# Optimization of Quadratic Functions

In this chapter we study the problem

(42)
$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \ \tfrac{1}{2}x^T H x + g^T x + \beta,$$

where $H \in \mathbb{R}^{n \times n}$ is symmetric, $g \in \mathbb{R}^n$, and $\beta \in \mathbb{R}$. It has already been observed that we may as well assume that $H$ is symmetric since

$$x^T H x = \tfrac{1}{2}x^T H x + \tfrac{1}{2}(x^T H x)^T = x^T \left[ \tfrac{1}{2}(H + H^T) \right] x,$$

where $\tfrac{1}{2}(H + H^T)$ is called the *symmetric part* of $H$. Therefore, in this chapter we assume that $H$ is symmetric. In addition, we have also noted that an objective function can always be shifted by a constant value without changing the solution set to the optimization problem. Therefore, we assume that $\beta = 0$ for most of our discussion. However, just as in the case of integration theory where it is often helpful to choose a particular constant of integration, in many applications there is a "natural" choice for $\beta$ that helps one interpret the problem as well as its solution.

The class of problems (42) is important for many reasons. Perhaps the most common instance of this problem is the linear least squares problem:

(43)
$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \ \tfrac{1}{2} \left\| Ax - b \right\|_2^2,$$

where $A \in \mathbb{R}^{m \times n}$, and $b \in \mathbb{R}^m$. By expanding the objective function in (43), we see that

(44)
$$\tfrac{1}{2} \left\| Ax - b \right\|_2^2 = \tfrac{1}{2}x^T(A^T A)x - (A^T b)^T x + \tfrac{1}{2} \left\| b \right\|_2^2 = \tfrac{1}{2}x^T H x + g^T x + \beta,$$

where $H = A^T A$, $g = -A^T b$, and $\beta = \tfrac{1}{2} \left\| b \right\|_2^2$. This connection to the linear least squares problem will be explored in detail later in this chapter. For the moment, we continue to exam the general problem (42). As in the case of the linear least squares problem, we begin by discussing characterizations of the solutions as well as their existence and uniqueness. In this discussion we try to follow the approach taken for the the linear least squares problem. However, in the case of (43), the matrix $H := A^T A$ and the vector $g = -A^T b$ possess special features that allowed us to establish very strong results on optimality conditions as well as on the existence and uniqueness of solutions. In the case of a general symmetric matrix $H$ and vector $g$ it is possible to obtain similar results, but there are some twists. Symmetric matrices have many special properties that can be exploited to help us achieve our goal. Therefore, we begin by recalling a few of these properties, specifically those related to eigenvalue decomposition.

## 1. Eigenvalue Decomposition of Symmetric Matrices

Given a matrix $A \in \mathbb{R}^{n \times n}$, we say that the scalar $\lambda$ is an eigenvalue of $A$ if there is a non-zero vector $x$ such that $Ax = \lambda x$, or equivalently, $\text{Null}(\lambda I - A) \neq \{0\}$. Observe that $\text{Null}(\lambda I - A) \neq \{0\}$ if and only if $(\lambda I - A)$ is singular, that is, $\det(\lambda I - A) = 0$. Consequently, $\lambda$ is an eigenvalue of $A$ if and only if $\det(\lambda I - A) = 0$. If we now think of $\lambda$ as a variable, this says that we can find all eigenvalues of $A$ by finding all roots of the equation $\det(\lambda I - A) = 0$. The function $p(\lambda) := \det(\lambda I - A)$ is easily seen to be a polynomial of degree $n$ in $\lambda$ which we call the *characteristic polynomial* of $A$. By the Fundamental Theorem of Algebra, we know that $p(\lambda)$ has $n$ roots over the complex numbers if we count the multiplicities of these roots. Hence, when we discuss eigenvalues and eigenvectors we are forced in the setting of complex numbers. For this reason we may as well assume that $A \in \mathbb{C}^{n \times n}$.

Working on $\mathbb{C}^n$ requires us to re-examine our notion of the Euclidean norm and its associated dot product. Recall that for a complex number $\zeta := x + iy$, with $x, y \in \mathbb{R}$ and $i := \sqrt{-1}$, the magnitude of $\zeta$ is given by $|\zeta| = \sqrt{\bar{\zeta}\zeta}$, where $\bar{\zeta} := x - iy$ is the complex conjugate of $\zeta$. If we now define the Euclidean norm of a vector $z \in \mathbb{C}^n$ to be the square root of the sum of the squares of magnitude of its components, then

$$\left\| z \right\|_2 = \sqrt{\sum_{k-1}^{n} |z_k|^2} = \sqrt{\sum_{k-1}^{n} \bar{z}_k z_k} = \sqrt{\bar{z}^T z} = \sqrt{z^* z},$$

where we define
$$z^* z = (\overline{z})^T z,$$
that is, $z^*$ takes $z$ to its *conjugate transpose*. When $z \in \mathbb{R}^n$, we have $z^* = z^T$, and we recover the usual formulas. With the $*$ operation, we can extend our notion of dot product (or, inner product) by writing
$$\langle z, y \rangle := z^* y \in \mathbb{C}.$$
When $z$ and $y$ are real vectors we recover usual notion of dot product for such vectors. Finally, for matrices $A \in \mathbb{C}^{n \times n}$, we define
$$A^* := \overline{A}^T,$$
that is, we conjugate every element of $A$ and then take the transpose. This notation is very helpful in a number of ways. For example, we have
$$\langle Ay, x \rangle = (Ay)^* x = y^* A^* x \quad \text{and} \quad \|Ax\|_2^2 = x^* A^* A x .$$
We call $A^*$ the *adjoint* of $A$.

Recall that a matrix $H \in \mathbb{R}^{n \times n}$ is said to be symmetric of $H^T = H$. By extension, we say that an matrix $Q \in \mathbb{C}^{n \times n}$ is *self-adjoint* if $Q^* = Q$. Thus, in particular, every real symmetric matrix is self adjoint. We have the following remarkable fact about self-adjoint matrices.

LEMMA 1.1. *If $Q \in \mathbb{C}^{n \times n}$ is self-adjoint, then $Q$ has only real eigenvalues. In particular, if $H$ is a real symmetric matrix, then $H$ has only real eigenvalues and for each such eigenvalue there is a real eigenvector. Moreover, if $(\lambda_1, v^1)$ and $(\lambda_2, v^2)$ are two eigenvalue-eigenvectors pairs for $H$ with $\lambda_1 \neq \lambda_2$, then $(v^1)^T v^2 = 0$.*

PROOF. Let $\lambda \in \mathbb{C}$ be an eigenvalue of $Q$. Then there is a non-zero eigenvector $x \in \mathbb{C}^n$ such that $Qx = \lambda x$. Therefore,
$$\lambda \|x\|_2^2 = \lambda x^* x = x^* Q x = x^* Q^* x = (x^* Q x)^* = (\lambda \|x\|_2^2)^* = \overline{\lambda} \|x\|_2^2,$$
so that $\lambda = \overline{\lambda}$ which can only occur if $\lambda$ is a real number.

If $H$ is real symmetrix, then it is self adjoint so all of its eigenvalues are real. If $\lambda$ is one such eigenvalue with associated eigenvector $z = x + \mathrm{i}y$ with $x, y \in \mathbb{R}^n$, then
$$Hx + \mathrm{i}Hy = Hz = \lambda z = \lambda x + \mathrm{i}\lambda y.$$
Consequently, $Hx = \lambda x$ and $Hy = \lambda y$ since both $Hx$ and $Hy$ are real vectors. Since $z \neq 0$, either $x$ or $y$ or both are non-zero, in any case we have a real eigenvector for $H$ corresponding to $\lambda$.

Next let $(\lambda_1, v^1)$ and $(\lambda_2, v^2)$ be eigenvalue-eigenvectors pairs for $H$ with $\lambda_1 \neq \lambda_2$. Then
$$\lambda_1 (v^1)^T v^2 = (Hv^1)^T v^2 = (v^1)^T H v^2 = \lambda_2 (v^1)^T v^2,$$
since $\lambda_1 \neq \lambda_2$, we must have $(v^1)^T v^2 = 0$. $\qquad\qquad\square$

Next, suppose $\lambda_1$ is an eigenvalue for the real symmetric matrix $H \in \mathbb{R}^{n \times n}$ and let the columns of the matrix $U_1 \in \mathbb{R}^{n \times k}$ form an orthonormal basis for the subspace $\mathrm{Null}(\lambda_1 I - H)$, where $k = \dim(\mathrm{Null}(\lambda_1 I - H)) \geq 1$. Let the columns of $U_2 \in \mathbb{R}^{n \times (n-k)}$ form an orthonormal basis for the subspace $\mathrm{Null}(\lambda_1 I - H)^\perp$ and set $\widetilde{U} = [U_1 \; U_2] \in \mathbb{R}^{n \times n}$. Then $\widetilde{U}^T \widetilde{U} = I$, that is, $\widetilde{U}$ is a *unitary matrix*. In particular, $\widetilde{U}^{-1} = \widetilde{U}^T$ and so $\widetilde{U}\widetilde{U}^T = I$ as well. We have the following relationships between $U_1$, $U_2$, and $H$:
$$HU_1 = \lambda_1 U_1, \quad U_1^T H U_1 = \lambda_1 U_1^T U_1 = \lambda_1 I_k \quad \text{and} \quad (U_1^T H U_2)^T = U_2^T H U_1 = \lambda_1 U_2^T U_1 = 0_{(n-k) \times k}.$$
Consequently,
$$(45) \qquad \widetilde{U}^T H \widetilde{U} = \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix} H \begin{bmatrix} U_1 & U_2 \end{bmatrix} = \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix} \begin{bmatrix} HU_1 & HU_2 \end{bmatrix} = \begin{bmatrix} U_1^T H U_1 & U_1^T H U_2 \\ U_2^T H U_1 & U_2^T H U_2 \end{bmatrix} = \begin{bmatrix} \lambda_1 I_k & 0 \\ 0 & U_2^T H U_2 \end{bmatrix},$$
and so
$$(46) \qquad H = \widetilde{U}\widetilde{U}^T H \widetilde{U}\widetilde{U}^T = \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{bmatrix} \lambda_1 I_k & 0 \\ 0 & U_2^T H U_2 \end{bmatrix} \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix}.$$
These observations provide the foundation for the following eigenvalue theorem for real symmetric matrices.

