$\|\mathbf{A}\tilde{\mathbf{x}}\|_1 = \sum_{i=1}^m \left| \sum_{k=1}^n a_{ik} \tilde{x}_k \right| = \sum_{i=1}^m |a_{ij}| = \max_k \sum_{i=1}^m |a_{ik}|.$$

## 4. Concepts from Geometry

### 4.1

$\Rightarrow$ : Let  $S = \{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}\}$  be a linear variety. Let  $\mathbf{x}, \mathbf{y} \in S$  and  $\alpha \in \mathbb{R}$ . Then,

$$\mathbf{A}(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) = \alpha\mathbf{A}\mathbf{x} + (1 - \alpha)\mathbf{A}\mathbf{y} = \alpha\mathbf{b} + (1 - \alpha)\mathbf{b} = \mathbf{b}.$$

Therefore,  $\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \in S$ .

$\Leftarrow$ : If  $S$  is empty, we are done. So, suppose  $\mathbf{x}_0 \in S$ . Consider the set  $S_0 = S - \mathbf{x}_0 = \{\mathbf{x} - \mathbf{x}_0 : \mathbf{x} \in S\}$ . Clearly, for all  $\mathbf{x}, \mathbf{y} \in S_0$  and  $\alpha \in \mathbb{R}$ , we have  $\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \in S_0$ . Note that  $\mathbf{0} \in S_0$ . We claim that  $S_0$  is a subspace. To see this, let  $\mathbf{x}, \mathbf{y} \in S_0$ , and  $\alpha \in \mathbb{R}$ . Then,  $\alpha\mathbf{x} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{0} \in S_0$ . Furthermore,  $\frac{1}{2}\mathbf{x} + \frac{1}{2}\mathbf{y} \in S_0$ , and therefore  $\mathbf{x} + \mathbf{y} \in S_0$  by the previous argument. Hence,  $S_0$  is a subspace. Therefore, by Exercise 3.13, there exists  $\mathbf{A}$  such that  $S_0 = \mathcal{N}(\mathbf{A}) = \{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{0}\}$ . Define  $\mathbf{b} = \mathbf{A}\mathbf{x}_0$ . Then,

$$\begin{aligned} S &= S_0 + \mathbf{x}_0 = \{\mathbf{y} + \mathbf{x}_0 : \mathbf{y} \in \mathcal{N}(\mathbf{A})\} \\ &= \{\mathbf{y} + \mathbf{x}_0 : \mathbf{A}\mathbf{y} = \mathbf{0}\} \\ &= \{\mathbf{y} + \mathbf{x}_0 : \mathbf{A}(\mathbf{y} + \mathbf{x}_0) = \mathbf{b}\} \\ &= \{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}\}. \end{aligned}$$

### 4.2

Let  $\mathbf{u}, \mathbf{v} \in \Theta = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| \leq r\}$ , and  $\alpha \in [0, 1]$ . Suppose  $\mathbf{z} = \alpha\mathbf{u} + (1 - \alpha)\mathbf{v}$ . To show that  $\Theta$  is convex, we need to show that  $\mathbf{z} \in \Theta$ , i.e.,  $\|\mathbf{z}\| \leq r$ . To this end,

$$\begin{aligned} \|\mathbf{z}\|^2 &= (\alpha\mathbf{u}^\top + (1 - \alpha)\mathbf{v}^\top)(\alpha\mathbf{u} + (1 - \alpha)\mathbf{v}) \\ &= \alpha^2\|\mathbf{u}\|^2 + 2\alpha(1 - \alpha)\mathbf{u}^\top\mathbf{v} + (1 - \alpha)^2\|\mathbf{v}\|^2. \end{aligned}$$

Since  $\mathbf{u}, \mathbf{v} \in \Theta$ , then  $\|\mathbf{u}\|^2 \leq r^2$  and  $\|\mathbf{v}\|^2 \leq r^2$ . Furthermore, by the Cauchy-Schwarz Inequality, we have  $\mathbf{u}^\top\mathbf{v} \leq \|\mathbf{u}\|\|\mathbf{v}\| \leq r^2$ . Therefore,

$$\|\mathbf{z}\|^2 \leq \alpha^2 r^2 + 2\alpha(1 - \alpha)r^2 + (1 - \alpha)^2 r^2 = r^2.$$

Hence,  $\mathbf{z} \in \Theta$ , which implies that  $\Theta$  is a convex set, i.e., the any point on the line segment joining  $\mathbf{u}$  and  $\mathbf{v}$  is also in  $\Theta$ .

### 4.3

Let  $\mathbf{u}, \mathbf{v} \in \Theta = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{b}\}$ , and  $\alpha \in [0, 1]$ . Suppose  $\mathbf{z} = \alpha\mathbf{u} + (1 - \alpha)\mathbf{v}$ . To show that  $\Theta$  is convex, we need to show that  $\mathbf{z} \in \Theta$ , i.e.,  $\mathbf{A}\mathbf{z} = \mathbf{b}$ . To this end,

$$\begin{aligned} \mathbf{A}\mathbf{z} &= \mathbf{A}(\alpha\mathbf{u} + (1 - \alpha)\mathbf{v}) \\ &= \alpha\mathbf{A}\mathbf{u} + (1 - \alpha)\mathbf{A}\mathbf{v}. \end{aligned}$$



Since  $\mathbf{u}, \mathbf{v} \in \Theta$ , then  $\mathbf{A}\mathbf{u} = \mathbf{b}$  and  $\mathbf{A}\mathbf{v} = \mathbf{b}$ . Therefore,

$$\mathbf{A}\mathbf{z} = \alpha\mathbf{b} + (1 - \alpha)\mathbf{b} = \mathbf{b},$$

and hence  $\mathbf{z} \in \Theta$ .

#### 4.4

Let  $\mathbf{u}, \mathbf{v} \in \Theta = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} \geq \mathbf{0}\}$ , and  $\alpha \in [0, 1]$ . Suppose  $\mathbf{z} = \alpha\mathbf{u} + (1 - \alpha)\mathbf{v}$ . To show that  $\Theta$  is convex, we need to show that  $\mathbf{z} \in \Theta$ , i.e.,  $\mathbf{z} \geq \mathbf{0}$ . To this end, write  $\mathbf{x} = [x_1, \dots, x_n]^\top$ ,  $\mathbf{y} = [y_1, \dots, y_n]^\top$ , and  $\mathbf{z} = [z_1, \dots, z_n]^\top$ . Then,  $z_i = \alpha x_i + (1 - \alpha)y_i$ ,  $i = 1, \dots, n$ . Since  $x_i, y_i \geq 0$ , and  $\alpha, 1 - \alpha \geq 0$ , we have  $z_i \geq 0$ . Therefore,  $\mathbf{z} \geq \mathbf{0}$ , and hence  $\mathbf{z} \in \Theta$ .

## 5. Elements of Calculus

### 5.1

Observe that

$$\|\mathbf{A}^k\| \leq \|\mathbf{A}^{k-1}\| \|\mathbf{A}\| \leq \|\mathbf{A}^{k-2}\| \|\mathbf{A}\|^2 \leq \dots \leq \|\mathbf{A}\|^k.$$

Therefore, if  $\|\mathbf{A}\| < 1$ , then  $\lim_{k \rightarrow \infty} \|\mathbf{A}^k\| = \mathbf{0}$  which implies that  $\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{0}$ .

### 5.2

For the case when  $\mathbf{A}$  has all real eigenvalues, the proof is simple. Let  $\lambda$  be the eigenvalue of  $\mathbf{A}$  with largest absolute value, and  $\mathbf{x}$  the corresponding (normalized) eigenvector, i.e.,  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  and  $\|\mathbf{x}\| = 1$ . Then,

$$\|\mathbf{A}\| \geq \|\mathbf{A}\mathbf{x}\| = \|\lambda\mathbf{x}\| = |\lambda|\|\mathbf{x}\| = |\lambda|,$$

which completes the proof for this case.

In general, the eigenvalues of  $\mathbf{A}$  and the corresponding eigenvectors may be complex. In this case, we proceed as follows (see [41]). Consider the matrix

$$\mathbf{B} = \frac{\mathbf{A}}{\|\mathbf{A}\| + \varepsilon},$$

where  $\varepsilon$  is a positive real number. We have

$$\|\mathbf{B}\| = \frac{\|\mathbf{A}\|}{\|\mathbf{A}\| + \varepsilon} < 1.$$

By Exercise 5.1,  $\mathbf{B}^k \rightarrow \mathbf{0}$  as  $k \rightarrow \infty$ , and thus by Lemma 5.1,  $|\lambda_i(\mathbf{B})| < 1$ ,  $i = 1, \dots, n$ . On the other hand, for each  $i = 1, \dots, n$ ,

$$\lambda_i(\mathbf{B}) = \frac{\lambda_i(\mathbf{A})}{\|\mathbf{A}\| + \varepsilon},$$

and thus

$$|\lambda_i(\mathbf{B})| = \frac{|\lambda_i(\mathbf{A})|}{\|\mathbf{A}\| + \varepsilon} < 1.$$

which gives

$$|\lambda_i(\mathbf{A})| < \|\mathbf{A}\| + \varepsilon.$$

Since the above arguments hold for any  $\varepsilon > 0$ , we have  $|\lambda_i(\mathbf{A})| \leq \|\mathbf{A}\|$ .

### 5.3

a.  $\nabla f(\mathbf{x}) = (\mathbf{a}\mathbf{b}^\top + \mathbf{b}\mathbf{a}^\top)\mathbf{x}.$

b.  $\mathbf{F}(\mathbf{x}) = \mathbf{a}\mathbf{b}^\top + \mathbf{b}\mathbf{a}^\top.$

## 5.4

---

We have

$$Df(\mathbf{x}) = [x_1/3, x_2/2],$$

and

$$\frac{d\mathbf{g}}{dt}(t) = \begin{bmatrix} 3 \\ 2 \end{bmatrix}.$$

By the chain rule,

$$\begin{aligned} \frac{d}{dt}F(t) &= Df(\mathbf{g}(t))\frac{d\mathbf{g}}{dt}(t) \\ &= [(3t+5)/3, (2t-6)/2] \begin{bmatrix} 3 \\ 2 \end{bmatrix} \\ &= 5t - 1. \end{aligned}$$

## 5.5

---

We have

$$Df(\mathbf{x}) = [x_2/2, x_1/2],$$

and

$$\frac{\partial \mathbf{g}}{\partial s}(s, t) = \begin{bmatrix} 4 \\ 2 \end{bmatrix}, \quad \frac{\partial \mathbf{g}}{\partial t}(s, t) = \begin{bmatrix} 3 \\ 1 \end{bmatrix}.$$

By the chain rule,

$$\begin{aligned} \frac{\partial}{\partial s}f(\mathbf{g}(s, t)) &= Df(\mathbf{g}(s, t))\frac{\partial \mathbf{g}}{\partial s}(s, t) \\ &= \frac{1}{2}[2s+t, 4s+3t] \begin{bmatrix} 4 \\ 2 \end{bmatrix} \\ &= 8s + 5t, \end{aligned}$$

and

$$\begin{aligned} \frac{\partial}{\partial t}f(\mathbf{g}(s, t)) &= Df(\mathbf{g}(s, t))\frac{\partial \mathbf{g}}{\partial t}(s, t) \\ &= \frac{1}{2}[2s+t, 4s+3t] \begin{bmatrix} 3 \\ 1 \end{bmatrix} \\ &= 5s + 3t. \end{aligned}$$

## 5.6

---

We have

$$Df(\mathbf{x}) = [3x_1^2x_2x_3^2 + x_2, x_1^3x_3^2 + x_1, 2x_1^3x_2x_3 + 1]$$

and

$$\frac{d\mathbf{x}}{dt}(t) = \begin{bmatrix} e^t + 3t^2 \\ 2t \\ 1 \end{bmatrix}.$$

By the chain rule,

$$\begin{aligned}
& \frac{d}{dt}f(\mathbf{x}(t)) \\
&= Df(\mathbf{x}(t))\frac{d\mathbf{x}}{dt}(t) \\
&= [3x_1(t)^2x_2(t)x_3(t)^2 + x_2(t), x_1(t)^3x_3(t)^2 + x_1(t), 2x_1(t)^3x_2(t)x_3(t) + 1] \begin{bmatrix} e^t + 3t^2 \\ 2t \\ 1 \end{bmatrix} \\
&= 12t(e^t + 3t^2)^3 + 2te^t + 6t^2 + 2t + 1.
\end{aligned}$$

### 5.7

Let  $\varepsilon > 0$  be given. Since  $\mathbf{f}(\mathbf{x}) = o(g(\mathbf{x}))$ , then

$$\lim_{\mathbf{x} \rightarrow \mathbf{0}} \frac{\|\mathbf{f}(\mathbf{x})\|}{g(\mathbf{x})} = 0.$$

Hence, there exists  $\delta > 0$  such that if  $\|\mathbf{x}\| < \delta$ , then

$$\frac{\|\mathbf{f}(\mathbf{x})\|}{g(\mathbf{x})} < \varepsilon,$$

which can be rewritten as

$$\|\mathbf{f}(\mathbf{x})\| \leq \varepsilon g(\mathbf{x}).$$

### 5.8

By Exercise 5.7, there exists  $\delta > 0$  such that if  $\|\mathbf{x}\| < \delta$ , then  $|o(g(\mathbf{x}))| < g(\mathbf{x})/2$ . Hence, if  $\|\mathbf{x}\| < \delta$ ,  $\mathbf{x} \neq \mathbf{0}$ , then

$$f(\mathbf{x}) \leq -g(\mathbf{x}) + |o(g(\mathbf{x}))| < -g(\mathbf{x}) + g(\mathbf{x})/2 = -\frac{1}{2}g(\mathbf{x}) < 0.$$

### 5.9

We have that

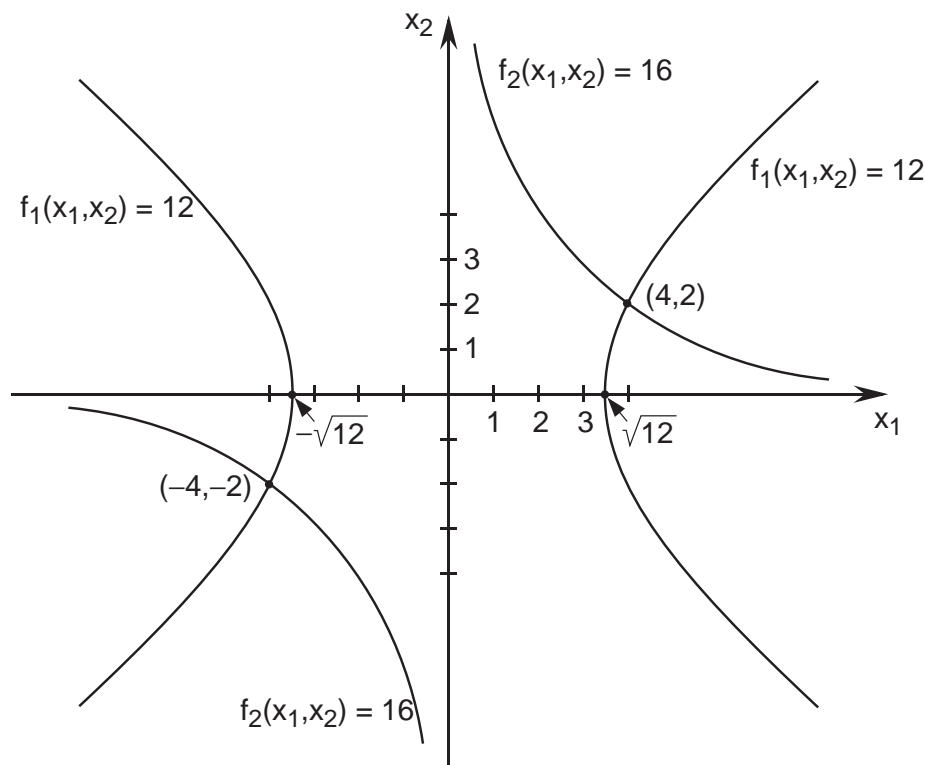
$$\{\mathbf{x} : f_1(\mathbf{x}) = 12\} = \{\mathbf{x} : x_1^2 - x_2^2 = 12\},$$

and

$$\{\mathbf{x} : f_2(\mathbf{x}) = 16\} = \{\mathbf{x} : x_2 = 8/x_1\}.$$

To find the intersection points, we substitute  $x_2 = 8/x_1$  into  $x_1^2 - x_2^2 = 12$  to get  $x_1^4 - 12x_1^2 - 64 = 0$ . Solving gives  $x_1^2 = 16, -4$ . Clearly, the only two possibilities for  $x_1$  are  $x_1 = +4, -4$ , from which we obtain  $x_2 = +2, -2$ . Hence, the intersection points are located at  $[4, 2]^\top$  and  $[-4, -2]^\top$ .

The level sets associated with  $f_1(x_1, x_2) = 12$  and  $f_2(x_1, x_2) = 16$  are shown as follows.



### 5.10

a. We have

$$f(\mathbf{x}) = f(\mathbf{x}_o) + Df(\mathbf{x}_o)(\mathbf{x} - \mathbf{x}_o) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_o)^\top D^2 f(\mathbf{x}_o)(\mathbf{x} - \mathbf{x}_o) + \dots$$

We compute

$$\begin{aligned} Df(\mathbf{x}) &= [e^{-x_2}, -x_1 e^{-x_2} + 1], \\ D^2 f(\mathbf{x}) &= \begin{bmatrix} 0 & -e^{-x_2} \\ -e^{-x_2} & x_1 e^{-x_2} \end{bmatrix}. \end{aligned}$$

Hence,

$$\begin{aligned} f(\mathbf{x}) &= 2 + [1, 0] \begin{bmatrix} x_1 - 1 \\ x_2 \end{bmatrix} + \frac{1}{2} [x_1 - 1, x_2] \begin{bmatrix} 0 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} x_1 - 1 \\ x_2 \end{bmatrix} + \dots \\ &= 1 + x_1 + x_2 - x_1 x_2 + \frac{1}{2} x_2^2 + \dots \end{aligned}$$

b. We compute

$$\begin{aligned} Df(\mathbf{x}) &= [4x_1^3 + 4x_1 x_2^2, 4x_1^2 x_2 + 4x_2^3], \\ D^2 f(\mathbf{x}) &= \begin{bmatrix} 12x_1^2 + 4x_2^2 & 8x_1 x_2 \\ 8x_1 x_2 & 4x_1^2 + 12x_2^2 \end{bmatrix}. \end{aligned}$$

Expanding  $f$  about the point  $\mathbf{x}_o$  yields

$$\begin{aligned} f(\mathbf{x}) &= 4 + [8, 8] \begin{bmatrix} x_1 - 1 \\ x_2 - 1 \end{bmatrix} + \frac{1}{2} [x_1 - 1, x_2 - 1] \begin{bmatrix} 16 & 8 \\ 8 & 16 \end{bmatrix} \begin{bmatrix} x_1 - 1 \\ x_2 - 1 \end{bmatrix} + \dots \\ &= 8x_1^2 + 8x_2^2 - 16x_1 - 16x_2 + 8x_1 x_2 + 12 + \dots \end{aligned}$$

c. We compute

$$\begin{aligned} Df(\mathbf{x}) &= [e^{x_1-x_2} + e^{x_1+x_2} + 1, -e^{x_1-x_2} + e^{x_1+x_2} + 1], \\ D^2f(\mathbf{x}) &= \begin{bmatrix} e^{x_1-x_2} + e^{x_1+x_2} & -e^{x_1-x_2} + e^{x_1+x_2} \\ -e^{x_1-x_2} + e^{x_1+x_2} & e^{x_1-x_2} + e^{x_1+x_2} \end{bmatrix}. \end{aligned}$$

Expanding  $f$  about the point  $\mathbf{x}_o$  yields

$$\begin{aligned} f(\mathbf{x}) &= 2 + 2e + [2e + 1, 1] \begin{bmatrix} x_1 - 1 \\ x_2 \end{bmatrix} + \frac{1}{2}[x_1 - 1, x_2] \begin{bmatrix} 2e & 0 \\ 0 & 2e \end{bmatrix} \begin{bmatrix} x_1 - 1 \\ x_2 \end{bmatrix} + \cdots \\ &= 1 + x_1 + x_2 + e(1 + x_1^2 + x_2^2) + \cdots. \end{aligned}$$

## 6. Basics of Unconstrained Optimization

### 6.1

a. In this case,  $\mathbf{x}^*$  is definitely not a local minimizer. To see this, note that  $\mathbf{d} = [1, -2]^\top$  is a feasible direction at  $\mathbf{x}^*$ . However,  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) = -1$ , which violates the FONC.

b. In this case,  $\mathbf{x}^*$  satisfies the FONC, and thus is possibly a local minimizer, but it is impossible to be definite based on the given information.

c. In this case,  $\mathbf{x}^*$  satisfies the SOSC, and thus is definitely a (strict) local minimizer.

d. In this case,  $\mathbf{x}^*$  is definitely not a local minimizer. To see this, note that  $\mathbf{d} = [0, 1]^\top$  is a feasible direction at  $\mathbf{x}^*$ , and  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) = 0$ . However,  $\mathbf{d}^\top \mathbf{F}(\mathbf{x}^*)\mathbf{d} = -1$ , which violates the SONC.

### 6.2

Because there are no constraints on  $x_1$  or  $x_2$ , we can utilize conditions for unconstrained optimization. To proceed, we first compute the function gradient and find the critical points, that is, the points that satisfy the FONC,

$$\nabla f(x_1, x_2) = \mathbf{0}.$$

The components of the gradient  $\nabla f(x_1, x_2)$  are

$$\frac{\partial f}{\partial x_1} = x_1^2 - 4 \quad \text{and} \quad \frac{\partial f}{\partial x_2} = x_2^2 - 16.$$

Thus there are four critical points:

$$\mathbf{x}^{(1)} = \begin{bmatrix} 2 \\ 4 \end{bmatrix}, \quad \mathbf{x}^{(2)} = \begin{bmatrix} 2 \\ -4 \end{bmatrix}, \quad \mathbf{x}^{(3)} = \begin{bmatrix} -2 \\ 4 \end{bmatrix}, \quad \text{and} \quad \mathbf{x}^{(4)} = \begin{bmatrix} -2 \\ -4 \end{bmatrix}.$$

We next compute the Hessian matrix of the function  $f$ :

$$\mathbf{F}(\mathbf{x}) = \begin{bmatrix} 2x_1 & 0 \\ 0 & 2x_2 \end{bmatrix}.$$

Note that  $\mathbf{F}(\mathbf{x}^{(1)}) > 0$  and therefore,  $\mathbf{x}^{(1)}$  is a strict local minimizer. Next,  $\mathbf{F}(\mathbf{x}^{(4)}) < 0$  and therefore,  $\mathbf{x}^{(4)}$  is a strict local maximizer. The Hessian is indefinite at  $\mathbf{x}^{(2)}$  and  $\mathbf{x}^{(3)}$  and so these points are neither maximizer nor minimizers.

### 6.3

Suppose  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ , and  $\mathbf{x}^* \in \Omega' \subset \Omega$ . Let  $\mathbf{x} \in \Omega'$ . Then,  $\mathbf{x} \in \Omega$  and therefore  $f(\mathbf{x}^*) \leq f(\mathbf{x})$ . Hence,  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega'$ .

### 6.4

Suppose  $\mathbf{x}^*$  is an interior point of  $\Omega$ . Therefore, there exists  $\varepsilon > 0$  such that  $\{\mathbf{y} : \|\mathbf{y} - \mathbf{x}^*\| < \varepsilon\} \subset \Omega$ . Since  $\mathbf{x}^*$  is a local minimizer of  $f$  over  $\Omega$ , there exists  $\varepsilon' > 0$  such that  $f(\mathbf{x}^*) \leq f(\mathbf{x})$  for all  $\mathbf{x} \in \{\mathbf{y} : \|\mathbf{y} - \mathbf{x}^*\| < \varepsilon'\}$ . Take  $\varepsilon'' = \min(\varepsilon, \varepsilon')$ . Then,  $\{\mathbf{y} : \|\mathbf{y} - \mathbf{x}^*\| < \varepsilon''\} \subset \Omega'$ , and  $f(\mathbf{x}^*) \leq f(\mathbf{x})$  for all  $\mathbf{x} \in \{\mathbf{y} : \|\mathbf{y} - \mathbf{x}^*\| < \varepsilon''\}$ . Thus,  $\mathbf{x}^*$  is a local minimizer of  $f$  over  $\Omega'$ .

To show that we cannot make the same conclusion if  $\mathbf{x}^*$  is not an interior point, let  $\Omega = \{0\}$ ,  $\Omega' = [-1, 1]$ , and  $f(x) = x$ . Clearly,  $0 \in \Omega$  is a local minimizer of  $f$  over  $\Omega$ . However,  $0 \in \Omega'$  is not a local minimizer of  $f$  over  $\Omega'$ .

### 6.5

a. The TONC is: if  $f''(0) = 0$ , then  $f'''(0) = 0$ . To prove this, suppose  $f''(0) = 0$ . Now, by the FONC, we also have  $f'(0) = 0$ . Hence, by Taylor's theorem,

$$f(x) = f(0) + \frac{f'''(0)}{3!}x^3 + o(x^3).$$

Since 0 is a local minimizer,  $f(x) \geq f(0)$  for all  $x$  sufficiently close to 0. Hence, for all such  $x$ ,

$$\frac{f'''(0)}{3!}x^3 \geq o(x^3).$$

Now, if  $x > 0$ , then

$$f'''(0) \geq 3! \frac{o(x^3)}{x^3},$$

which implies that  $f'''(0) \geq 0$ . On the other hand, if  $x < 0$ , then

$$f'''(0) \leq 3! \frac{o(x^3)}{x^3},$$

which implies that  $f'''(0) \leq 0$ . This implies that  $f'''(0) = 0$ , as required.

b. Let  $f(x) = -x^4$ . Then,  $f'(0) = 0$ ,  $f''(0) = 0$ , and  $f'''(0) = 0$ , which means that the FONC, SONC, and TONC are all satisfied. However, 0 is not a local minimizer:  $f(x) < 0$  for all  $x \neq 0$ .

c. The answer is yes. To see this, we first write

$$f(x) = f(0) + f'(0)x + \frac{f''(0)}{2}x^2 + \frac{f'''(0)}{3!}x^3.$$

Now, if the FONC is satisfied, then

$$f(x) = f(0) + \frac{f''(0)}{2}x^2 + \frac{f'''(0)}{3!}x^3.$$

Moreover, if the SONC is satisfied, then either (i)  $f''(0) > 0$  or (ii)  $f''(0) = 0$ . In the case (i), it is clear from the above equation that  $f(x) \geq f(0)$  for all  $x$  sufficiently close to 0 (because the third term on the right-hand side is  $o(x^2)$ ). In the case (ii), the TONC implies that  $f(x) = f(0)$  for all  $x$ . In either case,  $f(x) \geq f(0)$  for all  $x$  sufficiently close to 0. This shows that 0 is a local minimizer.

## 6.6

a. The TONC is: if  $f'(0) = 0$  and  $f''(0) = 0$ , then  $f'''(0) \geq 0$ . To prove this, suppose  $f'(0) = 0$  and  $f''(0) = 0$ . By Taylor's theorem, for  $x \geq 0$ ,

$$f(x) = f(0) + \frac{f'''(0)}{3!}x^3 + o(x^3).$$

Since 0 is a local minimizer,  $f(x) \geq f(0)$  for sufficiently small  $x \geq 0$ . Hence, for all  $x \geq 0$  sufficiently small,

$$f'''(0) \geq 3! \frac{o(x^3)}{x^3}.$$

This implies that  $f'''(0) \geq 0$ , as required.

b. Let  $f(x) = -x^4$ . Then,  $f'(0) = 0$ ,  $f''(0) = 0$ , and  $f'''(0) = 0$ , which means that the FONC, SONC, and TONC are all satisfied. However, 0 is not a local minimizer:  $f(x) < 0$  for all  $x > 0$ .

## 6.7

For convenience, let  $\mathbf{z}_0 = \mathbf{x}_0 + \arg \min_{\mathbf{x} \in \Omega} f(\mathbf{x})$ . Thus we want to show that  $\mathbf{z}_0 = \arg \min_{\mathbf{y} \in \Omega'} f(\mathbf{y})$ ; i.e., for all  $\mathbf{y} \in \Omega'$ ,  $f(\mathbf{y} - \mathbf{x}_0) \geq f(\mathbf{z}_0 - \mathbf{x}_0)$ . So fix  $\mathbf{y} \in \Omega'$ . Then,  $\mathbf{y} - \mathbf{x}_0 \in \Omega$ . Hence,

$$\begin{aligned} f(\mathbf{y} - \mathbf{x}_0) &\geq \min_{\mathbf{x} \in \Omega} f(\mathbf{x}) \\ &= f\left(\arg \min_{\mathbf{x} \in \Omega} f(\mathbf{x})\right) \\ &= f(\mathbf{z}_0 - \mathbf{x}_0), \end{aligned}$$

which completes the proof.

## 6.8

a. The gradient and Hessian of  $f$  are

$$\begin{aligned} \nabla f(\mathbf{x}) &= 2 \begin{bmatrix} 1 & 3 \\ 3 & 7 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 3 \\ 5 \end{bmatrix} \\ \mathbf{F}(\mathbf{x}) &= 2 \begin{bmatrix} 1 & 3 \\ 3 & 7 \end{bmatrix}. \end{aligned}$$

Hence,  $\nabla f([1, 1]^\top) = [11, 25]^\top$ , and  $\mathbf{F}([1, 1]^\top)$  is as shown above.

b. The direction of maximal rate of increase is the direction of the gradient. Hence, the directional derivative with respect to a unit vector in this direction is

$$\left( \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|} \right)^\top \nabla f(\mathbf{x}) = \frac{(\nabla f(\mathbf{x}))^\top \nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|} = \|\nabla f(\mathbf{x})\|.$$

At  $\mathbf{x} = [1, 1]^\top$ , we have  $\|\nabla f([1, 1]^\top)\| = \sqrt{11^2 + 25^2} = 27.31$ .

c. The FONC in this case is  $\nabla f(\mathbf{x}) = \mathbf{0}$ . Solving, we get

$$\mathbf{x} = \begin{bmatrix} 3/2 \\ -1 \end{bmatrix}.$$

The point above does not satisfy the SONC because the Hessian is not positive semidefinite (its determinant is negative).

## 6.9

a. A differentiable function  $f$  decreases most rapidly in the direction of the negative gradient. In our problem,

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix}^\top = \begin{bmatrix} 2x_1x_2 + x_2^3 & x_1^2 + 3x_1x_2^2 \end{bmatrix}^\top.$$

Hence, the direction of most rapid decrease is

$$-\nabla f(\mathbf{x}^{(0)}) = -\begin{bmatrix} 5 & 10 \end{bmatrix}^\top.$$

b. The rate of increase of  $f$  at  $\mathbf{x}^{(0)}$  in the direction  $-\nabla f(\mathbf{x}^{(0)})$  is

$$\nabla f(\mathbf{x}^{(0)})^\top \frac{-\nabla f(\mathbf{x}^{(0)})}{\|\nabla f(\mathbf{x}^{(0)})\|} = -\|\nabla f(\mathbf{x}^{(0)})\| = -\sqrt{125} = -5\sqrt{5}.$$

c. The rate of increase of  $f$  at  $\mathbf{x}^{(0)}$  in the direction  $\mathbf{d}$  is

$$\nabla f(\mathbf{x}^{(0)})^\top \frac{\mathbf{d}}{\|\mathbf{d}\|} = \begin{bmatrix} 5 & 10 \end{bmatrix} \begin{bmatrix} 3 \\ 4 \end{bmatrix} \frac{1}{5} = 11.$$

## 6.10

a. We can rewrite  $f$  as

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \begin{bmatrix} 4 & 4 \\ 4 & 2 \end{bmatrix} \mathbf{x} + \mathbf{x}^\top \begin{bmatrix} 3 \\ 4 \end{bmatrix} + 7.$$

The gradient and Hessian of  $f$  are

$$\begin{aligned} \nabla f(\mathbf{x}) &= \begin{bmatrix} 4 & 4 \\ 4 & 2 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \\ \mathbf{F}(\mathbf{x}) &= \begin{bmatrix} 4 & 4 \\ 4 & 2 \end{bmatrix}. \end{aligned}$$

Hence  $\nabla f([0, 1]^\top) = [7, 6]^\top$ . The directional derivative is

$$[1, 0]^\top \nabla f([0, 1]^\top) = 7.$$



b. The FONC in this case is  $\nabla f(\mathbf{x}) = \mathbf{0}$ . The only point satisfying the FONC is

$$\mathbf{x}^* = \frac{1}{4} \begin{bmatrix} -5 \\ 2 \end{bmatrix}.$$

The point above does not satisfy the SONC because the Hessian is not positive semidefinite (its determinant is negative). Therefore,  $f$  does not have a minimizer.

#### 6.11

a. Write the objective function as  $f(\mathbf{x}) = -x_2^2$ . In this problem the only feasible directions at  $\mathbf{0}$  are of the form  $\mathbf{d} = [d_1, 0]^\top$ . Hence,  $\mathbf{d}^\top \nabla f(\mathbf{0}) = 0$  for all feasible directions  $\mathbf{d}$  at  $\mathbf{0}$ .

b. The point  $\mathbf{0}$  is a local maximizer, because  $f(\mathbf{0}) = 0$ , while any feasible point  $\mathbf{x}$  satisfies  $f(\mathbf{x}) \leq 0$ .

The point  $\mathbf{0}$  is not a strict local maximizer because for any  $\mathbf{x}$  of the form  $\mathbf{x} = [x_1, 0]^\top$ , we have  $f(\mathbf{x}) = 0 = f(\mathbf{0})$ , and there are such points in any neighborhood of  $\mathbf{0}$ .

The point  $\mathbf{0}$  is not a local minimizer because for any point  $\mathbf{x}$  of the form  $\mathbf{x} = [x_1, x_1^2]^\top$  with  $x_1 > 0$ , we have  $f(\mathbf{x}) = -x_1^4 < 0$ , and there are such points in any neighborhood of  $\mathbf{0}$ . Since  $\mathbf{0}$  is not a local minimizer, it is also not a strict local minimizer.

#### 6.12

a. We have  $\nabla f(\mathbf{x}^*) = [0, 5]^\top$ . The only feasible directions at  $\mathbf{x}^*$  are of the form  $\mathbf{d} = [d_1, d_2]^\top$  with  $d_2 \geq 0$ . Therefore, for such feasible directions,  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) = 5d_2 \geq 0$ . Hence,  $\mathbf{x}^* = [0, 1]^\top$  satisfies the first order necessary condition.

b. We have  $\mathbf{F}(\mathbf{x}^*) = \mathbf{O}$ . Therefore, for any  $\mathbf{d}$ ,  $\mathbf{d}^\top \mathbf{F}(\mathbf{x}^*) \mathbf{d} \geq 0$ . Hence,  $\mathbf{x}^* = [0, 1]^\top$  satisfies the second order necessary condition.

c. Consider points of the form  $\mathbf{x} = [x_1, -x_1^2 + 1]^\top$ ,  $x_1 \in \mathbb{R}$ . Such points are in  $\Omega$ , and are arbitrarily close to  $\mathbf{x}^*$ . However, for such points  $\mathbf{x} \neq \mathbf{x}^*$ ,

$$f(\mathbf{x}) = 5(-x_1^2 + 1) = 5 - 5x_1^2 < 5 = f(\mathbf{x}^*).$$

Hence,  $\mathbf{x}^*$  is not a local minimizer.

#### 6.13

a. We have  $\nabla f(\mathbf{x}^*) = -[3, 0]^\top$ . The only feasible directions at  $\mathbf{x}^*$  are of the form  $\mathbf{d} = [d_1, d_2]^\top$  with  $d_1 \leq 0$ . Therefore, for such feasible directions,  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) = 3d_1 \geq 0$ . Hence,  $\mathbf{x}^* = [2, 0]^\top$  satisfies the first order necessary condition.

b. We have  $\mathbf{F}(\mathbf{x}^*) = \mathbf{O}$ . Therefore, for any  $\mathbf{d}$ ,  $\mathbf{d}^\top \mathbf{F}(\mathbf{x}^*) \mathbf{d} \geq 0$ . Hence,  $\mathbf{x}^* = [2, 0]^\top$  satisfies the second order necessary condition.

c. Yes,  $\mathbf{x}^*$  is a local minimizer. To see this, notice that any feasible point  $\mathbf{x} = [x_1, x_2]^\top \neq \mathbf{x}^*$  is such that  $x_1 < 2$ . Hence, for such points  $\mathbf{x} \neq \mathbf{x}^*$ ,

$$f(\mathbf{x}) = -3x_1 > -6 = f(\mathbf{x}^*).$$

In fact,  $\mathbf{x}^*$  is a strict local minimizer.

#### 6.14

a. We have  $\nabla f(\mathbf{x}) = [0, 1]$ , which is nonzero everywhere. Hence, no interior point satisfies the FONC. Moreover, any boundary point with a feasible direction  $\mathbf{d}$  such that  $d_2 < 0$  cannot satisfy the FONC, because for such a  $\mathbf{d}$ ,  $\mathbf{d}^\top \nabla f(\mathbf{x}) = d_2 < 0$ . By drawing a picture, it is easy to see that the only boundary point remaining is  $\mathbf{x}^* = [0, 1]^\top$ . For this point, any feasible direction satisfies  $d_2 \geq 0$ . Hence, for any feasible direction,  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) = d_2 \geq 0$ . Hence,  $\mathbf{x}^* = [0, 1]^\top$  satisfies the FONC, and is the only such point.

b. We have  $\mathbf{F}(\mathbf{x}) = \mathbf{O}$ . So any point (and in particular  $\mathbf{x}^* = [0, 1]^\top$ ) satisfies the SONC.

c. The point  $\mathbf{x}^* = [0, 1]^\top$  is not a local minimizer. To see this, consider points of the form  $\mathbf{x} = [\sqrt{1 - x_2^2}, x_2]^\top$  where  $x_2 \in [1/2, 1)$ . It is clear that such points are feasible, and are arbitrarily close to  $\mathbf{x}^* = [0, 1]^\top$ . However, for such points,  $f(\mathbf{x}) = x_2 < 1 = f(\mathbf{x}^*)$ .

---

**6.15**

a. We have  $\nabla f(\mathbf{x}^*) = [3, 0]^\top$ . The only feasible directions at  $\mathbf{x}^*$  are of the form  $\mathbf{d} = [d_1, d_2]^\top$  with  $d_1 \geq 0$ . Therefore, for such feasible directions,  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) = 3d_1 \geq 0$ . Hence,  $\mathbf{x}^* = [2, 0]^\top$  satisfies the first order necessary condition.

b. We have  $\mathbf{F}(\mathbf{x}^*) = \mathbf{O}$ . Therefore, for any  $\mathbf{d}$ ,  $\mathbf{d}^\top \mathbf{F}(\mathbf{x}^*) \mathbf{d} \geq 0$ . Hence,  $\mathbf{x}^* = [2, 0]^\top$  satisfies the second order necessary condition.

c. Consider points of the form  $\mathbf{x} = [-x_2^2 + 2, x_2]^\top$ ,  $x_2 \in \mathbb{R}$ . Such points are in  $\Omega$ , and could be arbitrarily close to  $\mathbf{x}^*$ . However, for such points  $\mathbf{x} \neq \mathbf{x}^*$ ,

$$f(\mathbf{x}) = 3(-x_2^2 + 2) = 6 - 6x_2^2 < 6 = f(\mathbf{x}^*).$$

Hence,  $\mathbf{x}^*$  is not a local minimizer.

---

**6.16**

a. We have  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ . Therefore, for any feasible direction  $\mathbf{d}$  at  $\mathbf{x}^*$ , we have  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) = 0$ . Hence,  $\mathbf{x}^*$  satisfies the first-order necessary condition.

b. We have

$$\mathbf{F}(\mathbf{x}^*) = \begin{bmatrix} 8 & 0 \\ 0 & -2 \end{bmatrix}.$$

Any feasible direction  $\mathbf{d}$  at  $\mathbf{x}^*$  has the form  $\mathbf{d} = [d_1, d_2]^\top$  where  $d_2 \leq 2d_1$ ,  $d_1, d_2 \geq 0$ . Therefore, for any feasible direction  $\mathbf{d}$  at  $\mathbf{x}^*$ , we have

$$\mathbf{d}^\top \mathbf{F}(\mathbf{x}^*) \mathbf{d} = 8d_1^2 - 2d_2^2 \geq 8d_1^2 - 2(2d_1)^2 = 0.$$

Hence,  $\mathbf{x}^*$  satisfies the second-order necessary condition.

c. We have  $f(\mathbf{x}^*) = 0$ . Any point of the form  $\mathbf{x} = [x_1, x_1^2 + 2x_1]^\top$ ,  $x_1 > 0$ , is feasible and has objective function value given by

$$f(\mathbf{x}) = 4x_1^2 - (x_1^2 + 2x_1)^2 = -(x_1^4 + 4x_1^3) < 0 = f(\mathbf{x}^*),$$

Moreover, there are such points in any neighborhood of  $\mathbf{x}^*$ . Therefore, the point  $\mathbf{x}^*$  is not a local minimizer.

---

**6.17**

a. We have  $\nabla f(\mathbf{x}^*) = [1/x_1^*, 1/x_2^*]^\top$ . If  $\mathbf{x}^*$  were an interior point, then  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ . But this is clearly impossible. Therefore,  $\mathbf{x}^*$  cannot possibly be an interior point.

b. We have  $\mathbf{F}(\mathbf{x}) = -\text{diag}[1/x_1^2, 1/x_2^2]$ , which is negative definite everywhere. Therefore, the second-order necessary condition is satisfied everywhere. (Note that because we have a maximization problem, *negative* definiteness is the relevant condition.)

---

**6.18**

Given  $x \in \mathbb{R}$ , let

$$f(x) = \sum_{i=1}^n (x - x_i)^2,$$

so that  $\bar{x}$  is the minimizer of  $f$ . By the FONC,

$$f'(\bar{x}) = 0,$$

and hence

$$\sum_{i=1}^n 2(\bar{x} - x_i) = 0,$$

which on solving gives

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i.$$

**6.19**

Let  $\theta_1$  be the angle from the horizontal to the bottom of the picture, and  $\theta_2$  the angle from the horizontal to the top of the picture. Then,  $\tan(\theta) = (\tan(\theta_2) - \tan(\theta_1))/(1 + \tan(\theta_2)\tan(\theta_1))$ . Now,  $\tan(\theta_1) = b/x$  and  $\tan(\theta_2) = (a+b)/x$ . Hence, the objective function that we wish to maximize is

$$f(x) = \frac{(a+b)/x - b/x}{1 + b(a+b)/x^2} = \frac{a}{x + b(a+b)/x}.$$

We have

$$f'(x) = -\frac{a^2}{(x + b(a+b)/x)^2} \left(1 - \frac{b(a+b)}{x^2}\right).$$

Let  $x^*$  be the optimal distance. Then, by the FONC, we have  $f'(x^*) = 0$ , which gives

$$\begin{aligned} 1 - \frac{b(a+b)}{(x^*)^2} &= 0 \\ \Rightarrow x^* &= \sqrt{b(a+b)}. \end{aligned}$$

**6.20**

The squared distance from the sensor to the baby's heart is  $1 + x^2$ , while the squared distance from the sensor to the mother's heart is  $1 + (2 - x)^2$ . Therefore, the signal to noise ratio is

$$f(x) = \frac{1 + (2 - x)^2}{1 + x^2}.$$

We have

$$\begin{aligned} f'(x) &= \frac{-2(2 - x)(1 + x^2) - 2x(1 + (2 - x)^2)}{(1 + x^2)^2} \\ &= \frac{4(x^2 - 2x - 1)}{(1 + x^2)^2}. \end{aligned}$$

By the FONC, at the optimal position  $x^*$ , we have  $f'(x^*) = 0$ . Hence, either  $x^* = 1 - \sqrt{2}$  or  $x^* = 1 + \sqrt{2}$ . From the figure, it easy to see that  $x^* = 1 - \sqrt{2}$  is the optimal position.

**6.21**

a. Let  $x$  be the decision variable. Write the total travel time as  $f(x)$ , which is given by

$$f(x) = \frac{\sqrt{1 + x^2}}{v_1} + \frac{\sqrt{1 + (d - x)^2}}{v_2}.$$

Differentiating the above expression, we get

$$f'(x) = \frac{x}{v_1 \sqrt{1 + x^2}} - \frac{d - x}{v_2 \sqrt{1 + (d - x)^2}}.$$

By the first order necessary condition, the optimal path satisfies  $f'(x^*) = 0$ , which corresponds to

$$\frac{x^*}{v_1 \sqrt{1 + (x^*)^2}} = \frac{d - x^*}{v_2 \sqrt{1 + (d - x^*)^2}},$$

or  $\sin \theta_1 / v_1 = \sin \theta_2 / v_2$ . Upon rearranging, we obtain the desired equation.

b. The second derivative of  $f$  is given by

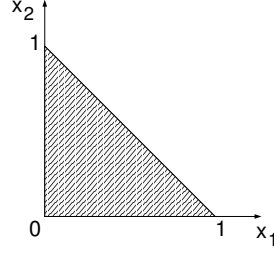
$$f''(x) = \frac{1}{v_1(1 + x^2)^{3/2}} + \frac{1}{v_2(1 + (d - x)^2)^{3/2}}.$$

Hence,  $f''(x^*) > 0$ , which shows that the second order sufficient condition holds.

---

**6.22**

a. We have  $f(\mathbf{x}) = U_1(x_1) + U_2(x_2)$  and  $\Omega = \{\mathbf{x} : x_1, x_2 \geq 0, x_1 + x_2 \leq 1\}$ . A picture of  $\Omega$  looks like:



b. We have  $\nabla f(\mathbf{x}) = [a_1, a_2]^\top$ . Because  $\nabla f(\mathbf{x}) \neq \mathbf{0}$ , for all  $\mathbf{x}$ , we conclude that no interior point satisfies the FONC. Next, consider any feasible point  $\mathbf{x}$  for which  $x_2 > 0$ . At such a point, the vector  $\mathbf{d} = [1, -1]^\top$  is a feasible direction. But then  $\mathbf{d}^\top \nabla f(\mathbf{x}) = a_1 - a_2 > 0$  which means that FONC is violated (recall that the problem is to *maximize*  $f$ ). So clearly the remaining candidates are those  $\mathbf{x}$  for which  $x_2 = 0$ . Among these, if  $x_1 < 1$ , then  $\mathbf{d} = [0, 1]^\top$  is a feasible direction, in which case we have  $\mathbf{d}^\top \nabla f(\mathbf{x}) = a_2 > 0$ . This leaves the point  $\mathbf{x} = [1, 0]^\top$ . At this point, any feasible direction  $\mathbf{d}$  satisfies  $d_1 \leq 0$  and  $d_2 \leq -d_1$ . Hence, for any feasible direction  $\mathbf{d}$ , we have

$$\mathbf{d}^\top \nabla f(\mathbf{x}) = d_1 a_1 + d_2 a_2 \leq d_1 a_1 + (-d_1) a_2 = d_1 (a_1 - a_2) \leq 0.$$

So, the only feasible point that satisfies the FONC is  $[1, 0]^\top$ .

c. We have  $\mathbf{F}(\mathbf{x}) = \mathbf{O} \leq 0$ . Hence, any point satisfies the SONC (again, recall that the problem is to *maximize*  $f$ ).

---

**6.23**

We have

$$\nabla f(\mathbf{x}) = \begin{bmatrix} 4(x_1 - x_2)^3 + 2x_1 - 2 \\ -4(x_1 - x_2)^3 - 2x_2 + 2 \end{bmatrix}.$$

Setting  $\nabla f(\mathbf{x}) = \mathbf{0}$  we get

$$\begin{aligned} 4(x_1 - x_2)^3 + 2x_1 - 2 &= 0 \\ -4(x_1 - x_2)^3 - 2x_2 + 2 &= 0. \end{aligned}$$

Adding the two equations, we obtain  $x_1 = x_2$ , and substituting back yields

$$x_1 = x_2 = 1.$$

Hence, the only point satisfying the FONC is  $[1, 1]^\top$ .

We have

$$\mathbf{F}(\mathbf{x}) = \begin{bmatrix} 12(x_1 - x_2)^2 + 2 & -12(x_1 - x_2)^2 \\ -12(x_1 - x_2)^2 & 12(x_1 - x_2)^2 - 2 \end{bmatrix}.$$

Hence

$$\mathbf{F}([1, 1]^\top) = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$$

Since  $\mathbf{F}([1, 1]^\top)$  is not positive semidefinite, the point  $[1, 1]^\top$  does not satisfy the SONC.

---

**6.24**

Suppose  $\mathbf{d}$  is a feasible direction at  $\mathbf{x}$ . Then, there exists  $\alpha_0 > 0$  such that  $\mathbf{x} + \alpha \mathbf{d} \in \Omega$  for all  $\alpha \in [0, \alpha_0]$ . Let  $\beta > 0$  be given. Then,  $\mathbf{x} + \alpha(\beta \mathbf{d}) \in \Omega$  for all  $\alpha \in [0, \alpha_0/\beta]$ . Since  $\alpha_0/\beta > 0$ , by definition  $\beta \mathbf{d}$  is also a feasible direction at  $\mathbf{x}$ .

---

**6.25**

$\Rightarrow$ : Suppose  $\mathbf{d}$  is feasible at  $\mathbf{x} \in \Omega$ . Then, there exists  $\alpha > 0$  such that  $\mathbf{x} + \alpha \mathbf{d} \in \Omega$ , that is,  $\mathbf{A}(\mathbf{x} + \alpha \mathbf{d}) = \mathbf{b}$ . Since  $\mathbf{A}\mathbf{x} = \mathbf{b}$  and  $\alpha \neq 0$ , we conclude that  $\mathbf{A}\mathbf{d} = \mathbf{0}$ .

$\Leftarrow$ : Suppose  $\mathbf{A}\mathbf{d} = \mathbf{0}$ . Then, for any  $\alpha \in [0, 1]$ , we have  $\alpha\mathbf{A}\mathbf{d} = \mathbf{0}$ . Adding this equation to  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , we obtain  $\mathbf{A}(\mathbf{x} + \alpha\mathbf{d}) = \mathbf{b}$ , that is,  $\mathbf{x} + \alpha\mathbf{d} \in \Omega$  for all  $\alpha \in [0, 1]$ . Therefore,  $\mathbf{d}$  is a feasible direction at  $\mathbf{x}$ .

**6.26**

The vector  $\mathbf{d} = [1, 1]^\top$  is a feasible direction at  $\mathbf{0}$ . Now,

$$\mathbf{d}^\top \nabla f(\mathbf{0}) = \frac{\partial f}{\partial x_1}(\mathbf{0}) + \frac{\partial f}{\partial x_2}(\mathbf{0}).$$

Since  $\nabla f(\mathbf{0}) \leq \mathbf{0}$  and  $\nabla f(\mathbf{0}) \neq \mathbf{0}$ , then

$$\mathbf{d}^\top \nabla f(\mathbf{0}) < 0.$$

Hence, by the FONC,  $\mathbf{0}$  is not a local minimizer.

**6.27**

We have  $\nabla f(\mathbf{x}) = \mathbf{c} \neq \mathbf{0}$ . Therefore, for any  $\mathbf{x} \in \overset{\circ}{\Omega}$ , we have  $\nabla f(\mathbf{x}) \neq \mathbf{0}$ . Hence, by Corollary 6.1,  $\mathbf{x} \in \overset{\circ}{\Omega}$  cannot be a local minimizer (and therefore it cannot be a solution).

**6.28**

The objective function is  $f(\mathbf{x}) = -c_1x_1 - c_2x_2$ . Therefore,  $\nabla f(\mathbf{x}) = [-c_1, -c_2]^\top \neq \mathbf{0}$  for all  $\mathbf{x}$ . Thus, by FONC, the optimal solution  $\mathbf{x}^*$  cannot lie in the interior of the feasible set. Next, for all  $\mathbf{x} \in L_1 \cup L_2$ ,  $\mathbf{d} = [1, 1]^\top$  is a feasible direction. Therefore,  $\mathbf{d}^\top \nabla f(\mathbf{x}) = -c_1 - c_2 < 0$ . Hence, by FONC, the optimal solution  $\mathbf{x}^*$  cannot lie in  $L_1 \cup L_2$ . Lastly, for all  $\mathbf{x} \in L_3$ ,  $\mathbf{d} = [1, -1]^\top$  is a feasible direction. Therefore,  $\mathbf{d}^\top \nabla f(\mathbf{x}) = c_2 - c_1 < 0$ . Hence, by FONC, the optimal solution  $\mathbf{x}^*$  cannot lie in  $L_3$ . Therefore, by elimination, the unique optimal feasible solution must be  $[1, 0]^\top$ .

**6.29**

a. We write

$$\begin{aligned} f(a, b) &= \frac{1}{n} \sum_{i=1}^n a^2 x_i^2 + b^2 + y_i^2 + 2x_i a b - 2x_i y_i a - 2y_i b \\ &= a^2 \left( \frac{1}{n} \sum_{i=1}^n x_i^2 \right) + b^2 + 2 \left( \frac{1}{n} \sum_{i=1}^n x_i \right) a b \\ &\quad - 2 \left( \frac{1}{n} \sum_{i=1}^n x_i y_i \right) a - 2 \left( \frac{1}{n} \sum_{i=1}^n y_i \right) b + \left( \frac{1}{n} \sum_{i=1}^n y_i^2 \right) \\ &= [a \ b] \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n x_i^2 & \frac{1}{n} \sum_{i=1}^n x_i \\ \frac{1}{n} \sum_{i=1}^n x_i y_i & 1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} \\ &\quad - 2 \left[ \frac{1}{n} \sum_{i=1}^n x_i y_i, \frac{1}{n} \sum_{i=1}^n y_i \right] \begin{bmatrix} a \\ b \end{bmatrix} + \frac{1}{n} \sum_{i=1}^n y_i^2 \\ &= \mathbf{z}^\top \mathbf{Q} \mathbf{z} - 2\mathbf{c}^\top \mathbf{z} + d, \end{aligned}$$

where  $\mathbf{z}$ ,  $\mathbf{Q}$ ,  $\mathbf{c}$  and  $d$  are defined in the obvious way.

b. If the point  $\mathbf{z}^* = [a^*, b^*]^\top$  is a solution, then by the FONC, we have  $\nabla f(\mathbf{z}^*) = 2\mathbf{Q}\mathbf{z}^* - 2\mathbf{c} = \mathbf{0}$ , which means  $\mathbf{Q}\mathbf{z}^* = \mathbf{c}$ . Now, since  $\overline{X^2} - (\overline{X})^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \overline{X})^2$ , and the  $x_i$  are not all equal, then  $\det \mathbf{Q} = \overline{X^2} - (\overline{X})^2 \neq 0$ . Hence,  $\mathbf{Q}$  is nonsingular, and hence

$$\mathbf{z}^* = \mathbf{Q}^{-1} \mathbf{c} = \frac{1}{\overline{X^2} - (\overline{X})^2} \begin{bmatrix} 1 & -\overline{X} \\ -\overline{X} & \overline{X^2} \end{bmatrix} \begin{bmatrix} \overline{XY} \\ \overline{Y} \end{bmatrix} = \begin{bmatrix} \frac{\overline{XY} - (\overline{X})(\overline{Y})}{\overline{X^2} - (\overline{X})^2} \\ \frac{(\overline{X^2})(\overline{Y}) - (\overline{X})(\overline{XY})}{\overline{X^2} - (\overline{X})^2} \end{bmatrix}.$$

Since  $\mathbf{Q} > 0$ , then by the SOSC, the point  $\mathbf{z}^*$  is a strict local minimizer. Since  $\mathbf{z}^*$  is the only point satisfying the FONC, then  $\mathbf{z}^*$  is the only local minimizer.

c. We have

$$a^* \overline{X} + b^* = \left( \frac{\overline{XY} - (\overline{X})(\overline{Y})}{\overline{X^2} - (\overline{X})^2} \right) \overline{X} + \frac{(\overline{X^2})(\overline{Y}) - (\overline{X})(\overline{XY})}{\overline{X^2} - (\overline{X})^2} = \overline{Y}.$$

**6.30**

Given  $\mathbf{x} \in \mathbb{R}^n$ , let

$$f(\mathbf{x}) = \frac{1}{p} \sum_{i=1}^p \|\mathbf{x} - \mathbf{x}^{(i)}\|^2$$

be the average squared error between  $\mathbf{x}$  and  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p)}$ . We can rewrite  $f$  as

$$\begin{aligned} f(\mathbf{x}) &= \frac{1}{p} \sum_{i=1}^p (\mathbf{x} - \mathbf{x}^{(i)})^\top (\mathbf{x} - \mathbf{x}^{(i)}) \\ &= \mathbf{x}^\top \mathbf{x} - 2 \left( \frac{1}{p} \sum_{i=1}^p \mathbf{x}^{(i)} \right)^\top \mathbf{x} + \frac{1}{p} \|\mathbf{x}^{(i)}\|^2. \end{aligned}$$

So  $f$  is a quadratic function. Since  $\bar{\mathbf{x}}$  is the minimizer of  $f$ , then by the FONC,  $\nabla f(\bar{\mathbf{x}}) = \mathbf{0}$ , i.e.,

$$2\bar{\mathbf{x}} - 2 \frac{1}{p} \sum_{i=1}^p \mathbf{x}^{(i)} = \mathbf{0}.$$

Hence, we get

$$\bar{\mathbf{x}} = \frac{1}{p} \sum_{i=1}^p \mathbf{x}^{(i)},$$

i.e.,  $\bar{\mathbf{x}}$  is just the average, or centroid, or center of gravity, of  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p)}$ .

The Hessian of  $f$  at  $\bar{\mathbf{x}}$  is

$$\mathbf{F}(\bar{\mathbf{x}}) = 2\mathbf{I}_n,$$

which is positive definite. Hence, by the SOSC,  $\bar{\mathbf{x}}$  is a strict local minimizer of  $f$  (in fact, it is a strict global minimizer because  $f$  is a convex quadratic function).

**6.31**

Fix any  $\mathbf{x} \in \Omega$ . The vector  $\mathbf{d} = \mathbf{x} - \mathbf{x}^*$  is feasible at  $\mathbf{x}^*$  (by convexity of  $\Omega$ ). By Taylor's formula, we have

$$f(\mathbf{x}) = f(\mathbf{x}^*) + \mathbf{d}^\top \nabla f(\mathbf{x}^*) + o(\|\mathbf{d}\|) \geq f(\mathbf{x}^*) + c\|\mathbf{d}\| + o(\|\mathbf{d}\|).$$

Therefore, for all  $\mathbf{x}$  sufficiently close to  $\mathbf{x}^*$ , we have  $f(\mathbf{x}) > f(\mathbf{x}^*)$ . Hence,  $\mathbf{x}^*$  is a strict local minimizer.

**6.32**

Since  $f \in \mathcal{C}^2$ ,  $\mathbf{F}(\mathbf{x}^*) = \mathbf{F}^\top(\mathbf{x}^*)$ . Let  $\mathbf{d} \neq \mathbf{0}$  be a feasible directions at  $\mathbf{x}^*$ . By Taylor's theorem,

$$f(\mathbf{x}^* + \mathbf{d}) - f(\mathbf{x}^*) = \frac{1}{2} \mathbf{d}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{d} + o(\|\mathbf{d}\|^2).$$

Using conditions a and b, we get

$$f(\mathbf{x}^* + \mathbf{d}) - f(\mathbf{x}^*) \geq c\|\mathbf{d}\|^2 + o(\|\mathbf{d}\|^2),$$

Therefore, for all  $\mathbf{d}$  such that  $\|\mathbf{d}\|$  is sufficiently small,

$$f(\mathbf{x}^* + \mathbf{d}) > f(\mathbf{x}^*),$$

and the proof is completed.

**6.33**

Necessity follows from the FONC. To prove sufficiency, we write  $f$  as

$$f(\mathbf{x}) = \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top \mathbf{Q} (\mathbf{x} - \mathbf{x}^*) - \frac{1}{2} \mathbf{x}^{*\top} \mathbf{Q} \mathbf{x}^*$$

where  $\mathbf{x}^* = \mathbf{Q}^{-1} \mathbf{b}$  is the unique vector satisfying the FONC. Clearly, since  $\frac{1}{2} \mathbf{x}^{*\top} \mathbf{Q} \mathbf{x}^*$  is a constant, and  $\mathbf{Q} > 0$ , then

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) = -\frac{1}{2} \mathbf{x}^{*\top} \mathbf{Q} \mathbf{x}^*,$$

and  $f(\mathbf{x}) = f(\mathbf{x}^*)$  if and only if  $\mathbf{x} = \mathbf{x}^*$ .

### 6.34

Write  $\mathbf{u} = [u_1, \dots, u_n]$ . We have

$$\begin{aligned} x_n &= ax_{n-1} + bu_n \\ &= a(ax_{n-2} + bu_{n-1}) + bu_n \\ &= a^2x_{n-2} + ab u_{n-1} + bu_n \\ &\vdots \\ &= a^n x_0 + a^{n-1}bu_1 + \dots + ab u_{n-1} + bu_n \\ &= \mathbf{c}^\top \mathbf{u}, \end{aligned}$$

where  $\mathbf{c} = [a^{n-1}b, \dots, ab, b]^\top$ . Therefore, the problem can be written as

$$\text{minimize } r\mathbf{u}^\top \mathbf{u} - q\mathbf{c}^\top \mathbf{u},$$

which is a positive definite quadratic in  $\mathbf{u}$ . The solution is therefore

$$\mathbf{u} = \frac{q}{2r}\mathbf{c},$$

or, equivalently,  $u_i = qa^{n-i}b/(2r)$ ,  $i = 1, \dots, n$ .

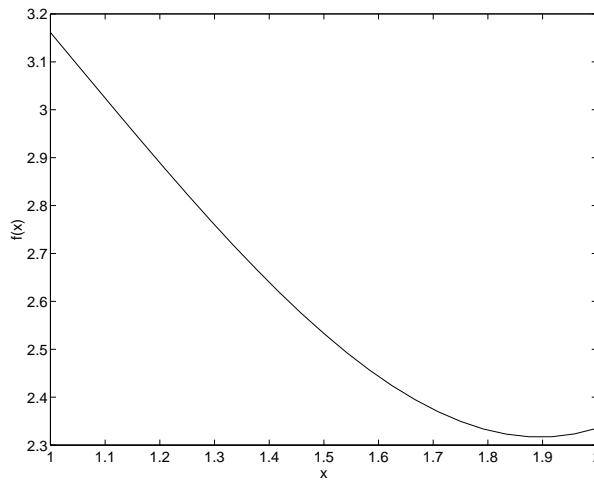
## 7. One Dimensional Search Methods

### 7.1

The range reduction factor for 3 iterations of the Golden Section method is  $((\sqrt{5}-1)/2)^3 = 0.236$ , while that of the Fibonacci method (with  $\varepsilon = 0$ ) is  $1/F_{3+1} = 0.2$ . Hence, if the desired range reduction factor is anywhere between 0.2 and 0.236 (e.g., 0.21), then the Golden Section method requires at least 4 iterations, while the Fibonacci method requires only 3. So, an example of a desired final uncertainty range is  $0.21 \times (8-5) = 0.63$ .

### 7.2

a. The plot of  $f(x)$  versus  $x$  is as below:



b. The number of steps needed for the Golden Section method is computed from the inequality:

$$0.61803^N \leq \frac{0.2}{2-1} \quad \Rightarrow \quad N \geq 3.34.$$

Therefore, the fewest possible number of steps is 4. Applying 4 steps of the Golden Section method, we end up with an uncertainty interval of  $[a_4, b_0] = [1.8541, 2.000]$ . The table with the results of the intermediate steps is displayed below:

Iteration $k$	$a_k$	$b_k$	$f(a_k)$	$f(b_k)$	New uncertainty interval
1	1.3820	1.6180	2.6607	2.4292	[1.3820,2]
2	1.6180	1.7639	2.4292	2.3437	[1.6180,2]
3	1.7639	1.8541	2.3437	2.3196	[1.7639,2]
4	1.8541	1.9098	2.3196	2.3171	[1.8541,2]

c. The number of steps needed for the Fibonacci method is computed from the inequality:

$$\frac{1 + 2\varepsilon}{F_{N+1}} \leq \frac{0.2}{2 - 1} \quad \Rightarrow \quad N \geq 4.$$

Therefore, the fewest possible number of steps is 4. Applying 4 steps of the Fibonacci method, we end up with an uncertainty interval of  $[a_4, b_0] = [1.8750, 2.000]$ . The table with the results of the intermediate steps is displayed below:

Iteration $k$	$\rho_k$	$a_k$	$b_k$	$f(a_k)$	$f(b_k)$	New unc. int.
1	0.3750	1.3750	1.6250	2.6688	2.4239	[1.3750,2]
2	0.4	1.6250	1.7500	2.4239	2.3495	[1.6250,2]
3	0.3333	1.7500	1.8750	2.3495	2.3175	[1.7500,2]
4	0.45	1.8750	1.8875	2.3175	2.3169	[1.8750,2]

d. We have  $f'(x) = 2x - 4 \sin x$ ,  $f''(x) = 2 - 4 \cos x$ . Hence, Newton's algorithm takes the form:

$$x^{(k+1)} = x^{(k)} - \frac{x^{(k)} - 2 \sin x^{(k)}}{1 - 2 \cos x^{(k)}}.$$

Applying 4 iterations with  $x^{(0)} = 1$ , we get  $x^{(1)} = -7.4727$ ,  $x^{(2)} = 14.4785$ ,  $x^{(3)} = 6.9351$ ,  $x^{(4)} = 16.6354$ . Apparently, Newton's method is not effective in this case.

### 7.3

---

a. We first create the M-file **f.m** as follows:

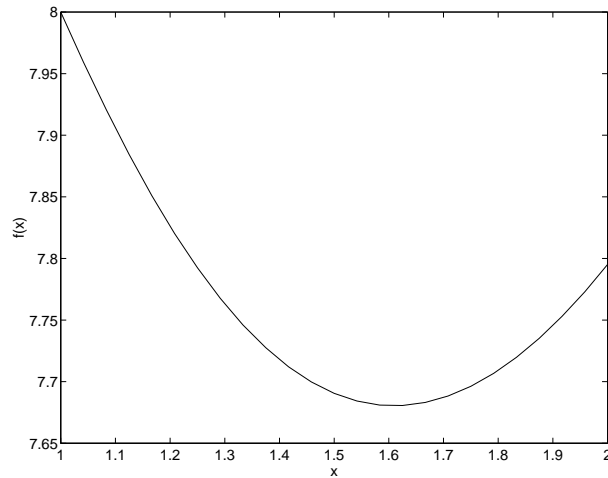
```
% f.m
function y=f(x)
y=8*exp(1-x)+7*log(x);
```

The MATLAB commands to plot the function are:

```
fplot('f',[1 2]);
xlabel('x');
ylabel('f(x)');
```

The resulting plot is as follows:





b. The MATLAB routine for the Golden Section method is:

```
%Matlab routine for Golden Section Search
```

```
left=1;
right=2;
uncert=0.23;
```

```
rho=(3-sqrt(5))/2;
```

```
N=ceil(log(uncert/(right-left))/log(1-rho)) %print N
```

```
lower='a';
a=left+(1-rho)*(right-left);
f_a=f(a);
for i=1:N,
    if lower=='a'
        b=a
        f_b=f_a
        a=left+rho*(right-left)
        f_a=f(a)
    else
        a=b
        f_a=f_b
        b=left+(1-rho)*(right-left)
        f_b=f(b)
    end %if
    if f_a<f_b
        right=b;
        lower='a'
    else
        left=a;
        lower='b'
    end %if
    New_Interval = [left,right]
end %for i
%-----
```

Using the above routine, we obtain  $N = 4$  and a final interval of  $[1.528, 1.674]$ . The table with the results of the intermediate steps is displayed below:

Iteration $k$	$a_k$	$b_k$	$f(a_k)$	$f(b_k)$	New uncertainty interval
1	1.382	1.618	7.7247	7.6805	[1.382,2]
2	1.618	1.764	7.6805	7.6995	[1.382,1.764]
3	1.528	1.618	7.6860	7.6805	[1.528,1.764]
4	1.618	1.674	7.6805	7.6838	[1.528,1.674]

c. The MATLAB routine for the Fibonacci method is:

```
%Matlab routine for Fibonacci Search technique
```

```
left=1;
right=2;
uncert=0.23;
epsilon=0.05;
```

```
F(1)=1;
F(2)=1;
```

```
N=0;
while F(N+2) < (1+2*epsilon)*(right-left)/uncert
    F(N+3)=F(N+2)+F(N+1);
    N=N+1;
end %while
```

```
N %print N
```

```
lower='a';
a=left+(F(N+1)/F(N+2))*(right-left);
f_a=f(a);
```

```
for i=1:N,
    if i~=N
        rho=1-F(N+2-i)/F(N+3-i)
    else
        rho=0.5-epsilon
    end %if
    if lower=='a'
        b=a
        f_b=f_a
        a=left+rho*(right-left)
        f_a=f(a)
    else
        a=b
        f_a=f_b
        b=left+(1-rho)*(right-left)
        f_b=f(b)
    end %if
    if f_a<f_b
        right=b;
        lower='a'
    else
        left=a;
        lower='b'
    end %if
    New_Interval = [left,right]
end %for i
%-----
```

Using the above routine, we obtain  $N = 3$  and a final interval of  $[1.58, 1.8]$ . The table with the results of the intermediate steps is displayed below:

Iteration $k$	$\rho_k$	$a_k$	$b_k$	$f(a_k)$	$f(b_k)$	New uncertainty interval
1	0.4	1.4	1.6	7.7179	7.6805	$[1.4, 2]$
2	0.333	1.6	1.8	7.6805	7.7091	$[1.4, 1.8]$
3	0.45	1.58	1.6	7.6812	7.6805	$[1.58, 1.8]$

#### 7.4

Now,  $\rho_k = 1 - \frac{F_{N-k+1}}{F_{N-k+2}}$ . Hence,

$$\begin{aligned}
 1 - \frac{\rho_k}{1 - \rho_k} &= 1 - \frac{1 - F_{N-k+1}/F_{N-k+2}}{F_{N-k+1}/F_{N-k+2}} \\
 &= 1 - \frac{F_{N-k+2} - F_{N-k+1}}{F_{N-k+1}} \\
 &= 1 - \frac{F_{N-k}}{F_{N-k+1}} \\
 &= \rho_{k+1}
 \end{aligned}$$

To show that  $0 \leq \rho_k \leq 1/2$ , we proceed by induction. Clearly  $\rho_1 = 1/2$  satisfies  $0 \leq \rho_1 \leq 1/2$ . Suppose  $0 \leq \rho_k \leq 1/2$ , where  $k \in \{1, \dots, N-1\}$ . Then,

$$\frac{1}{2} \leq 1 - \rho_k \leq 1$$

and hence

$$1 \leq \frac{1}{1 - \rho_k} \leq 2.$$

Therefore,

$$\frac{1}{2} \leq \frac{\rho_k}{1 - \rho_k} \leq 1.$$

Since  $\rho_{k+1} = 1 - \frac{\rho_k}{1 - \rho_k}$ , then

$$0 \leq \rho_{k+1} \leq \frac{1}{2}$$

as required.

#### 7.5

We proceed by induction. For  $k = 2$ , we have  $F_0 F_3 - F_1 F_2 = (1)(3) - (1)(2) = 1 = (-1)^2$ . Suppose  $F_{k-2} F_{k+1} - F_{k-1} F_k = (-1)^k$ . Then,

$$\begin{aligned}
 F_{k-1} F_{k+2} - F_k F_{k+1} &= F_{k-1} (F_{k+1} + F_k) - (F_{k-1} + F_{k-2}) F_{k+1} \\
 &= F_{k-1} F_k - F_{k-2} F_{k+1} \\
 &= -(-1)^k \\
 &= (-1)^{k+1}.
 \end{aligned}$$

#### 7.6

Define  $y_k = F_k$  and  $z_k = F_{k-1}$ . Then, we have

$$\begin{bmatrix} y_{k+1} \\ z_{k+1} \end{bmatrix} = \mathbf{A} \begin{bmatrix} y_k \\ z_k \end{bmatrix},$$

where

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix},$$

with initial condition

$$\begin{bmatrix} y_0 \\ z_0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

We can write

$$F_n = y_n = [1, 0] \begin{bmatrix} y_n \\ z_n \end{bmatrix} = [1, 0] \mathbf{A}^n \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

Since  $\mathbf{A}$  is symmetric, it can be diagonalized as

$$\mathbf{A} = \begin{bmatrix} \mathbf{u}^\top \\ \mathbf{v}^\top \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \begin{bmatrix} \mathbf{u} & \mathbf{v} \end{bmatrix}$$

where

$$\lambda_1 = \frac{1 + \sqrt{5}}{2}, \quad \lambda_2 = \frac{1 - \sqrt{5}}{2},$$

and

$$\mathbf{u} = -\frac{1}{5^{1/4}} \begin{bmatrix} \sqrt{2/(\sqrt{5}-1)} \\ \sqrt{(\sqrt{5}-1)/2} \end{bmatrix}, \quad \mathbf{v} = \frac{1}{5^{1/4}} \begin{bmatrix} \sqrt{(\sqrt{5}-1)/2} \\ -\sqrt{2/(\sqrt{5}-1)} \end{bmatrix}.$$

Therefore, we have

$$\begin{aligned} F_n &= \mathbf{u}^\top \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}^n \mathbf{u} \\ &= u_1^2 \lambda_1^n + u_2^2 \lambda_2^n \\ &= \frac{1}{\sqrt{5}} \left( \left( \frac{1 + \sqrt{5}}{2} \right)^{n+1} - \left( \frac{1 - \sqrt{5}}{2} \right)^{n+1} \right). \end{aligned}$$

## 7.7

The number  $\log(2)$  is the root of the equation  $g(x) = 0$ , where  $g(x) = \exp(x) - 2$ . The derivative of  $g$  is  $g'(x) = \exp(x)$ . Newton's method applied to this root finding problem is

$$x^{(k+1)} = x^{(k)} - \frac{\exp(x^{(k)}) - 2}{\exp(x^{(k)})} = x^{(k)} - 1 + 2\exp(-x^{(k)}).$$

Performing two iterations, we get  $x^{(1)} = 0.7358$  and  $x^{(2)} = 0.6940$ .

## 7.8

a. We compute  $g'(x) = 2e^x/(e^x + 1)^2$ . Therefore Newton's method of tangents for this problem takes the form

$$\begin{aligned} x^{(k+1)} &= x^{(k)} - \frac{(e^{x^{(k)}} - 1)/(e^{x^{(k)}} + 1)}{2e^{x^{(k)}}/(e^{x^{(k)}} + 1)^2} \\ &= x^{(k)} - \frac{(e^{2x^{(k)}} - 1)}{2e^{x^{(k)}}} \\ &= x^{(k)} - \sinh x^{(k)}. \end{aligned}$$

b. By symmetry, we need  $x^{(1)} = -x^{(0)}$  for cycling. Therefore,  $x^{(0)}$  must satisfy

$$-x^{(0)} = x^{(0)} - \sinh x^{(0)}.$$

The algorithm cycles if  $x^{(0)} = c$ , where  $c > 0$  is the solution to  $2c = \sinh c$ .

c. The algorithm converges to 0 if and only if  $|x^{(0)}| < c$ , where  $c$  is from part b.

### 7.9

The quadratic function that matches the given data  $x^{(k)}$ ,  $x^{(k-1)}$ ,  $x^{(k-2)}$ ,  $f(x^{(k)})$ ,  $f(x^{(k-1)})$ , and  $f(x^{(k-2)})$  can be computed by solving the following three linear equations for the parameters  $a$ ,  $b$ , and  $c$ :

$$a(x^{(k-i)})^2 + bx^{(k-i)} + c = f(x^{(k-i)}), \quad i = 0, 1, 2.$$

Then, the algorithm is given by  $x^{(k+1)} = -b/2a$  (so, in fact, we only need to find the *ratio* of  $a$  and  $b$ ). With some elementary algebra (e.g., using Cramer's rule without needing to calculate the determinant in the denominator), the algorithm can be written as:

$$x^{(k+1)} = \frac{\sigma_{12}f(x^{(k)}) + \sigma_{20}f(x^{(k-1)}) + \sigma_{01}f(x^{(k-2)})}{2(\delta_{12}f(x^{(k)}) + \delta_{20}f(x^{(k-1)}) + \delta_{01}f(x^{(k-2)}))}$$

where  $\sigma_{ij} = (x^{(k-i)})^2 - (x^{(k-j)})^2$  and  $\delta_{ij} = x^{(k-i)} - x^{(k-j)}$ .

### 7.10

a. A MATLAB routine for implementing the secant method is as follows.

```
function [x,v] = secant(g,xcurr,xnew,uncert);
%Matlab routine for finding root of g(x) using secant method
%
%  secant;
%  secant(g);
%  secant(g,xcurr,xnew);
%  secant(g,xcurr,xnew,uncert);
%
%  x=secant;
%  x=secant(g);
%  x=secant(g,xcurr,xnew);
%  x=secant(g,xcurr,xnew,uncert);
%
%  [x,v]=secant;
%  [x,v]=secant(g);
%  [x,v]=secant(g,xcurr,xnew);
%  [x,v]=secant(g,xcurr,xnew,uncert);
%
%The first variant finds the root of g(x) in the M-file g.m, with
%initial conditions 0 and 1, and uncertainty 10^(-5).
%The second variant finds the root of the function in the M-file specified
%by the string g, with initial conditions 0 and 1, and uncertainty 10^(-5).
%The third variant finds the root of the function in the M-file specified
%by the string g, with initial conditions specified by xcurr and xnew, and
%uncertainty 10^(-5).
%The fourth variant finds the root of the function in the M-file specified
%by the string g, with initial conditions specified by xcurr and xnew, and
%uncertainty specified by uncert.
%
%The next four variants returns the final value of the root as x.
%The last four variants returns the final value of the root as x, and
%the value of the function at the final value as v.

if nargin < 4
    uncert=10^(-5);
    if nargin < 3
        if nargin == 1
            xcurr=0;
            xnew=1;
        elseif nargin == 0
```

```

        g='g';
    else
        disp('Cannot have 2 arguments. ');
        return;
    end
end
end
end

g_curr=feval(g,xcurr);

while abs(xnew-xcurr)>xcurr*uncert,
    xold=xcurr;
    xcurr=xnew;
    g_old=g_curr;
    g_curr=feval(g,xcurr);
    xnew=(g_curr*xold-g_old*xcurr)/(g_curr-g_old);
end %while

%print out solution and value of g(x)
if nargout >= 1
    x=xnew;
    if nargout == 2
        v=feval(g,xnew);
    end
else
    final_point=xnew
    value=feval(g,xnew)
end %if
%-----

```

b. We get a solution of  $x = 0.0039671$ , with corresponding value  $g(x) = -9.908 \times 10^{-8}$ .

## 7.11

---

```

function alpha=linesearch_secant(grad,x,d)
%Line search using secant method

epsilon=10^(-4); %line search tolerance
max = 100; %maximum number of iterations
alpha_curr=0;
alpha=0.001;
dphi_zero=feval(grad,x)*d;
dphi_curr=dphi_zero;

i=0;
while abs(dphi_curr)>epsilon*abs(dphi_zero),
    alpha_old=alpha_curr;
    alpha_curr=alpha;
    dphi_old=dphi_curr;
    dphi_curr=feval(grad,x+alpha_curr*d)*d;
    alpha=(dphi_curr*alpha_old-dphi_old*alpha_curr)/(dphi_curr-dphi_old);
    i=i+1;
    if (i >= max) & (abs(dphi_curr)>epsilon*abs(dphi_zero)),
        disp('Line search terminating with number of iterations:');
        disp(i);
        break;
    end
end %while
%-----

```

## 7.12

a. We could carry out the bracketing using the one-dimensional function  $\phi_0(\alpha) = f(\mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)})$ , where  $\mathbf{d}^{(0)}$  is the negative gradient at  $\mathbf{x}^{(0)}$ , as described in Section 7.8. The decision variable would be  $\alpha$ . However, here we will directly represent the points in  $\mathbb{R}^2$  (which is equivalent, though unnecessary in general).

The uncertainty interval is calculated by the following procedure:

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \mathbf{x}, \quad \nabla f(\mathbf{x}) = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \mathbf{x}$$

Therefore,

$$\mathbf{d} = -\nabla f(\mathbf{x}^{(0)}) = - \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 0.8 \\ -0.25 \end{bmatrix} = \begin{bmatrix} -1.35 \\ -0.3 \end{bmatrix}$$

As the problem requires, we use  $\varepsilon = 0.075$ .

First, we begin calculating  $f(\mathbf{x}^{(0)})$  and  $\mathbf{x}^{(1)}$ :

$$f(\mathbf{x}^{(0)}) = f \left( \begin{bmatrix} 0.8 \\ -0.25 \end{bmatrix} \right) = 0.5025,$$

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \varepsilon \mathbf{d} = \begin{bmatrix} 0.8 \\ -0.25 \end{bmatrix} + 0.075 \begin{bmatrix} -1.35 \\ -0.3 \end{bmatrix} = \begin{bmatrix} 0.6987 \\ -0.2725 \end{bmatrix}$$

Then, we proceed as follows to find the uncertainty interval:

$$f(\mathbf{x}^{(1)}) = f \left( \begin{bmatrix} 0.6987 \\ -0.2725 \end{bmatrix} \right) = 0.3721$$

$$\mathbf{x}^{(2)} = \mathbf{x}^{(0)} + 2\varepsilon \mathbf{d} = \begin{bmatrix} 0.8 \\ -0.25 \end{bmatrix} + 0.15 \begin{bmatrix} -1.35 \\ -0.3 \end{bmatrix} = \begin{bmatrix} 0.5975 \\ -0.2950 \end{bmatrix}$$

$$f(\mathbf{x}^{(2)}) = f \left( \begin{bmatrix} 0.5975 \\ -0.2950 \end{bmatrix} \right) = 0.2678$$

$$\mathbf{x}^{(3)} = \mathbf{x}^{(0)} + 4\varepsilon \mathbf{d} = \begin{bmatrix} 0.8 \\ -0.25 \end{bmatrix} + 0.3 \begin{bmatrix} -1.35 \\ -0.3 \end{bmatrix} = \begin{bmatrix} 0.3950 \\ -0.3400 \end{bmatrix}$$

$$f(\mathbf{x}^{(3)}) = f \left( \begin{bmatrix} 0.3950 \\ -0.3400 \end{bmatrix} \right) = 0.1373$$

$$\mathbf{x}^{(4)} = \mathbf{x}^{(0)} + 8\varepsilon \mathbf{d} = \begin{bmatrix} 0.8 \\ -0.25 \end{bmatrix} + 0.6 \begin{bmatrix} -1.35 \\ -0.3 \end{bmatrix} = \begin{bmatrix} -0.0100 \\ -0.4300 \end{bmatrix}$$

$$f(\mathbf{x}^{(4)}) = f \left( \begin{bmatrix} -0.0100 \\ -0.4300 \end{bmatrix} \right) = 0.1893$$

Between  $f(\mathbf{x}^{(3)})$  and  $f(\mathbf{x}^{(4)})$  the function increases, which means that the minimizer must occur on the interval  $[\mathbf{x}^{(2)}, \mathbf{x}^{(4)}] = \left[ \begin{bmatrix} 0.5975 \\ -0.2950 \end{bmatrix}, \begin{bmatrix} -0.0100 \\ -0.4300 \end{bmatrix} \right]$ , with  $\mathbf{d} = \begin{bmatrix} -1.35 \\ -0.3 \end{bmatrix}$ .

