

First and second-order conditions in constrained optimisation

Anthony Horsley and Andrew J. Wrobel

Department of Economics
London School of Economics
Houghton Street
London WC2A 2AE
United Kingdom.
e-mail: LSEecon123@mac.com

1 December 2003 (revised 24 September 2004)
CDAM Research Report LSE-CDAM-2003-15

Abstract

We give a concise but complete and detailed exposition of the “classical” approach, based on directional variations, to the first-order and second-order conditions (FOC and SOCs) for finite-dimensional constrained optimisation with both equality and inequality constraints. Attention is paid to liminal constraints, which are active inequality constraints with zero Lagrange multipliers. The persistent assertion in economics texts that all active constraints can be treated like equality constraints is untrue with liminal constraints, and it gives a false “sufficient” SOC. Nor can liminal constraints be ignored like inactive constraints; this would give a false “necessary” SOC. Treating liminal constraints like equalities in the *necessary* SOC, or ignoring them in the *sufficient* SOC is not incorrect, but it gives weaker optimality criteria than the standard Kuhn-Tucker multiplier rules (although the resulting Strong Sufficient SOC does have a place in the directional-derivative results of solution-sensitivity analysis without strict complementarity). We also show that the square slack-variables method, which reduces inequalities to equalities, cannot deal properly with liminal constraints in SOCs.

Keywords: First-order conditions, second-order conditions, liminal constraints, binding constraints.

1991 *Mathematics Subject Classification:* Primary 90C30.

1991 *Journal of Economic Literature Classification:* C61.

1 Introduction

This is an exposition of the “classical” approach, based on directional variations, to the first-order and second-order conditions (FOC and SOCs) for constrained optimisation. The method consists in linearisation, and second-order expansion, of the objective and constraint functions. It is presented in the modern geometric language of tangent and normal cones, but it originated in the Euler-Lagrange calculus of variations. It was designed and developed for infinite-dimensional problems long before being systematically applied to finite-dimensional problems (which, because of their size in practical applications, became of real interest only with the advent of high-speed computers). Though only such problems, with finite numbers of variables and constraints, are considered here explicitly, the framework is general.

The cone method we use has also a more recent, and more powerful, “nonclassical” variant which consists in linearising only *after* mapping the variables to the values of the nonlinear functions in question (objective and constraints), i.e., a suitable cone is constructed in their codomain (instead of the domain as in the “classical” approach). This gives a framework for both optimal control and the calculus of variations that improves on the classical results. We intend to present this in detail elsewhere, and here the image method is only noted, as are the augmented Lagrangian and the penalty method (Section 11).

The problem is therefore linearised at the outset: after the preliminary Section 2, in which we introduce the relevant functions and their derivatives, we discuss the tangent and normal cones to the constraint set in Section 3. In Section 4 we linearise the constraint functions and define the associated linearisation cone. This equals the tangent cone if constraints are regular. Section 5 provides workable criteria for regularity in terms of the constraint gradients, viz., linear independence and the significantly less stringent Mangasarian-Fromovitz Constraint Qualification (MFCQ) a.k.a. normality (which for convex programmes reduces essentially to Slater’s Condition). Other implications of linear independence and normality are also noted. Section 6 gives examples of irregularity.

After this preparation, the immediate goal is to state the Lagrange multiplier rules in readily applicable forms. This is done first for the case of equality constraints only (Section 7), and then for the case of both equality and inequality constraints (Section 8). The formulations assume regularity but, provided that the constraints are linearly independent (or just normal), no explicit reference to the tangent cone is needed in applying the multiplier rules. Being simpler, the pure equality-constrained case is presented separately and first, but it is of course subsumed in the more general results with both kinds of constraint, i.e., in the Kuhn-Tucker multiplier rules.

Although this is standard material, even otherwise excellent economics texts persistently give faulty SOCs by overlooking the possibility that inequality constraints may be liminal, i.e., may have zero Lagrange multipliers. In particular, they give a “sufficient” second-order condition that is in fact *insufficient*: see [7] for a counterexample. In addi-

tion, some seem to regard that “result”, which is *false* with liminal inequalities, as such an obvious extension from the equality-constrained case as not to require a proper proof [12, p. 466]. Other texts “prove” the false “theorem” [3, Theorem II.3.4 (p. 38)]. These serve as references for later texts, such as [13, E.1.16 (ii)], which also refer to mathematical sources that could have been consulted for a correct formulation of the SOC itself, as well as for a better method of proof.

As we point out in Section 9, the failure of [3] to deal correctly with inequality constraints has its roots in the method employed, which consists in converting inequalities into equalities by introducing square slack variables. Attributed to Valentine in [5, p. 39] but actually even older, the method is effective with FOCs—and it “is used sparingly” in, e.g., [5, pp. 261–262]. But as a way of obtaining SOCs, this approach is inherently limited to the case of no liminal constraints. It offers no option but to exclude these by assumption, which is known as strict complementarity. Recognising this would have at least produced true, though not the best, results. If strict complementarity were indispensable for further analysis, then arguably not much would have been lost. But since the work started in [1] and completed in [8], this is no longer justified. Earlier, strict complementarity used to be the standard assumption for a sensitivity analysis of the optimal solution, and it is indeed necessary if the solution and its multipliers are to have ordinary derivatives with respect to the problem’s parameters: see, e.g., [4, Theorems 2.4.4 and 3.2.2] or [8, Theorem 1]. But although the solution and its multipliers are usually nondifferentiable without strict complementarity, they are still *directionally* differentiable: see, e.g., [4, Theorem 2.4.5] or [8, Theorems 3 and 4]. This suffices for most purposes, as is noted in [8, Section 3].

Proofs of necessity and sufficiency of the standard FOC and SOCs are deferred to Section 10, where the multiplier rules of Sections 7 and 8 are derived from more general results with an abstract constraint set. The Abstract Necessary FOC is put into a multiplier form by applying Farkas’ Lemma, which is a separation argument extending, to the case of inequalities, the purely algebraic Factorisation Lemma. These are given in the Appendix, along with a further extension known as Motzkin’s Lemma (which is used here only to reformulate the MFCQ, but is also of interest in multi-objective optimisation).

Throughout, we follow closely the exposition of Hestenes [5, Chapter 1] and [6, Chapters 3 and 4], identifying those statements in the sources that correspond to ours. The material is selected and arranged to give a concise but complete and detailed account, and extensive explanation is added to facilitate a thorough understanding of the method, its techniques and results. Application of the theory requires little other than facility with the Jacobian and Hessian matrices (reviewed briefly in Section 2, along with the multivariate Taylor expansion to second order), and with a determinantal or eigenvalue test of definiteness for quadratic forms subject to linear restrictions.

2 Problem formulation and preliminaries

The single-objective optimisation problem in question is taken to have been oriented to maximisation over an intersection of sublevel and level sets. Here, the number of decision variables, n , is assumed to be finite. So is the number of constraints, $m + l$: there are m equality and l inequality constraints. Neither the constraint set nor the maximand need be convex. The problem, then, is one of finite-dimensional nonconvex programming.

The equality-constraint functions are h_e for $e = 1, 2, \dots, m \geq 0$ (there may be no equality constraints at all). The inequality-constraint functions are g_i for $i = 1, \dots, l$ (where $l \geq 0$). Each of these functions is assumed to be defined and twice continuously differentiable on an open set $D \subseteq \mathbb{R}^n$. So the constraint set is

$$C = \{x \in D : h(x) = 0, g(x) \leq 0\}. \quad (1)$$

The maximand, f , is also taken to be of class C^2 (twice continuously differentiable) on D .

In matrix multiplication, the n -tuple of *decision variables* $x = (x_1, \dots, x_n)$ is regarded as a column; its transpose is a row $x^T = [x_1, \dots, x_n]$. In other words, the “default” arrangement of any tuple is as a column. This applies equally to the vector-valued constraint maps $h = (h_1, \dots, h_m)$ and $g = (g_1, \dots, g_l)$, which map \mathbb{R}^n into \mathbb{R}^m and \mathbb{R}^l , respectively. (The term “function” is used to mean a scalar-valued map.)

The *Jacobian matrix* of a map $g: \mathbb{R}^n \rightarrow \mathbb{R}^l$ is the $l \times n$ matrix of partial derivatives

$$Dg(x) = \left[\frac{\partial g_i}{\partial x_j} \right]_{i=1}^l \begin{matrix} n \\ j=1 \end{matrix}$$

i.e., its i -th row is the transposed gradient vector $\nabla g_i(x)^T$.

For an l -row λ^T and an $l \times n$ matrix A , the matrix product $\lambda^T A$ is a linear combination of the rows of A . So $\lambda^T Dg(x)$ is a linear combination of the gradients $(\nabla g_i)_{i=1}^l$.

The matrix product Av , where v is an n -column, is similarly a linear combination of the columns of A . But its entries can also be viewed as the scalar products of v and the rows of A ; for example, $Dg(x)v$ has $\nabla g_i(x)^T v$ as its i -th entry. The scalar product $p^T v$ is denoted also by $p \cdot v$.

The *right kernel* of A is the null space of the linear operation $v \mapsto Av$, i.e.,

$$\ker(A \cdot) = \{v \in \mathbb{R}^n : Av = 0\}.$$

Likewise, the left kernel is $\ker(\cdot A) = \{\lambda \in \mathbb{R}^l : \lambda^T A = 0\}$.

The *Hessian matrix* of a C^2 -function $L: \mathbb{R}^n \rightarrow \mathbb{R}$ is the symmetric square ($n \times n$) matrix of second partial derivatives

$$D^2L(x) = \left(\frac{\partial^2 L}{\partial x_r \partial x_s} (x) \right)_{r=1}^n \begin{matrix} n \\ s=1 \end{matrix}.$$

When H is an $n \times n$ matrix, the function $v \mapsto v^T H v$ is called a *quadratic form* on \mathbb{R}^n .

The Euclidean norm (the length) of a $v \in \mathbb{R}^n$ is¹

$$\|v\| = (v \cdot v)^{1/2} = \left(\sum_{j=1}^n (v_j)^2 \right)^{1/2}.$$

When a vector is thought of as an increment to x , it may be denoted by Δx (which must be understood as a single symbol).

By definition, a function L is twice differentiable at a point x if²

$$\lim_{v \rightarrow 0} \frac{L(x+v) - L(x) - \nabla L(x) \cdot v}{\|v\|} = 0$$

$$\lim_{v \rightarrow 0} \frac{L(x+v) - L(x) - \nabla L(x) \cdot v - \frac{1}{2} v^T D^2 L(x) v}{\|v\|^2} = 0.$$

A sequence $(x(k))_{k=1}^{\infty}$ in \mathbb{R}^n converges to x from a *direction* Δx if there is a sequence of positive scalars $\delta(k)$ such that

$$\delta(k) \rightarrow 0 \text{ and } \frac{x(k) - x}{\delta(k)} \rightarrow \Delta x \text{ as } k \rightarrow \infty.$$

For $\Delta x = 0$, this merely means that $x(k) \rightarrow x$ as $k \rightarrow \infty$. For $\Delta x \neq 0$, $x(k) \rightarrow x$ from the direction Δx if and only if

$$x(k) \rightarrow x \text{ and } \frac{x(k) - x}{\|x(k) - x\|} \rightarrow \frac{\Delta x}{\|\Delta x\|} \text{ as } k \rightarrow \infty. \quad (2)$$

(See, e.g., [6, p. 204].) If $x(k) \rightarrow x$ from Δx then (by substituting $x(k) - x$ for the v above)

$$\frac{\nabla L(x) \cdot \Delta x}{\|\Delta x\|} = \lim_{k \rightarrow \infty} \frac{L(x(k)) - L(x)}{\|x(k) - x\|} \quad (3)$$

$$\frac{1}{2} \frac{\Delta x^T D^2 L(x) \Delta x}{\|\Delta x\|^2} = \lim_{k \rightarrow \infty} \frac{L(x(k)) - L(x) - \nabla L(x) \cdot (x(k) - x)}{\|x(k) - x\|^2}. \quad (4)$$

3 Tangent and normal cones

A *cone* at the origin 0 of a vector space is a nonempty subset that is closed under scaling³—i.e., it is a set $K \neq \emptyset$ such that $\alpha K \subseteq K$ for each $\alpha \in \mathbb{R}_+$.⁴ The *cone generated*

¹All norms on \mathbb{R}^n are equivalent, i.e., define the same convergence concept.

²This means twice Fréchet differentiability at x . It follows from the existence and continuity of second partial derivatives on a neighbourhood of x .

³By this definition, a cone need not be convex or pointed. In some literature, a *convex* set closed under scaling is called a *wedge*, and a *pointed wedge* is called a *cone*.

⁴Then $x + K$ is a cone at x .

by a set X (a.k.a. the conical hull of X) is the smallest cone containing X . It is the union of all rays from the origin through the points of X ; i.e., $\text{cone } X = \bigcup_{\alpha \in \mathbb{R}_+} \alpha X$ if $X \neq \emptyset$. Additionally $\text{cone } \emptyset := \{0\}$, the zero cone. The symbol $:=$ means that the left-hand side (l.h.s.) is by definition equal to the right-hand side (r.h.s.).

The *convex hull* of X is the smallest convex set containing X ; it is denoted by $\text{conv } X$. The *convex cone generated* by a set X is the smallest convex cone containing X ; it equals $\text{cone conv } X = \text{conv cone } X$.

The *linear span* (a.k.a. linear hull) of X is the smallest linear space containing X ; it is denoted by $\text{span } X$.

A cone K is *line-free* if $K \cap (-K) = \{0\}$, i.e., if it contains no straight line through the origin. When K is a convex cone, $K \cap (-K)$ is the largest linear space contained in K ; it is called the lineality space of K . When it is zero, K is called *pointed* (a.k.a. salient).

Definition 1 *A convex cone K is pointed if $K \cap (-K) = \{0\}$.*

The usual concept of a tangent to a curve is adequate for most applications. It has, however, a generalisation that is useful when C contains sequences convergent to x but no curve issuing from x .

Definition 2 *A nonzero vector Δx is a sequential tangent to a set C , at a point $x \in C$, if there exists a sequence $(x(k))_{k=1}^{\infty}$ in $C \setminus \{x\}$ that converges to x from the direction Δx , i.e., satisfies (2). “Sequential tangent” is henceforth abbreviated to “tangent”.*

Definition 3 *A vector Δx is a curvilinear tangent to a set C , at a point $x \in C$, if there exists a parameterised curve \tilde{x} , defined on $[0, \bar{\epsilon}]$ for some $\bar{\epsilon} > 0$, such that: $\tilde{x}(0) = x$, $\tilde{x}(\epsilon) \in C$ for every $\epsilon \in [0, \bar{\epsilon}]$, and $(d\tilde{x}/d\epsilon)(0) = \Delta x$.*

The *tangent cone*, to C at an $x \in C$, consists of all the tangent directions and the zero vector, i.e.,

$$T_x C := \text{cone} \{ \Delta x : \Delta x \text{ is a unit vector tangent to } C \text{ at } x \}.$$

In other words, a tangent vector to C at x is any limit of the directions of arbitrarily small displacements from x within C . A curvilinear tangent is a vector tangent to a curve in C that issues from the point in question, x .

Remark 4 *Every curvilinear tangent vector is a (sequential) tangent, i.e., $\Delta x \in T_x C$ if there is a curve \tilde{x} in C that meets the conditions of Definition 3.*

The tangent cone is always closed: see, e.g., [5, Lemma 1.8.1] or [6, Lemma 4.4.2]. In general, it may be nonconvex; but it is convex if C is. In the convex case there is, in a sense, little need for “proper” curves: the straight-line segments lying in C , and their limits, give all the tangent directions.

Definition 5 A vector Δx is a feasible direction (a.k.a. a linear tangent) at a point $x \in C$ if there exists a scalar $\epsilon > 0$ such that the segment joining x and $x + \epsilon\Delta x$ lies in C .

Feasible directions form a subcone of $T_x C$, denoted by

$$F_x C := \{\Delta x : [x, x + \epsilon\Delta x] \subseteq C \text{ for some } \epsilon > 0\}.$$

Remark 6 If C is convex then so is $F_x C$. Furthermore, $T_x C$ equals the closure of $F_x C$. So $T_x C$ is convex, too. Also, $C \subseteq x + F_x C \subseteq x + T_x C$ (a cone at x).⁵

An outward normal (a.k.a. exterior normal) is a vector whose scalar product with every tangent vector is zero or less.⁶

Definition 7 A vector p is normal to a set C at a point $x \in C$ if $p \cdot \Delta x \leq 0$ for every $\Delta x \in T_x C$.