THEOREM 1.1. *[Eigenvalue Decomposition for Symmetric Matrices] Let $H \in \mathbb{R}^{n \times n}$ be a real symmetric matrix. Then there is a unitary matrix $U$ such that*
$$H = U \Lambda U^T,$$
*where $\Lambda := \mathrm{diag}(\lambda_1, \lambda_2, \ldots, \lambda_n)$ with $\lambda_1, \lambda_2, \ldots, \lambda_n$ being the eigenvalues of $H$ repeated according to multiplicity.*

PROOF. We proceed by induction on the dimension. The result is trivially true for $n = 1$. Assume that the result is true for all dimensions $k < n$ with $n > 1$ and show it is true for all $n \times n$ symmetric matrices. Let $H \in \mathbb{R}^{n \times n}$ be symmetric and let $\lambda_1$ be any eigenvalue of $H$ with $k = \dim(\text{Null}(\lambda_1 I - H)) \geq 1$. Let $U_1 \in \mathbb{R}^{n \times k}$ and $U_2 \in \mathbb{R}^{n \times (n-k)}$ be as in (45) and (46) above. If $k = n$, the result follows from (45) so we can assume that $k < n$.

Since (45) is a similarity transformation of $H$, $\widetilde{U}^T H \widetilde{U}$ has the same characteristic polynomial as $H$:

$$\det(\lambda I_n - H) = (\lambda - \lambda_1)^k q(\lambda), \quad \text{where} \quad q(\lambda) = \det(\lambda I_{n-k} - U_2^T H U_2).$$

Therefore, the eigenvalues of $U_2^T H U_2$ are necessarily those of $H$ that are not equal to $\lambda_1$ and each has the same multiplicity as they have for $H$.

Apply the induction hypothesis to the $(n - k) \times (n - k)$ matrix $U_2^T H U_2$ to obtain a real unitary matrix $V \in \mathbb{R}^{(n-k) \times (n-k)}$ such that

$$U_2^T H U_2 = V \widetilde{\Lambda} V^T,$$

where $\widetilde{\Lambda} = \text{diag}(\mu_1, \mu_2, \ldots, \mu_{(n-k)})$ with $\mu_1, \mu_2, \ldots, \mu_{(n-k)}$ being the eigenvalues of $H$ that are not equal to $\lambda_1$ with each having the same multiplicity as they have for $H$. Then, by (46)

$$H = \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{bmatrix} \lambda_1 I_k & 0 \\ 0 & U_2^T H U_2 \end{bmatrix} \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix} = \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{bmatrix} I_k & 0 \\ 0 & V \end{bmatrix} \begin{bmatrix} \lambda_1 I_k & 0 \\ 0 & \widetilde{\Lambda} \end{bmatrix} \begin{bmatrix} I_k & 0 \\ 0 & V^T \end{bmatrix} \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix}.$$

The result is obtained by setting

$$U = \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{bmatrix} I_k & 0 \\ 0 & V \end{bmatrix} = \begin{bmatrix} U_1 & U_2 V \end{bmatrix}$$

and observing that $U^T U = I$. $\qquad \square$

One important consequence of this result is the following theorem

THEOREM 1.2. *[The Rayleigh-Ritz Theorem] Let the symmetric matrix $H \in \mathbb{R}^{n \times n}$ have smallest eigenvalue $\lambda_{min}(H)$ and largest eigenvalue $\lambda_{max}(H)$. Then, for all $u \in \mathbb{R}^n$,*

$$\lambda_{min}(H) \|u\|_2^2 \leq u^T H u \leq \lambda_{max}(H) \|u\|_2^2,$$

*with equality holding on the left for every eigenvector $u$ for $\lambda_{min}(H)$ and equality holding on the right for every eigenvector $u$ for $\lambda_{max}(H)$.*

PROOF. Let $H = U \Lambda U^T$ be the eigenvalue decomposition of $H$ in Theorem 1.1 with $\Lambda = \text{diag}(\lambda_1, \lambda_2, \ldots, \lambda_n)$. Then the columns of $U$ form an orthonormal basis for $\mathbb{R}^n$. Therefore, given any $u \in \mathbb{R}^n \setminus \{0\}$, there is a $z \in \mathbb{R}^n \setminus \{0\}$ such that $u = Uz$. Hence

$$u^T H u = (Uz)^T U \Lambda U^T (Uz) = z^T \Lambda z = \sum_{j=1}^n \lambda_j z_j^2.$$

Clearly,

$$\lambda_{\min}(H) \|z\|_2^2 = \sum_{j=1}^n \lambda_{\min}(H) z_j^2 \leq \sum_{j=1}^n \lambda_j z_j^2 \leq \sum_{j=1}^n \lambda_{\max}(H) z_j^2 = \lambda_{\max}(H) \|z\|_2^2.$$

The result now follows since $\|z\|_2^2 = z^T z = z^T U^T U z = u^T u = \|u\|_2^2.$ $\qquad \square$

The following definition describes some important concepts associated with symmetric matrices that are important for optimization.

DEFINITION 1.1. *Let $H \in \mathbb{R}^{n \times n}$.*
(1) *$H$ is said to be positive definite if $x^T H x > 0$ for all $x \in \mathbb{R}^n \setminus \{0\}$.*
(2) *$H$ is said to be positive semi-definite if $x^T H x \geq 0$ for all $x \in \mathbb{R}^n$.*
(3) *$H$ is said to be negative definite if $x^T H x < 0$ for all $x \in \mathbb{R}^n \setminus \{0\}$.*
(4) *$H$ is said to be positive semi-definite if $x^T H x \leq 0$ for all $x \in \mathbb{R}^n$.*
(5) *$H$ is said to be indefinite if $H$ is none of the above.*

We denote the set of real $n \times n$ symmetric matrices by $\mathcal{S}^n$, the set of positive semi-definite real $n \times n$ symmetric matrices by $\mathcal{S}_+^n$, and the set of positive definite real $n \times n$ symmetric matrices by $\mathcal{S}_{++}^n$. It is easily seen that $\mathcal{S}^n$ is a vector space.

Theorem 1.2 provides necessary and sufficient conditions under which a symmetric matrix $H$ is positive/negative definite/semi-definite. For example, since $\lambda_{\min}(H) \|u\|_2^2 \leq u^T H u$ with equality when $u$ is an eigenvector associated

with $\lambda_{\min}(H)$, we have that $H$ is positive definite if and only if $\lambda_{\min}(H) > 0$. Similar results can be obtained for the other cases.

An additional property of positive semi-definite matrices is that they possess *square roots*. If $H \in \mathbb{R}^{n \times n}$ is symmetric and positive semi-definite, then Theorem 1.1 tells us that $H = U\Lambda U^T$, where $U$ is unitary and $\Lambda = \operatorname{diag}(\lambda_1, \ldots, \lambda_n)$ with $\lambda_i \geq 0$, $i = 1, \ldots, n$. If we define $\Lambda^{1/2} := \operatorname{diag}(\sqrt{\lambda_1}, \ldots, \sqrt{\lambda_n})$ and $H^{1/2} = U\Lambda^{1/2}U^T$, then $H = U\Lambda U^T = U\Lambda^{1/2}U^T = H^{1/2}H^{1/2}$, so $H^{1/2}$ provides a natural notion of the square root of a matrix. However, $H^{1/2}$ is not uniquely defined since we can always re-order the diagonal elements and their corresponding columns to produce the same effect. In addition, $H^{1/2}$ is always symmetric while in some instances choosing a non-symmetric square root may be beneficial. For example, if we consider the linear least squares problem (43), then $H = A^T A$. Should $A$ be considered a square root of $H$? In order to cover the full range of possible considerations, we make the following definition for the square root of a symmetric matrix.