MATLAB code to solve the problem is listed next.

```
% Coded by David Schwartzman
```

```

% In our case we have:
Q = [2 1; 1 2];
x0 = [0.8; -0.25];
e = 0.075;
f = zeros(1,10);
X = zeros(2,10);
x1 = x0;
d = -Q*x1;
f(1) = 0.5*(x1')*Q*x1;
for i=2:10
    X(:,i) = x1+e*d;
    f(i) = 0.5*X(:,i)'*Q*X(:,i);
    e = 2*e;
    if(f(i) > f(i-1))
        break;
    end
end
% The interval is defined by:
a = X(:,i-2);
b = X(:,i);
str = sprintf('The minimizer is located in: [a, b], where a = [%.4f; %.4f]...
    and b = [%.4f; %.4f]', a(1,1), a(2,1), b(1,1), b(2,1));
disp(str);

```

b. First, we determine the number of necessary iterations:

The initial uncertainty interval width is 0.6223. This width will be  $(0.6223)(0.618)^N$  after  $N$  stages. We choose  $N$  so that

$$(0.618)^N \leq \frac{0.01}{0.6223} \quad \Rightarrow N = 9$$

We show the first iteration of the algorithm; the rest are analogous and shown in the following table. From part a, we know that  $[a_0, b_0] = [x^{(2)}, x^{(4)}]$ , then:

$$[a_0, b_0] = \left[ \begin{bmatrix} 0.5975 \\ -0.2950 \end{bmatrix}, \begin{bmatrix} -0.0100 \\ -0.4300 \end{bmatrix} \right], \text{ with } f(a_0) = 0.2678, f(b_0) = 0.1893.$$

$$\begin{aligned}
 a_1 &= a_0 + \rho(b_0 - a_0) = [0.3655 \quad -0.3466]^T \\
 b_1 &= a_0 + (1 - \rho)(b_0 - a_0) = [0.2220 \quad -0.3784]^T \\
 f(a_1) &= 0.1270 \\
 f(b_1) &= 0.1085
 \end{aligned}$$

We can see  $f(a_1) > f(b_1)$ , hence the uncertainty interval is reduced to:

$$[a_1, b_1] = \left[ \begin{bmatrix} 0.3655 \\ -0.3466 \end{bmatrix}, \begin{bmatrix} -0.0100 \\ -0.4300 \end{bmatrix} \right]$$

So, calculating the norm of  $(b_0 - a_1)$ , we see that the uncertainty region width is now 0.38461.



Iteration	$a_k$	$b_k$	$f(a_k)$	$f(b_k)$	New Uncertainty Interval
1	$\begin{bmatrix} 0.3655 \\ -0.3466 \end{bmatrix}$	$\begin{bmatrix} 0.2220 \\ -0.3784 \end{bmatrix}$	0.1270	0.1085	$\left[ \begin{bmatrix} .3655 \\ -0.3466 \end{bmatrix}, \begin{bmatrix} 0.0100 \\ -0.4300 \end{bmatrix} \right]$
2	$\begin{bmatrix} 0.2220 \\ -0.3784 \end{bmatrix}$	$\begin{bmatrix} 0.1334 \\ -0.3981 \end{bmatrix}$	0.1085	0.1232	$\left[ \begin{bmatrix} .3655 \\ -0.3466 \end{bmatrix}, \begin{bmatrix} .1334 \\ -0.3981 \end{bmatrix} \right]$
3	$\begin{bmatrix} 0.2768 \\ -0.3663 \end{bmatrix}$	$\begin{bmatrix} 0.2220 \\ -0.3784 \end{bmatrix}$	0.1094	0.1085	$\left[ \begin{bmatrix} .2768 \\ -0.3663 \end{bmatrix}, \begin{bmatrix} .1334 \\ -0.3981 \end{bmatrix} \right]$
4	$\begin{bmatrix} 0.2220 \\ -0.3784 \end{bmatrix}$	$\begin{bmatrix} 0.1882 \\ -0.3860 \end{bmatrix}$	0.1085	0.1117	$\left[ \begin{bmatrix} .2768 \\ -0.3663 \end{bmatrix}, \begin{bmatrix} .1882 \\ -0.3860 \end{bmatrix} \right]$
5	$\begin{bmatrix} 0.2430 \\ -0.3738 \end{bmatrix}$	$\begin{bmatrix} 0.2220 \\ -0.3784 \end{bmatrix}$	0.1079	0.1085	$\left[ \begin{bmatrix} .2768 \\ -0.3663 \end{bmatrix}, \begin{bmatrix} .2220 \\ -0.3784 \end{bmatrix} \right]$
6	$\begin{bmatrix} 0.2559 \\ -0.3709 \end{bmatrix}$	$\begin{bmatrix} 0.2430 \\ -0.3738 \end{bmatrix}$	0.1081	0.1079	$\left[ \begin{bmatrix} 0.2559 \\ -0.3709 \end{bmatrix}, \begin{bmatrix} 0.2220 \\ -0.3784 \end{bmatrix} \right]$
7	$\begin{bmatrix} 0.2430 \\ -0.3738 \end{bmatrix}$	$\begin{bmatrix} 0.2350 \\ -0.3756 \end{bmatrix}$	0.1079	0.1080	$\left[ \begin{bmatrix} 0.2559 \\ -0.3709 \end{bmatrix}, \begin{bmatrix} 0.2350 \\ -0.3756 \end{bmatrix} \right]$
8	$\begin{bmatrix} 0.2479 \\ -0.3727 \end{bmatrix}$	$\begin{bmatrix} 0.2430 \\ -0.3738 \end{bmatrix}$	0.1080	0.1079	$\left[ \begin{bmatrix} 0.2479 \\ -0.3727 \end{bmatrix}, \begin{bmatrix} 0.2350 \\ -0.3756 \end{bmatrix} \right]$
9	$\begin{bmatrix} 0.2430 \\ -0.3738 \end{bmatrix}$	$\begin{bmatrix} 0.2399 \\ -0.3745 \end{bmatrix}$	0.1079	0.1079	$\left[ \begin{bmatrix} 0.2479 \\ -0.3727 \end{bmatrix}, \begin{bmatrix} 0.2399 \\ -0.3745 \end{bmatrix} \right]$

We can now see that the minimizer is located within  $\left[ \begin{bmatrix} 0.2479 \\ -0.3727 \end{bmatrix}, \begin{bmatrix} 0.2399 \\ -0.3745 \end{bmatrix} \right]$ , and its uncertainty interval width is 0.00819.

Matlab code used to perform calculations is listed next

```
% Coded by David Schwartzman
% To succesfully run this program, run
% the previous script to obtain a and b.
e = 0.01;
Q = [2 1; 1 2];
ro = 0.5*(3-sqrt(5));
% First we determine the number of necessary iterations.
d = norm(a-b); N = ceil( (log(e/d))/(log(1-ro)));
fa = 0.5*a'*Q*a;
fb = 0.5*b'*Q*b;
str1 = sprintf('Initial values:');
str2 = sprintf('a0 = [%f,%f].', a(1), a(2));
str3 = sprintf('b0 = [%f,%f].', b(1), b(2));
str4 = sprintf('f(a0) = %f.', fa);
str5 = sprintf('f(b0) = %f.', fb);
```

```

strn = sprintf('\n');
disp(strn);
disp(str1);disp(str2);
disp(str3);disp(str4);
disp(str5);
s = a + ro*(b-a);
t = a + (1-ro)*(b-a);
fs = 0.5*s'*Q*s;
ft = 0.5*t'*Q*t;
for i=1:N
    str1 = sprintf('Iteration number: %d', i);
    str2 = sprintf('a%d = [%.4f,%.4f].', i, s(1), s(2));
    str3 = sprintf('b%d = [%.4f,%.4f].', i, t(1), t(2));
    str4 = sprintf('f(a%d) = %.4f.', i, fs);
    str5 = sprintf('f(b%d) = %.4f.', i, ft);
    if (ft>fs)
        b = t;
        %fb = ft;
        t = s;
        ft = fs;
        s = a + ro*(b-a);
        fs = 0.5*s'*Q*s;
    else
        a = s;
        %fa = fs;
        s = t;
        fs = ft;
        t = a + (1-ro)*(b-a);
        ft = 0.5*t'*Q*t;
    end
    str6 = sprintf('New uncertainty interval: a%d = [%.4f,%.4f], ...
b%d = [%.4f,%.4f].', i, a(1), a(2), i, b(1), b(2));
    disp(strn);
    disp(str1)
    disp(str2)
    disp(str3)
    disp(str4)
    disp(str5)
    disp(str6)
end
% The interval where the minimizer is boxed in is given by:
an = a;
bn = b;
% We can return (an+bn)/2 as the minimizer.
min = (an+bn)/2;disp(strn);
str = sprintf('The minimizer x* is: [%.4f; %.4f]', min(1,1), min(2,1));
disp(str);

```

c. We need to determine the number of necessary iterations:

The initial uncertainty interval width is 0.6223. This width will be  $(0.6223)^{\frac{1+2\varepsilon}{F_{(N+1)}}}$ , where  $F_k$  is the k-th element of the Fibonacci sequence. We choose  $N$  so that

$$\frac{1+2\varepsilon}{F_{N+1}} \leq \frac{0.01}{0.6223} = 0.0161 \quad \Rightarrow F_{N+1} \geq \frac{1+2\varepsilon}{0.0161}$$

For  $\varepsilon = 0.05$ , we require  $F_{(N+1)} > 68.32$ , thus  $F_{(10)} = 89$  is enough, and we have  $N = 10 - 1$ , 9 iterations.

We show the first iteration of the algorithm; the rest are analogous and shown in the following table. From part a, we know that  $[a_0, b_0] = [x^{(2)}, x^{(4)}]$ , then:

$$[a_0, b_0] = \left[ \begin{bmatrix} 0.5975 \\ -0.2950 \end{bmatrix}, \begin{bmatrix} -0.0100 \\ -0.4300 \end{bmatrix} \right], \text{ with } f(a_0) = 0.2678, f(b_0) = 0.1893.$$

Recall that in the Fibonacci method,  $\rho_1 = 1 - \frac{F_N}{F_{N+1}} = 1 - \frac{55}{89} = 0.3820$ .

$$\begin{aligned} a_1 &= a_0 + \rho_1(b_0 - a_0) = [0.3654 \quad -0.3466]^T \\ b_1 &= a_0 + (1 - \rho_1)(b_0 - a_0) = [0.2221 \quad -0.3784]^T \\ f(a_1) &= 0.1270 \\ f(b_1) &= 0.1085 \end{aligned}$$

We can see  $f(a_1) > f(b_1)$ , hence the uncertainty interval is reduced to:

$$[a_1, b_0] = \left[ \begin{bmatrix} 0.3654 \\ -0.3466 \end{bmatrix}, \begin{bmatrix} -0.0100 \\ -0.4300 \end{bmatrix} \right]$$

So, calculating the norm of  $(b_0 - a_1)$ , we see that the uncertainty region width is now 0.38458.

$k$	$\rho_k$	$a_k$	$b_k$	$f(a_k)$	$f(b_k)$	New Uncertainty Interval
1	0.3820	$\begin{bmatrix} 0.3654 \\ -0.3466 \end{bmatrix}$	$\begin{bmatrix} 0.2221 \\ -0.3784 \end{bmatrix}$	0.1270	0.1085	$\left[ \begin{bmatrix} 0.3654 \\ -0.3466 \end{bmatrix}, \begin{bmatrix} -0.0100 \\ -0.4300 \end{bmatrix} \right]$
2	0.3818	$\begin{bmatrix} 0.2221 \\ -0.3784 \end{bmatrix}$	$\begin{bmatrix} 0.1333 \\ -0.3981 \end{bmatrix}$	0.1085	0.1232	$\left[ \begin{bmatrix} 0.3654 \\ -0.3466 \end{bmatrix}, \begin{bmatrix} 0.1333 \\ -0.3981 \end{bmatrix} \right]$
3	0.3824	$\begin{bmatrix} 0.2767 \\ -0.3663 \end{bmatrix}$	$\begin{bmatrix} 0.2221 \\ -0.3784 \end{bmatrix}$	0.1094	0.1085	$\left[ \begin{bmatrix} 0.2767 \\ -0.3663 \end{bmatrix}, \begin{bmatrix} 0.1333 \\ -0.3981 \end{bmatrix} \right]$

$k$	$\rho_k$	$a_k$	$b_k$	$f(a_k)$	$f(b_k)$	New Uncertainty Interval
4	0.3810	$\begin{bmatrix} 0.2221 \\ -0.3784 \end{bmatrix}$	$\begin{bmatrix} 0.1879 \\ -0.3860 \end{bmatrix}$	0.1085	0.1118	$\begin{bmatrix} \begin{bmatrix} 0.2767 \\ -0.3663 \end{bmatrix}, \begin{bmatrix} 0.1879 \\ -0.3860 \end{bmatrix} \end{bmatrix}$
5	0.3846	$\begin{bmatrix} 0.2426 \\ -0.3739 \end{bmatrix}$	$\begin{bmatrix} 0.2221 \\ -0.3784 \end{bmatrix}$	0.1079	0.1085	$\begin{bmatrix} \begin{bmatrix} 0.2767 \\ -0.3663 \end{bmatrix}, \begin{bmatrix} 0.2221 \\ -0.3784 \end{bmatrix} \end{bmatrix}$
6	0.3750	$\begin{bmatrix} 0.2562 \\ -0.3708 \end{bmatrix}$	$\begin{bmatrix} 0.2426 \\ -0.3739 \end{bmatrix}$	0.1082	0.1079	$\begin{bmatrix} \begin{bmatrix} 0.2562 \\ -0.3708 \end{bmatrix}, \begin{bmatrix} 0.2221 \\ -0.3784 \end{bmatrix} \end{bmatrix}$
7	0.4000	$\begin{bmatrix} 0.2426 \\ -0.3739 \end{bmatrix}$	$\begin{bmatrix} 0.2357 \\ -0.3754 \end{bmatrix}$	0.1079	0.1080	$\begin{bmatrix} \begin{bmatrix} 0.2562 \\ -0.3708 \end{bmatrix}, \begin{bmatrix} 0.2357 \\ -0.3754 \end{bmatrix} \end{bmatrix}$
8	0.3333	$\begin{bmatrix} 0.2494 \\ -0.3724 \end{bmatrix}$	$\begin{bmatrix} 0.2426 \\ -0.3739 \end{bmatrix}$	0.1080	0.1079	$\begin{bmatrix} \begin{bmatrix} 0.2494 \\ -0.3724 \end{bmatrix}, \begin{bmatrix} 0.2357 \\ -0.3754 \end{bmatrix} \end{bmatrix}$
9	0.4500	$\begin{bmatrix} 0.2426 \\ -0.3739 \end{bmatrix}$	$\begin{bmatrix} 0.2419 \\ -0.3740 \end{bmatrix}$	0.1079	0.1079	$\begin{bmatrix} \begin{bmatrix} 0.2494 \\ -0.3724 \end{bmatrix}, \begin{bmatrix} 0.2419 \\ -0.3740 \end{bmatrix} \end{bmatrix}$

We can now see that the minimizer is located within  $\begin{bmatrix} \begin{bmatrix} 0.2494 \\ -0.3724 \end{bmatrix}, \begin{bmatrix} 0.2419 \\ -0.3740 \end{bmatrix} \end{bmatrix}$ , and its uncertainty interval width is 0.00769.

Matlab copde used to perform calculations is listed next.

```
% Coded by David Schwartzman
% To succesfully run this program, run the first of the above scripts
% to obtain a and b.
% We take
e = 0.01;
Q = [2 1; 1 2];
% First determine the number of necessary iterations.
d = norm(a-b);
F_N_1 =(2*0.05+1)/(e/d);
F = zeros(1,20);
F(1) = 0;
F(2) = 1;
for i=1:20
    F(i+1) = F(i)+F(i+1);
    if(F(i+2)>= F_N_1)
        break;
    end
end
N = i-1;
ro = zeros(1, N+1);
for i = 1:N
    ro(i) = 1 - F(N+3-i)/F(N+4-i);
end
ro(N) = ro(N) - 0.05;
fa = 0.5*a'*Q*a;
fb = 0.5*b'*Q*b;
```

```

str1 = sprintf('Initial values:');
str2 = sprintf('a0 = [%.4f,%.4f].', a(1), a(2));
str3 = sprintf('b0 = [%.4f,%.4f].', b(1), b(2));
str4 = sprintf('f(a0) = %.4f.', fa);
str5 = sprintf('f(b0) = %.4f.', fb);
strn = sprintf('\n');
disp(strn);
disp(str1);
disp(str2);
disp(str3);
disp(str4);
disp(str5);
s = a + ro(1)*(b-a);
t = a + (1-ro(1))*(b-a);
fs = 0.5*s'*Q*s;
ft = 0.5*t'*Q*t;
for i=1:N
    str1 = sprintf('Iteration number: %d', i);
    str2 = sprintf('a%d = [%.4f,%.4f].', i, s(1), s(2));
    str3 = sprintf('b%d = [%.4f,%.4f].', i, t(1), t(2));
    str4 = sprintf('f(a%d) = %.4f.', i, fs);
    str5 = sprintf('f(b%d) = %.4f.', i, ft);
    if (ft>fs)
        b = t;
        t = s;
        ft = fs;
        s = a + ro(i+1)*(b-a);
        fs = 0.5*s'*Q*s;
    else
        a = s;
        s = t;
        fs = ft;
        t = a + (1-ro(i+1))*(b-a);
        ft = 0.5*t'*Q*t;
    end
    str6 = sprintf('New uncertainty interval: a%d = [%.4f,%.4f],...
        b%d = [%.4f,%.4f].', i, a(1), a(2), i, b(1), b(2));
    str7 = sprintf('Uncertainty interval width: %.5f', norm(a-b));
    disp(strn);
    disp(str1);
    disp(str2);
    disp(str3);
    disp(str4);
    disp(str5);
    disp(str6);
    disp(str7);
end
% The minimizer is boxed in the interval:
an = a;
bn = b;
% We can return (an+bn)/2 as the minimizer.
min = (an+bn)/2;
disp(strn);
str = sprintf('The minimizer x* is: [%.4f; %.4f]', min(1,1), min(2,1));

```

`disp(str);`

## 8. Gradient Methods

### 8.1

---

The function  $f$  is a quadratic and so we can represent it in standard form as

$$f = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{x} - \mathbf{x}^\top \begin{bmatrix} -1 \\ -1/2 \end{bmatrix} + 3 = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{x}^\top \mathbf{b} + c.$$

The first iteration is

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \alpha_0 \nabla f(\mathbf{x}^{(0)}).$$

To find  $\mathbf{x}^{(1)}$ , we need to compute  $\nabla f(\mathbf{x}^{(0)}) = \mathbf{g}^{(0)}$ . We have

$$\mathbf{g}^{(0)} = \mathbf{Q} \mathbf{x}^{(0)} - \mathbf{b} = \begin{bmatrix} 1 & 1/2 \end{bmatrix}^\top.$$

The step size,  $\alpha_0$ , can be computed as

$$\alpha_0 = \frac{\mathbf{g}^{(0)\top} \mathbf{g}^{(0)}}{\mathbf{g}^{(0)\top} \mathbf{Q} \mathbf{g}^{(0)}} = \frac{5}{6}.$$

Hence,

$$\mathbf{x}^{(1)} = -\alpha_0 \mathbf{g}^{(0)} = -\frac{5}{6} \begin{bmatrix} 1 \\ 1/2 \end{bmatrix} = -\begin{bmatrix} 5/6 \\ 5/12 \end{bmatrix}.$$

The second iteration is

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} - \alpha_1 \nabla f(\mathbf{x}^{(1)}),$$

where

$$\nabla f(\mathbf{x}^{(1)}) = \mathbf{g}^{(1)} = \mathbf{Q} \mathbf{x}^{(1)} - \mathbf{b} = \begin{bmatrix} 1/6 \\ -1/3 \end{bmatrix},$$

and

$$\alpha_1 = \frac{\mathbf{g}^{(1)\top} \mathbf{g}^{(1)}}{\mathbf{g}^{(1)\top} \mathbf{Q} \mathbf{g}^{(1)}} = \frac{5}{9}.$$

Hence,

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} - \alpha_1 \mathbf{g}^{(1)} = \begin{bmatrix} -5/6 \\ -5/12 \end{bmatrix} - \frac{5}{9} \begin{bmatrix} 1/6 \\ -1/3 \end{bmatrix} = \begin{bmatrix} -\frac{25}{27} \\ -\frac{25}{108} \end{bmatrix}.$$

The optimal solution is  $\mathbf{x}^* = [-1, -1/4]^\top$  obtained by solving the equation  $\mathbf{Q} \mathbf{x} = \mathbf{b}$ .

### 8.2

---

Let  $s$  be the order of convergence of  $\{\mathbf{x}^{(k)}\}$ . Suppose there exists  $c > 0$  such that for all  $k$  sufficiently large,

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \geq c \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p.$$

Hence, for all  $k$  sufficiently large,

$$\begin{aligned} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^s} &= \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p} \frac{1}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^{s-p}} \\ &\geq \frac{c}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^{s-p}}. \end{aligned}$$

Taking limits yields

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^s} \geq \frac{c}{\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^{s-p}}.$$

Since by definition  $s$  is the order of convergence,

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^s} < \infty.$$

Combining the above two inequalities, we get

$$\frac{c}{\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^{s-p}} < \infty.$$

Therefore, since  $\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\| = 0$ , we conclude that  $s \leq p$ , i.e., the order of convergence is at most  $p$ .

### 8.3

We use contradiction. Suppose  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$  and

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p} > 0$$

for some  $p < 1$ . We may assume that  $\mathbf{x}^{(k)} \neq \mathbf{x}^*$  for an infinite number of  $k$  (for otherwise, by convention, the ratio above is eventually 0). Fix  $\varepsilon > 0$ . Then, there exists  $K_1$  such that for all  $k \geq K_1$ ,

$$\frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p} > \varepsilon.$$

Dividing both sides by  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^{1-p}$ , we obtain

$$\frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|} > \frac{\varepsilon}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^{1-p}}.$$

Because  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$  and  $p < 1$ , we have  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^{1-p} \rightarrow 0$ . Hence, there exists  $K_2$  such that for all  $k \geq K_2$ ,  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^{1-p} < \varepsilon$ . Combining this inequality with the previous one yields

$$\frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|} > 1$$

for all  $k \geq \max(K_1, K_2)$ ; i.e.,

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| > \|\mathbf{x}^{(k)} - \mathbf{x}^*\|,$$

which contradicts the assumption that  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ .

### 8.4

a. The sequence converges to 0, because the exponent  $-2^{k^2}$  grows unboundedly negative as  $k \rightarrow \infty$ .

b. The order of convergence of  $\{x^{(k)}\}$  is  $\infty$ . To see this, we first write, for  $p \geq 1$ ,

$$\begin{aligned} \frac{|x^{(k+1)}|}{|x^{(k)}|^p} &= \frac{2^{-2^{(k+1)}^2}}{2^{-2^{k^2}p}} \\ &= 2^{-2^{k^2} + 2k+1 + p2^{k^2}} \\ &= 2^{-2^{k^2}(2^{2k+1} - p)}. \end{aligned}$$

But notice that the exponent  $-2^{k^2}(2^{2k+1} - p)$  grows unboundedly negative as  $k \rightarrow \infty$ , regardless of the value of  $p$ . Therefore, for any  $p$ ,

$$\lim_{k \rightarrow \infty} \frac{|x^{(k+1)}|}{|x^{(k)}|^p} = 0,$$

which means that the order of convergence is  $\infty$ .

## 8.5

---

a. We have

$$\begin{aligned}x^{(k)} &= ax^{(k-1)} \\&= a \cdot ax^{(k-2)} \\&= a^2 x^{(k-2)} \\&\vdots \\&= a^k x^{(0)}.\end{aligned}$$

Because  $0 < a < 1$ , we have  $a^k \rightarrow 0$ , and hence  $x^{(k)} \rightarrow 0$ .

b. Similarly, we have

$$\begin{aligned}y^{(k)} &= (y^{(k-1)})^b \\&= ((y^{(k-2)})^b)^b \\&= (y^{(k-2)})^{b^2} \\&\vdots \\&= (y^{(0)})^{b^k}.\end{aligned}$$

Because  $|y^{(0)}| < 1$  and  $b > 1$ , we have  $b^k \rightarrow \infty$  and hence  $y^{(k)} \rightarrow 0$ .

c. The order of convergence of  $\{x^{(k)}\}$  is 1 because

$$\lim_{k \rightarrow \infty} \frac{|x^{(k+1)}|}{|x^{(k)}|} = \lim_{k \rightarrow \infty} a = a,$$

and  $0 < a < \infty$ .

The order of convergence of  $\{y^{(k)}\}$  is  $b$  because

$$\lim_{k \rightarrow \infty} \frac{|y^{(k+1)}|}{|y^{(k)}|^b} = \lim_{k \rightarrow \infty} 1 = 1,$$

and  $0 < 1 < \infty$ .

d. Suppose  $|x^{(k)}| \leq c|x^{(0)}|$ . Using part a, we have  $a^k \leq c$ , which implies that  $k \geq \log(c)/\log(a)$ . So the smallest number of iterations  $k$  such that  $|x^{(k)}| \leq c|x^{(0)}|$  is  $\lceil \log(c)/\log(a) \rceil$  (the smallest integer not smaller than  $\log(c)/\log(a)$ ).

e. Suppose  $|y^{(k)}| \leq c|y^{(0)}|$ . Using part b, we have  $|y^{(0)}|^{b^k} \leq c|y^{(0)}|$ . Taking logs (twice) and rearranging, we have

$$k \geq \frac{1}{\log(b)} \log \left( 1 + \frac{\log(c)}{\log(|y^{(0)}|)} \right).$$

Denote the right-hand side by  $z$ . So the smallest number of iterations  $k$  such that  $|y^{(k)}| \leq c|y^{(0)}|$  is  $\lceil z \rceil$ .

f. Comparing the answer in part e with that of part d, we can see that as  $c \rightarrow 0$ , the answer in part d is  $\Omega(\log(c))$ , whereas the answer in part e is  $O(\log \log(c))$ . Hence, in the regime where  $c$  is very small, the number of iterations in part d (linear convergence) is (at least) exponentially larger than that in part e (superlinear convergence).

## 8.6

---

We have  $u_{k+1} = (1 - \rho)u_k$ , and  $u_k \rightarrow 0$ . Therefore,

$$\lim_{k \rightarrow \infty} \frac{|u_{k+1}|}{|u_k|} = 1 - \rho > 0$$

and thus the order of convergence is 1.



### 8.7

a. The value of  $x^*$  (in terms of  $a$ ,  $b$ , and  $c$ ) that minimizes  $f$  is  $x^* = b/a$ .

b. We have  $f'(x) = ax - b$ . Therefore, the recursive equation for the DDS algorithm is

$$x^{(k+1)} = x^{(k)} - \alpha(ax^{(k)} - b) = (1 - \alpha a)x^{(k)} + \alpha b.$$

c. Let  $\tilde{x} = \lim_{k \rightarrow \infty} x^{(k)}$ . Taking limits of both sides of  $x^{(k+1)} = x^{(k)} - \alpha(ax^{(k)} - b)$  (from part b), we get

$$\tilde{x} = \tilde{x} - \alpha(a\tilde{x} - b).$$

Hence, we get  $\tilde{x} = b/a = x^*$ .

d. To find the order of convergence, we compute

$$\begin{aligned} \frac{|x^{(k+1)} - b/a|}{|x^{(k)} - b/a|^p} &= \frac{|(1 - \alpha a)x^{(k)} + \alpha b - b/a|}{|x^{(k)} - b/a|^p} \\ &= \frac{|(1 - \alpha a)x^{(k)} - (1 - \alpha a)b/a|}{|x^{(k)} - b/a|^p} \\ &= |1 - \alpha a| |x^{(k)} - b/a|^{1-p}. \end{aligned}$$

Let  $z^{(k)} = |1 - \alpha a| |x^{(k)} - b/a|^{1-p}$ . Note that  $z^{(k)}$  converges to a finite nonzero number if and only if  $p = 1$  (if  $p < 1$ , then  $z^{(k)} \rightarrow 0$ , and if  $p > 1$ , then  $z^{(k)} \rightarrow \infty$ ). Therefore, the order of convergence of  $\{x^{(k)}\}$  is 1,

e. Let  $y^{(k)} = |x^{(k)} - b/a|$ . From part d, after some manipulation we obtain

$$y^{(k+1)} = |1 - \alpha a| y^{(k)} = |1 - \alpha a|^{k+1} y^{(0)}.$$

The sequence  $\{x^{(k)}\}$  converges (to  $b/a$ ) if and only if  $y^{(k)} \rightarrow 0$ . This holds if and only if  $|1 - \alpha a| < 1$ , which is equivalent to  $0 < \alpha < 2/a$ .

### 8.8

We rewrite  $f$  as  $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{b}^\top \mathbf{x}$ , where

$$\mathbf{Q} = \begin{bmatrix} 6 & 4 \\ 4 & 6 \end{bmatrix}$$

The characteristic polynomial of  $\mathbf{Q}$  is  $\lambda^2 - 12\lambda + 20$ . Hence, the eigenvalues of  $\mathbf{Q}$  are 2 and 10. Therefore, the largest range of values of  $\alpha$  for which the algorithm is globally convergent is  $0 < \alpha < 2/10$ .

### 8.9

a. We can write  $\mathbf{h}(\mathbf{x}) = \mathbf{Q} \mathbf{x} - \mathbf{b}$ , where  $\mathbf{b} = [-4, -1]^\top$  and

$$\mathbf{Q} = \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix}$$

is positive definite. Hence, the solution is

$$\mathbf{Q}^{-1} \mathbf{b} = \frac{1}{5} \begin{bmatrix} 3 & -2 \\ -2 & 3 \end{bmatrix} \begin{bmatrix} -4 \\ -1 \end{bmatrix} = \begin{bmatrix} -2 \\ 1 \end{bmatrix}.$$

b. By part a, the algorithm is a fixed-step-size gradient algorithm for a problem with gradient  $\mathbf{h}$ . The eigenvalues of  $\mathbf{Q}$  are 1 and 5. Hence, the largest range of values of  $\alpha$  such that the algorithm is globally convergent to the solution is  $0 < \alpha < 2/5$ .

c. The eigenvectors of  $\mathbf{Q}$  corresponding to eigenvalue 5 has the form  $c[1, 1]^\top$ , where  $c \in \mathbb{R}$ . Hence, to violate the descent property, we pick

$$\mathbf{x}^{(0)} = \mathbf{Q}^{-1} \mathbf{b} + c \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -3 \\ 0 \end{bmatrix}$$

where we choose  $c = -1$  so that  $\mathbf{x}^{(0)}$  has the specified form.

### 8.10

a. We have

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 3 & 1+a \\ 1+a & 3 \end{bmatrix} \mathbf{x} - \mathbf{x}^\top \begin{bmatrix} 1 \\ 1 \end{bmatrix} + b.$$

b. The unique global minimizer exists if and only if the Hessian is positive definite, which holds if and only if  $(1+a)^2 < 9$  (by Sylvester's criterion). Hence, the largest set of values of  $a$  and  $b$  such that the global minimizer of  $f$  exists is given by  $-4 < a < 2$  and  $b \in \mathbb{R}$  (unrestricted).

The minimizer is given by

$$\mathbf{x}^* = \frac{1}{9 - (1+a)^2} \begin{bmatrix} 3 & -(1+a) \\ -(1+a) & 3 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \frac{3 - (1+a)}{9 - (1+a)^2} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \frac{1}{4+a} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

c. The algorithm is a gradient algorithm with fixed step size  $\alpha = 2/5$ . The eigenvalues of the Hessian are (after some calculations)  $4+a$  and  $2-a$ . For global convergence, we need  $2/5 < 2/\lambda_{\max}$ , or  $\lambda_{\max} < 5$ , where  $\lambda_{\max} = \max(4+a, 2-a)$ . From this we deduce that  $-3 < a < 1$ . Hence, the largest set of values of  $a$  and  $b$  such that the algorithm is globally convergent is given by  $-3 < a < 1$  and  $b \in \mathbb{R}$  (unrestricted).

### 8.11

a. We have

$$\begin{aligned} f(x^{(k+1)}) &= (x^{(k+1)} - c)^2/2 \\ &= (x^{(k)} - \alpha_k(x^{(k)} - c) - c)^2/2 \\ &= (1 - \alpha_k)^2(x^{(k)} - c)^2/2 \\ &= (1 - \alpha_k)^2 f(x^{(k)}). \end{aligned}$$

b. We have  $x^{(k)} \rightarrow c$  if and only if  $f(x^{(k)}) \rightarrow 0$ . Hence, the algorithm is globally convergent if and only if  $f(x^{(k)}) \rightarrow 0$  for any  $x_0$ . From part a, we deduce that  $f(x^{(k)}) \rightarrow 0$  for any  $x_0$  if and only if  $\prod_{k=0}^{\infty} (1 - \alpha_k)^2 = 0$ . Because  $0 < \alpha < 1$ , this condition is equivalent to  $\prod_{k=0}^{\infty} (1 - \alpha_k) = 0$ , which holds if and only if

$$\sum_{k=0}^{\infty} \alpha_k = \infty.$$

### 8.12

The only local minimizer of  $f$  is  $x^* = 1/\sqrt{3}$ . Indeed, we have  $f'(x^*) = 0$  and  $f''(x^*) = 2\sqrt{3}$ . To find the largest range of values of  $\alpha$  such that the algorithm is locally convergent, we use a “linearization” argument: The algorithm is locally convergent if and only if the “linearized” algorithm  $x^{(k+1)} = x^{(k)} - \alpha f''(x^*)(x^{(k)} - x^*)$  is globally convergent. But the linearized algorithm is just a fixed step size algorithm applied to a quadratic with second derivative  $f''(x^*)$ . Therefore, the largest range of values of  $\alpha$  such that the algorithm is locally convergent is  $0 < \alpha < 2/f''(x^*) = 1/\sqrt{3}$ .

### 8.13

We use the formula from Lemma 8.1:

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

(we have  $V = f$  in this case). Using the expression for  $\gamma_k$ , we get, assuming  $x^{(k)} \neq 1$ ,

$$\gamma_k = \alpha 4 \cdot 2^{-k} (1 - \alpha 2^{-k}).$$

Hence,  $\gamma_k > 0$ , which means that  $f(x^{(k+1)}) < f(x^{(k)})$  if  $x^{(k)} \neq 1$  for  $k \geq 0$ . This implies that the algorithm has the descent property (for  $k \geq 0$ ).

We also note that

$$\sum_{k=0}^{\infty} \gamma_k = \alpha 4 \left( \sum_{k=0}^{\infty} 2^{-k} - \sum_{k=0}^{\infty} \alpha 4^{-k} \right) = \alpha 4 \left( 2 - \frac{4\alpha}{3} \right) < \inf ty.$$

Since  $\gamma_k > 0$  for all  $k \geq 0$ , we can apply the theorem given in class to deduce that the algorithm is not globally convergent.

#### 8.14

We have

$$|x^{(k+1)} - x^*| = |x^{(k)} - x^* - f'(x^{(k)})/f''(x^*)|.$$

By Taylor's Theorem,

$$f'(x^{(k)}) = f'(x^*) + f''(x^*)(x^{(k)} - x^*) + O(|x^{(k)} - x^*|^2).$$

Since  $f'(x^*) = 0$  by the FONC, we get

$$x^{(k)} - x^* - f'(x^{(k)})/f''(x^*) = O(|x^{(k)} - x^*|^2).$$

Combining the above with the first equation, we get

$$|x^{(k+1)} - x^*| = O(|x^{(k)} - x^*|^2),$$

which implies that the order of convergence is at least 2.

#### 8.15

a. The objective function is a quadratic that can be written as

$$f(x) = (\mathbf{a}x - \mathbf{b})^\top (\mathbf{a}x - \mathbf{b}) = \|\mathbf{a}\|^2 x^2 - 2\mathbf{a}^\top \mathbf{b}x + \|\mathbf{b}\|^2.$$

Hence, the minimizer is  $x^* = \mathbf{a}^\top \mathbf{b} / \|\mathbf{a}\|^2$ .

b. Note that  $f''(x) = 2\|\mathbf{a}\|^2$ . Thus, by the result for fixed step size gradient algorithms, the required largest range for  $\alpha$  is  $(0, 1/\|\mathbf{a}\|^2)$ .

#### 8.16

a. We have

$$\begin{aligned} f(x) = \|\mathbf{A}x - \mathbf{b}\|^2 &= (\mathbf{A}x - \mathbf{b})^\top (\mathbf{A}x - \mathbf{b}) \\ &= (\mathbf{x}^\top \mathbf{A}^\top - \mathbf{b}^\top)(\mathbf{A}x - \mathbf{b}) \\ &= \mathbf{x}^\top (\mathbf{A}^\top \mathbf{A})x - 2(\mathbf{A}^\top \mathbf{b})^\top x + \mathbf{b}^\top \mathbf{b} \end{aligned}$$

which is a quadratic function. The gradient is given by  $\nabla f(x) = 2(\mathbf{A}^\top \mathbf{A})x - 2(\mathbf{A}^\top \mathbf{b})$  and the Hessian is given by  $\mathbf{F}(x) = 2(\mathbf{A}^\top \mathbf{A})$ .

b. The fixed step size gradient algorithm for solving the above optimization problem is given by

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \alpha(2(\mathbf{A}^\top \mathbf{A})\mathbf{x}^{(k)} - 2\mathbf{A}^\top \mathbf{b}) \\ &= \mathbf{x}^{(k)} - 2\alpha\mathbf{A}^\top (\mathbf{A}\mathbf{x}^{(k)} - \mathbf{b}). \end{aligned}$$

c. The largest range of values for  $\alpha$  such that the algorithm in part b converges to the solution of the problem is given by

$$0 < \alpha < \frac{2}{\lambda_{\max}(2\mathbf{A}^\top \mathbf{A})} = \frac{1}{4}.$$

#### 8.17

a. We use contraposition. Suppose an eigenvalue of  $\mathbf{A}$  is negative:  $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$ , where  $\lambda < 0$  and  $\mathbf{v}$  is a corresponding eigenvector. Choose  $\mathbf{x}^{(0)} = \mathbf{v} + \mathbf{x}^*$ . Then,

$$\mathbf{x}^{(1)} = \mathbf{v} + \mathbf{x}^* - \alpha(\mathbf{A}\mathbf{v} + \mathbf{A}\mathbf{x}^* + \mathbf{b}) = \mathbf{v} + \mathbf{x}^* - \alpha\lambda\mathbf{v},$$

and hence

$$\mathbf{x}^{(1)} - \mathbf{x}^* = (1 - \alpha\lambda)(\mathbf{x}^{(0)} - \mathbf{x}^*).$$

Since  $1 - \alpha\lambda > 1$ , we conclude that the algorithm is not globally monotone.

b. Note that the algorithm is identical to a fixed step size gradient algorithm applied to a quadratic with Hessian  $\mathbf{A}$ . The eigenvalues of  $\mathbf{A}$  are 1 and 5. Therefore, the largest range of values of  $\alpha$  for which the algorithm is globally convergent is  $0 < \alpha < 2/5$ .

### 8.18

The steepest descent algorithm applied to the quadratic function  $f$  has the form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)} = \mathbf{x}^{(k)} - \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}} \mathbf{g}^{(k)}.$$

$\Rightarrow$ : If  $\mathbf{x}^{(1)} = \mathbf{Q}^{-1} \mathbf{b}$ , then

$$\mathbf{Q}^{-1} \mathbf{b} = \mathbf{x}^{(0)} - \alpha_0 \mathbf{g}^{(0)}.$$

Rearranging the above yields

$$\mathbf{Q} \mathbf{x}^{(0)} - \mathbf{b} = \alpha_0 \mathbf{Q} \mathbf{g}^{(0)}.$$

Since  $\mathbf{g}^{(0)} = \mathbf{Q} \mathbf{x}^{(0)} - \mathbf{b} \neq \mathbf{0}$ , we have

$$\mathbf{Q} \mathbf{g}^{(0)} = \frac{1}{\alpha_0} \mathbf{g}^{(0)}$$

which means that  $\mathbf{g}^{(0)}$  is an eigenvector of  $\mathbf{Q}$  (with corresponding eigenvalue  $1/\alpha_0$ ).

$\Leftarrow$ : By assumption,  $\mathbf{Q} \mathbf{g}^{(0)} = \lambda \mathbf{g}^{(0)}$ , where  $\lambda \in \mathbb{R}$ . We want to show that  $\mathbf{Q} \mathbf{x}^{(1)} = \mathbf{b}$ . We have

$$\begin{aligned} \mathbf{Q} \mathbf{x}^{(1)} &= \mathbf{Q} \left( \mathbf{x}^{(0)} - \frac{\mathbf{g}^{(0)\top} \mathbf{g}^{(0)}}{\mathbf{g}^{(0)\top} \mathbf{Q} \mathbf{g}^{(0)}} \mathbf{g}^{(0)} \right) \\ &= \mathbf{Q} \mathbf{x}^{(0)} - \frac{1}{\lambda} \frac{\mathbf{g}^{(0)\top} \mathbf{g}^{(0)}}{\mathbf{g}^{(0)\top} \mathbf{g}^{(0)}} \mathbf{Q} \mathbf{g}^{(0)} \\ &= \mathbf{Q} \mathbf{x}^{(0)} - \mathbf{g}^{(0)} \\ &= \mathbf{b}. \end{aligned}$$

### 8.19

a. Possible. Pick  $f$  such that  $\lambda_{\max} \geq 2\lambda_{\min}$  and  $\mathbf{x}^{(0)}$  such that  $\mathbf{g}^{(0)}$  is an eigenvector of  $\mathbf{Q}$  with eigenvalue  $\lambda_{\min}$ . Then,

$$\alpha_0 = \frac{\mathbf{g}^{(0)\top} \mathbf{g}^{(0)}}{\mathbf{g}^{(0)\top} \mathbf{Q} \mathbf{g}^{(0)}} = \frac{1}{\lambda_{\min}} \geq \frac{2}{\lambda_{\max}}.$$

b. Not possible. Indeed, using Rayleigh's inequality,

$$\alpha_0 = \frac{\mathbf{g}^{(0)\top} \mathbf{g}^{(0)}}{\mathbf{g}^{(0)\top} \mathbf{Q} \mathbf{g}^{(0)}} \leq \frac{1}{\lambda_{\min}}.$$

### 8.20

a. We rewrite  $f$  as

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{b}^\top \mathbf{x} - 22,$$

where

$$\mathbf{Q} = \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -3 \\ 1 \end{bmatrix}.$$

The eigenvalues of  $\mathbf{Q}$  are 1 and 5. Therefore, the range of values of the step size for which the algorithm converges to the minimizer is  $0 < \alpha < 2/5$ .

b. An eigenvector of  $\mathbf{Q}$  corresponding to the eigenvalue 5 is  $\mathbf{v} = [1, 1]^\top / 5$ . We have  $\mathbf{x}^* = \mathbf{Q}^{-1}\mathbf{b} = [-11, 9]^\top / 5$ . Hence, an initial condition that results in the algorithm diverging is

$$\mathbf{x}^{(0)} = \mathbf{x}^* + \mathbf{v} = \begin{bmatrix} -2 \\ 2 \end{bmatrix}.$$

---

### 8.21

In both cases, we compute the Hessian  $\mathbf{Q}$  of  $f$ , and find its largest eigenvalue  $\lambda_{\max}$ . Then the range we seek is  $0 < \alpha < 2/\lambda_{\max}$ .

a. In this case,

$$\mathbf{Q} = \begin{bmatrix} 6 & 4 \\ 4 & 6 \end{bmatrix},$$

with eigenvalues 2 and 10. Hence, the answer is  $0 < \alpha < 1/5$ .

b. In this case, again we have

$$\mathbf{Q} = \begin{bmatrix} 6 & 4 \\ 4 & 6 \end{bmatrix},$$

with eigenvalues 2 and 10. Hence, the answer is  $0 < \alpha < 1/5$ .

---

### 8.22

For the given algorithm we have

$$\gamma_k = \beta(2 - \beta) \left( \frac{(\mathbf{g}^{(k)})^\top \mathbf{g}^{(k)})^2}{(\mathbf{g}^{(k)})^\top \mathbf{Q} \mathbf{g}^{(k)}) (\mathbf{g}^{(k)})^\top \mathbf{Q}^{-1} \mathbf{g}^{(k)})} \right)$$

If  $0 < \beta < 2$ , then  $\beta(2 - \beta) > 0$ , and by Lemma 8.2,

$$\gamma_k \geq \beta(2 - \beta) \left( \frac{\lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbf{Q})} \right) > 0$$

which implies that  $\sum_{k=0}^{\infty} \gamma_k = \infty$ . Hence, by Theorem 8.1,  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$  for any  $\mathbf{x}^{(0)}$ .

If  $\beta \leq 0$  or  $\beta \geq 2$ , then  $\beta(2 - \beta) \leq 0$ , and by Lemma 8.2,

$$\gamma_k \leq \beta(2 - \beta) \left( \frac{\lambda_{\max}(\mathbf{Q})}{\lambda_{\min}(\mathbf{Q})} \right) < 0.$$

By Lemma 8.1,  $V(\mathbf{x}^{(k)}) \geq V(\mathbf{x}^{(0)})$ . Hence, if  $\mathbf{x}^{(0)} \neq \mathbf{x}^*$ , then  $\{V(\mathbf{x}^{(k)})\}$  does not converge to 0, and consequently  $\mathbf{x}^{(k)}$  does not converge to  $\mathbf{x}^*$ .

---

### 8.23

By Lemma 8.1,  $V(\mathbf{x}^{(k+1)}) = (1 - \gamma_k)V(\mathbf{x}^{(k)})$  for all  $k$ . Note that the algorithm has a descent property if and only if  $V(\mathbf{x}^{(k+1)}) < V(\mathbf{x}^{(k)})$  whenever  $\mathbf{g}^{(k)} \neq \mathbf{0}$ . Clearly, whenever  $\mathbf{g}^{(k)} \neq \mathbf{0}$ ,  $V(\mathbf{x}^{(k+1)}) < V(\mathbf{x}^{(k)})$  if and only if  $1 - \gamma_k < 1$ . The desired result follows immediately.

---

### 8.24

We have

$$\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} = \alpha_k \mathbf{d}^{(k)}$$

and hence

$$\langle \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}, \nabla f(\mathbf{x}^{(k+1)}) \rangle = \alpha_k \langle \mathbf{d}^{(k)}, \nabla f(\mathbf{x}^{(k+1)}) \rangle.$$

Now, let  $\phi_k(\alpha) = f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$ . Since  $\alpha_k$  minimizes  $\phi_k$ , then by the FONC,  $\phi'_k(\alpha_k) = 0$ . By the chain rule,  $\phi'_k(\alpha) = \mathbf{d}^{(k)\top} \nabla f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$ . Hence,

$$0 = \phi'_k(\alpha_k) = \mathbf{d}^{(k)\top} \nabla f(\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}) = \langle \mathbf{d}^{(k)}, \nabla f(\mathbf{x}^{(k+1)}) \rangle,$$

and so

$$\langle \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}, \nabla f(\mathbf{x}^{(k+1)}) \rangle = 0.$$

## 8.25

A simple MATLAB routine for implementing the steepest descent method is as follows.

```
function [x,N]=steep_desc(grad,xnew,options);
%   STEEP_DESC('grad',x0);
%   STEEP_DESC('grad',x0,OPTIONS);
%
%   x = STEEP_DESC('grad',x0);
%   x = STEEP_DESC('grad',x0,OPTIONS);
%
%   [x,N] = STEEP_DESC('grad',x0);
%   [x,N] = STEEP_DESC('grad',x0,OPTIONS);
%
%The first variant finds the minimizer of a function whose gradient
%is described in grad (usually an M-file: grad.m), using a gradient
%descent algorithm with initial point x0. The line search used in the
%secant method.
%The second variant allows a vector of optional parameters to
%defined. OPTIONS(1) controls how much display output is given; set
%to 1 for a tabular display of results, (default is no display:0).
%OPTIONS(2) is a measure of the precision required for the final point.
%OPTIONS(3) is a measure of the precision required of the gradient.
%OPTIONS(14) is the maximum number of iterations.
%For more information type HELP FOPTIONS.
%
%The next two variants returns the value of the final point.
%The last two variants returns a vector of the final point and the
%number of iterations.

if nargin ~= 3
    options = [];
    if nargin ~= 2
        disp('Wrong number of arguments. ');
        return;
    end
end

if length(options) >= 14
    if options(14)==0
        options(14)=1000*length(xnew);
    end
else
    options(14)=1000*length(xnew);
end

clc;
format compact;
format short e;

options = foptions(options);
print = options(1);
epsilon_x = options(2);
epsilon_g = options(3);
max_iter=options(14);

for k = 1:max_iter,
```

```

xcurr=xnew;
g_curr=feval(grad,xcurr);
if norm(g_curr) <= epsilon_g
    disp('Terminating: Norm of gradient less than');
    disp(epsilon_g);
    k=k-1;
    break;
end %if

alpha=linesearch_secant(grad,xcurr,-g_curr);

xnew = xcurr-alpha*g_curr;

if print,
    disp('Iteration number k =')
    disp(k); %print iteration index k
    disp('alpha =');
    disp(alpha); %print alpha
    disp('Gradient = ');
    disp(g_curr); %print gradient
    disp('New point =');
    disp(xnew); %print new point
end %if

if norm(xnew-xcurr) <= epsilon_x*norm(xcurr)
    disp('Terminating: Norm of difference between iterates less than');
    disp(epsilon_x);
    break;
end %if

if k == max_iter
    disp('Terminating with maximum number of iterations');
end %if
end %for

if nargout >= 1
    x=xnew;
    if nargout == 2
        N=k;
    end
else
    disp('Final point =');
    disp(xnew);
    disp('Number of iterations =');
    disp(k);
end %if
%-----

```

To apply the above MATLAB routine to the function in Example 8.1, we need the following M-file to specify the gradient.

```

function y=g(x)

y=[4*(x(1)-4).^3; 2*(x(2)-3); 16*(x(3)+5).^3];

```

We applied the algorithm as follows:

```

>> options(2) = 10^(-6);
>> options(3) = 10^(-6);

```

```
>> steep_desc('g', [-4;5;1], options)
```

```
Terminating: Norm of gradient less than
1.0000e-06
Final point =
4.0022e+00 3.0000e+00 -4.9962e+00
Number of iterations =
25
```

As we can see above, we obtained the final point  $[4.002, 3.000, -4.996]^\top$  after 25 iterations. The value of the objective function at the final point is  $7.2 \times 10^{-10}$ .

## 8.26

The algorithm terminated after 9127 iterations. The final point was  $[0.99992, 0.99983]^\top$ .

## 9. Newton's Method

### 9.1

a. We have  $f'(x) = 4(x - x_0)^3$  and  $f''(x) = 12(x - x_0)^2$ . Hence, Newton's method is represented as

$$x^{(k+1)} = x^{(k)} - \frac{x^{(k)} - x_0}{3},$$

which upon rewriting becomes

$$x^{(k+1)} - x_0 = \frac{2}{3} (x^{(k)} - x_0)$$

b. From part a,  $y^{(k)} = |x^{(k)} - x_0| = (2/3)|x^{(k-1)} - x_0| = (2/3)y^{(k-1)}$ .

c. From part b, we see that  $y^{(k)} = (2/3)^k y^{(0)}$  and therefore  $y^{(k)} \rightarrow 0$ . Hence  $x^{(k)} \rightarrow x_0$  for any  $x^{(0)}$ .

d. From part b, we have

$$\lim_{k \rightarrow \infty} \frac{|x^{(k+1)} - x_0|}{|x^{(k)} - x_0|} = \lim_{k \rightarrow \infty} \frac{2}{3} = \frac{2}{3} > 0$$

and hence the order of convergence is 1.

e. The theorem assumes that  $f''(x^*) \neq 0$ . However, in this problem,  $x^* = x_0$ , and  $f''(x^*) = 0$ .

### 9.2

a. We have

$$|x^{(k+1)} - x^*| = |x^{(k)} - x^* - \alpha_k f'(x^{(k)})|.$$

By Taylor's theorem applied to  $f'$ ,

$$f'(x^{(k)}) = f'(x^*) + f''(x^*)(x^{(k)} - x^*) + o(|x^{(k)} - x^*|).$$

Since  $f'(x^*) = 0$  by the FONC, we get

$$\begin{aligned} x^{(k)} - x^* - \alpha_k f'(x^{(k)}) &= (1 - \alpha_k f''(x^*))(x^{(k)} - x^*) + \alpha_k o(|x^{(k)} - x^*|) \\ &= o(|x^{(k)} - x^*|) + \alpha_k o(|x^{(k)} - x^*|) \\ &= (1 + \alpha_k) o(|x^{(k)} - x^*|). \end{aligned}$$

Because  $\{\alpha_k\}$  converges, it is bounded, and so  $(1 + \alpha_k) o(|x^{(k)} - x^*|) = o(|x^{(k)} - x^*|)$ . Combining the above with the first equation, we get

$$|x^{(k+1)} - x^*| = o(|x^{(k)} - x^*|),$$

which implies that the order of convergence is superlinear.

b. In the secant algorithm, if  $x^{(k)} \rightarrow x^*$ , then  $(f'(x^{(k)}) - f'(x^{(k-1)}))/(x^{(k)} - x^{(k-1)}) \rightarrow f''(x^*)$ . Since the secant algorithm has the form  $x^{(k+1)} = x^{(k)} - \alpha_k f'(x^{(k)})$  with  $\alpha_k = (x^{(k)} - x^{(k-1)})/(f'(x^{(k)}) - f'(x^{(k-1)}))$ , we



deduce that  $\alpha_k \rightarrow 1/f''(x^*)$ . Hence, if we apply the secant algorithm to a function  $f \in \mathcal{C}^2$ , and it converges to a local minimizer  $x^*$  such that  $f''(x^*) \neq 0$ , then the order of convergence is superlinear.

### 9.3

a. We compute  $f'(x) = 4x^{1/3}/3$  and  $f''(x) = 4x^{-2/3}/9$ . Therefore Newton's algorithm for this problem takes the form

$$x^{(k+1)} = x^{(k)} - \frac{4(x^{(k)})^{1/3}/3}{4(x^{(k)})^{-2/3}/9} = -2x^{(k)}.$$

b. From part a, we have  $x^{(k)} = 2^k x^{(0)}$ . Therefore, as long as  $x^{(0)} \neq 0$ , the sequence  $\{x^{(k)}\}$  does not converge to 0.

### 9.4

a. Clearly  $f(\mathbf{x}) \geq 0$  for all  $\mathbf{x}$ . We have

$$\begin{aligned} f(\mathbf{x}) = 0 & \Leftrightarrow x_2 - x_1^2 = 0 \quad \text{and} \quad 1 - x_1 = 0 \\ & \Leftrightarrow \mathbf{x} = [1, 1]^\top. \end{aligned}$$

Hence,  $f(\mathbf{x}) > f([1, 1]^\top)$  for all  $\mathbf{x} \neq [1, 1]^\top$ , and therefore  $[1, 1]^\top$  is the unique global minimizer.

b. We compute

$$\begin{aligned} \nabla f(\mathbf{x}) &= \begin{bmatrix} 400x_1^3 - 400x_1x_2 + 2x_1 - 2 \\ 200(x_2 - x_1^2) \end{bmatrix} \\ \mathbf{F}(\mathbf{x}) &= \begin{bmatrix} 1200x_1^2 - 400x_2 + 2 & -400x_1 \\ -400x_1 & 200 \end{bmatrix}. \end{aligned}$$

To apply Newton's method we use the inverse of the Hessian, which is

$$\mathbf{F}(\mathbf{x})^{-1} = \frac{1}{80000(x_1^2 - x_2) + 400} \begin{bmatrix} 200 & 400x_1 \\ 400x_1 & 1200x_1^2 - 400x_2 + 2 \end{bmatrix}.$$

Applying two iterations of Newton's method, we have  $\mathbf{x}^{(1)} = [1, 0]^\top$ ,  $\mathbf{x}^{(2)} = [1, 1]^\top$ . Therefore, in this particular case, the method converges in two steps! We emphasize, however, that this fortuitous situation is by no means typical, and is highly dependent on the initial condition.

c. Applying the gradient algorithm  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})$  with a fixed step size of  $\alpha_k = 0.05$ , we obtain  $\mathbf{x}^{(1)} = [0.1, 0]^\top$ ,  $\mathbf{x}^{(2)} = [0.17, 0.1]^\top$ .

### 9.5

If  $\mathbf{x}^{(0)} = \mathbf{x}^*$ , we are done. So, assume  $\mathbf{x}^{(0)} \neq \mathbf{x}^*$ . Since the standard Newton's method reaches the point  $\mathbf{x}^*$  in one step, we have

$$\begin{aligned} f(\mathbf{x}^*) &= f(\mathbf{x}^{(0)} + \mathbf{Q}^{-1}\mathbf{g}^{(0)}) \\ &= \min f(\mathbf{x}) \\ &\leq f(\mathbf{x}^{(0)} + \alpha\mathbf{Q}^{-1}\mathbf{g}^{(0)}) \end{aligned}$$

for any  $\alpha \geq 0$ . Hence,

$$\alpha_0 = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(0)} + \alpha\mathbf{Q}^{-1}\mathbf{g}^{(0)}) = 1.$$

Hence, in this case, the modified Newton's algorithm is equivalent to the standard Newton's algorithm, and thus  $\mathbf{x}^{(1)} = \mathbf{x}^*$ .

## 10. Conjugate Direction Methods

### 10.1

We proceed by induction to show that for  $k = 0, \dots, n-1$ , the set  $\{\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k)}\}$  is  $\mathbf{Q}$ -conjugate. We assume that  $\mathbf{d}^{(i)} \neq \mathbf{0}$ ,  $i = 1, \dots, k$ , so that  $\mathbf{d}^{(i)\top} \mathbf{Q} \mathbf{d}^{(i)} \neq 0$  and the algorithm is well defined.

For  $k = 0$ , the statement trivially holds. So, assume that the statement is true for  $k < n-1$ , i.e.,  $\{\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k)}\}$  is  $\mathbf{Q}$ -conjugate. We now show that  $\{\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k+1)}\}$  is  $\mathbf{Q}$ -conjugate. For this, we need only to show that for each  $j = 0, \dots, k$ , we have  $\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(j)} = 0$ . To this end,

$$\begin{aligned} \mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(j)} &= \left( \mathbf{p}^{(k+1)\top} - \sum_{i=0}^k \frac{\mathbf{p}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(i)}}{\mathbf{d}^{(i)\top} \mathbf{Q} \mathbf{d}^{(i)}} \mathbf{d}^{(i)\top} \right) \mathbf{Q} \mathbf{d}^{(j)} \\ &= \mathbf{p}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(j)} - \sum_{i=0}^k \frac{\mathbf{p}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(i)}}{\mathbf{d}^{(i)\top} \mathbf{Q} \mathbf{d}^{(i)}} \mathbf{d}^{(i)\top} \mathbf{Q} \mathbf{d}^{(j)}. \end{aligned}$$

By the induction hypothesis,  $\mathbf{d}^{(i)\top} \mathbf{Q} \mathbf{d}^{(j)} = 0$  for  $i \neq j$ . Therefore

$$\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(j)} = \mathbf{p}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(j)} - \frac{\mathbf{p}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(j)}}{\mathbf{d}^{(j)\top} \mathbf{Q} \mathbf{d}^{(j)}} \mathbf{d}^{(j)\top} \mathbf{Q} \mathbf{d}^{(j)} = 0.$$

In the above, we have assumed that the vectors  $\mathbf{d}^{(k)}$  are nonzero (so that  $\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)} \neq 0$  and the algorithm is well defined). To prove that this assumption holds, we use induction to show that  $\mathbf{d}^{(k)}$  is a (nonzero) linear combination of  $\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(k)}$  (which immediately implies that  $\mathbf{d}^{(k)}$  is nonzero because of the linear independence of  $\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(k)}$ ).

For  $k = 0$ , we have  $\mathbf{d}^{(0)} = \mathbf{p}^{(0)}$  by definition. Assume that the result holds for  $k < n-1$ ; i.e.,  $\mathbf{d}^{(k)} = \sum_{j=0}^k \alpha_j^{(k)} \mathbf{p}^{(j)}$ , where the coefficients  $\alpha_j^{(k)}$  are not all zero. Consider  $\mathbf{d}^{(k+1)}$ :

$$\begin{aligned} \mathbf{d}^{(k+1)} &= \mathbf{p}^{(k+1)} - \sum_{i=0}^k \beta_i \mathbf{d}^{(i)} \\ &= \mathbf{p}^{(k+1)} - \sum_{i=0}^k \beta_i \sum_{j=0}^i \alpha_j^{(i)} \mathbf{p}^{(j)} \\ &= \mathbf{p}^{(k+1)} - \sum_{j=0}^k \sum_{i=j}^k \beta_i \alpha_j^{(i)} \mathbf{p}^{(j)}. \end{aligned}$$

So, clearly  $\mathbf{d}^{(k+1)}$  is a nonzero linear combination of  $\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(k+1)}$ .

### 10.2

Let  $k \in \{0, \dots, n-1\}$  and  $\phi_k(\alpha) = f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$ . By the chain rule, we have

$$\phi'_k(\alpha_k) = \nabla f(\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)})^\top \mathbf{d}^{(k)} = \mathbf{g}^{(k+1)\top} \mathbf{d}^{(k)}.$$

Since  $\mathbf{g}^{(k+1)\top} \mathbf{d}^{(k)} = 0$ , we have  $\phi'_k(\alpha_k) = 0$ . Note that

$$\phi_k(\alpha) = \frac{1}{2} \left( \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)} \right) \alpha^2 + \mathbf{g}^{(k)\top} \mathbf{d}^{(k)} \alpha + \text{constant}.$$

As  $\phi$  is a quadratic function of  $\alpha$  (with positive coefficient in the quadratic term), we conclude that  $\alpha_k = \arg \min_{\alpha} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$ .

Note that since  $\mathbf{g}^{(k)\top} \mathbf{d}^{(k)} \neq 0$  is the coefficient of the linear term in  $\phi_k$ , we have  $\alpha_k \neq 0$ . For  $i \in \{0, \dots, k-1\}$ , we have

$$\begin{aligned} \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(i)} &= \frac{1}{\alpha_k} (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})^\top \mathbf{Q} \mathbf{d}^{(i)} \\ &= \frac{1}{\alpha_k} (\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)})^\top \mathbf{d}^{(i)} \\ &= \frac{1}{\alpha_k} (\mathbf{g}^{(k+1)\top} \mathbf{d}^{(i)} - \mathbf{g}^{(k)\top} \mathbf{d}^{(i)}) \\ &= 0 \end{aligned}$$

by assumption, which completes the proof.

### 10.3

From the conjugate gradient algorithm we have

$$\mathbf{d}^{(k)} = -\mathbf{g}^{(k)} + \frac{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k-1)}}{\mathbf{d}^{(k-1)\top} \mathbf{Q} \mathbf{d}^{(k-1)}} \mathbf{d}^{(k-1)}.$$

Premultiplying the above by  $\mathbf{d}^{(k)\top} \mathbf{Q}$  and using the fact that  $\mathbf{d}^{(k)}$  and  $\mathbf{d}^{(k-1)}$  are  $\mathbf{Q}$ -conjugate, yields

$$\begin{aligned} \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)} &= -\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)} + \frac{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k-1)}}{\mathbf{d}^{(k-1)\top} \mathbf{Q} \mathbf{d}^{(k-1)}} \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k-1)} \\ &= -\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}. \end{aligned}$$

### 10.4

a. Since  $\mathbf{Q}$  is symmetric, then there exists a set of vectors  $\{\mathbf{d}^{(1)}, \dots, \mathbf{d}^{(n)}\}$  such that  $\mathbf{Q} \mathbf{d}^{(i)} = \lambda_i \mathbf{d}^{(i)}$ ,  $i = 1, \dots, n$ , and  $\mathbf{d}^{(i)\top} \mathbf{d}^{(j)} = 0$ ,  $j \neq i$ , where the  $\lambda_i$  are (real) eigenvalues of  $\mathbf{Q}$ . Therefore, if  $i \neq j$ , we have  $\mathbf{d}^{(i)\top} \mathbf{Q} \mathbf{d}^{(j)} = \mathbf{d}^{(i)\top} (\lambda_j \mathbf{d}^{(j)}) = \lambda_j (\mathbf{d}^{(i)\top} \mathbf{d}^{(j)}) = 0$ . Hence the set  $\{\mathbf{d}^{(1)}, \dots, \mathbf{d}^{(n)}\}$  is  $\mathbf{Q}$ -conjugate.

b. Define  $\lambda_i = (\mathbf{d}^{(i)\top} \mathbf{Q} \mathbf{d}^{(i)}) / (\mathbf{d}^{(i)\top} \mathbf{d}^{(i)})$ . Let

$$\mathbf{D} = \begin{bmatrix} \mathbf{d}^{(1)\top} \\ \vdots \\ \mathbf{d}^{(n)\top} \end{bmatrix}.$$

Since  $\mathbf{Q}$  is positive definite and the set  $\{\mathbf{d}^{(1)}, \dots, \mathbf{d}^{(n)}\}$  is  $\mathbf{Q}$ -conjugate, then by Lemma 10.1, the set is also linearly independent. Hence,  $\mathbf{D}$  is nonsingular. By  $\mathbf{Q}$ -conjugacy, we have that for all  $i \neq j$ ,  $\mathbf{d}^{(i)\top} \mathbf{Q} \mathbf{d}^{(j)} = 0$ . By assumption, we have  $\mathbf{d}^{(i)\top} \lambda_j \mathbf{d}^{(j)} = \lambda_j \mathbf{d}^{(i)\top} \mathbf{d}^{(j)} = 0$ . Hence,  $\mathbf{d}^{(i)\top} \mathbf{Q} \mathbf{d}^{(j)} = \lambda_j \mathbf{d}^{(i)\top} \mathbf{d}^{(j)}$ . Moreover, for each  $i = 1, \dots, n$ , we have  $\mathbf{d}^{(i)\top} \mathbf{Q} \mathbf{d}^{(i)} = \mathbf{d}^{(i)\top} \lambda_i \mathbf{d}^{(i)} = \lambda_i \mathbf{d}^{(i)\top} \mathbf{d}^{(i)}$ . We can write the above conditions in matrix form:

$$\mathbf{D} \mathbf{Q} \mathbf{d}^{(i)} = \mathbf{D} (\lambda_i \mathbf{d}^{(i)}).$$

Since  $\mathbf{D}$  is nonsingular, then we have

$$\mathbf{Q} \mathbf{d}^{(i)} = \lambda_i \mathbf{d}^{(i)},$$

which completes the proof.

### 10.5

We have

$$\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k+1)} = \gamma_k \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k+1)} + \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}.$$

Hence, in order to have  $\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k+1)} = 0$ , we need

$$\gamma_k = -\frac{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k+1)}}.$$

### 10.6

We use induction. For  $k = 0$ , we have

$$\mathbf{d}^{(0)} = a_0 \mathbf{g}^{(0)} = -a_0 \mathbf{b} \in \mathcal{V}_1.$$

Moreover,  $\mathbf{x}^{(0)} = \mathbf{0} \in \mathcal{V}_0$ . Hence, the proposition is true at  $k = 0$ . Assume it is true at  $k$ . To show that it is also true at  $k + 1$ , note first that

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}.$$

Because  $\mathbf{x}^{(k)} \in \mathcal{V}_k \subset \mathcal{V}_{k+1}$  and  $\mathbf{d}^{(k)} \in \mathcal{V}_{k+1}$  by the induction hypothesis, we deduce that  $\mathbf{x}^{(k+1)} \in \mathcal{V}_{k+1}$ . Moreover,

$$\begin{aligned}\mathbf{d}^{(k+1)} &= a_k \mathbf{g}^{(k+1)} + b_k \mathbf{d}^{(k)} \\ &= a_k (\mathbf{Q}\mathbf{x}^{(k+1)} - \mathbf{b}) + b_k \mathbf{d}^{(k)}.\end{aligned}$$

But because  $\mathbf{x}^{(k+1)} \in \mathcal{V}_{k+1}$ ,  $\mathbf{Q}\mathbf{x}^{(k+1)} - \mathbf{b} \in \mathcal{V}_{k+2}$ . Moreover,  $\mathbf{d}^{(k)} \in \mathcal{V}_{k+1} \subset \mathcal{V}_{k+2}$ . Hence,  $\mathbf{d}^{(k+1)} \in \mathcal{V}_{k+2}$ . This completes the induction proof.

b. The conjugate gradient algorithm is an instance of the algorithm given in the question. By the “expanding subspace” theorem, we can say that in the conjugate gradient algorithm (with  $\mathbf{x}^{(0)} = \mathbf{0}$ ), at each  $k$ ,  $\mathbf{x}^{(k)}$  is the global minimizer of  $f$  on the Krylov subspace  $\mathcal{V}_k$ . Note that for all  $k \geq n$ ,  $\mathcal{V}_{k+1} = \mathcal{V}_k$ , because of the Cayley-Hamilton theorem, which allows us to express  $\mathbf{Q}^n$  as a linear combination of  $\mathbf{I}, \mathbf{Q}, \dots, \mathbf{Q}^{n-1}$ .

### 10.7

Expanding  $\phi(\mathbf{a})$  yields

$$\begin{aligned}\phi(\mathbf{a}) &= \frac{1}{2}(\mathbf{x}_0 + \mathbf{D}\mathbf{a})^\top \mathbf{Q}(\mathbf{x}_0 + \mathbf{D}\mathbf{a}) - (\mathbf{x}_0 + \mathbf{D}\mathbf{a})^\top \mathbf{b} \\ &= \frac{1}{2}\mathbf{a}^\top (\mathbf{D}^\top \mathbf{Q} \mathbf{D}) \mathbf{a} + \mathbf{a}^\top (\mathbf{D}^\top \mathbf{Q} \mathbf{x}_0 - \mathbf{D}^\top \mathbf{b}) + \left( \frac{1}{2}\mathbf{x}_0^\top \mathbf{Q} \mathbf{x}_0 - \mathbf{x}_0^\top \mathbf{b} \right).\end{aligned}$$

Clearly  $\phi$  is a quadratic function on  $\mathbb{R}^r$ . It remains to show that the matrix in the quadratic term,  $\mathbf{D}^\top \mathbf{Q} \mathbf{D}$ , is positive definite. Since  $\mathbf{Q} > 0$ , for any  $\mathbf{a} \in \mathbb{R}^r$ , we have

$$\mathbf{a}^\top (\mathbf{D}^\top \mathbf{Q} \mathbf{D}) \mathbf{a} = (\mathbf{D}\mathbf{a})^\top \mathbf{Q}(\mathbf{D}\mathbf{a}) \geq 0$$

and

$$\mathbf{a}^\top (\mathbf{D}^\top \mathbf{Q} \mathbf{D}) \mathbf{a} = (\mathbf{D}\mathbf{a})^\top \mathbf{Q}(\mathbf{D}\mathbf{a}) = 0$$

if and only if  $\mathbf{D}\mathbf{a} = \mathbf{0}$ . Since  $\text{rank } \mathbf{D} = r$ ,  $\mathbf{D}\mathbf{a} = \mathbf{0}$  if and only if  $\mathbf{a} = \mathbf{0}$ . Hence, the matrix  $\mathbf{D}^\top \mathbf{Q} \mathbf{D}$  is positive definite.

### 10.8

a. Let  $0 \leq k \leq n-1$  and  $0 \leq i \leq k$ . Then,

$$\begin{aligned}\mathbf{g}^{(k+1)T} \mathbf{g}^{(i)} &= \mathbf{g}^{(k+1)T} (\beta_{i-1} \mathbf{d}^{(i-1)} - \mathbf{d}^{(i)}) \\ &= \beta_{i-1} \mathbf{g}^{(k+1)T} \mathbf{d}^{(i-1)} - \mathbf{g}^{(k+1)T} \mathbf{d}^{(i)} \\ &= 0\end{aligned}$$

by Lemma 10.2.

b. Let  $0 \leq k \leq n-1$  and  $0 \leq i \leq k-1$ . Then,

$$\begin{aligned}\mathbf{g}^{(k+1)T} \mathbf{Q} \mathbf{g}^{(i)} &= (\beta_k \mathbf{d}^{(k)} - \mathbf{d}^{(k+1)})^\top \mathbf{Q} (\beta_{i-1} \mathbf{d}^{(i-1)} - \mathbf{d}^{(i)}) \\ &= \beta_k \beta_{i-1} \mathbf{d}^{(k)T} \mathbf{Q} \mathbf{d}^{(i-1)} - \beta_k \mathbf{d}^{(k)T} \mathbf{Q} \mathbf{d}^{(i)} - \beta_{i-1} \mathbf{d}^{(k+1)T} \mathbf{Q} \mathbf{d}^{(i-1)} + \mathbf{d}^{(k+1)T} \mathbf{Q} \mathbf{d}^{(i)} \\ &= 0\end{aligned}$$

by  $\mathbf{Q}$ -conjugacy of  $\mathbf{d}^{(k+1)}$ ,  $\mathbf{d}^{(k)}$ ,  $\mathbf{d}^{(i)}$  and  $\mathbf{d}^{(i-1)}$  (note that the iteration indices here are all distinct).

### 10.9

We represent  $f$  as

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 5 & -3 \\ -3 & 2 \end{bmatrix} \mathbf{x} - \mathbf{x}^\top \begin{bmatrix} 0 \\ 1 \end{bmatrix} - 7.$$

The conjugate gradient algorithm is based on the following formulas:

$$\begin{aligned}\mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}, \quad \alpha_k = -\frac{\mathbf{g}^{(k)\top} \mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}} \\ \mathbf{d}^{(k+1)} &= -\mathbf{g}^{(k+1)} + \beta_k \mathbf{d}^{(k)}, \quad \beta_k = \frac{\mathbf{g}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}}.\end{aligned}$$

We have,

$$\mathbf{d}^{(0)} = \mathbf{g}^{(0)} = \mathbf{Q} \mathbf{x}^{(0)} - \mathbf{b} = -\mathbf{b} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}.$$

We then proceed to compute

$$\alpha_0 = -\frac{\mathbf{g}^{(0)\top} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} = -\frac{\begin{bmatrix} 0 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ -1 \end{bmatrix}}{\begin{bmatrix} 0 & -1 \end{bmatrix} \begin{bmatrix} 5 & -3 \\ -3 & 2 \end{bmatrix} \begin{bmatrix} 0 \\ -1 \end{bmatrix}} = -\frac{1}{2}.$$

Hence,

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1/2 \end{bmatrix}.$$

We next proceed by evaluating the gradient of the objective function at  $\mathbf{x}^{(1)}$ ,

$$\mathbf{g}^{(1)} = \mathbf{Q} \mathbf{x}^{(1)} - \mathbf{b} = \begin{bmatrix} 5 & -3 \\ -3 & 2 \end{bmatrix} \begin{bmatrix} 0 \\ 1/2 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} -3/2 \\ 0 \end{bmatrix}.$$

Because the gradient is nonzero, we can proceed with the next step where we compute

$$\beta_0 = \frac{\mathbf{g}^{(1)\top} \mathbf{Q} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} = \frac{\begin{bmatrix} -3/2 & 0 \end{bmatrix} \begin{bmatrix} 5 & -3 \\ -3 & 2 \end{bmatrix} \begin{bmatrix} 0 \\ -1 \end{bmatrix}}{\begin{bmatrix} 0 & -1 \end{bmatrix} \begin{bmatrix} 5 & -3 \\ -3 & 2 \end{bmatrix} \begin{bmatrix} 0 \\ -1 \end{bmatrix}} = -\frac{9}{4}.$$

Hence, the direction  $\mathbf{d}^{(1)}$  is

$$\mathbf{d}^{(1)} = -\mathbf{g}^{(1)} + \beta_0 \mathbf{d}^{(0)} = \begin{bmatrix} 3/2 \\ 0 \end{bmatrix} - \frac{9}{4} \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 3/2 \\ 9/4 \end{bmatrix}.$$

It is easy to verify that the directions  $\mathbf{d}^{(0)}$  and  $\mathbf{d}^{(1)}$  are  $\mathbf{Q}$ -conjugate. Indeed,

$$\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(1)} = \begin{bmatrix} 0 & -1 \end{bmatrix} \begin{bmatrix} 5 & -3 \\ -3 & 2 \end{bmatrix} \begin{bmatrix} 3/2 \\ 9/4 \end{bmatrix} = 0.$$

## 10.10

a. We have  $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{b}^\top \mathbf{x}$  where

$$\mathbf{Q} = \begin{bmatrix} 5 & 2 \\ 2 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 3 \\ 1 \end{bmatrix}.$$

b. Since  $f$  is a quadratic function on  $\mathbb{R}^2$ , we need to perform only two iterations. For the first iteration we compute

$$\begin{aligned} \mathbf{d}^{(0)} &= -\mathbf{g}^{(0)} = [3, 1]^\top \\ \alpha_0 &= \frac{5}{29} \\ \mathbf{x}^{(1)} &= [0.51724, 0.17241]^\top \\ \mathbf{g}^{(1)} &= [-0.06897, 0.20690]^\top. \end{aligned}$$

For the second iteration we compute

$$\begin{aligned} \beta_0 &= 0.0047534 \\ \mathbf{d}^{(1)} &= [0.08324, -0.20214]^\top \\ \alpha_1 &= 5.7952 \\ \mathbf{x}^{(2)} &= [1.000, -1.000]^\top. \end{aligned}$$

c. The minimizer is given by  $\mathbf{x}^* = \mathbf{Q}^{-1}\mathbf{b} = [1, -1]^\top$ , which agrees with part b.

### 10.11

A MATLAB routine for the conjugate gradient algorithm with options for different formulas of  $\beta_k$  is:

```
function [x,N]=conj_grad(grad,xnew,options);
%   CONJ_GRAD('grad',x0);
%   CONJ_GRAD('grad',x0,OPTIONS);
%
%   x = CONJ_GRAD('grad',x0);
%   x = CONJ_GRAD('grad',x0,OPTIONS);
%
%   [x,N] = CONJ_GRAD('grad',x0);
%   [x,N] = CONJ_GRAD('grad',x0,OPTIONS);
%
%The first variant finds the minimizer of a function whose gradient
%is described in grad (usually an M-file: grad.m), using initial point
%x0.
%The second variant allows a vector of optional parameters to be
%defined:
%OPTIONS(1) controls how much display output is given; set
%to 1 for a tabular display of results, (default is no display: 0).
%OPTIONS(2) is a measure of the precision required for the final point.
%OPTIONS(3) is a measure of the precision required of the gradient.
%OPTIONS(5) specifies the formula for beta:
%   0=Powell;
%   1=Fletcher-Reeves;
%   2=Polak-Ribiere;
%   3=Hestenes-Stiefel.
%OPTIONS(14) is the maximum number of iterations.
%For more information type HELP FOPTIONS.
%
%The next two variants return the value of the final point.
%The last two variants return a vector of the final point and the
%number of iterations.

if nargin ~= 3
    options = [];
    if nargin ~= 2
        disp('Wrong number of arguments. ');
        return;
    end
```

```

end

numvars = length(xnew);
if length(options) >= 14
    if options(14)==0
        options(14)=1000*numvars;
    end
else
    options(14)=1000*numvars;
end

clc;
format compact;
format short e;

options = foptions(options);
print = options(1);
epsilon_x = options(2);
epsilon_g = options(3);
max_iter=options(14);

g_curr=feval(grad,xnew);
if norm(g_curr) <= epsilon_g
    disp('Terminating: Norm of initial gradient less than');
    disp(epsilon_g);
    return;
end %if

d=-g_curr;
reset_cnt = 0;
for k = 1:max_iter,

    xcurr=xnew;
    alpha=linesearch_secant(grad,xcurr,d);
    %alpha=-(d'*g_curr)/(d'*Q*d);
    xnew = xcurr+alpha*d;

    if print,
        disp('Iteration number k =')
        disp(k); %print iteration index k
        disp('alpha =');
        disp(alpha); %print alpha
        disp('Gradient = ');
        disp(g_curr'); %print gradient
        disp('New point =');
        disp(xnew'); %print new point
    end %if

    if norm(xnew-xcurr) <= epsilon_x*norm(xcurr)
        disp('Terminating: Norm of difference between iterates less than');
        disp(epsilon_x);
        break;
    end %if

    g_old=g_curr;
    g_curr=feval(grad,xnew);

    if norm(g_curr) <= epsilon_g
        disp('Terminating: Norm of gradient less than');

```

```

        disp(epsilon_g);
        break;
    end %if

    reset_cnt = reset_cnt+1;
    if reset_cnt == 3*numvars
        d=-g_curr;
        reset_cnt = 0;
    else
        if options(5)==0 %Powell
            beta = max(0,(g_curr'*(g_curr-g_old))/(g_old'*g_old));
        elseif options(5)==1 %Fletcher-Reeves
            beta = (g_curr'*g_curr)/(g_old'*g_old);
        elseif options(5)==2 %Polak-Ribiere
            beta = (g_curr'*(g_curr-g_old))/(g_old'*g_old);
        else %Hestenes-Stiefel
            beta = (g_curr'*(g_curr-g_old))/(d'*(g_curr-g_old));
        end %if
        d=-g_curr+beta*d;
    end

    if print,
        disp('New beta =');
        disp(beta);
        disp('New d =');
        disp(d);
    end

    if k == max_iter
        disp('Terminating with maximum number of iterations');
    end %if
end %for

if nargout >= 1
    x=xnew;
    if nargout == 2
        N=k;
    end
else
    disp('Final point =');
    disp(xnew);
    disp('Number of iterations =');
    disp(k);
end %if
%-----

```

We created the following M-file, `g.m`, for the gradient of Rosenbrock's function:

```

function y=g(x)
y=[-400*(x(2)-x(1).^2)*x(1)-2*(1-x(1)), 200*(x(2)-x(1).^2)]';

```

We tested the above routine as follows:

```

>> options(2)=10^(-7);
>> options(3)=10^(-7);
>> options(14)=100;
>> options(5)=0;
>> conj_grad('g',[-2;2],options);

```

Terminating: Norm of difference between iterates less than



```

1.0000e-07
Final_Point =
1.0000e+00 1.0000e+00
Number_of_iteration =
8
>> options(5)=1;
>> conj_grad('g',[-2;2],options);

Terminating: Norm of difference between iterates less than
1.0000e-07
Final_Point =
1.0000e+00 1.0000e+00
Number_of_iteration =
10
>> options(5)=2;
>> conj_grad('g',[-2;2],options);

Terminating: Norm of difference between iterates less than
1.0000e-07
Final_Point =
1.0000e+00 1.0000e+00
Number_of_iteration =
8
>> options(5)=3;
>> conj_grad('g',[-2;2],options);

Terminating: Norm of difference between iterates less than
1.0000e-07
Final_Point =
1.0000e+00 1.0000e+00
Number_of_iteration =
8

```

The reader is cautioned *not* to draw any conclusions about the superiority or inferiority of any of the formulas for  $\beta_k$  based only on the above single numerical experiment.