The normal cone, to C at x , is

$$N_x C := \{p : p \text{ is normal to } C \text{ at } x\} = \{p : p \cdot \Delta x \leq 0 \text{ for every } \Delta x \in T_x C\}.$$

This is a case of polarity for cones.

Definition 8 For any cone $K \subseteq \mathbb{R}^n$, its polar cone is⁷

$$K^\circ := \{p \in \mathbb{R}^n : p \cdot \Delta x \leq 0 \text{ for every } \Delta x \in K\}.$$

When K is a linear subspace, K° is equal to the orthogonal complement K^\perp .

In these terms, the normal cone is polar to the tangent cone, i.e., $N_x C := (T_x C)^\circ$. Like any polar cone, $N_x C$ is always convex and closed. Finally, note that the normal cone to $T_x C$ at 0 is also $N_x C$ (since $T_0 T_x C = T_x C$).

⁵This is in [6, Lemma 4.4.3].

⁶Geometrically, an outward normal forms an obtuse angle (an unoriented angle of at least $\pi/2$ radians) to every tangent vector.

⁷Its negative, $-K^\circ$, is also known as the dual of K .

4 Linearisation of constraints and their regularity

Linearisation of a set C , around a point $\bar{x} \in C$, consists in replacing C by $\bar{x} + \mathbb{T}_{\bar{x}}C$, i.e., the condition $x \in C$ is replaced by: $x - \bar{x} \in \mathbb{T}_{\bar{x}}C$. Here, “replacing” means that local optimality of \bar{x} on the constraint set C is to be characterised by a necessary or sufficient condition which is to hold for every vector in $\mathbb{T}_{\bar{x}}C$.

If such an abstract condition is to expand into a multiplier rule, the tangent cone $\mathbb{T}_{\bar{x}}C$ must be described in terms of the gradients of the constraint functions. To this end, the constraints $h(x) = 0$ and $g(x) \leq 0$ are linearised to

$$Dh(x)(x - \bar{x}) = 0 \quad \text{and} \quad \nabla g_i(\bar{x}) \cdot (x - \bar{x}) \leq 0 \text{ for each } i \text{ with } g_i(\bar{x}) = 0 \quad (5)$$

(any inactive inequality constraint, $g_i(\bar{x}) < 0$, is irrelevant to local linearisation around \bar{x}). The linearised a.k.a. tangential constraints (5) are met whenever the increment $\Delta x = x - \bar{x}$ is a tangent vector, but in general they can also be met by some nontangent vectors. In other words, with

$$A(\bar{x}) := \{i : g_i(\bar{x}) = 0\} \quad (6)$$

denoting the set of all the *active* inequality constraints, $\mathbb{T}_{\bar{x}}C$ is always contained in the convex cone

$$L_{\bar{x}}(h, g) := \{\Delta x : Dh(\bar{x})\Delta x = 0, \nabla g_i(\bar{x}) \cdot \Delta x \leq 0 \text{ for every } i \in A(\bar{x})\}. \quad (7)$$

In the case of no inequality constraints, $g = \emptyset$ formally, and the cone $L_{\bar{x}}(h, \emptyset)$ is actually a linear space.

Since $L_{\bar{x}}(h, g)$ consists of the increments $\Delta x = x - \bar{x}$ satisfying the linearised constraints (5), it is called the *linearisation cone*. It can be larger than the tangent cone, even when the latter is also convex. For example, in the pure equality-constrained case, L_x is always a linear space, but \mathbb{T}_x can be a “proper” cone, i.e., not a linear space (Example 19 below). However, the two cones should be equal if multiplier rules are to hold—and this condition is known as regularity.

Lemma 9 $\mathbb{T}_x C \subseteq L_x(h, g)$.⁸

Proof. Take any nonzero $\Delta x \in \mathbb{T}_x C$, scale it to unit length and take a sequence $x(k)$ in C converging to x from the direction Δx . Then

$$Dh(x)\Delta x = \lim_k \frac{h(x(k)) - h(x)}{\|x(k) - x\|} = 0$$

⁸This is in, e.g., [5, p.35, lines 12–14], [6, pp. 221–222] and [10, 5.2.12].

since $h(x(k)) = h(x)$ for each k . Similarly, if $g_i(x) = 0$ then

$$\nabla g_i(x) \cdot \Delta x = \lim_k \frac{g_i(x(k)) - g_i(x)}{\|x(k) - x\|} \leq 0$$

since $g_i(x(k)) \leq 0$ for each k . ■

Definition 10 A point $x \in C$ is regular for the representation of the constraint set C by the constraint functions h and g (abbreviated to: x is regular for h and g) if $T_x C = L_x(h, g)$.

Regularity is also known as Abadie's Constraint Qualification (ACQ) and as the Basic CQ. The somewhat stronger requirement that every vector from $L_x(h, g)$ be a *curvilinear* tangent is known as Kuhn-Tucker regularity or KTCQ; it holds for linearly independent constraints (Lemma 13 below). A constraint with a vanishing gradient may fail the KTCQ and still be regular.⁹

For the first-order multiplier rule, all that matters is that the normal cone $N_x C$ be expressed in terms of the constraint gradients. Regularity ensures this because the polar of $L_x(h, g)$ is always equal to the convex cone generated by the constraint gradients (without regard to sign in the case of equalities), i.e.,

$$\begin{aligned} L_x(h, g)^\circ &= \text{cone conv}(\{\pm \nabla h_e(x) : e = 1, \dots, m\} \cup \{\nabla g_i(x) : i \in A(x)\}) \\ &= \text{span}\{\nabla h_e(x) : e = 1, \dots, m\} + \text{cone conv}\{\nabla g_i(x) : i \in A(x)\} \end{aligned} \quad (8)$$

by Farkas' Lemma (Lemma 34). And if x is regular then

$$N_x C := (T_x C)^\circ = L_x(h, g)^\circ. \quad (9)$$

Known as *quasi-regularity*, this is exactly what is needed for the Kuhn-Tucker first-order multiplier rule to hold. The condition is somewhat weaker than regularity because the two polars in (9) are equal if and only if $L_x(h, g)$ equals the closed convex hull of $T_x C$, and this can be larger than $T_x C$ itself. (It is larger if, and only if, $T_x C$ is nonconvex. For such an example, see Section 18.)

The simplest nonlinear example of (Kuhn-Tucker) regularity is a single constraint with a nonzero gradient; this generalises to regularity of linearly independent constraints and further to regularity of so-called normal constraints (Lemmas 13 and 16 below).

⁹For a one-variable example with a C^2 -constraint, take either h or g equal to $x^5 \sin(1/x)$ for $x \neq 0$, and to 0 for $x = 0$. Its derivative at $x = 0$ vanishes, and this point fails the KTCQ because the constraint set contains no interval $[-\epsilon, 0]$ or $[0, \epsilon]$ with $\epsilon > 0$. But the constraint is nevertheless regular, since the (sequential) tangent cone is $T_0 = \mathbb{R} = L_0$.

Example 11 (Single regular equality constraint) *The tangent cone to a hypersurface C with the locus equation $h(x) = 0$, at a point \bar{x} with $\nabla h(\bar{x}) \neq 0$, is the hyperplane*

$$\mathbb{T}_{\bar{x}}C = \{\Delta x : \nabla h(\bar{x}) \cdot \Delta x = 0\}. \quad (10)$$

The normal cone is the straight line

$$\mathbb{N}_{\bar{x}}C = \text{span}\{\nabla h(\bar{x})\}. \quad (11)$$

Example 12 (Single regular inequality constraint) *The tangent cone to a sublevel set $C = \{x : g(x) \leq 0\}$, at a point \bar{x} with $g(\bar{x}) = 0$ and $\nabla g(\bar{x}) \neq 0$, is the half-space*

$$\mathbb{T}_{\bar{x}}C = \{\Delta x : \nabla g(\bar{x}) \cdot \Delta x \leq 0\}. \quad (12)$$

The normal cone is the ray (half-line)

$$\mathbb{N}_{\bar{x}}C = \text{cone}\{\nabla g(\bar{x})\}. \quad (13)$$

So, with nonzero constraint gradients, regularity means that the tangent and normal cones can be expressed in terms of the tangent and normal cones to the individual constraints. Specifically, the tangent cone equals the intersection (7) of the hyperplanes (10) and half-spaces (12) which are tangent to the individual constraints (Examples 11 and 12). And the normal cone is the (algebraic) sum (8) of the straight lines (11) and outward rays (13) which are normal to the individual constraints.

Comment: Regularity can be thought of as commutativity of two operations, viz., the linearisation and the definition of a constraint set in terms of constraint functions: starting from the functions (h, g) , one can first define the set C and then linearise it to the cone $\mathbb{T}_x C$. The alternative is to linearise h and g first, and then define the cone $\mathbb{L}_x(h, g)$. Either way, the result is the same—in the regular case.

5 Linear independence and Mangasarian-Fromovitz normality as conditions for regularity

Constraints are regular if their gradients are linearly independent. Known as LICQ, this condition can be restated as a rank condition on the constraints' Jacobian—viz., that the Jacobian matrix of the equality and active inequality constraint functions be of full row-rank. For brevity, denote the number of active inequality constraints by

$$a(x) := \text{card } A(x)$$

and let $g_{A(x)}$ mean the (finite) sequence of the active constraint functions; then $Dg_{A(x)}(x)$ is the corresponding $a(x) \times n$ submatrix of $Dg(x)$.¹⁰

¹⁰Formally, $g_{A(x)} = (g_{i_k})_{k=1}^{a(x)}$, and the k -th row of $Dg_{A(x)}(x)$ is $\nabla g_{i_k}(x)^T$, where i_k increases with k and $A(x) = \{i_k : k = 1, 2, \dots, a(x)\}$.

Lemma 13 (Linearisation under LICQ) *If $h(x) = 0$ and $g(x) \leq 0$, and the vectors $(\nabla h_e(x))_{e=1}^m$ and $(\nabla g_i(x))_{i \in \Lambda(x)}$ are linearly independent or, equivalently,*

$$\text{rank} \begin{bmatrix} Dh(x) \\ Dg_{\Lambda(x)}(x) \end{bmatrix} = m + a(x) \quad (14)$$

(i.e., this matrix has a full row rank), then:¹¹

1. Every $\Delta x \in L_x(h, g)$ is a curvilinear tangent vector at x to the set C given by (1), i.e., the Kuhn-Tucker Constraint Qualification holds at x .¹²

In other words, there is a curve $\tilde{x}: [0, \bar{\epsilon}] \rightarrow C$ (for some $\bar{\epsilon} > 0$) with $\tilde{x}(0) = x$ and $(d\tilde{x}/d\epsilon)(0) = \Delta x$. This implies that $h_e(\tilde{x}(\epsilon)) = 0$ for each e and $\epsilon \leq \bar{\epsilon}$, i.e., $h \circ \tilde{x} = 0$. In addition, \tilde{x} can be chosen so that $g_i(\tilde{x}(\epsilon)) = \epsilon \nabla g_i(x) \cdot \Delta x$ for each $i \in \Lambda(x)$ and each $\epsilon \leq \bar{\epsilon}$, i.e., $g_{\Lambda(x)}(\tilde{x}(\epsilon)) = \epsilon D(g_{\Lambda(x)})(x) \Delta x$.¹³ If h and g are of class C^k (for an integer $k \geq 1$), then \tilde{x} can be chosen to be of class C^k also.

2. So $L_x(h, g) = T_x C$ (i.e., x is regular for the representation of C by h and g).

Proof. Part 2 follows from Part 1 by Remark 4.

For Part 1, $\tilde{x}(\epsilon)$ can be obtained from $x + \epsilon \Delta x$ by adding a correction term $\tilde{c}(\epsilon)$, which is generally needed because $h(x + \epsilon \Delta x)$ is not exactly 0 and $g(x + \epsilon \Delta x)$ is not linear in ϵ .¹⁴ (Nor is $g_i(x + \epsilon \Delta x)$ always nonpositive if $g_i(x) = 0$ and $\nabla g_i(x) \cdot \Delta x = 0$; although if $\nabla g_i(x) \cdot \Delta x < 0$ then $g_i(x + \epsilon \Delta x) < 0$ for small enough $\epsilon > 0$.) Fix any basis for \mathbb{R}^n ; like any vector, $\tilde{c}(\epsilon)$ is a linear combination of the basis vectors. To minimise the number of nonzero coefficients, choose a subsequence $E(1), \dots, E(m + a(x))$ of as many basis vectors as there are active inequality and equality constraints in such a way that the $(m + a(x))$ -square matrix

$$J = \begin{bmatrix} [\nabla h_e(x) \cdot E(k)]_{e=1}^m & \begin{matrix} m+a(x) \\ k=1 \\ m+a(x) \end{matrix} \\ [\nabla g_i(x) \cdot E(k)]_{i \in \Lambda(x)} & \begin{matrix} m+a(x) \\ k=1 \end{matrix} \end{bmatrix} = \begin{bmatrix} Dh(x) \\ Dg_{\Lambda(x)}(x) \end{bmatrix} E$$

is nonsingular; then seek a correction of the form

$$\tilde{c}(\epsilon) = \sum_{k=1}^{m+a(x)} E(k) \tilde{r}_k(\epsilon) = E \tilde{r}(\epsilon)$$

¹¹This is in [5, Lemma 1.10.1]; it is also a part of [6, Theorem 4.10.3]. For equality constraints only, it is also in [6, Theorem 3.5.1] and [10, 5.1.7]; this is known as the Tangent Space Theorem.

¹²The KTCQ is further discussed in, e.g., [13, 1.D].

¹³In other words, any active constraint function g_i decreases along the curve \tilde{x} at a constant nonnegative rate, viz., $-\nabla g_i(x) \cdot \Delta x$ per unit of the curve's parameter.

¹⁴The correction is *not* needed if the constraints are locally linear, i.e., of the form $h(x) = Bx + \text{const.}$, etc.—so linear constraints are always regular, even when they are linearly dependent.

where the columns $E(1), E(2), \dots$ are concatenated in the $n \times (m + a(x))$ matrix E with entries $E_{jk} := E(k)_j$. (If the unit coordinate vectors of \mathbb{R}^n are chosen as the basis, then J is simply any maximal nonsingular submatrix of the Jacobian of $(h, g_{A(x)})$, the matrix E consists of the corresponding columns of the $n \times n$ unit matrix $I(n)$, and $c = Er$ is just r but with $n - m - a(x)$ zeros inserted at the right places.)

By the Implicit Function Theorem (in, e.g., [5, Theorem 1.7.1] and [6, Theorem 3.7.1]), the system of $m + a(x)$ nonlinear equations for r

$$h(x + \epsilon \Delta x + Er) = 0 \quad (15)$$

$$g_{A(x)}(x + \epsilon \Delta x + Er) - \epsilon Dg_{A(x)}(x) \Delta x = 0 \quad (16)$$

determines $(\tilde{r}_1, \tilde{r}_2, \dots)$ as smooth functions of ϵ , on an interval $[-\bar{\epsilon}, \bar{\epsilon}]$ for some $\bar{\epsilon} > 0$, with $\tilde{r}(0) = 0$. This is because, at $(\epsilon, r) = (0, 0)$, the Jacobian w.r.t. r of the l.h.s. of (15)–(16) is the nonsingular matrix J . It also follows that

$$\frac{d\tilde{r}}{d\epsilon}(0) = -J^{-1} \begin{bmatrix} Dh(x) \Delta x \\ Dg_{A(x)}(x) \Delta x - Dg_{A(x)}(x) \Delta x \end{bmatrix} = -J^{-1} \begin{bmatrix} 0 \\ 0 \end{bmatrix} = 0$$

and so the curve

$$\tilde{x}(\epsilon) := x + \epsilon \Delta x + E\tilde{r}(\epsilon)$$

meets all the requirements: $\tilde{x}(0) = x$

$$\frac{d\tilde{x}}{d\epsilon}(0) = \Delta x + E \frac{d\tilde{r}}{d\epsilon}(0) = \Delta x$$

and $h(\tilde{x}(\epsilon)) = 0$ as well as $g_{A(x)}(\tilde{x}(\epsilon)) = \epsilon Dg_{A(x)}(x) \cdot \Delta x \leq 0$ by (15)–(16), for every positive $\epsilon \leq \bar{\epsilon}$. ■

For inequality constraints, the linear independence condition can be weakened—without loss of regularity—to that of *positive* independence, i.e., nonexistence of a vanishing linear combination with positive coefficients (Definition 37). This is equivalent to the existence of a vector that meets the linearised inequality constraints *strictly* (when substituted for $x - \bar{x}$ in (5)): see Lemma 36. So the weaker CQ, known as the Mangasarian-Fromovitz *normality*, has two equivalent forms.