DEFINITION 1.2. *[Square Roots of Positive Semi-Definite Matrices] Let $H \in \mathbb{R}^{n \times n}$ be a symmetric positive semi-definite matrix. We say that the matrix $L \in \mathbb{R}^{n \times n}$ is a square root of $H$ if $H = LL^T$.*

## 2. Optimality Properties of Quadratic Functions

Recall that for the linear least squares problem, we were able to establish a necessary and sufficient condition for optimality, namely the normal equations, by working backward from a known solution. We now try to apply this same approach to quadratic functions, in particular, we try to extend the derivation in (31) to the objective function in (47). Suppose $\overline{x}$ is a local solution to the quadratic optimization problem

$$(47) \qquad \underset{x \in \mathbb{R}^n}{\operatorname{minimize}} \tfrac{1}{2} x^T H x + g^T x,$$

where $H \in \mathbb{R}^{n \times n}$ is symmetric and $g \in \mathbb{R}^n$, i.e., there is an $\epsilon > 0$ such that

$$(48) \qquad \tfrac{1}{2}\overline{x}^T H \overline{x} + g^T \overline{x} \leq \tfrac{1}{2} x^T H x + g^T x \quad \forall\, x \in \overline{x} + \epsilon \mathbb{B}_2,$$

where $\overline{x} + \epsilon \mathbb{B}_2 := \{\overline{x} + \epsilon u \,|\, u \in \mathbb{B}_2\}$ and $\mathbb{B}_2 := u\|u\|_2 \leq 1$ (hence, $\overline{x} + \epsilon \mathbb{B}_2 = \{x \,|\, \|\overline{x} - x\|_2 \leq \epsilon\}$). Note that, for all $x \in \mathbb{R}^n$,

$$
\begin{aligned}
\overline{x}^T H \overline{x} &= (x + (\overline{x} - x))^T H(x + (\overline{x} - x)) \\
&= x^T H x + 2x^T H(\overline{x} - x) + (\overline{x} - x)^T H(\overline{x} - x) \\
(49) \qquad &= x^T H x + 2(\overline{x} + (x - \overline{x}))^T H(\overline{x} - x) + (\overline{x} - x)^T H(\overline{x} - x) \\
&= x^T H x + 2\overline{x}^T H(\overline{x} - x) + 2(x - \overline{x})^T H(\overline{x} - x) + (\overline{x} - x)^T H(\overline{x} - x) \\
&= x^T H x + 2\overline{x}^T H(\overline{x} - x) - (\overline{x} - x)^T H(\overline{x} - x).
\end{aligned}
$$

Therefore, for all $x \in \overline{x} + \epsilon \mathbb{B}_2$,

$$
\begin{aligned}
\tfrac{1}{2}\overline{x}^T H \overline{x} + g^T \overline{x} &= (\tfrac{1}{2} x^T H x + g^T x) + (H\overline{x} + g)^T(\overline{x} - x) - \tfrac{1}{2}(\overline{x} - x)^T H(\overline{x} - x) \\
&\geq (\tfrac{1}{2}\overline{x}^T H \overline{x} + g^T \overline{x}) + (H\overline{x} + g)^T(\overline{x} - x) - \tfrac{1}{2}(\overline{x} - x)^T H(\overline{x} - x), \qquad \text{(since } \overline{x} \text{ is a local solution)}
\end{aligned}
$$

and so

$$(50) \qquad \tfrac{1}{2}(\overline{x} - x)^T H(\overline{x} - x) \geq (H\overline{x} + g)^T(\overline{x} - x) \qquad \forall\, x \in \overline{x} + \epsilon \mathbb{B}_2.$$

Let $0 \leq t \leq \epsilon$ and $v \in \mathbb{B}_2$ and define $x = \overline{x} + tv \in \overline{x} + \epsilon \mathbb{B}_2$. If we plug $x = \overline{x} + tv$ into (50), then

$$(51) \qquad \frac{t^2}{2} v^T H v \geq -t(H\overline{x} + g)^T v.$$

Dividing this expression by $t > 0$ and taking the limit as $t \downarrow 0$ tells us that

$$0 \leq (H\overline{x} + g)^T v \quad \forall\, v \in \mathbb{B}_2,$$

which implies that $H\overline{x} + g = 0$. Plugging this information back into (51) gives

$$\frac{t^2}{2} v^T H v \geq 0 \quad \forall\, v \in \mathbb{B}_2.$$

Dividing by $t^2/2$ for $t \neq 0$ tells us that

$$v^T H v \geq 0 \quad \forall\, v \in \mathbb{B}_2$$

or equivalently, that $H$ is positive semi-definite. These observations motivate the following theorem.

THEOREM 2.1. *[Existence and Uniqueness in Quadratic Optimization] Let $H \in \mathbb{R}^{n \times n}$ and $g \in \mathbb{R}^n$ be as in (47).*

(1) *A local solution to the problem (47) exists if and only if $H$ is positive semi-defnite and there exists a solution $\overline{x}$ to the equation $Hx + g = 0$ in which case $\overline{x}$ is a local solution to (47).*

(2) *If $\overline{x}$ is a local solution to (47), then it is a global solution to (47).*

(3) *The problem (47) has a unique global solution if and only if $H$ is positive definite in which case this solution is given by $\overline{x} = -H^{-1}g$.*

(4) *If either $H$ is not positive semi-definite or there is no solution to the equation $Hx + g = 0$ (or both), then*

$$-\infty = \inf_{x \in \mathbb{R}^n} \tfrac{1}{2}x^T H x + g^T x \ .$$

PROOF. (1) We have already shown that if a local solution $\overline{x}$ to (47) exists, then $H\overline{x} + g = 0$ and $H$ is positive semi-definite. On the other hand, suppose that $H$ is positive semi-definite and $\overline{x}$ is a solution to $Hx + g = 0$. Then, for all $x \in \mathbb{R}^n$, we can interchange the roles of $x$ and $\overline{x}$ in the second line of (49) to obtain

$$x^T H x = \overline{x}^T H \overline{x} + 2\overline{x}^T H(x - \overline{x}) + (x - \overline{x})^T H(x - \overline{x}).$$

Hence, for all $x \in \mathbb{R}^n$,

$$\tfrac{1}{2}x^T H x + g^T x = \tfrac{1}{2}\overline{x}^T H \overline{x} + g^T \overline{x} + (H\overline{x} + g)^T(x - \overline{x}) + \tfrac{1}{2}(x - \overline{x})^T H(x - \overline{x}) \geq \tfrac{1}{2}\overline{x}^T H \overline{x} + g^T \overline{x} \ ,$$

since $H\overline{x} + g = 0$ and $H$ is positive semi-definite. That is, $\overline{x}$ is a global solution to (47) and hence a local solution.

(2) Suppose $\overline{x}$ is a local solution to (47) so that, by Part (1), $H$ is positive semi-definite and $H\overline{x} + g = 0$, and there is an $\epsilon > 0$ such that (48) holds. Next observe that, for all $x, y \in \mathbb{R}^n$ and $\lambda \in \mathbb{R}$, we have

$$((1 - \lambda)x + \lambda y)^T H((1 - \lambda)x + \lambda y) - (1 - \lambda)x^T H x - \lambda y^T H y$$
$$= (1 - \lambda)^2 x^T H x + 2\lambda(1 - \lambda)x^T H y + \lambda^2 y^2 H y - (1 - \lambda)x^T H x - \lambda y^T H y$$
$$= -\lambda(1 - \lambda)x^T H x + 2\lambda(1 - \lambda)x^T H y - \lambda(1 - \lambda)y^T H y$$
$$= -\lambda(1 - \lambda)(x - y)^T H(x - y),$$

or equivalently,

$$(52) \qquad ((1 - \lambda)x + \lambda y)^T H((1 - \lambda)x + \lambda y) \leq (1 - \lambda)x^T H x + \lambda y^T H y - \lambda(1 - \lambda)(x - y)^T H(x - y).$$

Since $H$ is positive semi-definite, this implies that

$$(53) \qquad ((1 - \lambda)x + \lambda y)^T H((1 - \lambda)x + \lambda y) \leq (1 - \lambda)x^T H x + \lambda y^T H y \quad \forall \lambda \in [0, 1].$$

If $\overline{x}$ is not a global solution, then there is an $\hat{x}$ such that $f(\hat{x}) < f(\overline{x})$, where $f(x) := \tfrac{1}{2}x^T H x + g^T x$. By (48), we must have $\|\overline{x} - \hat{x}\|_2 > \epsilon$. Set $\lambda := \frac{\epsilon}{2\|\overline{x} - \hat{x}\|_2}$ so that $0 < \lambda < 1$, and define $x_\lambda := (1 - \lambda)\overline{x} + \lambda\hat{x} = \overline{x} + \lambda(\hat{x} - \overline{x})$ so that $x_\lambda \in \overline{x} + \epsilon\mathbb{B}_2$. But then, by (53),

$$f(x_\lambda) \leq (1 - \lambda)f(\overline{x}) + \lambda f(\hat{x}) < (1 - \lambda)f(\overline{x}) + \lambda f(\overline{x}) = f(\overline{x}),$$

which contradicts (48). Hence, no such $\hat{x}$ can exist so that $\overline{x}$ is a global solution to (47).

(3) If (47) has a unique global solution $\overline{x}$, then $\overline{x}$ must be the unique solution to the equation $Hx + g = 0$. This can only happen if $H$ is invertible. Hence, $H$ is invertible and positive semi-definite which implies that $H$ is positive definite. On the other hand, if $H$ is positive definite, then it is positive semi-definite and there is a unique solution to the equation $Hx + g = 0$, i.e., (49) has a unique global solution.

(4) The result follows if we can show that $f(x) := \tfrac{1}{2}x^T H x + g^T x$ is unbounded below when either $H$ is not positive semi-definite of there is no solution to the equation $Hx + g = 0$ (or both). Let us first suppose that $H$ is not positive semi-definite, or equivalently, $\lambda_{\min}(H) < 0$. Let $u \in \mathbb{R}^n$ be an eigenvector associated with the eigenvalue $\lambda_{\min}(H)$ with $\|u\|_2 = 1$. Then, for $x := tu$ with $t > 0$, we have $f(tu) = \lambda_{\min}(H)\frac{t^2}{2} + tg^T u \xrightarrow{t\uparrow\infty} -\infty$ since $\lambda_{\min}(H) < 0$, so $f$ is unbounded below.

Next suppose that there is no solution to the equation $Hx + g = 0$. In particular, $g \neq 0$ and $g \notin \text{Ran}(H) = \text{Null}(H)^\perp$. Then the orthogonal projection of $g$ onto $\text{Null}(H)$ cannot be zero: $\hat{g} := P_{\text{Null}(H)}(g) \neq 0$. Hence, for $x := -t\hat{g}$ with $t > 0$, we have $f(-t\hat{g}) = -t\|\hat{g}\|_2^2 \xrightarrow{t\uparrow\infty} -\infty$, so again $f$ is unbounded below. $\square$

The identity (52) is a very powerful tool in the analysis of quadratic functions. It was the key tool in showing that every local solution to (47) is necessarily a global solution. It is also remarkable, that a local solution exists if and only if $H$ is positive definite and there is a solution to the equation $Hx + g = 0$. We now show how these results can be extended to problems with linear equality constraints.

