

An introduction to algorithms for nonlinear optimization^{1,2}

Nicholas I. M. Gould^{3,4,5} and Sven Leyffer^{6,7,8}

ABSTRACT

We provide a concise introduction to modern methods for solving nonlinear optimization problems. We consider both linesearch and trust-region methods for unconstrained minimization, interior-point methods for problems involving inequality constraints, and SQP methods for those involving equality constraints. Theoretical as well as practical aspects are emphasised. We conclude by giving a personal view of some of the most significant papers in the area, and a brief guide to on-line resources.

¹ These notes were written to accompany a course on optimization we gave at the 10th EPSRC/LMS Numerical Analysis Summer School at the University of Durham, 15-19th July 2002. They are heavily influenced by a similar course given by the first author to M.Sc. students at Oxford University in Trinity Term, 2001.

² Some of this material appeared, in expanded form, in the book “Trust-region methods” [SIAM, Philadelphia, 2000] by Andy Conn, Philippe Toint and the first author. We are grateful to both Andy and Philippe for their contributions.

³ Computational Science and Engineering Department, Rutherford Appleton Laboratory, Chilton, Oxfordshire, OX11 0QX, England, EU. Email: n.gould@rl.ac.uk

⁴ Current reports available from “<http://www.numerical.rl.ac.uk/reports/reports.html>”.

⁵ This work was supported in part by the EPSRC grant GR/R46641

⁶ Department of Mathematics, University of Dundee, DD1 4HN, Scotland.
Email: sleyffer@maths.dundee.ac.uk

⁷ Address from 1st September 2002: Mathematics and Computer Science Division,
Argonne National Laboratory, 9700 S. Cass Avenue, Argonne, IL 60439, USA.
Email: leyffer@mcs.anl.gov

⁸ Current reports available from “<http://www-unix.mcs.anl.gov/~leyffer/>”.

Computational Science and Engineering Department
Atlas Centre
Rutherford Appleton Laboratory
Oxfordshire OX11 0QX
December 17, 2002.

Contents

INTRODUCTION	1
1 OPTIMALITY CONDITIONS AND WHY THEY ARE IMPORTANT	2
1.1 Optimization problems	2
1.2 Notation	2
1.3 Lipschitz continuity and Taylor's theorem	3
1.4 Optimality conditions	4
1.5 Optimality conditions for unconstrained minimization	5
1.6 Optimality conditions for constrained minimization	5
2 LINESEARCH METHODS FOR UNCONSTRAINED OPTIMIZATION	8
2.1 Linesearch methods	8
2.2 Practical linesearch methods	9
2.3 Convergence of generic linesearch methods	12
2.4 Method of steepest descent	13
2.5 More general descent methods	15
3 TRUST-REGION METHODS FOR UNCONSTRAINED OPTIMIZATION	20
3.1 Linesearch vs. trust-region methods	20
3.2 Trust-region models	20
3.3 Basic trust-region method	21
3.4 Basic convergence of trust-region methods	23
3.5 Solving the trust-region subproblem	26
3.6 Solving the large-scale problem	30
4 INTERIOR-POINT METHODS FOR INEQUALITY CONSTRAINED OPTIMIZATION	32
4.1 Merit functions for constrained minimization	33
4.2 The logarithmic barrier function for inequality constraints	33
4.3 A basic barrier-function algorithm	34
4.4 Potential difficulties	35
4.5 A different perspective: perturbed optimality conditions	37
4.6 A practical primal-dual method	40
5 SQP METHODS FOR EQUALITY CONSTRAINED OPTIMIZATION	41
5.1 Newton's method for first-order optimality	41
5.2 The Sequential Quadratic Programming iteration	42
5.3 Linesearch SQP methods	44
5.4 Trust-region SQP methods	47
CONCLUSION	53
APPENDIX A - SEMINAL BOOKS AND PAPERS	53
APPENDIX B - OPTIMIZATION RESOURCES ON THE WORLD-WIDE-WEB	60
APPENDIX C - SKETCHES OF PROOFS	64

INTRODUCTION

The solution of nonlinear optimization problems—that is the minimization or maximization of an objective function involving unknown parameters/variables in which the variables may be restricted by constraints—is one of the core components of computational mathematics. Nature (and man) loves to optimize, and the world is far from linear. In his book on Applied Mathematics, the eminent mathematician Gill Strang opines that optimization, along with the solution of systems of linear equations, and of (ordinary and partial) differential equations, is one of the three cornerstones of modern applied mathematics. It is strange, then, that despite fathering many of the pioneers in the area, the United Kingdom (in particular, higher education) has turned away from the subject in favour of Strang’s other key areas. Witness the numbers of lectures in the EPSRC summer school over the past ten years (broadly 3 series of lectures on numerical linear algebra, 4 on ODEs, 17 on PDEs, 4 in other areas, but *none* in optimization in Summer Schools V-XI) and the relative paucity of undergraduate or postgraduate courses in the area.

It is timely, therefore, to be given the opportunity to be able to review developments in nonlinear optimization. The past 10 years have shown an incredible growth in the power and applicability of optimization techniques, fueled in part by the “interior-point revolution” of the late 1980s. We have purposely chosen not to consider linear or discrete problems, nor to describe important applications such as optimal control, partially because there are other experts better able to review these fields, and partly because they would all make exciting courses on their own. Indeed, the prospect of simply describing nonlinear optimization in five lectures is extremely daunting, and each of the subjects we shall describe could easily fill the whole course! Of course, there is a strong cross-fertilisation of ideas from discrete to linear to nonlinear optimization, just as there are strong influences both to and from other branches of numerical analysis and computational mathematics.

This article is partitioned in broadly the same way as the course on which it is based. Optimality conditions play a vital role in optimization, both in the identification of optima, and in the design of algorithms to find them. We consider these in Section 1. Sections 2 and 3 are concerned with the two main techniques for solving unconstrained optimization problems. Although it can be argued that such problems arise relatively infrequently in practice (nonlinear fitting being a vital exception), the underlying linesearch and trust-region ideas are so important that it is best to understand them first in their simplest setting. The remaining two sections cover the problems we really wish to solve, those involving constraints. We purposely consider inequality constraints (alone) in one and equality constraints (alone) in the other, since then the key ideas may be developed without the complication of treating both kinds of constraints at once. Of course, real methods cope with both, and suitable algorithms will be hybrids of the methods we have considered.

We make no apologies for mixing theory in with algorithms, since (most) good algorithms have good theoretical underpinnings. So as not to disturb the development in the main text, the proofs of stated theorems have been relegated to Appendix C. In addition, we do not provide citations in the main text, but have devoted Appendix A to an annotated bibliography of what we consider to be essential references in nonlinear optimization. Such a list is, by its nature, selective, but we believe that the given references form a corpus of seminal work in the area, which should be read by any student interested in pursuing a career in optimization.

Before we start, we feel that one key development during the last five years has done more to promote the use of optimization than possibly any other. This is NEOS, the Network Enabled Optimization Server, at Argonne National Laboratory and Northwestern University in Chicago, see <http://www-neos.mcs.anl.gov/neos>. Here, users are able to submit problems for remote

solution, without charge, by a large (and expanding) collection of the world's best optimization solvers, many of them being only available otherwise commercially. Further details of what may be found on the World-Wide-Web are given in Appendix B.

1 OPTIMALITY CONDITIONS AND WHY THEY ARE IMPORTANT

1.1 Optimization problems

As we have said optimization is concerned with the minimization or maximization of an objective function, say, $f(x)$. Since

$$\text{maximum } f(x) = - \text{minimum } (-f(x))$$

there is no loss in generality in concentrating in this article on minimization—throughout, minimization will take place with respect to an n -vector, x , of real unknowns. A bit of terminology here: the smallest value of f gives its *minimum*, while any (there may be more than one) corresponding values of x are a *minimizer*.

There are a number of important subclasses of optimization problems. The simplest is *unconstrained minimization*, where we aim to

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

where the *objective function* $f: \mathbb{R}^n \rightarrow \mathbb{R}$. One level up is *equality constrained minimization*, where now we try to

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) \quad \text{subject to} \quad c(x) = 0$$

where the *constraints* $c: \mathbb{R}^n \rightarrow \mathbb{R}^m$. For consistency we shall assume that $m \leq n$, for otherwise it is unlikely (but not impossible) that there is an x that satisfies all of the equality constraints. Another important problem is *inequality constrained minimization*, in which we aim to

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) \quad \text{subject to} \quad c(x) \geq 0$$

where $c: \mathbb{R}^n \rightarrow \mathbb{R}^m$ and now m may be larger than n . The most general problem involves both equality and inequality constraints—some inequalities may have upper as well as lower bounds—and may be further sub-classified depending on the nature of the constraints. For instance, some of the $c_i(x)$ may be linear (that is $c_i(x) = a_i^T x - b_i$ for some vector a_i and scalar b_i), some may be simple bounds on individual components of x (for example, $c_i(x) = x_i$), or some may result from a network (“flow in = flow out”).