Lemma 14 (Two forms of MFCQ a.k.a. normality) *Given a point $x \in C$, the following two conditions are equivalent to each other:*¹⁵

1. If

$$\mu^T Dh(x) + \lambda^T Dg_{A(x)}(x) = 0 \quad (17)$$

and $\lambda \geq 0$ (where $\lambda \in \mathbb{R}^{a(x)}$), then $\lambda = 0$ and $\mu = 0$.

¹⁵This is in [6, Theorem 4.10.4].

2. (a) The equality constraint gradients $(\nabla h_e(x))_{e=1}^m$ are linearly independent; and
 (b) there is a vector $v \in \mathbb{R}^n$ such that $Dg_{A(x)}(x)v \ll 0$ and $Dh(x)v = 0$, i.e.,

$$\nabla h_e(x) \cdot v = 0 \text{ for each } e \quad \text{and} \quad \nabla g_i(x) \cdot v < 0 \text{ for each } i \in A(x). \quad (18)$$

Proof. When $A(x) \neq \emptyset$ (i.e., there is an active inequality constraint), apply Lemma 35 (Motzkin’s Alternative) with $A = \emptyset$ —i.e., without the matrix A —to $Dh(x)$ and $Dg_{A(x)}(x)$ in place of the remaining matrices B and C (renaming ν to λ). This shows that Condition 2b (i.e., solubility of (65) with the above substitutions) is equivalent to nonexistence of a μ and $\lambda > 0$ meeting (17). Nonexistence of such μ and λ follows obviously from Condition 1 (since the latter means nonexistence of a nonzero pair μ and $\lambda \geq 0$ meeting (17)).

Furthermore, when Condition 2b holds, multiplication of (17) by such a v shows that the λ in (17) must then be 0. This means that, under Condition 2b, Condition 2a is equivalent to Condition 1. So Condition 1 is equivalent to the conjunction of Conditions 2a and 2b.

This argument applies formally also when $A(x) = \emptyset$. But this case is actually trivial: if no inequality constraint is active, then Condition 2b holds vacuously, and Condition 1 reduces to Condition 2a. ■

Geometrically, the MFCQ means that: (i) the equality constraint gradients (∇h_e) are linearly independent, (ii) the convex cone generated by the inequality constraint gradients (∇g_i) is pointed, and (iii) the interior of its polar $L_x(\emptyset, g)$, which is nonempty by Lemma 36, contains a vector tangent to the equality-constraint manifold (i.e., a vector v tangent to each hypersurface with the locus equation $h_e = 0$).

That the MFCQ is weaker than the LICQ is obvious from its first form.

Corollary 15 (LICQ and MFCQ) *If the vectors $(\nabla h_e(x))_{e=1}^m$ and $(\nabla g_i(x))_{i \in A(x)}$ are linearly independent, then the Mangasarian-Fromovitz Constraint Qualification holds at x .¹⁶*

Proof. Linear independence of all the relevant gradients $(\nabla h_e$ and ∇g_i for each active i) obviously implies Condition 1 of Lemma 14. ■

But the MFCQ is strong enough to ensure regularity.

Lemma 16 (Linearisation under MFCQ) *The Mangasarian-Fromovitz Constraint Qualification, at an $x \in C$, implies that $L_x(h, g) = T_x C$ (i.e., that the point x is regular for the representation of the constraint set C by the functions h and g).¹⁷*

¹⁶This is a part of the first implication in [6, Theorem 4.10.3]: the conclusion “ $r = p$ ” means there the same as Condition 2b of Remark 14 here. Therefore, like the Proof of Remark 14, the derivation of MFCQ from LICQ in [6, p. 239, lines 6–12] uses a theorem of the alternative (Stiemke’s, which is Lemma 38).

¹⁷This is the second implication in [6, Theorem 4.10.3], where it is derived from [6, Theorem 3.5.1] and thus from the Implicit Function Theorem, like Lemma 16 here.

Proof. When there is no inequality constraint, i.e., $l = 0$, this is a special case of Lemma 13. When $l > 0$, the Proof of Lemma 13 has to be modified. To simplify the notation, one can assume that all the inequality constraints are active.

First consider any $\Delta x \in L_x(h, g)$ with $Dg(x) \Delta x \ll 0$ (in addition to $Dh(x) \Delta x = 0$).¹⁸ If actually there is no equality constraint, i.e., $m = 0$, then—to show that Δx is a tangent to C at x —it suffices to take the curve $\epsilon \mapsto x + \epsilon \Delta x$ (since $\nabla g_i(x) \cdot \Delta x < 0$ implies that $g_i(x + \epsilon \Delta x) < g_i(0) = 0$ for small $\epsilon > 0$ and each i). But if $m > 0$ then a correction term $\tilde{c}(\epsilon)$ must be added to ensure that the curve meets the equality constraints. This is done as in the Proof of Lemma 13: from any basis for \mathbb{R}^n , choose m vectors $E(1), \dots, E(m)$ such that the $m \times m$ square matrix

$$J := Dh(x) E = \left[\nabla h_e(x) \cdot E(k) \right]_{e=1}^m \quad \begin{matrix} m \\ k=1 \end{matrix}$$

is nonsingular, and seek a correction of the form

$$\tilde{c}(\epsilon) = \sum_{k=1}^m E(k) \tilde{r}_k(\epsilon) = E \tilde{r}(\epsilon)$$

where the columns $E(1), \dots, E(m)$ are concatenated in an $n \times m$ matrix with entries $E_{jk} := E(k)_j$. By the Implicit Function Theorem, the system of m nonlinear equations for r

$$h(x + \epsilon \Delta x + Er) = 0 \tag{19}$$

determines $(\tilde{r}_1, \dots, \tilde{r}_m)$ as smooth functions of ϵ , on an interval $[-\bar{\epsilon}, \bar{\epsilon}]$ for some $\bar{\epsilon} > 0$, with $\tilde{r}(0) = 0$. This is because the Jacobian w.r.t. r of the l.h.s. of (19) is the nonsingular matrix J . It also follows that

$$\frac{d\tilde{r}}{d\epsilon}(0) = -J^{-1} Dh(x) \Delta x = -J^{-1} 0 = 0$$

and so the curve

$$\tilde{x}(\epsilon) := x + \epsilon \Delta x + E \tilde{r}(\epsilon)$$

meets all the requirements: $\tilde{x}(0) = x$,

$$\frac{d\tilde{x}}{d\epsilon}(0) = \Delta x + E \frac{d\tilde{r}}{d\epsilon}(0) = \Delta x$$

and $\tilde{x}(\epsilon) \in C$ for every sufficiently small $\epsilon > 0$. This is because $h(\tilde{x}(\epsilon)) = 0$ by (19), and because for each i

$$\frac{d}{d\epsilon} (g_i \circ \tilde{x})(0) = \nabla g_i(x) \cdot \left(\Delta x + \frac{d\tilde{r}}{d\epsilon}(0) \right) = \nabla g_i(x) \cdot \Delta x < 0 \tag{20}$$

¹⁸The MFCQ means that the set of such vectors is nonempty, in which case it is the relative interior of L_x , i.e., its interior in $\ker(Dh(x) \cdot)$.

which implies that $g_i(\tilde{x}(\epsilon)) < g_i(\tilde{x}(0)) = 0$ for small enough $\epsilon > 0$.

Finally, given an arbitrary $\Delta x \in L_x(h, g)$, use the MFCQ to take a v satisfying (18). Then $\Delta x + \epsilon v \in T_x C$ for every $\epsilon > 0$ by the preceding argument (since $Dg(x)(\Delta x + \epsilon v) \ll 0$ and $Dh(x)(\Delta x + \epsilon v)v = 0$). As $\epsilon \searrow 0$, it follows that $\Delta x \in T_x C$ because $T_x C$ is closed. ■

As the Proof of Lemma 16 shows, the MFCQ implies the existence of a point satisfying all the inequality constraints strictly, i.e., a point x^S with

$$g_i(x^S) < 0 \text{ for each } i \quad \text{and} \quad h_e(x^S) = 0 \text{ for each } e. \quad (21)$$

In convex programming this is known as Slater's Condition (SCQ). It implies, conversely, that the MFCQ holds after discarding any dependent equality constraints.

Remark 17 *The Mangasarian-Fromovitz Constraint Qualification (at any $x \in C$) implies Slater's. For the converse, assume that g_i is convex and h_e is linear, for each i and e (and their domain D is a convex open set). Then:*

1. For every $x \in C$, (21) implies that the vector $v := x - x^S$ meets (18).
2. So if I is any maximal set of linearly independent equality constraints, i.e., $(\nabla h_e)_{e \in I}$ are linearly independent with the same span as the whole system $(\nabla h_e)_{e=1}^m$, then the MFCQ holds for the representation of C by $h_I := (\nabla h_e)_{e \in I}$ and g (at every $x \in C$).

Proof. That the MFCQ implies SCQ is shown in proving Lemma 16, after (20): for a small $\epsilon > 0$, $\tilde{x}(\epsilon)$ will do as x^S in (21).

For the converse, if $g_i(x) = 0$ then, by the gradient inequality for a convex function,

$$\nabla g_i(x) \cdot (x - x^S) \leq g_i(x) - g_i(x^S) = -g_i(x^S) < 0.$$

And if $h_e(x) = 0$ then $\nabla h_e \cdot (x - x^S) = h_e(x) - h_e(x^S) = 0$. This proves Part 1. Part 2 follows because h_I defines the same constraint set as h . ■

Comments:

1. These qualifications have implications other than regularity. The MFCQ means that there is no degenerate multiplier system in the Fritz John formulation of the necessary first-order condition, which assigns a multiplier $\nu \geq 0$ to the objective f . Such a generalised multiplier system exists always, also for irregular constraints: see, e.g., [5, p. 181] or [6, Theorem 6.6.2 and 6.10.3]. When $\nu > 0$, it can be set equal to 1 by scaling. When $\nu = 0$, the multiplier system is called *degenerate* (or *singular*). Nonexistence of a degenerate multiplier system is exactly the MFCQ (as is obvious from its first form, viz., Condition 1 of Lemma 14). Regularity, being a weaker condition, ensures only that a non-degenerate multiplier system exists (Theorems 22 and 25 below), without excluding the existence of a degenerate system. LICQ guarantees that a multiplier system exists, is not degenerate, and is unique after scaling ν to 1.

2. An equality constraint can be expressed as a pair of inequality constraints; this does not affect regularity (since the linearisation cone remains unchanged). Of course, neither LICQ nor MFCQ can hold once the constraints include a pair of opposite inequalities (with the same constraint function). The conversion of an equality to a pair of inequalities by itself gives rise to a degenerate multiplier system (with just two nonzero entries); this effect can be removed by strengthening the Fritz John FOC as in, e.g., [6, Theorem 5.7.1].
3. After linearisation, equalities can be replaced by inequalities without any disadvantage; this is done in the proofs of Appendix A.

6 Examples of irregular constraint representation

As the terminology and notation of Definition 10 make clear, regularity depends on a particular choice of the functions (h and g) representing the constraint set C , and not only on C itself. This is because the linearisation cone L_x does depend on h and g , although the tangent cone T_x depends only on C . The widely used abbreviation “ x is a regular point of C ” can mislead, since x can be regular for one representation of C but irregular for another. For example, a regular description of C (with a nonzero Jacobian matrix) can always be made irregular by squaring the equalities and cubing the inequalities.¹⁹ This produces constraint functions with zero gradients; such examples exist of course even with just one variable.

With two-variables, another zero-gradient example of irregularity can be constructed from a curve with a cusp: take the single constraint $0 = h(x_1, x_2) := x_1^3 + x_2^2$. Then

$$T_{(0,0)} \{(x_1, x_2) : x_1^3 + x_2^2 = 0\} = \mathbb{R}_- \times \{0\} \neq \mathbb{R}^2 = L_{(0,0)}(h, \emptyset)$$

because $\nabla h(0) = 0$ (where $g = \emptyset$ means absence of inequality constraints)—i.e., the tangent cone is a half-line, but the linearisation cone is the whole plane. (So the constraint is not quasi-regular either: polars of the two cones are different, too.) The same is true when the same function is used in the single inequality constraint $0 \geq g(x_1, x_2) := x_1^3 + x_2^2$, i.e.,

$$T_{(0,0)} \{(x_1, x_2) : x_1^3 + x_2^2 \leq 0\} = \mathbb{R}_- \times \{0\} \neq \mathbb{R}^2 = L_{(0,0)}(\emptyset, g).$$

But the reverse inequality constraint, $g \geq 0$ or $-g \leq 0$, is regular because its tangent cone is

$$T_{(0,0)} \{(x_1, x_2) : x_1^3 + x_2^2 \geq 0\} = \mathbb{R}^2 = L_{(0,0)}(\emptyset, -g).$$

¹⁹The set C does not change when h and g are replaced by h^2 and g^3 . But at an x with $h(x) = 0$ and $g(x) = 0$, one has $D(h^2)(x) = 0$ and $D(g^3)(x) = 0$ by the Chain Rule. So $L_x(h^2, g^3)$ is the whole space \mathbb{R}^n , which cannot equal $T_x C$ because $T_x C \subseteq L_x(h, g) \neq \mathbb{R}^n$ (unless both $Dh(x) = 0$ and $Dg(x) = 0$).

This shows also that regularity does depend on whether the constraints functions serve as equality or inequality constraints, and on the sense of any inequality. By contrast, the LICQ does not distinguish between equality and active inequality constraints.

Another constraint function with a zero gradient gives an example of quasi-regularity without regularity: take the single constraint $0 = h(x_1, x_2) := x_1^2 - x_2^2$. Then T_0C is C itself (i.e., two intersecting straight lines), but $L_0(h, \emptyset)$ is the whole plane (since $\nabla h(0) = 0$). However, their polars are both equal to the zero cone. (A three-dimensional variant gives an example of quasi-regularity, without regularity, in which the tangent cone is convex but not polyhedral like the linearisation cone: take the single constraint $0 \geq g(x_1, x_2, x_3) := x_1^2 + x_2^2 - x_3^2$; then T_0C , equal to C , is the “ice-cream cone”.)

Examples of irregularity with *non-zero* constraint gradients require two variables and can be constructed from two curves tangent to each other.

Example 18 (Irregular inequality constraints) Define $g_1(x_1, x_2) = x_1^3 - x_2$ and $g_2(x_1, x_2) = x_1^3 + x_2$. (The constraint set $C = \{g \leq 0\}$ has a cusp.) At $(0, 0)$, the tangent cone is a half-line:

$$T_{(0,0)}C = \{(\Delta x_1, 0) : \Delta x_1 \leq 0\} = \mathbb{R}_- \times \{0\}.$$

But $\nabla g_2(0, 0) = (0, 1) = -\nabla g_1(0, 0)$, so the linearisation cone is a whole line:

$$L_{(0,0)}(\emptyset, g) = \mathbb{R} \times \{0\} \neq T_{(0,0)}C.$$

So the point $(0, 0)$ is not regular. (Nor is it quasi-regular: the normal cone $N = T^\circ$ is the half-space $\mathbb{R}_+ \times \mathbb{R}$, but the polar L° is the line $\{0\} \times \mathbb{R}$.)

In Example 18, both cones stay the same when just one of the inequalities is changed to an equality. Changing both inequalities into equalities gives an irregularity example with a zero tangent cone: if $h_1(x_1, x_2) = x_1^3 - x_2$ and $h_2(x_1, x_2) = x_1^3 + x_2$ then C is the single point $\{(0, 0)\}$, and $T_{(0,0)}C = \{(0, 0)\}$. But $L_{(0,0)}(h, \emptyset)$ is the line $\mathbb{R} \times \{0\}$, as before. So the point $(0, 0)$ is not regular or quasi-regular. (The normal cone $N = T^\circ$ is the half-space $\mathbb{R}_+ \times \mathbb{R}$, but the polar L° is the line $\{0\} \times \mathbb{R}$.)

A variant of Example 18 gives irregular equality-constraints with a half-line as the tangent cone (which is therefore a “proper” cone, not a linear space).

Example 19 (Irregular equality constraints) Define $h_1(x_1, x_2) = (x_1^+)^3 - x_2$ and $h_2(x_1, x_2) = (x_1^+)^3 + x_2$, where $x^+ := \max\{x, 0\}$. Then h is of class C^2 , and the constraint set is the half-line

$$C = \{(x_1, 0) : x_1 \leq 0\} = \mathbb{R}_- \times \{0\}.$$

The tangent and linearisation cones are the same as in Example 18: $T_{(0,0)}C$ is the half-line $\mathbb{R}_- \times \{0\}$, but $L_{(0,0)}(h, \emptyset)$ is the whole line $\mathbb{R} \times \{0\}$.