## 11. Quasi-Newton Methods

### 11.1

---

a. Let

$$\phi(\alpha) = f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}).$$

Then, using the chain rule, we obtain

$$\phi'(\alpha) = \mathbf{d}^{(k)\top} \nabla f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}).$$

Hence

$$\phi'(0) = \mathbf{d}^{(k)\top} \mathbf{g}^{(k)}.$$

Since  $\phi'$  is continuous, then, if  $\mathbf{d}^{(k)\top} \mathbf{g}^{(k)} < 0$ , there exists  $\bar{\alpha} > 0$  such that for all  $\alpha \in (0, \bar{\alpha}]$ ,  $\phi(\alpha) < \phi(0)$ , i.e.,  $f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) < f(\mathbf{x}^{(k)})$ .

b. By part a,  $\phi(\alpha) < \phi(0)$  for all  $\alpha \in (0, \bar{\alpha}]$ . Hence,

$$\alpha_k = \arg \min_{\alpha \geq 0} \phi(\alpha) \neq 0$$

which implies that  $\alpha_k > 0$ .

c. Now,

$$\mathbf{d}^{(k)\top} \mathbf{g}^{(k+1)} = \mathbf{d}^{(k)\top} \nabla f(\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}) = \phi'_k(\alpha_k).$$

Since  $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) > 0$ , we have  $\phi'_k(\alpha_k) = 0$ . Hence,  $\mathbf{g}^{(k+1)\top} \mathbf{d}^{(k)} = 0$ .

d.

i. We have  $\mathbf{d}^{(k)} = -\mathbf{g}^{(k)}$ . Hence,  $\mathbf{d}^{(k)\top} \mathbf{g}^{(k)} = -\|\mathbf{g}^{(k)}\|^2$ . If  $\mathbf{g}^{(k)} \neq \mathbf{0}$ , then  $\|\mathbf{g}^{(k)}\|^2 > 0$ , and hence  $\mathbf{d}^{(k)\top} \mathbf{g}^{(k)} < 0$ .

ii. We have  $\mathbf{d}^{(k)} = -\mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}$ . Since  $\mathbf{F}(\mathbf{x}^{(k)}) > 0$ , we also have  $\mathbf{F}(\mathbf{x}^{(k)})^{-1} > 0$ . Therefore,  $\mathbf{d}^{(k)\top} \mathbf{g}^{(k)} = -\mathbf{g}^{(k)\top} \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)} < 0$  if  $\mathbf{g}^{(k)} \neq \mathbf{0}$ .

iii. We have

$$\mathbf{d}^{(k)} = -\mathbf{g}^{(k)} + \beta_{k-1} \mathbf{d}^{(k-1)}.$$

Hence,

$$\mathbf{d}^{(k)\top} \mathbf{g}^{(k)} = -\|\mathbf{g}^{(k)}\|^2 + \beta_{k-1} \mathbf{d}^{(k-1)\top} \mathbf{g}^{(k)}.$$

By part c,  $\mathbf{d}^{(k-1)\top} \mathbf{g}^{(k)} = 0$ . Hence, if  $\mathbf{g}^{(k)} \neq \mathbf{0}$ , then  $\|\mathbf{g}^{(k)}\|^2 > 0$ , and

$$\mathbf{d}^{(k)\top} \mathbf{g}^{(k)} = -\|\mathbf{g}^{(k)}\|^2 < 0.$$

iv. We have  $\mathbf{d}^{(k)} = -\mathbf{H}_k \mathbf{g}^{(k)}$ . Therefore, if  $\mathbf{H}_k > 0$  and  $\mathbf{g}^{(k)} \neq \mathbf{0}$ , then  $\mathbf{d}^{(k)\top} \mathbf{g}^{(k)} = -\mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)} < 0$ .

e. Using the equation  $\nabla f(\mathbf{x}) = \mathbf{Q}\mathbf{x} - \mathbf{b}$ , we get

$$\begin{aligned} \mathbf{d}^{(k)\top} \mathbf{g}^{(k+1)} &= \mathbf{d}^{(k)\top} (\mathbf{Q}\mathbf{x}^{(k+1)} - \mathbf{b}) \\ &= \mathbf{d}^{(k)\top} (\mathbf{Q}(\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}) - \mathbf{b}) \\ &= \alpha_k \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)} + \mathbf{d}^{(k)\top} (\mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b}) \\ &= \alpha_k \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)} + \mathbf{d}^{(k)\top} \mathbf{g}^{(k)}. \end{aligned}$$

By part c,  $\mathbf{d}^{(k)\top} \mathbf{g}^{(k+1)} = 0$ , which implies

$$\alpha_k = -\frac{\mathbf{d}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}}.$$

## 11.2

Yes, because:

1. The search direction is of the form  $\mathbf{d}^{(k)} = \mathbf{H}_k \nabla f(\mathbf{x}^{(k)})$  for matrix  $\mathbf{H}_k = \mathbf{F}(\mathbf{x}^{(k)})^{-1}$ ;
2. The matrix  $\mathbf{H}_k = \mathbf{F}(\mathbf{x}^{(k)})^{-1}$  is symmetric for  $f \in \mathcal{C}^2$ ;
3. If  $f$  is quadratic, then the quasi-Newton condition is satisfied:  $\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}$ ,  $0 \leq i \leq k$ . To see this, note that if the Hessian is  $\mathbf{Q}$ , then  $\mathbf{Q} \Delta \mathbf{x}^{(i)} = \Delta \mathbf{g}^{(i)}$ . Multiplying both sides by  $\mathbf{H}_k = \mathbf{Q}^{-1}$ , we obtain the desired result.

## 11.3

a. We have

$$f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) = \frac{1}{2} (\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})^\top \mathbf{Q} (\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) - (\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})^\top \mathbf{b} + c.$$

Using the chain rule, we obtain

$$\frac{d}{d\alpha} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) = (\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})^\top \mathbf{Q} \mathbf{d}^{(k)} - \mathbf{d}^{(k)\top} \mathbf{b}.$$

Equating the above to zero and solving for  $\alpha$  gives

$$\left(\mathbf{x}^{(k)\top} \mathbf{Q} - \mathbf{b}^\top\right) \mathbf{d}^{(k)} = -\alpha \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}.$$

Taking into account that  $\mathbf{g}^{(k)\top} = \mathbf{x}^{(k)\top} \mathbf{Q} - \mathbf{b}^\top$  and that  $\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)} > 0$  for  $\mathbf{g}^{(k)} \neq \mathbf{0}$ , we obtain

$$\alpha_k = -\frac{\mathbf{g}^{(k)\top} \mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}} = \frac{\mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}}.$$

b. The matrix  $\mathbf{Q}$  is symmetric and positive definite; hence  $\alpha_k > 0$  if  $\mathbf{H}_k = \mathbf{H}_k^\top > 0$ .

#### 11.4

a. The appropriate choice is  $\mathbf{H} = \mathbf{F}(\mathbf{x}^*)^{-1}$ . To show this, we can apply the same argument as in the proof of the theorem on the convergence of Newton's method. (We won't repeat it here.)

b. Yes (provided we incorporate the usual step size). Indeed, if we apply the algorithm with the choice of  $\mathbf{H}$  in part a, then when applied to a quadratic with Hessian  $\mathbf{Q}$ , the algorithm uses  $\mathbf{H} = \mathbf{Q}^{-1}$ , which definitely satisfies the quasi-Newton condition. In fact, the algorithm then behaves just like Newton's algorithm.

#### 11.5

Our objective is to minimize the quadratic

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{x}^\top \mathbf{b} + c.$$

We first compute the gradient  $\nabla f$  and evaluate it at  $\mathbf{x}^{(0)}$ ,

$$\nabla f(\mathbf{x}^{(0)}) = \mathbf{g}^{(0)} = \mathbf{Q} \mathbf{x}^{(0)} - \mathbf{b} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}.$$

It is a non-zero vector, so we proceed with the first iteration. Let  $\mathbf{H}_0 = \mathbf{I}_2$ . Then,

$$\mathbf{d}^{(0)} = -\mathbf{H}_0 \mathbf{g}^{(0)} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}.$$

The step size  $\alpha_0$  is

$$\alpha_0 = -\frac{\mathbf{g}^{(0)\top} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} = -\frac{\begin{bmatrix} -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix}}{\begin{bmatrix} 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix}} = \frac{2}{3}.$$

Hence,

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = \begin{bmatrix} 2/3 \\ -2/3 \end{bmatrix}.$$

We evaluate the gradient  $\nabla f$  and evaluate it at  $\mathbf{x}^{(1)}$  to obtain

$$\nabla f(\mathbf{x}^{(1)}) = \mathbf{g}^{(1)} = \mathbf{Q} \mathbf{x}^{(1)} - \mathbf{b} = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 2/3 \\ -2/3 \end{bmatrix} - \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} -1/3 \\ -1/3 \end{bmatrix}.$$

It is a non-zero vector, so we proceed with the second iteration. We compute  $\mathbf{H}_1$ , where

$$\mathbf{H}_1 = \mathbf{H}_0 + \frac{(\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)}) (\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})^\top}{\Delta \mathbf{g}^{(0)\top} (\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})}.$$

To find  $\mathbf{H}_1$  we need to compute,

$$\Delta \mathbf{x}^{(0)} = \mathbf{x}^{(1)} - \mathbf{x}^{(0)} = \begin{bmatrix} 2/3 \\ -2/3 \end{bmatrix} \quad \text{and} \quad \Delta \mathbf{g}^{(0)} = \mathbf{g}^{(1)} - \mathbf{g}^{(0)} = \begin{bmatrix} 2/3 \\ -4/3 \end{bmatrix}.$$

Using the above, we determine,

$$\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)} = \begin{bmatrix} 0 \\ 2/3 \end{bmatrix} \quad \text{and} \quad \Delta \mathbf{g}^{(0)\top} (\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)}) = -\frac{8}{9}.$$

Then, we obtain

$$\begin{aligned} \mathbf{H}_1 &= \mathbf{H}_0 + \frac{(\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)}) (\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})^\top}{\Delta \mathbf{g}^{(0)\top} (\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})} \\ &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \frac{\begin{bmatrix} 0 & 0 \\ 0 & 4/9 \end{bmatrix}}{-8/9} \\ &= \begin{bmatrix} 1 & 0 \\ 0 & 1/2 \end{bmatrix} \end{aligned}$$

and

$$\mathbf{d}^{(1)} = -\mathbf{H}_1 \mathbf{g}^{(1)} = \begin{bmatrix} 1/3 \\ 1/6 \end{bmatrix}.$$

We next compute

$$\alpha_1 = -\frac{\mathbf{g}^{(1)\top} \mathbf{d}^{(1)}}{\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(1)}} = 1.$$

Therefore,

$$\mathbf{x}^{(2)} = \mathbf{x}^* = \mathbf{x}^{(1)} + \alpha_1 \mathbf{d}^{(1)} = \begin{bmatrix} 1 \\ -1/2 \end{bmatrix}.$$

Note that  $\mathbf{g}^{(2)} = \mathbf{Q} \mathbf{x}^{(2)} - \mathbf{b} = \mathbf{0}$  as expected.

---

### 11.6

We are guaranteed that the step size satisfies  $\alpha_k > 0$  if the search direction is in the descent direction, i.e., the search direction  $\mathbf{d}^{(k)} = -\mathbf{M}_k \nabla f(\mathbf{x}^{(k)})$  has strictly positive inner product with  $-\nabla f(\mathbf{x}^{(k)})$  (see Exercise 11.1). Thus, the condition on  $\mathbf{M}_k$  that guarantees  $\alpha_k > 0$  is  $\nabla f(\mathbf{x}^{(k)})^\top \mathbf{M}_k \nabla f(\mathbf{x}^{(k)}) > 0$ , which corresponds to  $1 + a > 0$ , or  $a > -1$ . (Note that if  $a \leq -1$ , the search direction is not in the descent direction, and thus we cannot guarantee that  $\alpha_k > 0$ .)

---

### 11.7

Let  $\mathbf{x} \in \mathbb{R}^n$ . Then

$$\begin{aligned} \mathbf{x}^\top \mathbf{H}_{k+1} \mathbf{x} &= \mathbf{x}^\top \mathbf{H}_k \mathbf{x} + \mathbf{x}^\top \left( \frac{(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)}) (\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})^\top}{\Delta \mathbf{g}^{(k)\top} (\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})} \right) \mathbf{x} \\ &= \mathbf{x}^\top \mathbf{H}_k \mathbf{x} + \frac{(\mathbf{x}^\top (\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)}))^2}{\Delta \mathbf{g}^{(k)\top} (\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})}. \end{aligned}$$

Note that since  $\mathbf{H}_k > 0$ , we have  $\mathbf{x}^\top \mathbf{H}_k \mathbf{x} > 0$ . Hence, if  $\Delta \mathbf{g}^{(k)\top} (\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)}) > 0$ , then  $\mathbf{x}^\top \mathbf{H}_{k+1} \mathbf{x} > 0$ .

---

### 11.8

The complement of the Rank One update equation is

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{(\Delta \mathbf{g}^{(k)} - \mathbf{B}_k \Delta \mathbf{x}^{(k)}) (\Delta \mathbf{g}^{(k)} - \mathbf{B}_k \Delta \mathbf{x}^{(k)})^\top}{\Delta \mathbf{x}^{(k)\top} (\Delta \mathbf{g}^{(k)} - \mathbf{B}_k \Delta \mathbf{x}^{(k)})}.$$

Using the matrix inverse formula, we get

$$\begin{aligned}
& \mathbf{B}_{k+1}^{-1} \\
&= \mathbf{B}_k^{-1} \\
&- \frac{\mathbf{B}_k^{-1}(\Delta \mathbf{g}^{(k)} - \mathbf{B}_k \Delta \mathbf{x}^{(k)})(\Delta \mathbf{g}^{(k)} - \mathbf{B}_k \Delta \mathbf{x}^{(k)}) \mathbf{B}_k^{-1}}{\Delta \mathbf{x}^{(k)\top}(\Delta \mathbf{g}^{(k)} - \mathbf{B}_k \Delta \mathbf{x}^{(k)}) + (\Delta \mathbf{g}^{(k)} - \mathbf{B}_k \Delta \mathbf{x}^{(k)})^\top \mathbf{B}_k^{-1}(\Delta \mathbf{g}^{(k)} - \mathbf{B}_k \Delta \mathbf{x}^{(k)})} \\
&= \mathbf{B}_k^{-1} + \frac{(\Delta \mathbf{x}^{(k)} - \mathbf{B}_k^{-1} \Delta \mathbf{g}^{(k)})(\Delta \mathbf{x}^{(k)} - \mathbf{B}_k^{-1} \Delta \mathbf{g}^{(k)})^\top}{\Delta \mathbf{g}^{(k)\top}(\Delta \mathbf{x}^{(k)} - \mathbf{B}_k^{-1} \Delta \mathbf{g}^{(k)})}.
\end{aligned}$$

Substituting  $\mathbf{H}_k$  for  $\mathbf{B}_k^{-1}$ , we get a formula identical to the Rank One update equation. This should not be surprising, since there is only one update equation involving a rank one correction that satisfies the quasi-Newton condition.

### 11.9

We first compute the gradient  $\nabla f$  and evaluate it at  $\mathbf{x}^{(0)}$ ,

$$\nabla f(\mathbf{x}^{(0)}) = \mathbf{g}^{(0)} = \mathbf{Q}\mathbf{x}^{(0)} - \mathbf{b} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}.$$

It is a nonzero vector, so we proceed with the first iteration. Let  $\mathbf{H}_0 = \mathbf{I}_2$ . Then,

$$\mathbf{d}^{(0)} = -\mathbf{H}_0 \mathbf{g}^{(0)} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}.$$

The step size  $\alpha_0$  is

$$\alpha_0 = -\frac{\mathbf{g}^{(0)\top} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} = -\frac{\begin{bmatrix} -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix}}{\begin{bmatrix} 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix}} = \frac{2}{3}.$$

Hence,

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = \begin{bmatrix} 2/3 \\ -2/3 \end{bmatrix}.$$

We evaluate the gradient  $\nabla f$  and evaluate it at  $\mathbf{x}^{(1)}$  to obtain

$$\nabla f(\mathbf{x}^{(1)}) = \mathbf{g}^{(1)} = \mathbf{Q}\mathbf{x}^{(1)} - \mathbf{b} = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 2/3 \\ -2/3 \end{bmatrix} - \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} -1/3 \\ -1/3 \end{bmatrix}.$$

It is a nonzero vector, so we proceed with the second iteration. We compute  $\mathbf{H}_1$ , where

$$\mathbf{H}_1 = \mathbf{H}_0 + \frac{(\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})(\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})^\top}{\Delta \mathbf{g}^{(0)\top}(\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})}.$$

To find  $\mathbf{H}_1$  we need

$$\Delta \mathbf{x}^{(0)} = \mathbf{x}^{(1)} - \mathbf{x}^{(0)} = \begin{bmatrix} 2/3 \\ -2/3 \end{bmatrix} \quad \text{and} \quad \Delta \mathbf{g}^{(0)} = \mathbf{g}^{(1)} - \mathbf{g}^{(0)} = \begin{bmatrix} 2/3 \\ -4/3 \end{bmatrix}.$$

Using the above, we determine,

$$\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)} = \begin{bmatrix} 0 \\ 2/3 \end{bmatrix} \quad \text{and} \quad \Delta \mathbf{g}^{(0)\top}(\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)}) = -\frac{8}{9}.$$

Then, we obtain

$$\begin{aligned}
\mathbf{H}_1 &= \mathbf{H}_0 + \frac{(\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)}) (\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})^\top}{\Delta \mathbf{g}^{(0)\top} (\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})} \\
&= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \frac{\begin{bmatrix} 0 & 0 \\ 0 & 4/9 \end{bmatrix}}{-8/9} \\
&= \begin{bmatrix} 1 & 0 \\ 0 & 1/2 \end{bmatrix}
\end{aligned}$$

and

$$\mathbf{d}^{(1)} = -\mathbf{H}_1 \mathbf{g}^{(1)} = \begin{bmatrix} 1/3 \\ 1/6 \end{bmatrix}.$$

Note that  $\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(1)} = 0$ , that is,  $\mathbf{d}^{(0)}$  and  $\mathbf{d}^{(1)}$  are  $\mathbf{Q}$ -conjugate.

---

#### 11.10

The calculations are similar until we get to the second step:

$$\begin{aligned}
\mathbf{H}_1 &= \begin{bmatrix} 1/2 & -1/2 \\ -1/2 & 1/2 \end{bmatrix} \\
\mathbf{d}^{(0)} &= \mathbf{0}.
\end{aligned}$$

So the algorithm gets stuck at this point, which illustrates that it doesn't work.

---

#### 11.11

a. Since  $f$  is quadratic, and  $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$ , then

$$\alpha_k = -\frac{\mathbf{g}^{(k)\top} \mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}}.$$

b. Now,  $\mathbf{d}^{(k)} = -\mathbf{H}_k \mathbf{g}^{(k)}$ , where  $\mathbf{H}_k = \mathbf{H}_k^\top > 0$ . Substituting this into the formula for  $\alpha_k$  in part a, yields

$$\alpha_k = \frac{\mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}} > 0.$$

---

#### 11.12

(Our solution to this problem is based on a solution that was furnished to us by Michael Mera, a student in ECE 580 at Purdue in Spring 2005.) To proceed, we recall the formula of Lemma 11.1,

$$(\mathbf{A} + \mathbf{u} \mathbf{v}^\top)^{-1} = \mathbf{A}^{-1} - \frac{(\mathbf{A}^{-1} \mathbf{u})(\mathbf{v}^\top \mathbf{A}^{-1})}{1 + \mathbf{v}^\top \mathbf{A}^{-1} \mathbf{u}}$$

for  $1 + \mathbf{v}^\top \mathbf{A}^{-1} \mathbf{u} \neq 0$ . Recall the definitions from the hint,

$$\mathbf{A}_0 = \mathbf{B}_k, \quad \mathbf{u}_0 = \frac{\Delta \mathbf{g}^{(k)}}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)}}, \quad \mathbf{v}_0^\top = \Delta \mathbf{g}^{(k)\top},$$

and

$$\mathbf{A}_1 = \mathbf{B}_k + \frac{\Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top}}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)}} = \mathbf{A}_0 + \mathbf{u}_0 \mathbf{v}_0^\top, \quad \mathbf{u}_1 = -\frac{\mathbf{B}_k \Delta \mathbf{x}^{(k)}}{\Delta \mathbf{x}^{(k)\top} \mathbf{B}_k \Delta \mathbf{x}^{(k)}},$$

and

$$\mathbf{v}_1^\top = \Delta \mathbf{x}^{(k)\top} \mathbf{B}_k.$$

Using the above notation, we represent  $\mathbf{B}_{k+1}$  as

$$\begin{aligned}\mathbf{B}_{k+1} &= \mathbf{A}_0 + \mathbf{u}_0 \mathbf{v}_0^\top + \mathbf{u}_1 \mathbf{v}_1^\top \\ &= \mathbf{A}_1 + \mathbf{u}_1 \mathbf{v}_1^\top.\end{aligned}$$

Applying to the above Lemma 11.1 gives

$$\begin{aligned}\mathbf{H}_{k+1}^{BFGS} &= (\mathbf{B}_{k+1})^{-1} \\ &= (\mathbf{A}_1 + \mathbf{u}_1 \mathbf{v}_1^\top)^{-1} \\ &= \mathbf{A}_1^{-1} - \frac{\mathbf{A}_1^{-1} \mathbf{u}_1 \mathbf{v}_1^\top \mathbf{A}_1^{-1}}{1 + \mathbf{v}_1^\top \mathbf{A}_1^{-1} \mathbf{u}_1}.\end{aligned}$$

Substituting into the above the expression for  $\mathbf{A}_1^{-1}$  yields

$$\mathbf{H}_{k+1}^{BFGS} = \mathbf{A}_0^{-1} - \frac{\mathbf{A}_0^{-1} \mathbf{u}_0 \mathbf{v}_0^\top \mathbf{A}_0^{-1}}{1 + \mathbf{v}_0^\top \mathbf{A}_0^{-1} \mathbf{u}_0} - \frac{\left( \mathbf{A}_0^{-1} - \frac{\mathbf{A}_0^{-1} \mathbf{u}_0 \mathbf{v}_0^\top \mathbf{A}_0^{-1}}{1 + \mathbf{v}_0^\top \mathbf{A}_0^{-1} \mathbf{u}_0} \right) \mathbf{u}_1 \mathbf{v}_1^\top \left( \mathbf{A}_0^{-1} - \frac{\mathbf{A}_0^{-1} \mathbf{u}_0 \mathbf{v}_0^\top \mathbf{A}_0^{-1}}{1 + \mathbf{v}_0^\top \mathbf{A}_0^{-1} \mathbf{u}_0} \right)}{1 + \mathbf{v}_1^\top \left( \mathbf{A}_0^{-1} - \frac{\mathbf{A}_0^{-1} \mathbf{u}_0 \mathbf{v}_0^\top \mathbf{A}_0^{-1}}{1 + \mathbf{v}_0^\top \mathbf{A}_0^{-1} \mathbf{u}_0} \right) \mathbf{u}_1}.$$

Note that  $\mathbf{A}_0 = \mathbf{B}_k$ . Hence,  $\mathbf{A}_0^{-1} = \mathbf{B}_k^{-1} = \mathbf{H}_k$ . Using this and the notation introduced at the beginning of the solution, we obtain

$$\begin{aligned}\mathbf{H}_{k+1}^{BFGS} &= \mathbf{H}_k - \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)} + \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} \\ &\quad - \frac{\left( \mathbf{H}_k - \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)} + \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} \right) \left( \frac{-\mathbf{B}_k \Delta \mathbf{x}^{(k)} \Delta \mathbf{x}^{(k)\top} \mathbf{B}_k}{\Delta \mathbf{x}^{(k)\top} \mathbf{B}_k \Delta \mathbf{x}^{(k)}} \right)}{1 + \Delta \mathbf{x}^{(k)\top} \mathbf{B}_k \left( \mathbf{H}_k - \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)} + \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} \right) \left( \frac{-\mathbf{B}_k \Delta \mathbf{x}^{(k)}}{\Delta \mathbf{x}^{(k)\top} \mathbf{B}_k \Delta \mathbf{x}^{(k)}} \right)} \\ &\quad \times \left( \mathbf{H}_k - \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)} + \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} \right).\end{aligned}$$

We next perform some multiplications taking into account that  $\mathbf{H}_k = \mathbf{B}_k^{-1}$  and hence

$$\mathbf{H}_k \mathbf{B}_k = \mathbf{B}_k \mathbf{H}_k = \mathbf{I}_n.$$

We obtain

$$\begin{aligned}\mathbf{H}_{k+1}^{BFGS} &= \mathbf{H}_k - \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)} + \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} \\ &\quad - \frac{\left( 1 - \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)} + \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} \right) (-\Delta \mathbf{x}^{(k)} \Delta \mathbf{x}^{(k)\top}) \left( 1 - \frac{\Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)} + \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} \right)}{\Delta \mathbf{x}^{(k)\top} \mathbf{B}_k \Delta \mathbf{x}^{(k)} + \Delta \mathbf{x}^{(k)\top} \left( \mathbf{B}_k - \frac{\Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top}}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)} + \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} \right) (-\Delta \mathbf{x}^{(k)})}.\end{aligned}$$

We proceed with our manipulations. We first perform multiplications by  $\Delta \mathbf{x}^{(k)}$  and  $\Delta \mathbf{x}^{(k)\top}$  to obtain

$$\begin{aligned}\mathbf{H}_{k+1}^{BFGS} &= \mathbf{H}_k - \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)} + \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} \\ &\quad - \frac{\left( \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)}}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)} + \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} - \Delta \mathbf{x}^{(k)} \right) \left( \Delta \mathbf{x}^{(k)\top} - \frac{\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)} + \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} \right)}{\Delta \mathbf{x}^{(k)\top} \mathbf{B}_k \Delta \mathbf{x}^{(k)} - \Delta \mathbf{x}^{(k)\top} \mathbf{B}_k \Delta \mathbf{x}^{(k)} + \frac{\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)}}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)} + \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}}}.\end{aligned}$$

Cancelling the terms in the denominator of the last term above and performing further multiplications gives

$$\begin{aligned}
H_{k+1}^{BFGS} = H_k & - \frac{H_k \Delta g^{(k)} \Delta g^{(k)\top} H_k}{\Delta g^{(k)\top} \Delta x^{(k)} + \Delta g^{(k)\top} H_k \Delta g^{(k)}} \\
& + \frac{H_k \Delta g^{(k)} (\Delta g^{(k)\top} \Delta x^{(k)}) (\Delta x^{(k)\top} \Delta g^{(k)}) \Delta g^{(k)\top} H_k}{\Delta g^{(k)\top} \Delta x^{(k)} + \Delta g^{(k)\top} H_k \Delta g^{(k)}} \\
& + \frac{(\Delta x^{(k)\top} \Delta g^{(k)}) (\Delta g^{(k)\top} \Delta x^{(k)})}{(\Delta x^{(k)\top} \Delta g^{(k)}) (\Delta g^{(k)\top} \Delta x^{(k)})} \\
& + \frac{\Delta x^{(k)} \Delta x^{(k)\top} (\Delta g^{(k)\top} \Delta x^{(k)} + \Delta g^{(k)\top} H_k \Delta g^{(k)})}{(\Delta x^{(k)\top} \Delta g^{(k)}) (\Delta g^{(k)\top} \Delta x^{(k)})} \\
& - \frac{H_k \Delta g^{(k)} (\Delta g^{(k)\top} \Delta x^{(k)}) \Delta x^{(k)\top} + \Delta x^{(k)} \Delta g^{(k)\top} H_k}{(\Delta x^{(k)\top} \Delta g^{(k)}) (\Delta g^{(k)\top} \Delta x^{(k)})}.
\end{aligned}$$

Further simplification of the third and the fifth terms on the right hand-side of the above equation gives

$$\begin{aligned}
H_{k+1}^{BFGS} = H_k & - \frac{H_k \Delta g^{(k)} \Delta g^{(k)\top} H_k}{\Delta g^{(k)\top} \Delta x^{(k)} + \Delta g^{(k)\top} H_k \Delta g^{(k)}} \\
& + \frac{H_k \Delta g^{(k)} \Delta g^{(k)\top} H_k}{\Delta g^{(k)\top} \Delta x^{(k)} + \Delta g^{(k)\top} H_k \Delta g^{(k)}} \\
& + \frac{\Delta x^{(k)} \Delta x^{(k)\top} (\Delta g^{(k)\top} \Delta x^{(k)} + \Delta g^{(k)\top} H_k \Delta g^{(k)})}{(\Delta x^{(k)\top} \Delta g^{(k)}) (\Delta g^{(k)\top} \Delta x^{(k)})} \\
& - \frac{H_k \Delta g^{(k)} \Delta x^{(k)\top} + \Delta x^{(k)} \Delta g^{(k)\top} H_k}{\Delta x^{(k)\top} \Delta g^{(k)}}.
\end{aligned}$$

Note that the second and the third terms cancel out each other. We then represent the fourth term in alternative manner to obtain

$$\begin{aligned}
H_{k+1}^{BFGS} = H_k & + \frac{\Delta x^{(k)} \Delta x^{(k)\top}}{\Delta x^{(k)\top} \Delta g^{(k)}} \left( 1 + \frac{\Delta g^{(k)\top} H_k \Delta g^{(k)}}{\Delta g^{(k)\top} \Delta x^{(k)}} \right) \\
& - \frac{H_k \Delta g^{(k)} \Delta x^{(k)\top} + \Delta x^{(k)} \Delta g^{(k)\top} H_k}{\Delta x^{(k)\top} \Delta g^{(k)}},
\end{aligned}$$

which is the desired BFGS update formula.

### 11.13

The first step for both algorithms is clearly the same, since in either case we have

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \alpha_0 \mathbf{g}^{(0)}.$$

For the second step,

$$\begin{aligned}
\mathbf{d}^{(1)} & = -H_1 \mathbf{g}^{(1)} \\
& = - \left( \mathbf{I}_n + \left( 1 + \frac{\Delta \mathbf{g}^{(0)\top} \Delta \mathbf{g}^{(0)}}{\Delta \mathbf{g}^{(0)\top} \Delta \mathbf{x}^{(0)}} \right) \frac{\Delta \mathbf{x}^{(0)} \Delta \mathbf{x}^{(0)\top}}{\Delta \mathbf{x}^{(0)\top} \Delta \mathbf{g}^{(0)}} \right. \\
& \quad \left. - \frac{\Delta \mathbf{g}^{(0)} \Delta \mathbf{x}^{(0)\top} + (\Delta \mathbf{g}^{(0)} \Delta \mathbf{x}^{(0)\top})^\top}{\Delta \mathbf{g}^{(0)\top} \Delta \mathbf{x}^{(0)}} \right) \mathbf{g}^{(1)} \\
& = -\mathbf{g}^{(1)} - \left( 1 + \frac{\Delta \mathbf{g}^{(0)\top} \Delta \mathbf{g}^{(0)}}{\Delta \mathbf{g}^{(0)\top} \Delta \mathbf{x}^{(0)}} \right) \frac{\Delta \mathbf{x}^{(0)} \Delta \mathbf{x}^{(0)\top} \mathbf{g}^{(1)}}{\Delta \mathbf{x}^{(0)\top} \Delta \mathbf{g}^{(0)}} \\
& \quad + \frac{\Delta \mathbf{g}^{(0)} \Delta \mathbf{x}^{(0)\top} \mathbf{g}^{(1)} + \Delta \mathbf{x}^{(0)} \Delta \mathbf{g}^{(0)\top} \mathbf{g}^{(1)}}{\Delta \mathbf{g}^{(0)\top} \Delta \mathbf{x}^{(0)}}.
\end{aligned}$$

Since the line search is exact, we have

$$\Delta \mathbf{x}^{(0)\top} \mathbf{g}^{(1)} = \alpha_0 \mathbf{d}^{(0)\top} \mathbf{g}^{(1)} = 0.$$



Hence,

$$\begin{aligned}
\mathbf{d}^{(1)} &= -\mathbf{g}^{(1)} + \left( \frac{\Delta \mathbf{g}^{(0)\top} \mathbf{g}^{(1)}}{\Delta \mathbf{g}^{(0)\top} \Delta \mathbf{x}^{(0)}} \right) \Delta \mathbf{x}^{(0)} \\
&= -\mathbf{g}^{(1)} + \left( \frac{\mathbf{g}^{(1)\top} \Delta \mathbf{g}^{(0)}}{\Delta \mathbf{g}^{(0)\top} \mathbf{d}^{(0)}} \right) \mathbf{d}^{(0)} \\
&= -\mathbf{g}^{(1)} + \beta_0 \mathbf{d}^{(0)}
\end{aligned}$$

where

$$\beta_0 = \frac{\mathbf{g}^{(1)\top} \Delta \mathbf{g}^{(0)}}{\mathbf{d}^{(0)\top} \Delta \mathbf{g}^{(0)}} = \frac{\mathbf{g}^{(1)\top} (\mathbf{g}^{(1)} - \mathbf{g}^{(0)})}{\mathbf{d}^{(0)\top} (\mathbf{g}^{(1)} - \mathbf{g}^{(0)})}$$

is the Hestenes-Stiefel update formula for  $\beta_0$ . Since  $\mathbf{d}^{(0)} = -\mathbf{g}^{(0)}$ , and  $\mathbf{g}^{(1)\top} \mathbf{g}^{(0)} = 0$ , we have

$$\beta_0 = \frac{\mathbf{g}^{(1)\top} (\mathbf{g}^{(1)} - \mathbf{g}^{(0)})}{\mathbf{g}^{(0)\top} \mathbf{g}^{(0)}},$$

which is the Polak-Ribiere formula. Applying  $\mathbf{g}^{(1)\top} \mathbf{g}^{(0)} = 0$  again, we get

$$\beta_0 = \frac{\mathbf{g}^{(1)\top} \mathbf{g}^{(1)}}{\mathbf{g}^{(0)\top} \mathbf{g}^{(0)}},$$

which is the Fletcher-Reeves formula.

#### 11.14

a. Suppose the three conditions hold whenever applied to a quadratic. We need to show that when applied to a quadratic, for  $k = 0, \dots, n-1$  and  $i = 0, \dots, k$ ,  $\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}$ . For  $i = k$ , we have

$$\begin{aligned}
\mathbf{H}_{k+1} \Delta \mathbf{g}^{(k)} &= \mathbf{H}_k \Delta \mathbf{g}^{(k)} + \mathbf{U}_k \Delta \mathbf{g}^{(k)} \quad \text{by condition 1} \\
&= \mathbf{H}_k \Delta \mathbf{g}^{(k)} + \Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)} \quad \text{by condition 2} \\
&= \Delta \mathbf{x}^{(k)},
\end{aligned}$$

as required. For the rest of the proof ( $i = 0, \dots, k-1$ ), we use induction on  $k$ .

For  $k = 0$ , there is nothing to prove (covered by the  $i = k$  case). So suppose the result holds for  $k-1$ . To show the result for  $k$ , first fix  $i \in \{0, \dots, k-1\}$ . We have

$$\begin{aligned}
\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} &= \mathbf{H}_k \Delta \mathbf{g}^{(i)} + \mathbf{U}_k \Delta \mathbf{g}^{(i)} \\
&= \Delta \mathbf{x}^{(i)} + \mathbf{U}_k \Delta \mathbf{g}^{(i)} \quad \text{by the induction hypothesis} \\
&= \Delta \mathbf{x}^{(i)} + \mathbf{a}^{(k)} \Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(i)} + \mathbf{b}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(i)} \quad \text{by condition 3.}
\end{aligned}$$

So it suffices to show that the second and third terms are both 0. For the second term,

$$\begin{aligned}
\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(i)} &= \Delta \mathbf{x}^{(k)\top} \mathbf{Q} \Delta \mathbf{x}^{(i)} \\
&= \alpha_k \alpha_i \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(i)} \\
&= 0
\end{aligned}$$

because of the induction hypothesis, which implies  $\mathbf{Q}$ -conjugacy (where  $\mathbf{Q}$  is the Hessian of the given quadratic). Similarly, for the third term,

$$\begin{aligned}
\Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(i)} &= \Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(i)} \quad \text{by the induction hypothesis} \\
&= \Delta \mathbf{x}^{(k)\top} \mathbf{Q} \Delta \mathbf{x}^{(i)} \\
&= \alpha_k \alpha_i \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(i)} \\
&= 0,
\end{aligned}$$

again because of the induction hypothesis, which implies  $\mathbf{Q}$ -conjugacy. This completes the proof.

b. All three algorithms satisfy the conditions in part a. Condition 1 holds, as described in class. Condition 2 is straightforward to check for all three algorithms. For the rank-one and DFP algorithms, this is shown in the book. For BFGS, some simple matrix algebra establishes that it holds. Condition 3 holds by appropriate definition of the vectors  $\mathbf{a}^{(k)}$  and  $\mathbf{b}^{(k)}$ . In particular, for the rank-one algorithm,

$$\mathbf{a}^{(k)} = \frac{(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})}{(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})^\top \Delta \mathbf{g}^{(k)}}, \quad \mathbf{b}^{(k)} = -\frac{(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})}{(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})^\top \Delta \mathbf{g}^{(k)}}.$$

For the DFP algorithm,

$$\mathbf{a}^{(k)} = \frac{\Delta \mathbf{x}^{(k)}}{\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)}}, \quad \mathbf{b}^{(k)} = -\frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)}}{\Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}}.$$

Finally, for the BFGS algorithm,

$$\mathbf{a}^{(k)} = \left(1 + \frac{\Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)}}\right) \frac{\Delta \mathbf{x}^{(k)}}{\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)}} - \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)}}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)}}, \quad \mathbf{b}^{(k)} = \frac{\Delta \mathbf{x}^{(k)}}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)}}.$$

### 11.15

a. Suppose we apply the algorithm to a quadratic. Then, by the quasi-Newton property of DFP, we have  $\mathbf{H}_{k+1}^{DFP} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}$ ,  $0 \leq i \leq k$ . The same holds for BFGS. Thus, for the given  $\mathbf{H}_k$ , we have for  $0 \leq i \leq k$ ,

$$\begin{aligned} \mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} &= \phi \mathbf{H}_{k+1}^{DFP} \Delta \mathbf{g}^{(i)} + (1 - \phi) \mathbf{H}_{k+1}^{BFGS} \Delta \mathbf{g}^{(i)} \\ &= \phi \Delta \mathbf{x}^{(i)} + (1 - \phi) \Delta \mathbf{x}^{(i)} \\ &= \Delta \mathbf{x}^{(i)}, \end{aligned}$$

which shows that the above algorithm is a quasi-Newton algorithm (and hence also a conjugate direction algorithm).

b. By Theorem 11.4 and the discussion on BFGS, we have  $\mathbf{H}_k^{DFP} > 0$  and  $\mathbf{H}_k^{BFGS} > 0$ . Hence, for any  $\mathbf{x} \neq \mathbf{0}$ ,

$$\mathbf{x}^\top \mathbf{H}_k \mathbf{x} = \phi \mathbf{x}^\top \mathbf{H}_k^{DFP} \mathbf{x} + (1 - \phi) \mathbf{x}^\top \mathbf{H}_k^{BFGS} \mathbf{x} > 0$$

since  $\phi$  and  $1 - \phi$  are nonnegative. Hence,  $\mathbf{H}_k > 0$ , from which we conclude that the algorithm has the descent property if  $\alpha_k$  is computed by line search (by Proposition 11.1).

### 11.16

To show the result, we will prove the following precise statement: In the quadratic case (with Hessian  $\mathbf{Q}$ ), suppose that  $\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \rho_i \Delta \mathbf{x}^{(i)}$ ,  $0 \leq i \leq k$ ,  $k \leq n - 1$ . If  $\alpha_i \neq 0$ ,  $0 \leq i \leq k$ , then  $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k+1)}$  are  $\mathbf{Q}$ -conjugate.

We proceed by induction. We begin with the  $k = 0$  case: that  $\mathbf{d}^{(0)}$  and  $\mathbf{d}^{(1)}$  are  $\mathbf{Q}$ -conjugate. Because  $\alpha_0 \neq 0$ , we can write  $\mathbf{d}^{(0)} = \Delta \mathbf{x}^{(0)} / \alpha_0$ . Hence,

$$\begin{aligned} \mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(0)} &= -\mathbf{g}^{(1)\top} \mathbf{H}_1 \mathbf{Q} \mathbf{d}^{(0)} \\ &= -\mathbf{g}^{(1)\top} \mathbf{H}_1 \frac{\mathbf{Q} \Delta \mathbf{x}^{(0)}}{\alpha_0} \\ &= -\mathbf{g}^{(1)\top} \frac{\mathbf{H}_1 \Delta \mathbf{g}^{(0)}}{\alpha_0} \\ &= -\mathbf{g}^{(1)\top} \frac{\rho_0 \Delta \mathbf{x}^{(0)}}{\alpha_0} \\ &= -\rho_0 \mathbf{g}^{(1)\top} \mathbf{d}^{(0)}. \end{aligned}$$

But  $\mathbf{g}^{(1)\top} \mathbf{d}^{(0)} = 0$  as a consequence of  $\alpha_0 > 0$  being the minimizer of  $\phi(\alpha) = f(\mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)})$ . Hence,  $\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(0)} = 0$ .

Assume that the result is true for  $k - 1$  (where  $k < n - 1$ ). We now prove the result for  $k$ , that is, that  $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k+1)}$  are  $\mathbf{Q}$ -conjugate. It suffices to show that  $\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(i)} = 0$ ,  $0 \leq i \leq k$ . Given  $i$ ,  $0 \leq i \leq k$ , using the same algebraic steps as in the  $k = 0$  case, and using the assumption that  $\alpha_i \neq 0$ , we obtain

$$\begin{aligned} \mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(i)} &= -\mathbf{g}^{(k+1)\top} \mathbf{H}_{k+1} \mathbf{Q} \mathbf{d}^{(i)} \\ &\vdots \\ &= -\rho_i \mathbf{g}^{(k+1)\top} \mathbf{d}^{(i)}. \end{aligned}$$

Because  $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k)}$  are  $\mathbf{Q}$ -conjugate by assumption, we conclude from the expanding subspace lemma (Lemma 10.2) that  $\mathbf{g}^{(k+1)\top} \mathbf{d}^{(i)} = 0$ . Hence,  $\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(i)} = 0$ , which completes the proof.

### 11.17

A MATLAB routine for the quasi-Newton algorithm with options for different formulas of  $\mathbf{H}_k$  is:

```
function [x,N]=quasi_newton(grad,xnew,H,options);
%   QUASI_NEWTON('grad',x0,H0);
%   QUASI_NEWTON('grad',x0,H0,OPTIONS);
%
%   x = QUASI_NEWTON('grad',x0,H0);
%   x = QUASI_NEWTON('grad',x0,H0,OPTIONS);
%
%   [x,N] = QUASI_NEWTON('grad',x0,H0);
%   [x,N] = QUASI_NEWTON('grad',x0,H0,OPTIONS);
%
%The first variant finds the minimizer of a function whose gradient
%is described in grad (usually an M-file: grad.m), using initial point
%x0 and initial inverse Hessian approximation H0.
%The second variant allows a vector of optional parameters to be
%defined:
%OPTIONS(1) controls how much display output is given; set
%to 1 for a tabular display of results, (default is no display: 0).
%OPTIONS(2) is a measure of the precision required for the final point.
%OPTIONS(3) is a measure of the precision required of the gradient.
%OPTIONS(5) specifies the formula for the inverse Hessian update:
%   0=Rank One;
%   1=DFP;
%   2=BFGS;
%OPTIONS(14) is the maximum number of iterations.
%For more information type HELP FOPTIONS.
%
%The next two variants return the value of the final point.
%The last two variants return a vector of the final point and the
%number of iterations.

if nargin ~= 4
    options = [];
    if nargin ~= 3
        disp('Wrong number of arguments.');
```

```
        return;
    end
end

numvars = length(xnew);
if length(options) >= 14
    if options(14)==0
        options(14)=1000*numvars;
    end
else
```

```

    options(14)=1000*numvars;
end

clc;
format compact;
format short e;

options = foptions(options);
print = options(1);
epsilon_x = options(2);
epsilon_g = options(3);
max_iter=options(14);

reset_cnt = 0;
g_curr=feval(grad,xnew);
if norm(g_curr) <= epsilon_g
    disp('Terminating: Norm of initial gradient less than');
    disp(epsilon_g);
    return;
end %if

d=-H*g_curr;
for k = 1:max_iter,

    xcurr=xnew;
    alpha=lineasearch_secant(grad,xcurr,d);
    xnew = xcurr+alpha*d;

    if print,
        disp('Iteration number k =')
        disp(k); %print iteration index k
        disp('alpha =');
        disp(alpha); %print alpha
        disp('Gradient = ');
        disp(g_curr'); %print gradient
        disp('New point =');
        disp(xnew'); %print new point
    end %if

    if norm(xnew-xcurr) <= epsilon_x*norm(xcurr)
        disp('Terminating: Norm of difference between iterates less than');
        disp(epsilon_x);
        break;
    end %if

    g_old=g_curr;
    g_curr=feval(grad,xnew);

    if norm(g_curr) <= epsilon_g
        disp('Terminating: Norm of gradient less than');
        disp(epsilon_g);
        break;
    end %if

    p=alpha*d;
    q=g_curr-g_old;

    reset_cnt = reset_cnt+1;
    if reset_cnt == 3*numvars

```

```

        d=-g_curr;
        reset_cnt = 0;
    else
        if options(5)==0 %Rank One
            q'*(p-H*q)
            H = H+(p-H*q)*(p-H*q)'/(q'*(p-H*q));
        elseif options(5)==1 %DFP
            H = H+p*p'/(p'*q)-(H*q)*(H*q)'/(q'*H*q);
        else %BFGS
            H = H+(1+q'*H*q/(q'*p))*p*p'/(p'*q)-(H*q*p'+(H*q*p'))/(q'*p);
        end %if
        d=-H*g_curr;
    end

    if print,
        disp('New H =');
        disp(H);
        disp('New d =');
        disp(d);
    end

    if k == max_iter
        disp('Terminating with maximum number of iterations');
    end %if
end %for

if nargout >= 1
    x=xnew;
    if nargout == 2
        N=k;
    end
else
    disp('Final point =');
    disp(xnew);
    disp('Number of iterations =');
    disp(k);
end %if
%-----

```

We created the following M-file, `g.m`, for the gradient of Rosenbrock's function:

```

function y=g(x)
y=[-400*(x(2)-x(1).^2)*x(1)-2*(1-x(1)), 200*(x(2)-x(1).^2)]';

```

We tested the above routine as follows:

```

>> options(2)=10^(-7);
>> options(3)=10^(-7);
>> options(14)=100;
>> x0=[-2;2];
>> H0=eye(2);
>> options(5)=0;
>> quasi_newton('g',x0,H0,options);

```

```

Terminating: Norm of difference between iterates less than
1.0000e-07
Final point =
1.0000e+00 1.0000e+00
Number of iterations =
8

```

```

>> options(5)=1;
>> quasi_newton('g',x0,H0,options);

Terminating: Norm of difference between iterates less than
1.0000e-07
Final point =
1.0000e+00 1.0000e+00
Number of iterations =
8
>> options(5)=2;
>> quasi_newton('g',x0,H0,options);

Terminating: Norm of difference between iterates less than
1.0000e-07
Final point =
1.0000e+00 1.0000e+00
Number of iterations =
8

```

The reader is again cautioned *not* to draw any conclusions about the superiority or inferiority of any of the formulas for  $\mathbf{H}_k$  based only on the above single numerical experiment.

#### 11.18

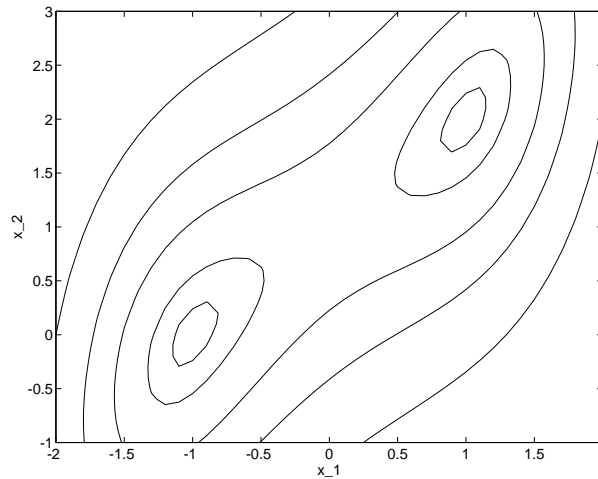
a. The plot of the level sets of  $f$  were obtained using the following MATLAB commands:

```

>> [X,Y]=meshdom(-2:0.1:2, -1:0.1:3);
>> Z=X.^4/4+Y.^2/2 -X.*Y +X-Y;
>> V=[-0.72, -0.6, -0.2, 0.5, 2];
>> contour(X,Y,Z,V)

```

The plot is depicted below:



b. With the initial condition  $[0, 0]^\top$ , the algorithm converges to  $[-1, 0]^\top$ , while with the initial condition  $[1.5, 1]^\top$ , the algorithm converges to  $[1, 2]^\top$ . These two points are the two strict local minimizers of  $f$  (as can be checked using the SOSC). The algorithm apparently converges to the minimizer “closer” to the initial point.

## 12. Solving $Ax = b$

### 12.1

---

Write the least squares cost in the usual notation  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$  where

$$\mathbf{A} = \begin{bmatrix} 3 \\ 5 \\ 6 \end{bmatrix}, \quad \mathbf{x} = [m], \quad \mathbf{b} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}.$$

The least squares estimate of the mass is

$$m^* = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} = \frac{31}{70}.$$

### 12.2

---

Write the least squares cost in the usual notation  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$  where

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 4 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} a \\ b \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix}.$$

The least squares estimate for  $[a, b]^\top$  is

$$\begin{aligned} \begin{bmatrix} a^* \\ b^* \end{bmatrix} &= (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} \\ &= \begin{bmatrix} 3 & 7 \\ 7 & 21 \end{bmatrix}^{-1} \begin{bmatrix} 12 \\ 31 \end{bmatrix} \\ &= \frac{1}{14} \begin{bmatrix} 21 & -7 \\ -7 & 3 \end{bmatrix}^{-1} \begin{bmatrix} 12 \\ 31 \end{bmatrix} \\ &= \frac{1}{14} \begin{bmatrix} 35 \\ 9 \end{bmatrix} \\ &= \begin{bmatrix} 5/2 \\ 9/14 \end{bmatrix}. \end{aligned}$$

### 12.3

---

a. We form

$$\mathbf{A} = \begin{bmatrix} 1^2/2 \\ 2^2/2 \\ 3^2/2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 5.00 \\ 19.5 \\ 44.0 \end{bmatrix}.$$

The least squares estimate of  $g$  is then given by

$$g = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} = 9.776.$$

b. We start with  $\mathbf{P}_0 = 0.040816$ , and  $\mathbf{x}^{(0)} = 9.776$ . We have  $\mathbf{a}_1 = 4^2/2 = 8$ , and  $\mathbf{b}^{(1)} = 78.5$ . Using the RLS formula, we get  $\mathbf{x}^{(1)} = 9.802$ , which is our updated estimate of  $g$ .

### 12.4

---

Let  $\mathbf{x} = [x_1, x_2, \dots, x_n]^\top$  and  $\mathbf{y} = [y_1, y_2, \dots, y_n]^\top$ . This least-squares estimation problem can be expressed as

$$\text{minimize } \|\alpha \mathbf{x} - \mathbf{y}\|^2,$$

with  $\alpha$  as the decision variable. Assuming that  $\mathbf{x} \neq \mathbf{0}$ , the solution is unique and is given by

$$\alpha^* = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{y} = \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}.$$

## 12.5

The least squares estimate of  $R$  is the least squares solution to

$$\begin{aligned} 1 \cdot R &= V_1 \\ &\vdots \\ 1 \cdot R &= V_n. \end{aligned}$$

Therefore, the least squares solution is

$$R^* = \left( \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \right)^{-1} \begin{bmatrix} V_1 \\ \vdots \\ V_n \end{bmatrix} = \frac{V_1 + \cdots + V_n}{n}.$$

## 12.6

We represent the data in the table and the decision variables  $a$  and  $b$  using the usual least squares matrix notation:

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 1 & 1 \\ 3 & 2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 6 \\ 4 \\ 5 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} a \\ b \end{bmatrix}.$$

The least squares estimate is given by

$$\mathbf{x}^* = \begin{bmatrix} a^* \\ b^* \end{bmatrix} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} = \begin{bmatrix} 11 & 9 \\ 9 & 9 \end{bmatrix}^{-1} \begin{bmatrix} 25 \\ 26 \end{bmatrix} = \frac{1}{18} \begin{bmatrix} 9 & -9 \\ -9 & 11 \end{bmatrix} \begin{bmatrix} 25 \\ 26 \end{bmatrix} = \begin{bmatrix} -1/2 \\ 61/18 \end{bmatrix}.$$

## 12.7

The problem can be formulated as a least-squares problem with

$$\mathbf{A} = \begin{bmatrix} 0.3 & 0.1 \\ 0.4 & 0.2 \\ 0.3 & 0.7 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 5 \\ 3 \\ 4 \end{bmatrix},$$

where the decision variable is  $\mathbf{x} = [x_1, x_2]$ , and  $x_1$  and  $x_2$  are the amounts of A and B, respectively. After some algebra, we obtain the solution:

$$\mathbf{x}^* = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} = \frac{1}{(0.34)(0.54) - (0.32)^2} \begin{bmatrix} 0.54 & -0.32 \\ -0.32 & 0.34 \end{bmatrix} \begin{bmatrix} 3.9 \\ 3.9 \end{bmatrix}.$$

Since we are only interest in the ratio of the first component of  $\mathbf{x}^*$  to the second component, we need only explicitly compute:

$$\text{Ratio} = \frac{0.54 - 0.32}{-0.32 + 0.34} = \frac{0.22}{0.02} = 11.$$

## 12.8

For each  $k$ , we can write

$$\begin{aligned} y_k &= ay_{k-1} + bu_k + v_k \\ &= a^2 y_{k-2} + abu_{k-1} + av_{k-1} + bu_k + v_k \\ &\vdots \\ &= a^{k-1} bu_1 + a^{k-2} bu_2 + \cdots + bu_k + a^{k-1} v_1 + a^{k-2} v_2 + \cdots + v_k \end{aligned}$$



Write  $\mathbf{u} = [u_1, \dots, u_n]^\top$ ,  $\mathbf{v} = [v_1, \dots, v_n]^\top$ , and  $\mathbf{y} = [y_1, \dots, y_n]^\top$ . Then,  $\mathbf{y} = \mathbf{C}\mathbf{u} + \mathbf{D}\mathbf{v}$ , where

$$\mathbf{C} = \begin{bmatrix} b & 0 & \cdots & 0 \\ ab & b & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ a^{n-1}b & \cdots & ab & b \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ a & 1 & \ddots & \vdots \\ \vdots & & \ddots & 0 \\ a^{n-1} & a^{n-2} & \cdots & 1 \end{bmatrix}.$$

Write  $\mathbf{b} = \mathbf{D}^{-1}\mathbf{y}$  and  $\mathbf{A} = \mathbf{D}^{-1}\mathbf{C}$  so that  $\mathbf{b} = \mathbf{A}\mathbf{u} + \mathbf{v}$ . Therefore, the linear least-squares estimate of  $\mathbf{u}$  given  $\mathbf{y}$  is

$$\mathbf{u}^* = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} = (\mathbf{C}^\top \mathbf{D}^{-\top} \mathbf{D}^{-1} \mathbf{C})^{-1} \mathbf{C}^\top \mathbf{D}^{-\top} \mathbf{D}^{-1} \mathbf{y}.$$

But  $\mathbf{C} = b\mathbf{D}$ . Hence,

$$\mathbf{u}^* = \frac{1}{b} \mathbf{D}^{-1} \mathbf{y} = \frac{1}{b} \begin{bmatrix} 1 & 0 & \cdots & 0 \\ -a & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & -a & 1 \end{bmatrix} \mathbf{y}.$$

(Notice that  $\mathbf{D}^{-1}$  has the simple form shown above.)

An alternative solution is first to define  $\mathbf{z} = [z_1, \dots, z_n]^\top$  by  $z_k = y_k - ay_{k-1}$ . Then, we have  $\mathbf{z} = b\mathbf{u} + \mathbf{v}$ . Therefore, the linear least-squares estimate of  $\mathbf{u}$  given  $\mathbf{y}$  (or, equivalently,  $\mathbf{z}$ ) is

$$\mathbf{u}^* = \frac{1}{b} \mathbf{z} = \frac{1}{b} \begin{bmatrix} 1 & 0 & \cdots & 0 \\ -a & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & -a & 1 \end{bmatrix} \mathbf{y}.$$

## 12.9

Define

$$\mathbf{X} = \begin{bmatrix} x_1 & 1 \\ \vdots & \vdots \\ x_p & 1 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_p \end{bmatrix}.$$

Since the  $x_i$  are not all equal, we have  $\text{rank } \mathbf{X} = 2$ . The objective function can be written as

$$f(a, b) = \left\| \mathbf{X} \begin{bmatrix} a \\ b \end{bmatrix} - \mathbf{y} \right\|^2.$$

Therefore, by Theorem 12.1 there exists a unique minimizer  $[a^*, b^*]^\top$  given by

$$\begin{aligned} \begin{bmatrix} a^* \\ b^* \end{bmatrix} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= \begin{bmatrix} \sum_{i=1}^p x_i^2 & \sum_{i=1}^p x_i \\ \sum_{i=1}^p x_i & p \end{bmatrix}^{-1} \begin{bmatrix} \sum_{i=1}^p x_i y_i \\ \sum_{i=1}^p y_i \end{bmatrix} \\ &= \begin{bmatrix} \overline{X^2} & \overline{X} \\ \overline{X} & 1 \end{bmatrix}^{-1} \begin{bmatrix} \overline{XY} \\ \overline{Y} \end{bmatrix} \\ &= \frac{1}{\overline{X^2} - (\overline{X})^2} \begin{bmatrix} 1 & -\overline{X} \\ -\overline{X} & \overline{X^2} \end{bmatrix} \begin{bmatrix} \overline{XY} \\ \overline{Y} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\overline{XY} - (\overline{X})(\overline{Y})}{\overline{X^2} - (\overline{X})^2} \\ \frac{(\overline{X^2})(\overline{Y}) - (\overline{X})(\overline{XY})}{\overline{X^2} - (\overline{X})^2} \end{bmatrix}. \end{aligned}$$

As we can see, the solution does not depend on  $\overline{Y^2}$ .

### 12.10

a. We wish to find  $\omega$  and  $\theta$  such that

$$\begin{aligned}\sin(\omega t_1 + \theta) &= y_1 \\ &\vdots \\ \sin(\omega t_p + \theta) &= y_p.\end{aligned}$$

Taking arcsin, we get the following system of linear equations:

$$\begin{aligned}\omega t_1 + \theta &= \arcsin y_1 \\ &\vdots \\ \omega t_p + \theta &= \arcsin y_p.\end{aligned}$$

b. We may write the system of linear equations in part a as  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , where

$$\mathbf{A} = \begin{bmatrix} t_1 & 1 \\ \vdots & \vdots \\ t_p & 1 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} \omega \\ \theta \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \arcsin y_1 \\ \vdots \\ \arcsin y_p \end{bmatrix}.$$

Since the  $t_i$  are not all equal, the first column of  $\mathbf{A}$  is not a scalar multiple of the second column. Therefore,  $\text{rank } \mathbf{A} = 2$ . Hence, the least squares solution is

$$\begin{aligned}\mathbf{x} &= (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} \\ &= \begin{bmatrix} \sum_{i=1}^p t_i^2 & \sum_{i=1}^p t_i \\ \sum_{i=1}^p t_i & p \end{bmatrix}^{-1} \begin{bmatrix} \sum_{i=1}^p t_i \arcsin y_i \\ \sum_{i=1}^p \arcsin y_i \end{bmatrix} \\ &= \begin{bmatrix} \overline{T^2} & \overline{T} \\ \overline{T} & 1 \end{bmatrix}^{-1} \begin{bmatrix} \overline{TY} \\ \overline{Y} \end{bmatrix} \\ &= \frac{1}{\overline{T^2} - (\overline{T})^2} \begin{bmatrix} 1 & -\overline{T} \\ -\overline{T} & \overline{T^2} \end{bmatrix} \begin{bmatrix} \overline{TY} \\ \overline{Y} \end{bmatrix} \\ &= \frac{1}{\overline{T^2} - (\overline{T})^2} \begin{bmatrix} \overline{TY} - (\overline{T})(\overline{Y}) \\ -(\overline{T})(\overline{TY}) + (\overline{T^2})(\overline{Y}) \end{bmatrix}.\end{aligned}$$

### 12.11

The given line can be expressed as the range of the matrix  $\mathbf{A} = [1, m]^\top$ . Let  $\mathbf{b} = [x_0, y_0]^\top$  be the given point. Therefore, the problem is a linear least squares problem of minimizing  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$ . The solution is given by

$$\mathbf{x}^* = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} = \frac{x_0 + my_0}{1 + m^2}.$$

Therefore, the point on the straight line that is closest to the given point  $[x_0, y_0]$  is given by  $[x^*, mx^*]^\top$ .

### 12.12

a. Write

$$\mathbf{A} = \begin{bmatrix} \mathbf{x}_1^\top & 1 \\ \vdots & \vdots \\ \mathbf{x}_p^\top & 1 \end{bmatrix} \in \mathbb{R}^{p \times (n+1)}, \quad \mathbf{z} = \begin{bmatrix} \mathbf{a} \\ c \end{bmatrix} \in \mathbb{R}^{n+1}, \quad \mathbf{b} = \begin{bmatrix} y_1 \\ \vdots \\ y_p \end{bmatrix} \in \mathbb{R}^p.$$

The objective function can then be written as  $\|\mathbf{A}\mathbf{z} - \mathbf{b}\|^2$ .

b. Let  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_p]^\top \in \mathbb{R}^{p \times n}$ , and  $\mathbf{e} = [1, \dots, 1]^\top \in \mathbb{R}^p$ . Then we may write  $\mathbf{A} = [\mathbf{X} \ \mathbf{e}]$ . The solution to the problem is  $(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$ . But

$$\mathbf{A}^\top \mathbf{A} = \begin{bmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{e} \\ \mathbf{e}^\top \mathbf{X} & p \end{bmatrix} = \begin{bmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{0} \\ \mathbf{0}^\top & p \end{bmatrix}$$

since  $\mathbf{X}^\top \mathbf{e} = \mathbf{x}_1 + \dots + \mathbf{x}_p = \mathbf{0}$  by assumption. Also,

$$\mathbf{A}^\top \mathbf{y} = \begin{bmatrix} \mathbf{X}^\top \mathbf{y} \\ \mathbf{e}^\top \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{e}^\top \mathbf{y} \end{bmatrix}$$

since  $\mathbf{X}^\top \mathbf{y} = y_1 \mathbf{x}_1 + \dots + y_p \mathbf{x}_p = \mathbf{0}$  by assumption. Therefore, the solution is given by

$$\mathbf{z}^* = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} = \begin{bmatrix} (\mathbf{X}^\top \mathbf{X})^{-1} & \mathbf{0} \\ \mathbf{0}^\top & 1/p \end{bmatrix} \begin{bmatrix} \mathbf{0} \\ \mathbf{e}^\top \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \frac{1}{p} \mathbf{e}^\top \mathbf{y} \end{bmatrix}.$$

The affine function of best fit is the constant function  $f(\mathbf{x}) = c$ , where

$$c = \frac{1}{p} \sum_{i=1}^p y_i.$$

### 12.13

a. Using the least squares formula, we have

$$\hat{\theta}_n = \left( [u_1, \dots, u_n] \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix} \right)^{-1} [u_1, \dots, u_n] \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \frac{\sum_{k=1}^n u_k y_k}{\sum_{k=1}^n u_k^2}.$$

b. Given  $u_k = 1$  for all  $k$ , we have

$$\hat{\theta}_n = \frac{1}{n} \sum_{k=1}^n y_k = \frac{1}{n} \sum_{k=1}^n (\theta + e_k) = \theta + \frac{1}{n} \sum_{k=1}^n e_k.$$

Hence,  $\hat{\theta}_n \rightarrow \theta$  if and only if  $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n e_k = 0$ .

### 12.14

We pose the problem as a least squares problem: minimize  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$  where  $\mathbf{x} = [a, b]^\top$ , and

$$\mathbf{A} = \begin{bmatrix} x_0 & 1 \\ x_1 & 1 \\ x_2 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}.$$

We have

$$\mathbf{A}^\top \mathbf{A} = \begin{bmatrix} \sum_{i=0}^2 x_i^2 & \sum_{i=0}^2 x_i \\ \sum_{i=0}^2 x_i & 3 \end{bmatrix}, \quad \mathbf{A}^\top \mathbf{b} = \begin{bmatrix} \sum_{i=0}^2 x_i x_{i+1} \\ \sum_{i=0}^2 x_{i+1} \end{bmatrix}.$$

Therefore, the least squares solution is

$$\begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \sum_{i=0}^2 x_i^2 & \sum_{i=0}^2 x_i \\ \sum_{i=0}^2 x_i & 3 \end{bmatrix}^{-1} \begin{bmatrix} \sum_{i=0}^2 x_i x_{i+1} \\ \sum_{i=0}^2 x_{i+1} \end{bmatrix} = \begin{bmatrix} 5 & 3 \\ 3 & 3 \end{bmatrix}^{-1} \begin{bmatrix} 18 \\ 11 \end{bmatrix} = \begin{bmatrix} 7/2 \\ 1/6 \end{bmatrix}.$$

**12.15**

We pose the problem as a least squares problem: minimize  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$  where  $\mathbf{x} = [a, b]^\top$ , and

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ h_1 & 0 \\ \vdots & \vdots \\ h_{n-1} & 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_n \end{bmatrix}$$

(note that  $h_0 = 0$ ). We have

$$\mathbf{A}^\top \mathbf{A} = \begin{bmatrix} \sum_{i=1}^{n-1} h_i^2 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{A}^\top \mathbf{b} = \begin{bmatrix} \sum_{i=1}^{n-1} h_i h_{i+1} \\ h_1 \end{bmatrix}.$$

The matrix  $\mathbf{A}^\top \mathbf{A}$  is nonsingular because we assume that at least one  $h_k$  is nonzero. Therefore, the least squares solution is

$$\begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^{n-1} h_i^2 & 0 \\ 0 & 1 \end{bmatrix}^{-1} \begin{bmatrix} \sum_{i=1}^{n-1} h_i h_{i+1} \\ h_1 \end{bmatrix} = \begin{bmatrix} (\sum_{i=1}^{n-1} h_i h_{i+1}) / (\sum_{i=1}^{n-1} h_i^2) \\ h_1 \end{bmatrix}.$$

**12.16**

We pose the problem as a least squares problem: minimize  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$  where  $\mathbf{x} = [a, b]^\top$ , and

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ s_1 & 1 \\ \vdots & \vdots \\ s_{n-1} & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{bmatrix}$$

(where we use  $s_0 = 0$ ). We have

$$\mathbf{A}^\top \mathbf{A} = \begin{bmatrix} \sum_{i=1}^{n-1} s_i^2 & \sum_{i=1}^{n-1} s_i \\ \sum_{i=1}^{n-1} s_i & n \end{bmatrix}, \quad \mathbf{A}^\top \mathbf{b} = \begin{bmatrix} \sum_{i=1}^{n-1} s_i s_{i+1} \\ \sum_{i=1}^n s_i \end{bmatrix}.$$

The matrix  $\mathbf{A}^\top \mathbf{A}$  is nonsingular because we assume that at least one  $s_k$  is nonzero. Therefore, the least squares solution is

$$\begin{aligned} \begin{bmatrix} a \\ b \end{bmatrix} &= \begin{bmatrix} \sum_{i=1}^{n-1} s_i^2 & \sum_{i=1}^{n-1} s_i \\ \sum_{i=1}^{n-1} s_i & n \end{bmatrix}^{-1} \begin{bmatrix} \sum_{i=1}^{n-1} s_i s_{i+1} \\ \sum_{i=1}^n s_i \end{bmatrix} \\ &= \frac{1}{n \sum_{i=1}^{n-1} s_i^2 - \left( \sum_{i=1}^{n-1} s_i \right)^2} \begin{bmatrix} n \sum_{i=1}^{n-1} s_i s_{i+1} - \sum_{i=1}^{n-1} s_i \sum_{i=1}^n s_i \\ - \sum_{i=1}^{n-1} s_i \sum_{i=1}^{n-1} s_i s_{i+1} + \sum_{i=1}^{n-1} s_i^2 \sum_{i=1}^n s_i \end{bmatrix}. \end{aligned}$$

**12.17**

This least-squares estimation problem can be expressed as

$$\text{minimize } \|\mathbf{a}\mathbf{x} - \mathbf{y}\|^2.$$

If  $\mathbf{x} = \mathbf{0}$ , then the problem has an infinite number of solutions: any  $a$  solves the problem. Assuming that  $\mathbf{x} \neq \mathbf{0}$ , the solution is unique and is given by

$$a^* = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{y} = \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}}.$$

**12.18**

The solution to this problem is the same as the solution to:

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|\mathbf{x} - \mathbf{b}\|^2 \\ & \text{subject to} && \mathbf{x} \in \mathcal{R}(\mathbf{A}). \end{aligned}$$

Substituting  $\mathbf{x} = \mathbf{A}\mathbf{y}$ , we see that this is simply a linear least squares problem with decision variable  $\mathbf{y}$ . The solution to the least squares problem is  $\mathbf{y}^* = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$ , which implies that the solution to the given problem is  $\mathbf{x}^* = \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$ .

**12.19**

We solve the problem using two different methods. The first method would be to use the Lagrange multiplier technique to solve the equivalent problem,

$$\begin{aligned} & \text{minimize} && \|\mathbf{x} - \mathbf{x}_0\|^2 \\ & \text{subject to} && h(\mathbf{x}) = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \mathbf{x} - 1 = 0, \end{aligned}$$

The lagrangian for the above problem has the form,

$$l(\mathbf{x}, \lambda) = x_1^2 + (x_2 + 3)^2 + x_3^2 + \lambda(x_1 + x_2 + x_3 - 1).$$

Applying the FONC gives

$$\nabla_{\mathbf{x}} l = \begin{bmatrix} 2x_1 + \lambda \\ 2x_2 + 6 + \lambda \\ 2x_3 + \lambda \end{bmatrix} \quad \text{and} \quad x_1 + x_2 + x_3 - 1 = 0.$$

Solving the above yields

$$\mathbf{x}^* = \begin{bmatrix} \frac{4}{3} \\ -\frac{5}{3} \\ \frac{4}{3} \end{bmatrix}.$$

The second approach is to use the well-known solution to the minimum norm problem. We first derive a general solution formula for the problem,

$$\begin{aligned} & \text{minimize} && \|\mathbf{x} - \mathbf{x}_0\| \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b}, \end{aligned}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m \leq n$ , and  $\text{rank } \mathbf{A} = m$ . To proceed, we first transform the above problem from the  $\mathbf{x}$  coordinates into the  $\mathbf{z} = \mathbf{x} - \mathbf{x}_0$  coordinates to obtain,

$$\begin{aligned} & \text{minimize} && \|\mathbf{z}\| \\ & \text{subject to} && \mathbf{A}\mathbf{z} = \mathbf{b} - \mathbf{A}\mathbf{x}_0. \end{aligned}$$

The solution to the above problem has the form,

$$\begin{aligned} \mathbf{z}^* &= \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} (\mathbf{b} - \mathbf{A}\mathbf{x}_0) \\ &= \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{b} - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{x}_0. \end{aligned}$$

Therefore, the solution to the original problem is

$$\begin{aligned} \mathbf{x}^* &= \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} (\mathbf{b} - \mathbf{A}\mathbf{x}_0) + \mathbf{x}_0 \\ &= \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{b} - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{x}_0 + \mathbf{x}_0 \\ &= \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{b} + \left( \mathbf{I}_n - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A} \right) \mathbf{x}_0. \end{aligned}$$

We substitute into the above formula the given numerical data to obtain

$$\begin{aligned}\mathbf{x}^* &= \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix} + \begin{bmatrix} \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} \\ -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} \end{bmatrix} \begin{bmatrix} 0 \\ -3 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} \frac{4}{3} \\ -\frac{5}{3} \\ \frac{4}{3} \end{bmatrix}.\end{aligned}$$

## 12.20

For each  $\mathbf{x} \in \mathbb{R}^n$ , let  $\mathbf{y} = \mathbf{x} - \mathbf{x}_0$ . Then, the original problem is equivalent to

$$\begin{aligned}&\text{minimize} && \|\mathbf{y}\| \\ &\text{subject to} && \mathbf{A}\mathbf{y} = \mathbf{b} - \mathbf{A}\mathbf{x}_0,\end{aligned}$$

in the sense that  $\mathbf{y}^*$  is a solution to the above problem if and only if  $\mathbf{x}^* = \mathbf{y}^* + \mathbf{x}_0$  is a solution to the original problem. By Theorem 12.2, the above problem has a unique solution given by

$$\mathbf{y}^* = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} (\mathbf{b} - \mathbf{A}\mathbf{x}_0) = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{b} - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{x}_0.$$

Therefore, the solution to the original problem is

$$\mathbf{x}^* = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{b} - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{x}_0 + \mathbf{x}_0 = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{b} + (\mathbf{I}_n - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}) \mathbf{x}_0.$$

Note that

$$\begin{aligned}\|\mathbf{x}^* - \mathbf{x}_0\| &= \|\mathbf{y}^*\| \\ &= \|\mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} (\mathbf{b} - \mathbf{A}\mathbf{x}_0)\| \\ &= \|\mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{b} - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{x}_0\|.\end{aligned}$$

## 12.21

The objective function of the given problem can be written as

$$f(\mathbf{x}) = \|\mathbf{B}\mathbf{x} - \mathbf{c}\|^2,$$

where

$$\mathbf{B} = \begin{bmatrix} \mathbf{A} \\ \vdots \\ \mathbf{A} \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_p \end{bmatrix}.$$

The solution is therefore

$$\mathbf{x}^* = (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{c} = (p\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top (\mathbf{b}_1 + \cdots + \mathbf{b}_p) = \frac{1}{p} \sum_{i=1}^p (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}_i = \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^*$$

Alternatively: Write

$$\|\mathbf{A}\mathbf{x} - \mathbf{b}_i\|^2 = \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} - 2\mathbf{x}^\top \mathbf{A}^\top \mathbf{b}_i + \|\mathbf{b}_i\|^2$$

Therefore, the given objective function can be written as

$$p\mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} - 2\mathbf{x}^\top \mathbf{A}^\top (\mathbf{b}_1 + \cdots + \mathbf{b}_p) + \|\mathbf{b}_1\|^2 + \cdots + \|\mathbf{b}_p\|^2.$$

The solution is therefore

$$\mathbf{x}^* = (p\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top (\mathbf{b}_1 + \cdots + \mathbf{b}_p) = \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^*$$

Note that the original problem can be written as the least squares problem

$$\text{minimize } \|\mathbf{Ax} - \mathbf{b}\|^2,$$

where

$$\mathbf{b} = \frac{\mathbf{b}_1 + \cdots + \mathbf{b}_p}{p}.$$

---

### 12.22

Write

$$\|\mathbf{Ax} - \mathbf{b}_i\|^2 = \mathbf{x}^\top \mathbf{A}^\top \mathbf{Ax} - 2\mathbf{x}^\top \mathbf{A}^\top \mathbf{b}_i + \|\mathbf{b}_i\|^2$$

Therefore, the given objective function can be written as

$$(\alpha_1 + \cdots + \alpha_p) \mathbf{x}^\top \mathbf{A}^\top \mathbf{Ax} - 2\mathbf{x}^\top \mathbf{A}^\top (\alpha_1 \mathbf{b}_1 + \cdots + \alpha_p \mathbf{b}_p) + \alpha_1 \|\mathbf{b}_1\|^2 + \cdots + \alpha_p \|\mathbf{b}_p\|^2.$$

The solution is therefore (by inspection)

$$\mathbf{x}^* = ((\alpha_1 + \cdots + \alpha_p) \mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top (\alpha_1 \mathbf{b}_1 + \cdots + \alpha_p \mathbf{b}_p) = \frac{1}{\alpha_1 + \cdots + \alpha_p} \sum_{i=1}^p \alpha_i \mathbf{x}_i^* = \sum_{i=1}^p \beta_i \mathbf{x}_i^*,$$

where  $\beta_i = \alpha_i / (\alpha_1 + \cdots + \alpha_p)$ .

Note that the original problem can be written as the least squares problem

$$\text{minimize } \|\mathbf{Ax} - \mathbf{b}\|^2,$$

where

$$\mathbf{b} = \frac{\alpha_1 \mathbf{b}_1 + \cdots + \alpha_p \mathbf{b}_p}{\alpha_1 + \cdots + \alpha_p}.$$

---

### 12.23

Let  $\mathbf{x}^* = \mathbf{A}^\top (\mathbf{AA}^\top)^{-1} \mathbf{b}$ . Suppose  $\mathbf{y}^*$  is a point in  $\mathcal{R}(\mathbf{A}^\top)$  that satisfies  $\mathbf{Ay}^* = \mathbf{b}$ . Then, there exists  $\mathbf{z}^* \in \mathbb{R}^m$  such that  $\mathbf{y}^* = \mathbf{A}^\top \mathbf{z}^*$ . Then, subtracting the equation  $\mathbf{A}(\mathbf{A}^\top (\mathbf{AA}^\top)^{-1} \mathbf{b}) = \mathbf{b}$  from the equation  $\mathbf{A}(\mathbf{A}^\top \mathbf{z}^*) = \mathbf{b}$ , we get

$$(\mathbf{AA}^\top)(\mathbf{z}^* - (\mathbf{AA}^\top)^{-1} \mathbf{b}) = \mathbf{0}.$$

Since  $\text{rank } \mathbf{A} = m$ ,  $\mathbf{AA}^\top$  is nonsingular. Therefore,  $\mathbf{z}^* - (\mathbf{AA}^\top)^{-1} \mathbf{b} = \mathbf{0}$ , which implies that

$$\mathbf{y}^* = \mathbf{A}^\top \mathbf{z}^* = \mathbf{A}^\top (\mathbf{AA}^\top)^{-1} \mathbf{b} = \mathbf{x}^*.$$

Hence,  $\mathbf{x}^* = \mathbf{A}^\top (\mathbf{AA}^\top)^{-1} \mathbf{b}$  is the only vector in  $\mathcal{R}(\mathbf{A}^\top)$  that satisfies  $\mathbf{Ax}^* = \mathbf{b}$ .

---

### 12.24

a. We have

$$\mathbf{x}^{(0)} = (\mathbf{A}_0^\top \mathbf{A}_0)^{-1} \mathbf{A}_0^\top \mathbf{b}^{(0)} = \mathbf{G}_0^{-1} \mathbf{A}_0^\top \mathbf{b}^{(0)}.$$

Similarly,

$$\mathbf{x}^{(1)} = (\mathbf{A}_1^\top \mathbf{A}_1)^{-1} \mathbf{A}_1^\top \mathbf{b}^{(1)} = \mathbf{G}_1^{-1} \mathbf{A}_1^\top \mathbf{b}^{(1)}.$$

b. Now,

$$\begin{aligned} \mathbf{G}_0 &= \begin{bmatrix} \mathbf{A}_1^\top & \mathbf{a}_1 \end{bmatrix} \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{a}_1^\top \end{bmatrix} \\ &= \mathbf{A}_1^\top \mathbf{A}_1 + \mathbf{a}_1 \mathbf{a}_1^\top \\ &= \mathbf{G}_1 + \mathbf{a}_1 \mathbf{a}_1^\top. \end{aligned}$$

Hence,

$$\mathbf{G}_1 = \mathbf{G}_0 - \mathbf{a}_1 \mathbf{a}_1^\top.$$

c. Using the Sherman-Morrison formula,

$$\begin{aligned}
\mathbf{P}_1 &= \mathbf{G}_1^{-1} \\
&= (\mathbf{G}_0 - \mathbf{a}_1 \mathbf{a}_1^\top)^{-1} \\
&= \mathbf{G}_0^{-1} - \frac{\mathbf{G}_0^{-1}(-\mathbf{a}_1) \mathbf{a}_1^\top \mathbf{G}_0^{-1}}{1 + (-\mathbf{a}_1)^\top \mathbf{G}_0^{-1} \mathbf{a}_1} \\
&= \mathbf{P}_0 + \frac{\mathbf{P}_0 \mathbf{a}_1 \mathbf{a}_1^\top \mathbf{P}_0}{1 - \mathbf{a}_1^\top \mathbf{P}_0 \mathbf{a}_1}.
\end{aligned}$$

d. We have

$$\begin{aligned}
\mathbf{A}_0^\top \mathbf{b}^{(0)} &= \mathbf{G}_0 \mathbf{G}_0^{-1} \mathbf{A}_0^\top \mathbf{b}^{(0)} \\
&= \mathbf{G}_0 \mathbf{x}^{(0)} \\
&= (\mathbf{G}_1 + \mathbf{a}_1 \mathbf{a}_1^\top) \mathbf{x}^{(0)} \\
&= \mathbf{G}_1 \mathbf{x}^{(0)} + \mathbf{a}_1 \mathbf{a}_1^\top \mathbf{x}^{(0)}.
\end{aligned}$$

e. Finally,

$$\begin{aligned}
\mathbf{x}^{(1)} &= \mathbf{G}_1^{-1} \mathbf{A}_1^\top \mathbf{b}^{(1)} \\
&= \mathbf{G}_1^{-1} (\mathbf{A}_1^\top \mathbf{b}^{(1)} + \mathbf{a}_1 b_1 - \mathbf{a}_1 b_1) \\
&= \mathbf{G}_1^{-1} (\mathbf{A}_0^\top \mathbf{b}^{(0)} - \mathbf{a}_1 b_1) \\
&= \mathbf{G}_1^{-1} (\mathbf{G}_1 \mathbf{x}^{(0)} + \mathbf{a}_1 \mathbf{a}_1^\top \mathbf{x}^{(0)} - \mathbf{a}_1 b_1) \\
&= \mathbf{x}^{(0)} - \mathbf{G}_1^{-1} \mathbf{a}_1 (b_1 - \mathbf{a}_1^\top \mathbf{x}^{(0)}) \\
&= \mathbf{x}^{(0)} - \mathbf{P}_1 \mathbf{a}_1 (b_1 - \mathbf{a}_1^\top \mathbf{x}^{(0)}).
\end{aligned}$$

The general RLS algorithm for removals of rows is:

$$\begin{aligned}
\mathbf{P}^{(k+1)} &= \mathbf{P}_k + \frac{\mathbf{P}_k \mathbf{a}_{k+1} \mathbf{a}_{k+1}^\top \mathbf{P}_k}{1 - \mathbf{a}_{k+1}^\top \mathbf{P}_k \mathbf{a}_{k+1}} \\
\mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \mathbf{P}_{k+1} \mathbf{a}_{k+1} (b_{k+1} - \mathbf{a}_{k+1}^\top \mathbf{x}^{(k)}).
\end{aligned}$$

## 12.25

Using the notation of the proof of Theorem 12.3, we can write

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mu (b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k)}) \frac{\mathbf{a}_{R(k)+1}}{\|\mathbf{a}_{R(k)+1}\|^2}.$$

Hence,

$$\mathbf{x}^{(k)} = \sum_{i=0}^{k-1} \left( \frac{\mu(2-\mu)}{\|\mathbf{a}_{R(i)+1}\|^2} (b_{R(i)+1} - \mathbf{a}_{R(i)+1}^\top \mathbf{x}^{(i)}) \right) \mathbf{a}_{R(i)+1}$$

which means that  $\mathbf{x}^{(k)}$  is in  $\text{span}[\mathbf{a}_1, \dots, \mathbf{a}_m] = \mathcal{R}(\mathbf{A}^\top)$ .

## 12.26

a. We claim that  $\mathbf{x}^*$  minimizes  $\|\mathbf{x} - \mathbf{x}^{(0)}\|$  subject to  $\{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}\}$  if and only if  $\mathbf{y}^* = \mathbf{x}^* - \mathbf{x}^{(0)}$  minimizes  $\|\mathbf{y}\|$  subject to  $\{\mathbf{A}\mathbf{y} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}\}$ .