1.2 Notation

It is convenient to introduce our most common notation and terminology at the outset. Suppose that $f(x)$ is at least twice continuously differentiable ($f \in \mathcal{C}^2$). We let $\nabla_x f(x)$ denote the vector of first partial derivatives, whose i -th component is $\partial f(x)/\partial x_i$. Similarly, the i, j -th component of the (symmetric) matrix $\nabla_{xx} f(x)$ is the second partial derivative $\partial^2 f(x)/\partial x_i \partial x_j$. We also write the usual Euclidean inner product between two p -vectors u and v as $\langle u, v \rangle \stackrel{\text{def}}{=} \sum_{i=1}^p u_i v_i$ (and mention, for those who care, that some but not all of what we have to say remains true in more general

Hilbert spaces!). We denote the set of points for which all the constraints are satisfied as \mathcal{C} , and say that any $x \in \mathcal{C}$ (resp. $x \notin \mathcal{C}$) is *feasible* (resp. *infeasible*).

With this in mind we define the *gradient* and *Hessian* (matrix) of the objective function f to be $g(x) \stackrel{\text{def}}{=} \nabla_x f(x)$ and $H(x) \stackrel{\text{def}}{=} \nabla_{xx} f(x)$, respectively. Likewise, the gradient and Hessian of the i -th constraint are $a_i(x) \stackrel{\text{def}}{=} \nabla_x c_i(x)$ and $H_i(x) \stackrel{\text{def}}{=} \nabla_{xx} c_i(x)$. The *Jacobian* (matrix) is

$$A(x) \stackrel{\text{def}}{=} \nabla_x c(x) \equiv \begin{pmatrix} a_1^T(x) \\ \vdots \\ a_m^T(x) \end{pmatrix}.$$

Finally, if y is a vector (of so-called *Lagrange multipliers*), the *Lagrangian* (function) is

$$\ell(x, y) \stackrel{\text{def}}{=} f(x) - \langle y, c(x) \rangle,$$

while its gradient and Hessian with respect to x are, respectively,

$$\begin{aligned} g(x, y) &\stackrel{\text{def}}{=} \nabla_x \ell(x, y) \equiv g(x) - \sum_{i=1}^m y_i a_i(x) \equiv g(x) - A^T(x)y \quad \text{and} \\ H(x, y) &\stackrel{\text{def}}{=} \nabla_{xx} \ell(x, y) \equiv H(x) - \sum_{i=1}^m y_i H_i(x). \end{aligned}$$

One last piece of notation: e_i is the i -th unit vector, while e is the vector of ones, and I is the (appropriately dimensioned) identity matrix.

1.3 Lipschitz continuity and Taylor's theorem

It might be argued that those who understand Taylor's theorem and have a basic grasp of linear algebra have all the tools they need to study continuous optimization—of course, this leaves aside all the beautiful mathematics needed to fully appreciate optimization in abstract settings, yet another future EPSRC summer school course, we hope!

Taylor's theorem(s) can most easily be stated for functions with Lipschitz continuous derivatives. Let \mathcal{X} and \mathcal{Y} open sets, let $F : \mathcal{X} \rightarrow \mathcal{Y}$, and let $\|\cdot\|_{\mathcal{X}}$ and $\|\cdot\|_{\mathcal{Y}}$ be norms on \mathcal{X} and \mathcal{Y} respectively. Then F is *Lipschitz continuous at* $x \in \mathcal{X}$ if there exists a function $\gamma(x)$ such that

$$\|F(z) - F(x)\|_{\mathcal{Y}} \leq \gamma(x)\|z - x\|_{\mathcal{X}}$$

for all $z \in \mathcal{X}$. Moreover F is *Lipschitz continuous throughout/in* \mathcal{X} if there exists a constant γ such that

$$\|F(z) - F(x)\|_{\mathcal{Y}} \leq \gamma\|z - x\|_{\mathcal{X}}$$

for all x and $z \in \mathcal{X}$. Lipschitz continuity relates (either locally or globally) the changes that occur in F to those that are permitted in x .

Armed with this, we have the following *Taylor* approximation results. The first suggests how good (or bad) a first-order (linear) or second-order (quadratic) Taylor series approximation to a scalar-valued function may be.

Theorem 1.1. Let \mathcal{S} be an open subset of \mathbb{R}^n , and suppose $f : \mathcal{S} \rightarrow \mathbb{R}$ is continuously differentiable throughout \mathcal{S} . Suppose further that $g(x)$ is Lipschitz continuous at x , with Lipschitz constant $\gamma^L(x)$ in some appropriate vector norm. Then, if the segment $x + \theta s \in \mathcal{S}$ for all $\theta \in [0, 1]$,

$$|f(x + s) - m^L(x + s)| \leq \frac{1}{2} \gamma^L(x) \|s\|^2, \text{ where}$$

$$m^L(x + s) = f(x) + \langle g(x), s \rangle.$$

If f is twice continuously differentiable throughout \mathcal{S} and $H(x)$ is Lipschitz continuous at x , with Lipschitz constant $\gamma^Q(x)$,

$$|f(x + s) - m^Q(x + s)| \leq \frac{1}{6} \gamma^Q(x) \|s\|^3, \text{ where}$$

$$m^Q(x + s) = f(x) + \langle g(x), s \rangle + \frac{1}{2} \langle s, H(x)s \rangle.$$

The second result compares how bad a first-order Taylor series approximation to a vector valued function might be.

Theorem 1.2. Let \mathcal{S} be an open subset of \mathbb{R}^n , and suppose $F : \mathcal{S} \rightarrow \mathbb{R}^m$ is continuously differentiable throughout \mathcal{S} . Suppose further that $\nabla_x F(x)$ is Lipschitz continuous at x , with Lipschitz constant $\gamma^L(x)$ in some appropriate vector norm and its induced matrix norm. Then, if the segment $x + \theta s \in \mathcal{S}$ for all $\theta \in [0, 1]$,

$$\|F(x + s) - M^L(x + s)\| \leq \frac{1}{2} \gamma^L(x) \|s\|^2,$$

where

$$M^L(x + s) = F(x) + \nabla_x F(x)s$$

1.4 Optimality conditions

Now is the time to come clean. It is very, very difficult to say anything about the solutions to the optimization problems given in Section 1.1. This is almost entirely because we are considering very general problems, for which there may be many local, often non-global, minimizers. There are two possible ways around this. We might choose to restrict the class of problems we allow, so that all local minimizers are global. But since this would rule out the vast majority of nonlinear problems that arise in practice, we instead choose to lower our sights, and only aim for local minimizers—there are methods that offer some guarantee of global optimality, but to date they are really restricted to small or very specially structured problems.

Formally, we still need to define what we mean by a local minimizer. A feasible point x_* is a *local* minimizer of $f(x)$ if there is an open neighbourhood \mathcal{N} of x_* such that $f(x_*) \leq f(x)$ for all $x \in \mathcal{C} \cap \mathcal{N}$. If there is an open neighbourhood \mathcal{N} of x_* such that $f(x_*) < f(x)$ for all $x \neq x_* \in \mathcal{C} \cap \mathcal{N}$, it is *isolated*.

While such definitions agree with our intuition, they are of very little use in themselves. What

we really need are optimality conditions. Optimality conditions are useful for three reasons. Firstly, they provide a means of guaranteeing that a candidate solution is indeed (locally) optimal—these are the so-called *sufficient conditions*. Secondly, they indicate when a point is not optimal—these are the *necessary conditions*. Finally they guide us in the design of algorithms, since lack of optimality indicates when we may improve our objective. We now give details.

1.5 Optimality conditions for unconstrained minimization

We first consider what we might deduce if we were fortunate enough to have found a local minimizer of $f(x)$. The following two results provide first- and second-order necessary optimality conditions (respectively).

Theorem 1.3. Suppose that $f \in C^1$, and that x_* is a local minimizer of $f(x)$. Then

$$g(x_*) = 0.$$

Theorem 1.4. Suppose that $f \in C^2$, and that x_* is a local minimizer of $f(x)$. Then $g(x_*) = 0$ and $H(x_*)$ is positive semi-definite, that is

$$\langle s, H(x_*)s \rangle \geq 0 \text{ for all } s \in \mathbb{R}^n.$$

But what if we have found a point that satisfies the above conditions? Is it a local minimizer? Yes, an isolated one, provided the following second-order sufficient optimality conditions are satisfied.

Theorem 1.5. Suppose that $f \in C^2$, that x_* satisfies the condition $g(x_*) = 0$, and that additionally $H(x_*)$ is positive definite, that is

$$\langle s, H(x_*)s \rangle > 0 \text{ for all } s \neq 0 \in \mathbb{R}^n.$$

Then x_* is an isolated local minimizer of f .

Notice how slim is the difference between these necessary and sufficient conditions.

1.6 Optimality conditions for constrained minimization

When constraints are present, things get more complicated. In particular, the geometry of the feasible region at (or near) to a minimizer plays a very subtle role. Consider a suspected minimizer x_* . We shall say that a constraint is *active* at x_* if and only if $c_i(x_*) = 0$. By necessity, equality constraints will be active, while determining which (if any) of the inequalities is active is probably the overriding concern in constrained optimization.