Adjoining a redundant constraint (a constraint implied by the other constraints) spoils linear independence. But it can never spoil regularity, and it can even regularise the constraints. Indeed, this is what happens when the redundant constraint $x_1 \leq 0$ is adjoined to Example 18 or Example 19. Of course, redundancy depends on whether the constraint is an equality or an inequality (and on the inequality's sense). For example, the single constraint $g_1(x_1, x_2) = x_1^2 + x_2 \leq 0$ is regular, and the system remains regular when the redundant constraint $x_2 \leq 0$ is adjoined. But the constraint $x_2 = 0$ is *not* redundant, and to adjoin it would make the point $(0, 0)$ irregular, with a single point as T but a whole line as L. (The same happens when the constraint $x_2 \geq 0$ is adjoined.) This gives another case in which regularity depends on whether the constraints functions serve as equality or inequality constraints (and on the sense of any inequality).

7 Multiplier rules for maxima with equality constraints only

In terms of the constraint map h with $g = \emptyset$, i.e., with no inequality constraints, the constraint set is $C = \{x : h(x) = 0\}$ and the cone $L_{\bar{x}}$ is actually a linear space. The FOC for a maximum at a regular point \bar{x} , given next, means in geometric terms that $\nabla f(\bar{x})$ is orthogonal to $L_{\bar{x}}(h, \emptyset)$, i.e., is a linear combination of $(\nabla h_e)_{e=1}^m$. Since $L_{\bar{x}}(h, \emptyset) = T_{\bar{x}}C$ by regularity, this means exactly that $\nabla f(\bar{x})$ is orthogonal to C .

Theorem 20 (Necessary FOC in Lagrange Multiplier Rule) *Assume that $h(\bar{x}) = 0$, and that \bar{x} is regular for h (with $g = \emptyset$). If \bar{x} is a local maximum point of f on C , then there exists a $\mu \in \mathbb{R}^m$ with*

$$\nabla f(\bar{x})^T = \mu^T Dh(\bar{x}) \quad (22)$$

or, in expanded form, $\nabla f(\bar{x}) = \sum_{e=1}^m \mu_e \nabla h_e(\bar{x})$. If the vectors $(\nabla h_e(\bar{x}))_{e=1}^m$ are linearly independent, i.e., $\text{rank } Dh(\bar{x}) = m$, then such a μ is unique.²⁰

Note that $\nabla_x(f - \mu \cdot h)$, which vanishes at \bar{x} , is the gradient of the *Lagrangian* a.k.a. *Lagrange function* L .²¹ Its additional arguments, μ , are called *Lagrange multipliers* (though the term is also used more narrowly to mean the particular multiplier values $\bar{\mu}$ in the FOC).

Definition 21 $L(\mu, x) := f(x) - \mu \cdot h(x)$ for every $\mu \in \mathbb{R}^m$ and $x \in D \subseteq \mathbb{R}^n$.

²⁰This is in [5, Theorem 1.9.1] and [6, Theorem 3.2.2 and p. 166, lines 4–7].

²¹In variational calculus, “Lagrange function” is a better name because “Lagrangian” means the maximand (or minimand).

Points of C meeting the Necessary FOC, called *stationary points* a.k.a. *critical points*, are examined by using the SOC. This consists in restricted definiteness of the Lagrangian's Hessian, i.e., its definiteness as a quadratic form on the subspace tangent to C . With no inequality constraints, the Lagrangian is equal to the objective function on the *whole* constraint set, i.e., $f(x) = L(\mu, x)$ for every μ and $x \in C$ (in brief, $f = L$ on C). This simplifies both formulation and proof of the SOCs.

Theorem 22 (SOCs in Lagrange Multiplier Rules) *Assume that \bar{x} is a stationary point supported by a multiplier $\bar{\mu} \in \mathbb{R}^m$, i.e., $h(\bar{x}) = 0$ and*

$$0 = \nabla_x L(\bar{\mu}, \bar{x})^T = \nabla f(\bar{x})^T - \bar{\mu}^T \text{D}h(\bar{x})$$

*or, in expanded form, $\nabla f(\bar{x}) = \sum_{e=1}^m \bar{\mu}_e \nabla h_e(\bar{x})$. Then:*²²

1. (Necessary SOC) *If \bar{x} is a local maximum point of f on C , and \bar{x} is regular for h (with $g = \emptyset$), then*

$$\Delta x^T \text{D}_{xx}^2 L(\bar{\mu}, \bar{x}) \Delta x \leq 0$$

for every Δx such that $\text{D}h(\bar{x}) \Delta x = 0$ (i.e., such that $\nabla h_e(\bar{x}) \cdot \Delta x = 0$ for each e).

2. (Sufficient SOC) *Conversely, if*

$$\Delta x^T \text{D}_{xx}^2 L(\bar{\mu}, \bar{x}) \Delta x < 0$$

for every nonzero Δx such that $\text{D}h(\bar{x}) \Delta x = 0$, then \bar{x} is a strict local maximum point of f on C . What is more, there exist numbers $\delta > 0$ and $\zeta > 0$ such that

$$f(x) < f(\bar{x}) - \zeta \|x - \bar{x}\|^2$$

for every $x \in C$ with $\|x - \bar{x}\| < \delta$ (i.e., for every x in the δ -neighbourhood of \bar{x} relative to C).²³

To use the SOCs, one needs a computational criterion of definiteness for a quadratic form on a linear subspace of \mathbb{R}^n . One such result, given next, extends Sylvester's Determinantal (Principal Minor) Test of unrestricted definiteness. It is stated in terms of the $n \times n$ symmetric matrix H representing the quadratic form and an $m \times n$ matrix B (of full row-rank) representing the subspace as $\ker(B \cdot)$.²⁴

²²This is in [5, Theorem 1.9.2], though without the tangent space spelt out, and in [6, Theorems 3.2.2 and 3.3.1].

²³One could, of course, equalise δ and ζ downwards to $\min\{\delta, \zeta\}$, but this is meaningless unless both can be measured in the same unit.

²⁴An alternative to using the extension is to change the variables, from x to a y , so that the subspace is spanned by a set of $n - m$ coordinate vectors (i.e., so that its locus equations become $\Delta y_{n-m+1} = 0, \dots, \Delta y_n = 0$ instead of $B\Delta x = 0$): then the original Sylvester's Criterion can be applied to the $(n - m)$ -square leading submatrix of the new matrix representing the quadratic.

Definition 23 The d -th leading submatrix of a square matrix M is the $d \times d$ matrix $[M_{rs}]_{r=1, s=1}^d$. Its determinant is the d -th leading minor.²⁵

Lemma 24 (Determinantal test of restricted definiteness) For a symmetric $n \times n$ matrix H and an $m \times n$ matrix B of rank m , the following conditions are equivalent to each other:

1. H is negative definite on the right kernel of B (i.e., $v^T H v < 0$ for every nonzero $v \in \mathbb{R}^n$ such that $Bv = 0$).
2. For each $d = 2m + 1, \dots, m + n$ the d -th leading minor of the $(m + n)$ -square matrix

$$\begin{bmatrix} 0 & B \\ B^T & H \end{bmatrix} \quad (23)$$

is of the sign $(-1)^{d-m}$ (i.e., the signs alternate and start from $(-1)^{m+1}$ or, equivalently, end with $(-1)^n$ for the determinant of the whole matrix (23)).

Proof. See, e.g., [2]. ■

Comments:

1. In Lemma 24, the ignored leading minors of (23) are those of dimension $d \leq 2m$. Those of dimensions $d \leq 2m - 1$ are zero. The one of dimension $2m$ is independent of H ; its sign is either $(-1)^m$ or 0.
2. For negative definiteness of H , on $\ker(B \cdot)$, the minor-sign sequence in Condition 2 of Lemma 24 must alternate and start from the correct sign $(-1)^{m+1}$ for dimension $d = 2m + 1$ or, equivalently, end with $(-1)^n$ for $d = m + n$. A constant sign equal to $(-1)^m$ indicates positive definiteness of H , on $\ker(B \cdot)$. If the sign alternates but starts from $(-1)^m$ and hence ends with $(-1)^{n-1}$, or it is constant but equal to $(-1)^{m+1}$, or it is neither constant nor alternating, then the form is *not* definite.²⁶
3. To verify the SOC for a constrained maximum (Theorem 22), the Determinantal Test (Lemma 24) can be applied to $H = D_{xx}^2 L(\bar{\mu}, \bar{x})$ with $B = -Dh(\bar{x})$. With this choice of sign, the compound matrix (23) is the *total* Hessian of L with respect to (μ, x) , i.e.,

$$\begin{bmatrix} 0 & B \\ B^T & H \end{bmatrix} = D^2 L = \begin{bmatrix} 0 & -Dh \\ -(Dh)^T & D_{xx}^2 L \end{bmatrix}. \quad (24)$$

²⁵Other matrices obtained by selecting a subset of rows and the corresponding subset of columns (with the same indices) are principal but not leading.

²⁶This is easiest to see in the unconstrained case with $n = 2$ (and $m = 0$). When the minor sign sequence is either $(+-)$ or $(--)$, this means that the eigenvalues of H are of opposite signs, i.e., the form H is *not* definite (and so, when $\nabla f(\bar{x}) = 0$ and $H = D_{xx}^2 f(\bar{x})$, the stationary point \bar{x} is a saddle point of f). When the sign sequence is either $(++)$ or $(-+)$, the eigenvalues are of the same sign, i.e., the form is definite (and so \bar{x} is a minimum or a maximum, respectively).

This is also known as the *bordered Hessian*—since the *partial* Hessian of L , w.r.t. x alone, is bordered by the constraint map’s Jacobian matrix (and its transpose).

4. The partial Hessian in (24) is that of L , *not* of f ; this matters when the constraints are nonlinear. For more explanation, see the Comment after the Proof of Theorem 22.
5. Of course, Lemma 24 can just as well be applied with $B = Dh(\bar{x})$ instead of $-Dh(\bar{x})$: the leading minors of the matrix (23) do not change when B is replaced by $-B$ (since this means changing the signs of m columns and m rows, in the d -th leading submatrix, for $d \geq 2m + 1$).
6. The SOC for a constrained maximum (Theorem 22) can also be verified by applying the determinantal test of *positive* definiteness to $-D_{xx}^2 L(\bar{\mu}, \bar{x})$. In other words, a stationary point \bar{x} (supported by $\bar{\mu}$) is a maximum if *all* the leading minors, of dimensions from $2m + 1$ to $m + n$, of the $(m + n)$ -square matrix

$$-D^2 L = \begin{bmatrix} 0 & Dh \\ (Dh)^T & -D_{xx}^2 L \end{bmatrix}$$

have the sign $(-1)^m$.

8 Multiplier rules for maxima with equality and inequality constraints

In terms of the constraint maps h and g , the constraint set is now

$$C = \{x : h(x) = 0, g(x) \leq 0\}$$

and in the FOC there are additional multipliers λ for the inequality constraints. These are always nonnegative, unlike the multipliers μ for the equality constraints, whose sign is *a priori* indefinite.

Theorem 25 (Necessary FOC in Kuhn-Tucker Multiplier Rule) *Assume that $h(\bar{x}) = 0$, $g(\bar{x}) \leq 0$, and that \bar{x} is regular for h and g . If \bar{x} is a local maximum point of f on C , then there exist a $\mu \in \mathbb{R}^m$ and a $\lambda \in \mathbb{R}^l$ with*

$$\lambda \geq 0 \tag{25}$$

$$\lambda \cdot g(\bar{x}) = 0 \tag{26}$$

$$\nabla f(\bar{x})^T = \mu^T Dh(\bar{x}) + \lambda^T Dg(\bar{x}) \tag{27}$$

i.e., in expanded form, $\nabla f(\bar{x}) = \sum_{e=1}^m \mu_e \nabla h_e(\bar{x}) + \sum_{i=1}^l \lambda_i \nabla g_i(\bar{x})$. If additionally the vectors $(\nabla h_e(\bar{x}))_{e=1}^m$ and $(\nabla g_i(\bar{x}))_{i \in A(\bar{x})}$ are linearly independent, i.e.,

$$\text{rank} \begin{bmatrix} Dh(\bar{x}) \\ Dg_{A(\bar{x})}(\bar{x}) \end{bmatrix} = m + \text{card } A(\bar{x})$$

then such a (μ, λ) is unique.²⁷

With inequalities, the constraints may have a corner point, i.e., a point at which the linearisation cone is pointed (and *a fortiori* the tangent cone is also line-free). For such points, there is a *sufficient* first-order condition.

Theorem 26 (Sufficient FOC in a multiplier rule) *Assume that \bar{x} is a stationary point supported by multipliers $\bar{\mu} \in \mathbb{R}^m$ and $\bar{\lambda} \in \mathbb{R}^l$, i.e., $h(\bar{x}) = 0$, $g(\bar{x}) \leq 0$, and*

$$\begin{aligned} 0 &\leq \bar{\lambda} \\ 0 &= \bar{\lambda} \cdot g(\bar{x}) \\ 0 &= \nabla_x L(\bar{\mu}, \bar{\lambda}, \bar{x})^\top = \nabla f(\bar{x})^\top - \bar{\mu}^\top Dh(\bar{x}) - \bar{\lambda}^\top Dg(\bar{x}) \end{aligned}$$

or, in expanded form, $\nabla f(\bar{x}) = \sum_{e=1}^m \bar{\mu}_e \nabla h_e(\bar{x}) + \sum_{i=1}^l \bar{\lambda}_i \nabla g_i(\bar{x})$. If additionally²⁸

$$\nabla f(\bar{x}) \cdot \Delta x \neq 0 \tag{28}$$

for every nonzero $\Delta x \in L_{\bar{x}}(h, g)$, then \bar{x} is a strict local maximum point of f on C . What is more, there exist numbers $\delta > 0$ and $\zeta > 0$ such that

$$f(x) \leq f(\bar{x}) - \zeta \|x - \bar{x}\| \tag{29}$$

for every $x \in C$ with $\|x - \bar{x}\| < \delta$.

But more usually the Sufficient FOC fails at a stationary point \bar{x} because, although $\nabla f(\bar{x}) \cdot \Delta x \geq 0$ for every $\Delta x \in L_{\bar{x}}$ by the Necessary FOC, the inequality is not always strict, i.e., $\nabla f(\bar{x}) \cdot \Delta x = 0$ for some nonzero $\Delta x \in L_{\bar{x}}$. Second-order conditions are then needed.

Definition 27 $L(\mu, \lambda, x) := f(x) - \mu \cdot h(x) - \lambda \cdot g(x)$ for every $(\mu, \lambda) \in \mathbb{R}^m \times \mathbb{R}^l$ and $x \in D \subseteq \mathbb{R}^n$.

²⁷This is in [5, Theorem 1.10.1], [6, Theorem 4.7.1] and [10, 5.2.18].

²⁸Since $L_{\bar{x}}(h, g)$ is a finitely generated convex cone, it suffices to assume (28) for each generator Δx . The assumption means that $\nabla f(\bar{x})$ lies in the interior of $L_{\bar{x}}(h, g)^\circ$, and it obviously implies that $L_{\bar{x}}(h, g)$ is pointed.

Unlike $\mu \cdot h$, the extra term $\lambda \cdot g$ of the Lagrangian L does not usually vanish on all of C . The condition $\lambda \cdot g(x) = 0$ is called *complementary slackness* (CS). Given a $\lambda \geq 0$, the set of those points of C which meet the CSC is

$$C_b(\lambda) := \{x \in C : \forall i (\lambda_i \neq 0 \Rightarrow g_i(x) = 0)\}. \quad (30)$$

Note that the Lagrangian is equal to the objective function only at those points which meet the constraints in a way consistent with the multiplier signs, i.e.,

$$L(\mu, \lambda, x) = f(x) \text{ if } x \in C_b(\lambda).$$

In contrast to the pure equality-constrained case, f does *not* equal L on all of C .

After solving the Necessary FOC system for (μ, λ, x) , it is known which inequality constraints are active and which of these are *binding*, i.e., have nonzero multipliers.²⁹ In the SOCs, binding constraints are treated just like equalities (and any inactive constraints are simply ignored). Active inequalities with *zero* multipliers are called *liminal* constraints; these are active but nonbinding (and hence are also known as “just active” or “degenerate active”). Except when there is just one such constraint, they cannot be ignored like inactive constraints.³⁰ Nor can liminal constraints be treated as equalities, like binding constraints. This has been persistently overlooked in economics texts: as a result, the so-called “sufficient” SOC of, e.g., [3, Theorems I.2.5 and II.3.4 (pp. 11 and 38)], [12, 19.8] and [13, 1.E.16 (ii)] is in fact *insufficient* (when there is a liminal constraint). See [7] for a counterexample. As for the error made in [3, Proof of Theorem II.3.4 (p. 38)], it is pinpointed at the end of our discussion of that method, in Section 9 before the final Comment.