To prove sufficiency, suppose  $\mathbf{y}^*$  minimizes  $\|\mathbf{y}\|$  subject to  $\{\mathbf{A}\mathbf{y} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}\}$ . Let  $\mathbf{x}^* = \mathbf{y}^* + \mathbf{x}^{(0)}$ . Consider any point  $\mathbf{x}_1 \in \{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}\}$ . Now,

$$\mathbf{A}(\mathbf{x}_1 - \mathbf{x}^{(0)}) = \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}.$$



Hence, by definition of  $\mathbf{y}^*$ ,

$$\|\mathbf{x}_1 - \mathbf{x}^{(0)}\| \geq \|\mathbf{y}^*\| = \|\mathbf{x}^* - \mathbf{x}^{(0)}\|.$$

Therefore  $\mathbf{x}^*$  minimizes  $\|\mathbf{x} - \mathbf{x}^{(0)}\|$  subject to  $\{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\}$ .

To prove necessity, suppose  $\mathbf{x}^*$  minimizes  $\|\mathbf{x} - \mathbf{x}^{(0)}\|$  subject to  $\{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\}$ . Let  $\mathbf{y}^* = \mathbf{x}^* - \mathbf{x}^{(0)}$ . Consider any point  $\mathbf{y}_1 \in \{\mathbf{y} : \mathbf{Ay} = \mathbf{b} - \mathbf{Ax}^{(0)}\}$ . Now,

$$\mathbf{A}(\mathbf{y}_1 + \mathbf{x}^{(0)}) = \mathbf{b}.$$

Hence, by definition of  $\mathbf{x}^*$ ,

$$\|\mathbf{y}_1\| = \|(\mathbf{y}_1 + \mathbf{x}^{(0)}) - \mathbf{x}^{(0)}\| \geq \|\mathbf{x}^* - \mathbf{x}^{(0)}\| = \|\mathbf{y}^*\|.$$

Therefore,  $\mathbf{y}^*$  minimizes  $\|\mathbf{y}\|$  subject to  $\{\mathbf{Ay} = \mathbf{b} - \mathbf{Ax}^{(0)}\}$ .

By Theorem 12.2, there exists a unique vector  $\mathbf{y}^*$  minimizing  $\|\mathbf{y}\|$  subject to  $\{\mathbf{Ay} = \mathbf{b} - \mathbf{Ax}^{(0)}\}$ . Hence, by the above claim, there exists a unique  $\mathbf{x}^*$  minimizing  $\|\mathbf{x} - \mathbf{x}^{(0)}\|$  subject to  $\{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\}$ .

b. Using the notation of the proof of Theorem 12.3, Kaczmarz's algorithm is given by

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mu(b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k)}) \mathbf{a}_{R(k)+1}.$$

Subtract  $\mathbf{x}^{(0)}$  from each side to give

$$(\mathbf{x}^{(k+1)} - \mathbf{x}^{(0)}) = (\mathbf{x}^{(k)} - \mathbf{x}^{(0)}) + \mu((b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(0)}) - \mathbf{a}_{R(k)+1}^\top (\mathbf{x}^{(k)} - \mathbf{x}^{(0)})) \mathbf{a}_{R(k)+1}.$$

Writing  $\mathbf{y}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^{(0)}$ , we get

$$\mathbf{y}^{(k+1)} = \mathbf{y}^{(k)} + \mu((b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(0)}) - \mathbf{a}_{R(k)+1}^\top \mathbf{y}^{(k)}) \mathbf{a}_{R(k)+1}.$$

Note that  $\mathbf{y}^{(0)} = \mathbf{0}$ . By Theorem 12.3, the sequence  $\{\mathbf{y}^{(k)}\}$  converges to the unique point  $\mathbf{y}^*$  that minimizes  $\|\mathbf{y}\|$  subject to  $\{\mathbf{Ay} = \mathbf{b} - \mathbf{Ax}^{(0)}\}$ . Hence  $\{\mathbf{x}^{(k)}\}$  converges to  $\mathbf{y}^* + \mathbf{x}^{(0)}$ . From the proof of part a,  $\mathbf{x}^* = \mathbf{y}^* + \mathbf{x}^{(0)}$  minimizes  $\|\mathbf{x} - \mathbf{x}^{(0)}\|$  subject to  $\{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\}$ . This completes the proof.

### 12.27

Following the proof of Theorem 12.3, assuming  $\|\mathbf{a}\| = 1$  without loss of generality, we arrive at

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2 = \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 - \mu(2 - \mu)(\mathbf{a}^\top (\mathbf{x}^{(k)} - \mathbf{x}^*))^2.$$

Since  $\mathbf{x}^{(k)}, \mathbf{x}^* \in \mathcal{R}(\mathbf{A}) = \mathcal{R}([\mathbf{a}^\top])$  by Exercise 12.25, we have  $\mathbf{x}^{(k)} - \mathbf{x}^* \in \mathcal{R}(\mathbf{A})$ . Hence, by the Cauchy-Schwarz inequality,

$$(\mathbf{a}^\top (\mathbf{x}^{(k)} - \mathbf{x}^*))^2 = \|\mathbf{a}\|^2 \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 = \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2,$$

since  $\|\mathbf{a}\| = 1$  by assumption. Thus, we obtain

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2 = (1 - \mu(2 - \mu)) \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 = \gamma^2 \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2$$

where  $\gamma = \sqrt{1 - \mu(2 - \mu)}$ . It is easy to check that  $0 \leq 1 - \mu(2 - \mu) < 1$  for all  $\mu \in (0, 2)$ . Hence,  $0 \leq \gamma < 1$ .

### 12.28

In Kaczmarz's algorithm with  $\mu = 1$ , we may write

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + (b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k)}) \frac{\mathbf{a}_{R(k)+1}}{\|\mathbf{a}_{R(k)+1}\|^2}.$$

Subtracting  $\mathbf{x}^*$  and premultiplying both sides by  $\mathbf{a}_{R(k)+1}^\top$  yields

$$\begin{aligned} \mathbf{a}_{R(k)+1}^\top (\mathbf{x}^{(k+1)} - \mathbf{x}^*) &= \mathbf{a}_{R(k)+1}^\top \left( \mathbf{x}^{(k)} - \mathbf{x}^* + (b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k)}) \frac{\mathbf{a}_{R(k)+1}}{\|\mathbf{a}_{R(k)+1}\|^2} \right) \\ &= \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k)} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^* + (b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k)}) \\ &= b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^* \\ &= 0. \end{aligned}$$

Substituting  $\mathbf{a}_{R(k)+1}^\top \mathbf{x}^* = b_{R(k)+1}$  yields the desired result.

### 12.29

We will prove this by contradiction. Suppose  $\mathbf{C}\mathbf{x}^*$  is not the minimizer of  $\|\mathbf{B}\mathbf{y} - \mathbf{b}\|^2$  over  $\mathbb{R}^r$ . Let  $\hat{\mathbf{y}}$  be the minimizer of  $\|\mathbf{B}\mathbf{y} - \mathbf{b}\|^2$  over  $\mathbb{R}^r$ . Then,  $\|\mathbf{B}\hat{\mathbf{y}} - \mathbf{b}\|^2 < \|\mathbf{B}\mathbf{C}\mathbf{x}^* - \mathbf{b}\|^2 = \|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|^2$ . Since  $\mathbf{C}$  is of full rank, there exists  $\hat{\mathbf{x}} \in \mathbb{R}^n$  such that  $\hat{\mathbf{y}} = \mathbf{C}\hat{\mathbf{x}}$ . Therefore,

$$\|\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}\|^2 = \|\mathbf{B}\mathbf{C}\hat{\mathbf{x}} - \mathbf{b}\|^2 = \|\mathbf{B}\hat{\mathbf{y}} - \mathbf{b}\|^2 < \|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|^2$$

which contradicts the assumption that  $\mathbf{x}^*$  is a minimizer of  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$  over  $\mathbb{R}^n$ .

### 12.30

a. Let  $\mathbf{A} = \mathbf{B}\mathbf{C}$  be a full rank factorization of  $\mathbf{A}$ . Now, we have  $\mathbf{A}^\dagger = \mathbf{C}^\dagger \mathbf{B}^\dagger$ , where  $\mathbf{B}^\dagger = (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top$  and  $\mathbf{C}^\dagger = \mathbf{C}^\top (\mathbf{C}\mathbf{C}^\top)^{-1}$ . On the other hand  $(\mathbf{A}^\top)^\dagger = (\mathbf{C}^\top \mathbf{B}^\top)^\dagger$ . Since  $\mathbf{A}^\top = \mathbf{C}^\top \mathbf{B}^\top$  is a full rank factorization of  $\mathbf{A}^\top$ , we have  $(\mathbf{A}^\top)^\dagger = (\mathbf{C}^\top \mathbf{B}^\top)^\dagger = (\mathbf{B}^\top)^\dagger (\mathbf{C}^\top)^\dagger$ . Therefore, to show that  $(\mathbf{A}^\top)^\dagger = (\mathbf{A}^\dagger)^\top$ , it is enough to show that

$$\begin{aligned} (\mathbf{B}^\top)^\dagger &= (\mathbf{B}^\dagger)^\top \\ (\mathbf{C}^\top)^\dagger &= (\mathbf{C}^\dagger)^\top. \end{aligned}$$

To this end, note that  $(\mathbf{B}^\top)^\dagger = \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1}$ , and  $(\mathbf{C}^\top)^\dagger = (\mathbf{C}\mathbf{C}^\top)^{-1} \mathbf{C}$ . On the other hand,  $(\mathbf{B}^\dagger)^\top = ((\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top)^\top = \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1}$ , and  $(\mathbf{C}^\dagger)^\top = (\mathbf{C}^\top (\mathbf{C}\mathbf{C}^\top)^{-1})^\top = (\mathbf{C}\mathbf{C}^\top)^{-1} \mathbf{C}$ , which completes the proof.

b. Note that  $\mathbf{A}^\dagger = \mathbf{C}^\dagger \mathbf{B}^\dagger$ , which is a full rank factorization of  $\mathbf{A}^\dagger$ . Therefore,  $(\mathbf{A}^\dagger)^\dagger = (\mathbf{B}^\dagger)^\dagger (\mathbf{C}^\dagger)^\dagger$ . Hence, to show that  $(\mathbf{A}^\dagger)^\dagger = \mathbf{A}$ , it is enough to show that

$$\begin{aligned} (\mathbf{B}^\dagger)^\dagger &= \mathbf{B} \\ (\mathbf{C}^\dagger)^\dagger &= \mathbf{C}. \end{aligned}$$

To this end, note that  $(\mathbf{B}^\dagger)^\dagger = ((\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top)^\dagger = \mathbf{B}$  since  $\mathbf{B}^\dagger$  is a full rank matrix. Similarly,  $(\mathbf{C}^\dagger)^\dagger = (\mathbf{C}^\top (\mathbf{C}\mathbf{C}^\top)^{-1})^\dagger = \mathbf{C}$  since  $\mathbf{C}^\dagger$  is a full rank matrix. This completes the proof.

### 12.31

$\Rightarrow$ : We prove properties 1–4 in turn.

1. This is immediate.

2. Let  $\mathbf{A} = \mathbf{B}\mathbf{C}$  be a full rank factorization of  $\mathbf{A}$ . We have  $\mathbf{A}^\dagger = \mathbf{C}^\dagger \mathbf{B}^\dagger$ , where  $\mathbf{B}^\dagger = (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top$  and  $\mathbf{C}^\dagger = \mathbf{C}^\top (\mathbf{C}\mathbf{C}^\top)^{-1}$ . Note that  $\mathbf{B}^\dagger \mathbf{B} = \mathbf{I}$  and  $\mathbf{C}\mathbf{C}^\dagger = \mathbf{I}$ . Now,

$$\begin{aligned} \mathbf{A}^\dagger \mathbf{A} \mathbf{A}^\dagger &= \mathbf{C}^\dagger \mathbf{B}^\dagger \mathbf{B} \mathbf{C} \mathbf{C}^\dagger \mathbf{B}^\dagger \\ &= \mathbf{C}^\dagger \mathbf{B}^\dagger \\ &= \mathbf{A}^\dagger. \end{aligned}$$

3. We have

$$\begin{aligned} (\mathbf{A} \mathbf{A}^\dagger)^\top &= (\mathbf{B} \mathbf{C} \mathbf{C}^\dagger \mathbf{B}^\dagger)^\top \\ &= (\mathbf{B} \mathbf{B}^\dagger)^\top \\ &= (\mathbf{B}^\dagger)^\top \mathbf{B}^\top \\ &= ((\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top)^\top \mathbf{B}^\top \\ &= \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \\ &= \mathbf{B} \mathbf{B}^\dagger \\ &= \mathbf{B} \mathbf{C} \mathbf{C}^\dagger \mathbf{B}^\dagger \\ &= \mathbf{A} \mathbf{A}^\dagger. \end{aligned}$$

4. We have

$$\begin{aligned}
(\mathbf{A}^\dagger \mathbf{A})^\top &= (\mathbf{C}^\dagger \mathbf{B}^\dagger \mathbf{B} \mathbf{C})^\top \\
&= (\mathbf{C}^\dagger \mathbf{C})^\top \\
&= \mathbf{C}^\top (\mathbf{C}^\dagger)^\top \\
&= \mathbf{C}^\top (\mathbf{C}^\top (\mathbf{C} \mathbf{C}^\top)^{-1})^\top \\
&= \mathbf{C}^\top (\mathbf{C} \mathbf{C}^\top)^{-1} \mathbf{C} \\
&= \mathbf{C}^\dagger \mathbf{C} \\
&= \mathbf{C}^\dagger \mathbf{B}^\dagger \mathbf{B} \mathbf{C} \\
&= \mathbf{A}^\dagger \mathbf{A}.
\end{aligned}$$

$\Leftarrow$ : By property 1, we immediately have  $\mathbf{A} \mathbf{A}^\dagger \mathbf{A} = \mathbf{A}$ . Therefore, it remains to show that there exist matrices  $\mathbf{U}$  and  $\mathbf{V}$  such that  $\mathbf{A}^\dagger = \mathbf{U} \mathbf{A}^\top$  and  $\mathbf{A}^\dagger = \mathbf{A}^\top \mathbf{V}$ .

For this, we note from property 2 that  $\mathbf{A}^\dagger = \mathbf{A}^\dagger \mathbf{A} \mathbf{A}^\dagger$ . But from property 3,  $\mathbf{A} \mathbf{A}^\dagger = (\mathbf{A} \mathbf{A}^\dagger)^\top = (\mathbf{A}^\dagger)^\top \mathbf{A}^\top$ . Hence,  $\mathbf{A}^\dagger = \mathbf{A}^\dagger (\mathbf{A}^\dagger)^\top \mathbf{A}^\top$ . Setting  $\mathbf{U} = \mathbf{A}^\dagger (\mathbf{A}^\dagger)^\top$ , we get that  $\mathbf{A}^\dagger = \mathbf{U} \mathbf{A}^\top$ .

Similarly, we note from property 4 that  $\mathbf{A}^\dagger \mathbf{A} = (\mathbf{A}^\dagger \mathbf{A})^\top = \mathbf{A}^\top (\mathbf{A}^\dagger)^\top$ . Substituting this back into property 2 yields  $\mathbf{A}^\dagger = \mathbf{A}^\dagger \mathbf{A} \mathbf{A}^\dagger = \mathbf{A}^\top (\mathbf{A}^\dagger)^\top \mathbf{A}^\dagger$ . Setting  $\mathbf{V} = (\mathbf{A}^\dagger)^\top \mathbf{A}^\dagger$  yields  $\mathbf{A}^\dagger = \mathbf{A}^\top \mathbf{V}$ . This completes the proof.

### 12.32

(Taken from [23, p. 24]) Let

$$\mathbf{A}_1 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix}, \quad \mathbf{A}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

We compute

$$\mathbf{A}_1^\dagger = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & -1 \end{bmatrix}, \quad \mathbf{A}_2^\dagger = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} = \mathbf{A}_2.$$

We have

$$\mathbf{A}_1 \mathbf{A}_2 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$$

which is a full rank factorization. Therefore,

$$(\mathbf{A}_1 \mathbf{A}_2)^\dagger = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1/2 & 1/2 \\ 0 & 0 & 0 \end{bmatrix}.$$

But

$$\mathbf{A}_2^\dagger \mathbf{A}_1^\dagger = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}.$$

Hence,  $(\mathbf{A}_1 \mathbf{A}_2)^\dagger \neq \mathbf{A}_2^\dagger \mathbf{A}_1^\dagger$ .

## 13. Unconstrained Optimization and Feedforward Neural Networks

### 13.1

a. The gradient of  $f$  is given by

$$\nabla f(\mathbf{w}) = -\mathbf{X}_d (\mathbf{y}_d - \mathbf{X}_d^\top \mathbf{w}).$$

b. The Conjugate Gradient algorithm applied to our training problem is:

1. Set  $k := 0$ ; select the initial point  $\mathbf{w}^{(0)}$ .
2.  $\mathbf{g}^{(0)} = -\mathbf{X}_d(\mathbf{y}_d - \mathbf{X}_d^\top \mathbf{w}^{(0)})$ . If  $\mathbf{g}^{(0)} = \mathbf{0}$ , stop, else set  $\mathbf{d}^{(0)} = -\mathbf{g}^{(0)}$ .
3.  $\alpha_k = -\frac{\mathbf{d}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{X}_d \mathbf{X}_d^\top \mathbf{d}^{(k)}}$
4.  $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \alpha_k \mathbf{d}^{(k)}$
5.  $\mathbf{g}^{(k+1)} = \mathbf{X}_d(\mathbf{y}_d - \mathbf{X}_d^\top \mathbf{w}^{(k+1)})$ . If  $\mathbf{g}^{(k+1)} = \mathbf{0}$ , stop.
6.  $\beta_k = \frac{\mathbf{g}^{(k+1)\top} \mathbf{X}_d \mathbf{X}_d^\top \mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{X}_d \mathbf{X}_d^\top \mathbf{d}^{(k)}}$
7.  $\mathbf{d}^{(k+1)} = -\mathbf{g}^{(k+1)} + \beta_k \mathbf{d}^{(k)}$
8. Set  $k := k + 1$ ; go to 3.

c. We form the matrix  $\mathbf{X}_d$  as

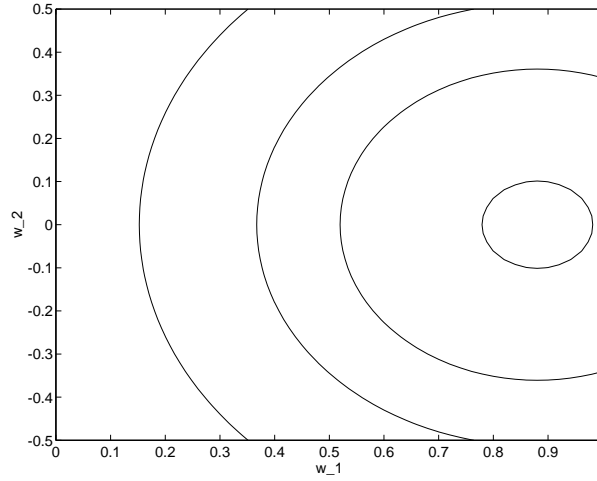
$$\mathbf{X}_d = \begin{bmatrix} -0.5 & -0.5 & -0.5 & 0 & 0 & 0 & 0.5 & 0.5 & 0.5 \\ -0.5 & 0 & 0.5 & -0.5 & 0 & 0.5 & -0.5 & 0 & 0.5 \end{bmatrix}$$

and the vector  $\mathbf{y}_d$  as

$$\mathbf{y}_d = [-0.42074, -0.47943, -0.42074, 0, 0, 0, 0.42074, 0.47943, 0.42074]^\top.$$

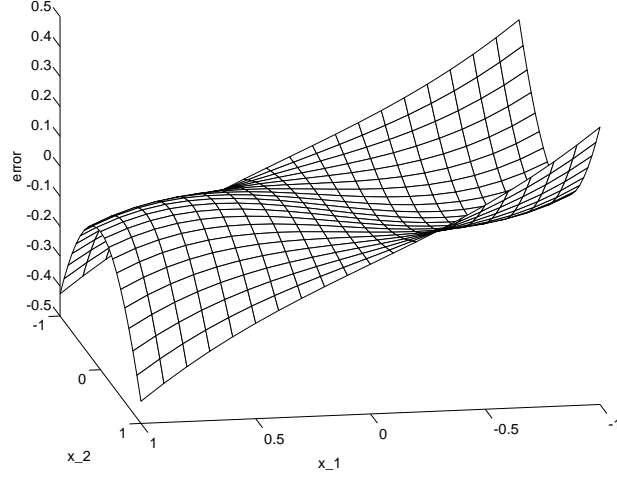
Running the Conjugate Gradient algorithm, we get a solution of  $\mathbf{w}^* = [0.8806, 0.000]^\top$ .

d. The level sets are shown in the figure below.



The solution in part c agrees with the level sets.

e. The plot of the error function is depicted below.



### 13.2

a. The expression we seek is

$$e_{k+1} = (1 - \mu)e_k.$$

To derive the above, we write

$$\begin{aligned} e_{k+1} - e_k &= y_d - \mathbf{x}_d^\top \mathbf{w}^{(k+1)} - (y_d - \mathbf{x}_d^\top \mathbf{w}^{(k)}) \\ &= -\mathbf{x}_d^\top (\mathbf{w}^{(k+1)} - \mathbf{w}^{(k)}). \end{aligned}$$

Substituting for  $\mathbf{w}^{(k+1)} - \mathbf{w}^{(k)}$  from the Widrow-Hoff algorithm yields

$$e_{k+1} - e_k = -\mu \mathbf{x}_d^\top \frac{e_k \mathbf{x}_d}{\mathbf{x}_d^\top \mathbf{x}_d} = -\mu e_k.$$

Hence,  $e_{k+1} = (1 - \mu)e_k$ .

b. For  $e_k \rightarrow 0$ , it is necessary and sufficient that  $|1 - \mu| < 1$ , which is equivalent to  $0 < \mu < 2$ .

### 13.3

a. The error satisfies

$$\mathbf{e}^{(k+1)} = (\mathbf{I}_p - \boldsymbol{\mu})\mathbf{e}^{(k)}.$$

To derive the above expression, we write

$$\begin{aligned} \mathbf{e}^{(k+1)} - \mathbf{e}^{(k)} &= \mathbf{y}_d - \mathbf{X}_d^\top \mathbf{w}^{(k+1)} - (\mathbf{y}_d - \mathbf{X}_d^\top \mathbf{w}^{(k)}) \\ &= -\mathbf{X}_d^\top (\mathbf{w}^{(k+1)} - \mathbf{w}^{(k)}). \end{aligned}$$

Substituting for  $\mathbf{w}^{(k+1)} - \mathbf{w}^{(k)}$  from the algorithm yields

$$\mathbf{e}^{(k+1)} - \mathbf{e}^{(k)} = -\mathbf{X}_d^\top \mathbf{X}_d (\mathbf{X}_d^\top \mathbf{X}_d)^{-1} \boldsymbol{\mu} \mathbf{e}^{(k)} = -\boldsymbol{\mu} \mathbf{e}^{(k)}.$$

Hence,  $\mathbf{e}^{(k+1)} = (\mathbf{I}_p - \boldsymbol{\mu})\mathbf{e}^{(k)}$ .

b. From part a, we see that  $\mathbf{e}^{(k)} = (\mathbf{I}_p - \boldsymbol{\mu})^k \mathbf{e}^{(0)}$ . Hence, by Lemma 5.1, a necessary and sufficient condition for  $\mathbf{e}^{(k)} \rightarrow \mathbf{0}$  for any  $\mathbf{e}^{(0)}$  is that all the eigenvalues of  $\mathbf{I}_p - \boldsymbol{\mu}$  must be located in the open unit circle. From Exercise 3.6, it follows that the above condition holds if and only if  $|1 - \lambda_i(\boldsymbol{\mu})| < 1$  for each eigenvalue  $\lambda_i(\boldsymbol{\mu})$  of  $\boldsymbol{\mu}$ . This is true if and only if  $0 < |\lambda_i(\boldsymbol{\mu})| < 2$  for each eigenvalue  $\lambda_i(\boldsymbol{\mu})$  of  $\boldsymbol{\mu}$ .

### 13.4

We modified the MATLAB routine of Exercise 8.25, by fixing the step size at a value  $\eta = 100$ . We need the following M-file for the gradient:

```

function y=Dfbp(w);

wh11=w(1);
wh21=w(2);
wh12=w(3);
wh22=w(4);
wo11=w(5);
wo12=w(6);

xd1=0; xd2=1; yd=1;

v1=wh11*xd1+wh12*xd2;
v2=wh21*xd1+wh22*xd2;
z1=sigmoid(v1);
z2=sigmoid(v2);
y1=sigmoid(wo11*z1+wo12*z2)
d1=(yd-y1)*y1*(1-y1);

y(1)=-d1*wo11*z1*(1-z1)*xd1;
y(2)=-d1*wo12*z2*(1-z2)*xd1;
y(3)=-d1*wo11*z1*(1-z1)*xd2;
y(4)=-d1*wo12*z2*(1-z2)*xd2;
y(5)=-d1*z1;
y(6)=-d1*z2;

y=y';

```

After 20 iterations of the backpropagation algorithm, we get the following weights:

$$\begin{aligned}
w_{11}^{o(20)} &= 2.883 \\
w_{12}^{o(20)} &= 3.194 \\
w_{11}^{h(20)} &= 0.1000 \\
w_{12}^{h(20)} &= 0.8179 \\
w_{21}^{h(20)} &= 0.3000 \\
w_{22}^{h(20)} &= 1.106.
\end{aligned}$$

The corresponding output of the network is  $y_1^{(20)} = 0.9879$ .

### 13.5

We used the following MATLAB routine:

```

function [x,N]=backprop(grad,xnew,options);
%   BACKPROP('grad',x0);
%   BACKPROP('grad',x0,OPTIONS);
%
%   x = BACKPROP('grad',x0);
%   x = BACKPROP('grad',x0,OPTIONS);
%
%   [x,N] = BACKPROP('grad',x0);
%   [x,N] = BACKPROP('grad',x0,OPTIONS);
%
%The first variant trains a net whose gradient
%is described in grad (usually an M-file: grad.m), using a backprop
%algorithm with initial point x0.
%The second variant allows a vector of optional parameters to
%defined. OPTIONS(1) controls how much display output is given; set
%to 1 for a tabular display of results, (default is no display: 0).

```

```

%OPTIONS(2) is a measure of the precision required for the final point.
%OPTIONS(3) is a measure of the precision required of the gradient.
%OPTIONS(14) is the maximum number of iterations.
%For more information type HELP FOPTIONS.
%
%The next two variants returns the value of the final point.
%The last two variants returns a vector of the final point and the
%number of iterations.

if nargin ~= 3
    options = [];
    if nargin ~= 2
        disp('Wrong number of arguments. ');
        return;
    end
end

if length(options) >= 14
    if options(14)==0
        options(14)=1000*length(xnew);
    end
else
    options(14)=1000*length(xnew);
end

clc;
format compact;
format short e;

options = foptions(options);
print = options(1);
epsilon_x = options(2);
epsilon_g = options(3);
max_iter=options(14);

for k = 1:max_iter,

    xcurr=xnew;

    g_curr=feval(grad,xcurr);

    if norm(g_curr) <= epsilon_g
        disp('Terminating: Norm of gradient less than');
        disp(epsilon_g);
        k=k-1;
        break;
    end %if

    alpha=10.0;

    xnew = xcurr-alpha*g_curr;

    if print,
        disp('Iteration number k =')
        disp(k); %print iteration index k
        disp('alpha =');
        disp(alpha); %print alpha
        disp('Gradient = ');
        disp(g_curr'); %print gradient
    end
end

```

```

        disp('New point =');
        disp(xnew'); %print new point
    end %if

    if norm(xnew-xcurr) <= epsilon_x*norm(xcurr)
        disp('Terminating: Norm of difference between iterates less than');
        disp(epsilon_x);
        break;
    end %if

    if k == max_iter
        disp('Terminating with maximum number of iterations');
    end %if
end %for

if nargout >= 1
    x=xnew;
    if nargout == 2
        N=k;
    end
else
    disp('Final point =');
    disp(xnew');
    disp('Number of iterations =');
    disp(k);
end %if
%-----

```

To apply the above routine, we need the following M-file for the gradient.

```

function y=grad(w,xd,yd);

wh11=w(1);
wh21=w(2);
wh12=w(3);
wh22=w(4);
wo11=w(5);
wo12=w(6);
t1=w(7);
t2=w(8);
t3=w(9);

xd1=xd(1); xd2=xd(2);

v1=wh11*xd1+wh12*xd2-t1;
v2=wh21*xd1+wh22*xd2-t2;
z1=sigmoid(v1);
z2=sigmoid(v2);
y1=sigmoid(wo11*z1+wo12*z2-t3);
d1=(yd-y1)*y1*(1-y1);

y(1)=-d1*wo11*z1*(1-z1)*xd1;
y(2)=-d1*wo12*z2*(1-z2)*xd1;
y(3)=-d1*wo11*z1*(1-z1)*xd2;
y(4)=-d1*wo12*z2*(1-z2)*xd2;
y(5)=-d1*z1;
y(6)=-d1*z2;
y(7)=d1*wo11*z1*(1-z1);
y(8)=d1*wo12*z2*(1-z2);
y(9)=d1;

```



```
y=y';
```

We applied our MATLAB routine as follows.

```
>> options(2)=10^(-7);
>> options(3)=10^(-7);
>> options(14)=10000;
>> w0=[0.1,0.3,0.3,0.4,0.4,0.6,0.1,0.1,-0.1]';
>> [wstar,N]=backprop('grad',w0,options)
```

Terminating with maximum number of iterations

```
wstar =
    -7.7771e+00
    -5.5932e+00
    -8.4027e+00
    -5.6384e+00
    -1.1010e+01
     1.0918e+01
    -3.2773e+00
    -8.3565e+00
     5.2606e+00
N =
    10000
```

As we can see from the above, the results coincide with Example 13.3. The table of the outputs of the trained network corresponding to the training input data is shown in Table 13.2.

## 14. Global Search Algorithms

### 14.1

---

The MATLAB program is as follows.

```
function [ output_args ] = nm_simplex( input_args )
%Nelder-Mead simplex method
%Based on the program by the Spring 2007 ECE580 student, Hengzhou Ding

disp ('We minimize a function using the Nelder-Mead method.')
disp ('There are two initial conditions.')
disp ('You can enter your own starting point.')
disp ('-----')

% disp('Select one of the starting points')
% disp ('[0.55;0.7] or [-0.9;-0.5]')
% x0=input('')
disp ('          ')
clear
close all;

disp('Select one of the starting points, or enter your own point')
disp('[0.55;0.7] or [-0.9;-0.5]')
disp('(Copy one of the above points and paste it at the prompt)')
x0=input('')

hold on
axis square
%Plot the contours of the objective function
[X1,X2]=meshgrid(-1:0.01:1);
```

```

Y=(X2-X1).^4+12.*X1.*X2-X1+X2-3;
[C,h] = contour(X1,X2,Y,20);
clabel(C,h);

% Initialize all parameters
lambda=0.1;
rho=1;
chi=2;
gamma=1/2;
sigma=1/2;
e1=[1 0]';
e2=[0 1]';
%x0=[0.55 0.7]';
%x0=[-0.9 -0.5]';

% Plot initial point and initialize the simplex
plot(x0(1),x0(2),'--*');
x(:,3)=x0;
x(:,1)=x0+lambda*e1;
x(:,2)=x0+lambda*e2;

while 1
% Check the size of simplex for stopping criterion
simpsize=norm(x(:,1)-x(:,2))+norm(x(:,2)-x(:,3))+norm(x(:,3)-x(:,1));
if(simpsize<1e-6)
break;
end
lastpt=x(:,3);
% Sort the simplex
x=sort_points(x,3);
% Reflection
centro=1/2*(x(:,1)+x(:,2));
xr=centro+rho*(centro-x(:,3));
% Accept condition
if(obj_fun(xr)>=obj_fun(x(:,1)) && obj_fun(xr)<obj_fun(x(:,2)))
x(:,3)=xr;
% Expand condition
elseif(obj_fun(xr)<obj_fun(x(:,1)))
xe=centro+rho*chi*(centro-x(:,3));
if(obj_fun(xe)<obj_fun(xr))
x(:,3)=xe;
else
x(:,3)=xr;
end

% Outside contraction or shrink
elseif(obj_fun(xr)>=obj_fun(x(:,2)) &&
obj_fun(xr)<obj_fun(x(:,3)))
xc=centro+gamma*rho*(centro-x(:,3));
if(obj_fun(xc)<obj_fun(x(:,3)))
x(:,3)=xc;
else
x=shrink(x,sigma);
end
% Inside contraction or shrink
else
xcc=centro-gamma*(centro-x(:,3));
if(obj_fun(xcc)<obj_fun(x(:,3)))
x(:,3)=xcc;

```

```

else
x=shrink(x,sigma);
end
end
% Plot the new point and connect
plot([lastpt(1),x(1,3)], [lastpt(2),x(2,3)], '--*');
end
% Output the final simplex (minimizer)
x(:,1)

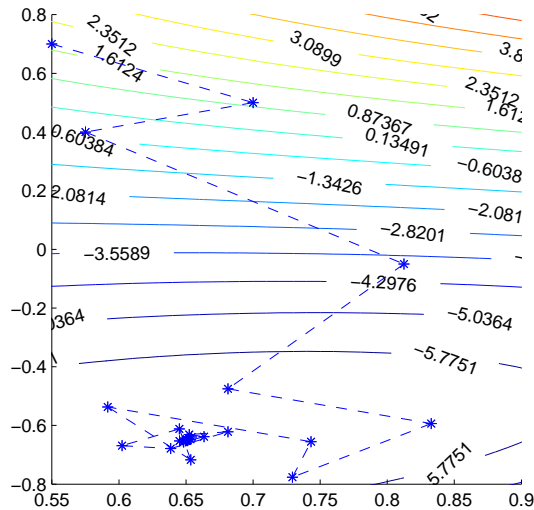
% obj_fun
function y = obj_fun(x)
y=(x(1)-x(2))4+12*x(1)*x(2)-x(1)+x(2)-3;

% sort_points
function y = sort_points(x,N)
for i=1:(N-1)
for j=1:(N-i)
if(obj_fun(x(:,j))>obj_fun(x(:,j+1)))
tmp=x(:,j);
x(:,j)=x(:,j+1);
x(:,j+1)=tmp;
end
end
end
y=x;

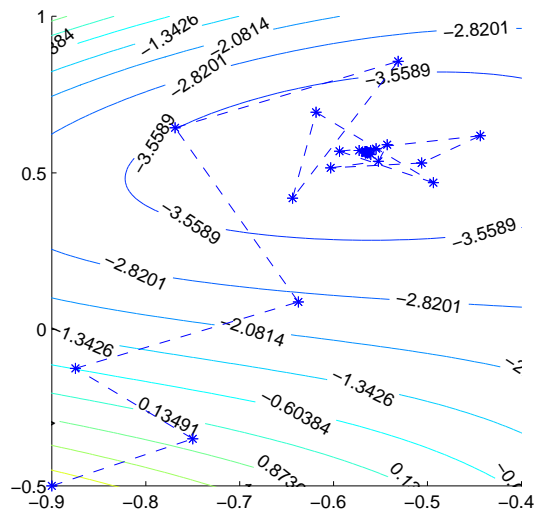
% shrink
function y = shrink(x,sigma)
x(:,2)=x(:,1)+sigma*(x(:,2)-x(:,1));
x(:,3)=x(:,1)+sigma*(x(:,3)-x(:,1));
y=x;

```

When we run the MATLAB code above with initial condition  $[0.55, 0.7]^\top$ , we obtain the following plot:



When we run the MATLAB code above with initial condition  $[-0.9, -0.5]^\top$ , we obtain the following plot:



Note that this function has two local minimizers. The algorithm terminates at these two minimizers with the two different initial conditions. This behavior depends on the value of  $\lambda$ , which determines the initial simplex. It is possible to reach both minimizers starting from the same initial point by using different values of  $\lambda$ . In the solution above, the initial simplex is “small”— $\lambda$  is just 0.1.

#### 14.2

A MATLAB routine for a naive random search algorithm is given by the M-file `rs_demo` shown below:

```
function [x,N]=random_search(funcname,xnew,options);
%Naive random search demo
% [x,N]=random_search(funcname,xnew,options);
% print = options(1);
% alpha = options(18);

if nargin ~= 3
    options = [];
    if nargin ~= 2
        disp('Wrong number of arguments. ');
        return;
    end
end

if length(options) >= 14
    if options(14)==0
        options(14)=1000*length(xnew);
    end
else
    options(14)=1000*length(xnew);
end
if length(options) < 18
    options(18)=1.0; %optional step size
end

format compact;
format short e;

options = foptions(options);
print = options(1);
epsilon_x = options(2);
```

```

epsilon_g = options(3);
max_iter=options(14);

alpha0 = options(18);

if funcname == 'f_r',
    ros_cnt
elseif funcname == 'f_p',
    pks_cnt;
end %if

if length(xnew) == 2
    plot(xnew(1),xnew(2),'o')
    text(xnew(1),xnew(2),'Start Point')
    xlower = [-2;-1];
    xupper = [2;3];
end

f_0=feval(funcname,xnew);
xbestcurr = xnew;
xbestold = xnew;
f_best=feval(funcname,xnew);
f_best=10^(sign(f_best))*f_best;

for k = 1:max_iter,

    xcurr=xbestcurr;
    f_curr=feval(funcname,xcurr);

    alpha = alpha0;
    xnew = xcurr + alpha*(2*rand(length(xcurr),1)-1);

    for i=1:length(xnew),
        xnew(i) = max(xnew(i),xlower(i));
        xnew(i) = min(xnew(i),xupper(i));
    end %for

    f_new=feval(funcname,xnew);

    if f_new < f_best,
        xbestold = xbestcurr;
        xbestcurr = xnew;
        f_best = f_new;
    end

    if print,
        disp('Iteration number k =')
        disp(k); %print iteration index k
        disp('alpha =');
        disp(alpha); %print alpha
        disp('New point =');
        disp(xnew'); %print new point
        disp('Function value =');
        disp(f_new); %print func value at new point
    end %if

    if norm(xnew-xbestold) <= epsilon_x*norm(xbestold)
        disp('Terminating: Norm of difference between iterates less than');
        disp(epsilon_x);
    end
end

```

```

        break;
    end %if

    pltpts(xbestcurr,xbestold);

    if k == max_iter
        disp('Terminating with maximum number of iterations');
    end %if
end %for

if nargout >= 1
    x=xnew;
    if nargout == 2
        N=k;
    end
else
    disp('Final point =');
    disp(xbestcurr);
    disp('Number of iterations =');
    disp(k);
end %if

```

A MATLAB routine for a simulated annealing algorithm is given by the M-file `sa_demo` shown below:

```

function [x,N]=simulated_annealing(funcname,xnew,options);
%Simulated annealing demo
% random_search(funcname,xnew,options);
% print = options(1);
% gamma = options(15);
% alpha = options(18);

if nargin ~= 3
    options = [];
    if nargin ~= 2
        disp('Wrong number of arguments. ');
        return;
    end
end

if length(options) >= 14
    if options(14)==0
        options(14)=1000*length(xnew);
    end
else
    options(14)=1000*length(xnew);
end

if length(options) < 15
    options(15)=5.0; %
end
if options(15)==0
    options(15)=5.0; %
end

if length(options) < 18
    options(18)=0.5; %optional step size
end

format compact;

```

```

format short e;

options = foptions(options);
print = options(1);
epsilon_x = options(2);
epsilon_g = options(3);
max_iter=options(14);

alpha = options(18);
gamma = options(15);
k0=2;

if funcname == 'f_r',
    ros_cnt
elseif funcname == 'f_p',
    pks_cnt;
end %if

if length(xnew) == 2
    plot(xnew(1),xnew(2),'o')
    text(xnew(1),xnew(2),'Start Point')
    xlower = [-2;-1];
    xupper = [2;3];
end

f_0=feval(funcname,xnew);
xbestcurr = xnew;
xbestold = xnew;
xcurr = xnew;
f_best=feval(funcname,xnew);
f_best=10^(sign(f_best))*f_best;

for k = 1:max_iter,

    f_curr=feval(funcname,xcurr);

    xnew = xcurr + alpha*(2*rand(length(xcurr),1)-1);

    for i=1:length(xnew),
        xnew(i) = max(xnew(i),xlower(i));
        xnew(i) = min(xnew(i),xupper(i));
    end %for

    f_new=feval(funcname,xnew);

    if f_new < f_curr,
        xcurr = xnew;
        f_curr = f_new;
    else
        cointoss = rand(1);
        Temp = gamma/log(k+k0);
        Probab = exp(-(f_new-f_curr)/Temp);
        if cointoss < Probab,
            xcurr = xnew;
            f_curr = f_new;
        end
    end

    if f_new < f_best,

```

```

    xbestold = xbestcurr;
    xbestcurr = xnew;
    f_best = f_new;
end

if print,
    disp('Iteration number k =')
    disp(k); %print iteration index k
    disp('alpha =');
    disp(alpha); %print alpha
    disp('New point =');
    disp(xnew'); %print new point
    disp('Function value =');
    disp(f_new); %print func value at new point
end %if

if norm(xnew-xbestold) <= epsilon_x*norm(xbestold)
    disp('Terminating: Norm of difference between iterates less
than');
    disp(epsilon_x);
    break;
end %if

pltpts(xbestcurr,xbestold);

if k == max_iter
    disp('Terminating with maximum number of iterations');
end %if
end %for

if nargout >= 1
    x=xnew;
    if nargout == 2
        N=k;
    end
else
    disp('Final point =');
    disp(xbestcurr');
    disp('Objective function value =');
    disp(f_best);
    disp('Number of iterations =');
    disp(k);
end %if

```

To use the above routines, we also need the following M-files:  
pltpts.m:

```

function out=pltpts(xnew,xcurr)

plot([xcurr(1),xnew(1)],[xcurr(2),xnew(2)],'r-',xnew(1),xnew(2),'o','Erasemode',
'none');
drawnow; % Draws current graph now
%pause(1)

out = [];

```

f\_p.m:

```

function y=f_p(x);

```



```

y=3*(1-x(1)).^2.*exp(-(x(1).^2)-(x(2)+1).^2) -
10.*(x(1)/5-x(1).^3-x(2).^5).*exp
(-(x(1).^2-x(2).^2) - exp(-(x(1)+1).^2-x(2).^2)/3;
y=-y;

```

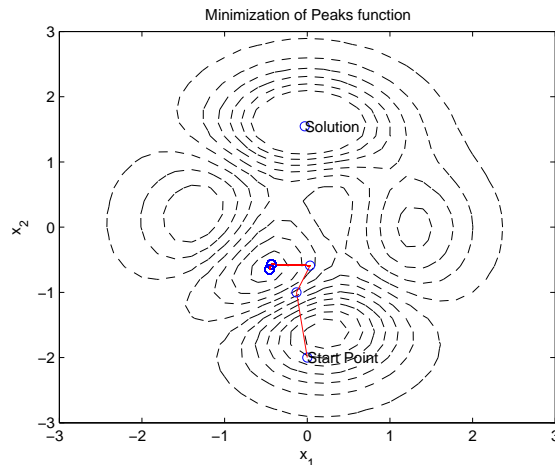
```

pks_cnt.m:

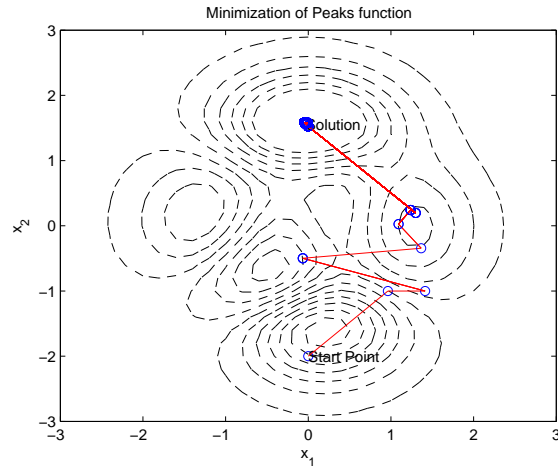
echo off
X = [-3:0.2:3]';
Y = [-3:0.2:3]';
[x,y]=meshgrid(X,Y) ;
func = 3*(1-x).^2.*exp(-x.^2-(y+1).^2) -
10.*(x/5-x.^3-y.^5).*exp(-x.^2-y.^2) -
exp(-(x+1).^2-y.^2)/3;
func = -func;
clf
levels = exp(-5:10);
levels = [-5:0.9:10];
contour(X,Y,func,levels,'k--')
xlabel('x_1')
ylabel('x_2')
title('Minimization of Peaks function')
drawnow;
hold on
plot(-0.0303,1.5455,'o')
text(-0.0303,1.5455,'Solution')

```

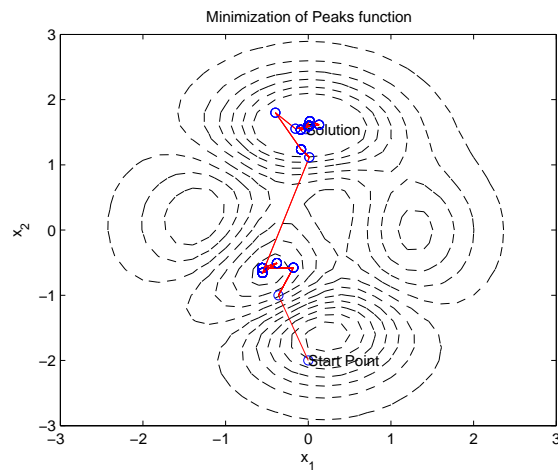
To run the naive random search algorithms, we first pick a value of  $\alpha = 0.5$ , which involves setting `options(18)=0.5`. We then use the command `rs_demo('f_p',[0;-2],options)`. The resulting plot of the algorithm trajectory is given below. As we can see, the algorithm is stuck at a local minimizer. (By running the algorithm several times, the reader can verify that this non-convergent behavior is typical.)



Next, we try  $\alpha = 1.5$ , which involves setting `options(18)=1.5`. We then use the command `rs_demo('f_p',[0;-2],options)` again, to obtain the plot shown below. This time, the algorithm reaches the global minimizer.



Finally, we again set  $\alpha = 0.5$ , using `options(18)=0.5`. We then run the simulated annealing code using `sa_demo('f_p',[0;-2],options)`. The algorithm can be seen to converge to the global minimizer, as plotted below.



### 14.3

A MATLAB routine for a particle swarm algorithm is:

```
% A particle swarm optimizer
% to find the minimum/maximum of the MATLABs' peaks function
% D---# of inputs to the function (dimension of problem)
clear
%Parameters
ps=10;
D=2;
ps_lb=-3;
ps_ub=3;
vel_lb=-1;
vel_ub=1;
iteration_n = 50;
range = [-3, 3, -3, 3]; % Range of the input variables
% Plot contours of peaks function
[x, y, z] = peaks;
%pcolor(x,y,z); shading interp; hold on;
```

```

    %contour(x, y, z, 20, 'r');
    mesh(x,y,z)
%    hold off;
    %colormap(gray);
    set(gca,'FontSize',14)
    axis([-3 3 -3 3 -9 9])
%axis square;
    xlabel('x_1','FontSize',14);
    ylabel('x_2','FontSize',14);
    zlabel('f(x_1,x_2)','FontSize',14);
    hold on

upper = zeros(iteration_n, 1);
average = zeros(iteration_n, 1);
lower = zeros(iteration_n, 1);

% initialize population of particles and their velocities at time
% zero,
% format of pos= (particle#, dimension)
% construct random population positions bounded by VR
% need to bound positions

ps_pos=ps_lb + (ps_ub-ps_lb).*rand(ps,D);

% need to bound velocities between -mv,mv

ps_vel=vel_lb + (vel_ub-vel_lb).*rand(ps,D);

% initial pbest positions

p_best = ps_pos;

% returns column of cost values (1 for each particle)

f1='3*(1-ps_pos(i,1))^2*exp(-ps_pos(i,1)^2-(ps_pos(i,2)+1)^2)';
f2='-10*(ps_pos(i,1)/5-ps_pos(i,1)^3-ps_pos(i,2)^5)*exp(-ps_pos(i,1)^2-ps_pos(i,2)^2)';
f3='-(1/3)*exp(-(ps_pos(i,1)+1)^2-ps_pos(i,2)^2)';

p_best_fit=zeros(ps,1);
for i=1:ps
    g1(i)=3*(1-ps_pos(i,1))^2*exp(-ps_pos(i,1)^2-(ps_pos(i,2)+1)^2);
    g2(i)=-10*(ps_pos(i,1)/5-ps_pos(i,1)^3-ps_pos(i,2)^5)*exp(-ps_pos(i,1)^2-ps_pos(i,2)^2);
    g3(i)=-(1/3)*exp(-(ps_pos(i,1)+1)^2-ps_pos(i,2)^2);
    p_best_fit(i)=g1(i)+g2(i)+g3(i);
end
p_best_fit;
hand_p3=plot3(ps_pos(:,1),ps_pos(:,2),p_best_fit','*k','markersize',15,'erase','xor');

% initial g_best

[g_best_val,g_best_idx] = max(p_best_fit);
%[g_best_val,g_best_idx] = min(p_best_fit); this is to minimize
g_best=ps_pos(g_best_idx,:);

% get new velocities, positions (this is the heart of the PSO
% algorithm)
for k=1:iteration_n
    for count=1:ps
        ps_vel(count,:) = 0.729*ps_vel(count,...           % prev vel

```

```

        +1.494*rand*(p_best(count,:)-ps_pos(count,:))...    % independent
        +1.494*rand*(g_best-ps_pos(count,:));              % social
    end
    ps_vel;
% update new position
    ps_pos = ps_pos + ps_vel;

%update p_best
for i=1:ps
    g1(i)=3*(1-ps_pos(i,1))^2*exp(-ps_pos(i,1)^2-(ps_pos(i,2)+1)^2);
    g2(i)=-10*(ps_pos(i,1)/5-ps_pos(i,1)^3-ps_pos(i,2)^5)*exp(-ps_pos(i,1)^2-ps_pos(i,2)^2);
    g3(i)=-(1/3)*exp(-(ps_pos(i,1)+1)^2-ps_pos(i,2)^2);
    ps_current_fit(i)=g1(i)+g2(i)+g3(i);

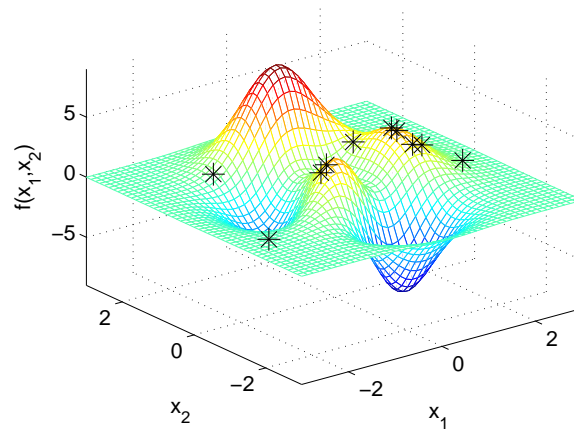
    if ps_current_fit(i)>p_best_fit(i)
        p_best_fit(i)=ps_current_fit(i);
        p_best(i,:)=ps_pos(i,:);
    end
end
p_best_fit;
%update g_best
[g_best_val,g_best_idx] = max(p_best_fit);
g_best=ps_pos(g_best_idx,:);

% Fill objective function vectors
    upper(k) = max(p_best_fit);
    average(k) = mean(p_best_fit);
lower(k) = min(p_best_fit);

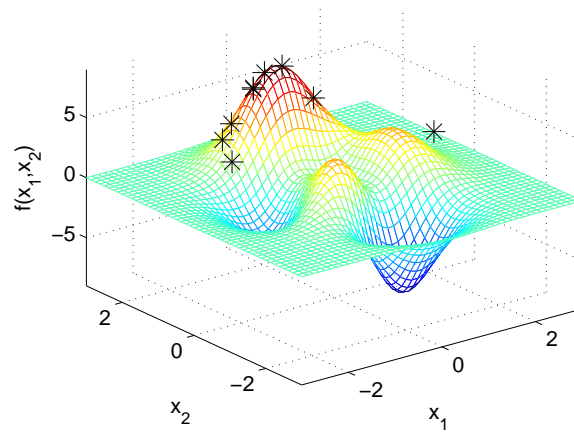
set(hand_p3,'xdata',ps_pos(:,1),'ydata',ps_pos(:,2),'zdata',ps_current_fit');
drawnow
pause
end
g_best
g_best_val
    figure;
x = 1:iteration_n;
plot(x, upper, 'o', x, average, 'x', x, lower, '*');
hold on;
plot(x, [upper average lower]);
hold off;
legend('Best', 'Average', 'Poorest');
xlabel('Iterations'); ylabel('Objective function value');

```

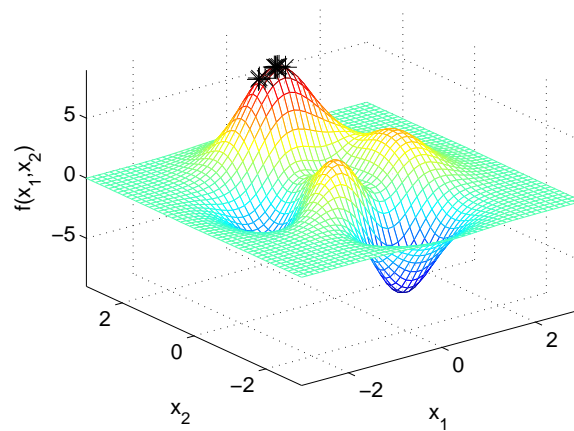
When we run the MATLAB code above, we obtain a plot of the initial set of particles, as shown below.



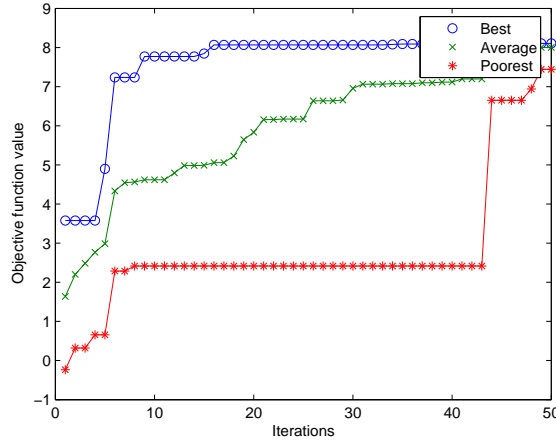
Then, after 30 iterations, we obtain:



Finally, after 50 iterations, we obtain:



A plot of the objective function values (best, average, and poorest) is shown below.



#### 14.4

- Expanding the right hand side of the second expression gives the desired result.
- Applying the algorithm, we get a binary representation of 11111001011, i.e.,

$$1995 = 2^{10} + 2^9 + 2^8 + 2^7 + 2^6 + 2^3 + 2^1 + 2^0.$$

- Applying the algorithm, we get a binary representation of 0.1011101, i.e.,

$$0.7265625 = 2^{-1} + 2^{-3} + 2^{-4} + 2^{-5} + 2^{-7}.$$

- We have  $19 = 2^4 + 2^1 + 2^0$ , i.e., the binary representation for 19 is 10011. For the fractional part, we need at least 7 bits to keep at least the same accuracy. We have  $0.95 = 2^{-1} + 2^{-2} + 2^{-3} + 2^{-4} + 2^{-7} + \dots$ , i.e., the binary representation is 0.1111001... Therefore, the binary representation of 19.95 with at least the same degree of accuracy is 10011.1111001.

#### 14.5

It suffices to prove the result for the case where only one symbol is swapped, since the general case is obtained by repeating the argument. We have two scenarios. First, suppose the symbol swapped is at a position corresponding to a don't care symbol in  $H$ . Clearly, after the swap, both chromosomes will still be in  $H$ . Second, suppose the symbol swapped is at a position corresponding to a fixed symbol in  $H$ . Since both chromosomes are in  $H$ , their symbols at that position must be identical. Hence, the swap does not change the chromosomes. This completes the proof.

#### 14.6

Consider a given chromosome in  $M(k) \cap H$ . The probability that it is chosen for crossover is  $q_c$ . If neither of its offsprings is in  $H$ , then at least one of the crossover points must be between the corresponding first and last fixed symbols of  $H$ . The probability of this is  $1 - (1 - \delta(H)/(L - 1))^2$ . To see this, note that the probability that each crossover point is not between the corresponding first and last fixed symbols is  $1 - \delta(H)/(L - 1)$ , and thus the probability that both crossover points are not between the corresponding first and last fixed symbols of  $H$  is  $(1 - \delta(H)/(L - 1))^2$ . Hence, the probability that the given chromosome is chosen for crossover and neither of its offsprings is in  $H$  is bounded above by

$$q_c \left( 1 - \left( 1 - \frac{\delta(H)}{L - 1} \right)^2 \right).$$

---

**14.7**

As for two-point crossover, the  $n$ -point crossover operation is a composition of  $n$  one-point crossover operations (i.e.,  $n$  one-point crossover operations in succession). The required result for this case is as follows.

**Lemma:**

Given a chromosome in  $M(k) \cap H$ , the probability that it is chosen for crossover and neither of its offsprings is in  $H$  is bounded above by

$$q_c \left( 1 - \left( 1 - \frac{\delta(H)}{L-1} \right)^n \right).$$

□

For the proof, proceed as in the solution of Exercise 14.6, replacing 2 by  $n$ .

---

**14.8**

```
function M=roulette_wheel(fitness);
%function M=roulette_wheel(fitness)
%fitness = vector of fitness values of chromosomes in population
%M = vector of indices indicating which chromosome in the
%   given population should appear in the mating pool

fitness = fitness - min(fitness); % to keep the fitness positive
if sum(fitness) == 0,
    disp('Population has identical chromosomes -- STOP');
    break;
else
    fitness = fitness/sum(fitness);
end
cum_fitness = cumsum(fitness);
for i = 1:length(fitness),
    tmp = find(cum_fitness-rand>0);
    M(i) = tmp(1);
end
```

---

**14.9**

```
% parent1, parent2 = two binary parent chromosomes (row vectors)

L = length(parent1);
crossover_pt = ceil(rand*(L-1));
offspring1 = [parent1(1:crossover_pt) parent2(crossover_pt+1:L)];
offspring2 = [parent2(1:crossover_pt) parent1(crossover_pt+1:L)];
```

---

**14.10**

```
% mating_pool = matrix of 0-1 elements; each row represents a chromosome
% p_m = probability of mutation
N = size(mating_pool,1);
L = size(mating_pool,2);
mutation_points = rand(N,L) < p_m;
new_population = xor(mating_pool,mutation_points);
```

---

**14.11**

A MATLAB routine for a genetic algorithm with binary encoding is:

```

function [winner,bestfitness] = ga(L,N,fit_func,options)
% function winner = GA(L,N,fit_func)
% Function call: GA(L,N,'f')
% L = length of chromosomes
% N = population size (must be an even number)
% f = name of fitness value function
%
%Options:
%print = options(1);
%selection = options(5);
%max_iter=options(14);
%p_c = options(18);
%p_m = p_c/100;
%
%Selection:
% options(5) = 0 for roulette wheel, 1 for tournament

clf;

if nargin ~= 4
    options = [];
    if nargin ~= 3
        disp('Wrong number of arguments.');
```

return;

end

end

```

if length(options) >= 14
    if options(14)==0
        options(14)=3*N;
    end
else
    options(14)=3*N;
end
if length(options) < 18
    options(18)=0.75; %optional crossover rate
end

%format compact;
%format short e;

options = foptions(options);
print = options(1);
selection = options(5);
max_iter=options(14);
p_c = options(18);
p_m = p_c/100;

P = rand(N,L)>0.5;
bestvaluesofar = 0;

%Initial evaluation
for i = 1:N,
    fitness(i) = feval(fit_func,P(i,:));
end
[bestvalue,best] = max(fitness);
if bestvalue > bestvaluesofar,
    bestsofar = P(best,:);
    bestvaluesofar = bestvalue;
end

```



```

end

for k = 1:max_iter,

    %Selection
    fitness = fitness - min(fitness); % to keep the fitness positive
    if sum(fitness) == 0,
        disp('Population has identical chromosomes -- STOP');
        disp('Number of iterations:');
        disp(k);
        for i = k:max_iter,
            upper(i)=upper(i-1);
            average(i)=average(i-1);
            lower(i)=lower(i-1);
        end
        break;
    else
        fitness = fitness/sum(fitness);
    end
    if selection == 0,
        %roulette-wheel
        cum_fitness = cumsum(fitness);
        for i = 1:N,
            tmp = find(cum_fitness-rand>0);
            m(i) = tmp(1);
        end
    else
        %tournament
        for i = 1:N,
            fighter1=ceil(rand*N);
            fighter2=ceil(rand*N);
            if fitness(fighter1)>fitness(fighter2),
                m(i) = fighter1;
            else
                m(i) = fighter2;
            end
        end
    end
    M = zeros(N,L);
    for i = 1:N,
        M(i,:) = P(m(i),:);
    end

    %Crossover
    Mnew = M;
    for i = 1:N/2
        ind1 = ceil(rand*N);
        ind2 = ceil(rand*N);
        parent1 = M(ind1,:);
        parent2 = M(ind2,:);
        if rand < p_c
            crossover_pt = ceil(rand*(L-1));
            offspring1 = [parent1(1:crossover_pt) parent2(crossover_pt+1:L)];
            offspring2 = [parent2(1:crossover_pt) parent1(crossover_pt+1:L)];
            Mnew(ind1,:) = offspring1;
            Mnew(ind2,:) = offspring2;
        end
    end
end

```

```

%Mutation
mutation_points = rand(N,L) < p_m;
P = xor(Mnew,mutation_points);

%Evaluation
for i = 1:N,
    fitness(i) = feval(fit_func,P(i,:));
end
[bestvalue,best] = max(fitness);
if bestvalue > bestvaluesofar,
    bestsofar = P(best,:);
    bestvaluesofar = bestvalue;
end

upper(k) = bestvalue;
average(k) = mean(fitness);
lower(k) = min(fitness);

end %for

if k == max_iter,
    disp('Algorithm terminated after maximum number of iterations:');
    disp(max_iter);
end

winner = bestsofar;
bestfitness = bestvaluesofar;

if print,
    iter = [1:max_iter]';
    plot(iter,upper,'o:',iter,average,'x-',iter,lower,'*--');
    legend('Best', 'Average', 'Worst');
    xlabel('Generations','FontSize',14);
    ylabel('Objective Function Value','FontSize',14);
    set(gca,'FontSize',14);
    hold off;
end

```

a. To run the routine, we create the following M-files.

```

function dec = bin2dec(bin,range);
%function dec = bin2dec(bin,range);
%Function to convert from binary (bin) to decimal (dec) in a given range

index = polyval(bin,2);
dec = index*((range(2)-range(1))/(2^length(bin)-1)) + range(1);

function y=f_manymax(x);

y=-15*(sin(2*x))^2-(x-2)^2+160;

function y=fit_func1(binchrom);
%1-D fitness function

f='f_manymax';
range=[-10,10];
x=bin2dec(binchrom,range);
y=feval(f,x);

```

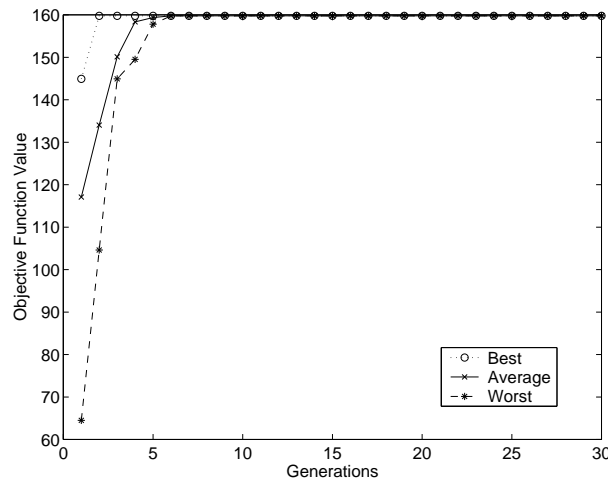
We use the following script to run the algorithm:

```

clear;
options(1)=1;
[x,y]=ga(8,10,'fit_func1',options);
f='f_manymax';
range=[-10,10];
disp('GA Solution:');
disp(bin2dec(x,range));
disp('Objective function value:');
disp(y);

```

Running the above algorithm, we obtain a solution of  $x^* = 1.6078$ , and an objective function value of 159.7640. The figure below shows a plot of the best, average, and worst solution from each generation of the population.



b. To run the routine, we create the following M-files (we also use the routine `bin2dec` from part a.

```

function y=f_peaks(x);

y=3*(1-x(1)).^2.*exp(-(x(1).^2)-(x(2)+1).^2) -
10.*(x(1)/5-x(1).^3-x(2).^5).*exp
(-x(1).^2-x(2).^2) - exp(-(x(1)+1).^2-x(2).^2)/3;

function y=fit_func2(binchrom);
%2-D fitness function

f='f_peaks';
xrange=[-3,3];
yrange=[-3,3];
L=length(binchrom);
x1=bin2dec(binchrom(1:L/2),xrange);
x2=bin2dec(binchrom(L/2+1:L),yrange);
y=feval(f,[x1,x2]);

```

We use the following script to run the algorithm:

```

clear;
options(1)=1;
[x,y]=ga(16,20,'fit_func2',options);
f='f_peaks';
xrange=[-3,3];

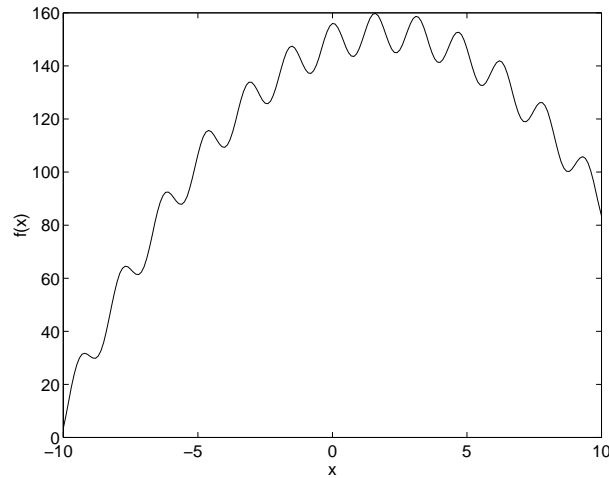
```

```

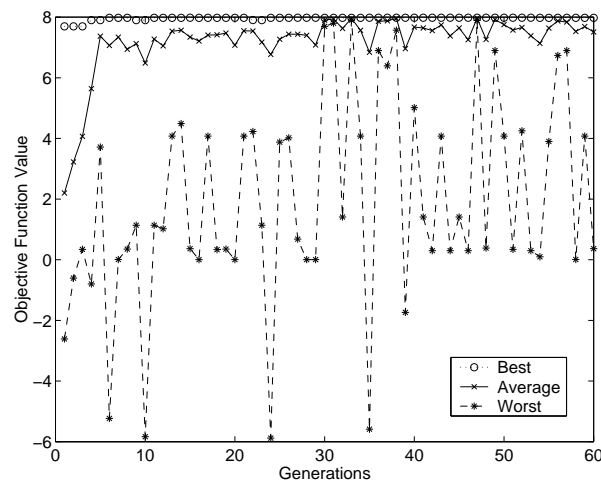
yrange=[-3,3];
L=length(x);
x1=bin2dec(x(1:L/2),xrange);
x2=bin2dec(x(L/2+1:L),yrange);
disp('GA Solution:');
disp([x1,x2]);
disp('Objective function value:');
disp(y);

```

A plot of the objective function is shown below.



Running the above algorithm, we obtain a solution of  $[-0.0353, 1.4941]^T$ , and an  $\mathbf{x}^* = [-0.0588, 1.5412]^T$ , and an objective function value of 7.9815. (Compare this solution with that of Example 14.3.) The figure below shows a plot of the best, average, and worst solution from each generation of the population.



#### 14.12

A MATLAB routine for a real-number genetic algorithm:

```

function [winner,bestfitness] = gar(Domain,N,fit_func,options)
% function winner = GAR(Domain,N,fit_func)
% Function call: GAR(Domain,N,'f')
% Domain = search space; e.g., [-2,2;-3,3] for the space [-2,2]x[-3,3]

```

```

% (number of rows of Domain = dimension of search space)
% N = population size (must be an even number)
% f = name of fitness value function
%
%Options:
%print = options(1);
%selection = options(5);
%max_iter=options(14);
%p_c = options(18);
%p_m = p_c/100;
%
%Selection:
% options(5) = 0 for roulette wheel, 1 for tournament

clf;

if nargin ~= 4
    options = [];
    if nargin ~= 3
        disp('Wrong number of arguments. ');
        return;
    end
end

if length(options) >= 14
    if options(14)==0
        options(14)=3*N;
    end
else
    options(14)=3*N;
end
if length(options) < 18
    options(18)=0.75; %optional crossover rate
end

%format compact;
%format short e;

options = foptions(options);
print = options(1);
selection = options(5);
max_iter=options(14);
p_c = options(18);
p_m = p_c/100;
n = size(Domain,1);
lowb = Domain(:,1)';
upb = Domain(:,2)';
bestvaluesofar = 0;
for i = 1:N,
    P(i,:) = lowb + rand(1,n).*(upb-lowb);
    %Initial evaluation
    fitness(i) = feval(fit_func,P(i,:));
end
[bestvalue,best] = max(fitness);
if bestvalue > bestvaluesofar,
    bestsofar = P(best,:);
    bestvaluesofar = bestvalue;
end

```

```

for k = 1:max_iter,

    %Selection
    fitness = fitness - min(fitness); % to keep the fitness positive
    if sum(fitness) == 0,
        disp('Population has identical chromosomes -- STOP');
        disp('Number of iterations:');
        disp(k);
        for i = k:max_iter,
            upper(i)=upper(i-1);
            average(i)=average(i-1);
            lower(i)=lower(i-1);
        end
        break;
    else
        fitness = fitness/sum(fitness);
    end
    if selection == 0,
        %roulette-wheel
        cum_fitness = cumsum(fitness);
        for i = 1:N,
            tmp = find(cum_fitness-rand>0);
            m(i) = tmp(1);
        end
    else
        %tournament
        for i = 1:N,
            fighter1=ceil(rand*N);
            fighter2=ceil(rand*N);
            if fitness(fighter1)>fitness(fighter2),
                m(i) = fighter1;
            else
                m(i) = fighter2;
            end
        end
    end
    end
    M = zeros(N,n);
    for i = 1:N,
        M(i,:) = P(m(i),:);
    end

    %Crossover
    Mnew = M;
    for i = 1:N/2
        ind1 = ceil(rand*N);
        ind2 = ceil(rand*N);
        parent1 = M(ind1,:);
        parent2 = M(ind2,:);
        if rand < p_c
            a = rand;
            offspring1 = a*parent1+(1-a)*parent2+(rand(1,n)-0.5).*(upb-lowb)/10;
            offspring2 = a*parent2+(1-a)*parent1+(rand(1,n)-0.5).*(upb-lowb)/10;
            %do projection
            for j = 1:n,
                if offspring1(j)<lowb(j),
                    offspring1(j)=lowb(j);
                elseif offspring1(j)>upb(j),
                    offspring1(j)=upb(j);
                end
            end
        end
    end
end

```

```

        if offspring2(j)<lowb(j),
            offspring2(j)=lowb(j);
        elseif offspring2(j)>upb(j),
            offspring2(j)=upb(j);
        end
    end
    Mnew(ind1,:) = offspring1;
    Mnew(ind2,:) = offspring2;
end
end

%Mutation
for i = 1:N,
if rand < p_m,
    a = rand;
    Mnew(i,:) = a*Mnew(i,:) + (1-a)*(lowb + rand(1,n).*(upb-lowb));
end
end

P = Mnew;

%Evaluation
for i = 1:N,
    fitness(i) = feval(fit_func,P(i,:));
end
[bestvalue,best] = max(fitness);
if bestvalue > bestvaluesofar,
    bestsofar = P(best,:);
    bestvaluesofar = bestvalue;
end

upper(k) = bestvalue;
average(k) = mean(fitness);
lower(k) = min(fitness);

end %for

if k == max_iter,
    disp('Algorithm terminated after maximum number of iterations:');
    disp(max_iter);
end

winner = bestsofar;
bestfitness = bestvaluesofar;

if print,
    iter = [1:max_iter]';
    plot(iter,upper,'o:',iter,average,'x-',iter,lower,'*--');
    legend('Best', 'Average', 'Worst');
    xlabel('Generations','FontSize',14);
    ylabel('Objective Function Value','FontSize',14);
    set(gca,'FontSize',14);
    hold off;
end
end

```

To run the routine, we create the following M-file for the given function.

```

function y=f_wave(x);

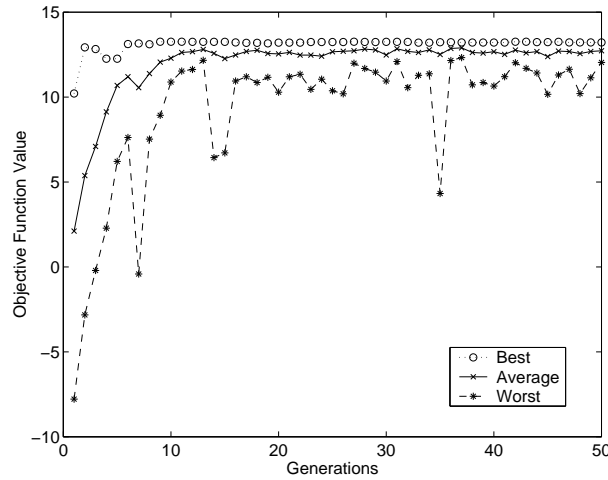
y=x(1)*sin(x(1)) + x(2)*sin(5*x(2));

```

We use the following script to run the algorithm:

```
options(1)=1;
options(14)=50;
[x,y]=gar([0,10;4,6],20,'f_wave',options);
disp('GA Solution:');
disp(x);
disp('Objective function value:');
disp(y);
```

Running the above algorithm, we obtain a solution of  $\mathbf{x}^* = [7.9711, 5.3462]^\top$ , and an objective function value of 13.2607. The figure below shows a plot of the best, average, and worst solution from each generation of the population.



Using the MATLAB function `fminunc` (from the Optimization Toolbox), we found the optimal point to be  $[7.9787, 5.3482]^\top$ , with objective function value 13.2612. We can see that this solution agrees with the solution obtained using our real-number genetic algorithm.



## 15. Introduction to Linear Programming

### 15.1

---

$$\begin{aligned}
 & \text{minimize} && -2x_1 - x_2 \\
 & \text{subject to} && x_1 + x_3 = 2 \\
 & && x_1 + x_2 + x_4 = 3 \\
 & && x_1 + 2x_2 + x_5 = 5 \\
 & && x_1, \dots, x_5 \geq 0
 \end{aligned}$$

### 15.2

---

We have

$$x_2 = ax_1 + bu_1 = a^2x_0 + abu_0 + bu_1 = a^2 + [ab, b]\mathbf{u}$$

where  $\mathbf{u} = [u_0, u_1]^\top$  is the decision variable. We can write the constraint as  $u_i \leq 1$  and  $u_i \geq -1$ . Hence, the problem is:

$$\begin{aligned}
 & \text{minimize} && a^2 + [ab, b]\mathbf{u} \\
 & \text{subject to} && -1 \leq u_i \leq 1, \quad i = 1, 2.
 \end{aligned}$$

Since  $a^2$  is a constant, we can remove it from the objective function without changing the solution. Introducing slack variables  $v_1, v_2, v_3, v_4$ , we obtain the standard form problem

$$\begin{aligned}
 & \text{minimize} && [ab, b]\mathbf{u} \\
 & \text{subject to} && u_0 + v_1 = 1 \\
 & && -u_0 + v_2 = 1 \\
 & && u_1 + v_3 = 1 \\
 & && -u_1 + v_4 = 1.
 \end{aligned}$$

### 15.3

---

Let  $x_i^+, x_i^- \geq 0$  be such that  $|x_i| = x_i^+ + x_i^-$ ,  $x_i = x_i^+ - x_i^-$ . Substituting into the original problem, we have

$$\begin{aligned}
 & \text{minimize} && c_1(x_1^+ + x_1^-) + c_2(x_2^+ + x_2^-) + \dots + c_n(x_n^+ + x_n^-) \\
 & \text{subject to} && \mathbf{A}(\mathbf{x}^+ - \mathbf{x}^-) = \mathbf{b} \\
 & && \mathbf{x}^+, \mathbf{x}^- \geq \mathbf{0},
 \end{aligned}$$

where  $\mathbf{x}^+ = [x_1^+, \dots, x_n^+]^\top$  and  $\mathbf{x}^- = [x_1^-, \dots, x_n^-]^\top$ . Rewriting, we get

$$\begin{aligned}
 & \text{minimize} && [\mathbf{c}^\top, \mathbf{c}^\top]\mathbf{z} \\
 & \text{subject to} && [\mathbf{A}, -\mathbf{A}]\mathbf{z} = \mathbf{b} \\
 & && \mathbf{z} \geq \mathbf{0},
 \end{aligned}$$

which is an equivalent linear programming problem in standard form.

Note that although the variables  $x_i^+$  and  $x_i^-$  in the solution are required to satisfy  $x_i^+x_i^- = 0$ , we do not need to explicitly include this in the constraint because any optimal solution to the above transformed problem automatically satisfies the condition  $x_i^+x_i^- = 0$ . To see this, suppose we have an optimal solution with both  $x_i^+ > 0$  and  $x_i^- > 0$ . In this case, note that  $c_i > 0$  (for otherwise we can add any arbitrary constant to both  $x_i^+$  and  $x_i^-$  and still satisfy feasibility, but decrease the objective function value). Then, by subtracting  $\min(x_i^+, x_i^-)$  from  $x_i^+$  and  $x_i^-$ , we have a new feasible point with lower objective function value, contradicting the optimality assumption. [See also M. A. Dahleh and I. J. Diaz-Bobillo, *Control of Uncertain Systems: A Linear Programming Approach*, Prentice Hall, 1995, pp. 189–190.]

## 15.4

Not every linear programming problem in standard form has a nonempty feasible set. Example:

$$\begin{array}{ll}\text{minimize} & x_1 \\ \text{subject to} & -x_1 = 1 \\ & x_1 \geq 0.\end{array}$$

Not every linear programming problem in standard form (even assuming a nonempty feasible set) has an optimal solution. Example:

$$\begin{array}{ll}\text{minimize} & -x_1 \\ \text{subject to} & x_2 = 1 \\ & x_1, x_2 \geq 0.\end{array}$$

## 15.5

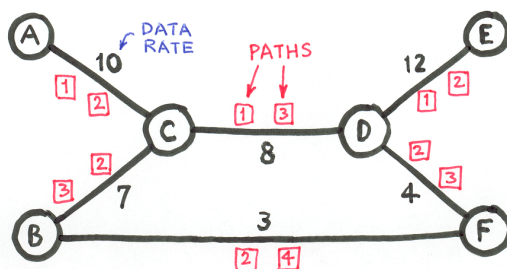
Let  $x_1$  and  $x_2$  represent the number of units to be shipped from A to C and to D, respectively, and  $x_3$  and  $x_4$  represent the number of units to be shipped from B to C and to D, respectively. Then, the given problem can be formulated as the following linear program:

$$\begin{array}{ll}\text{minimize} & x_1 + 2x_2 + 3x_3 + 4x_4 \\ \text{subject to} & x_1 + x_3 = 50 \\ & x_2 + x_4 = 60 \\ & x_1 + x_2 \leq 70 \\ & x_3 + x_4 \leq 80 \\ & x_1, x_2, x_3, x_4 \geq 0.\end{array}$$

Introducing slack variables  $x_5$  and  $x_6$ , we have the standard form problem

$$\begin{array}{ll}\text{minimize} & x_1 + 2x_2 + 3x_3 + 4x_4 \\ \text{subject to} & x_1 + x_3 = 50 \\ & x_2 + x_4 = 60 \\ & x_1 + x_2 + x_5 = 70 \\ & x_3 + x_4 + x_6 = 80 \\ & x_1, x_2, x_3, x_4, x_5, x_6 \geq 0.\end{array}$$

## 15.6



We can see that there are two paths from A to E ( $ACDE$  and  $ACBFDE$ ), and two paths from B to F ( $BCDF$  and  $BF$ ). Let  $x_1$  and  $x_2$  be the data rates for the two paths from A to E, respectively, and  $x_3$  and  $x_4$  the data rates for the two paths from B to F, respectively. The total revenue is then  $2(x_1 + x_2) + 3(x_3 + x_4)$ . For each link, we have a data rate constraint on the sum of all  $x_i$ s passing through that link. For example,

for link  $BC$ , there are two paths passing through it, with total data rate  $x_2 + x_3$ . Hence, the constraint for link  $BC$  is  $x_2 + x_3 \leq 7$ . Hence, the optimization problem is the following linear programming problem:

$$\begin{array}{ll} \text{maximize} & 2(x_1 + x_2) + 3(x_3 + x_4) \\ \text{subject to} & x_1 + x_2 \leq 10 \\ & x_1 + x_2 \leq 12 \\ & x_1 + x_3 \leq 8 \\ & x_2 + x_3 \leq 7 \\ & x_2 + x_3 \leq 4 \\ & x_2 + x_4 \leq 3 \\ & x_1, \dots, x_4 \geq 0. \end{array}$$

Converting this to standard form:

$$\begin{array}{ll} \text{minimize} & -2(x_1 + x_2) - 3(x_3 + x_4) \\ \text{subject to} & x_1 + x_2 + x_5 = 10 \\ & x_1 + x_2 + x_6 = 12 \\ & x_1 + x_3 + x_7 = 8 \\ & x_2 + x_3 + x_8 = 7 \\ & x_2 + x_3 + x_9 = 4 \\ & x_2 + x_4 + x_{10} = 3 \\ & x_1, \dots, x_{10} \geq 0. \end{array}$$

## 15.7

Let  $x_i \geq 0$ ,  $i = 1, \dots, 4$ , be the weight in pounds of item  $i$  to be used. Then, the total weight is  $x_1 + x_2 + x_3 + x_4$ . To satisfy the percentage content of fiber, fat, and sugar, and the total weight of 1000, we need

$$\begin{array}{rcl} 3x_1 + 8x_2 + 16x_3 + 4x_4 & = & 10(x_1 + x_2 + x_3 + x_4) \\ 6x_1 + 46x_2 + 9x_3 + 9x_4 & = & 2(x_1 + x_2 + x_3 + x_4) \\ 20x_1 + 5x_2 + 4x_3 + 0x_4 & = & 5(x_1 + x_2 + x_3 + x_4) \\ x_1 + x_2 + x_3 + x_4 & = & 1000 \end{array}$$

The total cost is  $2x_1 + 4x_2 + x_3 + 2x_4$ . Therefore, the problem is:

$$\begin{array}{ll} \text{minimize} & 2x_1 + 4x_2 + x_3 + 2x_4 \\ \text{subject to} & -7x_1 - 2x_2 + 6x_3 - 6x_4 = 0 \\ & 4x_1 + 44x_2 + 7x_3 + 7x_4 = 0 \\ & 15x_1 - x_3 - 5x_4 = 0 \\ & x_1 + x_2 + x_3 + x_4 = 1000 \\ & x_1, x_2, x_3, x_4 \geq 0 \end{array}$$

Alternatively, we could have simply replaced  $x_1 + x_2 + x_3 + x_4$  in the first three equality constraints above by 1000, to obtain:

$$\begin{array}{rcl} 3x_1 + 8x_2 + 16x_3 + 4x_4 & = & 10000 \\ 6x_1 + 46x_2 + 9x_3 + 9x_4 & = & 2000 \\ 20x_1 + 5x_2 + 4x_3 + 0x_4 & = & 5000 \\ x_1 + x_2 + x_3 + x_4 & = & 1000. \end{array}$$

Note that the only vector satisfying the above linear equations is  $[179, -175, 573, 422]^\top$ , which is not feasible. Therefore, the constraint does not have any any feasible points, which means that the problem does not have a solution.