## 3. Minimization of a Quadratic Function on an Affine Set

In this section we consider the problem

(54)
$$\begin{aligned} \text{minimize} \ & \tfrac{1}{2}x^T H x + g^T x \\ \text{subject to} \ & Ax = b, \end{aligned}$$

where $H \in \mathbb{R}^{n \times n}$ is symmetric, $g \in \mathbb{R}^n$, $A \in \mathbb{R}^{m \times n}$, and $b \in \mathbb{R}^m$. We assume that the system $Ax = b$ is consistent. That is, there exists $\hat{x} \in \mathbb{R}^n$ such that $A\hat{x} = b$ in which case

$$\{x \,|\, Ax = b\} = \hat{x} + \text{Null}(A).$$

Consequently, the problem (54) is of the form

(55)
$$\underset{x \in \hat{x} + S}{\text{minimize}} \ \tfrac{1}{2}x^T H x + g^T x \ ,$$

where $S$ is a subspace of $\mathbb{R}^n$. This representation of the problem shows that the problem (54) is trivial if $\text{Null}(A) = \{0\}$ since then the unique solution $\hat{x}$ to $Ax = b$ is the unique solution to (54). Hence, when considering the problem (54) it is always assumed that $\text{Null}(A) \neq \{0\}$, and furthermore, that $m < n$.

DEFINITION 3.1. *[Affine Sets] A subset $K$ of $\mathbb{R}^n$ is said to be an affine set if there exists a point $\hat{x} \in \mathbb{R}^n$ and a subspace $S \subset \mathbb{R}^n$ such that $K = \hat{x} + S = \{\hat{x} + u \,|\, u \in S\}$.*

We now develop necessary and sufficient optimality conditions for the problem (55), that is, for the minimization of a quadratic function over an affine set. For this we assume that we have a basis $v^1, v^2, \ldots, v^k$ for $S$ so that $\dim(S) = k$. Let $V \in \mathbb{R}^{n \times k}$ be the matrix whose columns are the vectors $v^1, v^2, \ldots, v^k$ so that $S = \text{Ran}(V)$. Then $\hat{x} + S = \{\hat{x} + Vz \,|\, z \in \mathbb{R}^k\}$. This allows us to rewrite the problem (55) as

(56)
$$\underset{z \in \mathbb{R}^k}{\text{minimize}} \ \tfrac{1}{2}(\hat{x} + Vz)^T H (\hat{x} + Vz) + g^T(\hat{x} + Vz) \ .$$

PROPOSITION 3.1. *Consider the two problems (55) and (56), where the columns of the matrix $V$ form a basis for the subspace $S$. The set of optimal solution to these problems are related as follows:*

$$\{\overline{x} \,|\, \overline{x} \text{ solves } (55)\} = \{\hat{x} + V\overline{z} \,|\, \overline{z} \text{ solves } (56)\} \ .$$

By expanding the objective function in (56), we obtain

$$\tfrac{1}{2}(\hat{x} + Vz)^T H(\hat{x} + Vz) + g^T(\hat{x} + Vz) = \tfrac{1}{2}z^T V^T H V z + (V^T(H\hat{x} + g))^T v + f(\hat{x}),$$

where $f(x) := \tfrac{1}{2}x^T H x + g^T x$. If we now set $\widehat{H} := V^T H V$, $\hat{g} := V^T(H\hat{x} + g)$, and $\beta := f(\hat{x})$, then problem (56) has the form of (42):

(57)
$$\underset{z \in \mathbb{R}^k}{\text{minimize}} \ \tfrac{1}{2}z^T \widehat{H} z + \hat{g}^T z \ ,$$

where, as usual, we have dropped the constant term $\beta = f(\hat{x})$. Since we have already developed necessary and sufficient conditions for optimality in this problem, we can use them to state similar conditions for the problem (55).

THEOREM 3.1. *[Optimization of Quadratics on Affine Sets]*
*Consider the problem (55).*
   *(1) A local solution to the problem (55) exists if and only if $u^T H u \geq 0$ for all $u \in S$ and there exists a vector $\overline{x} \in \hat{x} + S$ such that $H\overline{x} + g \in S^\perp$, in which case $\overline{x}$ is a local solution to (55).*
   *(2) If $\overline{x}$ is a local solution to (55), then it is a global solution.*
   *(3) The problem (55) has a unique global solution if and only if $u^T H u > 0$ for all $u \in S \setminus \{0\}$. Moreover, if $V \in \mathbb{R}^{n \times k}$ is any matrix such that $Ran(V) = S$ where $k = \dim(S)$, then a unique solution to (55) exists if and only if the matrix $V^T H V$ is positive definite in which case the unique solution $\overline{x}$ is given by*

$$\overline{x} = [I - V(V^T H V)^{-1}V^T H]\hat{x} - V(V^T H V)^{-1}V^T g \ .$$

*(4) If either there exists $\overline{u} \in S$ such that $\overline{u}^T H \overline{u} < 0$ or there does not exist $\overline{x} \in \hat{x} + S$ such that $H\overline{x} + g \in S^{\perp}$ (or both), then*

$$-\infty = \inf_{x \in \hat{x}+S} \tfrac{1}{2} x^T H x + g^T x \ .$$

PROOF. (1) By Proposition 3.1, a solution to (55) exists if and only if a solution to (56) exists. By Theorem 2.1, a solution to (56) exists if and only if $V^T H V$ is positive semi-definite and there is a solution $\overline{z}$ to the equation $V^T(H(\hat{x} + Vz) + g) = 0$ in which case $\overline{z}$ solves (56), or equivalently, by Proposition 3.1, $\overline{x} = \hat{x} + V\overline{z}$ solves (55). The condition that $V^T H V$ is positive semi-definite is equivalent to the statement that $z^T V^T H V z \geq 0$ for all $z \in \mathbb{R}^k$, or equivalently, $u^T H u \geq 0$ for all $u \in S$. The condition, $V^T(H(\hat{x} + V\overline{z}) + g) = 0$ is equivalent to $H\overline{x} + g \in \mathrm{Null}(V^T) = \mathrm{Ran}(V)^{\perp} = S^{\perp}$.

(2) This is an immediate consequence of Proposition 3.1 and Part (2) of Theorem 2.1.

(3) By Theorem 2.1, the problem (56) has a unique solution if and only if $V^T H V$ is positive definite in which case the solution is given by $\overline{z} = (V^T H V)^{-1} V^T (H\hat{x} + g)$. Note that $V^T H V$ is positive definite if and only if $u^T H u > 0$ for all $u \in S \setminus \{0\}$ which proves that this condition is necessary and sufficient. In addition, by Proposition 3.1, $\overline{x} = \hat{x} + V\overline{z} = [I - V(V^T H V)^{-1} V^T H]\hat{x} - V(V^T H V)^{-1} V^T g$ is the unique solution to (55).

(4) This follows the same pattern of proof using Part (4) of Theorem 2.1.  □

THEOREM 3.2. *[Optimization of Quadratics Subject to Linear Equality Constraints]*
*Consider the problem* (54).

(1) *A local solution to the problem* (54) *exists if and only if $u^T H u \geq 0$ for all $u \in \mathrm{Null}(A)$ and there exists a vector pair $(\overline{x}, \overline{y}) \in \mathbb{R}^n \times \mathbb{R}^m$ such that $H\overline{x} + A^T \overline{y} + g = 0$, in which case $\overline{x}$ is a local solution to* (55).
(2) *If $\overline{x}$ is a local solution to* (55), *then it is a global solution.*
(3) *The problem* (55) *has a unique global solution if and only if $u^T H u > 0$ for all $u \in \mathrm{Null}(A) \setminus \{0\}$.*
(4) *If $u^T H u > 0$ for all $u \in \mathrm{Null}(A) \setminus \{0\}$ and $\mathrm{rank}(A) = m$, the matrix*

$$M := \begin{bmatrix} H & A^T \\ A & 0 \end{bmatrix} \quad \text{is invertible, and the vector} \quad \begin{bmatrix} \overline{x} \\ \overline{y} \end{bmatrix} = M^{-1} \begin{bmatrix} -g \\ b \end{bmatrix}$$

*has $\overline{x}$ as the unique global solution to* (55).
(5) *If either there exists $\overline{u} \in \mathrm{Null}(A)$ such that $\overline{u}^T H \overline{u} < 0$ or there does not exist a vector pair $(\overline{x}, \overline{y}) \in \mathbb{R}^n \times \mathbb{R}^m$ such that $H\overline{x} + A^T \overline{y} + g = 0$ (or both), then*

$$-\infty = \inf_{x \in \hat{x}+S} \tfrac{1}{2} x^T H x + g^T x \ .$$

REMARK 3.1. *The condition that $\mathrm{rank}(A) = m$ in Part (4) of the theorem can always be satisfied by replacing $A$ by first row reducing $A$ to echelon form.*

PROOF. (1) Recall that $\mathrm{Null}(A)^{\perp} = \mathrm{Ran}(A^T)$. Hence, $w \in \mathrm{Null}(A)$ if and only if there exists $y \in \mathbb{R}^m$ such that $w = A^T y$. By setting $w = H\overline{x} + g$ the result follows from Part (1) of Theorem 3.1.

(2) Again, this is an immediate consequence of Proposition 3.1 and Part (2) of Theorem 2.1.

(3) This is just Part (3) of Theorem 3.1.

(4) Suppose $M \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$, then $Hx + A^T y = 0$ and $Ax = 0$. If we multiply $Hx + A^T y = 0$ on the left by $x^T$, we obtain $0 = x^T H x + x^T A^T y = x^T H x$ which implies that $x = 0$ since $x \in \mathrm{Null}(A)$. But then $A^T y = 0$, so that $y = 0$ since $\mathrm{rank}(A) = m$. Consequently, $\mathrm{Null}(M) = \{0\}$, i.e., $M$ is invertible. The result now follows from Part (1).