The set of all the binding constraints is denoted by

$$B(\lambda) := \{i : \lambda_i \neq 0\}.$$

A binding constraint is active: for $\lambda \geq 0$ and $x \in C$, the CSC means that $B(\lambda) \subseteq A(x)$. The set of all the liminal constraints is $A(x) \setminus B(\lambda)$. In this notation

$$C_b(\lambda) = \{x : h(x) = 0, g_{B(\lambda)}(x) = 0, g_{\setminus B(\lambda)} \leq 0\}$$

where $\setminus B(\lambda) = \{1, \dots, l\} \setminus B(\lambda)$. The SOCs are formulated as though the optimisation problem had this larger system of equality constraints, viz., $(h, g_{B(\lambda)}) = 0$.

²⁹In both mathematics and economics, the terms “active”, “effective”, “tight” and “binding” are used as synonyms, which is misleading: one feels that “binding” should mean more than merely “active”.

³⁰A single liminal constraint can be ignored because definiteness of a quadratic form on a half-space is equivalent to its definiteness on the whole space.

Theorem 28 (SOCs in Kuhn-Tucker Multiplier Rules) *Assume that \bar{x} is a stationary point supported by multipliers $\bar{\mu} \in \mathbb{R}^m$ and $\bar{\lambda} \in \mathbb{R}^l$, i.e., $h(\bar{x}) = 0$, $g(\bar{x}) \leq 0$, and*

$$\begin{aligned} 0 &\leq \bar{\lambda} \\ 0 &= \bar{\lambda} \cdot g(\bar{x}) \\ 0 &= \nabla_x L(\bar{\mu}, \bar{\lambda}, \bar{x})^\top = \nabla f(\bar{x})^\top - \bar{\mu}^\top \text{D}h(\bar{x}) - \bar{\lambda}^\top \text{D}g(\bar{x}) \end{aligned}$$

or, in expanded form, $\nabla f(\bar{x}) = \sum_{e=1}^m \bar{\mu}_e \nabla h_e(\bar{x}) + \sum_{i=1}^l \bar{\lambda}_i \nabla g_i(\bar{x})$. Then.³¹

1. (Necessary SOC) *If \bar{x} is a local maximum point of f on C , and \bar{x} is regular as a point of $C_b(\bar{\lambda})$, i.e., regular for $(h, g_{\text{B}(\bar{\lambda})})$ and $g_{\text{B}(\bar{\lambda})}$, then*

$$\Delta x^\top \text{D}_{xx}^2 L(\bar{\mu}, \bar{\lambda}, \bar{x}) \Delta x \leq 0 \quad (31)$$

for every Δx such that

$$\text{D}h(\bar{x}) \Delta x = 0, \text{ i.e., } \nabla h_e(\bar{x}) \cdot \Delta x = 0 \quad \text{for each } e \quad (32)$$

$$\text{D}g_{\text{B}(\bar{\lambda})}(\bar{x}) \Delta x = 0, \text{ i.e., } \nabla g_i(\bar{x}) \cdot \Delta x = 0 \quad \text{for each } i \text{ such that } \bar{\lambda}_i > 0 \quad (33)$$

$$\text{D}g_{\text{A}(\bar{x}) \setminus \text{B}(\bar{\lambda})}(\bar{x}) \Delta x \leq 0, \text{ i.e., } \nabla g_i(\bar{x}) \cdot \Delta x \leq 0 \quad \text{for } i \text{ with } \bar{\lambda}_i = 0 \text{ and } g_i(\bar{x}) = 0. \quad (34)$$

2. (Sufficient SOC) *Conversely, if*

$$\Delta x^\top \text{D}_{xx}^2 L(\bar{\mu}, \bar{\lambda}, \bar{x}) \Delta x < 0 \quad (35)$$

for every nonzero Δx meeting (32)–(34), then \bar{x} is a strict local maximum point of f on C . What is more, there exist numbers $\delta > 0$ and $\zeta > 0$ such that

$$f(x) \leq f(\bar{x}) - \zeta \|x - \bar{x}\|^2 \quad (36)$$

for every $x \in C$ with $\|x - \bar{x}\| < \delta$.

Comments (on the SO Multiplier Rules): These can be elucidated by reflecting on C_b 's tangent and linearisation cones, and on its relationship to C .

1. Conditions (32)–(34) mean exactly that Δx is in $L_{\bar{x}}\left(\left(h, g_{\text{B}(\bar{\lambda})}\right), g_{\text{B}(\bar{\lambda})}\right)$. This must equal $T_{\bar{x}}C_b(\bar{\lambda})$ if the Necessary SOC is indeed to hold at an optimum \bar{x} (supported by $\bar{\mu}$ and $\bar{\lambda}$).

³¹This is in [6, Theorems 4.7.4 and 4.7.5] and, with all the constraints assumed active, also in [5, Theorems 1.10.2 and 1.10.3].

2. In other words, \bar{x} must be regular for the representation of $C_b(\bar{\lambda})$ by $(h, g_{B(\bar{\lambda})})$ and $g_{\setminus B(\bar{\lambda})}$. This does *not* always follow from regularity of \bar{x} for the representation of C by h and g .³²
3. But linear independence of constraint gradients (14) implies both of the regularity conditions—since the LICQ does not distinguish between equalities and active inequalities. It then follows that $T_{\bar{x}}C_b(\bar{\lambda}) = T_{\bar{x}}C \cap \bigcap_{i \in B(\bar{\lambda})} \ker \nabla g_i(\bar{x})$.
4. Assume that the LICQ holds at \bar{x} , a stationary point supported by multipliers $\bar{\mu}$ and $\bar{\lambda} \geq 0$. If additionally $A(\bar{x}) = B(\bar{\lambda})$ then $C_b(\bar{\lambda})$ is, locally around \bar{x} , equal to the topological boundary of C relative to the manifold $\{x : h(x) = 0\}$.
5. Absence of liminal constraints, known also as *strict complementarity*, simplifies the use of the sufficient SOC: if all the active inequalities are binding, then (34) plays no part, and so the range for Δx , defined by the remaining Conditions (32)–(33), is actually a linear subspace (rather than a “proper” cone). So the Determinantal Test (Lemma 24) can be applied to the bordered Hessian of $(h, g_{B(\bar{\lambda})})$, with $m + \text{card } B(\bar{\lambda})$ in place of m . The same technique applies when there is just *one* liminal constraint: it can be ignored.³³

Comments (on strict complementarity):

1. In sensitivity analysis, liminal constraints have to be excluded by assumption if the optimal solution and its multipliers are to have ordinary derivatives with respect to the problem’s parameters: see, e.g., [4, Theorems 2.4.4 and 3.2.2] or [8, Theorem 1]. Without the strict complementarity assumption, the solution (and its multipliers) is usually nondifferentiable, but it is still *directionally* differentiable: see, e.g., [4, Theorem 2.4.5] or [8, Theorems 3 and 4].
2. Some topics involve liminal constraints of necessity. An example is the Le Chatelier Principle, discussed in, e.g., [11].

³²For example, take the constraints $g_1(x_1, x_2) = -x_1^3 + x_2 \leq 0$ and $g_2(x_1, x_2) = x_1^3 + x_2 \leq 0$ (with $h = \emptyset$) at $\bar{x} = (0, 0)$; these meet the MFCQ and are therefore regular. But they cease to be regular once one or both inequalities are turned into equality constraints, as the representation of $C_b(\bar{\lambda}_1, \bar{\lambda}_2)$ requires (if $\bar{\lambda} \neq 0$). For any maximand with $\nabla f(0, 0) = (0, 1)$, the point $(0, 0)$ is stationary, and every supporting $\bar{\lambda}$ is nonzero.

³³Matters complicate with two (or more) liminal constraints. In the two-variable case (and no other constraints), the question is whether $D_{xx}^2 L$ (equal to $D_{xx}^2 f$, since $\lambda_1 = 0 = \lambda_2$) is negative definite on a planar sector between the tangents to the two borders $\{g_1 = 0\}$ and $\{g_2 = 0\}$ of the sublevel sets. In addition to negative definiteness on the two extreme rays of the sector, one needs to verify that the form does not vanish anywhere in the sector.

9 Square slack reduction to equality constraints and its limitations

Each inequality constraint, $g_i(x) \leq 0$, on the decision variables $x = (x_1, \dots, x_n)$ can be converted to an equality $g_i(x) + s_i^2/2 = 0$ by introducing a variable whose square takes up any slack. Even though the maximand f does not depend on it, s_i is an extra decision variable: the converted problem is to be solved for both x and $s = (s_1, \dots, s_m)$. This artifice can be used to derive optimality conditions for the *original* inequality-constrained problem from those for the *converted*, purely equality-constrained problem. It is adequate as a way of deriving the Kuhn-Tucker Necessary FOC, although even this necessitates the use of both the FOC and the Necessary SOC (for the converted problem): the multiplier signs (25) do *not* follow from the FOC alone.

By replacing all inequalities with equalities, the method may also seem to somehow by-pass the complication of liminal (active but not binding) inequality constraints, but in fact it offers nothing in this respect. Applied to the equality-constrained converted problem, the second-order Lagrange multiplier rules translate into inferior results for the inequality-constrained original problem: instead of the standard Kuhn-Tucker second-order rules, they yield only a weaker necessary condition and a much stronger sufficient condition. Both consequences are undesirable, and widen the gap between the necessary and the sufficient conditions. The weaker necessary condition is so because, when there are liminal constraints, it asserts the Hessian's negative semidefiniteness on a smaller set of vectors Δx than that defined by (32)–(34) in the standard Necessary SOC. The “much stronger” sufficient condition rules out liminal constraints (i.e., it implies that $g_i(x)$ and the associated multiplier μ_i cannot both be zeros). It is simply the conjunction of strict complementarity and the standard Sufficient SOC (of Theorem 28). It is, pointlessly, even stronger than the so-called Strong Sufficient SOC.³⁴ The latter, though needlessly stringent as an optimality test,³⁵ does have a place in the directional-derivative results of solution-sensitivity analysis without strict complementarity: see, e.g., [4, Theorem 2.4.5] or [8, Theorems 3 and 4].

To derive the Kuhn-Tucker Necessary FOC in this way, and to see that it yields inferior *second-order* conditions, take for simplicity a problem with inequality constraints only, i.e., maximisation of $f(x)$ over x subject to $g(x) \leq 0$; even the one-variable, one-constraint case ($n = 1, l = 1$) can provide an illustration. After conversion, the problem is to maximise $f(x)$ over $x \in \mathbb{R}^n$ and $s \in \mathbb{R}^l$ subject to $g_i(x) + s_i^2/2 = 0$ for $i = 1, \dots, l$.

³⁴The Strong Sufficient SOC is like the the standard Sufficient SOC *except* for ignoring any liminal constraints, i.e., treating them just like inactive ones. This makes no difference if there is just one liminal constraint, but if there are two or more, to ignore even one of them is to strengthen the SOC: see [7].

³⁵It is put as an optimality test in, e.g., [10, 5.3.4].

Its Lagrangian is

$$\begin{aligned}\mathcal{L}(\mu, x, s) &= f(x) - \sum_{i=1}^l \mu_i \left(g_i(x) + \frac{s_i^2}{2} \right) = f(x) - \mu \cdot g(x) - \frac{1}{2} s^T \text{diag}(\mu) s \\ &= L(\mu, x) - \frac{1}{2} s^T \text{diag}(\mu) s\end{aligned}$$

where L is the Lagrangian for the *original* problem, and $\text{diag}(\mu)$ is the square matrix with the m -tuple μ as the diagonal and with zero off-diagonal entries. By Theorem 20, the Necessary FOC for the converted problem is

$$0 = -g_i(x) - \frac{s_i^2}{2} \quad \text{for each } i \quad (37)$$

$$0 = \nabla_x \mathcal{L} = \nabla f(x) - \sum_{i=1}^l \mu_i \nabla g_i(x) \quad (38)$$

$$0 = \nabla_s \mathcal{L} = -\text{diag}(\mu) s. \quad (39)$$

By (37), $g_i(x) = -s_i^2/2 \leq 0$. And (39) means that $\mu_i s_i = 0$ for each i —i.e., $\mu_i = 0$ or $g_i(x) = 0$ for each i . To show that $\mu \geq 0$ when x is a maximum point, use the Necessary SOC on the Lagrangian's Hessian

$$D_{(x,s)(x,s)}^2 \mathcal{L}(\mu, x, s) = \begin{bmatrix} D_{xx}^2 L(\mu, x) & 0 \\ 0 & -\text{diag}(\mu) \end{bmatrix}. \quad (40)$$

That is, by Part 1 of Theorem 22, $D_{(x,s)(x,s)}^2 \mathcal{L}$ is negative semidefinite on the linear space $\ker \left(\begin{bmatrix} Dg(x) & \text{diag}(s) \end{bmatrix} \cdot \right)$, i.e.,

$$\begin{aligned}0 &\geq \begin{bmatrix} \Delta x^T & \Delta s^T \end{bmatrix} \begin{bmatrix} D_{xx}^2 L(\mu, x) & 0 \\ 0 & -\text{diag}(\mu) \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta s \end{bmatrix} \\ &= \Delta x^T D_{xx}^2 L(\mu, x) \Delta x - \Delta s^T \text{diag}(\mu) \Delta s\end{aligned} \quad (41)$$

for every $(\Delta x, \Delta s)$ such that

$$\nabla g_i(x) \cdot \Delta x + s_i \Delta s_i = 0 \quad \text{for each } i. \quad (42)$$

For $\Delta x = 0$, this means that

$$0 \leq \sum_{i=1}^l \mu_i (\Delta s_i)^2 \quad \text{for every } \Delta s \text{ with } \sum_{i=1}^l s_i \Delta s_i = 0.$$

For any i with $s_i = 0$, the last condition does *not* restrict Δs_i , and it follows that $\mu_i \geq 0$. So $\mu \geq 0$ (since $\mu_i = 0$ if $s_i \neq 0$). This completes the derivation of the Kuhn-Tucker FOC

(for the original problem): it has been shown that $g(x) \leq 0$, $\mu \geq 0$, $\mu \cdot g(x) = 0$ and $\nabla f(x)^T = \mu^T Dg(x)$. The same argument is in [3, Proof of Theorem II.3.1 (pp. 34–36)].

The inferior second-order results obtainable in this way are derived next. Given a stationary point (x, s) and a supporting multiplier μ , which meet (37)–(39), recall the notation (6) for the set of active constraints (for the original problem), i.e.,

$$A(x) = \{i : g_i(x) = 0\} = \{i : s_i = 0\}.$$

The Necessary SOC for the converted problem—i.e., the negative semidefiniteness (41) under the restrictions (42) on $(\Delta x, \Delta s)$ —can be restated as

$$\Delta x^T D_{xx}^2 L(\mu, x) \Delta x \leq \sum_{i \in A(x)} \mu_i (\Delta s_i)^2 \quad \text{if } \nabla g_i(x) \cdot \Delta x + s_i \Delta s_i = 0 \text{ for each } i \quad (43)$$

since $\mu_i = 0$ for any $i \notin A(x)$. Furthermore, for any $i \notin A(x)$ its restriction on $(\Delta x, \Delta s)$ can be simply deleted from (43). This is because it can be solved for $\Delta s_i = -\nabla g_i(x) \cdot \Delta x / s_i$ given any Δx , and because the summation in (43) excludes any such i (i.e., the sum does not depend on Δs_i). So (43) is equivalent to

$$\Delta x^T D_{xx}^2 L(\mu, x) \Delta x \leq \sum_{i \in A(x)} \mu_i (\Delta s_i)^2 \quad \text{if } \nabla g_i(x) \cdot \Delta x = 0 \text{ for each } i \in A(x).$$

Finally, since each Δs_i in the above sum is unrestricted and $\mu \geq 0$, the Necessary SOC for the converted problem is equivalent to

$$\Delta x^T D_{xx}^2 L(\mu, x) \Delta x \leq 0 \quad \text{if } \nabla g_i(x) \cdot \Delta x = 0 \text{ for each } i \in A(x). \quad (44)$$

The same argument is in [3, Proof of Theorem II.3.3 (pp. 37–38)]. Note that (44) is weaker than the standard Kuhn-Tucker Necessary SOC for the original problem because it treats liminal constraints just like equality constraints: when i is liminal (i.e., $g_i(x) = 0$ and $\mu_i = 0$), the equality restriction on Δx in (44) is more stringent than the corresponding inequality in (34). For example, when there is just one liminal constraint, the standard formulation shows negative semidefiniteness on a space of one dimension higher than does (44). In the very simplest case of just one variable x and one constraint with $dg/dx \neq 0$, if $g(x) = 0$ and $\mu = 0$ then (44) becomes the tautology $0 \leq 0$ (since Δx can only be 0), whereas the standard Kuhn-Tucker Necessary SOC means then that $d^2 f/dx^2 \leq 0$.