### 15.8

The objective function is  $p_1 + \cdots + p_n$ . The constraint for the  $i$ th location is:  $g_{i,1}p_1 + \cdots + g_{i,n}p_n \geq P$ . Hence, the optimization problem is:

$$\begin{aligned} & \text{minimize} && p_1 + \cdots + p_n \\ & \text{subject to} && g_{i,1}p_1 + \cdots + g_{i,n}p_n \geq P, \quad i = 1, \dots, m \\ & && p_1, \dots, p_n \geq 0. \end{aligned}$$

By defining the notation  $\mathbf{G} = [g_{i,j}]$  ( $m \times n$ ),  $\mathbf{e}_n = [1, \dots, 1]^\top$  (with  $n$  components), and  $\mathbf{p} = [p_1, \dots, p_n]^\top$ , we can rewrite the problem as

$$\begin{aligned} & \text{minimize} && \mathbf{e}_n^\top \mathbf{p} \\ & \text{subject to} && \mathbf{G}\mathbf{p} \geq P\mathbf{e}_m \\ & && \mathbf{p} \geq \mathbf{0}. \end{aligned}$$

### 15.9

It is easy to check (using MATLAB, for example) that the matrix

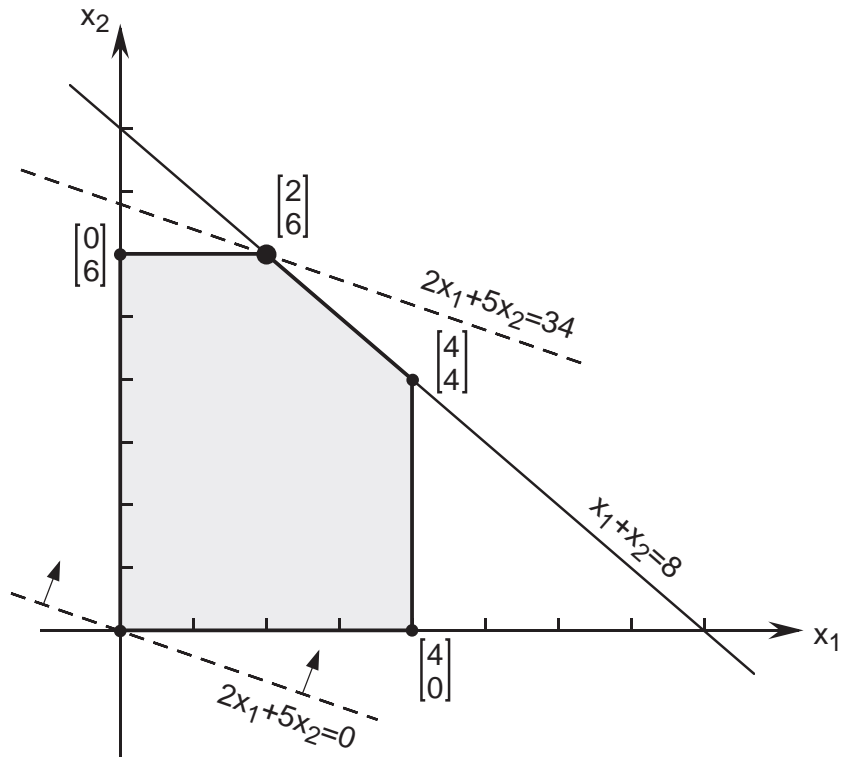
$$\mathbf{A} = \begin{bmatrix} 2 & -1 & 2 & -1 & 3 \\ 1 & 2 & 3 & 1 & 0 \\ 1 & 0 & -2 & 0 & -5 \end{bmatrix}$$

is of full rank (i.e.,  $\text{rank } \mathbf{A} = 3$ ). Therefore, the system has basic solutions. To find the basic solutions, we first select bases. Each basis consists of three linearly independent columns of  $\mathbf{A}$ . These columns correspond to basic variables of the basic solution. The remaining variables are nonbasic and are set to 0. The matrix  $\mathbf{A}$  has 5 columns; therefore, we have  $\binom{5}{3} = 10$  possible candidate basic solutions (corresponding to the 10 combinations of 3 columns out of 5). It turns out that all 10 combinations of 3 columns of  $\mathbf{A}$  are linearly independent. Therefore, we have 10 basic solutions. These are tabulated as follows:

Columns	Basic Solutions
1, 2, 3	$[-4/17, -80/17, 83/17, 0, 0]^\top$
1, 2, 4	$[-10, 49, 0, -83, 0]^\top$
1, 2, 5	$[105/31, 25/31, 0, 0, 83/31]^\top$
1, 3, 4	$[-12/11, 0, 49/11, -80/11, 0]^\top$
1, 3, 5	$[100/35, 0, 25/35, 0, 80/35]^\top$
1, 4, 5	$[65/18, 0, 0, 25/18, 49/18]^\top$
2, 3, 4	$[0, -6, 5, 2, 0]^\top$
2, 3, 5	$[0, -100/23, 105/23, 0, 4/23]^\top$
2, 4, 5	$[0, 13, 0, -21, 2]^\top$
3, 4, 5	$[0, 0, 65/19, -100/19, 12/19]^\top$

### 15.10

In the figure below, the shaded region corresponds to the feasible set. We then translate the line  $2x_1 + 5x_2 = 0$  across the shaded region until the line just touches the region at one point, and the line is as far as possible from the origin. The point of contact is the solution to the problem. In this case, the solution is  $[2, 6]^\top$ , and the corresponding cost is 34.



### 15.11

We use the following MATLAB commands:

```
>> f=[0,-10,0,-6,-20];
>> A=[1,-1,-1,0,0; 0,0,1,-1,-1];
>> b=[0;0];
>> vlb=zeros(5,1);
>> vub=[4;3;3;2;2];
>> x0=zeros(5,1);
>> neqcstr=2;
>> x=linprog(f,A,b,vlb,vub,x0,neqcstr)
```

x =

```
4.0000
2.0000
2.0000
0.0000
2.0000
```

The solution is  $[4, 2, 2, 0, 2]^T$ .

## 16. The Simplex Method

### 16.1

a. Performing a sequence of elementary row operations, we obtain

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & -1 & 3 & 2 \\ 2 & -1 & 3 & 0 & 1 \\ 3 & 1 & 2 & 3 & 3 \\ 1 & 2 & 3 & 1 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & -1 & 3 & 2 \\ 0 & -5 & 5 & -6 & -3 \\ 0 & -5 & 5 & -6 & -3 \\ 0 & 0 & 4 & -2 & -1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & -1 & 3 & 2 \\ 0 & -5 & 5 & -6 & -3 \\ 0 & 0 & 4 & -2 & -1 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} = \mathbf{B}.$$

Because elementary row operations do not change the rank of a matrix,  $\text{rank } \mathbf{A} = \text{rank } \mathbf{B}$ . Therefore  $\text{rank } \mathbf{A} = 3$ .

b. Performing a sequence of elementary row operations, we obtain

$$\begin{aligned} \mathbf{A} = \begin{bmatrix} 1 & \gamma & -1 & 2 \\ 2 & -1 & \gamma & 5 \\ 1 & 10 & -6 & 1 \end{bmatrix} &\rightarrow \begin{bmatrix} 1 & 10 & -6 & 1 \\ 1 & \gamma & -1 & 2 \\ 2 & -1 & \gamma & 5 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 10 & -6 & 1 \\ 0 & \gamma - 10 & 5 & 1 \\ 0 & -21 & \gamma + 12 & 3 \end{bmatrix} \\ &\rightarrow \begin{bmatrix} 1 & 1 & 10 & -6 \\ 0 & 1 & \gamma - 10 & 5 \\ 0 & 3 & -21 & \gamma + 12 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 1 & 10 & -6 \\ 0 & 1 & \gamma - 10 & 5 \\ 0 & 0 & -3(\gamma - 3) & \gamma - 3 \end{bmatrix} = \mathbf{B} \end{aligned}$$

Because elementary row operations do not change the rank of a matrix,  $\text{rank } \mathbf{A} = \text{rank } \mathbf{B}$ . Therefore  $\text{rank } \mathbf{A} = 3$  if  $\gamma \neq 3$  and  $\text{rank } \mathbf{A} = 2$  if  $\gamma = 3$ .

### 16.2

a.

$$\mathbf{A} = \begin{bmatrix} 3 & 1 & 0 & 1 \\ 6 & 2 & 1 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 4 \\ 5 \end{bmatrix}, \quad \mathbf{c} = [2, -1, -1, 0].$$

b. Pivoting the problem tableau about the elements (1, 4) and (2, 3), we obtain

$$\begin{array}{ccccc} 3 & 1 & 0 & 1 & 4 \\ 3 & 1 & 1 & 0 & 1 \\ 5 & 0 & 0 & 0 & 1 \end{array}$$

c. Basic feasible solution:  $\mathbf{x} = [0, 0, 1, 4]^\top$ ,  $\mathbf{c}^\top \mathbf{x} = -1$ .

d.  $[r_1, r_2, r_3, r_4] = [5, 0, 0, 0]$ .

e. Since the reduced cost coefficients are all  $\geq 0$ , the basic feasible solution in part c is optimal.

f. The original problem does indeed have a feasible solution, because the artificial problem has an optimal feasible solution with objective function value 0, as shown in the final phase I tableau.

g. Extract the submatrices corresponding to  $\mathbf{A}$  and  $\mathbf{b}$ , append the last row  $[\mathbf{c}^\top, 0]$ , and pivot about the (2, 1)th element to obtain

$$\begin{array}{ccccc} 0 & 0 & -1 & 1 & 3 \\ 1 & 1/3 & 1/3 & 0 & 1/3 \\ 0 & -5/3 & -5/3 & 0 & -2/3 \end{array}$$

### 16.3

The problem in standard form is:

$$\begin{array}{ll} \text{minimize} & -x_1 - x_2 - 3x_3 \\ \text{subject to} & x_1 + x_3 = 1 \\ & x_2 + x_3 = 2 \\ & x_1, x_2, x_3 \geq 0. \end{array}$$

We form the tableau for the problem:

$$\begin{array}{cccc} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 2 \\ -1 & -1 & -3 & 0 \end{array}$$

Performing necessary row operations, we obtain a tableau in canonical form:

$$\begin{array}{cccc} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 2 \\ 0 & 0 & -1 & 3 \end{array}$$

We pivot about the (1,3)th element to get:

$$\begin{array}{cccc} 1 & 0 & 1 & 1 \\ -1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 4 \end{array}$$

The reduced cost coefficients are all nonnegative. Hence, the current basic feasible solution is optimal:  $[0, 1, 1]^\top$ . The optimal cost is 4.

#### 16.4

The problem in standard form is:

$$\begin{array}{ll} \text{minimize} & -2x_1 - x_2 \\ \text{subject to} & x_1 + x_3 = 5 \\ & x_2 + x_4 = 7 \\ & x_1 + x_2 + x_5 = 9 \\ & x_1, \dots, x_5 \geq 0. \end{array}$$

We form the tableau for the problem:

$$\begin{array}{cccccc} 1 & 0 & 1 & 0 & 0 & 5 \\ 0 & 1 & 0 & 1 & 0 & 7 \\ 1 & 1 & 0 & 0 & 1 & 9 \\ -2 & -1 & 0 & 0 & 0 & 0 \end{array}$$

The above tableau is already in canonical form, and therefore we can proceed with the simplex procedure. We first pivot about the (1,1)th element, to get

$$\begin{array}{cccccc} 1 & 0 & 1 & 0 & 0 & 5 \\ 0 & 1 & 0 & 1 & 0 & 7 \\ 0 & 1 & -1 & 0 & 1 & 4 \\ 0 & -1 & 2 & 0 & 0 & 10 \end{array}$$

Next, we pivot about the (3,2)th element to get

$$\begin{array}{cccccc} 1 & 0 & 1 & 0 & 0 & 5 \\ 0 & 0 & 1 & 1 & -1 & 3 \\ 0 & 1 & -1 & 0 & 1 & 4 \\ 0 & 0 & 1 & 0 & 1 & 14 \end{array}$$

The reduced cost coefficients are all nonnegative. Hence, the optimal solution to the problem in standard form is  $[5, 4, 0, 3, 0]^\top$ . The corresponding optimal cost is  $-14$ .

---

**16.5**

a. Let  $\mathbf{B} = [\mathbf{a}_2, \mathbf{a}_1]$  represent the first two columns of  $\mathbf{A}$  ordered according to the basis corresponding to the given canonical tableau, and  $\mathbf{D}$  the second two columns. Then,

$$\mathbf{B}^{-1}\mathbf{D} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix},$$

Hence,

$$\mathbf{B} = \mathbf{D} \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}^{-1} = \begin{bmatrix} 3/2 & -1/2 \\ -2 & 1 \end{bmatrix}.$$

Hence,

$$\mathbf{A} = \begin{bmatrix} -1/2 & 3/2 & 0 & 1 \\ 1 & -2 & 1 & 0 \end{bmatrix}$$

An alternative approach is to realize that the canonical tableau is obtained from the problem tableau via elementary row operations. Therefore, we can obtain the entries of  $\mathbf{A}$  from the  $2 \times 4$  upper-left submatrix of the canonical tableau via elementary row operations also. Specifically, start with

$$\begin{bmatrix} 0 & 1 & 1 & 2 \\ 1 & 0 & 3 & 4 \end{bmatrix}$$

and then do two pivoting operations, one about  $(1, 4)$  and the other about  $(2, 3)$ .

b. The right-half of  $\mathbf{c}$  is given by

$$\mathbf{c}_D^\top = \mathbf{r}_D^\top + \mathbf{c}_B^\top \mathbf{B}^{-1} \mathbf{D} = [-1, 1] + [7, 8] \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = [-1, 1] + [31, 46] = [30, 47].$$

So  $\mathbf{c} = [8, 7, 30, 47]^\top$ .

c. First we calculate  $\mathbf{B}^{-1}\mathbf{b}$ , giving us the basic variable values:

$$\mathbf{B}^{-1}\mathbf{b} = \begin{bmatrix} 2 & 1 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} 5 \\ 6 \end{bmatrix} = \begin{bmatrix} 16 \\ 38 \end{bmatrix}.$$

Hence, the BFS is  $[38, 16, 0, 0]^\top$ .

d. The first two entries are 16 and 38, respectively. The last component is  $-\mathbf{c}_B^\top \mathbf{B}^{-1}\mathbf{b} = -(7(16) + 8(38)) = -416$ . Hence, the last column is the vector  $[16, 38, -416]^\top$ .

---

**16.6**

The columns in the constraint matrix  $\mathbf{A}$  corresponding to  $x_i^+$  and  $x_i^-$  are linearly dependent. Hence they cannot both enter a basis at the same time. This means that only one variable,  $x_i^+$  or  $x_i^-$ , can assume a nonnegative value; the nonbasic variable is necessarily zero.

---

**16.7**

a. From the given information, we have the  $4 \times 6$  canonical tableau

$$\begin{bmatrix} 1 & * & 0 & 0 & 1/2 & 1 \\ 0 & * & 1 & 0 & 0 & 2 \\ 0 & * & 0 & 1 & 0 & 3 \\ 0 & 1 & 0 & 0 & -1 & -6 \end{bmatrix}$$

Explanations:

- The given vector  $\mathbf{x}$  indicates that  $\mathbf{A}$  is  $3 \times 5$ .



- In the above tableau, we assume that the basis is  $[\mathbf{a}_1, \mathbf{a}_3, \mathbf{a}_4]$ , in this order. Other permutations of orders will result in interchanging rows among the first three rows of the tableau.
- The fifth column represents the coordinates of  $\mathbf{a}_5$  with respect to the basis  $[\mathbf{a}_1, \mathbf{a}_3, \mathbf{a}_4]$ . Because  $[-2, 0, 0, 0, 4]^\top$  lies in the nullspace of  $\mathbf{A}$ , we deduce that  $-2\mathbf{a}_1 + 4\mathbf{a}_5 = \mathbf{0}$ , which can be rewritten as  $\mathbf{a}_5 = (1/2)\mathbf{a}_1 + 0\mathbf{a}_3 + 0\mathbf{a}_4$ , and hence the coordinate vector is  $[1/2, 0, 0]^\top$ .

b. Let  $\mathbf{d}_0 = [-2, 0, 0, 0, 4]^\top$ . Then,  $\mathbf{A}\mathbf{d}_0 = \mathbf{0}$ . Therefore, the vector  $\mathbf{x}' = \mathbf{x} + \varepsilon\mathbf{d}_0$  also satisfies  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . Now,  $\mathbf{x}' = \mathbf{x} + \varepsilon\mathbf{d}_0 = [1 - 2\varepsilon, 0, 2, 3, 4\varepsilon]^\top$ . For  $\mathbf{x}'$  to be feasible, we must have  $\varepsilon \leq 1/2$ . Moreover, the objective function value of  $\mathbf{x}'$  is  $\mathbf{c}^\top \mathbf{x}' = z_0 + r_5 x'_5 = 6 - 4\varepsilon$ , where  $z_0$  is the objective function value of  $\mathbf{x}$ . So, if we pick any  $\varepsilon \in (0, 1/2]$ , then  $\mathbf{x}'$  will be a feasible solution with objective function value strictly less than 6. For example, with  $\varepsilon = 1/2$ ,  $\mathbf{x}' = [0, 0, 2, 3, 2]^\top$  is such a point. (We could also have obtained this solution by pivoting about the element (1, 5) in the tableau of part a.)

### 16.8

- a. The BFS is  $[6, 0, 7, 5, 0]^\top$ , with objective function value  $-8$ .
- b.  $\mathbf{r} = [0, 4, 0, 0, -4]^\top$ .
- c. Yes, because the 5th column has all negative entries.
- d. We pivot about the element (3, 2). The new canonical tableau is:

$$\begin{bmatrix} 0 & 0 & -1/3 & 1 & 0 & 8/3 \\ 1 & 0 & -2/3 & 0 & 0 & 4/3 \\ 0 & 1 & 1/3 & 0 & -1 & 7/3 \\ 0 & 0 & -4/3 & 0 & 0 & -4/3 \end{bmatrix}$$

- e. First note that based on the 5th column, the following point is feasible:

$$\mathbf{x} = \begin{bmatrix} 6 \\ 0 \\ 7 \\ 5 \\ 0 \end{bmatrix} + \epsilon \begin{bmatrix} 2 \\ 0 \\ 3 \\ 1 \\ 1 \end{bmatrix}.$$

Note that  $x_5 = \epsilon$ . Now, any solution of the form  $\mathbf{x} = [*, 0, *, *, \epsilon]^\top$  has an objective function value given by

$$z = z_0 + r_5 \epsilon$$

where  $z_0 = -8$  and  $r_5 = -4$  (from parts a and b). If  $z = -100$ , then  $\epsilon = 23$ . Hence, the following point has objective function value  $z = -100$ :

$$\mathbf{x} = \begin{bmatrix} 6 \\ 0 \\ 7 \\ 5 \\ 0 \end{bmatrix} + 23 \begin{bmatrix} 2 \\ 0 \\ 3 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 52 \\ 0 \\ 76 \\ 28 \\ 23 \end{bmatrix}.$$

- f. The entries of the 2nd column of the given canonical tableau are the coordinates of  $\mathbf{a}_2$  with respect to the basis  $\{\mathbf{a}_4, \mathbf{a}_1, \mathbf{a}_3\}$ . Therefore,

$$\mathbf{a}_2 = \mathbf{a}_4 + 2\mathbf{a}_1 + 3\mathbf{a}_3.$$

Therefore, the vector  $[2, -1, 3, 1, 0]^\top$  lies in the nullspace of  $\mathbf{A}$ . Similarly, using the entries of the 5th column, we deduce that  $[-2, 0, -3, -1, -1]^\top$  also lies in the nullspace of  $\mathbf{A}$ . These two vectors are linearly independent. Because  $\mathbf{A}$  has rank 3, the dimension of the nullspace of  $\mathbf{A}$  is 2. Hence, these two vectors form a basis for the nullspace of  $\mathbf{A}$ .

**16.9**

a. We can convert the problem to standard form by multiplying the objective function by  $-1$  and introducing a surplus variable  $x_3$ . We obtain:

$$\begin{array}{ll}\text{minimize} & x_1 + 2x_2 \\ \text{subject to} & x_2 - x_3 = 1 \\ & x_1, x_2, x_3 \geq 0.\end{array}$$

Note that we do not need to deal with the absence of the constraint  $x_2 \geq 0$  in the original problem, since  $x_2 \geq 1$  implies that  $x_2 \geq 0$  also. Had we used the rule of writing  $x_2 = u - v$  with  $u, v \geq 0$ , we obtain the standard form problem:

$$\begin{array}{ll}\text{minimize} & x_1 + 2u - 2v \\ \text{subject to} & u - v - x_3 = 1 \\ & x_1, u, v, x_3 \geq 0.\end{array}$$

b. For phase I, we set up the artificial problem tableau as:

$$\begin{array}{ccccc} 0 & 1 & -1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 \end{array}$$

Pivoting about element  $(1, 4)$ , we obtain the canonical tableau:

$$\begin{array}{ccccc} 0 & 1 & -1 & 1 & 1 \\ 0 & -1 & 1 & 0 & -1 \end{array}$$

Pivoting now about element  $(1, 2)$ , we obtain the next canonical tableau:

$$\begin{array}{ccccc} 0 & 1 & -1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 \end{array}$$

Hence, phase I terminates, and we use  $x_2$  as our initial basic variable for phase II.

For phase II, we set up the problem tableau as:

$$\begin{array}{ccccc} 0 & 1 & -1 & 1 & \\ 1 & 2 & 0 & 0 & \end{array}$$

Pivoting about element  $(1, 2)$ , we obtain

$$\begin{array}{ccccc} 0 & 1 & -1 & 1 & \\ 1 & 0 & 2 & -2 & \end{array}$$

Hence, the BFS  $[0, 1, 0]^\top$  is optimal, with objective function value 2. Therefore, the optimal solution to the original problem is  $[0, 1]^\top$  with objective function value  $-2$ .

**16.10**

a.  $[1, 0]^\top$

b.  $\begin{bmatrix} 1 & -1 & 1 \end{bmatrix}$

Note that the answer is not  $\begin{bmatrix} 1 & -1 & 1 \\ 0 & -1 & 1 \end{bmatrix}$ , which is the canonical tableau.

c. We choose  $q = 2$  because the only negative RCC value is  $r_2$ . However,  $y_{1,2} < 0$ . Therefore, the simplex algorithm terminates with the condition that the problem is unbounded.

d. Any vector of the form  $[x_1, x_1 - 1]^\top$ ,  $x_1 \geq 1$ , is feasible. Therefore the first component can take arbitrarily large (positive) values. Hence, the objective function, which is  $-x_1$ , can take arbitrarily negative values.

**16.11**

The problem in standard form is:

$$\begin{aligned}
 & \text{minimize} && x_1 + x_2 \\
 & \text{subject to} && x_1 + 2x_2 - x_3 = 3 \\
 & && 2x_1 + x_2 - x_4 = 3 \\
 & && x_1, x_2, x_3, x_4 \geq 0.
 \end{aligned}$$

We will use  $x_1$  and  $x_2$  as initial basic variables. Therefore, Phase I is not needed, and we immediately proceed with Phase II. The tableau for the problem is:

$$\begin{array}{ccccc}
 & \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{b} \\
 & 1 & 2 & -1 & 0 & 3 \\
 & 2 & 1 & 0 & -1 & 3 \\
 \mathbf{c}^\top & 1 & 1 & 0 & 0 & 0
 \end{array}$$

We compute

$$\begin{aligned}
 \boldsymbol{\lambda}^\top &= \mathbf{c}_B^\top \mathbf{B}^{-1} = [1, 1] \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}^{-1} = [1/3, 1/3], \\
 \mathbf{r}_D^\top &= \mathbf{c}_D^\top - \boldsymbol{\lambda}^\top \mathbf{D} = [0, 0] - [1/3, 1/3] \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix} = [1/3, 1/3] = [r_3, r_4].
 \end{aligned}$$

The reduced cost coefficients are all nonnegative. Hence, the solution to the standard form problem is  $[1, 1, 0, 0]^\top$ . Therefore, the solution to the original problem is  $[1, 1]^\top$ , and the corresponding cost is 2.

**16.12**

a. The problem in standard form is:

$$\begin{aligned}
 & \text{minimize} && 4x_1 + 3x_2 \\
 & \text{subject to} && 5x_1 + x_2 - x_3 = 11 \\
 & && 2x_1 + x_2 - x_4 = 8 \\
 & && x_1 + 2x_2 - x_5 = 7 \\
 & && x_1, \dots, x_5 \geq 0.
 \end{aligned}$$

We do not have an apparent basic feasible solution. Therefore, we will need to use the two phase method.

**Phase I:** We introduce artificial variables  $x_6, x_7, x_8$  and form the following tableau.

$$\begin{array}{cccccccccc}
 & \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{a}_5 & \mathbf{a}_6 & \mathbf{a}_7 & \mathbf{a}_8 & \mathbf{b} \\
 & 5 & 1 & -1 & 0 & 0 & 1 & 0 & 0 & 11 \\
 & 2 & 1 & 0 & -1 & 0 & 0 & 1 & 0 & 8 \\
 & 1 & 2 & 0 & 0 & -1 & 0 & 0 & 1 & 7 \\
 \mathbf{c}^\top & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0
 \end{array}$$

We then form the following revised tableau:

$$\begin{array}{c|ccc|c}
 \text{Variable} & & & & \mathbf{y}_0 \\
 \hline
 x_6 & 1 & 0 & 0 & 11 \\
 x_7 & 0 & 1 & 0 & 8 \\
 x_8 & 0 & 0 & 1 & 7
 \end{array}$$

We compute:

$$\begin{aligned}
 \boldsymbol{\lambda}^\top &= [1, 1, 1] \\
 \mathbf{r}_D^\top &= [r_1, r_2, r_3, r_4, r_5] = [-8, -4, 1, 1, 1].
 \end{aligned}$$

We form the augmented revised tableau by introducing  $\mathbf{y}_1 = \mathbf{B}^{-1}\mathbf{a}_1 = \mathbf{a}_1$ :

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$	$\mathbf{y}_1$
$x_6$	1	0	0	11	5
$x_7$	0	1	0	8	2
$x_8$	0	0	1	7	1

We now pivot about the first component of  $\mathbf{y}_1$  to get

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$
$x_1$	1/5	0	0	11/5
$x_7$	-2/5	1	0	18/5
$x_8$	-1/5	0	1	24/5

We compute

$$\begin{aligned}\boldsymbol{\lambda}^\top &= [-3/5, 1, 1] \\ \mathbf{r}_D^\top &= [r_2, r_3, r_4, r_5, r_6] = [-12/5, -3/5, 1, 1, 8/5].\end{aligned}$$

We bring  $\mathbf{y}_2 = \mathbf{B}^{-1}\mathbf{a}_2$  into the basis to get

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$	$\mathbf{y}_2$
$x_1$	1/5	0	0	11/5	1/5
$x_7$	-2/5	1	0	18/5	3/5
$x_8$	-1/5	0	1	24/5	9/5

We pivot about the third component of  $\mathbf{y}_2$  to get

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$
$x_1$	2/9	0	-1/9	5/3
$x_7$	-1/3	1	-1/3	2
$x_2$	-1/9	0	5/9	8/3

We compute

$$\begin{aligned}\boldsymbol{\lambda}^\top &= [-1/3, 1, -1/3] \\ \mathbf{r}_D^\top &= [r_3, r_4, r_5, r_6, r_8] = [-1/3, 1, -1/3, 4/3, 4/3].\end{aligned}$$

We bring  $\mathbf{y}_3 = \mathbf{B}^{-1}\mathbf{a}_3$  into the basis to get

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$	$\mathbf{y}_3$
$x_1$	2/9	0	-1/9	5/3	-2/9
$x_7$	-1/3	1	-1/3	2	1/3
$x_2$	-1/9	0	5/9	8/3	1/9

We pivot about the second component of  $\mathbf{y}_3$  to obtain

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$
$x_1$	0	2/3	-1/3	3
$x_3$	-1	3	-1	6
$x_2$	0	-1/3	2/3	2

We compute

$$\begin{aligned}\boldsymbol{\lambda}^\top &= [0, 0, 0] \\ \mathbf{r}_D^\top &= [r_4, r_5, r_6, r_7, r_8] = [0, 0, 1, 1, 1] \geq \mathbf{0}^\top.\end{aligned}$$

Thus, Phase I is complete, and the initial basic feasible solution is  $[3, 2, 6, 0, 0]^\top$ .

### Phase II

We form the tableau for the original problem:

	$\mathbf{a}_1$	$\mathbf{a}_2$	$\mathbf{a}_3$	$\mathbf{a}_4$	$\mathbf{a}_5$	$\mathbf{b}$
	5	1	-1	0	0	11
	2	1	0	-1	0	8
	1	2	0	0	-1	7
$\mathbf{c}^\top$	4	3	0	0	0	0

The initial revised tableau for Phase II is the final revised tableau for Phase I. We compute

$$\begin{aligned}\boldsymbol{\lambda}^\top &= [0, 5/3, 2/3] \\ \mathbf{r}_D^\top &= [r_4, r_5] = [5/3, 2/3] > \mathbf{0}^\top.\end{aligned}$$

Hence, the optimal solution to the original problem is  $[3, 2]^\top$ .

b. The problem in standard form is:

$$\begin{aligned}\text{minimize} \quad & -6x_1 - 4x_2 - 7x_3 - 5x_4 \\ \text{subject to} \quad & x_1 + 2x_2 + x_3 + 2x_4 + x_5 = 20 \\ & 6x_1 + 5x_2 + 3x_3 + 2x_4 + x_6 = 100 \\ & 3x_1 + 4x_2 + 9x_3 + 12x_4 + x_7 = 75 \\ & x_1, \dots, x_7 \geq 0.\end{aligned}$$

We have an apparent basic feasible solution:  $[0, 0, 0, 20, 100, 75]^\top$ , corresponding to  $\mathbf{B} = \mathbf{I}_3$ . We form the revised tableau corresponding to this basic feasible solution:

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$
$x_5$	1	0	0	20
$x_6$	0	1	0	100
$x_7$	0	0	1	75

We compute

$$\begin{aligned}\boldsymbol{\lambda}^\top &= [0, 0, 0] \\ \mathbf{r}_D^\top &= [r_1, r_2, r_3, r_4] = [-6, -4, -7, -5].\end{aligned}$$

We bring  $\mathbf{y}_2 = \mathbf{B}^{-1}\mathbf{a}_3 = \mathbf{a}_3$  into the basis to obtain

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$	$\mathbf{y}_3$
$x_5$	1	0	0	20	1
$x_6$	0	1	0	100	3
$x_7$	0	0	1	75	9

We pivot about the third component of  $\mathbf{y}_3$  to get

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$
$x_5$	1	0	-1/9	35/3
$x_6$	0	1	-1/3	75
$x_3$	0	0	1/9	25/3

We compute

$$\begin{aligned}\boldsymbol{\lambda}^\top &= [0, 0, -7/9] \\ \mathbf{r}_D^\top &= [r_1, r_2, r_4, r_7] = [-11/3, -8/9, 13/3, 7/9].\end{aligned}$$

We bring  $\mathbf{y}_1 = \mathbf{B}^{-1}\mathbf{a}_1$  into the basis to obtain

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$	$\mathbf{y}_1$
$x_5$	1	0	-1/9	35/3	2/3
$x_6$	0	1	-1/3	75	5
$x_3$	0	0	1/9	25/3	1/3

We pivot about the second component of  $\mathbf{y}_1$  to obtain

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$
$x_5$	1	-2/15	-1/15	5/3
$x_1$	0	1/5	-1/15	15
$x_3$	0	-1/15	2/15	10/3

We compute

$$\begin{aligned}\boldsymbol{\lambda}^\top &= [0, -11/15, -8/15] \\ \mathbf{r}_D^\top &= [r_2, r_4, r_6, r_7] = [27/15, 43/15, 11/15, 8/15] > \mathbf{0}^\top.\end{aligned}$$

The optimal solution to the original problem is therefore  $[15, 0, 10/3, 0]^\top$ .

#### 16.13

a. By inspection of  $\mathbf{r}^\top$ , we conclude that the basic variables are  $x_1, x_3, x_4$ , and the basis matrix is

$$\mathbf{B} = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}.$$

Since  $\mathbf{r}^\top \geq \mathbf{0}^\top$ , the basic feasible solution corresponding to the basis  $\mathbf{B}$  is optimal. This optimal basic feasible solution is  $[8, 0, 9, 7]^\top$ .

b. An optimal solution to the dual is given by

$$\boldsymbol{\lambda}^\top = \mathbf{c}_B^\top \mathbf{B}^{-1},$$

where  $\mathbf{c}_B^\top = [6, 4, 5]$ , and

$$\mathbf{B}^{-1} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}.$$

We obtain  $\boldsymbol{\lambda}^\top = [5, 6, 4]$ .

c. We have  $\mathbf{r}_D^\top = \mathbf{c}_D^\top - \boldsymbol{\lambda}^\top \mathbf{D}$ , where  $\mathbf{r}_D^\top = [1]$ ,  $\mathbf{c}_D^\top = [c_2]$ ,  $\boldsymbol{\lambda}^\top = [5, 6, 4]$ , and  $\mathbf{D} = [2, 1, 3]^\top$ . We get  $1 = c_2 - 10 - 6 - 12$ , which yields  $c_2 = 29$ .

#### 16.14

a. There are two basic feasible solutions:  $[1, 0]^\top$  and  $[0, 2]^\top$ .

b. The feasible set in  $\mathbb{R}^2$  for this problem is the line segment joining the two basic feasible solutions  $[1, 0]^\top$  and  $[0, 2]^\top$ . Therefore, if the problem has an optimal feasible solution that is not basic, then all points in the feasible set are optimal. For this, we need

$$\begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = \alpha \begin{bmatrix} 2 \\ 1 \end{bmatrix},$$

where  $\alpha \in \mathbb{R}$ .

c. Since all basic feasible solutions are optimal, the relative cost coefficients are all zero.

**16.15**

- a.  $2 - \alpha < 0$ ,  $\beta \leq 0$ ,  $\gamma \leq 0$ , and  $\delta$  anything.  
 b.  $2 - \alpha \geq 0$ ,  $\delta = -7$ ,  $\beta$  and  $\gamma$  anything.  
 c.  $2 - \alpha < 0$ ,  $\gamma > 0$ , either  $\beta \leq 0$  or  $5/\gamma \leq 4/\beta$ , and  $\delta$  anything.

**16.16**

a. The value of  $\alpha$  must be 0, because the objective function value is 0 (lower right corner), and  $\alpha$  is the value of an artificial variable.

The value of  $\beta$  must be 0, because it is the RCC value corresponding to a basic column.

The value of  $\gamma$  must be 2, because it must be a positive value. Otherwise, there is a feasible solution to the artificial problem with objective function value smaller than 0, which is impossible.

The value of  $\delta$  must be 0, because we must be able to bring the fourth column into the basis without changing the objective function value.

b. The given linear programming problem does indeed have a feasible solution:  $[0, 5, 6, 0]^\top$ . We obtain this by noticing that the right-most column is a linear combination of the second and third columns, with coefficients 5 and 6.

**16.17**

First, we convert the inequality constraint  $\mathbf{Ax} \geq \mathbf{b}$  into standard form. To do this, we introduce a variable  $\mathbf{w} \in \mathbb{R}^m$  of surplus variables to convert the inequality constraint into the following equivalent constraint:

$$[\mathbf{A}, -\mathbf{I}] \begin{bmatrix} \mathbf{x} \\ \mathbf{w} \end{bmatrix} = \mathbf{b}, \quad \mathbf{w} \geq \mathbf{0}.$$

Next, we introduce variables  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$  to replace the free variable  $\mathbf{x}$  by  $\mathbf{u} - \mathbf{v}$ . We then obtain the following equivalent constraint:

$$[\mathbf{A}, -\mathbf{A}, -\mathbf{I}] \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \\ \mathbf{w} \end{bmatrix} = \mathbf{b}, \quad \mathbf{u}, \mathbf{v}, \mathbf{w} \geq \mathbf{0}.$$

This form of the constraint is now in standard form. So we can now use Phase I of the simplex method to implement an algorithm to find a vectors  $\mathbf{u}$ ,  $\mathbf{v}$ , and  $\mathbf{w}$  satisfying the above constraint, if such exist, or to declare that none exists. If such exist, we output  $\mathbf{x} = \mathbf{u} - \mathbf{v}$ ; otherwise, we declare that no  $\mathbf{x}$  exists such that  $\mathbf{Ax} \geq \mathbf{b}$ . By construction, this algorithm is guaranteed to behave in the way specified by the question.

**16.18**

a. We form the tableau for the problem:

$$\begin{array}{cccccccc} 1 & 0 & 0 & 1/4 & -8 & -1 & 9 & 0 \\ 0 & 1 & 0 & 1/2 & -12 & -1/2 & 3 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & -3/4 & 20 & -1/2 & 6 & 0 \end{array}$$

The above tableau is already in canonical form, and therefore we can proceed with the simplex procedure. We first pivot about the (1, 4)th element, to get

$$\begin{array}{cccccccc} 4 & 0 & 0 & 1 & -32 & -4 & 36 & 0 \\ -2 & 1 & 0 & 0 & 4 & 3/2 & -15 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 3 & 0 & 0 & 0 & -4 & -7/2 & 33 & 0 \end{array}$$

Pivoting about the (2, 5)th element, we get

$$\begin{array}{cccccccc} -12 & 8 & 0 & 1 & 0 & 8 & -84 & 0 \\ -1/2 & 1/4 & 0 & 0 & 1 & 3/8 & -15/4 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 & -2 & 18 & 0 \end{array}$$

Pivoting about the (1,6)th element, we get

$$\begin{array}{cccccccc} -3/2 & 1 & 0 & 1/8 & 0 & 1 & -21/2 & 0 \\ 1/16 & -1/8 & 0 & -3/64 & 1 & 0 & 3/16 & 0 \\ 3/2 & -1 & 1 & -1/8 & 0 & 0 & 21/2 & 1 \\ -2 & 3 & 0 & 1/4 & 0 & 0 & -3 & 0 \end{array}$$

Pivoting about the (2,7)th element, we get

$$\begin{array}{cccccccc} 2 & -6 & 0 & -5/2 & 56 & 1 & 0 & 0 \\ 1/3 & -2/3 & 0 & -1/4 & 16/3 & 0 & 1 & 0 \\ -2 & 6 & 1 & 5/2 & -56 & 0 & 0 & 1 \\ -1 & 1 & 0 & -1/2 & 16 & 0 & 0 & 0 \end{array}$$

Pivoting about the (1,1)th element, we get

$$\begin{array}{cccccccc} 1 & -3 & 0 & -5/4 & 28 & 1/2 & 0 & 0 \\ 0 & 1/3 & 0 & 1/6 & -4 & -1/6 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & -2 & 0 & -7/4 & 44 & 1/2 & 0 & 0 \end{array}$$

Pivoting about the (2,2)th element, we get

$$\begin{array}{cccccccc} 1 & 0 & 0 & 1/4 & -8 & -1 & 9 & 0 \\ 0 & 1 & 0 & 1/2 & -12 & -1/2 & 3 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & -3/4 & 20 & -1/2 & 6 & 0 \end{array}$$

which is identical to the initial tableau. Therefore, cycling occurs.

b. We start with the initial tableau of part a, and pivot about the (1,4)th element to obtain

$$\begin{array}{cccccccc} 4 & 0 & 0 & 1 & -32 & -4 & 36 & 0 \\ -2 & 1 & 0 & 0 & 4 & 3/2 & -15 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 3 & 0 & 0 & 0 & -4 & -7/2 & 33 & 0 \end{array}$$

Pivoting about the (2,5)th element, we get

$$\begin{array}{cccccccc} -12 & 8 & 0 & 1 & 0 & 8 & -84 & 0 \\ -1/2 & 1/4 & 0 & 0 & 1 & 3/8 & -15/4 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 & -2 & 18 & 0 \end{array}$$

Pivoting about the (1,6)th element, we get

$$\begin{array}{cccccccc} -3/2 & 1 & 0 & 1/8 & 0 & 1 & -21/2 & 0 \\ 1/16 & -1/8 & 0 & -3/64 & 1 & 0 & 3/16 & 0 \\ 3/2 & -1 & 1 & -1/8 & 0 & 0 & 21/2 & 1 \\ -2 & 3 & 0 & 1/4 & 0 & 0 & -3 & 0 \end{array}$$

Pivoting about the (2,1)th element, we get

$$\begin{array}{cccccccc} 0 & -2 & 0 & -1 & 24 & 1 & -6 & 0 \\ 1 & -2 & 0 & -3/4 & 16 & 0 & 3 & 0 \\ 0 & 2 & 1 & 1 & -24 & 0 & 6 & 1 \\ 0 & -1 & 0 & -5/4 & 32 & 0 & 3 & 0 \end{array}$$



Pivoting about the (3,2)th element, we get

$$\begin{array}{cccccccc} 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1/4 & -8 & 0 & 9 & 1 \\ 0 & 1 & 1/2 & 1/2 & -12 & 0 & 3 & 1/2 \\ 0 & 0 & 1/2 & -3/4 & 20 & 0 & 6 & 1/2 \end{array}$$

Pivoting about the (3,4)th element, we get

$$\begin{array}{cccccccc} 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 1 & -1/2 & 3/4 & 0 & -2 & 0 & 15/2 & 3/4 \\ 0 & 2 & 1 & 1 & -24 & 0 & 6 & 1 \\ 0 & 3/2 & 5/4 & 0 & 2 & 0 & 21/2 & 5/4 \end{array}$$

The reduced cost coefficients are all nonnegative. Hence, the optimal solution to the problem is  $[3/4, 0, 0, 1, 0, 1, 0]^\top$ . The corresponding optimal cost is  $-5/4$ .

---

#### 16.19

a. We have

$$\mathbf{A}\mathbf{d}^{(0)} = \mathbf{A}(\mathbf{x}^{(1)} - \mathbf{x}^{(0)})/\alpha_0 = (\mathbf{b} - \mathbf{b})/\alpha_0 = 0.$$

Hence,  $\mathbf{d}^{(0)} \in \mathcal{N}(\mathbf{A})$ .

b. From our discussion of moving from one BFS to an adjacent BFS, we deduce that

$$\mathbf{d}^{(0)} = \begin{bmatrix} -\mathbf{y}_q \\ \mathbf{e}_{q-m} \end{bmatrix}.$$

In other words, the first  $m$  components of  $\mathbf{d}^{(0)}$  are  $-y_{1q}, \dots, -y_{mq}$ , and all the other components are 0 except the  $q$ th component, which is 1.

---

#### 16.20

The following is a MATLAB function that implements the simplex algorithm.

```
function [x,v]=simplex(c,A,b,v,options)
%           SIMPLEX(c,A,b,v);
%           SIMPLEX(c,A,b,v,options);
%
%           x = SIMPLEX(c,A,b,v);
%           x = SIMPLEX(c,A,b,v,options);
%
%           [x,v] = SIMPLEX(c,A,b,v);
%           [x,v] = SIMPLEX(c,A,b,v,options);
%
%SIMPLEX(c,A,b,v) solves the following linear program using the
%Simplex Method:
%  min c'x subject to Ax=b, x>=0,
%where [A b] is in canonical form, and v is the vector of indices of
%basic columns. Specifically, the v(i)-th column of A is the i-th
%standard basis vector.
%The second variant allows a vector of optional parameters to be
%defined:
%OPTIONS(1) controls how much display output is given; set
%to 1 for a tabular display of results (default is no display: 0).
%OPTIONS(5) specifies how the pivot element is selected;
%  0=choose the most negative relative cost coefficient;
%  1=use Bland's rule.
```

```

if nargin ~= 5
    options = [];
    if nargin ~= 4
        disp('Wrong number of arguments. ');
        return;
    end
end

format compact;
%format short e;

options = foptions(options);
print = options(1);

n=length(c);
m=length(b);

cB=c(v(:));

r = c'-cB'*A; %row vector of relative cost coefficients

cost = -cB'*b;

tabl=[A b;r cost];

if print,
    disp(' ');
    disp('Initial tableau:');
    disp(tabl);
end %if

while ones(1,n)*(r' >= zeros(n,1)) ~= n
    if options(5) == 0;
        [r_q,q] = min(r);
    else
        %Bland's rule
        q=1;
        while r(q) >= 0
            q=q+1;
        end
    end %if

    min_ratio = inf;
    p=0;
    for i=1:m,
        if tabl(i,q)>0
            if tabl(i,n+1)/tabl(i,q) < min_ratio
                min_ratio = tabl(i,n+1)/tabl(i,q);
                p = i;
            end %if
        end %if
    end %for
    if p == 0
        disp('Problem unbounded');
        break;
    end %if

    tabl=pivot(tabl,p,q);

```

```

    if print,
        disp('Pivot point:');
        disp([p,q]);
        disp('New tableau:');
        disp(tabl);
    end %if

    v(p) = q;
    r = tabl(m+1,1:n);
end %while

x=zeros(n,1);
x(v(:))=tabl(1:m,n+1);

```

The above function makes use of the following function that implements pivoting:

```

function Mnew=pivot(M,p,q)
%Mnew=pivot(M,p,q)
%Returns the matrix Mnew resulting from pivoting about the
%(p,q)th element of the given matrix M.

for i=1:size(M,1),
    if i==p
        Mnew(p,:)=M(p,:)/M(p,q);
    else
        Mnew(i,:)=M(i,:)-M(p,:)*(M(i,q)/M(p,q));
    end %if
end %for
%-----

```

We now apply the simplex algorithm to the problem in Example 16.2, as follows:

```

>> A=[1 0 1 0 0; 0 1 0 1 0; 1 1 0 0 1];
>> b=[4;6;8];
>> c=[-2;-5;0;0;0];
>> v=[3;4;5];
>> options(1)=1;
>> [x,v]=simplex(c,A,b,v,options);

```

Initial Tableau:

1	0	1	0	0	4
0	1	0	1	0	6
1	1	0	0	1	8
-2	-5	0	0	0	0

Pivot point:

2	2
---	---

New tableau:

1	0	1	0	0	4
0	1	0	1	0	6
1	0	0	-1	1	2
-2	0	0	5	0	30

Pivot point:

3	1
---	---

New tableau:

0	0	1	1	-1	2
0	1	0	1	0	6
1	0	0	-1	1	2
0	0	0	3	2	34

```
>> disp(x');

```

2	6	2	0	0
---	---	---	---	---

```
>> disp(v');
      3      2      1
```

As indicated above, the solution to the problem in standard form is  $[2, 6, 2, 0, 0]^T$ , and the objective function value is  $-34$ . The optimal cost for the original maximization problem is  $34$ .

### 16.21

The following is a MATLAB routine that implements the two-phase simplex method, using the MATLAB function from Exercise 16.20.

```
function [x,v]=tpsimplx(c,A,b,options)
%           TPSIMPLEX(c,A,b);
%           TPSIMPLEX(c,A,b,options);
%
%           x = TPSIMPLEX(c,A,b);
%           x = TPSIMPLEX(c,A,b,options);
%
%           [x,v] = TPSIMPLEX(c,A,b);
%           [x,v] = TPSIMPLEX(c,A,b,options);
%
%TPSIMPLEX(c,A,b) solves the following linear program using the
%two-phase simplex method:
% min c'x subject to Ax=b, x>=0.
%The second variant allows a vector of optional parameters to be
%defined:
%OPTIONS(1) controls how much display output is given; set
%to 1 for a tabular display of results (default is no display: 0).
%OPTIONS(5) specifies how the pivot element is selected;
% 0=choose the most negative relative cost coefficient;
% 1=use Bland's rule.

if nargin ~= 4
    options = [];
    if nargin ~= 3
        disp('Wrong number of arguments. ');
        return;
    end
end

clc;
format compact;
%format short e;

options = foptions(options);
print = options(1);

n=length(c);
m=length(b);

%Phase I
if print,
    disp(' ');
    disp('Phase I');
    disp('-----');
end

v=n*ones(m,1);
for i=1:m
    v(i)=v(i)+i;
```

```

end

[x,v]=simplex([zeros(n,1);ones(m,1)],[A eye(m)],b,v,options);

if all(v<=n),
    %Phase II
    if print
        disp(' ');
        disp('Phase II');
        disp('-----');
        disp('Basic columns:')
        disp(v')
    end

    %Convert [A b] into canonical augmented matrix
    Binv=inv(A(:,[v]));
    A=Binv*A;
    b=Binv*b;

    [x,v]=simplex(c,A,b,v,options);

    if print
        disp(' ');
        disp('Final solution:');
        disp(x');
    end
else
    %assumes nondegeneracy
    disp('Terminating: problem has no feasible solution.');
```

```

end
%-----
```

We now apply the above MATLAB routine to the problem in Example 16.5, as follows:

```

>> A=[1 1 0; 5 3 0 -1];
>> b=[4;8];
>> c=[-3;-5;0;0];
>> options(1)=1;
>> format rat;
>> tpsimplex(c,A,b,options);
```

```

Phase I
-----
```

Initial Tableau:

1	1	1	0	1	0	4
5	3	0	-1	0	1	8
-6	-4	-1	1	0	0	-12

Pivot point:

2 1

New tableau:

0	2/5	1	1/5	1	-1/5	12/5
1	3/5	0	-1/5	0	1/5	8/5
0	-2/5	-1	-1/5	0	6/5	-12/5

Pivot point:

1 3

New tableau:

0	2/5	1	1/5	1	-1/5	12/5
1	3/5	0	-1/5	0	1/5	8/5
0	*	0	*	1	1	*

```

Pivot point:
    2    2
New tableau:
   -2/3    0    1    1/3    1   -1/3    4/3
    5/3    1    0   -1/3    0    1/3    8/3
    *     0    0     *     1     1     *
Pivot point:
    1    4
New tableau:
   -2     0    3     1     3    -1     4
    1     1    1     0     1     0     4
    *     0    *     0     1     1     *

Basic columns:
    4     2

```

Phase II

-----

```

Initial Tableau:
   -2     0    3     1     4
    1     1    1     0     4
    2     0    5     0    20

```

```

Final solution:
    0     4     0     4

```

## 16.22

The following is a MATLAB function that implements the revised simplex algorithm.

```

function [x,v,Binv]=revsimp(c,A,b,v,Binv,options)
%           REVSIMP(c,A,b,v,Binv);
%           REVSIMP(c,A,b,v,Binv,options);
%
%           x = REVSIMP(c,A,b,v,Binv);
%           x = REVSIMP(c,A,b,v,Binv,options);
%
%           [x,v,Binv] = REVSIMP(c,A,b,v,Binv);
%           [x,v,Binv] = REVSIMP(c,A,b,v,Binv,options);
%
%REVSIMP(c,A,b,v,Binv) solves the following linear program using the
%revised simplex method:
%  min c'x  subject to Ax=b, x>=0,
%where v is the vector of indices of basic columns, and Binv is the
%inverse of the basis matrix. Specifically, the v(i)-th column of
%A is the i-th column of the basis vector.
%The second variant allows a vector of optional parameters to be
%defined:
%OPTIONS(1) controls how much display output is given; set
%to 1 for a tabular display of results (default is no display: 0).
%OPTIONS(5) specifies how the pivot element is selected;
%  0=choose the most negative relative cost coefficient;
%  1=use Bland's rule.

if nargin ~= 6
    options = [];
    if nargin ~= 5
        disp('Wrong number of arguments. ');
        return;
    end
end

```

```

end

format compact;
%format short e;

options = foptions(options);
print = options(1);

n=length(c);
m=length(b);

cB=c(v(:));
y0 = Binv*b;
lambdaT=cB'*Binv;
r = c'-lambdaT*A; %row vector of relative cost coefficients

if print,
    disp(' ');
    disp('Initial revised tableau [v B(-1) y0]:');
    disp([v Binv y0]);
    disp('Relative cost coefficients:');
    disp(r);
end %if

while ones(1,n)*(r' >= zeros(n,1)) ~= n
    if options(5) == 0;
        [r_q,q] = min(r);
    else
        %Bland's rule
        q=1;
        while r(q) >= 0
            q=q+1;
        end
    end %if

    yq = Binv*A(:,q);
    min_ratio = inf;
    p=0;
    for i=1:m,
        if yq(i)>0
            if y0(i)/yq(i) < min_ratio
                min_ratio = y0(i)/yq(i);
                p = i;
            end %if
        end %if
    end %for
    if p == 0
        disp('Problem unbounded');
        break;
    end %if

    if print,
        disp('Augmented revised tableau [v B(-1) y0 yq]:')
        disp([v Binv y0 yq]);
        disp('(p,q):');
        disp([p,q]);
    end

    augrevtabl=pivot([Binv y0 yq],p,m+2);

```

```

Binv=augrevtab1(:,1:m);
y0=augrevtab1(:,m+1);

v(p) = q;

cB=c(v(:));
lambdaT=cB'*Binv;
r = c'-lambdaT*A; %row vector of relative cost coefficients

if print,
    disp('New revised tableau [v B^(-1) y0]:');
    disp([v Binv y0]);
    disp('Relative cost coefficients:');
    disp(r);
end %if

end %while

x=zeros(n,1);
x(v(:))=y0;

```

The function makes use of the pivoting function in Exercise 16.20.

We now apply the simplex algorithm to the problem in Example 16.2, as follows:

```

>> A=[1 0 1 0 0; 0 1 0 1 0; 1 1 0 0 1];
>> b=[4;6;8];
>> c=[-2;-5;0;0;0];
>> v=[3;4;5];
>> Binv=eye(3);
>> options(1)=1;
>> [x,v,Binv]=rev_simp(c,A,b,v,Binv,options);

```

Initial revised tableau [v B<sup>(-1)</sup> y0]:

3	1	0	0	4
4	0	1	0	6
5	0	0	1	8

Relative cost coefficients:

-2	-5	0	0	0
----	----	---	---	---

Augmented revised tableau [v B<sup>(-1)</sup> y0 yq]:

3	1	0	0	4	0
4	0	1	0	6	1
5	0	0	1	8	1

(p,q):

2	2
---	---

New revised tableau [v B<sup>(-1)</sup> y0]:

3	1	0	0	4
2	0	1	0	6
5	0	-1	1	2

Relative cost coefficients:

-2	0	0	5	0
----	---	---	---	---

Augmented revised tableau [v B<sup>(-1)</sup> y0 yq]:

3	1	0	0	4	1
2	0	1	0	6	0
5	0	-1	1	2	1

(p,q):

3	1
---	---

New revised tableau [v B<sup>(-1)</sup> y0]:

3	1	1	-1	2
2	0	1	0	6
1	0	-1	1	2



```

Relative cost coefficients:
      0      0      0      3      2
>> disp(x');
      2      6      2      0      0
>> disp(v');
      3      2      1
>> disp(Binv);
      1      1     -1
      0      1      0
      0     -1      1

```

### 16.23

The following is a MATLAB routine that implements the two-phase revised simplex method, using the MATLAB function from Exercise 16.22.

```

function [x,v]=tprevsimp(c,A,b,options)
%           TPREVSIMP(c,A,b);
%           TPREVSIMP(c,A,b,options);
%
%           x = TPREVSIMP(c,A,b);
%           x = TPREVSIMP(c,A,b,options);
%
%           [x,v] = TPREVSIMP(c,A,b);
%           [x,v] = TPREVSIMP(c,A,b,options);
%
%TPREVSIMP(c,A,b) solves the following linear program using the
%two-phase revised simplex method:
% min c'x subject to Ax=b, x>=0.
%The second variant allows a vector of optional parameters to be
%defined:
%OPTIONS(1) controls how much display output is given; set
%to 1 for a tabular display of results (default is no display: 0).
%OPTIONS(5) specifies how the pivot element is selected;
% 0=choose the most negative relative cost coefficient;
% 1=use Bland's rule.

if nargin ~= 4
    options = [];
    if nargin ~= 3
        disp('Wrong number of arguments. ');
        return;
    end
end

clc;
format compact;
%format short e;

options = foptions(options);
print = options(1);

n=length(c);
m=length(b);

%Phase I
if print,
    disp(' ');
    disp('Phase I');
    disp('-----');

```

```

end

v=n*ones(m,1);
for i=1:m
    v(i)=v(i)+i;
end

[x,v,Binv]=rev_simp([zeros(n,1);ones(m,1)], [A eye(m)],b,v,eye(m),options);

%Phase II
if print
    disp(' ');
    disp('Phase II');
    disp('-----');
end

[x,v,Binv]=rev_simp(c,A,b,v,Binv,options);

if print
    disp(' ');
    disp('Final solution:');
    disp(x');
end
%-----

```

We now apply the above MATLAB routine to the problem in Example 16.5, as follows:

```

>> A=[4 2 -1 0; 1 4 0 -1];
>> b=[12;6];
>> c=[2;3;0;0];
>> options(1)=1;
>> format rat;
>> tprevsimp(c,A,b,options);

```

Phase I

-----

Initial revised tableau [v B<sup>(-1)</sup> y<sub>0</sub>]:

5	1	0	12
6	0	1	6

Relative cost coefficients:

-5	-6	1	1	0	0
----	----	---	---	---	---

Augmented revised tableau [v B<sup>(-1)</sup> y<sub>0</sub> y<sub>q</sub>]:

5	1	0	12	2
6	0	1	6	4

(p,q):

2	2
---	---

New revised tableau [v B<sup>(-1)</sup> y<sub>0</sub>]:

5	1	-1/2	9
2	0	1/4	3/2

Relative cost coefficients:

-7/2	0	1	-1/2	0	3/2
------	---	---	------	---	-----

Augmented revised tableau [v B<sup>(-1)</sup> y<sub>0</sub> y<sub>q</sub>]:

5	1	-1/2	9	7/2
2	0	1/4	3/2	1/4

(p,q):

1	1
---	---

New revised tableau [v B<sup>(-1)</sup> y<sub>0</sub>]:

1	2/7	-1/7	18/7
2	-1/14	2/7	6/7

Relative cost coefficients:

0    0    0    0    1    1

Phase II

-----

Initial revised tableau [v B<sup>-1</sup> y0]:

1    2/7   -1/7   18/7

2   -1/14   2/7   6/7

Relative cost coefficients:

\*    0    5/14   4/7

Final solution:

18/7   6/7   0    0

## 17. Duality

### 17.1

Since  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  are feasible, we have  $\mathbf{Ax} \geq \mathbf{b}$ ,  $\mathbf{x} \geq \mathbf{0}$ , and  $\boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top$ ,  $\boldsymbol{\lambda} \geq \mathbf{0}$ . Postmultiplying both sides of  $\boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top$  by  $\mathbf{x} \geq \mathbf{0}$  yields

$$\boldsymbol{\lambda}^\top \mathbf{Ax} \leq \mathbf{c}^\top \mathbf{x}.$$

Since  $\mathbf{Ax} \geq \mathbf{b}$  and  $\boldsymbol{\lambda}^\top \geq \mathbf{0}^\top$ , we have  $\boldsymbol{\lambda}^\top \mathbf{Ax} \geq \boldsymbol{\lambda}^\top \mathbf{b}$ . Hence,  $\boldsymbol{\lambda}^\top \mathbf{b} \leq \mathbf{c}^\top \mathbf{x}$ .

### 17.2

The primal problem is:

$$\begin{aligned} & \text{minimize} && \mathbf{e}_n^\top \mathbf{p} \\ & \text{subject to} && \mathbf{G}\mathbf{p} \geq \mathbf{P}\mathbf{e}_m \\ & && \mathbf{p} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{G} = [g_{i,j}]$ ,  $\mathbf{e}_n = [1, \dots, 1]^\top$  (with  $n$  components), and  $\mathbf{p} = [p_1, \dots, p_n]^\top$ . The dual of the problem is (using symmetric duality):

$$\begin{aligned} & \text{maximize} && \mathbf{P}\boldsymbol{\lambda}^\top \mathbf{e}_m \\ & \text{subject to} && \boldsymbol{\lambda}^\top \mathbf{G} \leq \mathbf{e}_n^\top \\ & && \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

### 17.3

a. We first transform the problem into standard form:

$$\begin{aligned} & \text{minimize} && -2x_1 - 3x_2 \\ & \text{subject to} && x_1 + 2x_2 + x_3 = 4 \\ & && 2x_1 + x_2 + x_4 = 5 \\ & && x_1, x_2, x_3, x_4 \geq 0. \end{aligned}$$

The initial tableau is:

$$\begin{array}{cccccc} 1 & 2 & 1 & 0 & 4 \\ 2 & 1 & 0 & 1 & 5 \\ -2 & -3 & 0 & 0 & 0 \end{array}$$

We now pivot about the (1,2)th element to get:

$$\begin{array}{cccccc} 1/2 & 1 & 1/2 & 0 & 2 \\ 3/2 & 0 & -1/2 & 1 & 3 \\ -1/2 & 0 & 3/2 & 0 & 6 \end{array}$$

Pivoting now about the (2, 1)th element gives:

$$\begin{array}{ccccc} 0 & 1 & 2/3 & -1/3 & 1 \\ 1 & 0 & -1/3 & 2/3 & 2 \\ 0 & 0 & 4/3 & 1/3 & 7 \end{array}$$

Thus, the solution to the standard form problem is  $x_1 = 2$ ,  $x_2 = 1$ ,  $x_3 = 0$ ,  $x_4 = 0$ . The solution to the original problem is  $x_1 = 2$ ,  $x_2 = 1$ .

b. The dual to the standard form problem is

$$\begin{array}{ll} \text{maximize} & 4\lambda_1 + 5\lambda_2 \\ \text{subject to} & \lambda_1 + 2\lambda_2 \leq -2 \\ & 2\lambda_1 + \lambda_2 \leq -3 \\ & \lambda_1, \lambda_2 \geq 0. \end{array}$$

From the discussion before Example 17.6, it follows that the solution to the dual is  $\lambda^\top = \mathbf{c}_I^\top - \mathbf{r}_I^\top = [-4/3, -1/3]$ .

#### 17.4

The dual problem is

$$\begin{array}{ll} \text{maximize} & 11\lambda_1 + 8\lambda_2 + 7\lambda_3 \\ \text{subject to} & 5\lambda_1 + 2\lambda_2 + \lambda_3 \leq 4 \\ & \lambda_1 + \lambda_2 + 2\lambda_3 \leq 3 \\ & \lambda_1, \lambda_2, \lambda_3 \geq 0. \end{array}$$

Note that we may arrive at the above in one of two ways: by applying the asymmetric form of duality, or by applying the symmetric form of duality to the original problem in standard form. From the solution of Exercise 16.11a, we have that the solution to the dual is  $\lambda^\top = \mathbf{c}_B^\top \mathbf{B}^{-1} = [0, 5/3, 2/3]$  (using the proof of the duality theorem).

#### 17.5

We represent the primal in the form

$$\begin{array}{ll} \text{minimize} & \mathbf{c}^\top \mathbf{x} \\ \text{subject to} & \mathbf{A}\mathbf{x} = \mathbf{b} \\ & \mathbf{x} \geq 0. \end{array}$$

The corresponding dual is

$$\begin{array}{ll} \text{maximize} & \lambda^\top \mathbf{b} \\ \text{subject to} & \lambda^\top \mathbf{A} \leq \mathbf{c}^\top, \end{array}$$

that is,

$$\begin{array}{ll} \text{maximize} & 2\lambda_1 + 7\lambda_2 + 3\lambda_3 \\ \text{subject to} & \begin{bmatrix} \lambda_1 & \lambda_2 & \lambda_3 \end{bmatrix} \begin{bmatrix} -2 & 1 & 1 & 0 & 0 \\ -1 & 2 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 \end{bmatrix} \leq \begin{bmatrix} -1 & -2 & 0 & 0 & 0 \end{bmatrix}. \end{array}$$

The solution to the dual can be obtained using the formula,  $\lambda^{*\top} = \mathbf{c}_B^\top \mathbf{B}^{-1}$ , where

$$\mathbf{c}_B^\top = \begin{bmatrix} -1 & -2 & 0 \end{bmatrix} \quad \text{and} \quad \mathbf{B} = \begin{bmatrix} -2 & 1 & 1 \\ -1 & 2 & 0 \\ 1 & 0 & 0 \end{bmatrix}.$$

Note that because the last element in  $\mathbf{c}_B^\top$  is zero, we do not need to calculate the last row of  $\mathbf{B}^{-1}$  when computing  $\boldsymbol{\lambda}^{*\top}$ , that is, these elements are “don’t care” elements that we denote using the asterisk. Hence,

$$\boldsymbol{\lambda}^{*\top} = \mathbf{c}_B^\top \mathbf{B}^{-1} = \begin{bmatrix} -1 & -2 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1/2 & 1/2 \\ * & * & * \end{bmatrix} = \begin{bmatrix} 0 & -1 & -2 \end{bmatrix}.$$

Note that

$$\mathbf{c}^\top \mathbf{x}^* = \boldsymbol{\lambda}^{*\top} \mathbf{b} = -13,$$

as expected.

### 17.6

a. Multiplying the objective function by  $-1$ , we see that the problem is of the form of the dual in the asymmetric form of duality. Therefore, the dual to the problem is of the form of the primal in the asymmetric form:

$$\begin{array}{ll} \text{minimize} & \boldsymbol{\lambda}^\top \mathbf{b} \\ \text{subject to} & \boldsymbol{\lambda}^\top \mathbf{A} = -\mathbf{c}^\top \\ & \boldsymbol{\lambda} \geq \mathbf{0} \end{array}$$

b. The given vector  $\mathbf{y}$  is feasible in the dual. Since  $\mathbf{b} = \mathbf{0}$ , any feasible point in the dual is optimal. Thus,  $\mathbf{y}$  is optimal in the dual, and the objective function value for  $\mathbf{y}$  is 0. Therefore, by the Duality Theorem, the primal also has an optimal feasible solution, and the corresponding objective function value is 0. Since the vector  $\mathbf{0}$  is feasible in the primal and has objective function value 0, the vector  $\mathbf{0}$  is a solution to the primal.

### 17.7

We introduce two sets of nonnegative variables:  $x_i^+ \geq 0$ ,  $x_i^- \geq 0$ ,  $i = 1, 2, \dots, 3$ . We can then represent the optimization problem in the form

$$\begin{array}{ll} \text{minimize} & (x_1^+ + x_1^-) + (x_2^+ + x_2^-) + (x_3^+ + x_3^-) \\ \text{subject to} & \begin{bmatrix} 1 & 1 & -1 & -1 & -1 & 1 \\ 0 & -1 & 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1^+ \\ x_2^+ \\ x_3^+ \\ x_1^- \\ x_2^- \\ x_3^- \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \\ & x_i^+ \geq 0, x_i^- \geq 0. \end{array}$$

We form the initial tableau,

$$\begin{array}{cccccc|c} x_1^+ & x_2^+ & x_3^+ & x_1^- & x_2^- & x_3^- & \mathbf{b} \\ 1 & 1 & -1 & -1 & -1 & 1 & 2 \\ 0 & -1 & 0 & 0 & 1 & 0 & 1 \\ \mathbf{c}^\top & 1 & 1 & 1 & 1 & 1 & 0 \end{array}$$

There is no apparent basic feasible solution. We add the second row to the first one to obtain,

$$\begin{array}{cccccc|c} x_1^+ & x_2^+ & x_3^+ & x_1^- & x_2^- & x_3^- & \mathbf{b} \\ 1 & 0 & -1 & -1 & 0 & 1 & 3 \\ 0 & -1 & 0 & 0 & 1 & 0 & 1 \\ \mathbf{c}^\top & 1 & 1 & 1 & 1 & 1 & 0 \end{array}$$

We next calculate the reduced cost coefficients,

$$\begin{array}{cccccc|c} x_1^+ & x_2^+ & x_3^+ & x_1^- & x_2^- & x_3^- & \mathbf{b} \\ 1 & 0 & -1 & -1 & 0 & 1 & 3 \\ 0 & -1 & 0 & 0 & 1 & 0 & 1 \\ \mathbf{c}^\top & 0 & 2 & 2 & 2 & 0 & -4 \end{array}$$

We have zeros under the basic columns. The reduced cost coefficients are all nonnegative. The optimal solution is,

$$\mathbf{x}^* = \begin{bmatrix} 3 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}^\top.$$

The optimal solution to the original problem is  $\mathbf{x}^* = [3, -1, 0]^\top$ .

The dual of the above linear program is

$$\begin{array}{ll} \text{maximize} & 2\lambda_1 + \lambda_2 \\ \text{subject to} & \begin{bmatrix} \lambda_1 & \lambda_2 \end{bmatrix} \begin{bmatrix} 1 & 1 & -1 & -1 & -1 & 1 \\ 0 & -1 & 0 & 0 & 1 & 0 \end{bmatrix} \leq \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix}. \end{array}$$

The optimal solution to the dual is

$$\begin{aligned} \boldsymbol{\lambda}^{*\top} &= \mathbf{c}_B^\top \mathbf{B}^{-1} \\ &= \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix}^{-1} \\ &= \begin{bmatrix} 1 & 2 \end{bmatrix}. \end{aligned}$$

---

## 17.8

a. The dual (asymmetric form) is

$$\begin{array}{ll} \text{maximize} & \lambda \\ \text{subject to} & \lambda a_i \leq 1, \quad i = 1, \dots, n. \end{array}$$

We can write the constraint as

$$\lambda \leq \min\{1/a_i : i = 1, \dots, n\} = 1/a_n.$$

Therefore, the solution to the dual problem is

$$\lambda = 1/a_n.$$

b. Duality Theorem: If the primal problem has an optimal solution, then so does the dual, and the optimal values of their respective objective functions are equal.

By the duality theorem, the primal has an optimal solution, and the optimal value of the objective function is  $1/a_n$ . The only feasible point in the primal with this objective function value is the basic feasible solution  $[0, \dots, 0, 1/a_n]^\top$ .

c. Suppose we start at a nonoptimal initial basic feasible solution,  $[0, \dots, 1/a_i, \dots, 0]^\top$ , where  $1 \leq i \leq n-1$ . The relative cost coefficient for the  $q$ th column,  $q \neq i$ , is

$$r_q = 1 - \frac{a_q}{a_i}.$$

Since  $a_n > a_j$  for any  $j \neq n$ ,  $r_q$  is the most negative relative cost coefficient if and only if  $q = n$ .

---

## 17.9

a. By asymmetric duality, the dual is given by

$$\begin{array}{ll} \text{minimize} & \lambda \\ \text{subject to} & \lambda \geq c_i, \quad i = 1, \dots, n. \end{array}$$

b. The constraint in part a implies that  $\lambda$  is feasible if and only if  $\lambda \geq c_4$ . Hence, the solution is  $\lambda^* = c_4$ .

c. By the duality theorem, the optimal objective function value for the given problem is  $c_4$ . The only solution that achieves this value is  $x_4^* = 1$  and  $x_i^* = 0$  for all  $i \neq 4$ .

---

**17.10**

a. The dual is

$$\begin{array}{ll} \text{minimize} & \boldsymbol{\lambda}^\top \mathbf{0} \\ \text{subject to} & \boldsymbol{\lambda}^\top \mathbf{A} \geq \mathbf{c}^\top \\ & \boldsymbol{\lambda} \geq \mathbf{0}. \end{array}$$

b. By the duality theorem, we conclude that the optimal value of the objective function is 0. The only vector satisfying  $\mathbf{x} \geq \mathbf{0}$  that has an objective function value of 0 is  $\mathbf{x} = \mathbf{0}$ . Therefore, the solution is  $\mathbf{x} = \mathbf{0}$ .

c. The constraint set contains only the vector  $\mathbf{0}$ . Any other vector  $\mathbf{x}$  satisfying  $\mathbf{x} \geq \mathbf{0}$  has at least one positive component, and consequently has a positive objective function value. But this contradicts the fact that the optimal solution has an objective function value of 0.

---

**17.11**

a. The artificial problem is:

$$\begin{array}{ll} \text{minimize} & [\mathbf{0}^\top, \mathbf{e}^\top] \mathbf{z} \\ \text{subject to} & [\mathbf{A}, \mathbf{I}] \mathbf{z} = \mathbf{b} \\ & \mathbf{z} \geq \mathbf{0}, \end{array}$$

where  $\mathbf{e} = [1, \dots, 1]^\top$  and  $\mathbf{z} = [\mathbf{x}^\top, \mathbf{y}^\top]^\top$ .

b. The dual to the artificial problem is:

$$\begin{array}{ll} \text{maximize} & \boldsymbol{\lambda}^\top \mathbf{b} \\ \text{subject to} & \boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{0}^\top \\ & \boldsymbol{\lambda}^\top \leq \mathbf{e}^\top. \end{array}$$

c. Suppose the given original linear programming problem has a feasible solution. By the FTLT, the original LP problem has a BFS. Then, by a theorem given in class, the artificial problem has an optimal feasible solution with  $\mathbf{y} = \mathbf{0}$ . Hence, by the Duality Theorem, the dual of the artificial problem also has an optimal feasible solution.

---

**17.12**

a. Possible. This situation arises if the primal is unbounded, which by the Weak Duality Lemma implies that the dual has no feasible solution.

b. Impossible, because the Duality Theorem requires that if the primal has an optimal feasible solution, then so does the dual.

c. Impossible, because the Duality Theorem requires that if the dual has an optimal feasible solution, then so does the primal. Also, the Weak Dual Lemma requires that if the primal is unbounded (i.e., has a feasible solution but no optimal feasible solution), then the dual must have no feasible solution.

---

**17.13**

To prove the result, we use Theorem 17.3 (Complementary Slackness). Since  $\boldsymbol{\mu} \geq \mathbf{0}$ , we have  $\mathbf{A}^\top \boldsymbol{\lambda} = \mathbf{c} - \boldsymbol{\mu} \leq \mathbf{c}$ . Hence,  $\boldsymbol{\lambda}$  is a feasible solution to the dual. Now,  $(\mathbf{c} - \mathbf{A}^\top \boldsymbol{\lambda})^\top \mathbf{x} = \boldsymbol{\mu}^\top \mathbf{x} = 0$ . Therefore, by Theorem 17.3,  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  are optimal for their respective problems.

---

**17.14**

To use the symmetric form of duality, we need to rewrite the problem as

$$\begin{array}{ll} \text{minimize} & -\mathbf{c}^\top (\mathbf{u} - \mathbf{v}), \\ \text{subject to} & -\mathbf{A}(\mathbf{u} - \mathbf{v}) \geq -\mathbf{b} \\ & \mathbf{u}, \mathbf{v} \geq \mathbf{0}, \end{array}$$

which we represent in the form

$$\begin{aligned} & \text{minimize} && [-\mathbf{c}^\top \ \mathbf{c}^\top] \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}, \\ & \text{subject to} && [-\mathbf{A} \ \mathbf{A}] \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} \geq -\mathbf{b} \\ & && \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} \geq \mathbf{0}. \end{aligned}$$

By the symmetric form of duality, the dual is:

$$\begin{aligned} & \text{maximize} && \boldsymbol{\lambda}^\top (-\mathbf{b}) \\ & \text{subject to} && \boldsymbol{\lambda}^\top [-\mathbf{A} \ \mathbf{A}] \leq [-\mathbf{c}^\top \ \mathbf{c}^\top] \\ & && \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

Note that for the constraint involving  $\mathbf{A}$ , we have

$$\begin{aligned} \boldsymbol{\lambda}^\top [-\mathbf{A} \ \mathbf{A}] \leq [-\mathbf{c}^\top \ \mathbf{c}^\top] & \Leftrightarrow -\boldsymbol{\lambda}^\top \mathbf{A} \leq -\mathbf{c}^\top \text{ and } \boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top \\ & \Leftrightarrow \boldsymbol{\lambda}^\top \mathbf{A} = \mathbf{c}^\top. \end{aligned}$$

Therefore, we can represent the dual as

$$\begin{aligned} & \text{minimize} && \boldsymbol{\lambda}^\top \mathbf{b} \\ & \text{subject to} && \boldsymbol{\lambda}^\top \mathbf{A} = \mathbf{c}^\top \\ & && \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

---

### 17.15

The corresponding dual can be written as:

$$\begin{aligned} & \text{maximize} && 3\lambda_1 + 3\lambda_2 \\ & \text{subject to} && \lambda_1 + 2\lambda_2 \leq 1 \\ & && 2\lambda_1 + \lambda_2 \leq 1 \\ & && \lambda_1, \lambda_2 \geq 0. \end{aligned}$$

To solve the dual, we refer back to the solution of Exercise 16.11. Using the idea of the proof of the duality theorem (Theorem 17.2), we obtain the solution to the dual as  $\boldsymbol{\lambda}^\top = \mathbf{c}_B^\top \mathbf{B}^{-1} = [1/3, 1/3]$ . The cost of the dual problem is 2, which verifies the duality theorem.

---

### 17.16

The dual to the above linear program (asymmetric form) is

$$\begin{aligned} & \text{maximize} && \boldsymbol{\lambda}^\top \mathbf{0} \\ & \text{subject to} && \boldsymbol{\lambda}^\top \mathbf{O} \leq \mathbf{c}^\top. \end{aligned}$$

The above dual problem has a feasible solution if and only if  $\mathbf{c} \geq \mathbf{0}$ . Since any feasible solution to the dual is also optimal, the dual has an optimal solution if and only if  $\mathbf{c} \geq \mathbf{0}$ . Therefore, by the duality theorem, the primal problem has a solution if and only if  $\mathbf{c} \geq \mathbf{0}$ .

If the solution to the dual exists, then the optimal value of the objective function in the primal is equal to that of the dual, which is clearly 0. In this case,  $\mathbf{0}$  is optimal, since  $\mathbf{c}^\top \mathbf{0} = 0$ .

---

### 17.17

Consider the primal problem

$$\begin{aligned} & \text{minimize} && \mathbf{0}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \geq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$



and its corresponding dual

$$\begin{array}{ll}\text{maximize} & \mathbf{y}^\top \mathbf{b} \\ \text{subject to} & \mathbf{y}^\top \mathbf{A} \leq \mathbf{0} \\ & \mathbf{y} \geq \mathbf{0}.\end{array}$$

$\Rightarrow$ : By assumption, there exists a feasible solution to the primal problem. Note that any feasible solution is also optimal, and has objective function value 0. Suppose  $\mathbf{y}$  satisfies  $\mathbf{A}^\top \mathbf{y} \leq \mathbf{0}$  and  $\mathbf{y} \geq \mathbf{0}$ . Then,  $\mathbf{y}$  is a feasible solution to the dual. Therefore, by the Weak Duality Lemma,  $\mathbf{b}^\top \mathbf{y} \leq 0$ .

$\Leftarrow$ : Note that the feasible region for the dual is nonempty, since  $\mathbf{0}$  is a feasible point. Also, by assumption,  $\mathbf{0}$  is an optimal solution, since any other feasible point  $\mathbf{y}$  satisfies  $\mathbf{b}^\top \mathbf{y} \leq \mathbf{b}^\top \mathbf{0} = 0$ . Hence, by the duality theorem, the primal problem has an (optimal) feasible solution.

#### 17.18

---

a. The dual is

$$\begin{array}{ll}\text{maximize} & \mathbf{y}^\top \mathbf{b} \\ \text{subject to} & \mathbf{y}^\top \mathbf{A} \leq \mathbf{0}^\top.\end{array}$$

b. The feasible set of the dual problem is always nonempty, because  $\mathbf{0}$  is clearly guaranteed to be feasible.

c. Suppose  $\mathbf{y}$  is feasible in the dual. Then, by assumption,  $\mathbf{b}^\top \mathbf{y} \leq 0$ . But the point  $\mathbf{0}$  is feasible and has objective function value 0. Hence,  $\mathbf{0}$  is optimal in the dual.

d. By parts b and c, the dual has an optimal feasible solution. Hence, by the duality theorem, the primal problem also has an optimal feasible solution.

e. By assumption, there exists a feasible solution to the primal problem. Note that any feasible solution in the primal has objective function value 0 (and hence so does the given solution). Suppose  $\mathbf{y}$  satisfies  $\mathbf{A}^\top \mathbf{y} \leq \mathbf{0}$ . Then,  $\mathbf{y}$  is a feasible solution to the dual. Therefore, by weak duality,  $\mathbf{b}^\top \mathbf{y} \leq 0$ .

#### 17.19

---

Consider the primal problem

$$\begin{array}{ll}\text{minimize} & \mathbf{y}^\top \mathbf{b} \\ \text{subject to} & \mathbf{y}^\top \mathbf{A} = \mathbf{0} \\ & \mathbf{y} \geq \mathbf{0}\end{array}$$

and its corresponding dual

$$\begin{array}{ll}\text{maximize} & \mathbf{0}^\top \mathbf{x} \\ \text{subject to} & \mathbf{A}\mathbf{x} \leq \mathbf{b}.\end{array}$$

$\Rightarrow$ : By assumption, there exists a feasible solution to the dual problem. Note that any feasible solution is also optimal, and has objective function value 0. Suppose  $\mathbf{y}$  satisfies  $\mathbf{A}^\top \mathbf{y} = \mathbf{0}$  and  $\mathbf{y} \geq \mathbf{0}$ . Then,  $\mathbf{y}$  is a feasible solution to the primal. Therefore, by the Weak Duality Lemma,  $\mathbf{b}^\top \mathbf{y} \geq 0$ .

$\Leftarrow$ : Note that the feasible region for the primal is nonempty, since  $\mathbf{0}$  is a feasible point. Also, by assumption,  $\mathbf{0}$  is an optimal solution, since any other feasible point  $\mathbf{y}$  satisfies  $\mathbf{b}^\top \mathbf{y} \geq \mathbf{b}^\top \mathbf{0} = 0$ . Hence, by the duality theorem, the dual problem has an (optimal) feasible solution.

#### 17.20

---

Let  $\mathbf{e} = [1, \dots, 1]^\top$ . Consider the primal problem

$$\begin{array}{ll}\text{minimize} & \mathbf{0}^\top \mathbf{x} \\ \text{subject to} & \mathbf{A}\mathbf{x} \leq -\mathbf{e}\end{array}$$

and its corresponding dual

$$\begin{array}{ll}\text{maximize} & \mathbf{e}^\top \mathbf{y} \\ \text{subject to} & \mathbf{y}^\top \mathbf{A} = \mathbf{0} \\ & \mathbf{y} \geq \mathbf{0}.\end{array}$$

$\Rightarrow$ : Suppose there exists  $\mathbf{Ax} < \mathbf{0}$ . Then, the vector  $\mathbf{x}' = \mathbf{x} / \min\{|\mathbf{Ax}_i|\}$  is a feasible solution to the primal problem. Note that any feasible solution is also optimal, and has objective function value 0. Suppose  $\mathbf{y}$  satisfies  $\mathbf{A}^\top \mathbf{y} = \mathbf{0}$ ,  $\mathbf{y} \geq \mathbf{0}$ . Then,  $\mathbf{y}$  is a feasible solution to the dual. Therefore, by the Weak Duality Lemma,  $\mathbf{e}^\top \mathbf{y} \leq 0$ . Since  $\mathbf{y} \geq \mathbf{0}$ , we conclude that  $\mathbf{y} = \mathbf{0}$ .

$\Leftarrow$ : Suppose  $\mathbf{0}$  is the only feasible solution to the dual. Then,  $\mathbf{0}$  is clearly also optimal. Hence, by the duality theorem, the primal problem has an (optimal) feasible solution  $\mathbf{x}$ . Since  $\mathbf{Ax} \leq -\mathbf{e}$  and  $-\mathbf{e} < \mathbf{0}$ , we get  $\mathbf{Ax} < \mathbf{0}$ .