In order to say anything about optimality, it is unfortunately necessary to rule out “nasty” local minimizers such as cusps on the constraint boundary. This requires that we have to ask that so-called

constraint qualifications hold—essentially these say that linear approximations to the constraints characterize all feasible perturbations about x_* and that perturbations which keep strongly active constraints strongly active (a *strongly* active constraint is one that will still be active if the data, and hence minimizer, is slightly perturbed) are completely characterized by their corresponding linearizations being forced to be active. Fortunately, such assumptions are automatically satisfied if the constraints are linear, or if the constraints that are active have independent gradients, and may actually be guaranteed in far weaker circumstances than these.

1.6.1 Optimality conditions for equality-constrained minimization

Given constraint qualifications, first- and second-order necessary optimality conditions for problems involving equality constraints are (respectively) as follows.

Theorem 1.6. Suppose that $f, c \in C^1$, and that x_* is a local minimizer of $f(x)$ subject to $c(x) = 0$. Then, so long as a first-order constraint qualification holds, there exist a vector of Lagrange multipliers y_* such that

$$\begin{aligned} c(x_*) &= 0 && \text{(primal feasibility)} \text{ and} \\ g(x_*) - A^T(x_*)y_* &= 0 && \text{(dual feasibility)}. \end{aligned}$$

Theorem 1.7. Suppose that $f, c \in C^2$, and that x_* is a local minimizer of $f(x)$ subject to $c(x) = 0$. Then, provided that first- and second-order constraint qualifications hold, there exist a vector of Lagrange multipliers y_* such that

$$\langle s, H(x_*, y_*)s \rangle \geq 0 \text{ for all } s \in \mathcal{N} \tag{1.1}$$

where

$$\mathcal{N} = \{s \in \mathbb{R}^n \mid A(x_*)s = 0\}.$$

Notice that there are two first-order optimality requirements: primal feasibility (the constraints are satisfied), and dual feasibility (the gradient of the objective function is expressible as a linear combination of the gradients of the constraints). It is not hard to anticipate that, just as in the unconstrained case, sufficient conditions occur when the requirement (1.1) is strengthened to $\langle s, H(x_*, y_*)s \rangle > 0$ for all $s \in \mathcal{N}$.

1.6.2 Optimality conditions for inequality-constrained minimization

Finally, when the problem involves inequality constraints, it is easy to imagine that only the constraints that are active at x_* play a role—the inactive constraints play no part in defining the minimizer—and indeed this is so. First- and second-order necessary optimality conditions are (respectively) as follows.

Theorem 1.8. Suppose that $f, c \in C^1$, and that x_* is a local minimizer of $f(x)$ subject to $c(x) \geq 0$. Then, provided that a first-order constraint qualification holds, there exist a vector of Lagrange multipliers y_* such that

$$\begin{aligned} c(x_*) &\geq 0 && \text{(primal feasibility),} \\ g(x_*) - A^T(x_*)y_* &= 0 \text{ and } y_* \geq 0 && \text{(dual feasibility) and} \\ c_i(x_*)[y_*]_i &= 0 && \text{(complementary slackness).} \end{aligned} \tag{1.2}$$

Theorem 1.9. Suppose that $f, c \in C^2$, and that x_* is a local minimizer of $f(x)$ subject to $c(x) \geq 0$. Then, provided that first- and second-order constraint qualifications hold, there exist a vector of Lagrange multipliers y_* for which primal/dual feasibility and complementary slackness requirements hold as well as

$$\langle s, H(x_*, y_*)s \rangle \geq 0 \text{ for all } s \in \mathcal{N}_+$$

where

$$\mathcal{N}_+ = \left\{ s \in \mathbb{R}^n \mid \begin{array}{l} \langle s, a_i(x_*) \rangle = 0 \text{ if } c_i(x_*) = 0 \text{ \& } [y_*]_i > 0 \text{ and} \\ \langle s, a_i(x_*) \rangle \geq 0 \text{ if } c_i(x_*) = 0 \text{ \& } [y_*]_i = 0 \end{array} \right\}. \tag{1.3}$$

See how dual feasibility now imposes an extra requirement, that the Lagrange multipliers be non-negative, while as expected there is an additional (complementary slackness) assumption that inactive constraints necessarily have zero Lagrange multipliers. Also notice that \mathcal{N}_+ , the set over which the Hessian of the Lagrangian is required to be positive semi-definite, may now be the intersection of a linear manifold and a cone, a particularly unpleasant set to work with.

The by-now obvious sufficient conditions also hold:

Theorem 1.10. Suppose that $f, c \in C^2$, and that x_* and a vector of Lagrange multipliers y_* satisfy (1.2) and

$$\langle s, H(x_*, y_*)s \rangle > 0$$

for all s in the set \mathcal{N}_+ given in (1.3). Then x_* is an isolated local minimizer of $f(x)$ subject to $c(x) \geq 0$.

2 LINESEARCH METHODS FOR UNCONSTRAINED OPTIMIZATION

In this and the next sections, we shall concentrate on the unconstrained minimization problem,

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

where the objective function $f: \mathbb{R}^n \rightarrow \mathbb{R}$. We shall assume that $f \in C^1$ (sometimes C^2) with Lipschitz continuous derivatives. Often in practice this assumption is violated, but nonetheless the methods converge (perhaps by good fortune) regardless.

Despite knowing how to characterise local minimizers of our problem, in practice it is rather unusual for us to be able to provide or compute an explicit minimizer. Instead, we would normally expect to fall back on a suitable iterative process. An *iteration* is simply a procedure whereby a sequence of points

$$\{x_k\}, \quad k = 1, 2, \dots$$

is generated, starting from some initial “guess” x_0 , with the overall aim of ensuring that (a subsequence) of the $\{x_k\}$ has favourable limiting properties. These might include that any limit generated satisfies first-order or, even better, second-order necessary optimality conditions.

Notice that we will not be able to guarantee that our iteration will converge to a global minimizer unless we know that f obeys very strong conditions, nor regrettably in general that any limit point is even a local minimizer (unless by chance it happens to satisfy second-order sufficiency conditions). What we normally do try to ensure is that, at the very least, the iteration is *globally* convergent, that is that (for at least) a subsequence of iterates $\{g(x_k)\}$ converges to zero. And our hope is that such a sequence converges at a reasonably fast asymptotic rate. These two preoccupations lie at the heart of computational optimization.

For brevity, in what follows, we shall write $f_k = f(x_k)$, $g_k = g(x_k)$ and $H_k = H(x_k)$.

2.1 Linesearch methods

Generically, linesearch methods work as follows. Firstly, a *search direction* p_k is calculated from x_k . This direction is required to be a *descent direction*, i.e.,

$$\langle p_k, g_k \rangle < 0 \quad \text{if } g_k \neq 0,$$

so that, for small steps along p_k , Taylor’s theorem (Theorem 1.1) guarantees that the objective function may be reduced. Secondly, a suitable *steplength* $\alpha_k > 0$ is calculated so that

$$f(x_k + \alpha_k p_k) < f_k.$$

The computation of α_k is the *linesearch*, and may itself be an iteration. Finally, given both search direction and steplength, the iteration concludes by setting

$$x_{k+1} = x_k + \alpha_k p_k.$$

Such a scheme sounds both natural and simple. But as with most simple ideas, it needs to be refined somewhat in order to become a viable technique. What might go wrong? Firstly, consider the example in Figure 2.1.

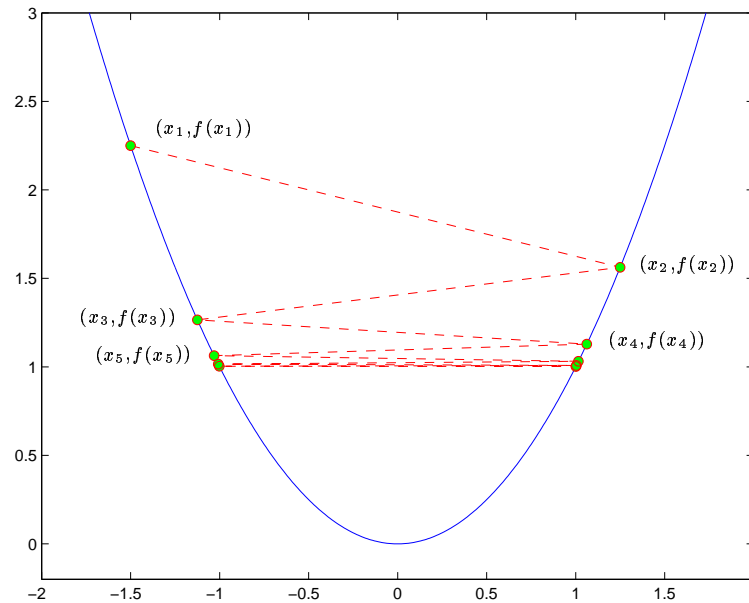


Figure 2.1: The objective function $f(x) = x^2$ and the iterates $x_{k+1} = x_k + \alpha_k p_k$ generated by the descent directions $p_k = (-1)^{k+1}$ and steps $\alpha_k = 2 + 3/2^{k+1}$ from $x_0 = 2$.

Here the search direction gives a descent direction, and the iterates oscillate from one side of the minimizer to the other. Unfortunately, the decrease per iteration is ultimately so small that the iterates converge to the pair ± 1 , neither of which is a stationary point. What has gone wrong? Simply the steps are too long relative to the amount of objective-function decrease that they provide.

Is this the only kind of failure? Unfortunately, no. For consider the example in Figure 2.2.