The Sufficient SOC for the converted problem is that (40) be negative definite under the restrictions (42), i.e., that

$$\Delta x^T D_{xx}^2 L(\mu, x) \Delta x \leq \sum_{i=1}^l \mu_i (\Delta s_i)^2 \quad \text{if } \begin{cases} \nabla g_i(x) \cdot \Delta x + s_i \Delta s_i = 0 \text{ for each } i \\ (\Delta x, \Delta s) \neq (0, 0) \end{cases} \quad (45)$$

This can be similarly restated (by recalling that $\mu_i = 0$ for any $i \notin A(x)$, and by solving again for $\Delta s_i = -\nabla g_i(x) \cdot \Delta x / s_i$, which is 0 if $\Delta x = 0$) as

$$\Delta x^T D_{xx}^2 L(\mu, x) \Delta x < \sum_{i \in A(x)} \mu_i (\Delta s_i)^2 \quad \text{if } \begin{cases} \nabla g_i(x) \cdot \Delta x = 0 \text{ for each } i \in A(x) \\ (\Delta x, \Delta s_{A(x)}) \neq (0, 0) \end{cases} . \quad (46)$$

Next, by considering three cases: (i) $\Delta x \neq 0$ but $\Delta s_A = 0$, (ii) $\Delta x = 0$ but $\Delta s_A \neq 0$, and (iii) both $\Delta x \neq 0$ and $\Delta s_A \neq 0$, (46) is shown to be equivalent to the conjunction of³⁶

$$\Delta x^T D_{xx}^2 L(\mu, x) \Delta x < 0 \quad \text{if } \Delta x \neq 0 \text{ and } \nabla g_i(x) \cdot \Delta x = 0 \text{ for each } i \in A(x) \quad (47)$$

$$\mu_i > 0 \quad \text{for each } i \in A(x) . \quad (48)$$

Without (48), Condition (47) is generally weaker than the standard Kuhn-Tucker Sufficient SOC (for the original problem) because, like (44), it treats liminal constraints just like equality constraints: when i is a liminal constraint (i.e., $g_i(x) = 0$ and $\mu_i = 0$), the equality restriction on Δx in (47) is more stringent than the corresponding inequality in (34). In other words, it is weaker because it requires the Hessian to be negative definite on a smaller set of nonzero vectors Δx than in (32)–(34). Under (48), which is strict complementarity, (47) becomes of course equivalent to the Kuhn-Tucker Sufficient SOC. Thus the method reproduces the standard sufficiency result—but *only* under the strong and extraneous assumption of strict complementarity.

Being weaker than the standard Sufficient SOC, Condition (47) is by itself insufficient for a local maximum—contrary to the assertion of [3, Theorems I.2.5 and II.3.4 (pp. 11 and 38)], which is reproduced in, e.g., [12, 19.8] and [13, 1.E.16 (ii)]. Although the argument in [3, Proof of Theorem II.3.4 (p. 38)] is largely the same as the above one, it goes wrong at the end [3, p. 38, line 3 f.b.]: in our notation, from $(\Delta x, \Delta s_{A(x)}) \neq (0, 0)$ it obviously does *not* follow that $\Delta x \neq 0$ (which *would* have given $\Delta x^T D_{xx}^2 L \Delta x < 0$). When $\Delta x = 0$, this argument must rely on positivity of $\sum_{i \in A(x)} \mu_i (\Delta s_i)^2$ for $\Delta s_{A(x)} \neq 0$ —and it necessitates the missing assumption of strict complementarity. This oversight produces a *false* “sufficient” SOC, which is in fact *insufficient* when there is a liminal constraint. When all inequality constraints are liminal,³⁷ its use can lead to a strict minimum point being misidentified as a strict maximum: see [7] for an example.

Comment: That the square-slack approach to inequalities cannot deal properly with liminal constraints is easiest to see in the one-variable, one-constraint case: the Sufficient SOC for the converted problem is that the 2×2 Hessian matrix $D_{(x,s)(x,s)}^2 \mathcal{L}$ be negative

³⁶Case (iii) does not add another condition because the inequality in (46) then follows from either of (47) and (48).

³⁷When there is a binding inequality, maximum and minimum points are distinguished from each other already by their multiplier signs in the FOCs: for a point to be stationary for both maximisation and minimisation, all the active inequality constraints must be liminal.

definite on $\ker((dg/dx, s) \cdot)$, i.e., on nonzero scalar multiples of the vector $(-s, dg/dx)$. This means that

$$0 > s^2 \left(\frac{d^2 f}{dx^2} - \mu \frac{d^2 g}{dx^2} \right) - \mu \left(\frac{dg}{dx} \right)^2 \quad (49)$$

which implies that $(\mu, s) \neq (0, 0)$, i.e., it implies strict complementarity. This can also be viewed as the simplest application of the Determinantal Test (Lemma 24): the expression in (49) equals

$$\det \begin{bmatrix} 0 & dg/dx & s \\ dg/dx & -d^2 f/dx^2 + \mu d^2 g/dx^2 & 0 \\ s & 0 & \mu \end{bmatrix}.$$

10 Maxima with abstract constraints and derivation of multiplier rules

The multiplier rules derive from the FOC and SOCs for a maximum on an arbitrary set C , which are given next.

Theorem 29 (Abstract Necessary FOC) *If \bar{x} is a local maximum point of f on C , then $\nabla f(\bar{x}) \cdot \Delta x \leq 0$ for every $\Delta x \in T_{\bar{x}}C$ (i.e., $\nabla f(\bar{x}) \in N_{\bar{x}}C$). So $\nabla f(\bar{x}) \cdot \Delta x = 0$ for every $\Delta x \in T_{\bar{x}}C \cap (-T_{\bar{x}}C)$.³⁸*

Proof. One can assume that $\|\Delta x\| = 1$. Take a sequence $(x(k))_{k=1}^{\infty}$ in $C \setminus \{\bar{x}\}$ converging to \bar{x} from the direction Δx , i.e., $(x(k) - \bar{x}) / \|x(k) - \bar{x}\| \rightarrow \Delta x$ as $k \rightarrow \infty$. Then $f(x(k)) \leq f(\bar{x})$ for every sufficiently large k , and so

$$0 \geq \lim_{k \rightarrow \infty} \frac{f(x(k)) - f(\bar{x})}{\|x(k) - \bar{x}\|} = \nabla f(\bar{x}) \cdot \Delta x$$

by (3). ■

Proof of Theorem 20 (NFOC in Lagrange Multiplier Rule, equalities only).

Here

$$C = \{x : h(x) = 0\}.$$

By Theorem 29, $\nabla f(\bar{x}) \cdot \Delta x \leq 0$ for $\Delta x \in T_{\bar{x}}C$. By the regularity assumption, $T_{\bar{x}}C$ equals $L_{\bar{x}}(h)$, which is a linear space (by its definition (7), since there are no inequality constraints). So actually

$$\nabla f(\bar{x}) \cdot \Delta x = 0 \quad \text{for every } \Delta x \in T_{\bar{x}}C = L_{\bar{x}}(h) := \ker(Dh(\bar{x}) \cdot).$$

To complete the proof, apply Lemma 33 (the Factorisation Lemma) to $p = \nabla f(\bar{x})$ and $B(e) = \nabla h_e(\bar{x})$. ■

³⁸This is in [5, Theorem 1.8.1] and in [6, Theorem 4.6.1].

Comment: So, without inequality constraints, conversion of the Abstract Necessary FOC into a Lagrange multiplier rule requires only the Factorisation Lemma, which is a purely algebraic result: it does not rely on separation of convex sets. The case of inequality constraints, dealt with next, does require a separation argument (viz., Farkas' Lemma).

Proof of Theorem 25 (NFOC in Kuhn-Tucker Multiplier Rule, with inequalities). Here

$$C = \{x : h(x) = 0, g(x) \leq 0\}.$$

By Theorem 29 and the regularity assumption,

$$\nabla f(\bar{x}) \cdot \Delta x \leq 0 \quad \text{for every } \Delta x \in T_{\bar{x}}C = L_{\bar{x}}(h, g).$$

To complete the proof, apply Lemma 34 (Farkas' Lemma) to $p = \nabla f(\bar{x})$ with $B = Dh_e(\bar{x})$ and $A = Dg_i(\bar{x})$. ■

Theorem 30 (Abstract Sufficient FOC) *If $\nabla f(\bar{x}) \cdot \Delta x < 0$ for every nonzero $\Delta x \in T_{\bar{x}}C$ (which implies that $T_{\bar{x}}C$ is line-free), then \bar{x} is a strict local maximum point of f on C .³⁹ What is more, there exist numbers $\delta > 0$ and $\zeta > 0$ such that*

$$f(x) \leq f(\bar{x}) - \zeta \|x - \bar{x}\|$$

for every $x \in C$ with $\|x - \bar{x}\| < \delta$.⁴⁰

Proof. Suppose contrarily that there is a sequence $(x(k))_{k=1}^{\infty}$ in C such that

$$x(k) \rightarrow \bar{x} \text{ as } k \rightarrow \infty \quad \text{and} \quad f(x(k)) > f(\bar{x}) - \frac{1}{k} \|x(k) - \bar{x}\| \text{ for each } k.$$

Then $x(k) \neq \bar{x}$, and since the unit sphere is compact, one can assume (by passing to a subsequence) that $(x(k) - \bar{x}) / \|x(k) - \bar{x}\|$ converges to some Δx (a unit vector). So, by (3),

$$\nabla f(x) \cdot \Delta x = \lim_{k \rightarrow \infty} \frac{f(x(k)) - f(\bar{x})}{\|x(k) - \bar{x}\|} \geq - \lim_{k \rightarrow \infty} \frac{1}{k} = 0$$

which contradicts the assumption. ■

Proof of Theorem 26 (SFOC, for corner points). For every $\Delta x \in L_{\bar{x}}(h, g)$

$$\nabla f(\bar{x}) \cdot \Delta x = \bar{\mu}^T Dh(\bar{x}) \Delta x + \bar{\lambda}^T Dg(\bar{x}) \Delta x = 0 + \bar{\lambda}^T Dg(\bar{x}) \Delta x \leq 0.$$

Since equality is excluded by the assumption (28), it follows that actually $\nabla f(\bar{x}) \cdot \Delta x < 0$ every nonzero $\Delta x \in L_{\bar{x}}(h, g)$, and *a fortiori* for every nonzero $\Delta x \in T_{\bar{x}}C$. So, by Theorem 30, \bar{x} is a strict local maximum point of f on $C = \{x : h(x) = 0, g(x) \leq 0\}$. ■

The next result will be applied to the Lagrangian L as a function of the decision variables x (with the multipliers $\bar{\mu}$ and $\bar{\lambda}$ fixed).

³⁹If, as in the case of regular constraints, $T_{\bar{x}}C$ is a finitely generated convex cone, then it obviously suffices that $\nabla f(\bar{x}) \cdot \Delta x < 0$ for each generator Δx .

⁴⁰This is in [6, Theorem 4.6.3].

Theorem 31 (Abstract SOCs) Assume that $\nabla L(\bar{x}) = 0$, for a real-valued function L that is twice differentiable at a point $\bar{x} \in C$. Then.⁴¹

1. (Necessary SOC) If \bar{x} is a local maximum point of L on C , then $\Delta x^T D^2 L(\bar{x}) \Delta x \leq 0$ for every $\Delta x \in T_{\bar{x}}C$.
2. (Sufficient SOC) Conversely, if $\Delta x^T D^2 L(\bar{x}) \Delta x < 0$ for every nonzero $\Delta x \in T_{\bar{x}}C$, then \bar{x} is a strict local maximum point of L on C . What is more, there exist numbers $\delta > 0$ and $\zeta > 0$ such that

$$L(x) \leq L(\bar{x}) - \zeta \|x - \bar{x}\|^2$$

for every $x \in C$ with $\|x - \bar{x}\| < \delta$.

Proof. For Part 1, after scaling Δx to unit length, take a sequence $(x(k))_{k=1}^{\infty}$ in $C \setminus \{\bar{x}\}$ converging to \bar{x} from the direction Δx . Then $L(x(k)) \leq L(\bar{x})$ for every sufficiently large k , and so

$$\begin{aligned} 0 &\geq \lim_{k \rightarrow \infty} \frac{L(x(k)) - L(\bar{x})}{\|x(k) - \bar{x}\|^2} = \lim_{k \rightarrow \infty} \frac{L(x(k)) - L(\bar{x}) - \nabla L(\bar{x}) \cdot (x(k) - \bar{x})}{\|x(k) - \bar{x}\|^2} \\ &= \frac{1}{2} \Delta x^T D^2 L(\bar{x}) \Delta x \end{aligned} \quad (50)$$

by (4).

For Part 2, suppose contrarily that there is a sequence $(x(k))_{k=1}^{\infty}$ in C such that

$$x(k) \rightarrow \bar{x} \text{ as } k \rightarrow \infty \quad \text{and} \quad L(x(k)) > L(\bar{x}) - \frac{1}{k} \|x(k) - \bar{x}\|^2 \text{ for each } k.$$

Then $x(k) \neq \bar{x}$, and since the unit sphere is compact, one can assume (by passing to a subsequence) that $(x(k) - \bar{x}) / \|x(k) - \bar{x}\|$ converges to some Δx (a unit vector). So, again by (4),

$$\begin{aligned} \frac{1}{2} \Delta x^T D^2 L(\bar{x}) \Delta x &= \lim_{k \rightarrow \infty} \frac{L(x(k)) - L(\bar{x}) - \nabla L(\bar{x}) \cdot (x(k) - \bar{x})}{\|x(k) - \bar{x}\|^2} \\ &= \lim_{k \rightarrow \infty} \frac{L(x(k)) - L(\bar{x})}{\|x(k) - \bar{x}\|^2} \geq - \lim_{k \rightarrow \infty} \frac{1}{k} = 0 \end{aligned} \quad (51)$$

which contradicts the assumption. ■

In the case of equality constraints only, to maximise f on C means exactly the same as to maximise $L(\bar{\mu}, \cdot)$ on C (since the two functions are simply equal on all of C).⁴²

⁴¹This is in [5, Theorem 1.8.2]; Part 1 is also in [6, Theorem 4.6.5].

⁴²In [5], this gives [5, Theorem 1.9.2] immediately from [5, Theorem 1.8.2].

So both of the above Abstract SOC's can be applied (to L) to prove the corresponding second-order multiplier rules without further ado.

Proof of Theorem 22 (SOC's as Multiplier Rules, equality constraints only).

With $\bar{\mu}$ fixed, abbreviate $L(\bar{\mu}, \cdot)$ to L ; then $\nabla L(\bar{x}) = 0$ by assumption. So Theorem 31 applies, and the result transcribes exactly into Theorem 22, since $L = f$ on $C = \{x : h(x) = 0\}$, and since $\ker(Dh(\bar{x}) \cdot) = L_{\bar{x}}(h) = T_{\bar{x}}C$ by regularity (which is needed only for Part 1, since $L_{\bar{x}}$ always contains $T_{\bar{x}}$ by Lemma 9). ■

Comment: The preceding proofs are next examined to explain why the SOC's must use the Hessian of the Lagrangian $L(\bar{\mu}, \cdot)$ and *not* that of the maximand f (unless all the constraints h are linear functions, in which case the two Hessians are equal).