### 17.21

a. Rewrite the primal as

$$\begin{aligned} & \text{minimize} && (-\mathbf{e})^\top \mathbf{x} \\ & \text{subject to} && (\mathbf{P} - \mathbf{I})^\top \mathbf{x} = \mathbf{0} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

By asymmetric duality, the dual is

$$\begin{aligned} & \text{maximize} && \boldsymbol{\lambda}^\top \mathbf{0} \\ & \text{subject to} && \boldsymbol{\lambda}^\top (\mathbf{P} - \mathbf{I})^\top \leq -\mathbf{e}^\top. \end{aligned}$$

b. To make the notation simpler, we rewrite the dual as:

$$\begin{aligned} & \text{maximize} && 0 \\ & \text{subject to} && (\mathbf{P} - \mathbf{I})\mathbf{y} \geq \mathbf{e}. \end{aligned}$$

Suppose the dual is feasible. Then, there exists a  $\mathbf{y}$  such that  $\mathbf{Py} \geq \mathbf{y} + \mathbf{e} > \mathbf{y}$ . Let  $y_i$  be the largest element of  $\mathbf{y}$ , and  $\mathbf{p}^{(i)\top}$  the  $i$ th row of  $\mathbf{P}$ . Then,  $\mathbf{p}^{(i)\top} \mathbf{y} > y_i$ . But, by definition of  $y_i$ ,  $\mathbf{y} \leq y_i \mathbf{e}$ . Hence,  $\mathbf{p}^{(i)\top} \mathbf{y} \leq y_i \mathbf{p}^{(i)\top} \mathbf{e} = y_i$ , which contradicts the inequality  $\mathbf{p}^{(i)\top} \mathbf{y} > y_i$ . Hence, the dual is not feasible.

c. The primal is certainly feasible, because  $\mathbf{0}$  is a feasible point. Therefore, by part b and strong duality, the primal must also be unbounded.

d. Because 0 is an achievable objective function value (it is the objective function value of  $\mathbf{0}$ ), and the problem is unbounded, we deduce that 1 is also achievable. Hence, there exists a feasible  $\mathbf{x}$  such that  $\mathbf{x}^\top \mathbf{e} = 1$ . This proves the desired result.

### 17.22

Write the LP problem

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{Ax} \geq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0} \end{aligned}$$

and the corresponding dual problem

$$\begin{aligned} & \text{maximize} && \boldsymbol{\lambda}^\top \mathbf{b} \\ & \text{subject to} && \boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top \\ & && \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

By a theorem on duality, if we can find feasible points  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  for the primal and dual, respectively, such that  $\mathbf{c}^\top \mathbf{x} = \boldsymbol{\lambda}^\top \mathbf{b}$ , then  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  are optimal for their respective problems. We can rewrite the previous set of relations as

$$\begin{bmatrix} -\mathbf{c}^\top & \mathbf{b}^\top \\ \mathbf{c}^\top & -\mathbf{b}^\top \\ \mathbf{A} & \mathbf{0} \\ \mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & -\mathbf{A}^\top \\ \mathbf{0} & \mathbf{I}_m \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\lambda} \end{bmatrix} \geq \begin{bmatrix} 0 \\ 0 \\ \mathbf{b} \\ \mathbf{0} \\ -\mathbf{c} \\ \mathbf{0} \end{bmatrix}.$$

Therefore, writing the above as  $\hat{\mathbf{A}}\mathbf{y} \geq \hat{\mathbf{b}}$ , where  $\hat{\mathbf{A}} \in \mathbb{R}^{(2m+2n+2) \times (m+n)}$  and  $\hat{\mathbf{b}} \in \mathbb{R}^{(2m+2n+2)}$ , we have that the first  $n$  components of  $\phi((2m+2n+2), (m+n), \hat{\mathbf{A}}, \hat{\mathbf{b}})$  is a solution to the given linear programming problem.

### 17.23

a. Consider the dual;  $\mathbf{b}$  does not appear in the constraint (but it does appear in the dual objective function). Thus, provided the level sets of the dual objective function do not exactly align with one of the faces of the constraint set (polyhedron), the optimal dual vector will not change if we perturb  $\mathbf{b}$  very slightly. Now, by the duality theorem,  $z(\mathbf{b}) = \boldsymbol{\lambda}^\top \mathbf{b}$ . Because  $\boldsymbol{\lambda}$  is constant in a neighborhood of  $\mathbf{b}$ , we deduce that  $\nabla z(\mathbf{b}) = \boldsymbol{\lambda}$ .

b. By part a, we deduce that the optimal objective function value will change by  $3\Delta b_1$ .

### 17.24

a. Weak duality lemma: if  $\mathbf{x}_0$  and  $\mathbf{y}_0$  are feasible points in the primal and dual, respectively, then  $f_1(\mathbf{x}_0) \geq f_2(\mathbf{y}_0)$ .

Proof: Because  $\mathbf{y}_0 \geq \mathbf{0}$  and  $\mathbf{A}\mathbf{x}_0 - \mathbf{b} \leq \mathbf{0}$ , we have  $\mathbf{y}_0^\top (\mathbf{A}\mathbf{x}_0 - \mathbf{b}) \leq 0$ . Therefore,

$$\begin{aligned} f_1(\mathbf{x}_0) &\geq f_1(\mathbf{x}_0) + \mathbf{y}_0^\top (\mathbf{A}\mathbf{x}_0 - \mathbf{b}) \\ &= \frac{1}{2} \mathbf{x}_0^\top \mathbf{x}_0 + \mathbf{y}_0^\top \mathbf{A}\mathbf{x}_0 - \mathbf{y}_0^\top \mathbf{b}. \end{aligned}$$

Now, we know that

$$\frac{1}{2} \mathbf{x}_0^\top \mathbf{x}_0 + \mathbf{y}_0^\top \mathbf{A}\mathbf{x}_0 \geq \frac{1}{2} \mathbf{x}^{*\top} \mathbf{x}^* + \mathbf{y}_0^\top \mathbf{A}\mathbf{x}^*,$$

where  $\mathbf{x}^* = -\mathbf{A}^\top \mathbf{y}_0$ . Hence,

$$\frac{1}{2} \mathbf{x}_0^\top \mathbf{x}_0 + \mathbf{y}_0^\top \mathbf{A}\mathbf{x}_0 \geq \frac{1}{2} \mathbf{y}_0^\top \mathbf{A}\mathbf{A}^\top \mathbf{y}_0 - \mathbf{y}_0^\top \mathbf{A}\mathbf{A}^\top \mathbf{y}_0 = -\frac{1}{2} \mathbf{y}_0^\top \mathbf{A}\mathbf{A}^\top \mathbf{y}_0.$$

Combining this with the above, we have

$$\begin{aligned} f_1(\mathbf{x}_0) &\geq -\frac{1}{2} \mathbf{y}_0^\top \mathbf{A}\mathbf{A}^\top \mathbf{y}_0 - \mathbf{y}_0^\top \mathbf{b} \\ &= f_2(\mathbf{y}_0). \end{aligned}$$

Alternatively, notice that

$$\begin{aligned} f_1(\mathbf{x}_0) - f_2(\mathbf{y}_0) &= \frac{1}{2} \mathbf{x}_0^\top \mathbf{x}_0 + \frac{1}{2} \mathbf{y}_0^\top \mathbf{A}\mathbf{A}^\top \mathbf{y}_0 + \mathbf{b}^\top \mathbf{y}_0 \\ &\geq \frac{1}{2} \mathbf{x}_0^\top \mathbf{x}_0 + \frac{1}{2} \mathbf{y}_0^\top \mathbf{A}\mathbf{A}^\top \mathbf{y}_0 + \mathbf{x}_0^\top \mathbf{A}^\top \mathbf{y}_0 \\ &= \frac{1}{2} \|\mathbf{x}_0 + \mathbf{A}^\top \mathbf{y}_0\|^2 \\ &\geq 0. \end{aligned}$$

b. Suppose  $f_1(\mathbf{x}_0) = f_2(\mathbf{y}_0)$  for feasible points  $\mathbf{x}_0$  and  $\mathbf{y}_0$ . Let  $\mathbf{x}$  be any feasible point in the primal. Then, by part a,  $f_1(\mathbf{x}) \geq f_2(\mathbf{y}_0) = f_1(\mathbf{x}_0)$ . Hence  $\mathbf{x}_0$  is optimal in the primal.

Similarly, let  $\mathbf{y}$  be any feasible point in the dual. Then, by part a,  $f_2(\mathbf{y}) \leq f_1(\mathbf{x}_0) = f_2(\mathbf{y}_0)$ . Hence  $\mathbf{y}_0$  is optimal in the dual.

## 18. Non-Simplex Methods

### 18.1

The following is a MATLAB function that implements the affine scaling algorithm.

```
function [x,N] = affscale(c,A,b,u,options);
%           AFFSCALE(c,A,b,u);
%           AFFSCALE(c,A,b,u,options);
```

```

%
%           x = AFFSCALE(c,A,b,u);
%           x = AFFSCALE(c,A,b,u,options);
%
%           [x,N] = AFFSCALE(c,A,b,u);
%           [x,N] = AFFSCALE(c,A,b,u,options);
%
%AFFSCALE(c,A,b,u) solves the following linear program using the
%affine scaling Method:
%   min c'x  subject to Ax=b, x>=0,
%where u is a strictly feasible initial solution.
%The second variant allows a vector of optional parameters to be
%defined:
%OPTIONS(1) controls how much display output is given; set
%to 1 for a tabular display of results (default is no display: 0).
%OPTIONS(2) is a measure of the precision required for the final point.
%OPTIONS(3) is a measure of the precision required cost value.
%OPTIONS(14) = max number of iterations.
%OPTIONS(18) = alpha.

if nargin ~= 5
    options = [];
    if nargin ~= 4
        disp('Wrong number of arguments. ');
        return;
    end
end
xnew=u;

if length(options) >= 14
    if options(14)==0
        options(14)=1000*length(xnew);
    end
else
    options(14)=1000*length(xnew);
end

%if length(options) < 18
    options(18)=0.99; %optional step size
%end

%clc;
format compact;
format short e;

options = foptions(options);
print = options(1);
epsilon_x = options(2);
epsilon_f = options(3);
max_iter=options(14);
alpha=options(18);

n=length(c);
m=length(b);

for k = 1:max_iter,

    xcurr=xnew;
    D = diag(xcurr);

```

```

Abar = A*D;
Pbar = eye(n) - Abar'*inv(Abar*Abar')*Abar;
d = -D*Pbar*D*c;
if d ~= zeros(n,1),
    nonzd = find(d<0);
    r = min(-xcurr(nonzd)./d(nonzd));
else
    disp('Terminating: d = 0');
    break;
end
xnew = xcurr+alpha*r*d;

if print,
    disp('Iteration number k =')
    disp(k); %print iteration index k
    disp('alpha_k =');
    disp(alpha*r); %print alpha_k
    disp('New point =');
    disp(xnew'); %print new point
end %if

if norm(xnew-xcurr) <= epsilon_x*norm(xcurr)
    disp('Terminating: Relative difference between iterates <');
    disp(epsilon_x);
    break;
end %if

if abs(c'*(xnew-xcurr)) < epsilon_f*abs(c'*xcurr),
    disp('Terminating: Relative change in objective function < ');
    disp(epsilon_f);
    break;
end %if

if k == max_iter
    disp('Terminating with maximum number of iterations');
end %if
end %for

if nargout >= 1
    x=xnew;
    if nargout == 2
        N=k;
    end
else
    disp('Final point =');
    disp(xnew');
    disp('Number of iterations =');
    disp(k);
end %if
%-----

```

We now apply the affine scaling algorithm to the problem in Example 16.2, as follows:

```

>> A=[1 0 1 0 0; 0 1 0 1 0; 1 1 0 0 1];
>> b=[4;6;8];
>> c=[-2;-5;0;0;0];
>> u=[2;3;2;3;3];
>> options(1)=0;
>> options(2)=10^(-7);
>> options(3)=10^(-7);

```

```
>> affscale(c,A,b,u,options);
Terminating: Relative difference between iterates <
1.0000e-07
Final point =
2.0000e+00 6.0000e+00 2.0000e+00 1.0837e-09 1.7257e-08
Number of iterations =
8
```

The result obtained after 8 iterations as indicated above agrees with the solution in Example 16.2:  $[2, 6, 2, 0, 0]^T$ .

## 18.2

The following is a MATLAB routine that implements the two-phase affine scaling method, using the MATLAB function from Exercise 18.1.

```
function [x,N]=tpaffscale(c,A,b,options)
% March 28, 2000
%
%           TPAFFSCALE(c,A,b);
%           TPAFFSCALE(c,A,b,options);
%
%           x = TPAFFSCALE(c,A,b);
%           x = TPAFFSCALE(c,A,b,options);
%
%           [x,N] = TPAFFSCALE(c,A,b);
%           [x,N] = TPAFFSCALE(c,A,b,options);
%
%TPAFFSCALE(c,A,b) solves the following linear program using the
%Two-Phase Affine Scaling Method:
% min c'x subject to Ax=b, x>=0.
%The second variant allows a vector of optional parameters to be
%defined:
%OPTIONS(1) controls how much display output is given; set
%to 1 for a tabular display of results (default is no display: 0).
%OPTIONS(2) is a measure of the precision required for the final point.
%OPTIONS(3) is a measure of the precision required cost value.
%OPTIONS(14) = max number of iterations.
%OPTIONS(18) = alpha.

if nargin ~= 4
    options = [];
    if nargin ~= 3
        disp('Wrong number of arguments. ');
        return;
    end
end

%clc;
format compact;
format short e;

options = foptions(options);
print = options(1);

n=length(c);
m=length(b);

%Phase I
if print,
    disp(' ');
```

```

    disp('Phase I');
    disp('-----');
end

u = rand(n,1);
v = b-A*u;
if v ~= zeros(m,1),
    u = affscale([zeros(1,n),1]',[A v],b,[u' 1]',options);
    u(n+1) = [];
end

if print
    disp(' ')
    disp('Initial condition for Phase II:')
    disp(u)
end

if u(n+1) < options(2),
    %Phase II
    u(n+1) = [];
    if print
        disp(' ');
        disp('Phase II');
        disp('-----');
        disp('Initial condition for Phase II:');
        disp(u);
    end
    [x,N]=affscale(c,A,b,u,options);
    if nargout == 0
        disp('Final point =');
        disp(x');
        disp('Number of iterations =');
        disp(N);
    end %if
else
    disp('Terminating: problem has no feasible solution.');
```

end

%-----

We now apply the above MATLAB routine to the problem in Example 16.5, as follows:

```

>> A=[1 1 1 0; 5 3 0 -1];
>> b=[4;8];
>> c=[-3;-5;0;0];
>> options(1)=0;
>> tpaffscale(c,A,b,options);
Terminating: Relative difference between iterates <
1.0000e-07
Terminating: Relative difference between iterates <
1.0000e-07
Final point =
4.0934e-09 4.0000e+00 9.4280e-09 4.0000e+00
Number of iterations =
7
```

The result obtained above agrees with the solution in Example 16.5:  $[0, 4, 0, 4]^T$ .

### 18.3

The following is a MATLAB routine that implements the affine scaling method applied to LP problems of the form given in the question by converting the given problem in Karmarkar's artificial form and then using the MATLAB function from Exercise 18.1.

```

function [x,N]=karaffscale(c,A,b,options)
%
%           KARAFFSCALE(c,A,b);
%           KARAFFSCALE(c,A,b,options);
%
%           x = KARAFFSCALE(c,A,b);
%           x = KARAFFSCALE(c,A,b,options);
%
%           [x,N] = KARAFFSCALE(c,A,b);
%           [x,N] = KARAFFSCALE(c,A,b,options);
%
%KARAFFSCALE(c,A,b) solves the following linear program using the
%Affine Scaling Method:
% min c'x subject to Ax>=b, x>=0.
%We use Karmarkar's artificial problem to convert the above problem into
%a form usable by the affine scaling method.
%The second variant allows a vector of optional parameters to be
%defined:
%OPTIONS(1) controls how much display output is given; set
%to 1 for a tabular display of results (default is no display: 0).
%OPTIONS(2) is a measure of the precision required for the final point.
%OPTIONS(3) is a measure of the precision required cost value.
%OPTIONS(14) = max number of iterations.
%OPTIONS(18) = alpha.

if nargin ~= 4
    options = [];
    if nargin ~= 3
        disp('Wrong number of arguments.');
```

return;

end

end

%clc;

format compact;

format short e;

options = foptions(options);

print = options(1);

n=length(c);

m=length(b);

%Convert to Karmarkar's artificial problem

x0 = ones(n,1);

l0 = ones(m,1);

u0 = ones(n,1);

v0 = ones(m,1);

AA = [

c' -b' zeros(1,n) zeros(1,m) (-c'\*x0+b'\*l0);

A zeros(m,m) zeros(m,n) -eye(m) (b-A\*x0+v0);

zeros(n,n) A' eye(n) zeros(n,m) (c-A'\*l0)

];

bb = [0; b; c];

cc = [zeros(2\*m+2\*n,1); 1];

y0 = [x0; l0; u0; v0; 1];

[y,N]=affscale(cc,AA,bb,y0,options);



```

if cc'*y <= options(3),
    x = y(1:n);
    if nargout == 0
        disp('Final point =');
        disp(x');
        disp('Final cost =');
        disp(c'*x);
        disp('Number of iterations =');
        disp(N);
    end %if
else
    disp('Terminating: problem has no optimal feasible solution.');
```

We now apply the above MATLAB routine to the problem in Example 15.15, as follows:

```

>> c=[-3;-5];
>> A=[-1 -5; -2 -1; -1 -1];
>> b=[-40;-20;-12];
>> options(2)=10^(-4);
>> karaffscale(c,A,b,options);
Terminating: Relative difference between iterates <
    1.0000e-04
Final point =
    5.1992e+00    6.5959e+00
Final cost =
   -4.8577e+01
Number of iterations =
    3
```

The solution from Example 15.15 is  $[5, 7]^\top$ . The accuracy of the result obtained above is disappointing. We believe that the inaccuracy here may be caused by our particularly simple numerical implementation of the affine scaling method. This illustrates the numerical issues that must be dealt with in any practically useful implementation of the affine scaling method.

#### 18.4

a. Suppose  $\mathbf{T}(\mathbf{x}) = \mathbf{T}(\mathbf{y})$ . Then,  $T_i(\mathbf{x}) = T_i(\mathbf{y})$  for  $i = 1, \dots, n+1$ . Note that for  $i = 1, \dots, n$ ,  $T_i(\mathbf{x}) = (x_i/a_i)T_{n+1}(\mathbf{x})$  and  $T_i(\mathbf{y}) = (y_i/a_i)T_{n+1}(\mathbf{y})$ . Therefore,

$$T_i(\mathbf{x}) = (x_i/a_i)T_{n+1}(\mathbf{x}) = T_i(\mathbf{y}) = (y_i/a_i)T_{n+1}(\mathbf{y}) = (y_i/a_i)T_{n+1}(\mathbf{x}),$$

which implies that  $x_i = y_i$ ,  $i = 1, \dots, n$ . Hence  $\mathbf{x} = \mathbf{y}$ .

b. Let  $\mathbf{y} \in \{\mathbf{x} \in \Delta : x_{n+1} > 0\}$ . Hence  $y_{n+1} > 0$ . Define  $\mathbf{x} = [x_1, \dots, x_n]^\top$  by  $x_i = a_i y_i / y_{n+1}$ ,  $i = 1, \dots, n$ .

Then,  $\mathbf{T}(\mathbf{x}) = \mathbf{y}$ . To see this, note that

$$T_{n+1}(\mathbf{x}) = \frac{1}{y_1/y_{n+1} + \dots + y_n/y_{n+1} + 1} = \frac{y_{n+1}}{y_1 + \dots + y_n + y_{n+1}} = y_{n+1}.$$

Also, for  $i = 1, \dots, n$ ,

$$T_i(\mathbf{x}) = (y_i/y_{n+1})T_{n+1}(\mathbf{x}) = y_i.$$

c. An immediate consequence of the solution to part b.

d. We have

$$T_{n+1}(\mathbf{a}) = \frac{1}{a_1/a_1 + \dots + a_n/a_n + 1} = \frac{1}{n+1},$$

and, for  $i = 1, \dots, n$ ,

$$T_i(\mathbf{a}) = (a_i/a_i)T_{n+1}(\mathbf{a}) = \frac{1}{n+1}.$$

e. Since  $\mathbf{y} = \mathbf{T}(\mathbf{x})$ , we have that for  $i = 1, \dots, n$ ,  $y_i = (x_i/a_i)y_{n+1}$ . Therefore,  $x'_i = y_i a_i = x_i y_{n+1}$ , which implies that  $\mathbf{x}' = y_{n+1} \mathbf{x}$ . Hence,  $\mathbf{A} \mathbf{x}' = y_{n+1} \mathbf{A} \mathbf{x} = \mathbf{b} y_{n+1}$ .

### 18.5

Let  $\mathbf{x} \in \mathbb{R}^n$ , and  $\mathbf{y} = \mathbf{T}(\mathbf{x})$ . Let  $\mathbf{a}_i$  be the  $i$ th column of  $\mathbf{A}$ ,  $i = 1, \dots, n$ . As in the hint, let  $\mathbf{A}'$  be given by

$$\mathbf{A}' = [a_1 \mathbf{a}_1, \dots, a_n \mathbf{a}_n, -\mathbf{b}].$$

Then,

$$\begin{aligned} \mathbf{A} \mathbf{x} = \mathbf{b} &\Leftrightarrow \mathbf{A} \mathbf{x} - \mathbf{b} = \mathbf{0} \\ &\Leftrightarrow [a_1, \dots, a_n, -\mathbf{b}] \begin{bmatrix} x_1 \\ \vdots \\ x_n \\ 1 \end{bmatrix} = \mathbf{0} \\ &\Leftrightarrow [a_1 \mathbf{a}_1, \dots, a_n \mathbf{a}_n, -\mathbf{b}] \begin{bmatrix} x_1/a_1 \\ \vdots \\ x_n/a_n \\ 1 \end{bmatrix} = \mathbf{0} \\ &\Leftrightarrow \mathbf{A}' \begin{bmatrix} (x_1/a_1)y_{n+1} \\ \vdots \\ (x_n/a_n)y_{n+1} \\ y_{n+1} \end{bmatrix} = \mathbf{0} \\ &\Leftrightarrow \mathbf{A}' \mathbf{y} = \mathbf{0}. \end{aligned}$$

### 18.6

The result follows from Exercise 18.5 by setting  $\mathbf{A} := \mathbf{c}^\top$  and  $\mathbf{b} := 0$ .

### 18.7

Consider the set  $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{e}^\top \mathbf{x} = 1, \mathbf{x} \geq \mathbf{0}, x_1 = 0\}$ , which can be written as  $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{A} \mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$ , where

$$\mathbf{A} = \begin{bmatrix} \mathbf{e}^\top \\ \mathbf{e}_1^\top \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 0 \end{bmatrix},$$

with  $\mathbf{e} = [1, \dots, 1]^\top$ ,  $\mathbf{e}_1 = [1, 0, \dots, 0]^\top$ . Let  $\mathbf{a}_0 = \mathbf{e}/n$ . By Exercise 12.20, the closest point on the set  $\{\mathbf{x} : \mathbf{A} \mathbf{x} = \mathbf{b}\}$  to the point  $\mathbf{a}_0$  is

$$\mathbf{x}^* = \mathbf{A}^\top (\mathbf{A} \mathbf{A}^\top)^{-1} (\mathbf{b} - \mathbf{A} \mathbf{a}_0) + \mathbf{a}_0 = \left[0, \frac{1}{n-1}, \dots, \frac{1}{n-1}\right]^\top.$$

Since  $\mathbf{x}^* \in \{\mathbf{x} : \mathbf{A} \mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}\} \subset \{\mathbf{x} : \mathbf{A} \mathbf{x} = \mathbf{b}\}$ , the point  $\mathbf{x}^*$  is also the closest point on the set  $\{\mathbf{x} : \mathbf{A} \mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$  to the point  $\mathbf{a}_0$ .

Let  $r = \|\mathbf{a}_0 - \mathbf{x}^*\|$ . Then, the sphere of radius  $r$  is inscribed in  $\Delta$ . Note that

$$r = \|\mathbf{a}_0 - \mathbf{x}^*\| = \frac{1}{\sqrt{n(n-1)}}.$$

Hence, the radius of the largest sphere inscribed in  $\Delta$  is larger than or equal to  $1/\sqrt{n(n-1)}$ . It remains to show that this largest radius is less than or equal to  $1/\sqrt{n(n-1)}$ . To this end, we show that this largest radius is less than or equal to  $1/\sqrt{n(n-1)} + \varepsilon$  for any  $\varepsilon > 0$ . For this, it suffices to show that the sphere of

radius  $1/\sqrt{n(n-1)} + \varepsilon$  is not inscribed in  $\Delta$ . To show this, consider the point

$$\begin{aligned}\mathbf{x} &= \mathbf{x}^* + \varepsilon \frac{\mathbf{x}^* - \mathbf{a}_0}{\|\mathbf{x}^* - \mathbf{a}_0\|} \\ &= \mathbf{x}^* + \varepsilon \sqrt{n(n-1)}(\mathbf{x}^* - \mathbf{a}_0) \\ &= \mathbf{x}^* + \varepsilon \left[ -\sqrt{\frac{n-1}{n}}, \frac{1}{\sqrt{n(n-1)}}, \dots, \frac{1}{\sqrt{n(n-1)}} \right]^\top.\end{aligned}$$

It is easy to verify that the point  $\mathbf{x}$  above is on the sphere of radius  $1/\sqrt{n(n-1)} + \varepsilon$ . However, clearly the first component of  $\mathbf{x}$  is negative. Therefore, the sphere of radius  $1/\sqrt{n(n-1)} + \varepsilon$  is not inscribed in  $\Delta$ . Our proof is thus completed.

### 18.8

We first consider the constraints. We claim that  $\mathbf{x} \in \Delta \Leftrightarrow \bar{\mathbf{x}} \in \Delta$ . To see this, note that if  $\mathbf{x} \in \Delta$ , then  $\mathbf{e}^\top \mathbf{x} = 1$  and hence

$$\mathbf{e}^\top \bar{\mathbf{x}} = \mathbf{e}^\top \mathbf{D}^{-1} \mathbf{x} / \mathbf{e}^\top \mathbf{D}^{-1} \mathbf{x} = 1,$$

which means that  $\bar{\mathbf{x}} \in \Delta$ . The same argument can be used for the converse. Next, we claim  $\mathbf{A}\mathbf{x} = \mathbf{0} \Leftrightarrow \mathbf{A}\mathbf{D}\bar{\mathbf{x}} = \mathbf{0}$ . To see this, write

$$\mathbf{A}\mathbf{x} = \mathbf{A}\mathbf{D}\mathbf{D}^{-1}\mathbf{x} = \mathbf{A}\mathbf{D}\bar{\mathbf{x}}(\mathbf{e}^\top \mathbf{D}^{-1}\mathbf{x}).$$

Since  $\mathbf{e}^\top \mathbf{D}^{-1}\mathbf{x} > 0$ , we have  $\mathbf{A}\mathbf{x} = \mathbf{0} \Leftrightarrow \mathbf{A}\mathbf{D}\bar{\mathbf{x}} = \mathbf{0}$ . Finally, we claim that if  $\mathbf{x}^*$  is an optimal solution to the original problem, then  $\bar{\mathbf{x}}^* = \mathbf{U}(\mathbf{x}^*)$  is an optimal solution to the transformed problem. To see this, recall that the problem in a Karmarkar's restricted problem, and hence by Assumption (B) we have  $\mathbf{c}^\top \mathbf{x}^* = 0$ . We now note that the minimum value of the objective function  $\mathbf{c}^\top \mathbf{D}\bar{\mathbf{x}}$  in the transformed problem is zero. This is because  $\mathbf{c}^\top \mathbf{D}\bar{\mathbf{x}} = \mathbf{c}^\top \mathbf{x} / \mathbf{e}^\top \mathbf{D}^{-1}\mathbf{x}$ , and  $\mathbf{e}^\top \mathbf{D}^{-1}\mathbf{x} > 0$ . Finally, we observe that at the point,  $\bar{\mathbf{x}}^* = \mathbf{U}(\mathbf{x}^*)$  the objective function value for the transformed problem is zero. Indeed,

$$\mathbf{c}^\top \mathbf{D}\bar{\mathbf{x}}^* = \mathbf{c}^\top \mathbf{D}\mathbf{D}^{-1}\mathbf{x}^* / \mathbf{e}^\top \mathbf{D}^{-1}\mathbf{x}^* = 0.$$

Therefore, the two problems are equivalent.

### 18.9

Let  $\mathbf{v} \in \mathbb{R}^{m+1}$  be such that  $\mathbf{v}^\top \mathbf{B} = \mathbf{0}^\top$ . We will show that

$$\mathbf{v}^\top \begin{bmatrix} \mathbf{A} \\ \mathbf{e}^\top \end{bmatrix} = \mathbf{0}^\top$$

and hence  $\mathbf{v} = \mathbf{0}$  by virtue of the assumption that

$$\text{rank} \begin{bmatrix} \mathbf{A} \\ \mathbf{e}^\top \end{bmatrix} = m + 1.$$

This in turn gives us the desired result.

To proceed, write  $\mathbf{v}$  as

$$\mathbf{v} = \begin{bmatrix} \mathbf{u} \\ v_{m+1} \end{bmatrix}$$

where  $\mathbf{u} \in \mathbb{R}^m$  constitute the first  $m$  components of  $\mathbf{v}$ . Then,

$$\mathbf{v}^\top \mathbf{B} = \mathbf{u}^\top \mathbf{A}\mathbf{D} + v_{m+1} \mathbf{e}^\top = \mathbf{0}^\top.$$

Postmultiplying the above by  $\mathbf{e}$ , and using the facts that  $\mathbf{D}\mathbf{e} = \mathbf{x}_0$ ,  $\mathbf{A}\mathbf{x}_0 = \mathbf{0}$ , and  $\mathbf{e}^\top \mathbf{e} = n$ , we get

$$\mathbf{u}^\top \mathbf{A}\mathbf{x}_0 + v_{m+1}n = v_{m+1}n = 0,$$

which implies that  $v_{m+1} = 0$ . Hence,  $\mathbf{u}^\top \mathbf{A} \mathbf{D} = \mathbf{0}^\top$ , which after postmultiplying by  $\mathbf{D}^{-1}$  gives  $\mathbf{u}^\top \mathbf{A} = \mathbf{0}^\top$ . Hence,

$$\mathbf{v}^\top \begin{bmatrix} \mathbf{A} \\ \mathbf{e}^\top \end{bmatrix} = \mathbf{0}^\top,$$

which implies that  $\mathbf{v} = \mathbf{0}$ . Hence,  $\text{rank } \mathbf{B} = m + 1$ .

### 18.10

We proceed by induction. For  $k = 0$ , the result is true because  $\mathbf{x}^{(0)} = \mathbf{a}_0$ . Now suppose that  $\mathbf{x}^{(k)}$  is a strictly interior point of  $\Delta$ . We first show that  $\bar{\mathbf{x}}^{(k+1)}$  is a strictly interior point. Now,

$$\bar{\mathbf{x}}^{(k+1)} = \mathbf{a}_0 - \alpha r \hat{\mathbf{c}}^{(k)}.$$

Then, since  $\alpha \in (0, 1)$  and  $\|\hat{\mathbf{c}}^{(k)}\| = 1$ , we have

$$\|\bar{\mathbf{x}}^{(k+1)} - \mathbf{a}_0\| \leq |\alpha r| \|\hat{\mathbf{c}}^{(k)}\| < r.$$

Since  $r$  is the radius of the largest sphere inscribed in  $\Delta$ ,  $\bar{\mathbf{x}}^{(k+1)}$  is a strictly interior point of  $\Delta$ . To complete the proof, we write

$$\mathbf{x}^{(k+1)} = \mathbf{U}_k^{-1}(\bar{\mathbf{x}}^{(k+1)}) = \frac{\mathbf{D}_k \bar{\mathbf{x}}^{(k+1)}}{\mathbf{e}^\top \mathbf{D}_k \bar{\mathbf{x}}^{(k+1)}}$$

We already know that  $\mathbf{x}^{(k+1)} \in \Delta$ . It therefore remains to show that it is strictly interior, i.e.,  $\mathbf{x}^{(k+1)} > \mathbf{0}$ . To see this, note that  $\mathbf{e}^\top \mathbf{D}_k \bar{\mathbf{x}}^{(k+1)} > 0$ . Furthermore, we can write

$$\mathbf{D}_k \bar{\mathbf{x}}^{(k+1)} = \begin{bmatrix} \bar{x}_1^{(k+1)} x_1^{(k)} \\ \vdots \\ \bar{x}_n^{(k+1)} x_n^{(k)} \end{bmatrix}.$$

Since  $\mathbf{x}^{(k)} = [x_1^{(k)}, \dots, x_n^{(k)}]^\top > \mathbf{0}$  by the induction hypothesis, and  $\bar{\mathbf{x}}^{(k+1)} = [\bar{x}_1^{(k+1)}, \dots, \bar{x}_n^{(k+1)}]^\top > \mathbf{0}$  by the above,  $\mathbf{x}^{(k+1)} > \mathbf{0}$  and hence it is a strictly interior point of  $\Delta$ .

## 19. Integer Linear Programming

### 19.1

The result follows from the simple observation that if  $\mathbf{M}$  is a submatrix of  $\mathbf{A}$ , then any submatrix of  $\mathbf{M}$  is also a submatrix of  $\mathbf{A}$ . Therefore, any property involving all submatrices of  $\mathbf{A}$  also applies to all submatrices of  $\mathbf{M}$ .

### 19.2

The result follows from the simple observation that any submatrix of  $\mathbf{A}^\top$  is the transpose of a submatrix of  $\mathbf{A}$ , and that the determinant of the transpose of a matrix equals the determinant of the original matrix.

### 19.3

The claim that  $\mathbf{A}$  is totally unimodular if  $[\mathbf{A}, \mathbf{I}]$  is totally unimodular follows from Exercise 19.1. To show the converse, suppose that  $\mathbf{A}$  is totally unimodular. We will show that any  $p \times p$  invertible submatrix of  $[\mathbf{A}, \mathbf{I}]$ ,  $p \leq \min(m, n)$ , has determinant  $\pm 1$ . We first note that any  $p \times p$  invertible submatrix of  $[\mathbf{A}, \mathbf{I}]$  that consists only of columns of  $\mathbf{A}$  has determinant  $\pm 1$  because  $\mathbf{A}$  is totally unimodular. Moreover, any  $p \times p$  invertible submatrix of  $\mathbf{I}$  has determinant 1.

Consider now a  $p \times p$  invertible submatrix of  $[\mathbf{A}, \mathbf{I}]$  composed of  $k$  columns of  $\mathbf{A}$  and  $p - k$  columns of  $\mathbf{I}$ . Without loss of generality, suppose that this submatrix is composed of the first  $p$  rows of  $[\mathbf{A}, \mathbf{I}]$ , the last  $k$  columns of  $\mathbf{A}$ , and the first  $p - k$  columns of  $\mathbf{I}$ . (This choice of rows and columns is without loss of generality because we can exchange rows and columns to arrive at this form, and each exchange only changes the sign of the determinant.) We now proceed as in the proof of Proposition 19.1.

### 19.4

The result follows from these properties of determinants: (1) that exchanging columns only changes the sign

of the determinant; (2) the determinant of a block triangular matrix is the product of the determinants of the diagonal blocks; and (3) the determinant of the identity matrix is 1. See also Exercise 2.4.

### 19.5

The vectors  $\mathbf{x}$  and  $\mathbf{z}$  together satisfy

$$\mathbf{A}\mathbf{x} + \mathbf{z} = \mathbf{b},$$

which means that  $\mathbf{z} = \mathbf{b} - \mathbf{A}\mathbf{x}$ . Because the right-hand side involves only integers,  $\mathbf{z}$  is an integer vector.

### 19.6

The following MATLAB code generates the figures.

```
%-----
% The vertices of the feasible set
x = [2/5 1; 2/5 -2/5]\[3; 1];
X = [0 0 x(1) 2.5];
Y = [0 3 x(2) 0];
fs=16; %FontSize
% Draw the feasible set for x1 x2 \in R.
vi = convhull(X,Y);
plot(X,Y, 'o');
axis on; axis equal;
axis([-0.2 4.2 -0.2 3.2]);
hold on
fill (X(vi), Y(vi), 'b','facealpha', 0.2);
text(.1,.5,['\fontsize{48}\Omega'],'position', [1.5 1.25])
set(gca,'FontSize',fs)
hold off
% The optimal solution has to be one of the extreme points
c = [-3 -4]';
% Draw the feasible set for the noninteger problem
figure
axis on; axis equal;
x = [-0.5:0.1:x(1)];
y1 = -x*0.4+3;
y2 = x-2.5;
fs=16; % Fontsize
plot(x,y1,'--b',x,y2,'--b','LineWidth',2);
axis([-0.2 4.2 -0.2 3.2]);
set(gca,'FontSize',fs)
hold on
X = [zeros(1,4) ones(1,3) 2*ones(1,3) 3];
Y = [0:max(floor(Y)) 0:max(floor(Y)-1) 0:max(floor(Y)-1) 1];
plot(X,Y,'b|s','LineWidth',2,...
      'MarkerEdgeColor','k',...
      'MarkerFaceColor','g',...
      'MarkerSize',10)
% Plot of the cost function
xc = [-1:0.5:5];
yc = -0.75*xc+(14/4)*ones(1,length(xc));
yc0 = -0.75*xc+(17.5/4)*ones(1,length(xc));
fs=16; % Fontsize
plot(xc, yc, 'r', xc, yc0, 'o-k', 'LineWidth',2);
set(gca,'FontSize',fs)
%text(.1,.5,['\fontsize{48}\Omega'],'position', [1.5 1.25])
[~,Xmin] = min(c'*[X; Y]);
str = sprintf('The maximizer is [%d, %d]' and the maximum is %.4f',...
              X(Xmin), Y(Xmin), -c'*[X(Xmin); Y(Xmin)]);
disp(str);
%-----
```

---

**19.7**

It suffices to show the following claim: If we introduce the equation

$$x_i + \sum_{j=m+1}^n \lfloor y_{ij} \rfloor x_j + x_{n+1} = \lfloor y_{i0} \rfloor$$

into the original constraint, then the result holds. The reason this suffices is that the Gomory cut is obtained by subtracting this equation from an equation obtained by elementary row operations on  $[\mathbf{A}, \mathbf{b}]$  (hence is equivalent to premultiplication by an invertible matrix).

To show the above claim, let  $x_{n+1}$  satisfy this constraint with an integer vector  $\mathbf{x}$ . Then,

$$x_i + \sum_{j=m+1}^n \lfloor y_{ij} \rfloor x_j + x_{n+1} = \lfloor y_{i0} \rfloor,$$

which implies that

$$x_{n+1} = \lfloor y_{i0} \rfloor - x_i - \sum_{j=m+1}^n \lfloor y_{ij} \rfloor x_j.$$

Because the right-hand side involves only integers,  $x_{n+1}$  is an integer.

---

**19.8**

If there is only one Gomory cut, then the result follows directly from Exercise 19.7. The general result follows by induction on the number of Gomory cuts, using Exercise 19.7 at each inductive step.

---

**19.9**

The result follows from Exercises 19.5 and 19.8.

---

**19.10**

The dual problem is:

$$\begin{aligned} &\text{minimize} && 3\lambda_1 + \lambda_2 \\ &\text{subject to} && \frac{2}{5}\lambda_1 + \frac{2}{5}\lambda_2 \geq 3 \\ &&& 1\lambda_1 - \frac{2}{5}\lambda_2 \geq 4 \\ &&& \lambda_1, \lambda_2 \geq 0 \\ &&& \lambda_1, \lambda_2 \in \mathbb{Z}. \end{aligned}$$

The problem is solved graphically using the same approach as in Example 19.5. We proceed by calculating the extreme points of the feasible set. We first assume that  $\lambda_1, \lambda_2 \in \mathbb{R}$ . The extreme points are calculated intersecting the given constraints, and they are:

$$\boldsymbol{\lambda}^{(1)} = [5, \frac{5}{2}]^\top, \quad \boldsymbol{\lambda}^{(2)} = [\frac{15}{2}, 0]^\top.$$

In Figure 24.10, we show the feasible set  $\Omega$  for the case when  $\lambda_1, \lambda_2 \in \mathbb{R}$ .

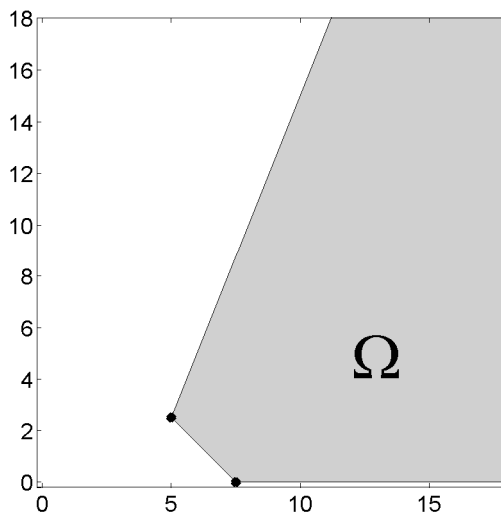
Next we sketch the feasible set for the case when  $\lambda_1, \lambda_2 \in \mathbb{Z}$  and solve the problem graphically. The graphical solution is depicted in Figure 24.11. We can see in Figure 24.11 that the optimal integer solution is

$$\mathbf{x}^* = [6, 2]^\top.$$

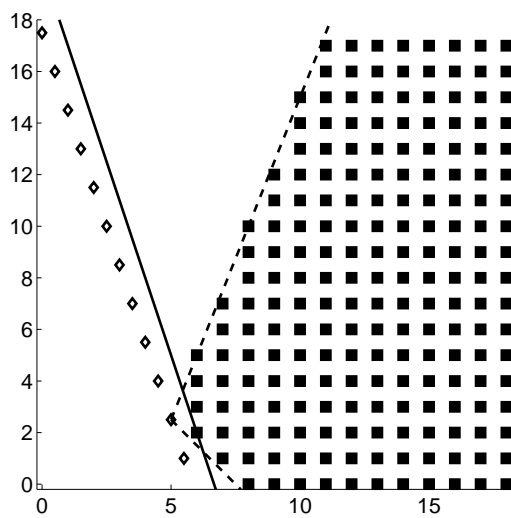
The following MATLAB code generates the figures.

```
%-----
% The vertices of the feasible set are:

x = [2/5 2/5; 1 -2/5] \ [3; 4];
x = [7.5 x(1) 20 20];
```



**Figure 24.10** Feasible set  $\Omega$  for the case when  $\lambda_1, \lambda_2 \in \mathbb{R}$  in Example 19.5.



**Figure 24.11** Real feasible set with  $\lambda_1, \lambda_2 \in \mathbb{Z}$ .

```

Y = [0 x(2) 0 40];
fs=16; % Fontsize
% Now we draw the set Omega, supposing x1 x2 \in R.
vi = convhull(X,Y);
plot(X(1:2),Y(1:2), 'rX', 'LineWidth',4);
axis on; axis equal;
axis([-0.2 18 -0.2 18]);
set(gca,'FontSize',fs)
%title('Feasible set supposing x1, x2 in R.','FontSize',14,'Fontname','Avant-garde');
hold on
fill (X(vi), Y(vi), 'b','facealpha', 0.2);
text(.1,.5,['\fontsize{48}\Omega'],'position', [12 5])
hold off

% We now the optimal solution has to be one of the extreme points.
c = [3 1]';

%Now we draw the real feasible set for the problem.
figure
axis on; axis equal;
axis([-0.2 18 -0.2 18]);
set(gca,'FontSize',fs)
%title('Feasible set and cost function','FontSize',14,'Fontname','Avant-garde');
hold on

X = [];
Y = [];

for i=1:18
    j=0;
    while ((j<=((i-4)*2.5)) && (j<18))
        if((j>=(7.5-i)))
            X = [X i];
            Y = [Y j];
        end
        j=j+1;
    end
end

plot(X,Y,'b|s','LineWidth',1,...
     'MarkerEdgeColor','k',...
     'MarkerFaceColor','g',...
     'MarkerSize',10)

x = [5:0.1:18];
y1 = (x-4)*2.5;
y2 = (7.5-x);

plot(x,y1,'--b',x,y2,'--b','LineWidth',2);
set(gca,'FontSize',fs)
%text(.1,.5,['\fontsize{48}\Omega'],'position', [12 5])

%Plot of the cost function at level 17.5
xc = [-1:0.5:18];
yc = (35/2)*ones(1,length(xc))-3*xc;
plot(xc, yc, 'dk', 'LineWidth',2);

%Plot of the cost function
xc = [-1:0.5:18];
yc = (20)*ones(1,length(xc))-3*xc;

```



```

plot(xc, yc, 'r', 'LineWidth',2);

[~,Xmin] = min(c'*[X; Y]);

str = sprintf('The minimizer is [%d, %d]'' and the maximum is %.4f]',...
    X(Xmin), Y(Xmin), c'*[X(Xmin); Y(Xmin)]);
disp(str);
%-----

```

## 20. Problems with Equality Constraints

### 20.1

The feasible set consists of the points

$$\mathbf{x} = \begin{bmatrix} 2 \\ a \end{bmatrix}, \quad a \leq -1.$$

We next find the gradients:

$$\nabla h(\mathbf{x}) = \begin{bmatrix} 2(x_1 - 2) \\ 0 \end{bmatrix} \quad \text{and} \quad \nabla g(\mathbf{x}) = \begin{bmatrix} 0 \\ 3(x_2 + 1)^2 \end{bmatrix}.$$

All feasible points are not regular because at the above points the gradients of  $h$  and  $g$  are not linearly independent. There are no regular points of the constraints.

### 20.2

a. As usual, let  $f$  be the objective function, and  $\mathbf{h}$  the constraint function. We form the Lagrangian  $l(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{h}(\mathbf{x})$ , and then find critical points by solving the following equations (Lagrange condition):

$$D_{\mathbf{x}}l(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{0}^\top, \quad D_{\boldsymbol{\lambda}}l(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{0}^\top.$$

We obtain

$$\begin{bmatrix} 2 & 2 & 0 & 1 & 4 \\ 2 & 6 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 5 \\ 1 & 2 & 0 & 0 & 0 \\ 4 & 0 & 5 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \lambda_1 \\ \lambda_2 \end{bmatrix} = \begin{bmatrix} -4 \\ -5 \\ -6 \\ 3 \\ 6 \end{bmatrix}.$$

The unique solution to the above system is

$$\mathbf{x}^* = [16/5, -1/10, -34/25]^\top, \quad \boldsymbol{\lambda}^* = [-27/5, -6/5]^\top.$$

Note that  $\mathbf{x}^*$  is a regular point. We now apply the SOSC. We compute

$$\mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{F}(\mathbf{x}^*) + [\boldsymbol{\lambda}^* \mathbf{H}(\mathbf{x}^*)] = \begin{bmatrix} 2 & 2 & 0 \\ 2 & 6 & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

The tangent plane is

$$\begin{aligned} T(\mathbf{x}^*) &= \left\{ \mathbf{y} \in \mathbb{R}^3 : \begin{bmatrix} 1 & 2 & 0 \\ 4 & 0 & 5 \end{bmatrix} \mathbf{y} = \mathbf{0} \right\} \\ &= \{a[-5/4, 5/8, 1]^\top : a \in \mathbb{R}\}. \end{aligned}$$

Let  $\mathbf{y} = a[-5/4, 5/8, 1]^\top \in T(\mathbf{x}^*)$ ,  $a \neq 0$ . We have

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{y} = \frac{75}{32} a^2 > 0.$$

Therefore,  $\mathbf{x}^*$  is a strict local minimizer.

b. The Lagrange condition for this problem is

$$\begin{aligned} 4 + 2\lambda x_1 &= 0 \\ 2x_2 + 2\lambda x_2 &= 0 \\ x_1^2 + x_2^2 - 9 &= 0. \end{aligned}$$

We have four points satisfying the Lagrange condition:

$$\begin{aligned}\mathbf{x}^{(1)} &= [3, 0]^\top, & \lambda^{(1)} &= -2/3 \\ \mathbf{x}^{(2)} &= [-3, 0]^\top, & \lambda^{(2)} &= 2/3 \\ \mathbf{x}^{(3)} &= [2, \sqrt{5}]^\top, & \lambda^{(3)} &= -1 \\ \mathbf{x}^{(4)} &= [2, -\sqrt{5}]^\top, & \lambda^{(4)} &= -1.\end{aligned}$$

Note that all four points  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(4)}$  are regular. We now apply the SOSC. We have

$$\begin{aligned}\mathbf{L}(\mathbf{x}, \lambda) &= \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix} + \lambda \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \\ T(\mathbf{x}) &= \{\mathbf{y} : [2x_1, 2x_2]\mathbf{y} = 0\}.\end{aligned}$$

For the first point, we have

$$\begin{aligned}\mathbf{L}(\mathbf{x}^{(1)}, \lambda^{(1)}) &= \begin{bmatrix} -4/3 & 0 \\ 0 & 2/3 \end{bmatrix} \\ T(\mathbf{x}^{(1)}) &= \{a[0, 1]^\top : a \in \mathbb{R}\}.\end{aligned}$$

Let  $\mathbf{y} = a[0, 1]^\top \in T(\mathbf{x}^{(1)})$ ,  $a \neq 0$ . Then

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^{(1)}, \lambda^{(1)}) \mathbf{y} = \frac{2}{3}a^2 > 0.$$

Hence,  $\mathbf{x}^{(1)}$  is a strict local minimizer.

For the second point, we have

$$\mathbf{L}(\mathbf{x}^{(2)}, \lambda^{(2)}) = \begin{bmatrix} 4/3 & 0 \\ 0 & 10/3 \end{bmatrix} > 0.$$

Hence,  $\mathbf{x}^{(2)}$  is a strict local minimizer.

For the third point, we have

$$\begin{aligned}\mathbf{L}(\mathbf{x}^{(3)}, \lambda^{(3)}) &= \begin{bmatrix} -2 & 0 \\ 0 & 0 \end{bmatrix} \\ T(\mathbf{x}^{(3)}) &= \{a[-\sqrt{5}, 2]^\top : a \in \mathbb{R}\}.\end{aligned}$$

Let  $\mathbf{y} = a[-\sqrt{5}, 2]^\top \in T(\mathbf{x}^{(3)})$ ,  $a \neq 0$ . Then

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^{(3)}, \lambda^{(3)}) \mathbf{y} = -10a^2 < 0.$$

Hence,  $\mathbf{x}^{(3)}$  is a strict local maximizer.

For the fourth point, we have

$$\begin{aligned}\mathbf{L}(\mathbf{x}^{(4)}, \lambda^{(4)}) &= \begin{bmatrix} -2 & 0 \\ 0 & 0 \end{bmatrix} \\ T(\mathbf{x}^{(4)}) &= \{a[\sqrt{5}, 2]^\top : a \in \mathbb{R}\}.\end{aligned}$$

Let  $\mathbf{y} = a[\sqrt{5}, 2]^\top \in T(\mathbf{x}^{(4)})$ ,  $a \neq 0$ . Then

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^{(4)}, \lambda^{(4)}) \mathbf{y} = -10a^2 < 0.$$

Hence,  $\mathbf{x}^{(4)}$  is a strict local maximizer.

c. The Lagrange condition for this problem is

$$\begin{aligned}x_2 + 2\lambda x_1 &= 0 \\x_1 + 8\lambda x_2 &= 0 \\x_1^2 + 4x_2^2 - 1 &= 0.\end{aligned}$$

We have four points satisfying the Lagrange condition:

$$\begin{aligned}\mathbf{x}^{(1)} &= [1/\sqrt{2}, -1/(2\sqrt{2})]^\top, & \lambda^{(1)} &= 1/4 \\ \mathbf{x}^{(2)} &= [-1/\sqrt{2}, 1/(2\sqrt{2})]^\top, & \lambda^{(2)} &= 1/4 \\ \mathbf{x}^{(3)} &= [1/\sqrt{2}, 1/(2\sqrt{2})]^\top, & \lambda^{(3)} &= -1/4 \\ \mathbf{x}^{(4)} &= [-1/\sqrt{2}, -1/(2\sqrt{2})]^\top, & \lambda^{(4)} &= -1/4.\end{aligned}$$

Note that all four points  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(4)}$  are regular. We now apply the SOSC. We have

$$\begin{aligned}\mathbf{L}(\mathbf{x}, \lambda) &= \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} + \lambda \begin{bmatrix} 2 & 0 \\ 0 & 8 \end{bmatrix}, \\ T(\mathbf{x}) &= \{\mathbf{y} : [2x_1, 8x_2]\mathbf{y} = 0\}.\end{aligned}$$

Note that

$$\begin{aligned}\mathbf{L}(\mathbf{x}, -1/4) &= \begin{bmatrix} -1/2 & 1 \\ 1 & -2 \end{bmatrix} \\ \mathbf{L}(\mathbf{x}, 1/4) &= \begin{bmatrix} 1/2 & 1 \\ 1 & 2 \end{bmatrix}.\end{aligned}$$

After standard manipulations, we conclude that the first two points are strict local maximizers, while the last two points are strict local minimizers.

### 20.3

We form the lagrangian

$$l(\mathbf{x}, \boldsymbol{\lambda}) = (\mathbf{a}^\top \mathbf{x})(\mathbf{b}^\top \mathbf{x}) + \lambda_1(x_1 + x_2) + \lambda_2(x_2 + x_3).$$

The Lagrange conditions take the form,

$$\begin{aligned}\nabla_{\mathbf{x}} l &= (\mathbf{a}\mathbf{b}^\top + \mathbf{b}\mathbf{a}^\top)\mathbf{x} + \begin{bmatrix} \nabla_{\mathbf{x}} h_1(\mathbf{x}) & \nabla_{\mathbf{x}} h_2(\mathbf{x}) \end{bmatrix} \boldsymbol{\lambda} \\ &= \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{bmatrix} \boldsymbol{\lambda} \\ &= \begin{bmatrix} x_2 + \lambda_1 \\ x_1 + x_3 + \lambda_1 + \lambda_2 \\ x_2 + \lambda_2 \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \\ \mathbf{h}(\mathbf{x}) &= \begin{bmatrix} x_1 + x_2 \\ x_2 + x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.\end{aligned}$$

It is easy to see that  $\mathbf{x}^* = \mathbf{0}$  and  $\boldsymbol{\lambda}^* = \mathbf{0}$  satisfy the Lagrange, FONC, conditions.

The Hessian of the lagrangian is

$$\mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{a}\mathbf{b}^\top + \mathbf{b}\mathbf{a}^\top = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

and the tangent space

$$T(\mathbf{x}^*) = \left\{ \mathbf{y} : \mathbf{y} = a \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}, a \in \mathbb{R} \right\}.$$

To verify if the critical point satisfies the SOSC, we evaluate

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{y} = -4a^2 < 0.$$

Thus the critical point is a strict local maximizer.

#### 20.4

By the Lagrange condition,  $\mathbf{x}^* = [x_1, x_2]^\top$  satisfies

$$\begin{aligned} x_1 + \lambda &= 0 \\ x_1 + 4 + 4\lambda &= 0. \end{aligned}$$

Eliminating  $\lambda$  we get

$$3x_1 - 4 = 0$$

which implies that  $x_1 = 4/3$ . Therefore,  $\nabla f(\mathbf{x}^*) = [4/3, 16/3]^\top$ .

#### 20.5

a. The Lagrange condition for this problem is:

$$\begin{aligned} 2(\mathbf{x}^* - \mathbf{x}_0) + 2\lambda^* \mathbf{x}^* &= \mathbf{0} \\ \|\mathbf{x}^*\|^2 &= 9, \end{aligned}$$

where  $\lambda^* \in \mathbb{R}$ . Rewriting the first equation we get  $(1 + \lambda^*)\mathbf{x}^* = \mathbf{x}_0$ , which when combined with the second equation gives two values for  $1 + \lambda^*$ :  $1 + \lambda_1^* = 2/3$  and  $1 + \lambda_2^* = -2/3$ . Hence there are two solutions to the Lagrange condition:  $\mathbf{x}^{*(1)} = (3/2)[1, \sqrt{3}]$ , and  $\mathbf{x}^{*(2)} = -(3/2)[1, \sqrt{3}]$ .

b. We have  $\mathbf{L}(\mathbf{x}^{*(i)}, \lambda_i^*) = (1 + \lambda_i^*)\mathbf{I}$ . To apply the SONC Theorem, we need to check regularity. This is easy, since the gradient of the constraint function at any point  $\mathbf{x}$  is  $2\mathbf{x}$ , which is nonzero at both the points in part a.

For the second point,  $1 + \lambda_2^* = -2/3$ , which implies that the point is not a local minimizer because the SONC does not hold.

On the other hand, the first point satisfies the SOSC (since  $1 + \lambda_1^* = 2/3$ ), which implies that it is a strict local minimizer.

#### 20.6

a. Let  $x_1$ ,  $x_2$ , and  $x_3$  be the dimensions of the closed box. The problem is

$$\begin{aligned} \text{minimize} \quad & 2(x_1x_2 + x_2x_3 + x_3x_1) \\ \text{subject to} \quad & x_1x_2x_3 = V. \end{aligned}$$

We denote  $f(\mathbf{x}) = 2(x_1x_2 + x_2x_3 + x_3x_1)$ , and  $h(\mathbf{x}) = x_1x_2x_3 - V$ . We have  $\nabla f(\mathbf{x}) = 2[x_2 + x_3, x_1 + x_3, x_1 + x_2]^\top$  and  $\nabla h(\mathbf{x}) = [x_2x_3, x_1x_3, x_1x_2]^\top$ . By the Lagrange condition, the dimensions of the box with minimum surface area satisfies

$$\begin{aligned} 2(b + c) + \lambda bc &= 0 \\ 2(a + c) + \lambda ac &= 0 \\ 2(a + b) + \lambda ab &= 0 \\ abc &= V, \end{aligned}$$

where  $\lambda \in \mathbb{R}$ .

b. Regularity of  $\mathbf{x}^*$  means  $\nabla h(\mathbf{x}^*) \neq 0$  (since there is only one scalar equality constraint in this case). Since  $\mathbf{x}^* = [a, b, c]^\top$  is a feasible point, we must have  $a, b, c \neq 0$  (for otherwise the volume will be 0). Hence,  $\nabla h(\mathbf{x}^*) \neq 0$ , which implies that  $\mathbf{x}^*$  is regular.

c. Multiplying the first equation by  $a$  and the second equation by  $b$ , and then subtracting the first from the second, we obtain:

$$c(a - b) = 0.$$

Since  $c \neq 0$  (see part b), we conclude that  $a = b$ . By a similar procedure on the second and third equations, we conclude that  $b = c$ . Hence, substituting into the fourth (constraint) equation, we obtain

$$a = b = c = V^{1/3},$$

with  $\lambda = -4V^{-1/3}$ .

d. The Hessian of the Lagrangian is given by

$$\mathbf{L}(\mathbf{x}^*, \lambda) = \begin{bmatrix} 0 & 2 + \lambda c & 2 + \lambda b \\ 2 + \lambda c & 0 & 2 + \lambda a \\ 2 + \lambda b & 2 + \lambda a & 0 \end{bmatrix} = \begin{bmatrix} 0 & -2 & -2 \\ -2 & 0 & -2 \\ -2 & -2 & 0 \end{bmatrix} = -2 \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}.$$

The matrix  $\mathbf{L}(\mathbf{x}^*, \lambda)$  is not positive definite (there are several ways to check this: we could use Sylvester's criterion, or we could compute the eigenvalues of  $\mathbf{L}(\mathbf{x}^*, \lambda)$ , which are  $2, 2, -4$ ). Therefore, we need to compute the tangent space  $T(\mathbf{x}^*)$ . Note that

$$Dh(\mathbf{x}^*) = \nabla h(\mathbf{x}^*)^\top = [bc, ac, ab] = V^{2/3}[1, 1, 1].$$

Hence,

$$T(\mathbf{x}^*) = \{\mathbf{y} : Dh(\mathbf{x}^*)\mathbf{y} = 0\} = \{\mathbf{y} : [1, 1, 1]\mathbf{y} = 0\} = \{\mathbf{y} : y_3 = -(y_1 + y_2)\}.$$

Let  $\mathbf{y} \in T(\mathbf{x}^*)$ ,  $\mathbf{y} \neq \mathbf{0}$ . Note that either  $y_1 \neq 0$  or  $y_2 \neq 0$ . We have,

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \lambda) \mathbf{y} = -2\mathbf{y}^\top \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} \mathbf{y} = -4(y_1 y_2 + y_1 y_3 + y_2 y_3).$$

Substituting  $y_3 = -(y_1 + y_2)$ , we obtain

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \lambda) \mathbf{y} = -4(y_1 y_2 - y_1(y_1 + y_2) - y_2(y_1 + y_2)) = 4(y_1^2 + y_1 y_2 + y_2^2) = 4\mathbf{z}^\top \mathbf{Q} \mathbf{z}$$

where  $\mathbf{z} = [y_1, y_2]^\top \neq \mathbf{0}$  and

$$\mathbf{Q} = \begin{bmatrix} 1 & 1/2 \\ 1/2 & 1 \end{bmatrix} > 0.$$

Therefore,  $\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \lambda) \mathbf{y} > 0$ , which shows that the SOSC is satisfied.

An alternative (simpler) calculation:

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \lambda) \mathbf{y} = -2\mathbf{y}^\top \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} \mathbf{y} = -2(y_1(y_2 + y_3) + y_2(y_1 + y_3) + y_3(y_1 + y_2)).$$

Substituting  $y_1 = -(y_2 + y_3)$ ,  $y_2 = -(y_1 + y_3)$ , and  $y_3 = -(y_1 + y_2)$  in the first, second, and third terms, respectively, we obtain

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \lambda) \mathbf{y} = 2(y_1^2 + y_2^2 + y_3^2) > 0.$$

## 20.7

a. We first compute critical points by applying the Lagrange conditions. These are:

$$\begin{aligned} 2x_1 + 2\lambda x_1 &= 0 \\ 6x_1 + 2\lambda x_2 &= 0 \\ 1 + 2\lambda x_3 &= 0 \\ x_1^2 + x_2^2 + x_3^2 - 16 &= 0. \end{aligned}$$

There are six points satisfying the Lagrange condition:

$$\begin{aligned} \mathbf{x}^{(1)} &= [\sqrt{63}/2, 0, 1/2]^\top, & \lambda^{(1)} &= -1 \\ \mathbf{x}^{(2)} &= [-\sqrt{63}/2, 0, 1/2]^\top, & \lambda^{(2)} &= -1 \\ \mathbf{x}^{(3)} &= [0, 0, 4]^\top, & \lambda^{(3)} &= -1/8 \\ \mathbf{x}^{(4)} &= [0, 0, -4]^\top, & \lambda^{(4)} &= 1/8 \\ \mathbf{x}^{(5)} &= [0, \sqrt{575}/6, 1/6]^\top, & \lambda^{(5)} &= -3 \\ \mathbf{x}^{(6)} &= [0, -\sqrt{575}/6, 1/6]^\top, & \lambda^{(6)} &= -3. \end{aligned}$$

All the above points are regular. We now apply second order conditions to establish their nature. For this, we compute

$$\mathbf{F}(\mathbf{x}) = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad \mathbf{H}(\mathbf{x}) = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix},$$

and

$$T(\mathbf{x}^*) = \{\mathbf{y} \in \mathbb{R}^3 : [2x_1, 2x_2, 2x_3]\mathbf{y} = 0\}.$$

For the first point, we have

$$\begin{aligned} \mathbf{L}(\mathbf{x}^{(1)}, \lambda^{(1)}) &= \begin{bmatrix} 0 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & -2 \end{bmatrix} \\ T(\mathbf{x}^{(1)}) &= \{[-a/\sqrt{63}, b, a]^\top : a, b \in \mathbb{R}\}. \end{aligned}$$

Let  $\mathbf{y} = [-a/\sqrt{63}, b, a]^\top \in T(\mathbf{x}^{(1)})$ , where  $a$  and  $b$  are not both zero. Then

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^{(1)}, \lambda^{(1)}) \mathbf{y} = 4b^2 - 2a^2 \begin{cases} > 0 & \text{if } |a| < b\sqrt{2} \\ = 0 & \text{if } |a| = b\sqrt{2} \\ < 0 & \text{if } |a| > b\sqrt{2} \end{cases}.$$

From the above, we see that  $\mathbf{x}^{(1)}$  does not satisfy the SONC. Therefore,  $\mathbf{x}^{(1)}$  cannot be an extremizer. Performing similar calculations for  $\mathbf{x}^{(2)}$ , we conclude that  $\mathbf{x}^{(2)}$  cannot be an extremizer either.

For the third point, we have

$$\begin{aligned} \mathbf{L}(\mathbf{x}^{(3)}, \lambda^{(3)}) &= \begin{bmatrix} 7/4 & 0 & 0 \\ 0 & 23/4 & 0 \\ 0 & 0 & -1/4 \end{bmatrix} \\ T(\mathbf{x}^{(3)}) &= \{[a, b, 0]^\top : a, b \in \mathbb{R}\}. \end{aligned}$$

Let  $\mathbf{y} = [a, b, 0]^\top \in T(\mathbf{x}^{(3)})$ , where  $a$  and  $b$  are not both zero. Then

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^{(3)}, \lambda^{(3)}) \mathbf{y} = \frac{7}{4}a^2 + \frac{23}{4}b^2 > 0.$$

Hence,  $\mathbf{x}^{(3)}$  is a strict local minimizer. Performing similar calculations for the remaining points, we conclude that  $\mathbf{x}^{(4)}$  is a strict local minimizer, and  $\mathbf{x}^{(5)}$  and  $\mathbf{x}^{(6)}$  are both strict local maximizers.

b. The Lagrange condition for the problem is:

$$\begin{aligned} 2x_1 + \lambda(6x_1 + 4x_2) &= 0 \\ 2x_2 + \lambda(4x_1 + 12x_2) &= 0 \\ 3x_2^2 + 4x_1x_2 + 6x_2^2 - 140 &= 0. \end{aligned}$$

We represent the first two equations as

$$\begin{bmatrix} 2 + 6\lambda & 4\lambda \\ 4\lambda & 2 + 12\lambda \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

From the constraint equation, we note that  $\mathbf{x} = [0, 0]^\top$  cannot satisfy the Lagrange condition. Therefore, the determinant of the above matrix must be zero. Solving for  $\lambda$  yields two possible values:  $-1/7$  and  $-1/2$ . We then have four points satisfying the Lagrange condition:

$$\begin{aligned} \mathbf{x}^{(1)} &= [2, 4]^\top, & \lambda^{(1)} &= -1/7 \\ \mathbf{x}^{(2)} &= [-2, -4]^\top, & \lambda^{(2)} &= -1/7 \\ \mathbf{x}^{(3)} &= [-2\sqrt{14}, \sqrt{14}]^\top, & \lambda^{(3)} &= -1/2 \\ \mathbf{x}^{(4)} &= [2\sqrt{14}, -\sqrt{14}]^\top, & \lambda^{(4)} &= -1/2. \end{aligned}$$

Applying the SOSC, we conclude that  $\mathbf{x}^{(1)}$  and  $\mathbf{x}^{(2)}$  are strict local minimizers, and  $\mathbf{x}^{(3)}$  and  $\mathbf{x}^{(4)}$  are strict local maximizers.

## 20.8

a. We can represent the problem as

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}) \\ &\text{subject to} && h(\mathbf{x}) = \mathbf{0}, \end{aligned}$$

where  $f(\mathbf{x}) = 2x_1 + 3x_2 - 4$ , and  $h(\mathbf{x}) = x_1x_2 - 6$ . We have  $Df(\mathbf{x}) = [2, 3]$ , and  $Dh(\mathbf{x}) = [x_2, x_1]$ . Note that  $\mathbf{0}$  is not a feasible point. Therefore, any feasible point is regular. If  $\mathbf{x}^*$  is a local extremizer, then by the Lagrange multiplier theorem, there exists  $\lambda^* \in \mathbb{R}$  such that  $Df(\mathbf{x}^*) + \lambda^* Dh(\mathbf{x}^*) = \mathbf{0}^\top$ , or

$$\begin{aligned} 2 + \lambda^* x_2^* &= 0 \\ 3 + \lambda^* x_1^* &= 0. \end{aligned}$$

Solving, we get two possible extremizers:  $\mathbf{x}^{(1)} = [3, 2]^\top$ , with corresponding Lagrange multiplier  $\lambda^{(1)} = -1$ , and  $\mathbf{x}^{(2)} = [-3, 2]^\top$ , with corresponding Lagrange multiplier  $\lambda^{(2)} = 1$ .

b. We have  $\mathbf{F}(\mathbf{x}) = \mathbf{O}$ , and

$$\mathbf{H}(\mathbf{x}) = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}.$$

First, consider the point  $\mathbf{x}^{(1)} = [3, 2]^\top$ , with corresponding Lagrange multiplier  $\lambda^{(1)} = -1$ . We have

$$\mathbf{L}(\mathbf{x}^{(1)}, \lambda^{(1)}) = - \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix},$$

and

$$T(\mathbf{x}^{(1)}) = \{\mathbf{y} : [2, 3]\mathbf{y} = \mathbf{0}\} = \{\alpha[-3, 2]^\top : \alpha \in \mathbb{R}\}.$$

Let  $\mathbf{y} = \alpha[-3, 2]^\top \in T(\mathbf{x}^{(1)})$ ,  $\alpha \neq 0$ . We have

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^{(1)}, \lambda^{(1)}) \mathbf{y} = 12\alpha^2 > 0.$$



Therefore, by the SOSC,  $\mathbf{x}^{(1)} = [3, 2]^\top$  is a strict local minimizer.

Next, consider the point  $\mathbf{x}^{(2)} = -[3, 2]^\top$ , with corresponding Lagrange multiplier  $\lambda^{(2)} = 1$ . We have

$$\mathbf{L}(\mathbf{x}^{(2)}, \lambda^{(2)}) = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}.$$

and

$$T(\mathbf{x}^{(2)}) = \{\mathbf{y} : -[2, 3]\mathbf{y} = \mathbf{0}\} = \{\alpha[-3, 2]^\top : \alpha \in \mathbb{R}\} = T(\mathbf{x}^{(1)}).$$

Let  $\mathbf{y} = \alpha[-3, 2]^\top \in T(\mathbf{x}^{(2)})$ ,  $\alpha \neq 0$ . We have

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^{(2)}, \lambda^{(2)}) \mathbf{y} = -12\alpha^2 < 0.$$

Therefore, by the SOSC,  $\mathbf{x}^{(2)} = -[3, 2]^\top$  is a strict local maximizer.

c. Note that  $f(\mathbf{x}^{(1)}) = 8$ , while  $f(\mathbf{x}^{(2)}) = -16$ . Therefore,  $\mathbf{x}^{(1)}$ , although a strict local minimizer, is not a global minimizer. Likewise,  $\mathbf{x}^{(2)}$ , although a strict local maximizer, is not a global maximizer.

## 20.9

We observe that  $f(x_1, x_2)$  is a ratio of two quadratic functions, that is, we can represent  $f(x_1, x_2)$  as

$$f(x_1, x_2) = \frac{\mathbf{x}^\top \mathbf{Q} \mathbf{x}}{\mathbf{x}^\top \mathbf{P} \mathbf{x}}.$$

Therefore, if a point  $\mathbf{x}$  is a maximizer of  $f(x_1, x_2)$  then so is any nonzero multiple of this point because

$$\frac{(t\mathbf{x})^\top \mathbf{Q}(t\mathbf{x})}{(t\mathbf{x})^\top \mathbf{P}(t\mathbf{x})} = \frac{t^2 \mathbf{x}^\top \mathbf{Q} \mathbf{x}}{t^2 \mathbf{x}^\top \mathbf{P} \mathbf{x}} = \frac{\mathbf{x}^\top \mathbf{Q} \mathbf{x}}{\mathbf{x}^\top \mathbf{P} \mathbf{x}}.$$

Thus any nonzero multiple of a solution is also a solution. To proceed, represent the original problem in an equivalent form,

$$\begin{aligned} & \text{maximize} && \mathbf{x}^\top \mathbf{Q} \mathbf{x} = 18x_1^2 - 8x_1x_2 + 12x_2^2 \\ & \text{subject to} && \mathbf{x}^\top \mathbf{P} \mathbf{x} = 2x_1^2 + 2x_2^2 = 1. \end{aligned}$$

Thus, we wish to maximize  $f(x_1, x_2) = 18x_1^2 - 8x_1x_2 + 12x_2^2$  subject to the equality constraint,  $h(x_1, x_2) = 1 - 2x_1^2 - 2x_2^2 = 0$ . We apply the Lagrange's method to solve the problem. We form the Lagrangian function,

$$l(\mathbf{x}, \lambda) = f + \lambda h,$$

compute its gradient and find critical points. We have,

$$\begin{aligned} \nabla_{\mathbf{x}} l &= \nabla_{\mathbf{x}} \left( \mathbf{x}^\top \begin{bmatrix} 18 & -4 \\ -4 & 12 \end{bmatrix} \mathbf{x} + \lambda \left( 1 - \mathbf{x}^\top \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{x} \right) \right) \\ &= 2 \begin{bmatrix} 18 & -4 \\ -4 & 12 \end{bmatrix} \mathbf{x} - 2\lambda \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{x} \\ &= \mathbf{0}. \end{aligned}$$

We represent the above in an equivalent form,

$$\left( \lambda \mathbf{I}_2 - \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}^{-1} \begin{bmatrix} 18 & -4 \\ -4 & 12 \end{bmatrix} \right) \mathbf{x} = \mathbf{0}.$$

That is, solving the problem is being reduced to solving an eigenvalue-eigenvector problem,

$$\left( \lambda \mathbf{I}_2 - \begin{bmatrix} 9 & -2 \\ -2 & 6 \end{bmatrix} \right) \mathbf{x} = \begin{bmatrix} \lambda - 9 & 2 \\ 2 & \lambda - 6 \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

The characteristic polynomial is

$$\det \begin{bmatrix} \lambda - 9 & 2 \\ 2 & \lambda - 6 \end{bmatrix} = \lambda^2 - 15\lambda + 50 = (\lambda - 5)(\lambda - 10).$$

The eigenvalues are 5 and 10. Because we are interested in finding a maximizer, we conclude that the value of the maximized function is 10, while the corresponding maximizer corresponds to an appropriately scaled, to satisfy the constraint, eigenvector of this eigenvalue. An eigenvector can easily be found by taking any nonzero column of the adjoint matrix of

$$10\mathbf{I}_2 - \begin{bmatrix} 9 & -2 \\ -2 & 6 \end{bmatrix}.$$

Performing simple manipulations gives

$$\text{adj} \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} = \begin{bmatrix} 4 & -2 \\ -2 & 1 \end{bmatrix}.$$

Thus,

$$\sqrt{0.1} \begin{bmatrix} -2 \\ 1 \end{bmatrix}$$

is a maximizer for the equivalent problem. Any multiple of the above vector is a solution of the original maximization problem.

---

#### 20.10

We use the technique of Example 20.8. First, we write the objective function in the form  $\mathbf{x}^\top \mathbf{Q} \mathbf{x}$ , where

$$\mathbf{Q} = \mathbf{Q}^\top = \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix}.$$

The characteristic polynomial of  $\mathbf{Q}$  is  $\lambda^2 - 6\lambda + 5$ , and the eigenvalues of  $\mathbf{Q}$  are 1 and 5. The solutions to the problem are the unit length eigenvectors of  $\mathbf{Q}$  corresponding to the eigenvalue 5, which are  $\pm[1, 1]^\top/\sqrt{2}$ .

---

#### 20.11

Consider the problem

$$\begin{array}{ll} \text{minimize} & \|\mathbf{A}\mathbf{x}\|^2 \\ \text{subject to} & \|\mathbf{x}\|^2 = 1. \end{array}$$

The optimal objective function value of this problem is the smallest value that  $\|\mathbf{y}\|^2$  can take. The above can be solved easily using Lagrange multipliers. The Lagrange conditions are

$$\begin{array}{rcl} \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} - \lambda \mathbf{x}^\top & = & \mathbf{0}^\top \\ 1 - \mathbf{x}^\top \mathbf{x} & = & 0. \end{array}$$

The first equation can be rewritten as  $\mathbf{A}^\top \mathbf{A} \mathbf{x} = \lambda \mathbf{x}$ , which implies that  $\lambda$  is an eigenvalue of  $\mathbf{A}^\top \mathbf{A}$ . Moreover, premultiplying by  $\mathbf{x}^\top$  yields  $\mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} = \lambda \mathbf{x}^\top \mathbf{x} = \lambda$ , which indicates that the Lagrange multiplier is equal to the optimal objective function value. Hence, the range of values that  $\|\mathbf{y}\| = \|\mathbf{A}\mathbf{x}\|$  can take is 1 to  $\sqrt{20}$ .

---

#### 20.12

Consider the following optimization problem (we need to use squared norms to make the functions differentiable):

$$\begin{array}{ll} \text{minimize} & -\|\mathbf{A}\mathbf{x}\|^2 \\ \text{subject to} & \|\mathbf{x}\|^2 = 1. \end{array}$$

As usual, write  $f(\mathbf{x}) = -\|\mathbf{Ax}\|^2$  and  $h(\mathbf{x}) = \|\mathbf{x}\|^2 - 1$ . We have  $\nabla f(\mathbf{x}) = -2\mathbf{A}^\top \mathbf{Ax}$  and  $\nabla h(\mathbf{x}) = 2\mathbf{x}$ . Note that all feasible solutions are regular. Let  $\mathbf{x}^*$  be an optimal solution. Note that the optimal value of the objective function is  $f(\mathbf{x}^*) = -\|\mathbf{A}\|_2^2$ . The Lagrange condition for the above problem is:

$$\begin{aligned} -2\mathbf{A}^\top \mathbf{Ax}^* + \lambda^*(2\mathbf{x}^*) &= 0 \\ \|\mathbf{x}^*\|^2 &= 1. \end{aligned}$$

From the first equation, we see that

$$\mathbf{A}^\top \mathbf{Ax}^* = \lambda^* \mathbf{x}^*,$$

which implies that  $\lambda^*$  is an eigenvalue of  $\mathbf{A}^\top \mathbf{A}$ , and  $\mathbf{x}^*$  is the corresponding eigenvector. Premultiplying the above equation by  $\mathbf{x}^{*\top}$  and combining the result with the constraint equation, we obtain

$$\lambda^* = \mathbf{x}^{*\top} \mathbf{A}^\top \mathbf{Ax}^* = \|\mathbf{Ax}^*\|^2 = -f(\mathbf{x}^*) = \|\mathbf{A}\|_2^2.$$

Therefore, because  $\mathbf{x}^*$  minimizes  $f(\mathbf{x}^*)$ , we deduce that  $\lambda^*$  must be the largest eigenvalue of  $\mathbf{A}^\top \mathbf{A}$ ; i.e.,  $\lambda^* = \lambda_1$ . Therefore,

$$\|\mathbf{A}\|_2 = \sqrt{\lambda_1}.$$

---

### 20.13

Let  $h(\mathbf{x}) = 1 - \mathbf{x}^\top \mathbf{Px} = 0$ . Let  $\mathbf{x}_0$  be such that  $h(\mathbf{x}_0) = 0$ . Then,  $\mathbf{x}_0 \neq \mathbf{0}$ . For  $\mathbf{x}_0$  to be a regular point, we need to show that  $\{\nabla h(\mathbf{x}_0)\}$  is a linearly independent set, i.e.,  $\nabla h(\mathbf{x}_0) \neq \mathbf{0}$ . Now,  $\nabla h(\mathbf{x}) = -2\mathbf{Px}$ . Since  $\mathbf{P}$  is nonsingular, and  $\mathbf{x}_0 \neq \mathbf{0}$ , then  $\nabla h(\mathbf{x}_0) = -2\mathbf{Px}_0 \neq \mathbf{0}$ .

---

### 20.14

Note that the point  $[1, 1]^\top$  is a regular point. Applying the Lagrange multiplier theorem gives

$$\begin{aligned} a + 2\lambda^* &= 0 \\ b + 2\lambda^* &= 0. \end{aligned}$$

Hence,  $a = b$ .

---

### 20.15

a. Denote the solution by  $[x_1^*, x_2^*]^\top$ . The Lagrange condition for this problem has the form

$$\begin{aligned} x_2^* - 2 + 2\lambda^* x_1^* &= 0 \\ x_1^* - 2\lambda^* x_2^* &= 0 \\ (x_1^*)^2 - (x_2^*)^2 &= 0. \end{aligned}$$

From the first and third equations it follows that  $x_1^*, x_2^* \neq 0$ . Then, combining the first and second equations, we obtain

$$\lambda^* = \frac{2 - x_2^*}{2x_1^*} = \frac{x_1^*}{2x_2^*}$$

which implies that  $2x_2^* - (x_2^*)^2 = (x_1^*)^2$ . Hence,  $x_2^* = 1$ , and by the third Lagrange equation,  $(x_1^*)^2 = 1$ . Thus, the only two points satisfying the Lagrange condition are  $[1, 1]^\top$  and  $[-1, 1]^\top$ . Note that both points are regular.

b. Consider the point  $\mathbf{x}^* = [-1, 1]^\top$ . The corresponding Lagrange multiplier is  $\lambda^* = -1/2$ . The Hessian of the Lagrangian is

$$\mathbf{L}(\mathbf{x}^*, \lambda^*) = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix} = \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix}.$$

The tangent plane is given by

$$T(\mathbf{x}^*) = \{\mathbf{y} : [-2, -2]\mathbf{y} = 0\} = \{[a, -a]^\top : a \in \mathbb{R}\}.$$

Let  $\mathbf{y} \in T(\mathbf{x}^*)$ ,  $\mathbf{y} \neq \mathbf{0}$ . Then,  $\mathbf{y} = [a, -a]^\top$  for some  $a \neq 0$ . We have  $\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \lambda^*) \mathbf{y} = -2a^2 < 0$ . Hence, SONC does not hold in this case, and therefore  $\mathbf{x}^* = [-1, 1]^\top$  cannot be local minimizer. In fact, the point is a strict local maximizer.

c. Consider the point  $\mathbf{x}^* = [1, 1]^\top$ . The corresponding Lagrange multiplier is  $\lambda^* = 1/2$ . The Hessian of the Lagrangian is

$$\mathbf{L}(\mathbf{x}^*, \lambda^*) = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

The tangent plane is given by

$$T(\mathbf{x}^*) = \{\mathbf{y} : [2, -2]\mathbf{y} = 0\} = \{[a, a]^\top : a \in \mathbb{R}\}.$$

Let  $\mathbf{y} \in T(\mathbf{x}^*)$ ,  $\mathbf{y} \neq \mathbf{0}$ . Then,  $\mathbf{y} = [a, a]$  for some  $a \neq 0$ . We have  $\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \lambda^*) \mathbf{y} = 2a^2 > 0$ . Hence, by the SOSC, the point  $\mathbf{x}^* = [1, 1]^\top$  is a strict local minimizer.