Now the iterates approach the minimizer from one side, but the stepsizes are so small that each iterate falls woefully short of the minimizer, and ultimately converge to the non-stationary value 1.

So now we can see that a simple-minded linesearch method can fail if the linesearch allows steps that are either too long or too short relative to the amount of decrease that might be obtained with a well-chosen step.

2.2 Practical linesearch methods

In the early days, it was often suggested that α_k should be chosen to minimize $f(x_k + \alpha p_k)$. This is known as an *exact* linesearch. In most cases, exact linesearches prove to be both very expensive—they are essentially univariate minimizations—and most definitely not cost effective, and are consequently rarely used nowadays.

Modern linesearch methods prefer to use *inexact* linesearches, which are guaranteed to pick steps that are neither too long nor too short. In addition, they aim to pick a “useful” initial “guess” for each stepsize so as to ensure fast asymptotic convergence—we will return to this when we discuss Newton’s method. The main contenders amongst the many possible inexact linesearches are the so-called “backtracking- Armijo” and the “Armijo-Goldstein” varieties. The former are extremely easy to implement, and form the backbone of most Newton-like linesearch methods. The latter are particularly important when using secant quasi-Newton methods (see Section 2.5.3), but alas we do

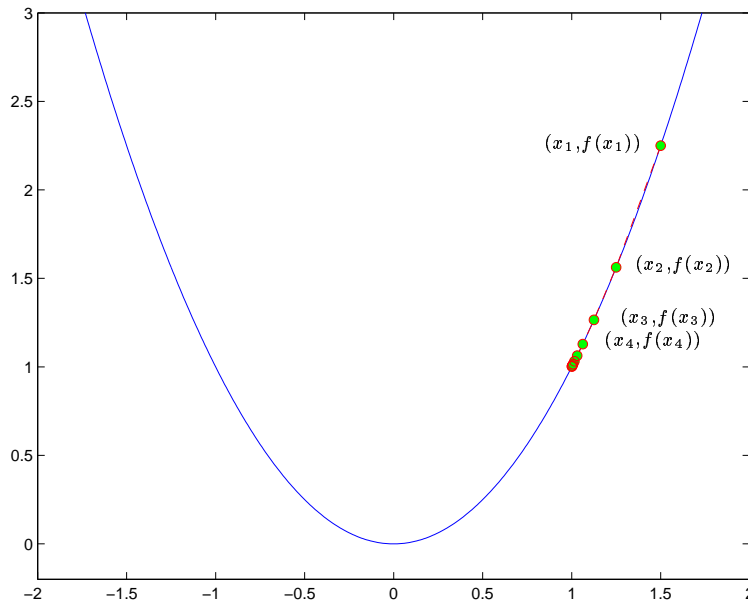


Figure 2.2: The objective function $f(x) = x^2$ and the iterates $x_{k+1} = x_k + \alpha_k p_k$ generated by the descent directions $p_k = -1$ and steps $\alpha_k = 1/2^{k+1}$ from $x_0 = 2$.

not have space to describe them here.

Here is a basic *backtracking* linesearch to find α_k :

Given $\alpha_{\text{init}} > 0$ (e.g., $\alpha_{\text{init}} = 1$),
 let $\alpha^{(0)} = \alpha_{\text{init}}$ and $l = 0$.
 Until $f(x_k + \alpha^{(l)} p_k) < f_k$
 set $\alpha^{(l+1)} = \tau \alpha^{(l)}$, where $\tau \in (0, 1)$ (e.g., $\tau = \frac{1}{2}$)
 and increase l by 1.
 Set $\alpha_k = \alpha^{(l)}$.

Notice that the backtracking strategy prevents the step from getting too small, since the first allowable value stepsize of the form $\alpha_{\text{init}} \tau^i$, $i = 0, 1, \dots$ is accepted. However, as it stands, there is still no mechanism for preventing too large steps relative to decrease in f . What is needed is a tighter requirement than simply that $f(x_k + \alpha^{(l)} p_k) < f_k$. Such a role is played by the Armijo condition.

The *Armijo condition* is that the steplength be asked to give slightly more than simply decrease in f . The actual requirement is that

$$f(x_k + \alpha_k p_k) \leq f(x_k) + \alpha_k \beta \langle p_k, g_k \rangle$$

for some $\beta \in (0, 1)$ (e.g., $\beta = 0.1$ or even $\beta = 0.0001$)—this requirement is often said to give *sufficient decrease*. Observe that, since $\langle p_k, g_k \rangle < 0$, the longer the step, the larger the required decrease in

f. The range of permitted values for the stepsize is illustrated in Figure 2.3.

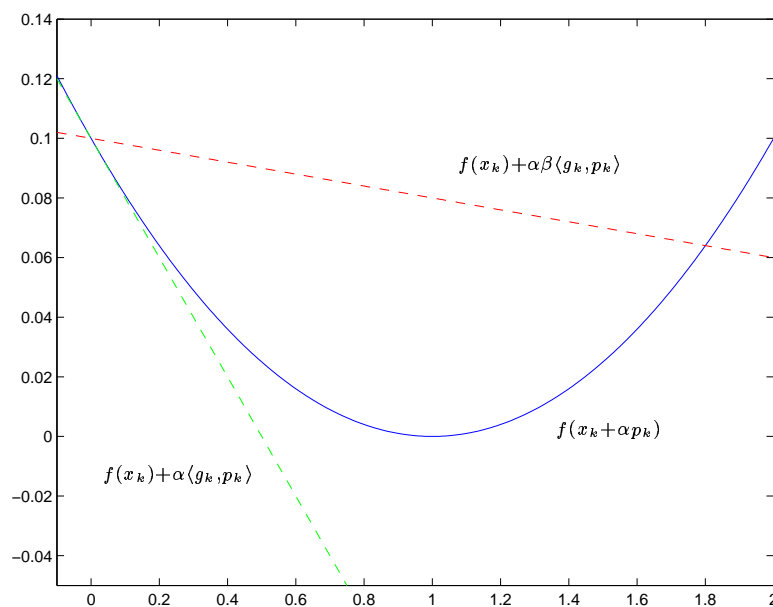


Figure 2.3: A steplength of anything up to 1.8 is permitted for this example, in the case where $\beta = 0.2$.

The Armijo condition may then be inserted into our previous backtracking scheme to give the aptly-named *Backtracking-Armijo* linesearch:

Given $\alpha_{\text{init}} > 0$ (e.g., $\alpha_{\text{init}} = 1$),
 let $\alpha^{(0)} = \alpha_{\text{init}}$ and $l = 0$.
 Until $f(x_k + \alpha^{(l)} p_k) \leq f(x_k) + \alpha^{(l)} \beta \langle p_k, g_k \rangle$
 set $\alpha^{(l+1)} = \tau \alpha^{(l)}$, where $\tau \in (0, 1)$ (e.g., $\tau = \frac{1}{2}$)
 and increase l by 1.
 Set $\alpha_k = \alpha^{(l)}$.

Of course, it is one thing to provide likely-sounding rules to control stepsize selection, but another to be sure that they have the desired effect. Indeed, can we even be sure that there are points which satisfy the Armijo condition? Yes, for we have

Theorem 2.1. Suppose that $f \in C^1$, that $g(x)$ is Lipschitz continuous with Lipschitz constant $\gamma(x)$, that $\beta \in (0, 1)$ and that p is a descent direction at x . Then the Armijo condition

$$f(x + \alpha p) \leq f(x) + \alpha \beta \langle p, g(x) \rangle$$

is satisfied for all $\alpha \in [0, \alpha_{\max}(x, p)]$, where

$$\alpha_{\max}(x, p) = \frac{2(\beta - 1) \langle p, g(x) \rangle}{\gamma(x) \|p\|_2^2}.$$

Note that since $\gamma(x)$ is rarely known, the theorem does not provide a recipe for computing $\alpha_{\max}(x, p)$, merely a guarantee that there is such a suitable value. The numerator in $\alpha_{\max}(x, p)$ corresponds to the slope and the denominator to the curvature term. It can be interpreted as follows: If the curvature term is large, then the admissible range of α is small. Similarly, if the projected gradient along the search direction is large, then the range of admissible α is larger.

It then follows that the Backtracking-Armijo linesearch can be guaranteed to terminate with a suitably modest stepsize.

Corollary 2.2. Suppose that $f \in C^1$, that $g(x)$ is Lipschitz continuous with Lipschitz constant γ_k at x_k , that $\beta \in (0, 1)$ and that p_k is a descent direction at x_k . Then the stepsize generated by the backtracking-Armijo linesearch terminates with

$$\alpha_k \geq \min \left(\alpha_{\text{init}}, \frac{2\tau(\beta - 1) \langle p_k, g_k \rangle}{\gamma_k \|p_k\|_2^2} \right).$$

Again, since γ_k is rarely known, the corollary does not give a practical means for computing α_k , just an assurance that there is a suitable value. Notice that the stepsize is certainly not too large, since it is bounded above by α_{\max} , and can only be small when $\langle p, g(x) \rangle / \|p\|_2^2$ is. This will be the key to the successful termination of generic linesearch methods.