(5) By Part (1), this is just a restatement of Theorem 3.1 Part (4).  □

The vector $\overline{y}$ appearing in this Theorem is call a *Lagrange multiplier* vector. Lagrange multiplier vectors play an essential role in constrained optimization and lie at the heart of what is called *duality theory*. This theory is more fully developed in Chapter **??**.

We now study how one might check when $H$ is positive semi-definite as well as solving the equation $Hx + g = 0$ when $H$ is positive semi-definite.

## 4. The Principal Minor Test for Positive Definiteness

Let $H \in \mathcal{S}^n$. We wish to obtain a test of when $H$ is positive definite without having to compute its eigenvalue decomposition. First note that $H_{ii} = e_i^T H e_i$, so that $H$ can be positive definite only if $H_{ii} > 0$. This is only a "sanity check" for whether a matrix is positive definite. That is, if any diagonal element of $H$ is not positive, then $H$ cannot be positive definite. In this section we develop a necessary and sufficient condition for $H$ to be positive definite that makes use of the determinant. We begin with the following lemma.

LEMMA 4.1. *Let $H \in \mathcal{S}^n$, $u \in \mathbb{R}^n$, and $\alpha \in \mathbb{R}$, and consider the block matrix*

$$\hat{H} := \begin{bmatrix} H & u \\ u^T & \alpha \end{bmatrix} \in \mathcal{S}^{(n+1)} .$$

(1) *The matrix $\hat{H}$ is positive semi-definite if and only if $H$ is positive semi-definite and there exists a vector $z \in \mathbb{R}^n$ such that $u = Hz$ and $\alpha \geq z^T H z$.*
(2) *The matrix $\hat{H}$ is positive definite if and only if $H$ is positive definite and $\alpha > u^T H^{-1} u$.*

PROOF. (1) Suppose H is positive semi-definite, and there exists $z$ such that $u = Hz$ and $\alpha \geq z^T H z$. Then for any $\hat{x} = \begin{bmatrix} x \\ x_n \end{bmatrix}$ where $x_n \in \mathbb{R}$ and $x \in \mathbb{R}^n$, we have

$$
\begin{aligned}
\hat{x}^T \hat{H} \hat{x} &= x^T H x + 2x^T H x_n z + x_n^2 \alpha \\
&= (x + x_n z)^T H(x + x_n z) + x_n^2(\alpha - z^T H z) \geq 0.
\end{aligned}
$$

Hence, $\hat{H}$ is positive semi-definite.

Conversely, suppose that $\hat{H}$ is positive semi-definite. Write $u = u_1 + u_2$ where $u_1 \in \text{Ran}(H)$ and $u_2 \in \text{Ran}(H)^{\perp} = \text{Null}(H)$, so that there is a $z \in \mathbb{R}^n$ such that $u_1 = Hz$. Then, for all $\hat{x} = \begin{pmatrix} x \\ x_n \end{pmatrix} \in \mathbb{R}^{(n+1)}$,

$$
\begin{aligned}
0 \leq \hat{x}^T \hat{H} \hat{x} &= x^T H x + 2x_n u^T x + \alpha x_n^2 \\
&= x^T H x + 2x_n (u_1 + u_2)^T x + \alpha x_n^2 \\
&= x^T H x + 2x_n z^T H x + x_n^2 z^T H z + x_n^2(\alpha - z^T H z) + 2x_n u_2^T x \\
&= (x + x_n z)^T H(x + x_n z) + x_n^2(\alpha - z^T H z) + 2x_n u_2^T x.
\end{aligned}
$$

Taking $x_n = 0$ tells us that $H$ is positive semi-definite, and taking $\hat{x} = \begin{pmatrix} -t u_2 \\ 1 \end{pmatrix}$ for $t \in \mathbb{R}$ gives

$$\alpha - 2t \left\| u_2 \right\|_2^2 \geq 0 \quad \text{for all } t \in \mathbb{R},$$

which implies that $u_2 = 0$. Finally, taking $\hat{x} = \begin{pmatrix} -z \\ 1 \end{pmatrix}$, tells us that $z^T H z \leq \alpha$ which proves the result.

(2) The proof follows the pattern of Part (1) but now we can take $z = H^{-1} u$. □

If the matrix $H$ is invertible, we can apply a kind of block Gaussian elimination to the matrix $\hat{H}$ in the lemma to obtain a matrix with block upper triangular structure:

$$\begin{bmatrix} I & 0 \\ (-H^{-1}u)^T & 1 \end{bmatrix} \begin{bmatrix} H & u \\ u^T & \alpha \end{bmatrix} = \begin{bmatrix} H & u \\ 0 & (\alpha - u^T H^{-1} u) \end{bmatrix} .$$

One consequence of this relationship is that

(58)
$$
\begin{aligned}
\det \begin{bmatrix} H & u \\ u^T & \alpha \end{bmatrix} &= \det \begin{bmatrix} I & 0 \\ (-H^{-1}u)^T & 1 \end{bmatrix} \det \begin{bmatrix} H & u \\ u^T & \alpha \end{bmatrix} \\
&= \det \left( \begin{bmatrix} I & 0 \\ (-H^{-1}u)^T & 1 \end{bmatrix} \begin{bmatrix} H & u \\ u^T & \alpha \end{bmatrix} \right) \\
&= \det \begin{bmatrix} H & u \\ 0 & (\alpha - u^T H^{-1} u) \end{bmatrix} \\
&= \det(H)(\alpha - u^T H^{-1} u).
\end{aligned}
$$

We use this determinant identity in conjunction with the previous lemma to establish a test for whether a matrix is positive definite based on determinants. The test requires us to introduce the following elementary definition.

DEFINITION 4.1. *[Principal Minors] The kth principal minor of a matrix $B \in \mathbb{R}^{n \times n}$ is the determinant of the upper left–hand corner $k \times k$–submatrix of $B$ for $1 \leq k \leq n$.*

PROPOSITION 4.1. *[The Principal Minor Test] Let $H \in \mathcal{S}^n$. Then $H$ is positive definite if and only if each of its principal minors is positive.*

PROOF. The proof proceeds by induction on the dimension $n$ of $H$. The result is clearly true for $n = 1$. We now assume the result is true for $1 \leq k \leq n$ and show it is true for dimension $n + 1$. Write

$$H := \begin{bmatrix} \hat{H} & u \\ u^T & \alpha \end{bmatrix}.$$

Then Lemma 4.1 tells us that $H$ is positive definite if and only if $\hat{H}$ is positive definite and $\alpha > u^T \hat{H}^{-1} u$. By the induction hypothesis, $\hat{H}$ is positive definite if and only if all of its principal minors are positive. If we now combine this with the expression (58), we get that $H$ is positive definite if and only if all principal minors of $\hat{H}$ are positive and, by (58), $\det(H) = \det(\hat{H})(\alpha - u^T \hat{H}^{-1} u) > 0$, or equivalently, all principal minors of $H$ are positive. □

This result only applies to positive definite matrices, and does not provide insight into how to solve linear equations involving $H$ such as $Hx + g = 0$. These two issues can be addressed through the Cholesky factorization.

## 5. The Cholesky Factorizations

We now consider how one might solve a quadratic optimization problem. Recall that a solution only exists when $H$ is positive semi-definite and there is a solution to the equation $Hx + g = 0$. Let us first consider solving the equation when $H$ is positive definite. We use a procedure similar to the LU factorization but which also takes advantage of symmetry.

Suppose

$$H = \begin{bmatrix} \alpha_1 & h_1^T \\ h_1 & \widetilde{H}_1 \end{bmatrix}, \quad \text{where} \quad \widetilde{H}_1 \in \mathcal{S}^n.$$

Note that $\alpha_1 = e_1^T H e_1 > 0$ since $H$ is positive definite (if $\alpha_1 \leq 0$, then $H$ cannot be positive definite), so there is no need to apply a permutation. Multiply $H$ on the left by the Gaussian elimination matrix for the first column, we obtain

$$L_1^{-1} H = \begin{bmatrix} 1 & 0 \\ -\frac{h_1}{\alpha_1} & I \end{bmatrix} \begin{bmatrix} \alpha_1 & h_1^T \\ h_1 & \widetilde{H}_1 \end{bmatrix} = \begin{bmatrix} \alpha_1 & h_1^T \\ 0 & \widetilde{H}_1 - \alpha^{-1} h_1 h_1^T \end{bmatrix}.$$

By symmetry, we have

$$L_1^{-1} H L_1^{-T} = \begin{bmatrix} \alpha_1 & h_1^T \\ 0 & \widetilde{H}_1 - \alpha^{-1} h_1 h_1^T \end{bmatrix} \begin{bmatrix} 1 & -\frac{h_1^T}{\alpha_1} \\ 0 & I \end{bmatrix} = \begin{bmatrix} \alpha_1 & 0 \\ 0 & \widetilde{H}_1 - \alpha^{-1} h_1 h_1^T \end{bmatrix}.$$

Set $H_1 = \widetilde{H}_1 - \alpha^{-1} h_1 h_1^T$. Observe that for every non-zero vector $v \in \mathbb{R}^{(n-1)}$,

$$v^T H_1 v = \begin{pmatrix} 0 \\ v \end{pmatrix}^T \begin{bmatrix} \alpha_1 & 0 \\ 0 & H_1 \end{bmatrix} \begin{pmatrix} 0 \\ v \end{pmatrix} = \left( L_1^{-T} \begin{pmatrix} 0 \\ v \end{pmatrix} \right)^T H \left( L_1^{-T} \begin{pmatrix} 0 \\ v \end{pmatrix} \right) > 0,$$

which shows that $H_1$ is positive definite. Decomposing $H_1$ as we did $H$ gives

$$H_1 = \begin{bmatrix} \alpha_2 & h_2^T \\ h_2 & \widetilde{H}_2 \end{bmatrix}, \quad \text{where} \quad \widetilde{H}_2 \in \mathcal{S}^{(n-1)}.$$

Again, $\alpha_2 > 0$ since $H_1$ is positive definite (if $\alpha_2 \leq 0$, then $H$ cannot be positive definite). Hence, can repeat the reduction process for $H_1$. But if at any stage we discover and $\alpha_i \leq 0$, then we terminate, since $H$ cannot be positive definite.