#### 20.16

a. The point  $\mathbf{x}^*$  is the solution to the optimization problem

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{0}. \end{aligned}$$

Since  $\text{rank } \mathbf{A} = m$ , any feasible point is regular. By the Lagrange multiplier theorem, there exists  $\boldsymbol{\lambda}^* \in \mathbb{R}^m$  such that

$$(\mathbf{x}^* - \mathbf{x}_0)^\top - \boldsymbol{\lambda}^{*\top} \mathbf{A} = \mathbf{0}^\top.$$

Postmultiplying both sides by  $\mathbf{x}^*$  and using the fact that  $\mathbf{A}\mathbf{x}^* = \mathbf{0}$ , we get

$$(\mathbf{x}^* - \mathbf{x}_0)^\top \mathbf{x}^* = 0.$$

b. From part a, we have

$$\mathbf{x}^* - \mathbf{x}_0 = \mathbf{A}^\top \boldsymbol{\lambda}^*.$$

Premultiplying both sides by  $\mathbf{A}$  we get

$$-\mathbf{A}\mathbf{x}_0 = (\mathbf{A}\mathbf{A}^\top) \boldsymbol{\lambda}^*$$

from which we conclude that  $\boldsymbol{\lambda}^* = -(\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{x}_0$ . Hence,

$$\mathbf{x}^* = \mathbf{x}_0 + \mathbf{A}^\top \boldsymbol{\lambda}^* = \mathbf{x}_0 - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{x}_0 = (\mathbf{I}_n - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}) \mathbf{x}_0.$$

#### 20.17

a. The Lagrange condition is (omitting all superscript-\* for convenience):

$$\begin{aligned} (\mathbf{A}\mathbf{x} - \mathbf{b})^\top \mathbf{A} + \boldsymbol{\lambda}^\top \mathbf{C} &= \mathbf{0}^\top \\ \mathbf{C}\mathbf{x} &= \mathbf{d}. \end{aligned}$$

For simplicity, write  $\mathbf{Q} = \mathbf{A}^\top \mathbf{A}$ , which is positive definite. From the first equation, we have

$$\mathbf{x} = \mathbf{Q}^{-1} \mathbf{A}^\top \mathbf{b} - \mathbf{Q}^{-1} \mathbf{C}^\top \boldsymbol{\lambda}.$$

Multiplying both sides by  $\mathbf{C}$  and using the second equation, we have

$$\mathbf{d} = \mathbf{C}\mathbf{Q}^{-1} \mathbf{A}^\top \mathbf{b} - \mathbf{C}\mathbf{Q}^{-1} \mathbf{C}^\top \boldsymbol{\lambda},$$

from which we obtain

$$\boldsymbol{\lambda} = (\mathbf{C}\mathbf{Q}^{-1} \mathbf{C}^\top)^{-1} (\mathbf{C}\mathbf{Q}^{-1} \mathbf{A}^\top \mathbf{b} - \mathbf{d}).$$

Substituting back into the equation for  $\mathbf{x}$ , we obtain

$$\mathbf{x} = \mathbf{Q}^{-1} \mathbf{A}^\top \mathbf{b} - \mathbf{Q}^{-1} \mathbf{C}^\top (\mathbf{C}\mathbf{Q}^{-1} \mathbf{C}^\top)^{-1} (\mathbf{C}\mathbf{Q}^{-1} \mathbf{A}^\top \mathbf{b} - \mathbf{d}).$$

b. Rewrite the objective function as

$$\frac{1}{2} \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} - \mathbf{b}^\top \mathbf{A} \mathbf{x} + \frac{1}{2} \|\mathbf{b}\|^2.$$

As before, write  $\mathbf{Q} = \mathbf{A}^\top \mathbf{A}$ . Completing the squares and setting  $\mathbf{y} = \mathbf{x} - \mathbf{Q}^{-1} \mathbf{A}^\top \mathbf{b}$ , the objective function can be written as

$$\frac{1}{2} \mathbf{y}^\top \mathbf{Q} \mathbf{y} + \text{const.}$$

Hence, the problem can be converted to the equivalent QP:

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \mathbf{y}^\top \mathbf{Q} \mathbf{y} \\ & \text{subject to} && \mathbf{C} \mathbf{y} = \mathbf{d} - \mathbf{C} \mathbf{Q}^{-1} \mathbf{A}^\top \mathbf{b}. \end{aligned}$$

The solution to this QP is

$$\mathbf{y}^* = \mathbf{Q}^{-1} \mathbf{C}^\top (\mathbf{C} \mathbf{Q}^{-1} \mathbf{C}^\top)^{-1} (\mathbf{d} - \mathbf{C} \mathbf{Q}^{-1} \mathbf{A}^\top \mathbf{b}).$$

Hence, the solution to the original problem is:

$$\begin{aligned} \mathbf{x}^* &= \mathbf{Q}^{-1} \mathbf{A}^\top \mathbf{b} + \mathbf{Q}^{-1} \mathbf{C}^\top (\mathbf{C} \mathbf{Q}^{-1} \mathbf{C}^\top)^{-1} (\mathbf{d} - \mathbf{C} \mathbf{Q}^{-1} \mathbf{A}^\top \mathbf{b}) \\ &= (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} + (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{C}^\top (\mathbf{C} (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{C}^\top)^{-1} (\mathbf{d} - \mathbf{C} (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}), \end{aligned}$$

which agrees with the solution obtained in part a.

---

#### 20.18

Write  $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{c}^\top \mathbf{x} + d$  (actually, we could have ignored  $d$ ) and  $\mathbf{h}(\mathbf{x}) = \mathbf{b} - \mathbf{A} \mathbf{x}$ . We have

$$Df(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q} - \mathbf{c}^\top, \quad Dh(\mathbf{x}) = -\mathbf{A}.$$

The Lagrange condition is

$$\begin{aligned} \mathbf{x}^{*\top} \mathbf{Q} - \mathbf{c}^\top - \boldsymbol{\lambda}^{*\top} \mathbf{A} &= \mathbf{0}^\top \\ \mathbf{b} - \mathbf{A} \mathbf{x}^* &= \mathbf{0}. \end{aligned}$$

From the first equation we get

$$\mathbf{x}^* = \mathbf{Q}^{-1} (\mathbf{A}^\top \boldsymbol{\lambda}^* + \mathbf{c}).$$

Multiplying both sides by  $\mathbf{A}$  and using the second equation (constraint), we get

$$(\mathbf{A} \mathbf{Q}^{-1} \mathbf{A}^\top) \boldsymbol{\lambda}^* + \mathbf{A} \mathbf{Q}^{-1} \mathbf{c} = \mathbf{b}.$$

Since  $\mathbf{Q} > 0$  and  $\mathbf{A}$  is of full rank, we can write

$$\boldsymbol{\lambda}^* = (\mathbf{A} \mathbf{Q}^{-1} \mathbf{A}^\top)^{-1} (\mathbf{b} - \mathbf{A} \mathbf{Q}^{-1} \mathbf{c}).$$

Hence,

$$\mathbf{x}^* = \mathbf{Q}^{-1} \mathbf{c} + \mathbf{Q}^{-1} \mathbf{A}^\top (\mathbf{A} \mathbf{Q}^{-1} \mathbf{A}^\top)^{-1} (\mathbf{b} - \mathbf{A} \mathbf{Q}^{-1} \mathbf{c}).$$

Alternatively, we could have rewritten the given problem in our usual quadratic programming form with variable  $\mathbf{y} = \mathbf{x} - \mathbf{Q}^{-1} \mathbf{c}$ .

---

#### 20.19

Clearly, we have  $\mathcal{M} = \mathcal{R}(\mathbf{B})$ , i.e.,  $\mathbf{y} \in \mathcal{M}$  if and only if there exists  $\mathbf{x} \in \mathbb{R}^m$  such that  $\mathbf{y} = \mathbf{B} \mathbf{x}$ . Hence

$$\begin{aligned} & \mathbf{L} \text{ is positive semidefinite on } \mathcal{M} \\ \Leftrightarrow & \text{ for all } \mathbf{y} \in \mathcal{M}, \quad \mathbf{y}^\top \mathbf{L} \mathbf{y} \geq 0 \\ \Leftrightarrow & \text{ for all } \mathbf{x} \in \mathbb{R}^m, \quad (\mathbf{B} \mathbf{x})^\top \mathbf{L} (\mathbf{B} \mathbf{x}) \geq 0 \\ \Leftrightarrow & \text{ for all } \mathbf{x} \in \mathbb{R}^m, \quad \mathbf{x}^\top (\mathbf{B}^\top \mathbf{L} \mathbf{B}) \mathbf{x} \geq 0 \\ \Leftrightarrow & \text{ for all } \mathbf{x} \in \mathbb{R}^m, \quad \mathbf{x}^\top \mathbf{L}_{\mathcal{M}} \mathbf{x} \geq 0 \\ \Leftrightarrow & \mathbf{L}_{\mathcal{M}} \geq 0. \end{aligned}$$

For positive definiteness, the same argument applies, with  $\geq$  replaced by  $>$ .

## 20.20

a. By simple manipulations, we can write

$$x_2 = a^2 x_0 + abu_0 + bu_1.$$

Therefore, the problem is

$$\begin{aligned} & \text{minimize} && \frac{1}{2}(u_0^2 + u_1^2) \\ & \text{subject to} && a^2 x_0 + abu_0 + bu_1 = 0. \end{aligned}$$

Alternatively, we may use a vector notation: writing  $\mathbf{u} = [u_0, u_1]^\top$ , we have

$$\begin{aligned} & \text{minimize} && f(\mathbf{u}) \\ & \text{subject to} && h(\mathbf{u}) = 0, \end{aligned}$$

where  $f(\mathbf{u}) = \frac{1}{2}\|\mathbf{u}\|^2$ , and  $h(\mathbf{u}) = a^2 x_0 + [ab, b]\mathbf{u}$ . Since the vector  $\nabla h(\mathbf{u}) = [ab, b]^\top$  is nonzero for any  $\mathbf{u}$ , then any feasible point is regular. Therefore, by the Lagrange multiplier theorem, there exists  $\lambda^* \in \mathbb{R}$  such that

$$\begin{aligned} u_0^* + \lambda^* ab &= 0 \\ u_1^* + \lambda^* b &= 0 \\ a^2 x_0 + abu_0^* + bu_1^* &= 0. \end{aligned}$$

We have three linear equations in three unknowns, that upon solving yields

$$u_0^* = -\frac{a^3 x_0}{b(1+a^2)}, \quad u_1^* = -\frac{a^2 x_0}{b(1+a^2)}.$$

b. The Hessians of  $f$  and  $h$  are  $\mathbf{F}(\mathbf{u}) = \mathbf{I}_2$  ( $2 \times 2$  identity matrix) and  $\mathbf{H}(\mathbf{u}) = \mathbf{O}$ , respectively. Hence, the Hessian of the Lagrangian is  $\mathbf{L}(\mathbf{u}^*, \lambda^*) = \mathbf{I}_2$ , which is positive definite. Therefore,  $\mathbf{u}^*$  satisfies the SOSC, and is therefore a strict local minimizer.

## 20.21

Letting  $\mathbf{z} = [x_2, u_1, u_2]^\top$ , the objective function is  $\mathbf{z}^\top \mathbf{Q} \mathbf{z}$ , where

$$\mathbf{Q} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1/2 & 0 \\ 0 & 0 & 1/3 \end{bmatrix}.$$

The linear constraint on  $\mathbf{z}$  is obtained by writing

$$x_2 = 2x_1 + u_2 = 2(2 + u_1) + u_2,$$

which can be written as  $\mathbf{A} \mathbf{z} = b$ , where

$$\mathbf{A} = [1, -2, -1], \quad b = 4.$$

Hence, using the method of Section 20.6, the solution is

$$\mathbf{z}^* = \mathbf{Q}^{-1} \mathbf{A}^\top (\mathbf{A} \mathbf{Q}^{-1} \mathbf{A}^\top)^{-1} b = \begin{bmatrix} 1 \\ -4 \\ -3 \end{bmatrix} \cdot (12)^{-1} \cdot 4 = \begin{bmatrix} 1/3 \\ -4/3 \\ -1 \end{bmatrix}.$$

Thus, the optimal controls are  $u_1^* = -4/3$  and  $u_2^* = -1$ .

**20.22**

The composite input vector is

$$\mathbf{u} = \begin{bmatrix} u_0 & u_1 & u_2 \end{bmatrix}^\top.$$

The performance index  $J$  is  $J = \frac{1}{2} \mathbf{u}^\top \mathbf{u}$ . To obtain the constraint  $\mathbf{A}\mathbf{u} = b$ , where  $\mathbf{A} \in \mathbb{R}^{1 \times 3}$ , we proceed as follows. First, we write

$$\begin{aligned} x_2 &= x_1 + 2u_1 \\ &= x_0 + 2u_0 + 2u_1. \end{aligned}$$

Using the above, we obtain

$$\begin{aligned} x_3 &= 9 \\ &= x_2 + 2u_2 \\ &= x_0 + 2u_0 + 2u_1 + 2u_2. \end{aligned}$$

We represent the above in the format  $\mathbf{A}\mathbf{u} = b$  as follows

$$\begin{bmatrix} 2 & 2 & 2 \end{bmatrix} \begin{bmatrix} u_0 \\ u_1 \\ u_2 \end{bmatrix} = 6.$$

Thus we formulated the problem of finding the optimal control sequence as a constrained optimization problem

$$\begin{aligned} &\text{minimize} \quad \frac{1}{2} \mathbf{u}^\top \mathbf{u} \\ &\text{subject to} \quad \mathbf{A}\mathbf{u} = b. \end{aligned}$$

To solve the above problem, we form the Lagrangian

$$l(\mathbf{u}, \lambda) = \frac{1}{2} \mathbf{u}^\top \mathbf{u} + \lambda (\mathbf{A}\mathbf{u} - b),$$

where  $\lambda$  is the Lagrange multiplier. Applying the Lagrange first-order condition yields

$$\mathbf{u} + \mathbf{A}^\top \lambda = \mathbf{0} \quad \text{and} \quad \mathbf{A}\mathbf{u} = b.$$

From the first of the above conditions, we calculate,  $\mathbf{u} = -\mathbf{A}^\top \lambda$ . Substituting the above into the second of the Lagrange conditions gives

$$\lambda = -\left(\mathbf{A}\mathbf{A}^\top\right)^{-1} b.$$

Combining the last two equations, we obtain a closed-form formula for the optimal input sequence

$$\mathbf{u} = \mathbf{A}^\top \left(\mathbf{A}\mathbf{A}^\top\right)^{-1} b.$$

In our problem,

$$\mathbf{u} = \begin{bmatrix} u_0 \\ u_1 \\ u_2 \end{bmatrix} = \frac{b}{\left(\mathbf{A}\mathbf{A}^\top\right)} \mathbf{A}^\top = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

## 21. Problems With Inequality Constraints

## 21.1

---

a. We form the Lagrangian function,

$$l(\mathbf{x}, \mu) = x_1^2 + 4x_2^2 + \mu(4 - x_1^2 - 2x_2^2).$$

The KKT conditions take the form,

$$\begin{aligned} D_{\mathbf{x}}l(\mathbf{x}, \mu) &= \begin{bmatrix} 2x_1 - 2\mu x_1 & 8x_2 - 4\mu x_2 \end{bmatrix} = \mathbf{0}^\top \\ \mu(4 - x_1^2 - 2x_2^2) &= 0 \\ \mu &\geq 0 \\ 4 - x_1^2 - 2x_2^2 &\leq 0. \end{aligned}$$

From the first of the above equality, we obtain

$$(1 - \mu)x_1 = 0 \quad (2 - \mu)x_2 = 0.$$

We first consider the case when  $\mu = 0$ . Then, we obtain the point,  $\mathbf{x}^{(1)} = \mathbf{0}$ , which does not satisfy the constraints.

The next case is when  $\mu = 1$ . Then we have to have  $x_2 = 0$  and using  $\mu(4 - x_1^2 - 2x_2^2) = 0$  gives

$$\mathbf{x}^{(2)} = \begin{bmatrix} 2 \\ 0 \end{bmatrix} \quad \text{and} \quad \mathbf{x}^{(3)} = -\begin{bmatrix} 2 \\ 0 \end{bmatrix}.$$

For the case when  $\mu = 2$ , we have to have  $x_1 = 0$  and we get

$$\mathbf{x}^{(4)} = \begin{bmatrix} 0 \\ \sqrt{2} \end{bmatrix} \quad \text{and} \quad \mathbf{x}^{(5)} = -\begin{bmatrix} 0 \\ \sqrt{2} \end{bmatrix}.$$

b. The Hessian of  $l$  is

$$\mathbf{L} = \begin{bmatrix} 2 & 0 \\ 0 & 8 \end{bmatrix} + \mu \begin{bmatrix} -2 & 0 \\ 0 & -4 \end{bmatrix}.$$

When  $\mu = 1$ ,

$$\mathbf{L} = \begin{bmatrix} 0 & 0 \\ 0 & 4 \end{bmatrix}.$$

We next find the subspace

$$\begin{aligned} \tilde{T} &= T = \{\mathbf{y} : \begin{bmatrix} \pm 4 & 0 \end{bmatrix} \mathbf{y} = 0\} \\ &= \{\mathbf{y} = a \begin{bmatrix} 0 & 1 \end{bmatrix}^\top : a \in \mathbb{R}\}. \end{aligned}$$

We then check for positive definiteness of  $\mathbf{L}$  on  $\tilde{T}$ ,

$$\mathbf{y}^\top \mathbf{L} \mathbf{y} = a^2 \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & 4 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = 4a^2 > 0.$$

Hence,  $\mathbf{x}^{(2)}$  and  $\mathbf{x}^{(3)}$  satisfy the SOSC to be strict local minimizers.

When  $\mu = 2$ ,

$$\mathbf{L} = \begin{bmatrix} -2 & 0 \\ 0 & 0 \end{bmatrix},$$

and

$$T = \{\mathbf{y} = a \begin{bmatrix} 1 & 0 \end{bmatrix}^\top : a \in \mathbb{R}\}.$$



We have

$$\mathbf{y}^\top \mathbf{L} \mathbf{y} = -2a^2 < 0.$$

Thus,  $\mathbf{x}^{(4)}$  and  $\mathbf{x}^{(5)}$  do not satisfy the SONC to be minimizers.

In summary, only  $\mathbf{x}^{(2)}$  and  $\mathbf{x}^{(3)}$  are strict local minimizers.

## 21.2

a. We first find critical points by applying the Karush-Kuhn-Tucker conditions, which are

$$\begin{aligned} 2x_1 - 2 - 2\mu_1 x_1 + 5\mu_2 &= 0 \\ 2x_2 - 10 + \frac{1}{5}\mu_1 + \frac{1}{2}\mu_2 &= 0 \\ \mu_1 \left( \frac{1}{5}x_2 - x_1^2 \right) + \mu_2 \left( 5x_1 + \frac{1}{2}x_2 - 5 \right) &= 0 \\ \boldsymbol{\mu} &\geq 0. \end{aligned}$$

We have to check four possible combinations.

**Case 1:** ( $\mu_1 = 0, \mu_2 = 0$ ) Solving the first and second Karush-Kuhn-Tucker equations yields  $\mathbf{x}^{(1)} = [1, 5]^\top$ . However, this point is not feasible and is therefore not a candidate minimizer.

**Case 2:** ( $\mu_1 > 0, \mu_2 = 0$ ) We have two possible solutions:

$$\begin{aligned} \mathbf{x}^{(2)} &= [-0.98, 4.8]^\top & \mu_1^{(2)} &= 2.02 \\ \mathbf{x}^{(3)} &= [-0.02, 0]^\top & \mu_1^{(3)} &= 50. \end{aligned}$$

Both  $\mathbf{x}^{(2)}$  and  $\mathbf{x}^{(3)}$  satisfy the constraints, and are therefore candidate minimizers.

**Case 3:** ( $\mu_1 = 0, \mu_2 > 0$ ) Solving the corresponding equation yields:

$$\mathbf{x}^{(4)} = [0.5050, 4.9505]^\top \quad \mu_1^{(4)} = 0.198.$$

The point  $\mathbf{x}^{(4)}$  is not feasible, and hence is not a candidate minimizer.

**Case 4:** ( $\mu_1 > 0, \mu_2 > 0$ ) We have two solutions:

$$\begin{aligned} \mathbf{x}^{(5)} &= [0.732, 2.679]^\top & \boldsymbol{\mu}^{(5)} &= [13.246, 3.986]^\top \\ \mathbf{x}^{(6)} &= [-2.73, 37.32]^\top & \boldsymbol{\mu}^{(6)} &= [188.8, -204]^\top. \end{aligned}$$

The point  $\mathbf{x}^{(5)}$  is feasible, but  $\mathbf{x}^{(6)}$  is not.

We are left with three candidate minimizers:  $\mathbf{x}^{(2)}$ ,  $\mathbf{x}^{(3)}$ , and  $\mathbf{x}^{(5)}$ . It is easy to check that they are regular. We now check if each satisfies the second order conditions. For this, we compute

$$\mathbf{L}(\mathbf{x}, \boldsymbol{\mu}) = \begin{bmatrix} 2 - 2\mu_1 & 0 \\ 0 & 2 \end{bmatrix}.$$

For  $\mathbf{x}^{(2)}$ , we have

$$\begin{aligned} \mathbf{L}(\mathbf{x}^{(2)}, \boldsymbol{\mu}^{(2)}) &= \begin{bmatrix} -2.04 & 0 \\ 0 & 2 \end{bmatrix} \\ \tilde{T}(\mathbf{x}^{(2)}) &= \{a[-0.1021, 1]^\top : a \in \mathbb{R}\}. \end{aligned}$$

Let  $\mathbf{y} = a[-0.1021, 1]^\top \in \tilde{T}(\mathbf{x}^{(2)})$  with  $a \neq 0$ . Then

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^{(2)}, \boldsymbol{\mu}^{(2)}) \mathbf{y} = 1.979a^2 > 0.$$

Thus, by the SOSC,  $\mathbf{x}^{(2)}$  is a strict local minimizer.

For  $\mathbf{x}^{(3)}$ , we have

$$\begin{aligned} \mathbf{L}(\mathbf{x}^{(3)}, \boldsymbol{\mu}^{(3)}) &= \begin{bmatrix} -97.958 & 0 \\ 0 & 2 \end{bmatrix} \\ \tilde{T}(\mathbf{x}^{(3)}) &= \{a[-4.898, 1]^\top : a \in \mathbb{R}\}. \end{aligned}$$

Let  $\mathbf{y} = a[-4.898, 1]^\top \in \tilde{T}(\mathbf{x}^{(3)})$  with  $a \neq 0$ . Then,

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^{(3)}, \boldsymbol{\mu}^{(3)}) \mathbf{y} = -2347.9a^2 < 0.$$

Thus,  $\mathbf{x}^{(3)}$  does not satisfy the SOSC. In fact, in this case, we have  $T(\mathbf{x}^{(3)}) = \tilde{T}(\mathbf{x}^{(3)})$ , and hence  $\mathbf{x}^{(3)}$  does not satisfy the SONC either. We conclude that  $\mathbf{x}^{(3)}$  is not a local minimizer. We can easily check that  $\mathbf{x}^{(3)}$  is not a local maximizer either.

For  $\mathbf{x}^{(5)}$ , we have

$$\begin{aligned} \mathbf{L}(\mathbf{x}^{(5)}, \boldsymbol{\mu}^{(5)}) &= \begin{bmatrix} -24.4919 & 0 \\ 0 & 2 \end{bmatrix} \\ \tilde{T}(\mathbf{x}^{(5)}) &= \{\mathbf{0}\}. \end{aligned}$$

The SOSC is trivially satisfied, and therefore  $\mathbf{x}^{(5)}$  is a strict local minimizer.

b. The Karush-Kuhn-Tucker conditions are:

$$\begin{aligned} 2x_1 - \mu_1 - \mu_3 &= 0 \\ 2x_2 - \mu_2 - \mu_3 &= 0 \\ -x_1 &\leq 0 \\ -x_2 &\leq 0 \\ -x_1 - x_2 + 5 &\leq 0 \\ -\mu_1 x_1 - \mu_2 x_2 + \mu_3(-x_1 - x_2 + 5) &= 0 \\ \mu_1, \mu_2, \mu_3 &\geq 0. \end{aligned}$$

It is easy to verify that the only combination of Karush-Kuhn-Tucker multipliers resulting in a feasible point is  $\mu_1 = \mu_2 = 0$ ,  $\mu_3 > 0$ . For this case, we obtain  $\mathbf{x}^* = [2.5, 2.5]^\top$ ,  $\boldsymbol{\mu}^* = [0, 0, 5]^\top$ . We have

$$\mathbf{L}(\mathbf{x}^*, \boldsymbol{\mu}^*) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} > 0.$$

Hence,  $\mathbf{x}^*$  is a strict local minimizer (in fact, the only one for this problem).

c. The Karush-Kuhn-Tucker conditions are:

$$\begin{aligned} 2x_1 + 6x_2 - 4 + 2\mu_1 x_1 + 2\mu_2 &= 0 \\ 6x_1 - 2 + 2\mu_1 - 2\mu_2 &= 0 \\ x_1^2 + 2x_2 - 1 &\leq 0 \\ 2x_1 - 2x_2 - 1 &\leq 0 \\ \mu_1(x_1^2 + 2x_2 - 1) + \mu_2(2x_1 - 2x_2 - 1) &= 0 \\ \mu_1, \mu_2 &\geq 0. \end{aligned}$$

It is easy to verify that the only combination of Karush-Kuhn-Tucker multipliers resulting in a feasible point is  $\mu_1 = 0$ ,  $\mu_2 > 0$ . For this case, we obtain  $\mathbf{x}^* = [9/14, 2/14]^\top$ ,  $\boldsymbol{\mu}^* = [0, 13/14]^\top$ . We have

$$\begin{aligned} \mathbf{L}(\mathbf{x}^*, \boldsymbol{\mu}^*) &= \begin{bmatrix} 2 & 6 \\ 6 & 0 \end{bmatrix} \\ \tilde{T}(\mathbf{x}^*) &= \{a[1, 1] : a \in \mathbb{R}\}. \end{aligned}$$

Let  $\mathbf{y} = a[1, 1] \in \tilde{T}(\mathbf{x}^*)$  with  $a \neq 0$ . Then

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \boldsymbol{\mu}^*) \mathbf{y} = 14a^2 > 0.$$

Hence,  $\mathbf{x}^*$  is a strict local minimizer (in fact, the only one for this problem).

### 21.3

The Karush-Kuhn-Tucker conditions are:

$$\begin{aligned}
2x_1 + 2\lambda x_1 + 2\lambda x_2 + 2\mu x_1 &= 0 \\
2x_2 + 2\lambda x_1 + 2\lambda x_2 - \mu &= 0 \\
x_1^2 + 2x_1x_2 + x_2^2 - 1 &= 0 \\
x_1^2 - x_2 &\leq 0 \\
\mu(x_1^2 - x_2) &= 0 \\
\mu &\geq 0.
\end{aligned}$$

We have two cases to consider.

**Case 1:** ( $\mu > 0$ ) Substituting  $x_2 = x_1^2$  into the third equation and combining the result with the first two yields two possible points:

$$\begin{aligned}
\mathbf{x}^{(1)} &= [-1.618, 0.618]^\top & \mu^{(1)} &= -3.7889 \\
\mathbf{x}^{(2)} &= [2.618, 0.382]^\top & \mu^{(2)} &= -0.2111.
\end{aligned}$$

Note that the resulting  $\mu$  values violate the condition  $\mu > 0$ . Hence, neither of the points are minimizers (although they are candidates for maximizers).

**Case 2:** ( $\mu = 0$ ) Subtracting the second equation from the first yields  $x_1 = x_2$ , which upon substituting into the third equation gives two possible points:

$$\mathbf{x}^{(3)} = [1/2, 1/2]^\top, \quad \mathbf{x}^{(4)} = [-1/2, 1/2]^\top.$$

Note that  $\mathbf{x}^{(4)}$  is not a feasible point, and is therefore not a candidate minimizer.

Therefore, the only remaining candidate is  $\mathbf{x}^{(3)}$ , with corresponding  $\lambda^{(3)} = -1/2$  (and  $\mu = 0$ ). We now check second order conditions. We have

$$\begin{aligned}
\mathbf{L}(\mathbf{x}^{(3)}, 0, \lambda^{(3)}) &= \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \\
\tilde{T}(\mathbf{x}^{(3)}) &= \{a[1, -1]^\top : a \in \mathbb{R}\}.
\end{aligned}$$

Let  $\mathbf{y} = a[1, -1]^\top \in \tilde{T}(\mathbf{x}^{(3)})$  with  $a \neq 0$ . Then

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^{(3)}, 0, \lambda^{(3)}) \mathbf{y} = 4a^2 > 0.$$

Therefore, by the SOSC,  $\mathbf{x}^{(3)}$  is a strict local minimizer.

### 21.4

The optimization problem is:

$$\begin{aligned}
&\text{minimize} && \mathbf{e}_n^\top \mathbf{p} \\
&\text{subject to} && \mathbf{G}\mathbf{p} \geq P\mathbf{e}_m \\
&&& \mathbf{p} \geq \mathbf{0},
\end{aligned}$$

where  $\mathbf{G} = [g_{i,j}]$ ,  $\mathbf{e}_n = [1, \dots, 1]^\top$  (with  $n$  components), and  $\mathbf{p} = [p_1, \dots, p_n]^\top$ . The KKT condition for this problem is:

$$\begin{aligned}
\mathbf{e}_n^\top - \mu_1^\top \mathbf{G} - \mu_2^\top &= \mathbf{0}^\top \\
\mu_1^\top (P\mathbf{e}_m - \mathbf{G}\mathbf{p}) - \mu_2^\top \mathbf{p} &= 0 \\
\mathbf{G}\mathbf{p} &\geq P\mathbf{e}_m \\
\mu_1, \mu_2, \mathbf{p} &\geq \mathbf{0}.
\end{aligned}$$

**21.5**

a. We have  $f(\mathbf{x}) = x_2 - (x_1 - 2)^3 + 3$  and  $g(\mathbf{x}) = 1 - x_2$ . Hence,  $\nabla f(\mathbf{x}) = [-3(x_1 - 2)^2, 1]^\top$  and  $\nabla g(\mathbf{x}) = [0, -1]^\top$ . The KKT condition is

$$\begin{aligned}\mu &\geq 0 \\ -3(x_1 - 2)^2 &= 0 \\ 1 - \mu &= 0 \\ \mu(1 - x_2) &= 0 \\ 1 - x_2 &\leq 0.\end{aligned}$$

The only solution to the above conditions is  $\mathbf{x}^* = [2, 1]^\top$ ,  $\mu^* = 1$ .

To check if  $\mathbf{x}^*$  is regular, we note that the constraint is active. We have  $\nabla g(\mathbf{x}^*) = [0, -1]^\top$ , which is nonzero. Hence,  $\mathbf{x}^*$  is regular.

b. We have

$$\mathbf{L}(\mathbf{x}^*, \mu^*) = \mathbf{F}(\mathbf{x}^*) + \mu^* \mathbf{G}(\mathbf{x}^*) = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}.$$

Hence, the point  $\mathbf{x}^*$  satisfies the SONC.

c. Since  $\mu^* > 0$ , we have  $\tilde{T}(\mathbf{x}^*, \mu^*) = T(\mathbf{x}^*) = \{\mathbf{y} : [0, -1]\mathbf{y} = 0\} = \{\mathbf{y} : y_2 = 0\}$ , which means that  $\tilde{T}$  contains nonzero vectors. Hence, the SOSOC does not hold at  $\mathbf{x}^*$ .

**21.6**

a. Write  $f(\mathbf{x}) = x_2$ ,  $g(\mathbf{x}) = -(x_2 + (x_1 - 1)^2 - 3)$ . We have  $\nabla f(\mathbf{x}) = [0, 1]^\top$  and  $\nabla g(\mathbf{x}) = [-2(x_1 - 1), -1]^\top$ . The KKT conditions are:

$$\begin{aligned}\mu &\geq 0 \\ -2\mu(x_1 - 1) &= 0 \\ 1 - \mu &= 0 \\ \mu(x_2 + (x_1 - 1)^2 - 3) &= 0 \\ x_2 + (x_1 - 1)^2 + 3 &\leq 0.\end{aligned}$$

From the third equation we get  $\mu = 1$ . The second equation then gives  $x_1 = 1$ , and the fourth equation gives  $x_2 = 3$ . Therefore, the only point that satisfies the KKT condition is  $\mathbf{x}^* = [1, 3]^\top$ , with a KKT multiplier of  $\mu^* = 1$ .

b. Note that the constraint  $x_2 + (x_1 - 1)^2 + 3 \geq 0$  is active at  $\mathbf{x}^*$ . We have  $T(\mathbf{x}^*) = \{\mathbf{y} : [0, -1]\mathbf{y} = 0\} = \{\mathbf{y} : y_2 = 0\}$ , and  $N(\mathbf{x}^*) = \{\mathbf{y} : \mathbf{y} = [0, -1]z, z \in \mathbb{R}\} = \{\mathbf{y} : y_1 = 0\}$ . Because  $\mu^* > 0$ , we have  $\tilde{T}(\mathbf{x}^*) = T(\mathbf{x}^*) = \{\mathbf{y} : y_2 = 0\}$ .

c. We have

$$\mathbf{L}(\mathbf{x}^*, \mu^*) = \mathbf{O} + 1 \begin{bmatrix} -2 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} -2 & 0 \\ 0 & 0 \end{bmatrix}.$$

From part b,  $T(\mathbf{x}^*) = \{\mathbf{y} : y_2 = 0\}$ . Therefore, for any  $\mathbf{y} \in T(\mathbf{x}^*)$ ,  $\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \mu^*) \mathbf{y} = -2y_1^2 \leq 0$ , which means that  $\mathbf{x}^*$  does not satisfy the SONC.

**21.7**

a. We need to consider two optimization problems. We first consider the minimization problem

$$\begin{aligned}\text{minimize} \quad & (x_1 - 2)^2 + (x_2 - 1)^2 \\ \text{subject to} \quad & x_1^2 - x_2 \leq 0 \\ & x_1 + x_2 - 2 \leq 0 \\ & -x_1 \leq 0.\end{aligned}$$

Then, we form the Lagrangian function

$$l(\mathbf{x}, \boldsymbol{\mu}) = (x_1 - 2)^2 + (x_2 - 1)^2 + \mu_1(x_1^2 - x_2) + \mu_2(x_1 + x_2 - 2) + \mu_3(-x_1).$$

The KKT condition takes the form

$$\begin{aligned}\nabla_{\mathbf{x}} l(\mathbf{x}, \boldsymbol{\mu}) &= \begin{bmatrix} 2(x_1 - 2) + 2\mu_1 x_1 + \mu_2 - \mu_3 & 2(x_2 - 1) - \mu_1 + \mu_2 \end{bmatrix} = \mathbf{0}^T \\ \mu_1(x_1^2 - x_2) &= 0 \\ \mu_2(x_1 + x_2 - 2) &= 0 \\ \mu_3(-x_1) &= 0 \\ \mu_i &\geq 0.\end{aligned}$$

The point  $\mathbf{x}^* = \mathbf{0}$  satisfies the above conditions for  $\mu_1 = -2$ ,  $\mu_2 = 0$ , and  $\mu_3 = -4$ . Thus the point  $\mathbf{x}^*$  does not satisfy the KKT conditions for minimum.

We next consider the maximization problem

$$\begin{aligned}\text{minimize} \quad & -(x_1 - 2)^2 - (x_2 - 1)^2 \\ \text{subject to} \quad & x_1^2 - x_2 \leq 0 \\ & x_1 + x_2 - 2 \leq 0 \\ & -x_1 \leq 0.\end{aligned}$$

The Lagrangian function for the above problem is,

$$l(\mathbf{x}, \boldsymbol{\mu}) = -(x_1 - 2)^2 - (x_2 - 1)^2 + \mu_1(x_1^2 - x_2) + \mu_2(x_1 + x_2 - 2) + \mu_3(-x_1).$$

The KKT condition takes the form

$$\begin{aligned}\nabla_{\mathbf{x}} l(\mathbf{x}, \boldsymbol{\mu}) &= \begin{bmatrix} -2(x_1 - 2) + 2\mu_1 x_1 + \mu_2 - \mu_3 & -2(x_2 - 1) - \mu_1 + \mu_2 \end{bmatrix} = \mathbf{0}^\top \\ \mu_1(x_1^2 - x_2) &= 0 \\ \mu_2(x_1 + x_2 - 2) &= 0 \\ \mu_3(-x_1) &= 0 \\ \mu_i &\geq 0.\end{aligned}$$

The point  $\mathbf{x}^* = \mathbf{0}$  satisfies the above conditions for  $\mu_1 = 2$ ,  $\mu_2 = 0$ , and  $\mu_3 = 4$ . Hence, the point  $\mathbf{x}^*$  satisfies the KKT conditions for maximum.

b. We next compute the Hessian, with respect to  $\mathbf{x}$ , of the lagrangian to obtain

$$\mathbf{L} = \mathbf{F} + \mu_1^* \mathbf{G}_1 = \begin{bmatrix} -2 & 0 \\ 0 & -2 \end{bmatrix} + \begin{bmatrix} 4 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$$

which is indefinite on  $\mathbb{R}^2$ . We next find the subspace

$$\begin{aligned}\tilde{T} &= \left\{ \mathbf{y} : \begin{bmatrix} \nabla g_1(\mathbf{x}^*)^\top \\ \nabla g_3(\mathbf{x}^*)^\top \end{bmatrix} \mathbf{y} = \mathbf{0} \right\} = \left\{ \mathbf{y} : \begin{bmatrix} 0 & -1 \\ -1 & 0 \end{bmatrix} \mathbf{y} = \mathbf{0} \right\} \\ &= \{\mathbf{0}\},\end{aligned}$$

That is,  $\tilde{T}$  is a trivial subspace that consists only of the zero vector. Thus the SOSC for  $\mathbf{x}^*$  to be a strict local maximizer is trivially satisfied.

## 21.8

a. Write  $h(\mathbf{x}) = x_1 - x_2$ ,  $g(\mathbf{x}) = -x_1$ . We have  $Df(\mathbf{x}) = [x_2^2, 2x_1x_2]$ ,  $Dh(\mathbf{x}) = [1, -1]$ ,  $Dg(\mathbf{x}) = [-1, 0]$ .

Note that all feasible points are regular. The KKT condition is:

$$\begin{aligned} x_2^2 + \lambda - \mu &= 0 \\ 2x_1x_2 - \lambda &= 0 \\ \mu x_1 &= 0 \\ \mu &\geq 0 \\ x_1 - x_2 &= 0 \\ x_1 &\geq 0. \end{aligned}$$

We first try  $x_1 = x_1^* = 0$  (active inequality constraint). Substituting and manipulating, we have the solution  $x_1^* = x_2^* = 0$  with  $\mu^* = 0$ , which is a legitimate solution. If we then try  $x_1 = x_1^* > 0$  (inactive inequality constraint), we find that there is no consistent solution to the KKT condition. Thus, there is only one point satisfying the KKT condition:  $\mathbf{x}^* = \mathbf{0}$ .

c. The tangent space at  $\mathbf{x}^* = \mathbf{0}$  is given by

$$T(\mathbf{0}) = \{\mathbf{y} : [1, -1]\mathbf{y} = 0, [-1, 0]\mathbf{y} = 0\} = \{\mathbf{0}\}.$$

Therefore, the SONC holds for the solution in part a.

d. We have

$$\mathbf{L}(\mathbf{x}, \lambda, \mu) = \begin{bmatrix} 0 & 2x_2 \\ 2x_2 & 2x_1 \end{bmatrix}.$$

Hence, at  $\mathbf{x}^* = \mathbf{0}$ , we have  $\mathbf{L}(\mathbf{0}, 0, 0) = \mathbf{O}$ . Since the active constraint at  $\mathbf{x}^* = \mathbf{0}$  is degenerate, we have

$$\tilde{T}(\mathbf{0}, 0) = \{\mathbf{y} : [1, -1]\mathbf{y} = 0\},$$

which is nontrivial. Hence, for any nonzero vector  $\mathbf{y} \in \tilde{T}(\mathbf{0}, 0)$ , we have  $\mathbf{y}^\top \mathbf{L}(\mathbf{0}, 0, 0)\mathbf{y} = 0 \not> 0$ . Thus, the SOSC does not hold for the solution in part a.

## 21.9

a. The KKT condition for the problem is:

$$\begin{aligned} (\mathbf{A}\mathbf{x} - \mathbf{b})^\top \mathbf{A} + \lambda \mathbf{e}^\top - \boldsymbol{\mu}^\top &= \mathbf{0}^\top \\ \boldsymbol{\mu}^\top \mathbf{x} &= 0 \\ \boldsymbol{\mu} &\geq \mathbf{0} \\ \mathbf{e}^\top \mathbf{x} - 1 &= 0 \\ \mathbf{x} &\geq \mathbf{0} \end{aligned}$$

where  $\mathbf{e} = [1, \dots, 1]^\top$ .

b. A feasible point  $\mathbf{x}^*$  is regular in this problem if the vectors  $\mathbf{e}$ ,  $\mathbf{e}_i$ ,  $i \in J(\mathbf{x}^*)$  are linearly independent, where  $J(\mathbf{x}^*) = \{i : x_i^* = 0\}$  and  $\mathbf{e}_i$  is the vector with 0 in all components except the  $i$ th component, which is 1.

In this problem, all feasible points are regular. To see this, note that  $\mathbf{0}$  is not feasible. Therefore, any feasible point results in the set  $J(\mathbf{x}^*)$  having fewer than  $n$  elements, which implies that the vectors  $\mathbf{e}$ ,  $\mathbf{e}_i$ ,  $i \in J(\mathbf{x}^*)$  are linearly independent.

## 21.10

By the KKT Theorem, there exists  $\mu^* \geq 0$  such that

$$\begin{aligned} (\mathbf{x}^* - \mathbf{x}_0) + \mu^* \nabla g(\mathbf{x}^*) &= \mathbf{0} \\ \mu^* g(\mathbf{x}^*) &= 0. \end{aligned}$$

Premultiplying both sides of the first equation by  $(\mathbf{x}^* - \mathbf{x}_0)^\top$ , we obtain

$$\|\mathbf{x}^* - \mathbf{x}_0\|^2 + \mu^* (\mathbf{x}^* - \mathbf{x}_0)^\top \nabla g(\mathbf{x}^*) = 0.$$

Since  $\|\mathbf{x}^* - \mathbf{x}_0\|^2 > 0$  (because  $g(\mathbf{x}_0) > 0$ ) and  $\mu^* \geq 0$ , we deduce that  $(\mathbf{x}^* - \mathbf{x}_0)^\top \nabla g(\mathbf{x}^*) < 0$  and  $\mu^* > 0$ . From the second KKT condition above, we conclude that  $g(\mathbf{x}^*) = 0$ .

### 21.11

a. By inspection, we guess the point  $[2, 2]^\top$  (drawing a picture may help).

b. We write  $f(\mathbf{x}) = (x_1 - 3)^2 + (x_2 - 4)^2$ ,  $g_1(\mathbf{x}) = -x_1$ ,  $g_2(\mathbf{x}) = -x_2$ ,  $g_3(\mathbf{x}) = x_1 - 2$ ,  $g_4(\mathbf{x}) = x_2 - 2$ ,  $\mathbf{g} = [g_1, g_2, g_3, g_4]^\top$ . The problem becomes

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{g}(\mathbf{x}) \leq \mathbf{0}. \end{aligned}$$

We now check the SOSC for the point  $\mathbf{x}^* = [2, 2]^\top$ . We have two active constraints:  $g_3, g_4$ . Regularity holds, since  $\nabla g_3(\mathbf{x}^*) = [1, 0]^\top$  and  $\nabla g_4(\mathbf{x}^*) = [0, 1]^\top$ . We have  $\nabla f(\mathbf{x}^*) = [-2, -4]^\top$ . We need to find a  $\mu^* \in \mathbb{R}^4$ ,  $\mu^* \geq \mathbf{0}$ , satisfying FONC. From the condition  $\mu^{*\top} \mathbf{g}(\mathbf{x}^*) = 0$ , we deduce that  $\mu_1^* = \mu_2^* = 0$ . Hence,  $Df(\mathbf{x}^*) + \mu^{*\top} D\mathbf{g}(\mathbf{x}^*) = \mathbf{0}^\top$  if and only if  $\mu^* = [0, 0, 2, 4]^\top$ . Now,

$$\mathbf{F}(\mathbf{x}^*) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \quad [\mu^* \mathbf{G}(\mathbf{x}^*)] = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}.$$

Hence

$$\mathbf{L}(\mathbf{x}^*, \mu^*) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

which is positive definite on  $\mathbb{R}^2$ . Hence, SOSC is satisfied, and  $\mathbf{x}^*$  is a strict local minimizer.

### 21.12

The KKT condition is

$$\begin{aligned} \mathbf{x}^\top \mathbf{Q} + \mu^\top \mathbf{A} &= \mathbf{0}^\top \\ \mu^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) &= 0 \\ \mu &\geq \mathbf{0} \\ \mathbf{A}\mathbf{x} - \mathbf{b} &\leq \mathbf{0}. \end{aligned}$$

Postmultiplying the first equation by  $\mathbf{x}$  gives

$$\mathbf{x}^\top \mathbf{Q}\mathbf{x} + \mu^\top \mathbf{A}\mathbf{x} = 0.$$

We note from the second equation that  $\mu^\top \mathbf{A}\mathbf{x} = \mu^\top \mathbf{b}$ . Hence,

$$\mathbf{x}^\top \mathbf{Q}\mathbf{x} + \mu^\top \mathbf{b} = 0.$$

Since  $\mathbf{Q} > 0$ , the first term is nonnegative. Also, the second term is nonnegative because  $\mu \geq \mathbf{0}$  and  $\mathbf{b} \geq \mathbf{0}$ . Hence, we conclude that both terms must be zero. Because  $\mathbf{Q} > 0$ , we must have  $\mathbf{x} = \mathbf{0}$ .

*Aside:* Actually, we can deduce that the only solution to the KKT condition must be  $\mathbf{0}$ , as follows. The problem is convex; thus, the only points satisfying the KKT condition are global minimizers. However, we see that  $\mathbf{0}$  is a feasible point, and is the only point for which the objective function value is 0. Further, the objective function is bounded below by 0. Hence,  $\mathbf{0}$  is the only global minimizer.

### 21.13

a. We have one scalar equality constraint with  $h(\mathbf{x}) = [c, d]^\top \mathbf{x} - e$  and two scalar inequality constraints with  $\mathbf{g}(\mathbf{x}) = -\mathbf{x}$ . Hence, there exists  $\mu^* \in \mathbb{R}^2$  and  $\lambda^* \in \mathbb{R}$  such that

$$\begin{aligned} \mu^* &\geq \mathbf{0} \\ a + c\lambda^* - \mu_1^* &= 0 \\ b + d\lambda^* - \mu_2^* &= 0 \\ \mu^{*\top} \mathbf{x}^* &= 0 \\ cx_1^* + dx_2^* &= e \\ \mathbf{x}^* &\geq \mathbf{0}. \end{aligned}$$

b. Because  $\mathbf{x}^*$  is a basic feasible solution, and the equality constraint precludes the point  $\mathbf{0}$ , exactly one of the inequality constraints is active. The vectors  $\nabla h(\mathbf{x}^*) = [c, d]^\top$  and  $\nabla g_1 = [1, 0]^\top$  are linearly independent. Similarly, the vectors  $\nabla h(\mathbf{x}^*) = [c, d]^\top$  and  $\nabla g_2 = [0, 1]^\top$  are linearly independent. Hence,  $\mathbf{x}^*$  must be regular.

c. The tangent space is given by

$$\begin{aligned} T(\mathbf{x}^*) &= \{\mathbf{y} \in \mathbb{R}^n : Dh(\mathbf{x}^*)\mathbf{y} = \mathbf{0}, Dg_j(\mathbf{x}^*)\mathbf{y} = 0, j \in J(\mathbf{x}^*)\} \\ &= \mathcal{N}(\mathbf{M}), \end{aligned}$$

where  $\mathbf{M}$  is a matrix with the first row equal to  $Dh(\mathbf{x}^*) = [c, d]$ , and the second row is either  $Dg_1 = [1, 0]$  or  $Dg_2 = [0, 1]$ . But, as we have seen in part b,  $\text{rank } \mathbf{M} = 2$ . Hence,  $T(\mathbf{x}^*) = \{\mathbf{0}\}$ .

d. Recall that we can take  $\boldsymbol{\mu}^*$  to be the relative cost coefficient vector (i.e., the KKT conditions are satisfied with  $\boldsymbol{\mu}^*$  being the relative cost coefficient vector). If the relative cost coefficients of all nonbasic variables are strictly positive, then  $\mu_j^* > 0$  for all  $j \in J(\mathbf{x}^*)$ . Hence,  $\tilde{T}(\mathbf{x}^*, \boldsymbol{\mu}^*) = T(\mathbf{x}^*) = \{\mathbf{0}\}$ , which implies that  $L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) > 0$  on  $\tilde{T}(\mathbf{x}^*, \boldsymbol{\mu}^*)$ . Hence, the SOSC is satisfied.

#### 21.14

Let  $\mathbf{x}^*$  be a solution. Since  $\mathbf{A}$  is of full rank,  $\mathbf{x}^*$  is regular. The KKT Theorem states that  $\mathbf{x}^*$  satisfies:

$$\begin{aligned} \boldsymbol{\mu}^* &\geq \mathbf{0} \\ \mathbf{c}^\top + \boldsymbol{\mu}^{*\top} \mathbf{A} &= \mathbf{0} \\ \boldsymbol{\mu}^{*\top} \mathbf{A} \mathbf{x}^* &= 0. \end{aligned}$$

If we postmultiply the second equation by  $\mathbf{x}^*$  and subtract the third from the result, we get

$$\mathbf{c}^\top \mathbf{x}^* = 0.$$

#### 21.15

a. We can write the LP as

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}) \\ &\text{subject to} && \mathbf{h}(\mathbf{x}) = \mathbf{0}, \\ &&& \mathbf{g}(\mathbf{x}) \leq \mathbf{0}, \end{aligned}$$

where  $f(\mathbf{x}) = \mathbf{c}^\top \mathbf{x}$ ,  $\mathbf{h}(\mathbf{x}) = \mathbf{A}\mathbf{x} - \mathbf{b}$ , and  $\mathbf{g}(\mathbf{x}) = -\mathbf{x}$ . Thus, we have  $Df(\mathbf{x}) = \mathbf{c}^\top$ ,  $D\mathbf{h}(\mathbf{x}) = \mathbf{A}$ , and  $D\mathbf{g}(\mathbf{x}) = -\mathbf{I}$ . The Karush-Kuhn-Tucker conditions for the above problem have the form: if  $\mathbf{x}^*$  is a local minimizer, then there exists  $\boldsymbol{\lambda}^*$  and  $\boldsymbol{\mu}^*$  such that

$$\begin{aligned} \boldsymbol{\mu}^* &\geq \mathbf{0} \\ \mathbf{c}^\top + \boldsymbol{\lambda}^{*\top} \mathbf{A} - \boldsymbol{\mu}^{*\top} &= \mathbf{0}^\top \\ \boldsymbol{\mu}^{*\top} \mathbf{x}^* &= 0. \end{aligned}$$

b. Let  $\mathbf{x}^*$  be an optimal feasible solution. Then,  $\mathbf{x}^*$  satisfies the Karush-Kuhn-Tucker conditions listed in part a. Since  $\boldsymbol{\mu}^* \geq \mathbf{0}$ , then from the second condition in part a, we obtain  $(-\boldsymbol{\lambda}^*)^\top \mathbf{A} \leq \mathbf{c}^\top$ . Hence,  $\bar{\boldsymbol{\lambda}} = -\boldsymbol{\lambda}^*$  is a feasible solution to the dual (see Chapter 17). Postmultiplying the second condition in part a by  $\mathbf{x}^*$ , we have

$$0 = \mathbf{c}^\top \mathbf{x}^* + \boldsymbol{\lambda}^{*\top} \mathbf{A} \mathbf{x}^* - \boldsymbol{\mu}^{*\top} \mathbf{x}^* = \mathbf{c}^\top \mathbf{x}^* + \boldsymbol{\lambda}^{*\top} \mathbf{b}$$

which gives

$$\mathbf{c}^\top \mathbf{x}^* = \bar{\boldsymbol{\lambda}}^\top \mathbf{b}.$$

Hence,  $\bar{\boldsymbol{\lambda}}$  achieves the same objective function value for the dual as  $\mathbf{x}^*$  for the primal.

c. From part a, we have  $\boldsymbol{\mu}^{*\top} = \mathbf{c}^\top - \bar{\boldsymbol{\lambda}}^\top \mathbf{A}$ . Substituting this into  $\boldsymbol{\mu}^{*\top} \mathbf{x}^* = 0$  yields the desired result.



**21.16**

By definition of  $J(\mathbf{x}^*)$ , we have  $g_i(\mathbf{x}^*) < 0$  for all  $i \notin J(\mathbf{x}^*)$ . Since by assumption  $g_i$  is continuous for all  $i$ , there exists  $\varepsilon > 0$  such that  $g_i(\mathbf{x}) < 0$  for all  $i \notin J(\mathbf{x}^*)$  and all  $\mathbf{x}$  in the set  $B = \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\| < \varepsilon\}$ . Let  $S_1 = \{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}, g_j(\mathbf{x}) \leq 0, j \in J(\mathbf{x}^*)\}$ . We claim that  $S \cap B = S_1 \cap B$ . To see this, note that clearly  $S \cap B \subset S_1 \cap B$ . To show that  $S_1 \cap B \subset S \cap B$ , suppose  $\mathbf{x} \in S_1 \cap B$ . Then, by definition of  $S_1$  and  $B$ , we have  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ ,  $g_j(\mathbf{x}) \leq 0$  for all  $j \in J(\mathbf{x}^*)$ , and  $g_i(\mathbf{x}) < 0$  for all  $i \notin J(\mathbf{x}^*)$ . Hence,  $\mathbf{x} \in S \cap B$ .

Since  $\mathbf{x}^*$  is a local minimizer of  $f$  over  $S$ , and  $S \cap B \subset S$ ,  $\mathbf{x}^*$  is also a local minimizer of  $f$  over  $S \cap B = S_1 \cap B$ . Hence, we conclude that  $\mathbf{x}^*$  is a regular local minimizer of  $f$  on  $S_1$ . Note that  $S' \subset S_1$ , and  $\mathbf{x}^* \in S'$ . Therefore,  $\mathbf{x}^*$  is a regular local minimizer of  $f$  on  $S'$ .

**21.17**

Write  $f(\mathbf{x}) = x_1^2 + x_2^2$ ,  $g_1(\mathbf{x}) = x_1^2 - x_2 - 4$ ,  $g_2(\mathbf{x}) = x_2 - x_1 - 2$ , and  $\mathbf{g} = [g_1, g_2]^\top$ . We have  $\nabla f(\mathbf{x}) = [2x_1, 2x_2]^\top$ ,  $\nabla g_1(\mathbf{x}) = [2x_1, -1]^\top$ ,  $\nabla g_2(\mathbf{x}) = [-1, 1]^\top$ , and  $D^2 f(\mathbf{x}) = \text{diag}[2, 2]$ . We compute

$$\nabla f(\mathbf{x}) + \boldsymbol{\mu}^\top \nabla \mathbf{g}(\mathbf{x}) = [2x_1 + 2\mu_1 x_1 - \mu_2, 2x_2 - \mu_1 + \mu_2]^\top.$$

We use the FONC to find critical points. Rewriting  $\nabla f(\mathbf{x}) + \boldsymbol{\mu}^\top \nabla \mathbf{g}(\mathbf{x}) = \mathbf{0}$ , we obtain

$$x_1 = \frac{\mu_2}{2 + 2\mu_1}, \quad x_2 = \frac{\mu_1 - \mu_2}{2}.$$

We also use  $\boldsymbol{\mu}^\top \mathbf{g}(\mathbf{x}) = 0$  and  $\boldsymbol{\mu} \geq \mathbf{0}$ , giving

$$\mu_1(x_1^2 - x_2 - 4) = 0, \quad \mu_2(x_2 - x_1 - 2) = 0.$$

The vector  $\boldsymbol{\mu}$  has two components; therefore, we try four different cases.

**Case 1:** ( $\mu_1 > 0$ ,  $\mu_2 > 0$ ) We have

$$x_1^2 - x_2 - 4 = 0, \quad x_2 - x_1 - 2 = 0.$$

We obtain two solutions:  $\mathbf{x}^{(1)} = [-2, 0]^\top$  and  $\mathbf{x}^{(2)} = [3, 5]^\top$ . For  $\mathbf{x}^{(1)}$ , the two FONC equations give  $\mu_1 = \mu_2$  and  $-2(2 + 2\mu_1) = \mu_1$ , which yield  $\mu_1 = \mu_2 = -4/5$ . This is not a legitimate solution since we require  $\boldsymbol{\mu} \geq \mathbf{0}$ . For  $\mathbf{x}^{(2)}$ , the two FONC equations give  $\mu_1 - \mu_2 = 10$  and  $3(2 + 2\mu_1) = \mu_2$ , which yield  $\boldsymbol{\mu} = [-16/5, -66/5]$ . Again, this is not a legitimate solution.

**Case 2:** ( $\mu_1 = 0$ ,  $\mu_2 > 0$ ) We have

$$x_2 - x_1 - 2 = 0, \quad x_1 = \frac{\mu_2}{2}, \quad x_2 = -\frac{\mu_2}{2}.$$

Hence,  $x_1 = -x_2$ , and thus  $\mathbf{x} = [-1, 1]$ ,  $\mu_2 = -2$ . This is not a legitimate solution since we require  $\boldsymbol{\mu} \geq \mathbf{0}$ .

**Case 3:** ( $\mu_1 > 0$ ,  $\mu_2 = 0$ ) We have

$$x_1^2 - x_2 - 4 = 0, \quad x_1 = 0, \quad x_2 = \frac{\mu_1}{2}.$$

Therefore,  $x_2 = -4$ ,  $\mu_1 = -8$ , and again we don't have a legitimate solution.

**Case 4:** ( $\mu_1 = 0$ ,  $\mu_2 = 0$ ) We have  $x_1 = x_2 = 0$ , and all constraints are inactive. This is a legitimate candidate for the minimizer. We now apply the SOSC. Note that since the candidate is an interior point of the constraint set, the SOSC for the problem is equivalent to the SOSC for unconstrained optimization. The Hessian matrix  $D^2 f(\mathbf{x}) = \text{diag}[2, 2]$  is symmetric and positive definite. Hence, by the SOSC, the point  $\mathbf{x}^* = [0, 0]^\top$  is the strict local minimizer (in fact, it is easy to see that it is a global minimizer).

**21.18**

Write  $f(\mathbf{x}) = x_1^2 + x_2^2$ ,  $g_1(\mathbf{x}) = -x_1 + x_2^2 + 4$ ,  $g_2(\mathbf{x}) = x_1 - 10$ , and  $\mathbf{g} = [g_1, g_2]^\top$ . We have  $\nabla f(\mathbf{x}) = [2x_1, 2x_2]^\top$ ,  $\nabla g_1(\mathbf{x}) = [-1, 2x_2]^\top$ ,  $\nabla g_2(\mathbf{x}) = [1, 0]^\top$ ,  $D^2 f(\mathbf{x}) = \text{diag}[2, 2]$ ,  $D^2 g_1(\mathbf{x}) = \text{diag}[0, 2]$ , and  $D^2 g_2(\mathbf{x}) = \mathbf{O}$ . We compute

$$\nabla f(\mathbf{x}) + \boldsymbol{\mu}^\top \nabla \mathbf{g}(\mathbf{x}) = [2x_1 - \mu_1 + \mu_2, 2x_2 + 2\mu_1 x_2]^\top.$$

We use the FONC to find critical points. Rewriting  $\nabla f(\mathbf{x}) + \boldsymbol{\mu}^\top \nabla \mathbf{g}(\mathbf{x}) = \mathbf{0}$ , we obtain

$$x_1 = \frac{\mu_1 - \mu_2}{2}, \quad x_2(1 + \mu_1) = 0.$$

Since we require  $\boldsymbol{\mu} \geq \mathbf{0}$ , we deduce that  $x_2 = 0$ . Using  $\boldsymbol{\mu}^\top \mathbf{g}(\mathbf{x}) = 0$  gives

$$\mu_1(-x_1 + 4) = 0, \quad \mu_2 = 0.$$

We are left with two cases.

**Case 1:** ( $\mu_1 > 0, \mu_2 = 0$ ) We have  $-x_1 + 4 = 0$ , and  $\mu_1 = 8$ , which is a legitimate candidate.

**Case 2:** ( $\mu_1 = 0, \mu_2 = 0$ ) We have  $x_1 = x_2 = 0$ , which is not a legitimate candidate, since it is not a feasible point.

We now apply SOSC to our candidate  $\mathbf{x} = [4, 0]^\top$ ,  $\boldsymbol{\mu} = [8, 0]^\top$ . Now,

$$\mathbf{L}([4, 0]^\top, [8, 0]^\top) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} + 8 \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 18 \end{bmatrix},$$

which is positive definite on all of  $\mathbb{R}^2$ . The point  $[4, 0]^\top$  is clearly regular. Hence, by the SOSC,  $\mathbf{x}^* = [4, 0]$  is a strict local minimizer.

### 21.19

Write  $f(\mathbf{x}) = x_1^2 + x_2^2$ ,  $g_1(\mathbf{x}) = -x_1 - x_2^2 + 4$ ,  $g_2(\mathbf{x}) = 3x_2 - x_1$ ,  $g_3(\mathbf{x}) = -3x_2 - x_1$  and  $\mathbf{g} = [g_1, g_2, g_3]^\top$ . We have  $\nabla f(\mathbf{x}) = [2x_1, 2x_2]^\top$ ,  $\nabla g_1(\mathbf{x}) = [-1, -2x_2]^\top$ ,  $\nabla g_2(\mathbf{x}) = [-1, 3]^\top$ ,  $\nabla g_3(\mathbf{x}) = [-1, -3]^\top$ ,  $D^2 f(\mathbf{x}) = \text{diag}[2, 2]$ ,  $D^2 g_1(\mathbf{x}) = \text{diag}[0, -2]$ , and  $D^2 g_2(\mathbf{x}) = D^2 g_3(\mathbf{x}) = \mathbf{O}$ .

From the figure, we see that the two candidates are  $\mathbf{x}^{(1)} = [3, 1]$  and  $\mathbf{x}^{(2)} = [3, -1]$ . Both points are easily verified to be regular.

For  $\mathbf{x}^{(1)}$ , we have  $\mu_3 = 0$ . Now,

$$Df(\mathbf{x}^{(1)}) + \boldsymbol{\mu}^\top D\mathbf{g}(\mathbf{x}^{(1)}) = [6 - \mu_1 - \mu_2, 2 - 2\mu_1 + 3\mu_2] = \mathbf{0}^\top,$$

which yields  $\mu_1 = 4, \mu_2 = 2$ . Now,  $\tilde{T}(\mathbf{x}^{(1)}) = \{\mathbf{0}\}$ . Therefore, any matrix is positive definite on  $\tilde{T}(\mathbf{x}^{(1)})$ . Hence, by the SOSC,  $\mathbf{x}^{(1)}$  is a strict local minimizer.

For  $\mathbf{x}^{(2)}$ , we have  $\mu_2 = 0$ . Now,

$$Df(\mathbf{x}^{(1)}) + \boldsymbol{\mu}^\top D\mathbf{g}(\mathbf{x}^{(1)}) = [6 - \mu_1 - \mu_3, -2 + 2\mu_1 - 3\mu_3] = \mathbf{0}^\top,$$

which yields  $\mu_1 = 4, \mu_3 = 2$ . Now, again we have  $\tilde{T}(\mathbf{x}^{(2)}) = \{\mathbf{0}\}$ . Therefore, any matrix is positive definite on  $\tilde{T}(\mathbf{x}^{(2)})$ . Hence, by the SOSC,  $\mathbf{x}^{(2)}$  is a strict local minimizer.

### 21.20

a. Write  $f(\mathbf{x}) = 3x_1$  and  $g(\mathbf{x}) = 2 - x_1 - x_2^2$ . We have  $\nabla f(\mathbf{x}^*) = [3, 0]^\top$  and  $\nabla g(\mathbf{x}^*) = [-1, 0]^\top$ . Hence, letting  $\mu^* = 3$ , we have  $\nabla f(\mathbf{x}^*) + \mu^* \nabla g(\mathbf{x}^*) = \mathbf{0}$ . Note also that  $\mu^* \geq 0$  and  $\mu^* g(\mathbf{x}^*) = 0$ . Hence,  $\mathbf{x}^* = [2, 0]^\top$  satisfies the KKT (first order necessary) condition.

b. We have  $\mathbf{F}(\mathbf{x}^*) = \mathbf{O}$  and  $\mathbf{G}(\mathbf{x}^*) = \text{diag}[0, -2]$ . Hence,  $\mathbf{L}(\mathbf{x}^*, \mu^*) = \mathbf{O} + 3 \text{diag}[0, -2] = \text{diag}[0, -6]$ . Also,  $T(\mathbf{x}^*) = \{\mathbf{y} : [-1, 0]\mathbf{y} = 0\} = \{\mathbf{y} : y_1 = 0\}$ . Hence,  $\mathbf{x}^* = [2, 0]^\top$  does not satisfy the second order necessary condition.

c. No. Consider points of the form  $\mathbf{x} = [-x_2^2 + 2, x_2]^\top$ ,  $x_2 \in \mathbb{R}$ . Such points are feasible, and could be arbitrarily close to  $\mathbf{x}^*$ . However, for such points  $\mathbf{x} \neq \mathbf{x}^*$ ,

$$f(\mathbf{x}) = 3(-x_2^2 + 2) = 6 - 6x_2^2 < 6 = f(\mathbf{x}^*).$$

Hence,  $\mathbf{x}^*$  is not a local minimizer.

### 21.21

The KKT condition for the problem is

$$\begin{aligned} \boldsymbol{\mu}^* &\geq \mathbf{0} \\ \mathbf{x}^* + \lambda^* \mathbf{a} - \boldsymbol{\mu}^* &= \mathbf{0} \\ \boldsymbol{\mu}^{*\top} \mathbf{x}^* &= 0. \end{aligned}$$

Premultiplying the second KKT condition above by  $\boldsymbol{\mu}^{*\top}$  and using the third condition, we get

$$\lambda^* \boldsymbol{\mu}^{*\top} \mathbf{a} = \|\boldsymbol{\mu}^*\|^2.$$

Also, premultiplying the second KKT condition above by  $\mathbf{x}^{*\top}$  and using the feasibility condition  $\mathbf{a}^\top \mathbf{x}^* = b$ , we get

$$\lambda^* = -\frac{\|\mathbf{x}^*\|^2}{b} < 0.$$

We conclude that  $\boldsymbol{\mu}^* = \mathbf{0}$ . For if not, the equation  $\lambda^* \boldsymbol{\mu}^{*\top} \mathbf{a} = \|\boldsymbol{\mu}^*\|^2$  implies that  $\boldsymbol{\mu}^{*\top} \mathbf{a} < 0$ , which contradicts  $\boldsymbol{\mu}^* \geq \mathbf{0}$  and  $\mathbf{a} \geq \mathbf{0}$ .

Rewriting the second KKT condition with  $\boldsymbol{\mu}^* = \mathbf{0}$  yields

$$\mathbf{x}^* = -\lambda^* \mathbf{a}.$$

Using the feasibility condition  $\mathbf{a}^\top \mathbf{x}^* = b$ , we get

$$\mathbf{x}^* = \mathbf{a} \frac{b}{\|\mathbf{a}\|^2}.$$

## 21.22

a. Suppose  $(x_1^*)^2 + (x_2^*)^2 < 1$ . Then, the point  $\mathbf{x}^* = [x_1^*, x_2^*]^\top$  lies in the interior of the constraint set  $\Omega = \{\mathbf{x} : \|\mathbf{x}\|^2 \leq 1\}$ . Hence, by the FONC for unconstrained optimization, we have that  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ , where  $f(\mathbf{x}) = \|\mathbf{x} - [a, b]^\top\|^2$  is the objective function. Now,  $\nabla f(\mathbf{x}^*) = 2(\mathbf{x}^* - [a, b]^\top) = \mathbf{0}$ , which implies that  $\mathbf{x}^* = [a, b]^\top$  which violates the assumption  $(x_1^*)^2 + (x_2^*)^2 < 1$ .

b. First, we need to show that  $\mathbf{x}^*$  is a regular point. For this, note that if we write the constraint as  $g(\mathbf{x}) = \|\mathbf{x}\|^2 - 1 \leq 0$ , then  $\nabla g(\mathbf{x}^*) = 2\mathbf{x}^* \neq \mathbf{0}$ . Therefore,  $\mathbf{x}^*$  is a regular point. Hence, by the Karush-Kuhn-Tucker theorem, there exists  $\mu \in \mathbb{R}$ ,  $\mu \geq 0$ , such that

$$\nabla f(\mathbf{x}^*) + \mu \nabla g(\mathbf{x}^*) = \mathbf{0},$$

which gives

$$\mathbf{x}^* = \frac{1}{1 + \mu} \begin{bmatrix} a \\ b \end{bmatrix}.$$

Hence,  $\mathbf{x}^*$  is unique, and we can write  $x_1^* = \alpha a$ ,  $x_2^* = \alpha b$ , where  $\alpha = 1/(1 + \mu) \geq 0$ .

c. Using part b and the fact that  $\|\mathbf{x}^*\| = 1$ , we get  $\|\mathbf{x}^*\|^2 = \alpha^2 \|[a, b]^\top\|^2 = 1$ , which gives  $\alpha = 1/\sqrt{a^2 + b^2}$ .

## 21.23

a. The Karush-Kuhn-Tucker conditions for this problem are

$$\begin{aligned} 2x_1^* + \mu^* \exp(x_1^*) &= 0 \\ 2(x_2^* + 1) - \mu^* &= 0 \\ \mu^* (\exp(x_1^*) - x_2^*) &= 0 \\ \exp(x_1^*) &\leq x_2^* \\ \mu^* &\geq 0. \end{aligned}$$

b. From the second equation in part a, we obtain  $\mu^* = 2(x_2^* + 1)$ . Since  $x_2^* \geq \exp(x_1^*) > 0$ , then  $\mu^* > 0$ . Hence, by the third equation in part a, we obtain  $x_2^* = \exp(x_1^*)$ .

c. Since  $\mu^* = 2(x_2^* + 1) = 2(\exp(x_1^*) + 1)$ , then by the first equation in part a, we have

$$2x_1^* + 2(\exp(x_1^*) + 1) \exp(x_1^*) = 0$$

which implies

$$x_1^* = -(\exp(2x_1^*) + \exp(x_1^*)).$$

Since  $\exp(x_1^*), \exp(2x_1^*) > 0$ , then  $x_1^* < 0$ , and hence  $\exp(x_1^*), \exp(2x_1^*) < 1$ . Therefore,  $x_1^* > -2$ .

### 21.24

a. We rewrite the problem as

$$\begin{array}{ll} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & g(\mathbf{x}) \leq 0, \end{array}$$

where  $f(\mathbf{x}) = \mathbf{c}^\top \mathbf{x}$  and  $g(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|^2 - 1$ . Hence,  $\nabla f(\mathbf{x}) = \mathbf{c}$  and  $\nabla g(\mathbf{x}) = \mathbf{x}$ . Note that  $\mathbf{x}^* \neq \mathbf{0}$  (for otherwise it would not be feasible), and therefore it is a regular point. By the KKT theorem, there exists  $\mu^* \leq 0$  such that  $\mathbf{c} = \mu^* \mathbf{x}^*$  and  $\mu^* g(\mathbf{x}^*) = 0$ . Since  $\mathbf{c} \neq \mathbf{0}$ , we must have  $\mu^* \neq 0$ . Therefore,  $g(\mathbf{x}^*) = 0$ , which implies that  $\|\mathbf{x}^*\|^2 = 2$ .

b. From part a, we have  $\alpha^2 \|\mathbf{e}\|^2 = 2$ . Since  $\|\mathbf{e}\|^2 = n$ , we have  $\alpha = \sqrt{2/n}$ .

To find  $\mathbf{c}$ , we use

$$4 = \mathbf{c}^\top \mathbf{x}^* + 8 = \mu^* \|\mathbf{x}^*\|^2 + 8 = 2\mu^* + 8,$$

and thus  $\mu^* = -2$ . Hence,  $\mathbf{c} = -2\alpha \mathbf{e} = -\left(2\sqrt{2/n}\right) \mathbf{e}$ .

### 21.25

We can represent the equivalent problem as

$$\begin{array}{ll} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & g(\mathbf{x}) \leq 0, \end{array}$$

where  $g(\mathbf{x}) = \frac{1}{2}\|\mathbf{h}(\mathbf{x})\|^2$ . Note that

$$\nabla g(\mathbf{x}) = D\mathbf{h}(\mathbf{x})^\top \mathbf{h}(\mathbf{x}).$$

Therefore, the KKT condition is:

$$\begin{array}{rcl} \mu^* & \geq & 0 \\ \nabla f(\mathbf{x}^*) + \mu^* D\mathbf{h}(\mathbf{x}^*)^\top \mathbf{h}(\mathbf{x}^*) & = & \mathbf{0} \\ \mu^* \|\mathbf{h}(\mathbf{x}^*)\| & = & 0. \end{array}$$

Note that for a feasible point  $\mathbf{x}^*$ , we have  $\mathbf{h}(\mathbf{x}^*) = \mathbf{0}$ . Therefore, the KKT condition becomes

$$\begin{array}{rcl} \mu^* & \geq & 0 \\ \nabla f(\mathbf{x}^*) & = & \mathbf{0}. \end{array}$$

Note that  $\nabla g(\mathbf{x}^*) = \mathbf{0}$ . Therefore, any feasible point  $\mathbf{x}^*$  is not regular. Hence, the KKT theorem cannot be applied in this case. This should be clear, since obviously  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  is not necessary for optimality in general.

## 22. Convex Optimization Problems

### 22.1

The given function is a quadratic, which we represent in the form

$$f = \mathbf{x}^\top \begin{bmatrix} -1 & -\alpha & 1 \\ -\alpha & -1 & -2 \\ 1 & -2 & -5 \end{bmatrix} \mathbf{x}.$$

A quadratic function is concave if and only if it is negative-semidefinite. Equivalently, if and only if its negative is positive-semidefinite. On the other hand, a symmetric matrix is positive semidefinite if and only

if all its principal minors, not just the leading principal minors, are nonnegative. Thus we will determine the range of the parameter  $\alpha$  for which

$$-f = \mathbf{x}^\top \begin{bmatrix} 1 & \alpha & -1 \\ \alpha & 1 & 2 \\ -1 & 2 & 5 \end{bmatrix} \mathbf{x}.$$

is positive-semidefinite. It is easy to see that the three first-order principal minors (diagonal elements of  $\mathbf{F}$ ) are all positive. There are three second-order principal minors. Only one of them, the leading principal minor, is a function of the parameter  $\alpha$ ,

$$\det \mathbf{F}(1:1, 1:1) = \det \begin{bmatrix} 1 & \alpha \\ \alpha & 1 \end{bmatrix} = 1 - \alpha^2.$$

The above second-order leading principal minor is nonnegative if and only if

$$\alpha \in [-1, 1].$$

The other second-order principal minors are

$$\det \mathbf{F}(1:2, 1:2) \quad \text{and} \quad \det \mathbf{F}(2:3, 2:3),$$

and they are positive. There is only one third-order principal minor,  $\det \mathbf{F}$ , where

$$\begin{aligned} \det \mathbf{F} &= \det \begin{bmatrix} 1 & 2 \\ 2 & 5 \end{bmatrix} - \alpha \det \begin{bmatrix} \alpha & -1 \\ 2 & 5 \end{bmatrix} - \det \begin{bmatrix} \alpha & -1 \\ 1 & 2 \end{bmatrix} \\ &= 1 - \alpha(5\alpha + 2) - (2\alpha + 1) \\ &= 1 - 5\alpha^2 - 2\alpha - 2\alpha - 1 \\ &= -5\alpha^2 - 4\alpha. \end{aligned}$$

The third-order principal minor is nonnegative if and only if,  $-5\alpha^2 - 4\alpha \geq 0$ , that is, if and only if

$$\alpha \in [-4/5, 0].$$

Combining this with  $\alpha \in [-1, 1]$  from above, we conclude that the function  $f$  is negative-semidefinite, equivalently, the quadratic function  $f$  is concave, if and only if

$$\alpha \in [-4/5, 0].$$

---

## 22.2

We have

$$\begin{aligned} \phi(\alpha) &= \frac{1}{2}(\mathbf{x} + \alpha \mathbf{d})^\top \mathbf{Q}(\mathbf{x} + \alpha \mathbf{d}) - (\mathbf{x} + \alpha \mathbf{d})^\top \mathbf{b} \\ &= \frac{1}{2}(\mathbf{d}^\top \mathbf{Q} \mathbf{d})\alpha^2 + \mathbf{d}^\top (\mathbf{Q} \mathbf{x} - \mathbf{b})\alpha + \left( \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \mathbf{x}^\top \mathbf{b} \right). \end{aligned}$$

This is a quadratic function of  $\alpha$ . Since  $\mathbf{Q} > 0$ , then

$$\frac{d^2 \phi}{d\alpha^2}(\alpha) = \mathbf{d}^\top \mathbf{Q} \mathbf{d} > 0$$

and hence by Theorem 22.5,  $\phi$  is strictly convex.

---

## 22.3

Write  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q} \mathbf{x}$ , where

$$\mathbf{Q} = \frac{1}{2} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}.$$

Let  $\mathbf{x}, \mathbf{y} \in \Omega$ . Then,  $\mathbf{x} = [a_1, ma_1]^\top$  and  $\mathbf{y} = [a_2, ma_2]^\top$  for some  $a_1, a_2 \in \mathbb{R}$ . By Proposition 22.1, it is enough to show that  $(\mathbf{y} - \mathbf{x})^\top \mathbf{Q}(\mathbf{y} - \mathbf{x}) \geq 0$ . By substitution,

$$(\mathbf{y} - \mathbf{x})^\top \mathbf{Q}(\mathbf{y} - \mathbf{x}) = m(a_2 - a_1)^2 \geq 0,$$

which completes the proof.

#### 22.4

Let  $\mathbf{x}, \mathbf{y} \in \Omega$  and  $\alpha \in (0, 1)$ . Then,  $h(\mathbf{x}) = h(\mathbf{y}) = c$ . By convexity of  $\Omega$ ,  $h(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) = c$ . Therefore,

$$h(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) = \alpha h(\mathbf{x}) + (1 - \alpha)h(\mathbf{y})$$

and so  $h$  is convex over  $\Omega$ . We also have

$$-h(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) = \alpha(-h(\mathbf{x})) + (1 - \alpha)(-h(\mathbf{y})),$$

which shows that  $-h$  is convex, and thus  $h$  is concave.

#### 22.5

At  $x = 0$ , for  $\xi \in [-1, 1]$ , we have, for all  $y \in \mathbb{R}$ ,

$$|y| \geq |0| + \xi(y - 0) = \xi y.$$

Thus, in this case any  $\xi$  in the interval  $[-1, 1]$  is a subgradient of  $f$  at  $x = 0$ .

At  $x = 1$ ,  $\xi = 1$  is the only subgradient of  $f$ , because, for all  $y \in \mathbb{R}$ ,

$$|y| \geq 1 + \xi(y - 1) = y.$$

#### 22.6

Let  $\mathbf{x}, \mathbf{y} \in \Omega$  and  $\alpha \in (0, 1)$ . For convenience, write  $\bar{f} = \max\{f_1, \dots, f_\ell\} = \max_i f_i$ . We have

$$\begin{aligned} \bar{f}(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) &= \max_i f_i(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \\ &\leq \max_i (\alpha f_i(\mathbf{x}) + (1 - \alpha)f_i(\mathbf{y})) \quad \text{by convexity of each } f_i \\ &\leq \alpha \max_i f_i(\mathbf{x}) + (1 - \alpha) \max_i f_i(\mathbf{y}) \quad \text{by property of max} \\ &= \alpha \bar{f}(\mathbf{x}) + (1 - \alpha)\bar{f}(\mathbf{y}) \end{aligned}$$

which implies that  $\bar{f}$  is convex.

#### 22.7

$\Rightarrow$ : This is true by definition.

$\Leftarrow$ : Let  $\mathbf{d} \in \mathbb{R}^n$  be given. We want to show that  $\mathbf{d}^\top \mathbf{Q} \mathbf{d} \geq 0$ . Now, fix some vector  $\mathbf{x} \in \Omega$ . Since  $\Omega$  is open, there exists  $\alpha \neq 0$  such that  $\mathbf{y} = \mathbf{x} - \alpha \mathbf{d} \in \Omega$ . By assumption,

$$0 \leq (\mathbf{y} - \mathbf{x})^\top \mathbf{Q}(\mathbf{y} - \mathbf{x}) = \alpha^2 \mathbf{d}^\top \mathbf{Q} \mathbf{d}$$

which implies that  $\mathbf{d}^\top \mathbf{Q} \mathbf{d} \geq 0$ .

#### 22.8

Yes, the problem is a convex optimization problem.

First we show that the objective function  $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$  is convex. We write

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top (\mathbf{A}^\top \mathbf{A}) \mathbf{x} - (\mathbf{b}^\top \mathbf{A}) \mathbf{x} + \text{constant}$$

which is a quadratic function with Hessian  $\mathbf{A}^\top \mathbf{A}$ . Since the Hessian  $\mathbf{A}^\top \mathbf{A}$  is positive semidefinite, the objective function  $f$  is convex.

Next we show that the constraint set is convex. Consider two feasible points  $\mathbf{x}$  and  $\mathbf{y}$ , and let  $\lambda \in (0, 1)$ . Then,  $\mathbf{x}$  and  $\mathbf{y}$  satisfy  $\mathbf{e}^\top \mathbf{x} = 1, \mathbf{x} \geq \mathbf{0}$  and  $\mathbf{e}^\top \mathbf{y} = 1, \mathbf{y} \geq \mathbf{0}$ , respectively. We have

$$\mathbf{e}^\top (\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) = \lambda \mathbf{e}^\top \mathbf{x} + (1 - \lambda) \mathbf{e}^\top \mathbf{y} = \lambda + (1 - \lambda) = 1.$$

Moreover, each component of  $\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}$  is given by  $\lambda x_i + (1 - \lambda) y_i$ , which is nonnegative because every term here is nonnegative. Hence,  $\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}$  is a feasible point, which shows that the constraint set is convex.

## 22.9

We need to show that  $\Omega$  is a convex set, and  $f$  is a convex function on  $\Omega$ .

To show that  $\Omega$  is a convex set, we need to show that for any  $\mathbf{y}, \mathbf{z} \in \Omega$  and  $\alpha \in (0, 1)$ , we have  $\alpha \mathbf{y} + (1 - \alpha) \mathbf{z} \in \Omega$ . Let  $\mathbf{y}, \mathbf{z} \in \Omega$  and  $\alpha \in (0, 1)$ . Thus,  $y_1 = y_2 \geq 0$  and  $z_1 = z_2 \geq 0$ . Hence,

$$\mathbf{x} = \alpha \mathbf{y} + (1 - \alpha) \mathbf{z} = \begin{bmatrix} \alpha y_1 + (1 - \alpha) z_1 \\ \alpha y_2 + (1 - \alpha) z_2 \end{bmatrix}$$

Now,

$$x_1 = \alpha y_1 + (1 - \alpha) z_1 = \alpha y_2 + (1 - \alpha) z_2 = x_2,$$

and since  $\alpha, 1 - \alpha \geq 0$ ,

$$x_1 \geq 0.$$

Hence,  $\mathbf{x} \in \Omega$  and therefore  $\Omega$  is convex.