1. Theorem 31 *cannot* be applied to f instead of $L(\bar{\mu}, \cdot)$ because $\nabla f(\bar{x})$ does not vanish (unlike $\nabla_x L(\bar{\mu}, \bar{x})$, which vanishes by the FOC).
2. And Theorem 31 does depend on the assumption that $\nabla L(\bar{x}) = 0$: this *cannot* be weakened to

$$\nabla L(\bar{x}) \cdot \Delta x = 0 \quad \text{for } \Delta x \in T_{\bar{x}}C. \quad (52)$$

When $T_{\bar{x}}C$ is a linear space, the weaker condition (52) is exactly the FOC for \bar{x} to be a maximum point of L (and hence it would hold in Part 1 of Theorem 31 even if $\nabla L(\bar{x})$ were nonzero). But (52) would not suffice (for the Proof of Theorem 31 if $\nabla L(\bar{x})$ were nonzero) because the term $\nabla L(\bar{x}) \cdot (x(k) - \bar{x})$ in (50) and (51) would not vanish (unless the constraints are linear, in which case C and $\bar{x} + T_{\bar{x}}C$ are locally identical, so $x(k) - \bar{x} \in T_{\bar{x}}C$ and hence $\nabla L(\bar{x}) \cdot (x(k) - \bar{x}) = 0$ purely by (52)). Nor would this term generally vanish in the limit after the division, i.e., it would not be of higher order than $\|x(k) - \bar{x}\|^2$ (although it would still be of higher order than $\|x(k) - \bar{x}\|$, purely by (3) and (52)).

Even with inequality constraints, the Necessary Abstract SOC is still applicable enough to prove the corresponding multiplier rule. It is applied, to $L(\bar{\mu}, \bar{\lambda}, \cdot)$, on the smaller set $C_b(\bar{\lambda})$ —on which f and L are still equal. Elsewhere on the constraint set, $f \leq L$. So a maximum point \bar{x} for f , which always lies in C_b , need *not* be a maximum point for L on C (since $L(\bar{x}) = f(\bar{x}) \geq f(x) \leq L(x)$ for $x \in C$). But it is, of course, a maximum point for L on C_b .

By contrast, the above Sufficient Abstract SOC is not readily applicable if the binding inequality constraints are indeed to be treated just like equality constraints.⁴³ It is to achieve this important objective that the Hessian's negative definiteness is assumed,

⁴³There are at least two advantages from being able to handle the binding constraints like equality constraints. First, the sufficient SOC matches the necessary SOC (i.e., is obtained just by making the definiteness strict). Second, when there is no liminal constraint, the negative definiteness condition is to be verified on a linear subspace (rather than on a cone).

in Part 2 of Theorem 28, only for those $\Delta x \in T_{\bar{x}}C$ with $\nabla g_i(\bar{x}) \cdot \Delta x = 0$ for each binding constraint i . Since Part 2 of Theorem 31 requires negative definiteness for every $\Delta x \in T_{\bar{x}}C$, it does not apply as it stands. What is needed, then, is a variant (of Part 2 of Theorem 31) that, from the weaker premise, draws the weaker but adequate conclusion that \bar{x} is a (local) maximum point for f , though not necessarily for L (on C).⁴⁴ Such a variant is given next.

Theorem 32 (Modified Abstract Sufficient SOC) *In addition to $\nabla L(\bar{x}) = 0$, for a real-valued L that is twice differentiable at $\bar{x} \in C$, assume that $L \geq f$ on C and $L(\bar{x}) = f(\bar{x})$. If, furthermore,*

$$\Delta x^T D^2 L(\bar{x}) \Delta x < 0$$

for every nonzero $\Delta x \in T_{\bar{x}}C$ such that $\nabla(L - f)(\bar{x}) \Delta x = 0$ (i.e., $\nabla f(\bar{x}) \cdot \Delta x = 0$), then \bar{x} is a (strict) local maximum point of f on C . What is more, there exist numbers $\delta > 0$ and $\zeta > 0$ such that

$$f(x) \leq f(\bar{x}) - \zeta \|x - \bar{x}\|^2$$

*for every $x \in C$ with $\|x - \bar{x}\| < \delta$.*⁴⁵

Proof. Suppose contrarily that there is a sequence $(x(k))_{k=1}^{\infty}$ in C such that

$$x(k) \rightarrow \bar{x} \text{ as } k \rightarrow \infty \text{ and } f(x(k)) > f(\bar{x}) - \frac{1}{k} \|x(k) - \bar{x}\|^2 \text{ for every } k.$$

Then $x(k) \neq \bar{x}$; and one can assume that $(x(k) - \bar{x}) / \|x(k) - \bar{x}\|$ converges to a unit vector Δx . So, by (4),

$$\begin{aligned} \frac{1}{2} \Delta x^T D^2 L(\bar{x}) \Delta x &= \lim_{k \rightarrow \infty} \frac{L(x(k)) - L(\bar{x})}{\|x(k) - \bar{x}\|^2} \geq \limsup_{k \rightarrow \infty} \frac{f(x(k)) - f(\bar{x})}{\|x(k) - \bar{x}\|^2} \\ &\geq \liminf_{k \rightarrow \infty} \frac{f(x(k)) - f(\bar{x})}{\|x(k) - \bar{x}\|^2} \geq - \lim_{k \rightarrow \infty} \frac{1}{k} = 0. \end{aligned}$$

This also shows that the sequence $(L - f)(x(k)) / \|x(k) - \bar{x}\|^2$ is bounded. Since $L(\bar{x}) = f(\bar{x})$, it follows that

$$0 = \lim_{k \rightarrow \infty} \frac{(L - f)(x(k))}{\|x(k) - \bar{x}\|} = \lim_{k \rightarrow \infty} \frac{(L - f)(x(k)) - (L - f)(\bar{x})}{\|x(k) - \bar{x}\|} = \nabla(L - f)(\bar{x}) \cdot \Delta x$$

by (3). This, together with $\Delta x^T D^2 L(\bar{x}) \Delta x \geq 0$, contradicts the assumption. ■

⁴⁴Optimality of \bar{x} for L on C would imply its optimality for f on C , from $f(\bar{x}) = L(\bar{x}) \geq L(x) \geq f(x)$.

⁴⁵This is in [6, Theorem 4.6.4], and implicitly in [5, Proof of Theorem 1.10.3].

Proof of Theorem 28 (SOCs as Multiplier Rules, with inequality constraints). Here

$$C = \{x : h(x) = 0, g(x) \leq 0\}.$$

With $\bar{\mu}$ and $\bar{\lambda}$ fixed, abbreviate $L(\bar{\mu}, \bar{\lambda}, \cdot)$ to L ; then $L = f$ on the set $C_b(\bar{\lambda})$ given by (30). Also, $\nabla L(\bar{x}) = 0$ by assumption.

For Part 1, since \bar{x} is a (local) maximum point of f on C , it is *a fortiori* a maximum point of L on $C_b(\bar{\lambda})$. So Theorem 31 applies, with C_b in place of C . This gives the semi-definiteness (31) of $D_{xx}^2 L$ for every $\Delta x \in T_{\bar{x}} C_b(\bar{\lambda})$, which equals $L_{\bar{x}}\left(\left(h, g_{B(\bar{\lambda})}\right), g_{\setminus B(\bar{\lambda})}\right)$ by the regularity assumption.

For Part 2, note first that $L \geq f$ on C and $L(\bar{x}) = f(\bar{x})$. So, to apply Theorem 32, one needs only to verify the definiteness (35) of $D_{xx}^2 L$ for every nonzero $\Delta x \in T_{\bar{x}} C$ such that $\nabla(L - f)(\bar{x}) \cdot \Delta x = 0$. Since $\bar{\lambda}_i = 0$ for $i \notin B(\bar{\lambda})$, and since $Dh(\bar{x}) \cdot \Delta x = 0$,

$$\nabla(L - f)(\bar{x}) \cdot \Delta x = \sum_{i \in B(\bar{\lambda})} \bar{\lambda}_i \nabla g_i(\bar{x}) \cdot \Delta x.$$

In other words, it suffices to verify (35) for every nonzero Δx in the cone

$$T_{\bar{x}} C \cap \bigcap_{i \in B(\bar{\lambda})} \ker \nabla g_i(\bar{x}) = T_{\bar{x}} C \cap \ker \left(Dg_{B(\bar{\lambda})}(\bar{x}) \cdot \right).$$

Since, by Lemma 9, this is contained in

$$L_{\bar{x}}(h, g) \cap \ker \left(Dg_{B(\bar{\lambda})}(\bar{x}) \cdot \right) = L_{\bar{x}}\left(\left(h, g_{B(\bar{\lambda})}\right), g_{\setminus B(\bar{\lambda})}\right) \quad (53)$$

the requirement is met by the assumption for Part 2 (of Theorem 28)—which is exactly that (35) holds for every nonzero Δx in the last cone in (53).⁴⁶ So Theorem 32 applies; this yields the inequality (36), which shows that \bar{x} is a strict maximum point. ■

11 Remarks on image linearisation as another approach

The above derivation of the FOCs consists in linearising the objective f and the inequality constraints g as functions of the decision variables, x (assuming for simplicity that there are no equality constraints). For a *small* increment Δx to the point in question, \bar{x} , this guarantees a close approximation to the functions' values and thus, for regular

⁴⁶Under the regularity assumption of Part 1, $L_{\bar{x}}\left(\left(h, g_B\right), g_{\setminus B}\right) = T_{\bar{x}} C_b$. This is not needed for Part 2, though.

constraints, an adequate approximation of the original problem by optimisation of the linear functional $p \cdot = \nabla f(\bar{x}) \cdot$ over the cone $L_{\bar{x}}(g)$, which is equal to the tangent at \bar{x} to the constraint set $C = \{g \leq 0\}$. When \bar{x} is a corner point of C (i.e., $L_{\bar{x}}$ is a pointed cone), this actually gives a *sufficient* FOC: if \bar{x} is a *unique* optimum for the linearised problem, then it is also a local optimum for the original problem (Theorems 26 and 30). Of more importance is the corresponding *necessary* FOC: if \bar{x} is any local optimum for the original problem, then it is also an optimum for the linearised problem (Theorem 29). This is converted into a multiplier rule (Theorem 25) by applying Farkas' Lemma, which can be proved by separating a point from a closed convex cone with a hyperplane. The separation argument has two versions dual to each other. In the First Proof of Lemma 34, the point $p = \nabla f(\bar{x})$ is separated from the relevant cone by means of a $v \in \mathbb{R}^n$, i.e., a vector of decision variables with the interpretation of an increment to \bar{x} .⁴⁷ In the Second Proof of Lemma 34, an orthant with its vertex at $(1, 0) \in \mathbb{R} \times \mathbb{R}^l$, where l is the number of constraints, is separated from the range of the *linearised* objective and constraint functions, i.e., from the image of \mathbb{R}^n under the map $v \mapsto (p \cdot v, Av)$ with $p = \nabla f(\bar{x})$ and $A = Dg(\bar{x})$. This time, the sets being separated lie in the space of values of the objective and constraint functions, and they are separated by means of a vector $(-1, \bar{\lambda})$ that gives the multipliers $\bar{\lambda} \in \mathbb{R}^l$.

The latter separation argument can also be set up by linearising only *after* mapping the variables space (here, \mathbb{R}^n or a subset D) to the values space (here, \mathbb{R}^{1+l}) by means of the nonlinear objective and constraint functions (f, g) . The image set $\{(f(x), g(x)) : x \in D\}$ is then linearised, around a point $(f(\bar{x}), g(\bar{x}))$, by taking either the whole tangent cone, if it is convex, or a suitably chosen convex subcone of the tangent: see [5, Chapter 4] or [6, Chapter 6].⁴⁸ By dealing primarily with the functions' values rather than their variables, this method allows any variation of \bar{x} , say to $\tilde{x}(\epsilon)$, that results in small changes of the relevant values (f and g). This is particularly useful when the decision variables form an infinite-dimensional space X , since it allows also the use of a nondirectional variation $\bar{x} - x(\epsilon)$ that may converge to zero in some topology on X , as $\epsilon \searrow 0$, but not necessarily from any particular direction in the specified space X . When X is a space of absolutely continuous functions, one well-known example is Weierstrass' needle-like variation. Used originally to derive his necessary condition for a strong extremum, the method was later adapted for a derivation of Pontriagin's Maximum Principle. Thus the image approach provides a unifying framework for the classical calculus of variations and the more recent optimal control theory: see [5, Chapters 4 ff.].

Another approach to multiplier rules consists in convexifying the Lagrangian by augmenting it with quadratic penalty terms. This leads to the principle of constraint removal and to an algorithm using both penalties and multipliers: see [6, Chapter 5].

⁴⁷The separated cone and point (p) lie in what is the dual parameter space when the linearised problem is regarded as the primal, in the duality framework of linear or convex programming.

⁴⁸Even the case of a convex image set captures some nonconvex programmes: see [6, Example 6.6.1].

A Theorems of the Alternative

The linear algebra tool for converting the abstract FOC into a multiplier rule is the Factorisation Lemma with a “signed” extension known as Farkas’ Lemma. For multi-objective optimisation,⁴⁹ there are “multi-vector” extensions of Farkas’ Lemma, known as Tucker’s and Motzkin’s Lemmas. The latter is used here for a different purpose, viz., to show the equivalence of two forms of the Mangasarian-Fromovitz Constraint Qualification (Lemma 14).

The Factorisation Lemma can be stated as a criterion for a nonhomogeneous system of linear equations to have a solution—i.e., for a given vector p to be a linear combination of, say, the rows of a given matrix B (when the system’s variables μ are arranged in a row to the left of B). For the system $p^T = \mu^T B$ to be soluble for μ , it is obviously necessary that the system $Bv = 0$ and $p^T v \neq 0$ have no solution for v . The point is that this condition is also sufficient. Thus the lemma is a *theorem of the alternative*—stating that always either one system or the other has a solution (but obviously never both).

Lemma 33 (Factorisation) *Given a vector $p \in \mathbb{R}^n$ and an $m \times n$ (real) matrix B , exactly one of the following two systems has a solution: either*

$$p^T = \mu^T B \quad \text{for some } \mu \tag{54}$$

or

$$\begin{cases} p^T v \neq 0 \\ Bv = 0 \end{cases} \quad \text{for some } v \tag{55}$$

(but not both (54) and (55)).

In other words, with $B_{e\bullet}$ denoting the e -th row of B , there exists a $\mu \in \mathbb{R}^m$ such that

$$p^T = \sum_{e=1}^m \mu_e B_{e\bullet} \tag{56}$$

if (and only if)

$$\ker(B\cdot) \subseteq \ker(p^T\cdot). \tag{57}$$

Comments:

1. In (57), the data are viewed as linear operations, $v \mapsto Bv$ and $v \mapsto p^T v$. The inclusion between their kernels is not only necessary but also sufficient for (56) to be met by some μ . (In the Proof of Theorem 20), this is applied to the maximand’s gradient ∇f as p and the constraints’ Jacobian Dh as B .)

⁴⁹See, e.g., [10, (7.5.14)].

2. In the language of linear operations, whose composition corresponds to matrix multiplication, Lemma 33 means that, given a linear functional $p: \mathbb{R}^n \rightarrow \mathbb{R}$ and a linear map $B: \mathbb{R}^n \rightarrow \mathbb{R}^m$, the inclusion $\ker(B) \subseteq \ker(p)$ is sufficient (and obviously necessary) for p to factorise into the composition, $\mu \circ B$, of B and some linear functional $\mu: \mathbb{R}^m \rightarrow \mathbb{R}$ (hence the lemma's name). The lemma extends to linear operations between any linear spaces; in particular, p may be vector-valued.⁵⁰

Farkas' Lemma gives a similar criterion for a (nonhomogeneous) system of linear equations to have a *nonnegative* solution (for λ)—i.e., for a given vector p to be a linear combination with nonnegative coefficients of a given set of vectors, say the rows of a given matrix A . For the system $p^T = \lambda^T A$ to be soluble for λ , it is obviously necessary that the system of homogeneous linear inequalities $Av \leq 0$ and $p^T v > 0$ have no solution for v . The point is that this condition is also sufficient. Lemma 33 can be deduced from this by rewriting the equality $Bv = 0$ as a pair of opposite inequalities ($Bv \leq 0$ and $Bv \geq 0$). Indeed, as stated next, Farkas' Lemma contains the Factorisation Lemma.

Lemma 34 (Farkas' Alternative) *Given a vector $p \in \mathbb{R}^n$, an $m \times n$ matrix B and an $l \times n$ matrix A (both real), exactly one of the following two systems has a solution: either*

$$p^T = \mu^T B + \lambda^T A \quad \text{for some } \mu \text{ and } \lambda \geq 0 \quad (58)$$

or

$$\begin{cases} p^T v > 0 \\ Bv = 0 \\ Av \leq 0 \end{cases} \quad \text{for some } v \quad (59)$$

(but not both (58) and (59)).⁵¹

In other words, with $A_{i\bullet}$ denoting the i -th row of A , there exists a $\mu \in \mathbb{R}^m$ and a nonnegative $\lambda \in \mathbb{R}_+^l$ such that

$$p^T = \sum_{e=1}^m \mu_e B_{e\bullet} + \sum_{i=1}^l \lambda_i A_{i\bullet} \quad (60)$$

if (and only if) for every $v \in \mathbb{R}^n$

$$(Bv = 0 \text{ and } Av \leq 0) \Rightarrow p \cdot v \leq 0. \quad (61)$$

First Proof. To start with, one can assume that $m = 0$, i.e., one can omit B (replacing $\mu^T B$ and Bv by zeros in (58) and (59)). Since both systems cannot be

⁵⁰Also, the Factorisation Lemma has its counterparts for other algebraic structures—groups, rings, etc. It is also known as the Homomorphism Theorem or Sard's Quotient Theorem.