2.3 Convergence of generic linesearch methods

In order to tie all of the above together, we first need to state our Generic Linesearch Method:

Given an initial guess x_0 , let $k = 0$
 Until convergence:
 Find a descent direction p_k at x_k .
 Compute a stepsize α_k using a
 backtracking-Armijo linesearch along p_k .
 Set $x_{k+1} = x_k + \alpha_k p_k$, and increase k by 1.

It is then quite straightforward to apply Corollary 2.2 to deduce the following very general convergence result.

Theorem 2.3. Suppose that $f \in C^1$ and that g is Lipschitz continuous on \mathbb{R}^n . Then, for the iterates generated by the Generic Linesearch Method,

either

$$g_l = 0 \text{ for some } l \geq 0$$

or

$$\lim_{k \rightarrow \infty} f_k = -\infty$$

or

$$\lim_{k \rightarrow \infty} \min \left(|\langle p_k, g_k \rangle|, \frac{|\langle p_k, g_k \rangle|}{\|p_k\|_2} \right) = 0.$$

In words, either we find a first-order stationary point in a finite number of iterations, or we encounter a sequence of iterates for which the objective function is unbounded from below, or the slope (or a normalized slope) along the search direction converges to zero. While the first two of these possibilities are straightforward and acceptable consequences, the latter is perhaps not. For one thing, it certainly does not say that the gradient converges to zero, that is the iterates may not ultimately be first-order critical, since it might equally occur if the search direction and gradient tend to be mutually orthogonal. Thus we see that simply requiring that p_k be a descent direction is not a sufficiently demanding requirement. We will return to this shortly, but first we consider *the* archetypical globally convergent algorithm, the method of steepest descent.

2.4 Method of steepest descent

We have just seen that the Generic Linesearch Method may not succeed if the search direction becomes orthogonal to the gradient. Is there a direction for which this is impossible? Yes, when the search direction is the descent direction

$$p_k = -g_k,$$

the so-called *steepest-descent* direction—the epithet is appropriate since this direction solves the problem

$$\underset{p \in \mathbb{R}^n}{\text{minimize}} \quad m_k^L(x_k + p) \stackrel{\text{def}}{=} f_k + \langle p, g_k \rangle \quad \text{subject to} \quad \|p\|_2 = \|g_k\|_2,$$

and thus gives the greatest possible reduction in a first-order model of the objective function for a step whose length is specified. Global convergence follows immediately from Theorem 2.3.

Theorem 2.4. Suppose that $f \in C^1$ and that g is Lipschitz continuous on \mathbb{R}^n . Then, for the iterates generated by the Generic Linesearch Method using the steepest-descent direction,

either

$$g_l = 0 \text{ for some } l \geq 0$$

or

$$\lim_{k \rightarrow \infty} f_k = -\infty$$

or

$$\lim_{k \rightarrow \infty} g_k = 0.$$

As we mentioned above, this theorem suggests that steepest descent really is the archetypical globally convergent method, and in practice many other methods resort to steepest descent when they run into trouble. However, the method is not scale invariant, as re-scaling variables can lead to widely different “steepest-descent” directions. Even worse, as we can see in Figure 2.4, convergence

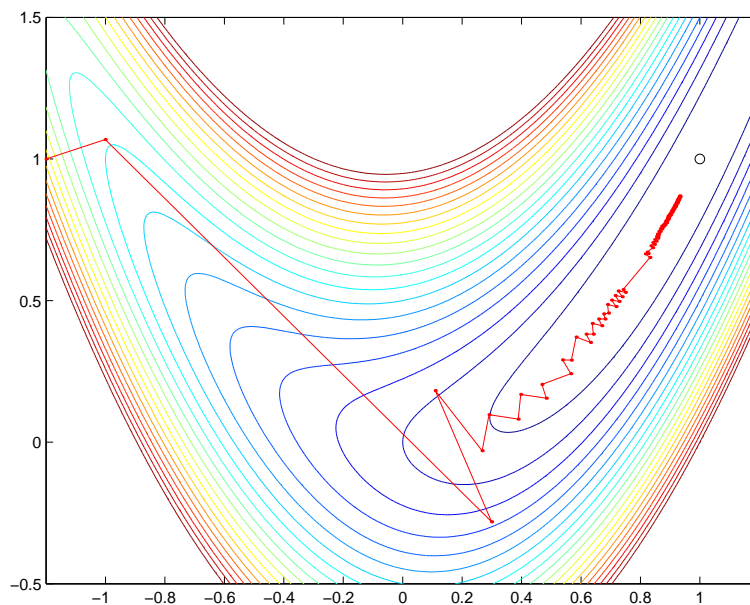


Figure 2.4: Contours for the objective function $f(x, y) = 10(y - x^2)^2 + (x - 1)^2$, and the iterates generated by the Generic Linesearch steepest-descent method.

may be (and actually almost always is) very slow in theory, while numerically convergence sometimes does not occur at all as the iteration stagnates. In practice, steepest-descent is all but worthless in most cases. The figure exhibits quite typical behaviour in which the iterates repeatedly oscillate

from one side of a objective function “valley” to the other. All of these phenomena may be attributed to a lack of attention to problem curvature when building the search direction. We now turn to methods that try to avoid this defect.

2.5 More general descent methods

2.5.1 Newton and Newton-like methods

Let B_k be a symmetric, positive definite matrix. Then it is trivial to show that the search direction p_k for which

$$B_k p_k = -g_k$$

is a descent direction. In fact, this direction solves the direction-finding problem

$$\underset{p \in \mathbb{R}^n}{\text{minimize}} \quad m_k^Q(x_k + p) \stackrel{\text{def}}{=} f_k + \langle p, g_k \rangle + \frac{1}{2} \langle p, B_k p \rangle, \quad (2.1)$$

where $m_k^Q(x_k + p)$ is a quadratic approximation to the objective function at x_k .

Of particular interest is the possibility that $B_k = H_k$, for in this case $m_k^Q(x_k + p)$ gives a second-order Taylor’s approximation to $f(x_k + p)$. The resulting direction for which

$$H_k p_k = -g_k$$

is known as the *Newton* direction, and any method which uses it is a Newton method. But notice that the Newton direction is only guaranteed to be useful in a linesearch context if the Hessian H_k is positive definite, for otherwise p_k might turn out to be an ascent direction.

It is also worth saying that while one can motivate such Newton-like methods from the prospective of minimizing a local second-order model of the objective function, one could equally argue that they aim to find a zero of a local first-order model

$$g(x_k + p) \approx g_k + B_k p_k$$

of its gradient. So long as B_k remains “sufficiently” positive definite, we can make precisely the same claims for these second-order methods as for those based on steepest descent.

Theorem 2.5. Suppose that $f \in C^1$ and that g is Lipschitz continuous on \mathbb{R}^n . Then, for the iterates generated by the Generic Linesearch Method using the Newton or Newton-like direction,

either

$$g_l = 0 \text{ for some } l \geq 0$$

or

$$\lim_{k \rightarrow \infty} f_k = -\infty$$

or

$$\lim_{k \rightarrow \infty} g_k = 0$$

provided that the eigenvalues of B_k are uniformly bounded and bounded away from zero.

Indeed, one can regard such methods as “scaled” steepest descent, but they have the advantage that they can be made scale invariant for suitable B_k , and crucially, as we see in Figure 2.5, their

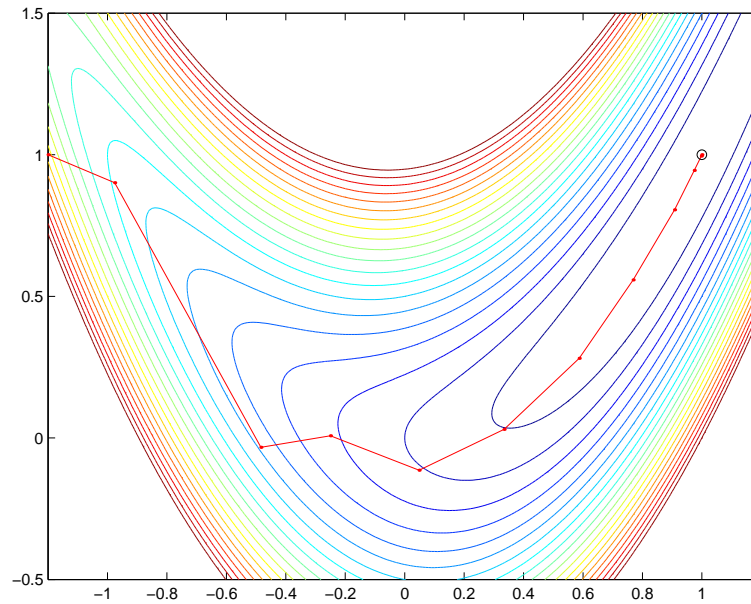


Figure 2.5: Contours for the objective function $f(x, y) = 10(y - x^2)^2 + (x - 1)^2$, and the iterates generated by the Generic Linesearch Newton method.

convergence is often significantly faster than steepest descent. In particular, in the case of the Newton direction, the Generic Linesearch method will usually converge very rapidly indeed.