To show that  $f$  is convex on  $\Omega$ , we need to show that for any  $\mathbf{y}, \mathbf{z} \in \Omega$  and  $\alpha \in [0, 1]$ ,  $f(\alpha \mathbf{y} + (1 - \alpha) \mathbf{z}) \leq \alpha f(\mathbf{y}) + (1 - \alpha) f(\mathbf{z})$ . Let  $\mathbf{y}, \mathbf{z} \in \Omega$  and  $\alpha \in [0, 1]$ . Thus,  $y_1 = y_2 \geq 0$  and  $z_1 = z_2 \geq 0$ , so that  $f(\mathbf{y}) = y_1^3$  and  $f(\mathbf{z}) = z_1^3$ . Also,  $\alpha^3 \leq \alpha$  and  $(1 - \alpha)^3 \leq (1 - \alpha)$ . We have

$$\begin{aligned} f(\alpha \mathbf{y} + (1 - \alpha) \mathbf{z}) &= (\alpha y_1 + (1 - \alpha) z_1)^3 \\ &= \alpha^3 y_1^3 + (1 - \alpha)^3 z_1^3 + 3\alpha^2 x_1^2 (1 - \alpha) y_1 + 3\alpha x_1 (1 - \alpha)^2 y_1^2 \\ &\leq \alpha y_1^3 + (1 - \alpha) z_1^3 + \max(y_1, z_1)(\alpha^3 - \alpha + (1 - \alpha)^3 - (1 - \alpha)) \\ &\quad + 3\alpha^2 (1 - \alpha) + 3\alpha (1 - \alpha)^2 \\ &= \alpha y_1^3 + (1 - \alpha) z_1^3 \\ &= \alpha f(\mathbf{y}) + (1 - \alpha) f(\mathbf{z}). \end{aligned}$$

Hence,  $f$  is convex.

## 22.10

Since the problem is a convex optimization problem, we know for sure that any point of the form  $\alpha \mathbf{y} + (1 - \alpha) \mathbf{z}$ ,  $\alpha \in (0, 1)$ , is a global minimizer. However, any other point may or may not be a minimizer. Hence, the largest set of points  $G \subset \Omega$  for which we can be sure that every point in  $G$  is a global minimizer, is given by

$$G = \{\alpha \mathbf{y} + (1 - \alpha) \mathbf{z} : 0 \leq \alpha \leq 1\}.$$

## 22.11

a. Let  $f$  be the objective function and  $\Omega$  the constraint set. Consider the set  $\Gamma = \{\mathbf{x} \in \Omega : f(\mathbf{x}) \leq 1\}$ . This set contains all three of the given points. Moreover, by Lemma 22.1,  $\Gamma$  is convex. Now, if we take the average of the first two points (which is a convex combination of them), the resulting point  $(1/2)[1, 0, 0]^\top + (1/2)[0, 1, 0]^\top = (1/2)[1, 1, 0]^\top$  is in  $\Gamma$ , because  $\Gamma$  is convex. Similarly, the point  $(2/3)(1/2)[1, 1, 0]^\top + (1/3)[0, 0, 1]^\top = (1/3)[1, 1, 1]^\top$  is also in  $\Gamma$ , because  $\Gamma$  is convex. Hence, the objective function value of  $(1/3)[1, 1, 1]^\top$  must be  $\leq 1$ .

b. If the three points are all global minimizers, then the point  $(1/3)[1, 1, 1]^\top$ , which must cannot have higher objective function value than the given three points (by part a), must also be a global minimizer.

## 22.12

a. The Lagrange condition for the problem is given by:

$$\begin{aligned} \mathbf{x}^\top \mathbf{Q} + \boldsymbol{\lambda}^\top \mathbf{A} &= \mathbf{0} \\ \mathbf{A} \mathbf{x} &= \mathbf{b}. \end{aligned}$$

From the first equation above, we obtain

$$\mathbf{x} = \mathbf{Q}^{-1} \mathbf{A}^\top \boldsymbol{\lambda}.$$

Applying the second equation (constraint on  $\mathbf{x}$ ), we have

$$(\mathbf{A}\mathbf{Q}^{-1}\mathbf{A}^\top)\boldsymbol{\lambda} = \mathbf{b}.$$

Since  $\text{rank } \mathbf{A} = m$ , the matrix  $\mathbf{A}\mathbf{Q}^{-1}\mathbf{A}^\top$  is invertible. Therefore, the only solution to the Lagrange condition is

$$\mathbf{x} = \mathbf{Q}^{-1}\mathbf{A}^\top(\mathbf{A}\mathbf{Q}^{-1}\mathbf{A}^\top)^{-1}\mathbf{b}.$$

b. The point in part a above is a global minimizer because the problem is a convex optimization problem (by problem 1, the constraint set is convex; the objective function is convex because its Hessian,  $\mathbf{Q}$ , is positive definite).

### 22.13

By Theorem 22.4, for all  $\mathbf{x} \in \Omega$ , we have

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) + Df(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*).$$

Substituting  $Df(\mathbf{x}^*)$  from the equation  $Df(\mathbf{x}^*) + \sum_{j \in J(\mathbf{x}^*)} \mu_j^* \mathbf{a}_j^\top = \mathbf{0}^\top$  into the above inequality yields

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) - \sum_{j \in J(\mathbf{x}^*)} \mu_j^* \mathbf{a}_j^\top (\mathbf{x} - \mathbf{x}^*).$$

Observe that for each  $j \in J(\mathbf{x}^*)$ ,

$$\mathbf{a}_j^\top \mathbf{x}^* + b_j = 0,$$

and for each  $\mathbf{x} \in \Omega$ ,

$$\mathbf{a}_j^\top \mathbf{x} + b_j \geq 0.$$

Hence, for each  $j \in J(\mathbf{x}^*)$ ,

$$\mathbf{a}_j^\top (\mathbf{x} - \mathbf{x}^*) \geq 0.$$

Since  $\mu_j^* \leq 0$ , we get

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) - \sum_{j \in J(\mathbf{x}^*)} \mu_j^* \mathbf{a}_j^\top (\mathbf{x} - \mathbf{x}^*) \geq f(\mathbf{x}^*)$$

and the proof is completed.

### 22.14

a. Let  $\Omega = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} \geq b\}$ ,  $\mathbf{x}_1, \mathbf{x}_2 \in \Omega$ , and  $\lambda \in [0, 1]$ . Then,  $\mathbf{a}^\top \mathbf{x}_1 \geq b$  and  $\mathbf{a}^\top \mathbf{x}_2 \geq b$ . Therefore,

$$\begin{aligned} \mathbf{a}^\top (\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) &= \lambda \mathbf{a}^\top \mathbf{x}_1 + (1 - \lambda) \mathbf{a}^\top \mathbf{x}_2 \\ &\geq \lambda b + (1 - \lambda) b \\ &= b \end{aligned}$$

which means that  $\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in \Omega$ . Hence,  $\Omega$  is a convex set.

b. Rewrite the problem as

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}) \\ &\text{subject to} && g(\mathbf{x}) \leq 0 \end{aligned}$$

where  $f(\mathbf{x}) = \|\mathbf{x}\|^2$  and  $g(\mathbf{x}) = b - \mathbf{a}^\top \mathbf{x}$ . Now,  $\nabla g(\mathbf{x}) = -\mathbf{a} \neq \mathbf{0}$ . Therefore, any feasible point is regular. By the Karush-Kuhn-Tucker theorem, there exists  $\mu^* \geq 0$  such that

$$\begin{aligned} 2\mathbf{x}^* - \mu^* \mathbf{a} &= \mathbf{0} \\ \mu^* (b - \mathbf{a}^\top \mathbf{x}^*) &= 0. \end{aligned}$$

Since  $\mathbf{x}^*$  is a feasible point, then  $\mathbf{x}^* \neq \mathbf{0}$ . Therefore, by the first equation, we see that  $\mu^* \neq 0$ . The second equation then implies that  $b - \mathbf{a}^\top \mathbf{x}^* = 0$ .



c. By the first Karush-Kuhn-Tucker equation, we have  $\mathbf{x}^* = \mu^* \mathbf{a}/2$ . Since  $\mathbf{a}^\top \mathbf{x}^* = b$ , then  $\mu^* \mathbf{a}^\top \mathbf{a}/2 = \mathbf{a}^\top \mathbf{x}^* = b$ , and therefore  $\mu^* = 2b/\|\mathbf{a}\|^2$ . Since  $\mathbf{x}^* = \mu^* \mathbf{a}/2$  then  $\mathbf{x}^*$  is uniquely given by  $\mathbf{x}^* = b\mathbf{a}/\|\mathbf{a}\|^2$ .

### 22.15

a. Let  $f(\mathbf{x}) = \mathbf{c}^\top \mathbf{x}$  and  $\Omega = \{\mathbf{x} : \mathbf{x} \geq \mathbf{0}\}$ . Suppose  $\mathbf{x}, \mathbf{y} \in \Omega$ , and  $\alpha \in (0, 1)$ . Then,  $\mathbf{x}, \mathbf{y} \geq \mathbf{0}$ . Hence,  $\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \geq \mathbf{0}$ , which means  $\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \in \Omega$ . Furthermore,

$$\mathbf{c}^\top (\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) = \alpha\mathbf{c}^\top \mathbf{x} + (1 - \alpha)\mathbf{c}^\top \mathbf{y}.$$

Therefore,  $f$  is convex. Hence, the problem is a convex programming problem.

b.  $\Rightarrow$ : We use contraposition. Suppose  $c_i < 0$  for some  $i$ . Let  $\mathbf{d} = [0, \dots, 1, \dots, 0]$ , where 1 appears in the  $i$ th component. Clearly  $\mathbf{d}$  is a feasible direction for any point  $\mathbf{x}^* \geq \mathbf{0}$ . However,  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) = \mathbf{d}^\top \mathbf{c} = c_i < 0$ . Therefore, the FONC does not hold, and any point  $\mathbf{x}^* \geq \mathbf{0}$  cannot be a minimizer.

$\Leftarrow$ : Suppose  $\mathbf{c} \geq \mathbf{0}$ . Let  $\mathbf{x}^* = \mathbf{0}$ , and  $\mathbf{d}$  a feasible direction at  $\mathbf{x}^*$ . Then,  $\mathbf{d} \geq \mathbf{0}$ . Hence,  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) \geq \mathbf{0}$ . Therefore, by Theorem 22.7,  $\mathbf{x}^*$  is a solution.

The above also proves that if a solution exists, then  $\mathbf{0}$  is a solution.

c. Write  $\mathbf{g}(\mathbf{x}) = -\mathbf{x}$  so that the constraint can be expressed as  $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}$ .

$\Rightarrow$ : We have  $D\mathbf{g}(\mathbf{x}) = -\mathbf{I}$ , which has full rank. Therefore, any point is regular. Suppose a solution  $\mathbf{x}^*$  exists. Then, by the KKT theorem, there exists  $\mu^* \geq \mathbf{0}$  such that  $\mathbf{c}^\top - \mu^{*\top} = \mathbf{0}^\top$  and  $\mu^{*\top} \mathbf{x}^* = 0$ . Hence,  $\mathbf{c} = \mu^* \geq \mathbf{0}$ .

$\Leftarrow$ : Suppose  $\mathbf{c} \geq \mathbf{0}$ . Let  $\mathbf{x}^* = \mathbf{0}$  and  $\mu^* = \mathbf{c}$ . Then,  $\mu^* \geq \mathbf{0}$ ,  $\mathbf{c}^\top - \mu^{*\top} = \mathbf{0}^\top$ , and  $\mu^{*\top} \mathbf{x}^* = 0$ , i.e., the KKT condition is satisfied. By part a,  $\mathbf{x}^*$  is a solution to the problem.

The above also proves that if a solution exists, then  $\mathbf{0}$  is a solution.

### 22.16

a. The standard form problem is

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

which can be written as

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{h}(\mathbf{x}) = \mathbf{0}, \\ & && \mathbf{g}(\mathbf{x}) \leq \mathbf{0}, \end{aligned}$$

where  $f(\mathbf{x}) = \mathbf{c}^\top \mathbf{x}$ ,  $\mathbf{h}(\mathbf{x}) = \mathbf{A}\mathbf{x} - \mathbf{b}$ , and  $\mathbf{g}(\mathbf{x}) = -\mathbf{x}$ . Thus, we have  $Df(\mathbf{x}) = \mathbf{c}^\top$ ,  $D\mathbf{h}(\mathbf{x}) = \mathbf{A}$ , and  $D\mathbf{g}(\mathbf{x}) = -\mathbf{I}$ . The Karush-Kuhn-Tucker conditions for the above problem has the form:

$$\begin{aligned} \mu^* & \geq \mathbf{0} \\ \mathbf{c}^\top + \lambda^{*\top} \mathbf{A} - \mu^{*\top} & = \mathbf{0}^\top \\ \mu^{*\top} \mathbf{x}^* & = 0. \end{aligned}$$

b. The Karush-Kuhn-Tucker conditions are sufficient for optimality in this case because the problem is a convex optimization problem, i.e., the objective function is a convex function, and the feasible set is a convex set.

c. The dual problem is

$$\begin{aligned} & \text{maximize} && \lambda^\top \mathbf{b} \\ & \text{subject to} && \lambda^\top \mathbf{A} \leq \mathbf{c}^\top. \end{aligned}$$

c. Let

$$\mu^* = \mathbf{c}^\top - \lambda^{*\top} \mathbf{A}.$$

Since  $\lambda^*$  is feasible for the dual, we have  $\mu^* \geq \mathbf{0}$ . Rewriting the above equation, we get

$$\mathbf{c}^\top - \lambda^{*\top} \mathbf{A} - \mu^* = \mathbf{0}.$$

The Complementary Slackness condition  $(\mathbf{c}^\top - \lambda^{*\top} \mathbf{A})\mathbf{x}^* = \mathbf{0}$  can be written as  $\mu^{*\top} \mathbf{x}^* = \mathbf{0}$ . Therefore, the Karush-Kuhn-Tucker conditions hold. By part b,  $\mathbf{x}^*$  is optimal.

## 22.17

a. We can treat  $\mathbf{s}^{(1)}$  and  $\mathbf{s}^{(2)}$  as vectors in  $\mathbb{R}^n$ . We have

$$S_a = \{\mathbf{s} : \mathbf{s} = x_1 \mathbf{s}^{(1)} + x_2 \mathbf{s}^{(2)}, x_1, x_2 \in \mathbb{R}, \quad s_i \geq a, i = 1, \dots, n\}.$$

Let  $\mathbf{a} = [a, \dots, a]^\top$ . The optimization problem is:

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2}(x_1^2 + x_2^2) \\ \text{subject to} \quad & x_1 \mathbf{s}^{(1)} + x_2 \mathbf{s}^{(2)} \geq \mathbf{a}. \end{aligned}$$

b. The KKT conditions are:

$$\begin{aligned} x_1 - \mu^\top \mathbf{s}^{(1)} &= 0 \\ x_2 - \mu^\top \mathbf{s}^{(2)} &= 0 \\ \mu^\top (x_1 \mathbf{s}^{(1)} + x_2 \mathbf{s}^{(2)} - \mathbf{a}) &= \mathbf{0} \\ \mu &\geq \mathbf{0} \\ x_1 \mathbf{s}^{(1)} + x_2 \mathbf{s}^{(2)} &\geq \mathbf{a}. \end{aligned}$$

c. Yes, because the Hessian of the Lagrangian is  $\mathbf{I}$  (identity), which is positive definite.

d. Yes, this is a convex optimization problem. The objective function is quadratic, with identity Hessian (hence positive definite). The constraint set is of the form  $\mathbf{A}\mathbf{x} \geq \mathbf{a}$ , and hence is a linear variety.

## 22.18

a. We first show that the set of probability vectors

$$\Omega = \{\mathbf{q} \in \mathbb{R}^n : q_1 + \dots + q_n = 1, \quad q_i > 0, \quad i = 1, \dots, n\}$$

is a convex set. Let  $\mathbf{y}, \mathbf{z} \in \Omega$ , so  $y_1 + \dots + y_n = 1$ ,  $y_i > 0$ ,  $z_1 + \dots + z_n = 1$ , and  $z_i > 0$ . Let  $\alpha \in (0, 1)$  and  $\mathbf{x} = \alpha \mathbf{y} + (1 - \alpha) \mathbf{z}$ . We have

$$\begin{aligned} x_1 + \dots + x_n &= \alpha y_1 + (1 - \alpha) z_1 + \dots + \alpha y_n + (1 - \alpha) z_n \\ &= \alpha(y_1 + \dots + y_n) + (1 - \alpha)(z_1 + \dots + z_n) \\ &= \alpha + (1 - \alpha) \\ &= 1. \end{aligned}$$

Also, because  $y_i > 0$ ,  $z_i > 0$ ,  $\alpha > 0$ , and  $1 - \alpha > 0$ , we conclude that  $x_i > 0$ . Thus,  $\mathbf{x} \in \Omega$ , which shows that  $\Omega$  is convex.

b. We next show that the function  $f$  is a convex function on  $\Omega$ . For this, we compute

$$\mathbf{F}(\mathbf{q}) = \begin{bmatrix} \frac{p_1}{q_1^2} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{p_n}{q_n^2} \end{bmatrix},$$

which shows that  $\mathbf{F}(\mathbf{q}) > 0$  for all  $\mathbf{q}$  in the open set  $\{\mathbf{q} : q_i > 0, \quad i = 1, \dots, n\}$ , which contains  $\Omega$ . Therefore,  $f$  is convex on  $\Omega$ .

c. Fix a probability vector  $\mathbf{p}$ . Consider the optimization problem

$$\begin{aligned} & \text{minimize} && p_1 \log \left( \frac{p_1}{x_1} \right) + \cdots + p_n \log \left( \frac{p_n}{x_n} \right) \\ & \text{subject to} && x_1 + \cdots + x_n = 1 \\ & && x_i > 0, \quad i = 1, \dots, n \end{aligned}$$

By parts a and b, the problem is a convex optimization problem. We ignore the constraint  $x_i > 0$  and write down the Lagrange conditions for the equality-constraint problem:

$$\begin{aligned} -\frac{p_i}{x_i^*} + \lambda^* &= 0, \quad i = 1, \dots, n \\ x_1^* + \cdots + x_n^* &= 1. \end{aligned}$$

Rewrite the first set of equations as  $x_i^* = p_i/\lambda^*$ . Combining this with the constraint and the fact that  $p_1 + \cdots + p_n = 1$ , we obtain  $\lambda^* = 1$ , which means that  $x_i^* = p_i$ . Therefore, the unique global minimizer is  $\mathbf{x}^* = \mathbf{p}$ .

Note that  $f(\mathbf{x}^*) = 0$ . Hence, we conclude that  $f(\mathbf{x}) \geq 0$  for all  $\mathbf{x} \in \Omega$ . Moreover,  $f(\mathbf{x}) = 0$  if and only if  $\mathbf{x}^* = 0$ . This proves the required result.

d. Given two probability vectors  $\mathbf{p}$  and  $\mathbf{q}$ , the number

$$D(\mathbf{p}, \mathbf{q}) = p_1 \log \left( \frac{p_1}{q_1} \right) + \cdots + p_n \log \left( \frac{p_n}{q_n} \right)$$

is called the *relative entropy* (or *Kullback-Liebler divergence*) between  $\mathbf{p}$  and  $\mathbf{q}$ . It is used in information theory to measure the “distance” between two probability vectors. The result of part c justifies the use of  $D$  as a measure of “distance” (although  $D$  is not a *metric* because it is not symmetric).

## 22.19

We claim that a solution exists, and that it is unique. To prove the first claim, choose  $\varepsilon > 0$  such that there exists  $\mathbf{x} \in \Omega$  satisfying  $\|\mathbf{x} - \mathbf{z}\| < \varepsilon$ . Consider the modified problem

$$\begin{aligned} & \text{minimize} && \|\mathbf{x} - \mathbf{z}\| \\ & \text{subject to} && \mathbf{x} \in \Omega \cap \{\mathbf{y} : \|\mathbf{y} - \mathbf{z}\| \leq \varepsilon\}. \end{aligned}$$

If this modified problem has a solution, then clearly so does the original problem. The objective function here is continuous, and the constraint set is closed and bounded. Hence, by Weierstrass’s Theorem, a solution to the problem exists.

Let  $f$  be the objective function. Next, we show that  $f$  is convex (and hence the problem is a convex optimization problem). Let  $\mathbf{x}, \mathbf{y} \in \Omega$  and  $\alpha \in (0, 1)$ . Then,

$$\begin{aligned} f(\alpha\mathbf{x} + (1-\alpha)\mathbf{y}) &= \|\alpha\mathbf{x} + (1-\alpha)\mathbf{y} - \mathbf{z}\| \\ &= \|\alpha(\mathbf{x} - \mathbf{z}) + (1-\alpha)(\mathbf{y} - \mathbf{z})\| \\ &\leq \alpha\|\mathbf{x} - \mathbf{z}\| + (1-\alpha)\|\mathbf{y} - \mathbf{z}\| \\ &= \alpha f(\mathbf{x}) + (1-\alpha)f(\mathbf{y}), \end{aligned}$$

which shows that  $f$  is convex.

To prove uniqueness, let  $\mathbf{x}_1$  and  $\mathbf{x}_2$  be solutions to the problem. Then, by convexity,  $\mathbf{x}_3 = (\mathbf{x}_1 + \mathbf{x}_2)/2$  is also a solution. But

$$\begin{aligned} \|\mathbf{x}_3 - \mathbf{z}\| &= \left\| \frac{\mathbf{x}_1 + \mathbf{x}_2}{2} - \mathbf{z} \right\| \\ &= \left\| \frac{\mathbf{x}_1 - \mathbf{z}}{2} + \frac{\mathbf{x}_2 - \mathbf{z}}{2} \right\| \\ &\leq \frac{1}{2}(\|\mathbf{x}_1 - \mathbf{z}\| + \|\mathbf{x}_2 - \mathbf{z}\|) \\ &= \|\mathbf{x}_3 - \mathbf{z}\|, \end{aligned}$$

from which we conclude that the triangle inequality above holds with equality, implying that  $\mathbf{x}_1 - \mathbf{z} = \alpha(\mathbf{x}_2 - \mathbf{z})$  for some  $\alpha \geq 0$ . Because  $\|\mathbf{x}_1 - \mathbf{z}\| = \|\mathbf{x}_2 - \mathbf{z}\| = \alpha\|\mathbf{x}_1 - \mathbf{z}\|$ , we have  $\alpha = 1$ . From this, we obtain  $\mathbf{x}_1 = \mathbf{x}_2$ , which proves uniqueness.

## 22.20

a. Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\mathbf{B} \in \mathbb{R}^{n \times n}$  be symmetric and  $\mathbf{A} \geq 0$ ,  $\mathbf{B} \geq 0$ . Fix  $\alpha \in (0, 1)$ ,  $\mathbf{x} \in \mathbb{R}^n$ , and let  $\mathbf{C} = \alpha\mathbf{A} + (1 - \alpha)\mathbf{B}$ . Then,

$$\begin{aligned}\mathbf{x}^\top \mathbf{C} \mathbf{x} &= \mathbf{x}^\top [\alpha\mathbf{A} + (1 - \alpha)\mathbf{B}] \mathbf{x} \\ &= \alpha \mathbf{x}^\top \mathbf{A} \mathbf{x} + (1 - \alpha) \mathbf{x}^\top \mathbf{B} \mathbf{x}.\end{aligned}$$

Since  $\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$ ,  $\mathbf{x}^\top \mathbf{B} \mathbf{x} \geq 0$ , and  $\alpha, (1 - \alpha) > 0$  by assumption, then  $\mathbf{x}^\top \mathbf{C} \mathbf{x} \geq 0$ , which proves the required result.

b. We first show that the constraint set  $\Omega = \{\mathbf{x} : \mathbf{F}_0 + \sum_{j=1}^n x_j \mathbf{F}_j \geq 0\}$  is convex. So, let  $\mathbf{x}, \mathbf{y} \in \Omega$  and  $\alpha \in (0, 1)$ . Let  $\mathbf{z} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{y}$ . Then,

$$\begin{aligned}\mathbf{F}_0 + \sum_{j=1}^n z_j \mathbf{F}_j &= \mathbf{F}_0 + \sum_{j=1}^n [\alpha x_j + (1 - \alpha)y_j] \mathbf{F}_j \\ &= \mathbf{F}_0 + \alpha \sum_{j=1}^n x_j \mathbf{F}_j + (1 - \alpha) \sum_{j=1}^n y_j \mathbf{F}_j \\ &= \alpha [\mathbf{F}_0 + \sum_{j=1}^n x_j \mathbf{F}_j] + (1 - \alpha) [\mathbf{F}_0 + \sum_{j=1}^n y_j \mathbf{F}_j].\end{aligned}$$

By assumption, we have

$$\begin{aligned}\mathbf{F}_0 + \sum_{j=1}^n x_j \mathbf{F}_j &\geq 0 \\ \mathbf{F}_0 + \sum_{j=1}^n y_j \mathbf{F}_j &\geq 0.\end{aligned}$$

By part c, we conclude that

$$\mathbf{F}_0 + \sum_{j=1}^n z_j \mathbf{F}_j \geq 0,$$

which implies that  $\mathbf{z} \in \Omega$ .

To show that the objective function  $f(\mathbf{x}) = \mathbf{c}^\top \mathbf{x}$  is convex on  $\Omega$ , let  $\mathbf{x}, \mathbf{y} \in \Omega$  and  $\alpha \in (0, 1)$ . Then,

$$\begin{aligned}f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) &= \mathbf{c}^\top (\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \\ &= \alpha \mathbf{c}^\top \mathbf{x} + (1 - \alpha) \mathbf{c}^\top \mathbf{y} \\ &= \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y})\end{aligned}$$

which shows that  $f$  is convex.

c. The objective function is already in the required form. To rewrite the constraint, let  $a_{i,j}$  be the  $(i, j)$ th entry of  $\mathbf{A}$ ,  $i = 1, \dots, m$ ,  $j = 1, \dots, n$ . Then, the constraint  $\mathbf{A}\mathbf{x} \geq \mathbf{b}$  can be written as

$$a_{i,1}x_1 + a_{i,2}x_2 + \dots + a_{i,n}x_n \geq b_i, \quad i = 1, \dots, m$$

Now form the diagonal matrices

$$\begin{aligned}\mathbf{F}_0 &= \text{diag}\{-b_1, \dots, -b_m\} \\ \mathbf{F}_j &= \text{diag}\{a_{1,j}, \dots, a_{m,j}\}, \quad j = 1, \dots, n.\end{aligned}$$

Note that a diagonal matrix is positive semidefinite if and only if every diagonal element is nonnegative. Hence, the constraint  $\mathbf{Ax} \geq \mathbf{b}$  can be written as  $\mathbf{F}_0 + \sum_{j=1}^n x_j \mathbf{F}_j \geq 0$ . The left hand side is a diagonal matrix, and the  $i$ th diagonal element is simply  $-b_i + a_{i,1}x_1 + a_{i,2}x_2 + \cdots + a_{i,n}x_n$ .

## 22.21

a. We have

$$\Omega = \{\mathbf{x} : x_1 + \cdots + x_n = 1; x_1, \dots, x_n > 0; x_1 \geq 2x_i, i = 2, \dots, n\}.$$

So let  $\mathbf{x}, \mathbf{y} \in \Omega$  and  $\alpha \in (0, 1)$ . Consider  $\mathbf{z} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{y}$ . We have

$$\begin{aligned} z_1 + \cdots + z_n &= \alpha x_1 + (1 - \alpha)y_1 + \cdots + \alpha x_n + (1 - \alpha)y_n \\ &= \alpha(x_1 + \cdots + x_n) + (1 - \alpha)(y_1 + \cdots + y_n) \\ &= \alpha + 1 - \alpha \\ &= 1. \end{aligned}$$

Moreover, for each  $i$ , because  $x_i > 0$ ,  $y_i > 0$ ,  $\alpha > 0$  and  $1 - \alpha > 0$ , we have  $z_i > 0$ . Finally, for each  $i$ ,

$$z_1 = \alpha x_1 + (1 - \alpha)y_1 \geq \alpha 2x_i + (1 - \alpha)2y_i = 2z_i.$$

Hence,  $\mathbf{z} \in \Omega$ , which implies that  $\Omega$  is convex.

b. We first show that the negative of the objective function is convex. For this, we will compute its Hessian, which turns out to be a diagonal matrix with  $i$ th diagonal entry  $1/x_i^2$ , which is strictly positive. Hence, the Hessian is positive definite, which implies that the negative of the objective function is convex.

Combining the above with part a, we conclude that the problem is a convex optimization problem. Hence, the FONC (for set constraints) is necessary and sufficient. Let  $\mathbf{x}$  be a given allocation. The FONC at  $\mathbf{x}$  is  $\mathbf{d}^\top \nabla f(\mathbf{x}) \geq 0$  for all feasible directions  $\mathbf{d}$  at  $\mathbf{x}$ . But because  $\Omega$  is convex, the FONC can be written as  $(\mathbf{y} - \mathbf{x})^\top \nabla f(\mathbf{x}) \geq 0$  for all  $\mathbf{y} \in \Omega$ . Computing  $\nabla f(\mathbf{x})$  for  $f(\mathbf{x}) = -\sum_{i=1}^n \log(x_i)$ , we get the proportional fairness condition.

## 22.22

a. We rewrite the problem into a minimization problem by multiplying the objective function by  $-1$ . Thus, the new objective function is the sum of the functions  $-U_i$ . Because each  $U_i$  is concave,  $-U_i$  is convex, and hence their sum is convex.

To show that the constraint set  $\Omega = \{\mathbf{x} : \mathbf{e}^\top \mathbf{x} \leq C\}$  (where  $\mathbf{e} = [1, \dots, 1]^\top$ ) is convex, let  $\mathbf{x}_1, \mathbf{x}_2 \in \Omega$ , and  $\lambda \in [0, 1]$ . Then,  $\mathbf{e}^\top \mathbf{x}_1 \leq C$  and  $\mathbf{e}^\top \mathbf{x}_2 \leq C$ . Therefore,

$$\begin{aligned} \mathbf{e}^\top (\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) &= \lambda \mathbf{e}^\top \mathbf{x}_1 + (1 - \lambda) \mathbf{e}^\top \mathbf{x}_2 \\ &\leq \lambda C + (1 - \lambda) C \\ &= C \end{aligned}$$

which means that  $\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in \Omega$ . Hence,  $\Omega$  is a convex set.

b. Because the problem is a convex optimization problem, the following KKT condition is necessary and sufficient for  $\mathbf{x}^*$  to be a global minimizers:

$$\begin{aligned} \mu^* &\geq 0 \\ -U'_i(x_i^*) + \mu^* &= 0, \quad i = 1, \dots, n \\ \mu^* \left( \sum_{i=1}^n x_i^* - C \right) &= 0 \\ \sum_{i=1}^n x_i^* &\leq C. \end{aligned}$$

Note that because  $-U_i(x_i) + \mu^* x_i$  is a convex function of  $x_i$ , the second line above can be written as  $x_i^* = \arg \max_x (U_i(x) - \mu^* x)$ .

c. Because each  $U_i$  is concave and increasing, we conclude that  $\sum_{i=1}^n x_i^* = C$ ; for otherwise we could increase some  $x_i^*$  and hence  $U_i(x_i^*)$  and also  $\sum_{i=1}^n U_i(x_i^*)$ , contradicting the optimality of  $\mathbf{x}^*$ .

### 22.23

First note that the optimization problem that you construct cannot be a convex problem (for otherwise, the FONC implies that  $\mathbf{x}^*$  is a global minimizer, which then implies that the SONC holds). Let  $f(\mathbf{x}) = x_2$ ,  $g(\mathbf{x}) = -(x_2 + x_1^2)$ , and  $\mathbf{x}^* = \mathbf{0}$ . Then,  $\nabla f(\mathbf{x}^*) = [0, 1]^\top$ . Any feasible direction  $\mathbf{d}$  at  $\mathbf{x}^*$  is of the form  $\mathbf{d} = [d_1, d_2]^\top$  with  $d_2 \geq 0$ . Hence,  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) \geq 0$ , which shows that the FONC holds.

Because  $\nabla g(\mathbf{x}^*) = -[0, 1]$  (so  $\mathbf{x}^*$  is regular), we see that if  $\mu^* = 1$ , then  $\nabla f(\mathbf{x}^*) + \mu^* \nabla g(\mathbf{x}^*) = \mathbf{0}$ , and so the KKT condition holds.

Because  $\mathbf{F}(\mathbf{x}^*) = \mathbf{O}$ , the SONC for set constraint  $\Omega$  holds. However,

$$\mathbf{L}(\mathbf{x}^*, \mu^*) = \mathbf{O} + \mu^* \begin{bmatrix} -2 & 0 \\ 0 & 0 \end{bmatrix} < \mathbf{0}$$

and  $T(\mathbf{x}^*) = \{\mathbf{y} : y_2 = 0\}$ , which shows that the SONC for inequality constraint  $g(\mathbf{x}) \leq 0$  does not hold.

### 22.24

a. Let  $\mathbf{x}_0$  and  $\boldsymbol{\mu}_0$  be feasible points in the primal and dual, respectively. Then,  $\mathbf{g}(\mathbf{x}_0) \leq \mathbf{0}$  and  $\boldsymbol{\mu}_0 \geq \mathbf{0}$ , and so  $\boldsymbol{\mu}_0^\top \mathbf{g}(\mathbf{x}_0) \leq 0$ . Hence,

$$\begin{aligned} f(\mathbf{x}_0) &\geq f(\mathbf{x}_0) + \boldsymbol{\mu}_0^\top \mathbf{g}(\mathbf{x}_0) \\ &= l(\mathbf{x}_0, \boldsymbol{\mu}_0) \\ &\geq \min_{\mathbf{x} \in \mathbb{R}^n} l(\mathbf{x}, \boldsymbol{\mu}_0) \\ &= q(\boldsymbol{\mu}_0). \end{aligned}$$

b. Suppose  $f(\mathbf{x}_0) = q(\boldsymbol{\mu}_0)$  for feasible points  $\mathbf{x}_0$  and  $\boldsymbol{\mu}_0$ . Let  $\mathbf{x}$  be any feasible point in the primal. Then, by part a,  $f(\mathbf{x}) \geq q(\boldsymbol{\mu}_0) = f(\mathbf{x}_0)$ . Hence  $\mathbf{x}_0$  is optimal in the primal.

Similarly, let  $\boldsymbol{\mu}$  be any feasible point in the dual. Then, by part a,  $q(\boldsymbol{\mu}) \leq f(\mathbf{x}_0) = q(\boldsymbol{\mu}_0)$ . Hence  $\boldsymbol{\mu}_0$  is optimal in the dual.

c. Let  $\mathbf{x}^*$  be optimal in the primal. Then, by the KKT Theorem, there exists  $\boldsymbol{\mu}^* \in \mathbb{R}^m$  such that

$$\begin{aligned} \nabla_x l(\mathbf{x}^*, \boldsymbol{\mu}^*) &= (Df(\mathbf{x}^*) + \boldsymbol{\mu}^{*\top} D\mathbf{g}(\mathbf{x}^*))^\top = \mathbf{0} \\ \boldsymbol{\mu}^{*\top} \mathbf{g}(\mathbf{x}^*) &= \mathbf{0} \\ \boldsymbol{\mu}^* &\geq \mathbf{0}. \end{aligned}$$

Therefore,  $\boldsymbol{\mu}^*$  is feasible in the dual. Further, we note that  $l(\cdot, \boldsymbol{\mu}^*)$  is a convex function (because  $f$  is convex,  $\boldsymbol{\mu}^{*\top} \mathbf{g}$  is convex being the sum of convex functions, and hence  $l$  is the sum of two convex functions). Hence, we have  $l(\mathbf{x}^*, \boldsymbol{\mu}^*) = \min_{\mathbf{x} \in \mathbb{R}^n} l(\mathbf{x}, \boldsymbol{\mu}^*)$ . Therefore,

$$\begin{aligned} q(\boldsymbol{\mu}^*) &= \min_{\mathbf{x} \in \mathbb{R}^n} l(\mathbf{x}, \boldsymbol{\mu}^*) \\ &= l(\mathbf{x}^*, \boldsymbol{\mu}^*) \\ &= f(\mathbf{x}^*) + \boldsymbol{\mu}^{*\top} \mathbf{g}(\mathbf{x}^*) \\ &= f(\mathbf{x}^*). \end{aligned}$$

By part b,  $\boldsymbol{\mu}^*$  is optimal in the dual.

### 22.25

a. The Schur complement of  $\mathbf{M}(1, 1)$  is

$$\begin{aligned} \Delta_{11} &= \mathbf{M}(2:3, 2:3) - \mathbf{M}(2:3, 1) \mathbf{M}(1, 1)^{-1} \mathbf{M}(1, 2:3) \\ &= \begin{bmatrix} 1 & 2 \\ 2 & 5 \end{bmatrix} - \begin{bmatrix} \gamma \\ -1 \end{bmatrix} \begin{bmatrix} \gamma & -1 \end{bmatrix} \\ &= \begin{bmatrix} 1 - \gamma^2 & 2 + \gamma \\ 2 + \gamma & 4 \end{bmatrix}. \end{aligned}$$

b. The Schur complement of  $\mathbf{M}(2: !3, 2: !3)$  is

$$\begin{aligned}
\Delta_{22} &= \mathbf{M}(1, 1) - \mathbf{M}(1, 2: !3) \mathbf{M}(2: !3, 2: 3)^{-1} \mathbf{M}(2: !3, 1) \\
&= 1 - \begin{bmatrix} \gamma & -1 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 2 & 5 \end{bmatrix}^{-1} \begin{bmatrix} \gamma \\ -1 \end{bmatrix} \\
&= 1 - \begin{bmatrix} \gamma & -1 \end{bmatrix} \begin{bmatrix} 5 & -2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} \gamma \\ -1 \end{bmatrix} \\
&= -\gamma(5\gamma + 4).
\end{aligned}$$

## 22.26

Let

$$\mathbf{P} = \begin{bmatrix} x_1 & x_2 \\ x_2 & x_3 \end{bmatrix}$$

and let

$$\mathbf{P}_1 = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{P}_2 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \text{and} \quad \mathbf{P}_3 = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}.$$

Then, we can represent the Lyapunov inequality  $\mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} < 0$  as

$$\begin{aligned}
\mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} &= x_1 (\mathbf{A}^\top \mathbf{P}_1 + \mathbf{P}_1 \mathbf{A}) + x_2 (\mathbf{A}^\top \mathbf{P}_2 + \mathbf{P}_2 \mathbf{A}) \\
&\quad + x_3 (\mathbf{A}^\top \mathbf{P}_3 + \mathbf{P}_3 \mathbf{A}) \\
&= -x_1 \mathbf{F}_1 - x_2 \mathbf{F}_2 - x_3 \mathbf{F}_3 \\
&< 0,
\end{aligned}$$

where

$$\mathbf{F}_i = -\mathbf{A}^\top \mathbf{P}_i - \mathbf{P}_i \mathbf{A}, \quad i = 1, 2, 3.$$

Equivalently,

$$\mathbf{P} = \mathbf{P}^\top \quad \text{and} \quad \mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} < 0$$

if and only if

$$\mathbf{F}(\mathbf{x}) = x_1 \mathbf{F}_1 + x_2 \mathbf{F}_2 + x_3 \mathbf{F}_3 > 0.$$

## 22.27

The quadratic inequality

$$\mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} + \mathbf{P} \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^\top \mathbf{P} < 0$$

can be equivalently represented by the following LMI:

$$\begin{bmatrix} -\mathbf{R} & \mathbf{B}^\top \mathbf{P} \\ \mathbf{P} \mathbf{B} & \mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} \end{bmatrix} < 0,$$

or as the following LMI:

$$\begin{bmatrix} \mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} & \mathbf{P} \mathbf{B} \\ \mathbf{B}^\top \mathbf{P} & -\mathbf{R} \end{bmatrix} < 0.$$

It is easy to verify, using Schur complements, that the above two LMIs are equivalent to the following quadratic inequality:

$$\begin{bmatrix} \mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} + \mathbf{P} \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^\top \mathbf{P} & \mathbf{O} \\ \mathbf{O} & -\mathbf{R} \end{bmatrix} < 0.$$

## 22.28

The MATLAB code is as follows:

```
A = [-0.9501    -0.4860    -0.4565  
     -0.2311    -0.8913    -0.0185  
     -0.6068    -0.7621    -0.8214];  
setlmis([]);  
P=lmivar(1,[3 1]);  
lmiterm([1 1 1 P],1,A,'s')  
lmiterm([2 1 1 0],0.1)  
lmiterm([-2 1 1 P],1,1)  
lmiterm([3 1 1 P],1,1)  
lmiterm([-3 1 1 0],1)  
lmis=getlmis;  
[tmin,xfas]=feasp(lmis);  
P=dec2mat(lmis,xfas,P)
```

## 23. Algorithms for Constrained Optimization

### 23.1

- a. By drawing a simple picture, it is easy to see that  $\Pi[\mathbf{x}] = \mathbf{x}/\|\mathbf{x}\|$ , provided  $\mathbf{x} \neq \mathbf{0}$ .  
b. By inspection, we see that the solutions are  $[0, 1]^\top$  and  $[0, -1]^\top$ . (Or use Rayleigh's inequality.)  
c. Now,

$$\mathbf{x}^{(k+1)} = \Pi[\mathbf{x}^{(k)} + \alpha \nabla f(\mathbf{x}^{(k)})] = \beta_k(\mathbf{x}^{(k)} + \alpha \mathbf{Q}\mathbf{x}^{(k)}) = \beta_k(\mathbf{I} + \alpha \mathbf{Q})\mathbf{x}^{(k)},$$

where  $\beta_k = 1/\|(\mathbf{I} + \alpha \mathbf{Q})\mathbf{x}^{(k)}\|$ . For the particular given form of  $\mathbf{Q}$ , we have

$$\begin{aligned}x_1^{(k+1)} &= \beta_k(1 + \alpha)x_1^{(k)} \\x_2^{(k+1)} &= \beta_k(1 + 2\alpha)x_2^{(k)}.\end{aligned}$$

Hence,

$$\mathbf{y}^{(k+1)} = \left( \frac{1 + \alpha}{1 + 2\alpha} \right) \mathbf{y}^{(k)}.$$

- d. Assuming  $x_2^{(0)} \neq 0$ ,  $\mathbf{y}^{(0)}$  is well defined. Hence, by part c, we can write

$$\mathbf{y}^{(k)} = \left( \frac{1 + \alpha}{1 + 2\alpha} \right)^k \mathbf{y}^{(0)}.$$

Because  $\alpha > 0$ ,

$$\frac{1 + \alpha}{1 + 2\alpha} < 1,$$

which implies that  $\mathbf{y}^{(k)} \rightarrow \mathbf{0}$ . But

$$\begin{aligned}1 &= \|\mathbf{x}^{(k)}\| \\&= \sqrt{(x_1^{(k)})^2 + (x_2^{(k)})^2} \\&= \sqrt{(x_2^{(k)})^2 \left( \frac{(x_1^{(k)})^2}{(x_2^{(k)})^2} + 1 \right)} \\&= |x_2^{(k)}| \sqrt{(y^{(k)})^2 + 1},\end{aligned}$$



which implies that

$$|x_2^{(k)}| = \frac{1}{\sqrt{(y^{(k)})^2 + 1}}.$$

Because  $y^{(k)} \rightarrow 0$ , we have  $|x_2^{(k)}| \rightarrow 1$ . By the expression for  $x_2^{(k+1)}$  in part c, we see that the sign of  $x_2^{(k)}$  does not change with  $k$ . Hence, we deduce that either  $x_2^{(k)} \rightarrow 1$  or  $x_2^{(k)} \rightarrow -1$ . This also implies that  $x_1^{(k)} \rightarrow 0$ . Hence,  $\mathbf{x}^{(k)}$  converges to a solution to the problem.

e. If  $x_2^{(0)} = 0$ , then  $x_2^{(k)} = 0$  for all  $k$ , which means that  $x_1^{(k)} = 1$  or  $-1$  for all  $k$ . In this case, the algorithm is stuck at the initial condition  $[1, 0]^\top$  or  $[-1, 0]^\top$  (which are in fact the minimizers).

### 23.2

a. Yes. To show: Suppose that  $\mathbf{x}^{(k)}$  is a global minimizer of the given problem. Then, for all  $\mathbf{x} \in \Omega$ ,  $\mathbf{x} \neq \mathbf{x}^{(k)}$ , we have  $\mathbf{c}^\top \mathbf{x} \geq \mathbf{c}^\top \mathbf{x}^{(k)}$ . Rewriting, we obtain  $\mathbf{c}^\top (\mathbf{x} - \mathbf{x}^{(k)}) \geq 0$ . Recall that

$$\begin{aligned} \Pi[\mathbf{x}^{(k)} - \nabla f(\mathbf{x}^{(k)})] &= \arg \min_{\mathbf{x} \in \Omega} \|\mathbf{x} - (\mathbf{x}^{(k)} - \nabla f(\mathbf{x}^{(k)}))\|^2 \\ &= \arg \min_{\mathbf{x} \in \Omega} \|\mathbf{x} - \mathbf{x}^{(k)} + \mathbf{c}\|^2. \end{aligned}$$

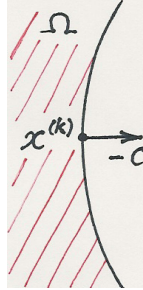
But, for any  $\mathbf{x} \in \Omega$ ,  $\mathbf{x} \neq \mathbf{x}^{(k)}$ ,

$$\begin{aligned} \|\mathbf{x} - \mathbf{x}^{(k)} + \mathbf{c}\|^2 &= \|\mathbf{x} - \mathbf{x}^{(k)}\|^2 + \|\mathbf{c}\|^2 + 2\mathbf{c}^\top (\mathbf{x} - \mathbf{x}^{(k)}) \\ &> \|\mathbf{c}\|^2, \end{aligned}$$

where we used the facts that  $\|\mathbf{x} - \mathbf{x}^{(k)}\|^2 > 0$  and  $\mathbf{c}^\top (\mathbf{x} - \mathbf{x}^{(k)}) \geq 0$ . On the other hand,  $\|\mathbf{x}^{(k)} - \mathbf{x}^{(k)} + \mathbf{c}\|^2 = \|\mathbf{c}\|^2$ . Hence,

$$\mathbf{x}^{(k+1)} = \Pi[\mathbf{x}^{(k)} - \nabla f(\mathbf{x}^{(k)})] = \mathbf{x}^{(k)}.$$

b. No. Counterexample:



### 23.3

a. Suppose  $\mathbf{x}^{(k)}$  satisfies the FONC. Then,  $\nabla f(\mathbf{x}^{(k)}) = \mathbf{0}$ . Hence,  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$ . Conversely, suppose  $\mathbf{x}^{(k)}$  does not satisfy the FONC. Then,  $\nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$ . Hence,  $\alpha_k > 0$ , and so  $\mathbf{x}^{(k+1)} \neq \mathbf{x}^{(k)}$ .

b. Case (i): Suppose  $\mathbf{x}^{(k)}$  is a corner point. Without loss of generality, take  $\mathbf{x}^{(k)} = [1, 1]^\top$ . (We can do this because any other corner point can be mapped to this point by changing variables  $x_i$  to  $-x_i$  as appropriate.) Note that any feasible direction  $\mathbf{d}$  at  $\mathbf{x}^{(k)} = [1, 1]^\top$  satisfies  $\mathbf{d} \leq \mathbf{0}$ . Therefore,

$$\begin{aligned} \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} &\Leftrightarrow -\nabla f(\mathbf{x}^{(k)}) \geq \mathbf{0} \\ &\Leftrightarrow \mathbf{d}^\top \nabla f(\mathbf{x}^{(k)}) \geq 0 \text{ for all feasible } \mathbf{d} \text{ at } \mathbf{x}^{(k)} \\ &\Leftrightarrow \mathbf{x}^{(k)} \text{ satisfies FONC.} \end{aligned}$$

Case (ii): Suppose  $\mathbf{x}^{(k)}$  is not a corner point (i.e., is an edge point). Without loss of generality, take  $\mathbf{x}^{(k)} \in \{\mathbf{x} : x_1 = 1, -1 < x_2 < 1\}$ . (We can do this because any other edge point can be mapped to this point by changing variables  $x_i$  to  $-x_i$  as appropriate.) Note that any feasible direction  $\mathbf{d}$  at  $\mathbf{x}^{(k)} \in \{\mathbf{x} : x_1 = 1, -1 < x_2 < 1\}$  satisfies  $d_1 \leq 0$ . Therefore,

$$\begin{aligned} \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} &\Leftrightarrow -\nabla f(\mathbf{x}^{(k)}) = [a, 0]^\top, \quad a > 0 \\ &\Leftrightarrow \mathbf{d}^\top \nabla f(\mathbf{x}^{(k)}) \geq 0 \text{ for all feasible } \mathbf{d} \text{ at } \mathbf{x}^{(k)} \\ &\Leftrightarrow \mathbf{x}^{(k)} \text{ satisfies FONC.} \end{aligned}$$

**23.4**

By definition of  $\Pi$ , we have

$$\begin{aligned}\Pi[\mathbf{x}_0 + \mathbf{y}] &= \arg \min_{\mathbf{x} \in \Omega} \|\mathbf{x} - (\mathbf{x}_0 + \mathbf{y})\| \\ &= \arg \min_{\mathbf{x} \in \Omega} \|(\mathbf{x} - \mathbf{x}_0) - \mathbf{y}\|.\end{aligned}$$

By Exercise 6.7, we can write

$$\arg \min_{\mathbf{x} \in \Omega} \|(\mathbf{x} - \mathbf{x}_0) - \mathbf{y}\| = \mathbf{x}_0 + \arg \min_{\mathbf{z} \in \mathcal{N}(\mathbf{A})} \|\mathbf{z} - \mathbf{y}\|.$$

The term  $\arg \min_{\mathbf{z} \in \mathcal{N}(\mathbf{A})} \|\mathbf{z} - \mathbf{y}\|$  is simply the orthogonal projection of  $\mathbf{y}$  onto  $\mathcal{N}(\mathbf{A})$ . By Exercise 6.7, we have

$$\arg \min_{\mathbf{z} \in \mathcal{N}(\mathbf{A})} \|\mathbf{z} - \mathbf{y}\| = \mathbf{P}\mathbf{y},$$

where  $\mathbf{P} = \mathbf{I} - \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}\mathbf{A}$ . Hence,

$$\Pi[\mathbf{x}_0 + \mathbf{y}] = \mathbf{x}_0 + \mathbf{P}\mathbf{y}.$$

**23.5**

Since  $\alpha_k \geq 0$  is a minimizer of  $\phi_k(\alpha) = f(\mathbf{x}^{(k)} - \alpha\mathbf{P}\mathbf{g}^{(k)})$ , we apply the FONC to  $\phi_k(\alpha)$  to obtain

$$\phi'_k(\alpha) = (\mathbf{x}^{(k)} - \alpha\mathbf{P}\mathbf{g}^{(k)})^\top \mathbf{Q}(-\mathbf{P}\mathbf{g}^{(k)}) - \mathbf{b}^\top(-\mathbf{P}\mathbf{g}^{(k)}).$$

Therefore,  $\phi'_k(\alpha) = 0$  if  $\alpha \mathbf{g}^{(k)\top} \mathbf{P}\mathbf{Q}\mathbf{P}\mathbf{g}^{(k)} = (\mathbf{x}^{(k)\top} \mathbf{Q} - \mathbf{b}^\top) \mathbf{P}\mathbf{g}^{(k)}$ . But

$$\mathbf{x}^{(k)\top} \mathbf{Q} - \mathbf{b}^\top = \mathbf{g}^{(k)\top}.$$

Hence

$$\alpha_k = \frac{\mathbf{g}^{(k)\top} \mathbf{P}\mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{P}\mathbf{Q}\mathbf{P}\mathbf{g}^{(k)}}.$$

**23.6**

By Exercise 23.5, the projected steepest descent algorithm applied to this problem takes the form

$$\begin{aligned}\mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \mathbf{P}\mathbf{x}^{(k)} \\ &= (\mathbf{I}_n - \mathbf{P})\mathbf{x}^{(k)} \\ &= \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}\mathbf{A}\mathbf{x}^{(k)}.\end{aligned}$$

If  $\mathbf{x}^{(0)} \in \{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}\}$ , then  $\mathbf{A}\mathbf{x}^{(0)} = \mathbf{b}$ , and hence

$$\mathbf{x}^{(1)} = \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}\mathbf{b}$$

which solves the problem (see Section 12.3).

**23.7**

a. Define

$$\phi_k(\alpha) = f(\mathbf{x}^{(k)} - \alpha\mathbf{P}\nabla f(\mathbf{x}^{(k)}))$$

By the Chain Rule,

$$\begin{aligned}\phi'_k(\alpha_k) &= \frac{d\phi_k}{d\alpha}(\alpha_k) \\ &= -(\nabla f(\mathbf{x}^{(k)} - \alpha_k\mathbf{P}\nabla f(\mathbf{x}^{(k)})))^\top \mathbf{P}\nabla f(\mathbf{x}^{(k)}) \\ &= -(\nabla f(\mathbf{x}^{(k+1)}))^\top \mathbf{P}\nabla f(\mathbf{x}^{(k)}).\end{aligned}$$

Since  $\alpha_k$  minimizes  $\phi_k$ ,  $\phi'_k(\alpha_k) = 0$ , and thus  $\mathbf{g}^{(k+1)\top} \mathbf{P}\mathbf{g}^{(k)} = 0$ .

b. We have  $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} = -\alpha_k \mathbf{P} \mathbf{g}^{(k)}$  and  $\mathbf{x}^{(k+2)} - \mathbf{x}^{(k+1)} = -\alpha_{k+1} \mathbf{P} \mathbf{g}^{(k+1)}$ . Therefore,

$$\begin{aligned} (\mathbf{x}^{(k+2)} - \mathbf{x}^{(k+1)})^\top (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) &= \alpha_{k+1} \alpha_k \mathbf{g}^{(k+1)\top} \mathbf{P}^\top \mathbf{P} \mathbf{g}^{(k)} \\ &= \alpha_{k+1} \alpha_k \mathbf{g}^{(k+1)\top} \mathbf{P} \mathbf{g}^{(k)} \\ &= 0 \end{aligned}$$

by part a, and the fact that  $\mathbf{P} = \mathbf{P}^\top = \mathbf{P}^2$ .

### 23.8

a. minimize  $f(\mathbf{x}) + \gamma P(\mathbf{x})$ .

b. Suppose  $\mathbf{x}^\gamma \notin \Omega$ . Then,  $P(\mathbf{x}^\gamma) > 0$  by definition of  $P$ . Because  $\mathbf{x}^\gamma$  is a global minimizer of the unconstrained problem, we have

$$\begin{aligned} f(\mathbf{x}^\gamma) + \gamma P(\mathbf{x}^\gamma) &\leq f(\mathbf{x}^*) + \gamma P(\mathbf{x}^*) \\ &= f(\mathbf{x}^*), \end{aligned}$$

which implies that

$$f(\mathbf{x}^\gamma) \leq f(\mathbf{x}^*) - \gamma P(\mathbf{x}^\gamma) < f(\mathbf{x}^*).$$

### 23.9

We use the penalty method. First, we construct the unconstrained objective function with penalty parameter  $\gamma$ :

$$f(\mathbf{x}) = x_1^2 + 2x_2^2 + \gamma(x_1 + x_2 - 3)^2.$$

Because  $f$  is a quadratic with positive definite quadratic term, it is easy to find its minimizer:

$$\mathbf{x}_\gamma = \frac{1}{1 + 2/(3\gamma)} \begin{bmatrix} 2 \\ 1 \end{bmatrix}.$$

For example, we can obtain the above by solving the FONC:

$$\begin{aligned} 2(1 + \gamma)x_1 + 2\gamma x_2 - 6\gamma &= 0 \\ 2\gamma x_1 + 2(2 + \gamma)x_2 - 6\gamma &= 0. \end{aligned}$$

Now letting  $\gamma \rightarrow \infty$ , we obtain

$$\mathbf{x}^* = \begin{bmatrix} 2 \\ 1 \end{bmatrix}.$$

(It is easy to verify, using other means, that this is indeed the correct solution.)

### 23.10

Using the penalty method, we construct the unconstrained problem

$$\text{minimize } x + \gamma(\max(a - x, 0))^2$$

To find the solution to the above problem, we use the FONC. It is easy to see that the solution  $x^*$  satisfies  $x^* < a$ . The derivative of the above objective function in the region  $x < a$  is  $1 + 2\gamma(x - a)$ . Thus, by the FONC, we have  $x^* = a - 1/(2\gamma)$ . Since the true solution is at  $a$ , the difference is  $1/(2\gamma)$ . Therefore, for  $1/(2\gamma) \leq \varepsilon$ , we need  $\gamma \geq 1/(2\varepsilon)$ . The smallest such  $\gamma$  is  $1/(2\varepsilon)$ .

### 23.11

a. We have

$$\frac{1}{2} \|\mathbf{x}\|^2 + \gamma \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 1 + 2\gamma & 2\gamma \\ 2\gamma & 1 + 2\gamma \end{bmatrix} \mathbf{x} - \mathbf{x}^\top \begin{bmatrix} 2\gamma \\ 2\gamma \end{bmatrix} + \gamma.$$

The above is a quadratic with positive definite Hessian. Therefore, the minimizer is

$$\begin{aligned}\mathbf{x}_\gamma^* &= \begin{bmatrix} 1+2\gamma & 2\gamma \\ 2\gamma & 1+2\gamma \end{bmatrix}^{-1} \begin{bmatrix} 2\gamma \\ 2\gamma \end{bmatrix} \\ &= \frac{1}{2+1/2\gamma} \begin{bmatrix} 1 \\ 1 \end{bmatrix}.\end{aligned}$$

Hence,

$$\lim_{\gamma \rightarrow \infty} \mathbf{x}_\gamma^* = \frac{1}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

The solution to the original constrained problem is (see Section 12.3)

$$\mathbf{x}^* = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{b} = \frac{1}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

b. We represent the objective function of the associated unconstrained problem as

$$\frac{1}{2} \|\mathbf{x}\|^2 + \gamma \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 = \frac{1}{2} \mathbf{x}^\top (\mathbf{I}_n + 2\gamma \mathbf{A}^\top \mathbf{A}) \mathbf{x} - \mathbf{x}^\top (2\gamma \mathbf{A}^\top \mathbf{b}) + \gamma \mathbf{b}^\top \mathbf{b}.$$

The above is a quadratic with positive definite Hessian. Therefore, the minimizer is

$$\begin{aligned}\mathbf{x}_\gamma^* &= (\mathbf{I}_n + 2\gamma \mathbf{A}^\top \mathbf{A})^{-1} (2\gamma \mathbf{A}^\top \mathbf{b}) \\ &= \left( \frac{1}{2\gamma} \mathbf{I}_n + \mathbf{A}^\top \mathbf{A} \right)^{-1} \mathbf{A}^\top \mathbf{b}.\end{aligned}$$

Let  $\mathbf{A} = \mathbf{U} [\mathbf{S} \quad \mathbf{O}] \mathbf{V}^\top$  be the singular value decomposition of  $\mathbf{A}$ . For simplicity, denote  $\varepsilon = 1/2\gamma$ . We have

$$\begin{aligned}\mathbf{x}_\gamma^* &= (\varepsilon \mathbf{I}_n + \mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} \\ &= \left( \varepsilon \mathbf{I}_n + \mathbf{V} \begin{bmatrix} \mathbf{S} \\ \mathbf{O} \end{bmatrix} \mathbf{U}^\top \mathbf{U} [\mathbf{S} \quad \mathbf{O}] \mathbf{V}^\top \right)^{-1} \mathbf{A}^\top \mathbf{b} \\ &= \left( \varepsilon \mathbf{I}_n + \mathbf{V} \begin{bmatrix} \mathbf{S}^2 & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix} \mathbf{V}^\top \right)^{-1} \mathbf{A}^\top \mathbf{b} \\ &= \mathbf{V} \begin{bmatrix} \varepsilon \mathbf{I}_m + \mathbf{S}^2 & \mathbf{O} \\ \mathbf{O} & \varepsilon \mathbf{I}_{n-m} \end{bmatrix}^{-1} \mathbf{V}^\top \mathbf{A}^\top \mathbf{b}.\end{aligned}$$

Note that

$$\mathbf{V}^\top \mathbf{A}^\top = \begin{bmatrix} \mathbf{S} \\ \mathbf{O} \end{bmatrix} \mathbf{U}^\top.$$

Also,

$$\begin{bmatrix} \varepsilon \mathbf{I}_m + \mathbf{S}^2 & \mathbf{O} \\ \mathbf{O} & \varepsilon \mathbf{I}_{n-m} \end{bmatrix}^{-1} = \begin{bmatrix} (\varepsilon \mathbf{I}_m + \mathbf{S}^2)^{-1} & \mathbf{O} \\ \mathbf{O} & \frac{1}{\varepsilon} \mathbf{I}_{n-m} \end{bmatrix},$$

where  $(\varepsilon \mathbf{I}_m + \mathbf{S}^2)^{-1}$  is diagonal. Hence,

$$\begin{aligned}
\mathbf{x}_\gamma^* &= (\varepsilon \mathbf{I}_n + \mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} \\
&= \mathbf{V} \begin{bmatrix} (\varepsilon \mathbf{I}_m + \mathbf{S}^2)^{-1} & \mathbf{O} \\ \mathbf{O} & \frac{1}{\varepsilon} \mathbf{I}_{n-m} \end{bmatrix} \begin{bmatrix} \mathbf{S} \\ \mathbf{O} \end{bmatrix} \mathbf{U}^\top \\
&= \mathbf{V} \begin{bmatrix} \mathbf{S} \\ \mathbf{O} \end{bmatrix} (\varepsilon \mathbf{I}_m + \mathbf{S}^2)^{-1} \mathbf{U}^\top \\
&= \mathbf{V} \begin{bmatrix} \mathbf{S} \\ \mathbf{O} \end{bmatrix} \mathbf{U}^\top \mathbf{U} (\varepsilon \mathbf{I}_m + \mathbf{S}^2)^{-1} \mathbf{U}^\top \\
&= \mathbf{A}^\top \mathbf{U} (\varepsilon \mathbf{I}_m + \mathbf{S}^2)^{-1} \mathbf{U}^\top.
\end{aligned}$$

Note that as  $\gamma \rightarrow \infty$ ,  $\varepsilon \rightarrow 0$ , and

$$\mathbf{U} (\varepsilon \mathbf{I}_m + \mathbf{S}^2)^{-1} \mathbf{U}^\top \rightarrow \mathbf{U} (\mathbf{S}^2)^{-1} \mathbf{U}^\top.$$

But,

$$\mathbf{U} (\mathbf{S}^2)^{-1} \mathbf{U}^\top = (\mathbf{U} \mathbf{S}^2 \mathbf{U}^\top)^{-1} = (\mathbf{A} \mathbf{A}^\top)^{-1}.$$

Therefore,

$$\mathbf{x}_\gamma^* \rightarrow \mathbf{A}^\top (\mathbf{A} \mathbf{A}^\top)^{-1} \mathbf{b} = \mathbf{x}^*.$$

## 24. Multi-Objective Optimization

### 24.1

---

The MATLAB code is as follows:

```

function multi_op
%MULTI_OP, illustrates multi-objective optimization.

clear
clc
disp ('')
disp ('This is a demo illustrating multi-objective optimization.')
disp ('The numerical example is a modification of the example')
disp ('from the 2002 book by A. Osyczka,')
disp ('Example 5.1 on pages 101--105')
disp ('-----')
disp ('Select the population size denoted POPSIZE, for example, 50.')
disp ('')
POPSIZE=input('Population size POPSIZE = ');
disp ('-----')
disp ('Select the number of iterations denoted NUMITER; e.g., 10.')
disp ('')
NUMITER=input('Number of iterations NUMITER = ');
disp ('')
disp ('-----')

% Main
for i = 1:NUMITER
    fprintf('Working on Iteration %.0f...\n',i)
    xmat = genxmat(POPSIZE);
    if i~=1
        for j = 1:length(xR)

```

```

        xmat = [xmat;xR{j}];
    end
end
[xR,fR] = Select_P(xmat);
fprintf('Number of Pareto solutions: %.0f\n',length(fR))

end
disp ('')
disp ('-----')

fprintf(' Pareto solutions \n')
celldisp(xR)
disp ('')
disp ('-----')

fprintf(' Objective vector values \n')
celldisp(fR)
xlabel('f_1','FontSize',16)
ylabel('f_2','FontSize',16)
title('Pareto optimal front','FontSize',16)
set(gca,'FontSize',16)
grid
for i=1:length(xR)
    xx(i)=xR{i}(1);
    yy(i)=xR{i}(2);
end
XX=[xx; yy];
figure
axis([1 7 5 10])
hold on
for i=1:size(XX,2)
    plot(XX(1,i),XX(2,i),'marker','o','markersize',6)
end
xlabel('x_1','FontSize',16)
ylabel('x_2','FontSize',16)
title('Pareto optimal solutions','FontSize',16)
set(gca,'FontSize',16)
grid
hold off
figure
axis([-2 10 2 13])
hold on
plot([2 6],[5 5],'marker','o','markersize',6)
plot([6 6],[5 9],'marker','o','markersize',6)
plot([2 6],[9 9],'marker','o','markersize',6)
plot([2 2],[5 9],'marker','o','markersize',6)
for i=1:size(XX,2)
    plot(XX(1,i),XX(2,i),'marker','x','markersize',10)
end
x1=-2:.2:10;
x2=2:.2:13;
[X1, X2]=meshgrid(x1,x2);
Z1=-X1.^2 - X2;
v=[0 -5 -7 -10 -15 -20 -30 -40 -60];
cs1=contour(X1,X2,Z1,v);
clabel(cs1)
Z2=X1+X2.^2;
v2=[20 25 35 40 60 80 100 120];
cs2=contour(X1,X2,Z2,v2);

```

```

xlabel(cs2)
xlabel('x_1','FontSize',16)
ylabel('x_2','FontSize',16)
title('Level sets of f_1 and f_2, and Pareto optimal
points','FontSize',16)
set(gca,'FontSize',16)
grid
hold off

function xmat0 = genxmat(POPSIZE)
xmat0 = rand(POPSIZE,2);
xmat0(:,1) = xmat0(:,1)*4+2;
xmat0(:,2) = xmat0(:,2)*4+5;

function [xR,fR] = Select_P(xmat)

% Declaration
J = size(xmat,1);

% Init
Rset = [1];
j = 1;
isstep7 = 0;

% Step 1
x{1} = xmat(1,:);
f{1} = evalfcn(x{1});

% Step 2
while j < J
    j = j+1;

    % Step 3
    r = 1;
    rdel = [];
    q = 0;
    R = length(Rset);

    for k = 1:size(xmat,1)
        x{k} = xmat(k,:);
        f{k} = evalfcn(x{k});
    end

    % Step 4
    while 1
        %for r=1:R
        if all(f{j}<f{Rset(r)})
            q = q+1;
            rdel = [rdel r];
        else

            % Step 5
            if all(f{j}>=f{Rset(r)})
                break
            end
        end
    end

    % Step 6

```

```

        r=r+1;
        if r > R
            isstep7 = 1;
            break
        end
    end
end

% Step 7
if isstep7 == 1
    isstep7 = 0;
    if (q~=0)
        Rset(rdel) = [];
        Rset = [Rset j];
    else

        %Step 8
        Rset = [Rset j];
    end
end

for k = 1:size(xmat,1)
    x{k} = xmat(k,:);
    f{k} = evalfcn(x{k});
end

R = length(Rset);
end

% Return the Pareto solution.
for i = 1:length(Rset)
    xR{i} = x{Rset(i)};
    fR{i} = f{Rset(i)};
end

x1 = [];
y1 = [];
x2 = [];
y2 = [];
for k = 1:size(xmat,1)
    if ismember(k,Rset)
        x1 = [x1 f{k}(1)];
        y1 = [y1 f{k}(2)];
    else
        x2 = [x2 f{k}(1)];
        y2 = [y2 f{k}(2)];
    end
end
end
%newplot
plot(x1,y1,'xr',x2,y2,'.b')
drawnow

function y = f1(x)
% y = x(1)^2+x(2);
% The above function is the original function in the Osyczka's 2002
% book,
% (Example 5.1, page 101).
% Its negative makes a much more interesting example.

```



```

y = -(x(1)^2+x(2));

function y = f2(x)
y = x(1)+x(2)^2;

function y = evalfcn(x)
y(1) = f1(x);
y(2) = f2(x);

```

## 24.2

a. We proceed using contraposition. Assume that  $\mathbf{x}^*$  is not Pareto optimal. Therefore, there exists a point  $\hat{\mathbf{x}} \in \Omega$  such that

$$f_i(\hat{\mathbf{x}}) \leq f_i(\mathbf{x}^*) \quad \text{for all } i = 1, 2, \dots, \ell$$

and for some  $j$ ,  $f_j(\hat{\mathbf{x}}) < f_j(\mathbf{x}^*)$ . Since  $\mathbf{c} > \mathbf{0}$ ,

$$\mathbf{c}^\top \mathbf{f}(\mathbf{x}^*) > \mathbf{c}^\top \mathbf{f}(\hat{\mathbf{x}}),$$

which implies that  $\mathbf{x}^*$  is not a global minimizer for the weighted-sum problem.

For the converse, consider the following counterexample:  $\Omega = \{\mathbf{x} \in \mathbb{R}^2 : \|\mathbf{x}\| \geq 1, \mathbf{x} \geq \mathbf{0}\}$  and  $\mathbf{f}(\mathbf{x}) = [x_1, x_2]^\top$ . It is easy to see that the Pareto front is  $\{\mathbf{x} : \|\mathbf{x}\| = 1, \mathbf{x} \geq \mathbf{0}\}$  (i.e., the part of the unit circle in the nonnegative quadrant). So  $\mathbf{x}^* = (1/\sqrt{2})[1, 1]^\top$  is a Pareto minimizer. However, there is no  $\mathbf{c} > \mathbf{0}$  such that  $\mathbf{x}^*$  is a global minimizer of the weighted-sum problem. To see this, fix  $\mathbf{c} > \mathbf{0}$  (assuming  $c_1 \leq c_2$  without loss of generality) and consider the objective function value  $f(\mathbf{x}^*) = (c_1 + c_2)/\sqrt{2}$  for the weighted-sum problem. Now, the point  $\mathbf{x}_0 = [1, 0]^\top$  is also a feasible point. Moreover  $f(\mathbf{x}_0) = c_1 \leq (c_1 + c_2)/2 \leq f(\mathbf{x}^*)$ . So  $\mathbf{x}^*$  is not a global minimizer of the weighted-sum problem.

b. We proceed using contraposition. Assume that  $\mathbf{x}^*$  is not Pareto optimal. Therefore, there exists a point  $\hat{\mathbf{x}} \in \Omega$  such that

$$f_i(\hat{\mathbf{x}}) \leq f_i(\mathbf{x}^*) \quad \text{for all } i = 1, 2, \dots, \ell$$

and for some  $j$ ,  $f_j(\hat{\mathbf{x}}) < f_j(\mathbf{x}^*)$ . By assumption, for all  $i = 1, \dots, \ell$ ,  $f_i(\mathbf{x}^*) \geq 0$ , which implies that

$$\sum_{i=1}^{\ell} (f_i(\mathbf{x}^*))^p > \sum_{i=1}^{\ell} (f_i(\hat{\mathbf{x}}))^p$$

(because  $p > 0$ ). Hence,  $\mathbf{x}^*$  is not a global minimizer for the minimum-norm problem.

For the converse, consider the following counterexample:  $\Omega = \{\mathbf{x} \in \mathbb{R}^2 : x_1 + 2x_2 \geq 2, \mathbf{x} \geq \mathbf{0}\}$  and  $\mathbf{f}(\mathbf{x}) = [x_1, x_2]^\top$ . It is easy to see that the Pareto front is  $\{\mathbf{x} : x_1 + 2x_2 = 2, \mathbf{x} \geq \mathbf{0}\}$ . So  $\mathbf{x}^* = [1, 1/2]^\top$  is a Pareto minimizer. However, there is no  $p > 0$  such that  $\mathbf{x}^*$  is a global minimizer of the minimum-norm problem. To see this, fix  $p > 0$  and consider the objective function value  $f(\mathbf{x}^*) = 1 + (1/2)^p$  for the minimum-norm problem. Now, the point  $\mathbf{x}_0 = [0, 1]^\top$  is also a feasible point. Moreover  $f(\mathbf{x}_0) = 1 \leq 1 + (1/2)^p = f(\mathbf{x}^*)$ . So  $\mathbf{x}^*$  is not a global minimizer of the minimum-norm problem.

c. For the first part, consider the following counterexample:  $\Omega = \{\mathbf{x} \in \mathbb{R}^2 : x_1 + x_2 \geq 2, \mathbf{x} \geq \mathbf{0}\}$  and  $\mathbf{f}(\mathbf{x}) = [x_1, x_2]^\top$ . The Pareto front is  $\{\mathbf{x} : x_1 + x_2 = 2, \mathbf{x} \geq \mathbf{0}\}$ , and  $\mathbf{x}^* = [1/2, 3/4]^\top$  is a Pareto minimizer. But

$$\begin{aligned}
f(\mathbf{x}^*) &= \max\{f_1(\mathbf{x}^*), f_2(\mathbf{x}^*)\} \\
&= \max\{1/2, 3/4\} \\
&= 3/4.
\end{aligned}$$

However,  $\mathbf{x}_0 = [1, 1]^\top$  is also a feasible point, and  $f(\mathbf{x}_0) = 1 < f(\mathbf{x}^*)$ . Hence,  $\mathbf{x}^*$  is not a global minimizer of the minimax problem.

For the second part, suppose  $\Omega = \{\mathbf{x} \in \mathbb{R}^2 : x_1 \leq 2\}$  and  $\mathbf{f}(\mathbf{x}) = [x_1, 2]^\top$ . Then, for any  $\mathbf{x} \in \Omega$ ,  $\max\{f_1(\mathbf{x}), f_2(\mathbf{x})\} = 2$ . So any  $\mathbf{x}^* \in \mathbb{R}^2$  is a global minimizer of the minimax (single-objective) problem.

However, consider another point  $\hat{\mathbf{x}}$  such that  $\hat{x}_1 < x_1^*$ . Then,  $f_1(\hat{\mathbf{x}}) < f_1(\mathbf{x}^*)$  and  $f_2(\hat{\mathbf{x}}) = f_2(\mathbf{x}^*)$ . Hence,  $\mathbf{x}^*$  is not a Pareto minimizer.

In fact, in the above example, no Pareto minimizer exists. However, if we set  $\Omega = \{\mathbf{x} \in \mathbb{R}^2 : 1 \leq x_1 \leq 2\}$ , then the counterexample is still valid, but in this case any point of the form  $[1, x_2]^\top$  is a Pareto minimizer.

### 24.3

Let

$$f(\mathbf{x}) = \mathbf{c}^{*\top} \mathbf{f}(\mathbf{x})$$

where  $\mathbf{x} \in \Omega = \{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}\}$ . The function  $f$  is convex because all the functions  $f_i$  are convex and  $c_i^* > 0, i = 1, 2, \dots, \ell$ . We can represent the given first-order condition in the following form: for any feasible direction  $\mathbf{d}$  at  $\mathbf{x}^*$ , we have

$$\mathbf{d}^\top \nabla f(\mathbf{x}^*) \geq 0.$$

By Theorem 22.7, the point  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ . Therefore,

$$f(\mathbf{x}^*) \leq f(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \Omega.$$

That is,

$$\sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}^*) \leq \sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \Omega.$$

To finish the proof, we now assume that  $\mathbf{x}^*$  is not Pareto optimal and the above condition holds. We then proceed using the proof by contradiction. Because, by assumption,  $\mathbf{x}^*$  is not Pareto optimal, there exists a point  $\hat{\mathbf{x}} \in \Omega$  such that

$$f_i(\hat{\mathbf{x}}) \leq f_i(\mathbf{x}^*) \quad \text{for all } i = 1, 2, \dots, \ell$$

and for some  $j$ ,  $f_j(\hat{\mathbf{x}}) < f_j(\mathbf{x}^*)$ . Since for all  $i = 1, 2, \dots, \ell$ ,  $c_i^* > 0$ , we must have

$$\sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}^*) > \sum_{i=1}^{\ell} c_i^* f_i(\hat{\mathbf{x}}),$$

which contradicts the above condition,  $\sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}^*) \leq \sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x})$  for all  $\mathbf{x} \in \Omega$ . This completes the proof. (See also Exercise 24.2, part a.)

### 24.4

Let

$$f(\mathbf{x}) = \mathbf{c}^{*\top} \mathbf{f}(\mathbf{x})$$

where  $\mathbf{x} \in \Omega = \{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}\}$ . The function  $f$  is convex because all the functions  $f_i$  are convex and  $c_i^* > 0, i = 1, 2, \dots, \ell$ . We can represent the given Lagrange condition in the form

$$\begin{aligned} Df(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} D\mathbf{g}(\mathbf{x}^*) &= \mathbf{0}^\top \\ \mathbf{h}(\mathbf{x}^*) &= \mathbf{0}. \end{aligned}$$

By Theorem 22.8, the point  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ . Therefore,

$$f(\mathbf{x}^*) \leq f(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \Omega.$$

That is,

$$\sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}^*) \leq \sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \Omega.$$

To finish the proof, we now assume that  $\mathbf{x}^*$  is not Pareto optimal and the above condition holds. We then proceed using the proof by contradiction. Because, by assumption,  $\mathbf{x}^*$  is not Pareto optimal, there exists a point  $\hat{\mathbf{x}} \in \Omega$  such that

$$f_i(\hat{\mathbf{x}}) \leq f_i(\mathbf{x}^*) \quad \text{for all } i = 1, 2, \dots, \ell$$

and for some  $j$ ,  $f_j(\hat{\mathbf{x}}) < f_j(\mathbf{x}^*)$ . Since for all  $i = 1, 2, \dots, \ell$ ,  $c_i^* > 0$ , we must have

$$\sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}^*) > \sum_{i=1}^{\ell} c_i^* f_i(\hat{\mathbf{x}}),$$

which contradicts the above condition,  $\sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}^*) \leq \sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x})$  for all  $\mathbf{x} \in \Omega$ . This completes the proof. (See also Exercise 24.2, part a.)

#### 24.5

Let

$$f(\mathbf{x}) = \mathbf{c}^{*\top} \mathbf{f}(\mathbf{x})$$

where  $\mathbf{x} \in \Omega = \{\mathbf{x} : \mathbf{g}(\mathbf{x}) \leq \mathbf{0}\}$ . The function  $f$  is convex because all the functions  $f_i$  are convex and  $c_i^* > 0$ ,  $i = 1, 2, \dots, \ell$ . We can represent the given KKT condition in the form

$$\begin{aligned} \boldsymbol{\mu}^* &\geq \mathbf{0} \\ Df(\mathbf{x}^*) + \boldsymbol{\mu}^{*\top} D\mathbf{g}(\mathbf{x}^*) &= \mathbf{0}^\top \\ \boldsymbol{\mu}^{*\top} \mathbf{g}(\mathbf{x}^*) &= 0 \\ \mathbf{g}(\mathbf{x}^*) &\leq \mathbf{0}. \end{aligned}$$

By Theorem 22.9, the point  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ . Therefore,

$$f(\mathbf{x}^*) \leq f(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \Omega.$$

That is,

$$\sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}^*) \leq \sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \Omega.$$

To finish the proof, we now assume that  $\mathbf{x}^*$  is not Pareto optimal and the above condition holds. We then proceed using the proof by contradiction. Because, by assumption,  $\mathbf{x}^*$  is not Pareto optimal, there exists a point  $\hat{\mathbf{x}} \in \Omega$  such that

$$f_i(\hat{\mathbf{x}}) \leq f_i(\mathbf{x}^*) \quad \text{for all } i = 1, 2, \dots, \ell$$

and for some  $j$ ,  $f_j(\hat{\mathbf{x}}) < f_j(\mathbf{x}^*)$ . Since for all  $i = 1, 2, \dots, \ell$ ,  $c_i^* > 0$ , we must have

$$\sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}^*) > \sum_{i=1}^{\ell} c_i^* f_i(\hat{\mathbf{x}}),$$

which contradicts the above condition,  $\sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}^*) \leq \sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x})$  for all  $\mathbf{x} \in \Omega$ . This completes the proof. (See also Exercise 24.2, part a.)

#### 24.6

Let

$$f(\mathbf{x}) = \mathbf{c}^{*\top} \mathbf{f}(\mathbf{x})$$

where  $\mathbf{x} \in \Omega = \{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}, \mathbf{g}(\mathbf{x}) \leq \mathbf{0}\}$ . The function  $f$  is convex because all the functions  $f_i$  are convex and  $c_i^* > 0$ ,  $i = 1, 2, \dots, \ell$ . We can represent the given KKT-type condition in the form

$$\begin{aligned} \boldsymbol{\mu}^* &\geq \mathbf{0} \\ Df(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} D\mathbf{h}(\mathbf{x}^*) + \boldsymbol{\mu}^{*\top} D\mathbf{g}(\mathbf{x}^*) &= \mathbf{0}^\top \\ \boldsymbol{\mu}^{*\top} \mathbf{g}(\mathbf{x}^*) &= 0 \\ \mathbf{h}(\mathbf{x}^*) &= \mathbf{0} \\ \mathbf{g}(\mathbf{x}^*) &\leq \mathbf{0}. \end{aligned}$$

By Theorem 22.9, the point  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ . Therefore,

$$f(\mathbf{x}^*) \leq f(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \Omega.$$

That is,

$$\sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}^*) \leq \sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \Omega.$$

To finish the proof, we now assume that  $\mathbf{x}^*$  is not Pareto optimal and the above condition holds. We then proceed using the proof by contradiction. Because, by assumption,  $\mathbf{x}^*$  is not Pareto optimal, there exists a point  $\hat{\mathbf{x}} \in \Omega$  such that

$$f_i(\hat{\mathbf{x}}) \leq f_i(\mathbf{x}^*) \quad \text{for all } i = 1, 2, \dots, \ell$$

and for some  $j$ ,  $f_j(\hat{\mathbf{x}}) < f_j(\mathbf{x}^*)$ . Since for all  $i = 1, 2, \dots, \ell$ ,  $c_i^* > 0$ , we must have

$$\sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}^*) > \sum_{i=1}^{\ell} c_i^* f_i(\hat{\mathbf{x}}),$$

which contradicts the above condition,  $\sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x}^*) \leq \sum_{i=1}^{\ell} c_i^* f_i(\mathbf{x})$  for all  $\mathbf{x} \in \Omega$ . This completes the proof. (See also Exercise 24.2, part a.)

#### 24.7

The given minimax problem is equivalent to the problem given in the hint:

$$\begin{array}{ll} \text{minimize} & z \\ \text{subject to} & f_i(\mathbf{x}) - z \leq 0, \quad i = 1, 2. \end{array}$$

Suppose  $[\mathbf{x}^{*\top}, z^*]^\top$  is a local minimizer for the above problem (which is equivalent to  $\mathbf{x}^*$  being a local minimizer for the original problem). Then, by the KKT Theorem, there exists  $\boldsymbol{\mu}^{*\top} \geq \mathbf{0}$ , where  $\boldsymbol{\mu}^* \in \mathbb{R}^2$ , such that

$$\begin{aligned} [\mathbf{0}^\top, 1] + \boldsymbol{\mu}^{*\top} \begin{bmatrix} [\nabla f_1(\mathbf{x}^*)^\top, -1] \\ [\nabla f_2(\mathbf{x}^*)^\top, -1] \end{bmatrix} &= \mathbf{0}^\top \\ \boldsymbol{\mu}^{*\top} \begin{bmatrix} f_1(\mathbf{x}^*) - z^* \\ f_2(\mathbf{x}^*) - z^* \end{bmatrix} &= 0. \end{aligned}$$

Rewriting the first equation above, we get

$$\mu_1^* \nabla f_1(\mathbf{x}^*) + \mu_2^* \nabla f_2(\mathbf{x}^*) = \mathbf{0}, \quad \mu_1^* + \mu_2^* = 1.$$

Rewriting the second equation, we get

$$\mu_i^* (f_i(\mathbf{x}^*) - z^*) = 0, \quad i = 1, 2.$$

Suppose  $f_i(\mathbf{x}^*) < \max\{f_1(\mathbf{x}^*), f_2(\mathbf{x}^*)\}$ , where  $i \in \{1, 2\}$ . Then,  $z^* > f_i(\mathbf{x}^*)$ . Hence, by the above equation we conclude that  $\mu_i^* = 0$ .