If we can continue this process $n$ times, we will have constructed a lower triangular matrix

$$L := L_1 L_2 \cdots L_n \quad \text{such that} \quad L^{-1} H L^{-T} = D, \quad \text{where} \quad D := \text{diag}(\alpha_1, \alpha_2, \ldots, \alpha_n)$$

is a diagonal matrix with strictly positive diagonal entries. On the other hand, if at some point in the process we discover an $\alpha_i$ that is not positive, then $H$ cannot be positive definite and the process terminates. That is, this computational procedure simultaneously tests whether $H$ is positive definite as it tries to diagonalize $H$. We will call this process the *Cholesky diagonalization procedure*. It is used to establish the following factorization theorem.

THEOREM 5.1. *[The Cholesky Factorization] Let $H \in \mathcal{S}_+^n$ have rank $k$. Then there is a lower triangular matrix $L \in \mathbb{R}^{n \times k}$ such that $H = LL^T$. Moreover, if the rank of $H$ is $n$, then there is a positive diagonal matrix $D$ and a lower triangular matrix $\widetilde{L}$ with ones on it diagonal such that $H = \widetilde{L}D\widetilde{L}^T$.*

PROOF. Let the columns of the matrix $V_1 \in \mathbb{R}^{n \times k}$ be an orthonormal basis for $\text{Ran}(H)$ and the columns of $V_2 \in \mathbb{R}^{n \times (n-k)}$ be an orthonormal basis for $\text{Null}(H)$ and set $V = [V_1 \ V_2] \in \mathbb{R}^{n \times n}$. Then

$$V^T H V = \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix} H [V_1 \ V_2]$$

$$= \begin{bmatrix} V_1^T H V_1 & V_1^T H V_2 \\ V_2^T H V_1 & V_2^T H V_2 \end{bmatrix}$$

$$= \begin{bmatrix} V_1^T H V_1 & 0 \\ 0 & 0 \end{bmatrix}.$$

Since $\text{Ran}(H) = \text{Null}(H^T)^{\perp} = \text{Null}(H)^{\perp}$, $V_1 H V_1^T \in \mathbb{R}^{k \times k}$ is symmetric and positive definite. By applying the procedure described prior to the statement of the theorem, we construct a nonsingular lower triangular matrix $\widetilde{L} \in \mathbb{R}^{k \times k}$ and a diagonal matrix $D = \text{diag}(\alpha_1, \alpha_2, \ldots, \alpha_k)$, with $\alpha_i > 0$, $i = 1, \ldots, k$, such that $V_1 H V_1^T = \widetilde{L}D\widetilde{L}^T$. Set $\widehat{L} = \widetilde{L}D^{1/2}$ so that $V_1 H V_1^T = \widehat{L}\widehat{L}^T$. If $k = n$, taking $V = I$ proves the theorem by setting $L = \widehat{L}$. If $k < n$,

$$H = [V_1 \ V_2] \begin{bmatrix} \widehat{L}\widehat{L}^T & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix} = (V_1\widehat{L})(V_1\widehat{L})^T.$$

Let $(V_1\widehat{L})^T \in \mathbb{R}^{k \times n}$ have reduced QR factorization $(V_1\widehat{L})^T = QR$ (see Theorem 5.1). Since $\widehat{L}^T$ has rank $k$, $Q \in \mathbb{R}^{k \times k}$ is unitary and $R = [R_1 \ R_2]$ with $R_1 \in \mathbb{R}^{k \times k}$ nonsingular and $R_2 \in \mathbb{R}^{k \times (n-k)}$. Therefore,

$$H = (V_1\widehat{L})(V_1\widehat{L})^T = R^T Q^T Q R = R^T R.$$

The theorem follows by setting $L = R^T$.                                                                  $\square$

When $H$ is positive definite, the factorization $H = LL^T$ is called the Cholesky factorization of $H$, and when $\text{rank}(H) < n$ it is called the *generalized Cholesky factorization* of $H$. In the positive definite case, the Cholesky diagonalization procedure computes the Cholesky factorization of $H$. On the other hand, when $H$ is only positive semi-definite, the proof of the theorem provides a guide for obtaining the generalized Cholesky factorization.

## 5.1. Computing the Generalized Cholesky Factorization.

**Step 1:** Initiate the Cholesky diagonalization procedure. If the procedure successfully completes $n$ iterations, the Cholesky factorization has been obtained. Otherwise the procedure terminates at some iteration $k+1 < n$. If $\alpha_{k+1} < 0$, proceed no further since the matrix $H$ is not positive semi-definite. If $\alpha_{k+1} = 0$, proceed to Step 2.

**Step 2:** In Step 1, the factorization

$$\widehat{L}^{-1} H \widehat{L}^{-T} = \begin{bmatrix} \widehat{D} & 0 \\ 0 & \widehat{H} \end{bmatrix},$$

where

$$\widehat{L} = \begin{bmatrix} \widehat{L}_1 & 0 \\ \widehat{L}_2 & I_{(n-k)} \end{bmatrix}$$

with $\widehat{L}_1 \in \mathbb{R}^{k \times k}$ lower triangular with ones on the diagonal, $\widehat{D} = \text{diag}(\alpha_1, \alpha_2, \ldots, \alpha_k) \in \mathbb{R}^{k \times k}$ with $\alpha_i > 0$ $i = 1, \ldots, k$, and $\widehat{H} \in \mathbb{R}^{(n-k) \times (n-k)}$ with $\widehat{H}$ symmetric has a nontrivial null space. Let the full QR factorization of $\widehat{H}$ be given by

$$\widehat{H} = [U_1 \ U_2] \begin{bmatrix} R_1 & R_2 \\ 0 & 0 \end{bmatrix} = U \begin{bmatrix} R \\ 0 \end{bmatrix},$$

where
  - $U = [U_1 \ U_2] \in \mathbb{R}^{k \times k}$ is unitary,
  - the columns of $U_1 \in \mathbb{R}^{k \times k_1}$ form an orthonormal basis for $\text{Ran}(\widehat{H})$ with $k_1 = \text{rank}\left(\widehat{H}\right) < k$,
  - the columns of $U_2 \in \mathbb{R}^{k \times (k-k_1)}$ for an orthonormal basis for $\text{Null}(\widehat{H})$,
  - $R_1 \in \mathbb{R}^{k_1 \times k_1}$ is upper triangular and nonsingular,

- $R_2 \in \mathbb{R}^{k_1 \times (k-k_1)}$, and
- $R = [R_1 \ R_2] \in \mathbb{R}^{k_1 \times k}$.

Consequently,

$$
\begin{bmatrix} U_1^T \widehat{H} U_1 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} U_1^T \widehat{H} U_1 & U_1^T \widehat{H} U_2 \\ U_2^T \widehat{H} U_1 & U_2^T \widehat{H} U_2 \end{bmatrix}
$$
$$
= U^T \widehat{H} U
$$
$$
= \begin{bmatrix} R \\ 0 \end{bmatrix} [U_1 \ U_2]
$$
$$
= \begin{bmatrix} RU_1 & RU_2 \\ 0 & 0 \end{bmatrix},
$$

and so $RU_2 = 0$ and $U_1^T \widehat{H} U_1 = RU_1 \in \mathbb{R}^{k_1 \times k_1}$ is a nonsingular symmetric matrix.

Note that only the reduced QR factorization of $H = U_1 R$ is required since $U_1^T \widehat{H} U_1 = RU_1$.

**Step 4:** Initiate the Cholesky diagonalization procedure on $U_1^T \widehat{H} U_1$. If the procedure successfully completes $k_1$ iterations, the Cholesky factorization

$$
U_1^T \widehat{H} U_1 = \widehat{L}_3 \widehat{L}_3^T
$$

has been obtained. If this does not occur, the procedure terminates at some iteration $j < k_1$ with $\alpha_j < 0$ since $U_1^T \widehat{H} U_1$ is nonsingular. In this case, terminate the process since $H$ cannot be positive semi-definfite. Otherwise proceed to Step 5.