Theorem 2.6. Suppose that $f \in C^2$ and that H is Lipschitz continuous on \mathbb{R}^n . Then suppose that the iterates generated by the Generic Linesearch Method with $\alpha_{\text{init}} = 1$ and $\beta < \frac{1}{2}$, in which the search direction is chosen to be the Newton direction $p_k = -H_k^{-1}g_k$ whenever H_k is positive definite, has a limit point x_* for which $H(x_*)$ is positive definite. Then

- (i) $\alpha_k = 1$ for all sufficiently large k ,
- (ii) the entire sequence $\{x_k\}$ converges to x_* , and
- (iii) the rate is Q-quadratic, i.e, there is a constant $\kappa \geq 0$.

$$\lim_{k \rightarrow \infty} \frac{\|x_{k+1} - x_*\|_2}{\|x_k - x_*\|_2^2} \leq \kappa.$$

2.5.2 Modified-Newton methods

Of course, away from a local minimizer there is no reason to believe that H_k will be positive definite, so precautions need to be taken to ensure that Newton and Newton-like linesearch methods, for which B_k is (or is close to) H_k , satisfy the assumptions of the global convergence Theorem 2.5. If H_k is indefinite, it is usual to solve instead

$$(H_k + M_k)p_k = -g_k,$$

where M_k is chosen so that $H_k + M_k$ is “sufficiently” positive definite and $M_k = 0$ when H_k is itself “sufficiently” positive definite. This may be achieved in a number of ways.

Firstly, if H_k has the spectral (that is eigenvector-eigenvalue) decomposition $H_k = Q_k D_k Q_k^T$, then M_k may be chosen so that

$$H_k + M_k = Q_k \max(\epsilon I, |D_k|) Q_k^T$$

for some “small” ϵ . This will shift all the insufficiently positive eigenvalues by as little as possible as is needed to make the overall matrix positive definite. While such a decomposition may be too expensive to compute for larger problems, a second, cheaper alternative is to find (or estimate) the smallest (necessarily real!) eigenvalue, $\lambda_{\min}(H_k)$, of H_k , and to set

$$M_k = \max(0, \epsilon - \lambda_{\min}(H_k))I$$

so as to shift *all* the eigenvalues by just enough as to make the smallest “sufficiently” positive. While this is often tried in practice, in the worst case it may have the effect of over-emphasising one large, negative eigenvalue at the expense of the remaining small, positive ones, and in producing a direction which is essentially steepest descent. Finally, a good compromise is instead to attempt a Cholesky factorization of H_k , and to alter the generated factors if there is evidence that the factorization will otherwise fail. There are a number of so-called *Modified Cholesky* factorizations, each of which will obtain

$$H_k + M_k = L_k L_k^T,$$

where M_k is zero for sufficiently positive-definite H_k , and “not-unreasonably large” in all other cases.

2.5.3 Quasi-Newton methods

It was fashionable in the 1960s and 1970s to attempt to build suitable approximations B_k to the Hessian, H_k . Activity in this area has subsequently died down, possibly because people started to realize that computing exact second derivatives was not as onerous as they had previously contended, but these techniques are still of interest particularly when gradients are awkward to obtain (such as when the function values are simply given as the result of some other, perhaps hidden, computation). There are broadly two classes of what may be called quasi-Newton methods.

The first are simply based on estimating columns of H_k by *finite differences*. For example, we might use the approximation

$$(H_k)e_i \approx h^{-1}(g(x_k + he_i) - g_k) \stackrel{\text{def}}{=} (B_k)e_i$$

for some “small” scalar $h > 0$. The difficulty here is in choosing an appropriate value for h : too large a value gives inaccurate approximations, while a too small one leads to large numerical cancellation errors.

The second sort of quasi-Newton methods are known as *secant approximations*, and try to ensure the *secant condition*

$$B_{k+1}s_k = y_k, \text{ where } s_k = x_{k+1} - x_k \text{ and } y_k = g_{k+1} - g_k,$$

that would be true if $H(x)$ were constant, is satisfied. The secant condition gives a lot of flexibility, and among the many methods that have been discovered, the *Symmetric Rank-1* method, for which

$$B_{k+1} = B_k + \frac{(y_k - B_k s_k)(y_k - B_k s_k)^T}{\langle s_k, y_k - B_k s_k \rangle},$$

and the *BFGS* method, for which

$$B_{k+1} = B_k + \frac{y_k y_k^T}{\langle s_k, y_k \rangle} - \frac{B_k s_k s_k^T B_k}{\langle s_k, B_k s_k \rangle}$$

are the best known (and generally the best). Note that the former may give indefinite approximations (or even fail), while the latter is guaranteed to generate symmetric and positive definite matrices so long as B_0 is positive definite and $\langle s_k, y_k \rangle > 0$ (the last condition may be ensured by an appropriate ‘‘Goldstein’’ linesearch). Since both of these secant methods are based on low-rank updates, it is possible to keep the per-iteration linear algebraic requirements at a more modest level for such methods than is generally possible with Newton or finite-difference methods.

2.5.4 Conjugate-gradient and truncated-Newton methods

And what if the problem is large and matrix factorization is out of the question? We have already considered (and rejected) steepest-descent methods. Is there something between the simplicity of steepest descent and the power (but expense) of Newton-like methods? Fortunately, the answer is yes.

Suppose that instead of solving (2.1), we instead find our search direction as

$$p_k = (\text{approximate}) \arg \min_{p \in \mathbb{R}^n} q(p) = f_k + \langle p, g_k \rangle + \frac{1}{2} \langle p, B_k p \rangle,$$

where we assume that B_k is positive definite—the key word here is *approximate*. Suppose that instead of minimizing q over all $p \in \mathbb{R}^n$, we restrict p to lie in a (much) smaller subspace—of course if we do this we will not (likely) obtain the optimal value of q , but we might hope to obtain a good approximation with considerably less effort.

Let $D^i = (d^0 : \dots : d^{i-1})$ be any collection of i vectors, let

$$\mathcal{D}^i = \{p \mid p = D^i p_d \text{ for some } p_d \in \mathbb{R}^i\}$$

be the subspace spanned by D^i , and suppose that we choose to pick

$$p^i = \arg \min_{p \in \mathcal{D}^i} q(p).$$

Then immediately $D^{i T} g^i = 0$, where $g^i = B_k p^i + g_k$ is the gradient of q at p^i . More revealingly, since $p^{i-1} \in \mathcal{D}^i$, it follows that $p^i = p^{i-1} + D^i p_d^i$, where

$$\begin{aligned} p_d^i &= \arg \min_{p_d \in \mathbb{R}^i} \langle p_d, D^{i T} g^{i-1} \rangle + \frac{1}{2} \langle p_d, D^{i T} B_k D^i p_d \rangle \\ &= -(D^{i T} B_k D^i)^{-1} D^{i T} g^{i-1} = -\langle d^{i-1}, g^{i-1} \rangle (D^{i T} B_k D^i)^{-1} e^i. \end{aligned}$$

Hence

$$p^i = p^{i-1} - \langle d^{i-1}, g^{i-1} \rangle D^i (D^{i T} B_k D^i)^{-1} e^i. \quad (2.2)$$

All of this is true regardless of D^i . But now suppose that the members of \mathcal{D}^i are B_k -conjugate, that is to say that $\langle d_i, B_k d_j \rangle = 0$ for all $i \neq j$. If this is so (2.2) becomes

$$p^i = p^{i-1} + \alpha^{i-1} d^{i-1}, \quad \text{where } \alpha^{i-1} = -\frac{\langle d^{i-1}, g^{i-1} \rangle}{\langle d^{i-1}, B_k d^{i-1} \rangle}. \quad (2.3)$$

Thus so long as we can generate B_k -conjugate vectors, we can build up successively improving approximations to the minimize of q by solving a sequence of *one-dimensional* minimization problems—the relationship (2.3) may be interpreted as finding α^{i-1} to minimize $q(p^{i-1} + \alpha d^{i-1})$. But can we find suitable B_k -conjugate vectors?

Surprisingly perhaps, yes, it is easy. Since g^i is independent of \mathcal{D}^i , let

$$d^i = -g^i + \sum_{j=0}^{i-1} \beta^{ij} d^j$$

for some unknown β^{ij} . Then elementary manipulation (and a cool head) shows that if we choose β^{ij} so that d^i is B -conjugate to \mathcal{D}^i , we obtain the wonderful result that

$$\beta^{ij} = 0 \text{ for } j < i - 1, \text{ and } \beta^{i, i-1} \equiv \beta^{i-1} = \frac{\|g_i\|_2^2}{\|g_{i-1}\|_2^2}.$$

That is, almost all of the β^{ij} are zero! Summing all of this up, we arrive at the method of *conjugate gradients* (CG):

Given $p^0 = 0$, set $g^0 = g_k$, $d^0 = -g_k$ and $i = 0$.
 Until g^i is “small”, iterate:
 $\alpha^i = \|g^i\|_2^2 / \langle d^i, B d^i \rangle$
 $p^{i+1} = p^i + \alpha^i d^i$
 $g^{i+1} = g^i + \alpha^i B_k d^i$
 $\beta^i = \|g^{i+1}\|_2^2 / \|g^i\|_2^2$
 $d^{i+1} = -g^{i+1} + \beta^i d^i$
 and increase i by 1.

Important features are that $\langle d^j, g^{i+1} \rangle = 0$ and $\langle g^j, g^{i+1} \rangle = 0$ for all $j = 0, \dots, i$, and most particularly that $\langle p^i, g_k \rangle \leq \langle p^{i-1}, g_k \rangle < 0$ for $i = 1, \dots, n$, from which we see that *any* $p_k = p^i$ is a descent direction.