⁵¹This is in [5, Lemma 1.5.3], [6, Theorem 4.3.4], [9, pp. 11, 68–69, 115, 141] and [10, 5.2.3].

soluble (simultaneously), it suffices to choose either system and show that its insolubility implies solubility of the other. With either choice, this can be done by separation with a hyperplane. If (59) is chosen to be insoluble, suitable subsets of \mathbb{R}^{1+l} are separated. If (58) is chosen to be insoluble then p can be separated, by a $v \in \mathbb{R}^n$, from the range of $\cdot A$ on \mathbb{R}_+^l . The latter argument, detailed next, is slightly simpler.

By assumption, p does not lie in $\mathbb{R}_+^l A$, the image of \mathbb{R}_+^l under the linear map $\lambda \mapsto \lambda^T A$. Since this is a finitely generated convex cone,⁵² it is a closed set: see, e.g., [5, Lemma 1.5.5] or [6, Theorem 4.3.2].⁵³ Therefore, p can be separated strongly from $\mathbb{R}_+^l A$ by a hyperplane, i.e., there exists a $v \in \mathbb{R}^n$ and a scalar z such that

$$p^T v > z \geq \lambda^T A v \quad (62)$$

for every $\lambda \geq 0$. It follows that $z \geq 0$ (by setting $\lambda = 0$). Furthermore, z can be chosen to be 0 (since if $\lambda^T A v$ were positive, it could be made arbitrarily large by scaling λ up, and so it could not be bounded from above). So $\lambda^T A v \leq 0$ for every $\lambda \geq 0$, i.e., $A v \leq 0$. And $p^T v > 0$, as required.

The case of $m > 0$ is reduced to the case of $m = 0$ (with $2m + l$ instead of l) by rewriting $Bv = 0$ as $Bv \leq 0$ and $-Bv \leq 0$. This is because the existence of $\lambda' \geq 0$, $\lambda'' \geq 0$ and $\lambda \geq 0$ such that $p = (\lambda' - \lambda'')^T B + \lambda^T A$ is equivalent to (58), by setting $\mu = \lambda' - \lambda''$. ■

Comments:

1. Farkas' Lemma contains the Factorisation Lemma because if a linear functional has a semidefinite sign on a linear space then it actually vanishes on it. Here, this means that the condition $p \cdot v \leq 0$ for every $v \in \ker(B \cdot)$ is actually equivalent to the apparently stronger condition $p \cdot v = 0$ for $v \in \ker(B \cdot)$.

Second Proof of Lemma 34. This time, it is (59) that is assumed to be insoluble, and solubility of (58) is to be deduced. As before, one can set $m = 0$.

By assumption, the image of \mathbb{R}^n under the linear map $v \mapsto (p^T v, A v) \in \mathbb{R}^{1+l}$ is disjoint from the set

$$\mathbb{R}_{++} \times \mathbb{R}_-^l := \{(w, -u) : w > 0, u \geq 0\}.$$

This means that the image space $(p^T, A) \mathbb{R}^n$ is disjoint from the closed orthant $(1, 0) + (\mathbb{R}_+ \times \mathbb{R}_-^l)$. Equivalently

$$(1; 0, \dots, 0) \notin (p^T, A) \mathbb{R}^n + (\mathbb{R}_- \times \mathbb{R}_+^l).$$

⁵²It is the convex hull of the cone generated by $\{I_i \cdot A : i = 1, \dots, l\}$, the images under $\cdot A$ of the coordinate unit vectors $\begin{bmatrix} 0 & \dots & 1 & 0 & \dots \end{bmatrix}$.

⁵³What is more, a finitely generated convex cone is the same as a polyhedral cone (intersection of a finite number of half-spaces): see, e.g., [9, 4.6.1 and 4.6.2].

Being the (algebraic) sum of a linear space and an orthant, this set is a finitely generated convex cone, so it is closed: see, e.g., [5, Lemma 1.5.5] or [6, Theorem 4.3.2]. Therefore it can be separated strongly from $(1, 0)$ by a hyperplane perpendicular to a vector $(-\nu, \lambda)$, i.e., there exists a scalar ν , a $\lambda \in \mathbb{R}^l$ and a scalar z such that

$$-\nu < z \leq (-\nu p^T + \lambda^T A) v + (\nu w + \lambda \cdot u) \quad (63)$$

for every $v, u \geq 0$ and (scalar) $w \geq 0$. It follows that $z \leq 0$ (by setting $v = 0, w = 0$ and $u = 0$). Furthermore, z can be chosen to be 0 (since if the r.h.s. of the second inequality in (63) were negative, it could be made arbitrarily large in absolute value by scaling up v, w and u). It follows that $\lambda \geq 0$ (since $\lambda \cdot u \geq 0$ for every $u \geq 0$, by setting $v = 0$ and $w = 0$), and similarly $\nu \geq 0$ (by setting $v = 0$ and $u = 0$). Actually $\nu > 0$ (by the first inequality in (63)) and so, by scaling (ν, λ) , one can set $\nu = 1$. Then, finally,

$$p^T = \lambda^T A$$

(because if not, then the term $(-p^T + \lambda^T A) v$ in (63) could be made negative and arbitrarily large in absolute value by a choice of v). ■

Lemma 35 (Motzkin’s Alternative) *Given a $q \times n$ matrix C , an $m \times n$ matrix B and an $l \times n$ matrix A (all real), exactly one of the following two systems has a solution: either*

$$0 = \nu^T C + \mu^T B + \lambda^T A \quad \text{for some } \nu > 0 \text{ and } \mu \text{ and } \lambda \geq 0 \quad (64)$$

or⁵⁴

$$\begin{cases} Cv \ll 0 \\ Bv = 0 \\ Av \leq 0 \end{cases} \quad \text{for some } v \quad (65)$$

(but not both (64) and (65)).⁵⁵

First Proof. This can be proved like Farkas’ Lemma (which it contains): as in the Proof of Lemma 34, it suffices to choose either system and show that its insolubility implies solubility of the other; and one can assume that $m = 0$ (replacing $\mu^T B$ and Bv by zeros in (64) and (65)).

⁵⁴The symbols $<$ and \ll denote semistrict and strict vector inequalities, i.e., $<$ means “ \leq but \neq ”, whilst \ll means strict inequality for each pair of entries. Also, it is usually assumed that the matrix C is nonempty, i.e., that $q \geq 1$ (in addition to $n \geq 1$). Formally, this is unnecessary because when $q = 0$, there is no $\nu > 0$ (so (64) is insoluble, whilst (65) has $v = 0$ as a solution).

⁵⁵This is in [9, p. 135] and [10, 7.5.3]. In [10], Proof 3 needs a correction: instead of being scaled to have 1 as one of its entries, the y_1^* in [10, (7.5.6)] should be scaled to lie in a compact base for \mathbb{R}_+^q —such as the unit simplex in (66) here. This will make the set H in [10, (7.5.6)] convex and compact. As it stands, H has neither property, but both are needed for the strong separation argument.

Suppose that (64) is insoluble, i.e., the ranges of $\cdot C$ on $-(\mathbb{R}_+^q \setminus \{0\})$ and of $\cdot A$ on \mathbb{R}_+^l are disjoint. Equivalently, $\mathbb{R}_+^l A$ is disjoint from $-S_1^q C$, where

$$S_1^q := \left\{ \nu \in \mathbb{R}^q : \nu \geq 0, \sum_{j=1}^n \nu_j = 1 \right\} \quad (66)$$

is the unit simplex in \mathbb{R}_+^q (i.e., the image of \mathbb{R}_+^l under the linear map $\lambda \mapsto \lambda^T A$ is disjoint from the image of $-S_1^q$ under the linear map $\nu \mapsto \nu^T C$). Since S_1^q is compact and convex, so is its linear image; and therefore $-S_1^q C$ can be strongly separated, by a hyperplane, from the closed convex cone $\mathbb{R}_+^l A$. In other words, there exists a $v \in \mathbb{R}^n$ and a scalar z such that

$$-\nu^T C v > z \geq \lambda^T A v \quad (67)$$

for every $\nu \in S_1^q$ and $\lambda \geq 0$. It follows that $z \geq 0$ (by setting $\lambda = 0$). Furthermore, z can be chosen to be 0 (since if $\lambda^T A v$ were positive, it could be made arbitrarily large by scaling λ up). So $\lambda^T A v \leq 0$ for every $\lambda \geq 0$, i.e., $A v \leq 0$.

Similarly $C v \ll 0$ (i.e., $C v$ is strictly negative), since $\nu^T C v < 0$ for every $\nu \in S_1^q$ (or, equivalently, for every $\nu > 0$, i.e., for semipositive v). ■

Comments:

1. Like Farkas' Lemma, Motzkin's Lemma can also be proved by separation in the other space, viz., \mathbb{R}^{q+l} .

Second Proof of Lemma 35. This time, it is (65) that is assumed to be insoluble, and solubility of (64) is to be deduced. As before, one can set $m = 0$.

By assumption, the image of \mathbb{R}^n under the linear map $v \mapsto (Cv, Av) \in \mathbb{R}^{q+l}$ is disjoint from the set

$$\mathbb{R}_{--}^q \times \mathbb{R}_-^l := -\{(w, u) : w \gg 0, u \geq 0\}.$$

This means that the image space $(C, A)\mathbb{R}^n$ is disjoint from the closed orthant $(-1; 0) + (\mathbb{R}_+^q \times \mathbb{R}_-^l)$. Equivalently

$$(-1, \dots, -1; 0, \dots, 0) \notin (C, A)\mathbb{R}^n + (\mathbb{R}_+ \times \mathbb{R}_+^l).$$

Therefore this closed convex cone can be separated strongly from $(-1; 0)$ by a hyperplane (in $\mathbb{R} \times \mathbb{R}^l$), i.e., there exists a $\nu \in \mathbb{R}^q$, a $\lambda \in \mathbb{R}^l$ and a scalar z such that

$$-\sum_{k=1}^q \nu_k < z \leq (\nu^T C + \lambda^T A) v + (\nu \cdot w + \lambda \cdot u) \quad (68)$$

for every $v, u \geq 0$ and $w \geq 0$. It follows that $z \leq 0$ (by setting $v = 0, w = 0$ and $u = 0$). Furthermore, z can be chosen to be 0 (since if the r.h.s. of the second

inequality in (68) were negative, it could be made arbitrarily large in absolute value by scaling up v , w and u). So $\lambda \geq 0$ (by setting $v = 0$ and $w = 0$), and similarly $\nu \geq 0$ (by setting $v = 0$ and $u = 0$). And $\nu \neq 0$ (by the first inequality in (68)); so $\nu > 0$. Finally,

$$0 = \nu^T C + \lambda^T A$$

(because if not, then the term $(\nu^T C + \lambda^T A)v$ in (68) could be made negative and arbitrarily large in absolute value by a choice of v). ■

2. The proof by separation in $\mathbb{R}^{q+l} = \mathbb{R}^q \times \mathbb{R}^l$ (with $m = 0$) can also be split into two stages, as in [10, p. 170, Proof 2]: first, a separation argument in \mathbb{R}^q produces a ν , which is then be used to “scalarise” one dimension of C , i.e., to reduce Motzkin’s Lemma to Farkas’ Lemma with $p^T = -\nu^T C$.
3. Another theorem of the alternative, Tucker’s, is obtained “by swapping the strictness and semi-strictness” in Motzkin’s, i.e., by having $\nu \gg 0$ in (64) and $Cv < 0$ in (65)—which together imply that $\nu^T Cv < 0$, as before. See, e.g., [9] or [10, 7.5.7]. The two, Motzkin’s and Tucker’s, are combined in what is known as Slater’s Alternative (which contains all those given here) : see, e.g., [9] or [10, 7.5.11].

A special cases of Motzkin’s Lemma gives a criterion for a homogeneous system of linear equations to have a semipositive solution.

Lemma 36 (Gordan’s Alternative) *Given a $q \times n$ (real) matrix C , exactly one of the following two systems has a solution: either*

$$0 = \nu^T C \quad \text{for some } \nu > 0 \tag{69}$$

or

$$Cv \ll 0 \quad \text{for some } v \tag{70}$$

(but not both).⁵⁶

Proof. Set $l = 0$ and $m = 0$ in Lemma 35 (i.e., take both A and B to be empty). ■

Definition 37 *A set of vectors is positively independent if none of its semipositive linear combinations equals zero (i.e., if a nonnegative combination vanishes, then all its coefficients are zeros).*

By Gordan’s Lemma, a finite set of vectors $p(1), p(2), \dots, p(q)$ is positively independent if and only if there exists a v with $p(i) \cdot v < 0$ for each i .⁵⁷ Geometrically, a

⁵⁶This is in [9, pp. 71 and 137] and [10, pp. 171–172].

⁵⁷The rows of a matrix C are positively independent if and only if system (69) is insoluble for ν , i.e., if and only if system (70) has a solution v .

finite set of nonzero vectors is positively independent if and only if the convex cone they generate is pointed (i.e., line-free): see, e.g., [6, Exercise 4.3.13]. And a linear functional has a negative scalar product with each of a finite number of nonzero vectors if and only if it lies in the interior of the polar to the convex cone they generate. So, in geometric terms, Gordan’s Lemma means that a finitely generated convex cone is pointed if and only if its polar is solid (i.e., has a nonempty interior).

Another special case of Motzkin’s Lemma gives a similar criterion for a homogeneous system of linear equations to have a strictly positive solution.

Lemma 38 (Stiemke’s Alternative) *Given a $m \times n$ (real) matrix B , exactly one of the following two systems has a solution: either*

$$0 < \mu^T B \quad \text{for some } \mu \tag{71}$$

or

$$\begin{cases} v \gg 0 \\ Bv = 0 \end{cases} \quad \text{for some } v \tag{72}$$

(but not both).⁵⁸

Proof. Set $q = n$ and $C = -I(n)$ and an empty A (i.e., $l = 0$) in Lemma 35. ■

References

- [1] Bigelow, J. H., and N. E. Shapiro (1974): “Implicit function theorems for mathematical programming and for systems of inequalities”, *Mathematical Programming*, 6, 141–156.
- [2] Debreu, G. (1952): “Definite and semidefinite quadratic forms”, *Econometrica*, 20, 295–300.
- [3] El-Hodiri, M. A. (1971): *Constrained extrema: introduction to the differentiable case with economic applications*. (Lecture Notes in Operations Research and Mathematical Systems, vol. 56). Berlin-Heidelberg-New York: Springer.
- [4] Fiacco, A. V. (1983): *Introduction to sensitivity and stability analysis in nonlinear programming*. New York: Academic Press.
- [5] Hestenes, M. R. (1966): *Calculus of variations and optimal control theory*. New York-London-Sydney: Wiley.

⁵⁸This is in [6, Theorem 4.3.5] and, as a case of Tucker’s Lemma, in [9, pp. 138] and [10, p. 172].

- [6] Hestenes, M. R. (1975): *Optimization theory: the finite-dimensional case*. New York-London-Sydney: Wiley.
- [7] Horsley, A., and A. J. Wrobel (2003): “Liminal inequality constraints and second-order optimality conditions”, CDAM Research Report LSE-CDAM-2003-16.
- [8] Jittorntrum, A. V. (1984): “Solution point differentiability without strict complementarity in nonlinear programming”, *Mathematical Programming Study*, 21, 127–138.
- [9] Panik, M. J. (1993): *Fundamentals of convex analysis*. Dordrecht-Boston-London: Kluwer.
- [10] Ponstein, J. (1980): *Approaches to the theory of optimisation*. Cambridge-New York-Melbourne: Cambridge University Press.
- [11] Silberberg, E. (1971): “The Le Chatelier Principle as a corollary to a generalized envelope theorem”, *Journal of Economic Theory*, 3, 146–155.
- [12] Simon, C., and L. Blume (1994): *Mathematics for economists*. New York-London: Norton.
- [13] Takayama, A. (1985): *Mathematical economics*. Cambridge-London-New York: Cambridge University Press.