**Step 5:** We now have

$$
H = \begin{bmatrix} \widehat{L}_1 & 0 \\ \widehat{L}_2 & I_{(n-k)} \end{bmatrix} \begin{bmatrix} \widehat{D} & 0 \\ 0 & \widehat{H} \end{bmatrix} \begin{bmatrix} \widehat{L}_1^T & \widehat{L}_2^T \\ 0 & I_{(n-k)} \end{bmatrix}
$$
$$
= \begin{bmatrix} \widehat{L}_1 & 0 \\ \widehat{L}_2 & I_{(n-k)} \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & U \end{bmatrix} \begin{bmatrix} \widehat{D} & 0 \\ 0 & U^T \widehat{H} \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & U^T \end{bmatrix} \begin{bmatrix} \widehat{L}_1^T & \widehat{L}_2^T \\ 0 & I_{(n-k)} \end{bmatrix}
$$
$$
= \begin{bmatrix} \widehat{L}_1 & 0 & 0 \\ \widehat{L}_2 & U_1 & U_2 \end{bmatrix} \begin{bmatrix} \widehat{D} & 0 & 0 \\ 0 & \widehat{L}_3 \widehat{L}_3^T & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \widehat{L}_1^T & \widehat{L}_2^T \\ 0 & U_1^T \\ 0 & U_2^T \end{bmatrix}
$$
$$
= \begin{bmatrix} \widetilde{L}_1 \widehat{D}^{1/2} & 0 & 0 \\ \widehat{L}_2 \widehat{D}^{1/2} & U_1 \widehat{L}_3 & 0 \end{bmatrix} \begin{bmatrix} \widehat{D}^{1/2} \widetilde{L}_1^T & \widehat{D}^{1/2} \widehat{L}_2^T \\ 0 & \widehat{L}_3^T U_1^T \\ 0 & 0 \end{bmatrix}
$$
$$
= \begin{bmatrix} L_1 & 0 \\ L_2 & U_1 \widehat{L}_3 \end{bmatrix} \begin{bmatrix} L_1^T & L_2^T \\ 0 & \widehat{L}_3^T U_1^T \end{bmatrix},
$$

where $L_1 = \widetilde{L}_1 \widehat{D}^{1/2} \in \mathbb{R}^{k \times k}$ is lower triangular, $L_2 = \widehat{L}_2 \widehat{D}^{1/2} \in \mathbb{R}^{(n-k) \times k}$, and $U_1 \widehat{L}_3 \in \mathbb{R}^{(n-k) \times k_1}$. In particular, $k + k_1 = \text{rank}(H)$ since $L_1$ has rank $k$ and $U_1 \widehat{L}_3$ has rank $k_1$. Let $\widehat{L}_3^T U_1^T$ have QR factorization $\widehat{L}_3^T U_1^T = V L_3^T$, where $V \in \mathbb{R}^{k_1 \times k_1}$ is unitary and $L_3 \in \mathbb{R}^{k_1 \times (n-k)}$ is lower triangular. Then

$$
H = \begin{bmatrix} L_1 & 0 \\ L_2 & U_1 \widehat{L}_3 \end{bmatrix} \begin{bmatrix} L_1^T & L_2^T \\ 0 & \widehat{L}_3^T U_1^T \end{bmatrix} = \begin{bmatrix} L_1 & 0 \\ L_2 & L_3 V^T \end{bmatrix} \begin{bmatrix} L_1^T & L_2^T \\ 0 & V L_3^T \end{bmatrix} = \begin{bmatrix} L_1 & 0 \\ L_2 & L_3 \end{bmatrix} \begin{bmatrix} L_1^T & L_2^T \\ 0 & L_3^T \end{bmatrix},
$$

since $V^T V = I_{k_1}$. This is the generalized Cholesky factorization of $H$.

### 5.2. Computing Solutions to the Quadratic Optimization Problem via Cholesky Factorizations.

**Step 1:** Apply the procedure described in the previous section for computing the generalized Cholesky factorization of $H$. If it is determined that $H$ is not positive definite, then proceed no further since the problem (42) has no solution and the optimal value is $-\infty$.

**Step 2:** Step 1 provides us with the generalized Cholesky factorization for $H = LL^T$ with $L^T = [L_1^T \ L_2^T]$, where $L_1 \in \mathbb{R}^{k \times k}$ and $L_2 \in \mathbb{R}^{(n-k) \times k}$ with $k = \text{rank}(H)$. Write

$$
g = \begin{pmatrix} g_1 \\ g_2 \end{pmatrix},
$$

where $g_1 \in \mathbb{R}^k$ and $g_2 \in \mathbb{R}^{(n-k)}$. Since $\operatorname{Ran}(H) = \operatorname{Ran}(L)$, the system $Hx + g = 0$ is solvable if and only if $-g \in \operatorname{Ran}(L)$. That is, there exists $w \in \mathbb{R}^k$ such that $Lw = -g$, or equivalently,

$$L_1 w = -g_1 \quad \text{and} \quad L_2 w = -g_2.$$

Since $L_1$ is invertible, the system $L_1 w = -g_1$ has as its unique solution $\overline{w} = L_1^{-1} g_1$. Note that $\overline{w}$ is easy to compute by forward substitution since $L_1$ is lower triangular. Having $\overline{w}$ check to see if $L_2 \overline{w} = -g_2$. If this is not the case, then proceed no further, since the system $Hx + g = 0$ has no solution and so the optimal value in (42) is $-\infty$. Otherwise, proceed to Step 3.

**Step 3:** Use back substitution to solve the equation $L_1^T y = \overline{w}$ for $\overline{y} := L_1^{-T} \overline{w}$ and set

$$\overline{x} = \begin{pmatrix} \overline{y} \\ 0 \end{pmatrix}.$$

Then

$$H\overline{x} = LL^T \overline{x} = \begin{bmatrix} L_1 \\ L_2 \end{bmatrix} [L_1^T \ L_2^T] \begin{pmatrix} \overline{y} \\ 0 \end{pmatrix} = \begin{bmatrix} L_1 \\ L_2 \end{bmatrix} \overline{w} = -g \ .$$

Hence, $\overline{x}$ solves the equation $Hx + g = 0$ and so is an optimal solution to the quadratic optimization problem (42).

## 6. Linear Least Squares Revisited

We have already see that the least squares problem is a special case of the problem of minimizing a quadratic function. But what about the reverse? Part (4) of Theorem 2.1 tells us that, in general, the reverse cannot be true since the linear least squares problem always has a solution. But what about the case when the quadratic optimization problem has a solution? In this case the matrix $H$ is necessarily positive semi-definite and a solution to the system $Hx + g = 0$ exists. By Theorem 5.1, there is a lower triangular matrix $L \in \mathbb{R}^{n \times k}$, where $k = \operatorname{rank}(H)$, such that $H = LL^T$. Set $A := L^T$. In particular, this implies that $\operatorname{Ran}(H) = \operatorname{Ran}(L) = \operatorname{Ran}(A^T)$. Since $Hx + g = 0$, we know that $-g \in \operatorname{Ran}(H) = \operatorname{Ran}(A^T)$, and so there is a vector $b \in \mathbb{R}^k$ such that $-g = A^T b$. Consider the linear least squares problem

$$\min_{x \in \mathbb{R}^n} \tfrac{1}{2} \|Ax - b\|_2^2 \ .$$

As in (44), expand the objective in this problem to obtain

$$
\begin{aligned}
\tfrac{1}{2} \|Ax - b\|_2^2 &= \tfrac{1}{2} x^T (A^T A) x - (A^T b)^T x + \tfrac{1}{2} \|b\|_2^2 \\
&= \tfrac{1}{2} x^T LL^T x + g^T x + \beta \\
&= \tfrac{1}{2} x^T Hx + g^T x + \beta,
\end{aligned}
$$

where $\beta = \tfrac{1}{2} \|b\|_2^2$. We have just proved the following result.

PROPOSITION 6.1. *A quadratic optimization problem of the form* (42) *has an optimal solution if and only if it is equivalent to a linear least squares problem.*

## 7. The Conjugate Gradient Algorithm

The Cholesky factorization is an important and useful tool for computing solutions to the quadratic optimization problem, but it is too costly to be employed in many very large scale applications. In some applications, the matrix $H$ is too large to be stored or it is not available as a data structure. However, in these problems it is often the case that the matrix vector product $Hx$ can be obtained for a given vector $x \in \mathbb{R}^n$. This occurs, for example, in a signal processing applications. In this section, we develop an algorithm for solving the quadratic optimization problem (47) that only requires access to the matrix vector products $Hx$. Such an algorithm is called a *matrix free* method since knowledge the whole matrix $H$ is not required. In such cases the Cholesky factorization is inefficient to compute. The focus of this section is the study of the matrix free method known as the *conjugate gradient algorithm*. Throughout this section we assume that the matrix $H$ is positive definite.

**7.1. Conjugate Direction Methods.** Consider the problem (47) where it is known that $H$ is symmetric and positive definite. In this case it is possible to define a notion of *orthogonality* or *conjugacy* with respect to $H$.

DEFINITION 7.1 (Conjugacy). *Let* $H \in \mathcal{S}^n_{++}$*. We say that the vectors* $x, y \in \mathbb{R}^n \backslash \{0\}$ *are* $H$*-conjugate (or* $H$*-orthogonal) if* $x^T H y = 0$*.*

PROPOSITION 7.1. *[Conjugacy implies Linear Independence]*
*If* $H \in \mathcal{S}^n_{++}$ *and the set of nonzero vectors* $d^0, d^1, \ldots, d^k$ *are (pairwise)* $H$*-conjugate, then these vectors are linearly independent.*

PROOF. If $0 = \sum_{i=0}^{k} \mu_i d^i$, then for $\bar{i} \in \{0, 1, \ldots, k\}$

$$0 = (d^{\bar{i}})^T H [\sum_{i=0}^{k} \mu_i d^i] = \mu_{\bar{i}} (d^{\bar{i}})^T H d^{\bar{i}},$$

Hence $\mu_i = 0$ for each $i = 0, \ldots, k$. □

Let $x^0 \in \mathbb{R}^n$ and suppose that the vectors $d^0, d^1, \ldots, d^{k-1} \in \mathbb{R}^n$ are $H$-conjugate. Set $S = \text{Span}(d^0, d^1, \ldots, d^{k-1})$. Theorem 3.1 tells us that there is a unique optimal solution $\bar{x}$ to the problem $\min \left\{ \frac{1}{2} x^T H x_g^T x \,\middle|\, x \in x^0 + S \right\}$, and that $\bar{x}$ is uniquely identified by the condition $H\bar{x} + g \in S^\perp$, or equivalently, $0 = (d^j)^T (H\bar{x} + g)$, $j = 0, 1, \ldots, k-1$. Since $\bar{x} \in x_0 + S$, there are scalars $\mu_0, \ldots, \mu_{n-1}$ such that

(59) $$\bar{x} = x^0 + \mu_0 d^0 + \ldots + \mu_{k-1} d^{k-1},$$

and so, for each $j = 0, 1, \ldots, k-1$,

$$\begin{aligned}
0 &= (d^j)^T (H\bar{x} + g) \\
&= (d^j)^T \left( H(x^0 + \mu_0 d^0 + \ldots + \mu_{k-1} d^{k-1}) + g \right) \\
&= (d^j)^T (Hx^0 + g) + \mu_0 (d^j)^T H d^0 + \ldots + \mu_{k-1} (d^j)^T H d^{k-1} \\
&= (d^j)^T (Hx^0 + g) + \mu_j (d^j)^T H d^j \ .
\end{aligned}$$

Therefore,

(60) $$\mu_j = \frac{-(Hx^0 + g)^T (d^j)}{(d^j)^T H d^j} \quad j = 0, 1 \ldots, k-1 \ .$$

This observation motivates the following theorem.