In practice the above conjugate gradient iteration may be seen to offer a compromise between the steepest-descent direction (stopping when $i = 1$) and a Newton (-like) direction (stopping when $i = n$). For this reason, using such a curtailed conjugate gradient step within a linesearch (or trust-region) framework is often known as a *truncated*-Newton method. Frequently the size of g^i relative to g_k is used as a stopping criteria, a particularly popular rule being to stop the conjugate-gradient iteration when

$$\|g^i\| \leq \min(\|g_k\|^\omega, \eta) \|g_k\|,$$

where η and $\omega \in (0, 1)$, since then a faster-than-linear asymptotic convergence rate may be achieved if $B_k = H_k$.

3 TRUST-REGION METHODS FOR UNCONSTRAINED OPTIMIZATION

In this section, we continue to concentrate on the unconstrained minimization problem, and shall as before assume that the objective function is C^1 (sometimes C^2) with Lipschitz continuous derivatives.

3.1 Linesearch vs. trust-region methods

One might view linesearch methods as naturally “optimistic”. Fairly arbitrary search directions are permitted—essentially 50% of all possible directions give descent from a given point—while unruly behaviour is held in check via the linesearch. There is, however, another possibility, that more control is taken when choosing the search direction, with the hope that this will then lead to a higher probability that the (full) step really is useful for reducing the objective. This naturally “conservative” approach is the basis of trust-region methods.

As we have seen, linesearch methods pick a descent direction p_k , then pick a stepsize α_k to “reduce” $f(x_k + \alpha p_k)$ and finally accept $x_{k+1} = x_k + \alpha_k p_k$. *Trust-region* methods, by contrast, pick the overall step s_k to reduce a “model” of $f(x_k + s)$, and accept $x_{k+1} = x_k + s_k$ if the decrease predicted by the model is realised by $f(x_k + s_k)$. Since there is no guarantee that this will always be so, the fall-back mechanism is to set $x_{k+1} = x_k$, and to “refine” the model when the existing model produces a poor step. Thus, while a linesearch method recovers from a poor step by retreating along a parametric (usually linear) curve, a trust-region method recovers by reconsidering the whole step-finding procedure.

3.2 Trust-region models

It is natural to build a model of $f(x_k + s)$ by considering Taylor series approximations. Of particular interest are the *linear* model

$$m_k^L(s) = f_k + \langle s, g_k \rangle,$$

and the *quadratic* model

$$m_k^Q(s) = f_k + \langle s, g_k \rangle + \frac{1}{2} \langle s, B_k s \rangle,$$

where B_k is a symmetric approximation to the local Hessian matrix H_k . However, such models are far from perfect. In particular, the models are unlikely to resemble $f(x_k + s)$ if s is large. More seriously, the models may themselves be unbounded from below so that any attempts to minimize them may result in a large step. This defect will always occur for the linear model (unless $g_k = 0$), and also for the quadratic model if B_k is indefinite (and possibly if B_k is only positive semi-definite). Thus simply using a Taylor-series model is fraught with danger.

There is, fortunately, a simple and effective way around this conundrum. The idea is to prevent the model $m_k(s)$ from being unboundedness by imposing a *trust-region* constraint

$$\|s\| \leq \Delta_k,$$

for some “suitable” scalar *radius* $\Delta_k > 0$, on the step. This is a natural idea, since we know from Theorem 1.1 that we can improve the approximation error $|f(x_k + s) - m_k(s)|$ by restricting the allowable step. Thus our *trust-region subproblem* is to

$$\text{approximately minimize } m_k(s) \text{ subject to } \|s\| \leq \Delta_k, \\ s \in \mathbb{R}^n$$

and we shall choose s_k as approximate solution of this problem. In theory, it does not depend on which norm $\|\cdot\|$ we use (at least, in finite-dimensional spaces), but in practice it might!

For simplicity, we shall concentrate on the second-order (Newton-like) model

$$m_k(s) = m_k^Q(s) = f_k + \langle s, g_k \rangle + \frac{1}{2} \langle s, B_k s \rangle$$

and any (consistent) trust-region norm $\|\cdot\|$ for which

$$\kappa_s \|\cdot\| \leq \|\cdot\|_2 \leq \kappa_l \|\cdot\|$$

for some $\kappa_l \geq \kappa_s > 0$. Notice that the gradient of $m_k(s)$ at $s = 0$ coincides with the gradient of f at x_k , and also, unlike for linesearch methods, $B_k = H_k$ is always allowed. The vast majority of models use the ℓ_1 , ℓ_2 or ℓ_∞ norms on \mathbb{R}^n , and for these we have $\|\cdot\|_2 \leq \|\cdot\|_2 \leq \|\cdot\|_2$ (obviously!!), $n^{-\frac{1}{2}} \|\cdot\|_1 \leq \|\cdot\|_2 \leq \|\cdot\|_1$ and $\|\cdot\|_\infty \leq \|\cdot\|_2 \leq n \|\cdot\|_\infty$.

3.3 Basic trust-region method

Having decided upon a suitable model, we now turn to the trust-region algorithm itself. As we have suggested, we shall choose to “accept” $x_{k+1} = x_k + s_k$ whenever (a reasonable fraction of) the predicted model decrease $f_k - m_k(s_k)$ is realized by the actual decrease $f_k - f(x_k + s_k)$. We measure this by computing the ratio

$$\rho_k = \frac{f_k - f(x_k + s_k)}{f_k - m_k(s_k)}$$

of actual to predicted decrease, and accepting the trust-region step when ρ_k is not unacceptably smaller than 1.0. If the ratio is close to (or larger than) 1.0, there is good reason to believe that future step computations may well benefit from an increase in the trust-region radius, so we allow a radius increase in this case. If, by contrast, there is poor agreement between the actual and predicted decrease (and particularly, if f actually increases), the current step is poor and should be rejected. In this case, we reduce the trust-region radius to encourage a more suitable step at the next iteration.

We may summarize the basic trust-region method as follows:

Given $k = 0$, $\Delta_0 > 0$ and x_0 , until “convergence” do:
 Build the second-order model $m(s)$ of $f(x_k + s)$.
 “Solve” the trust-region subproblem to find s_k
 for which $m(s_k)$ “<” f_k and $\|s_k\| \leq \Delta_k$, and define

$$\rho_k = \frac{f_k - f(x_k + s_k)}{f_k - m_k(s_k)}.$$

If $\rho_k \geq \eta_v$ [*very successful*]

$$0 < \eta_v < 1$$

set $x_{k+1} = x_k + s_k$ and $\Delta_{k+1} = \gamma_i \Delta_k$.

$$\gamma_i \geq 1$$

Otherwise if $\rho_k \geq \eta_s$ then [*successful*]

$$0 < \eta_s \leq \eta_v < 1$$

set $x_{k+1} = x_k + s_k$ and $\Delta_{k+1} = \Delta_k$.

Otherwise [*unsuccessful*]

set $x_{k+1} = x_k$ and $\Delta_{k+1} = \gamma_d \Delta_k$.

$$0 < \gamma_d < 1$$

Increase k by 1.

Reasonable values might be $\eta_v = 0.9$ or 0.99 , $\eta_s = 0.1$ or 0.01 , $\gamma_i = 2$, and $\gamma_d = 0.5$. In practice, these parameters might even be allowed to vary (within reasonable limits) from iteration to iteration. In particular, there would seem to be little justification in increasing the trust region radius following a very successful iteration unless $\|s_k\| \approx \Delta_k$, nor in decreasing the radius by less than is required to “cut off” an unsuccessful s_k .

In practice, the trust-region radius is *not* increased for a very successful iterations, if the step is much shorter, say less than half the trust-region radius. There exist various schemes for choosing an initial trust-region radius. However, if the problem is well scaled, then $\Delta_0 = O(1)$ is reasonable. Poor scaling can affect the performance of trust-region methods. In practice it often suffices that the variables of the (scaled) problem have roughly the same order of magnitude.

It remains for us to decide what we mean by “solving” the trust-region subproblem. We shall see in Section 3.5 that (at least in the ℓ_2 -trust-region norm case) it is possible to find the (global) solution to the subproblem. However, since this may result in a considerable amount of work, we first seek “minimal” conditions under which we can guarantee convergence of the above algorithm to a first-order critical point.

We have already seen that steepest-descent linesearch methods have very powerful (theoretical) convergence properties. The same is true in the trust-region framework. Formally, at the very least, we shall require that we achieve as much reduction in the model as we would from an iteration of steepest descent. That is, if we define the *Cauchy* point as $s_k^c = -\alpha_k^c g_k$, where

$$\begin{aligned} \alpha_k^c &= \arg \min_{\alpha > 0} m_k(-\alpha g_k) \text{ subject to } \alpha \|g_k\| \leq \Delta_k \\ &= \arg \min_{0 < \alpha \leq \Delta_k / \|g_k\|} m_k(-\alpha g_k) \end{aligned} ,$$

we shall require that our step s_k satisfies

$$m_k(s_k) \leq m_k(s_k^c) \text{ and } \|s_k\| \leq \Delta_k. \quad (3.1)$$

Notice that the Cauchy point is extremely easy to find, since it merely requires that we minimize

the quadratic model along a line segment. In practice, we shall hope to—and can—do far better than this, but for now (3.1) suffices.