THEOREM 7.1. *[Expanding Subspace Theorem]*
*Consider the problem* (47) *with* $H \in \mathcal{S}^n_{++}$*, and set* $f(x) = \frac{1}{2} x^T H x + g^T x$*. Let* $\{d^i\}_{i=0}^{n-1}$ *be a sequence of nonzero* $H$*-conjugate vectors in* $\mathbb{R}^n$*. Then, for any* $x^0 \in \mathbb{R}^n$ *the sequence* $\{x^k\}$ *generated according to*

$$x^{k+1} := x^k + t_k d^k,$$

*with*

$$t_k := \arg\min\{f(x^k + td^k) : t \in \mathbb{R}\},$$

*has the property that* $f(x) = \frac{1}{2} x^T H x + g^T x$ *attains its minimum value on the affine set* $x^0 + \text{Span}\{d^0, \ldots, d^{k-1}\}$ *at the point* $x^k$*. In particular, if* $k = n$*, then* $x^n$ *is the unique global solution to the problem* (47)*.*

PROOF. Let us first compute the value of the $t_k$'s. For $k = 0, \ldots, k-1$, define $\varphi_k : \mathbb{R} \to \mathbb{R}$ by

$$\begin{aligned}
\varphi_k(t) &= f(x^k + td^k) \\
&= \frac{t^2}{2} (d^k)^T H d^k + t(g^k)^T d^k + f(x_k),
\end{aligned}$$

where $g^k = Hx^k + g$. Then, for $k = 0, \ldots, k-1$, $\varphi'_k(t) = t(d^k)^T H d^k + (g^k)^T d^k$ and $\varphi''_k(t) = (d^k)^T H d^k > 0$. Since $\varphi''_k(t) > 0$, our one dimensional calculus tells us that $\varphi_k$ attains its global minmimum value at the unique solution $t_k$ to the equation $\phi'_k(t) = 0$, i.e.,

$$t_k = -\frac{(g^k)^T d^k}{(d^k)^T H d^k}.$$

Therefore,

$$x^k = x^0 + t_0 d^0 + t_1 d^1 + \cdots + t_k d^k$$

with

$$t_k = -\frac{(g^k)^T d^k}{(d^k)^T H d_k}, \quad k = 0, 1, \ldots, k.$$

In the discussion preceding the theorem it was shown that if $\bar{x}$ is the solution to the problem

$$\min \left\{ f(x) \,\middle|\, x \in x^0 + \mathrm{Span}(d^0, d^1, \ldots, d^k) \right\},$$

then $\bar{x}$ is given by (59) and (60). Therefore, if we can now show that $\mu_j = t_j$, $j = 0, 1, \ldots, k$, then $\bar{x} = x_k$ proving the result. For each $j \in \{0, 1, \ldots, k\}$ we have

$$
\begin{aligned}
(g^j)^T d_j &= (Hx^j + g)^T d_j \\
&= \left(H(x^0 + t_0 d^0 + t_1 d^1 + \cdots + t_{j-1} d^{j-1}) + g\right)^T d^j \\
&= (Hx^0 + g)^T d^j + t_0 (d^0)^T H d^j + t_1 (d^1)^T H d^j + \cdots + t_{j-1}(d^{j-1})^T H d^j \\
&= (Hx^0 + g)^T d^j \\
&= (g^0)^T d_j .
\end{aligned}
$$

Therefore, for each $j \in \{0, 1, \ldots, k\}$,

$$t_j = \frac{-(g^j)^T d^j}{(d^j)^T H d^j} = \frac{-(g^0)^T d^j}{(d^j)^T Q d_j} = \mu_j,$$

which proves the result.                                                                                □

**7.2. The Conjugate Gradient Algorithm.** The major drawback of the Conjugate Direction Algorithm of the previous section is that it seems to require that a set of $H$-conjugate directions must be obtained before the algorithm can be implemented. This is in opposition to our working assumption that $H$ is so large that it cannot be kept in storage since any set of $H$-conjugate directions requires the same amount of storage as $H$. However, it is possible to generate the directions $d^j$ one at a time and then discard them after each iteration of the algorithm. One example of such an algorithm is the Conjugate Gradient Algorithm.

**The C-G Algorithm:**

**Initialization:** $x^0 \in \mathbb{R}^n$, $d^0 = -g^0 = -(Hx^0 + g)$.

For $k = 0, 1, 2, \ldots$

$$
\begin{aligned}
t_k &:= -(g^k)^T d^k / (d^k)^T H d^k \\
x^{k+1} &:= x^k + t_k d^k \\
g^{k+1} &:= Hx^{k+1} + g \qquad\qquad \text{(STOP if } g^{k+1} = 0) \\
\beta_k &:= (g^{k+1})^T H d^k / (d^k)^T H d^k \\
d^{k+1} &:= -g^{k+1} + \beta_k d^k \\
k &:= k + 1.
\end{aligned}
$$

THEOREM 7.2. *[*CONJUGATE GRADIENT THEOREM*]*
*The C-G algorithm is a conjugate direction method. If it does not terminate at $x^k$ (i.e. $g^k \neq 0$), then*

(1) *Span* $[g^0, g^1, \ldots, g^k] = span\ [g^0, Hg^0, \ldots, H^k g^0]$
(2) *Span* $[d^0, d^1, \ldots, d^k] = span\ [g^0, Hg^0, \ldots, H^k g^0]$
(3) $(d^k)^T Q d^i = 0$ *for* $i \leq k - 1$
(4) $t_k = (g^k)^T g^k / (d^k)^T H d^k$
(5) $\beta_k = (g^{k+1})^T g^{k+1} / (g^k)^T g^k$.

PROOF. We first prove (1)-(3) by induction. The results are clearly true for $k = 0$. Now suppose they are true for $k$, we show they are true for $k + 1$. First observe that

$$g^{k+1} = g_k + t_k H d^k$$

so that $g^{k+1} \in \mathrm{Span}[g^0, \ldots, H^{k+1} g^0]$ by the induction hypothesis on (1) and (2). Also $g^{k+1} \notin \mathrm{Span}\ [d^0, \ldots, d^k]$, otherwise, by Theorem 3.1 Part (1), $g^{k+1} = Hx^{k+1} + g = 0$ since the method is a conjugate direction method up to step $k$ by the induction hypothesis. Hence $g^{k+1} \notin \mathrm{Span}\ [g^0, \ldots, H^k g^0]$ and so $\mathrm{Span}\ [g^0, g^1, \ldots, g^{k+1}] = \mathrm{Span}\ [g^0, \ldots, H^{k+1} g^0]$, which proves (1).

To prove (2) write

$$d^{k+1} = -g^{k+1} + \beta_k d^k$$

so that (2) follows from (1) and the induction hypothesis on (2).

To see (3) observe that

$$(d^{k+1})^T H d^i = -(g^{k+1})^T H d^i + \beta_k (d^k)^T H d^i.$$

For $i = k$ the right hand side is zero by the definition of $\beta_k$. For $i < k$ both terms vanish. The term $(g^{k+1})^T H d^i = 0$ by Theorem 7.1 since $H d^i \in \mathrm{Span}[d^0, \ldots, d^k]$ by (1) and (2). The term $(d^k)^T H d^i$ vanishes by the induction hypothesis on (3).

To prove (4) write

$$-(g^k)^T d^k = (g^k)^T g^k - \beta_{k-1}(g^k)^T d^{k-1}$$

where $(g^k)^T d^{k-1} = 0$ by Theorem 7.1.

To prove (5) note that $(g^{k+1})^T g^k = 0$ by Theorem 7.1 because $g^k \in \mathrm{Span}[d^0, \ldots, d^k]$. Hence

$$(g^{k+1})^T H d^k = \frac{1}{t_k}(g^{k+1})^T[g^{k+1} - g^k] = \frac{1}{t_k}(g^{k+1})^T g^{k+1}.$$

Therefore,

$$\beta_k = \frac{1}{t_k}\frac{(g^{k+1})^T g^{k+1}}{(d^k)^T H d^k} = \frac{(g^{k+1})^T g^{k+1}}{(g^k)^T g^k}.$$

$\square$

**Remarks:**

(1) The C–G method is an example of a *descent method* since the values

$$f(x_0), \ f(x_1), \ \ldots, f(x_n)$$

form a decreasing sequence.

(2) It should be observed that due to the occurrence of round-off error the C-G algorithm is best implemented as an iterative method. That is, at the end of n steps, $x^n$ may not be the global optimal solution and the intervening directions $d^k$ may not be $H$-conjugate. Consequently, the algorithm is usually iterated until $\left\|g^k\right\|_2$ is sufficiently small. Due to the observations in the previous remark, this approach is guarenteed to continue to reduce the function value if possible since the overall method is a descent method. In this sense the C–G algorithm is self correcting.