Figure 3.1 illustrates the trust-region problem in four different situations. The contours of the original function are shown as dotted lines, while the contours of the trust-region model appear as solid lines with the ℓ_2 trust-region ball in bold. Clockwise from top left, the plots depict the following situations: first, a quadratic model with positive definite Hessian, next a linear model about the same point, the third plot shows a quadratic model with indefinite Hessian and the final plot is a quadratic model with positive definite Hessian whose minimizers lies outside the trust-region.

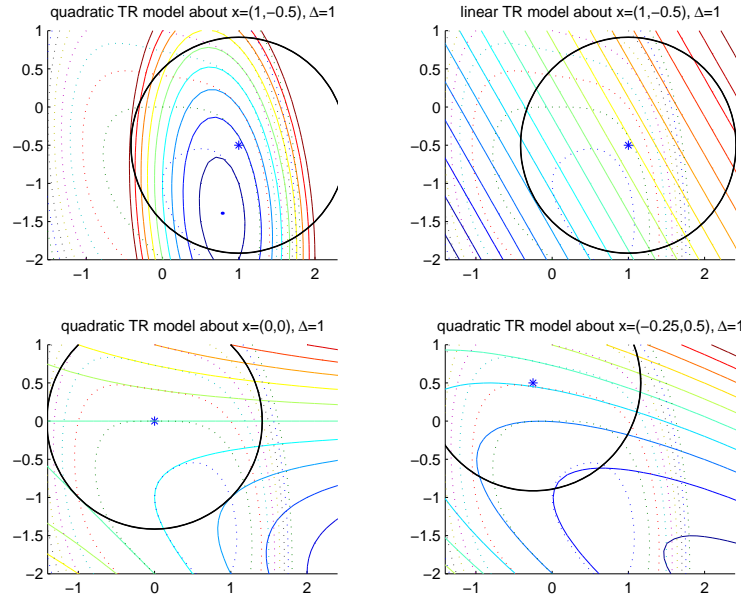


Figure 3.1: Trust-region models of $f(x) = x_1^4 + x_1x_2 + (1+x_2)^2$ about different points.

We now examine the convergence of this trust-region method.

3.4 Basic convergence of trust-region methods

The first thing to note is that we can guarantee a reasonable reduction in the model at the Cauchy point.

Theorem 3.1. If $m_k(s)$ is the second-order model and s_k^c is its Cauchy point within the trust-region $\|s\| \leq \Delta_k$, then

$$f_k - m_k(s_k^c) \geq \frac{1}{2} \|g_k\|_2 \min \left[\frac{\|g_k\|_2}{1 + \|B_k\|_2}, \kappa_s \Delta_k \right].$$

Observe that the guaranteed reduction depends on how large the current gradient is, and is also affected by the size of both the trust-region radius and the (inverse) of the Hessian.

Since our algorithm requires that the step does at least as well as the Cauchy point, we then have the following immediate corollary.

Corollary 3.2. If $m_k(s)$ is the second-order model, and s_k is an improvement on the Cauchy point within the trust-region $\|s\| \leq \Delta_k$,

$$f_k - m_k(s_k) \geq \frac{1}{2} \|g_k\|_2 \min \left[\frac{\|g_k\|_2}{1 + \|B_k\|_2}, \kappa_s \Delta_k \right].$$

This is a typical trust-region result, in that it relates the model reduction to a measure of the distance to optimality, in this case measured in terms of the norm of the gradient.

It is also necessary to say something about how much the model and the objective can vary. Since we are using a second-order model for which the first-two terms are exactly those from the Taylor's approximation, it is not difficult to believe that the difference between model and function will vary like the square of the norm of s_k , and indeed this is so.

Lemma 3.3. Suppose that $f \in C^2$, and that the true and model Hessians satisfy the bounds $\|H_k\|_2 \leq \kappa_h$ and $\|B_k\|_2 \leq \kappa_b$ for all k and some $\kappa_h \geq 1$ and $\kappa_b \geq 0$. Then

$$|f(x_k + s_k) - m_k(s_k)| \leq \kappa_d \Delta_k^2,$$

where $\kappa_d = \frac{1}{2} \kappa_l^2 (\kappa_h + \kappa_b)$, for all k .

Actually the result is slightly weaker than necessary since, for our purposes, we have chosen to replace $\|s_k\|$ by its (trust-region) bound Δ_k .

Armed with these bounds, we now arrive at a crucial result, namely that it will always be possible to make progress from a non-optimal point ($g_k \neq 0$).

Lemma 3.4. Suppose that $f \in C^2$, that the true and model Hessians satisfy the bounds $\|H_k\|_2 \leq \kappa_h$ and $\|B_k\|_2 \leq \kappa_b$ for all k and some $\kappa_h \geq 1$ and $\kappa_b \geq 0$, and that $\kappa_d = \frac{1}{2} \kappa_l^2 (\kappa_h + \kappa_b)$. Suppose furthermore that $g_k \neq 0$ and that

$$\Delta_k \leq \|g_k\|_2 \min \left(\frac{1}{\kappa_s (\kappa_h + \kappa_b)}, \frac{\kappa_s (1 - \eta_v)}{2 \kappa_d} \right).$$

Then iteration k is very successful and

$$\Delta_{k+1} \geq \Delta_k.$$

This result is fairly intuitive, since when the radius shrinks the model looks more and more like its

first-order Taylor's expansion (provided B_k is bounded) and thus ultimately there must be good local agreement between the model and objective functions.

The next result is a variation on its predecessor, and says that the radius is uniformly bounded away from zero if the same is true of the sequence of gradients, that is the radius will not shrink to zero at non-optimal points.

Lemma 3.5. Suppose that $f \in C^2$, that the true and model Hessians satisfy the bounds $\|H_k\|_2 \leq \kappa_h$ and $\|B_k\|_2 \leq \kappa_b$ for all k and some $\kappa_h \geq 1$ and $\kappa_b \geq 0$, and that $\kappa_d = \frac{1}{2}\kappa_l^2(\kappa_h + \kappa_b)$. Suppose furthermore that there exists a constant $\epsilon > 0$ such that $\|g_k\|_2 \geq \epsilon$ for all k . Then

$$\Delta_k \geq \kappa_\epsilon \stackrel{\text{def}}{=} \epsilon \gamma_d \min \left(\frac{1}{\kappa_s(\kappa_h + \kappa_b)}, \frac{\kappa_s(1 - \eta_v)}{2\kappa_d} \right)$$

for all k .

We may then deduce that if there are only a finite number of successful iterations, the iterates must be first-order optimal after the last of these.

Lemma 3.6. Suppose that $f \in C^2$, and that both the true and model Hessians remain bounded for all k . Suppose furthermore that there are only finitely many successful iterations. Then $x_k = x_*$ for all sufficiently large k and $g(x_*) = 0$.

Having ruled out this special (and highly unlikely) case, we then have our first global convergence result, namely that otherwise there is at least one sequence of gradients that converge to zero.

Theorem 3.7. Suppose that $f \in C^2$, and that both the true and model Hessians remain bounded for all k . Then either

$$g_l = 0 \text{ for some } l \geq 0$$

or

$$\lim_{k \rightarrow \infty} f_k = -\infty$$

or

$$\liminf_{k \rightarrow \infty} \|g_k\| = 0.$$

Is this all we can show? Is it possible for a second sub-sequence of gradients to stay bounded away from zero? Fortunately, no.

Corollary 3.8. Suppose that $f \in C^2$, and that both the true and model Hessians remain bounded for all k . Then either

$$g_l = 0 \text{ for some } l \geq 0$$

or

$$\lim_{k \rightarrow \infty} f_k = -\infty$$

or

$$\lim_{k \rightarrow \infty} g_k = 0.$$

Thus we have the highly-satisfying result that the gradients of the sequence $\{x_k\}$ generated by our algorithm converge to, or are all ultimately, zero. This does not mean that a subsequence of $\{x_k\}$ itself converges, but if it does, the limit is first-order critical.

It is also possible to show that an enhanced version of our basic algorithm converges to second-order critical points. To do so, we need to ensure that the Hessian of the model converges to that of the objective (as would obviously be the case if $B_k = H_k$), and that the step s_k has a significant component along the eigenvector corresponding to the most negative eigenvalue of B_k (if any). It is also possible to show that if $B_k = H_k$, if $\{x_k\}$ has a limit x_* for which $H(x_*)$ is positive definite, and if s_k is chosen to

$$\underset{s \in \mathbb{R}^n}{\text{minimize}} \quad m_k(s) \text{ subject to } \|s\| \leq \Delta_k, \quad (3.2)$$

the step Δ_k stays bounded away from zero, and thus the iteration ultimately becomes Newton's method (c.f. (2.1)).

In conclusion, we have seen that trust-region methods have a very rich underlying convergence theory. But so much for theory. We now turn to the outstanding practical issue, namely how one might hope to find a suitable step s_k . We will consider two possibilities, one that aims to get a very good approximation to (3.2), and a second, perhaps less ambitious method that is more geared towards large-scale computation.

3.5 Solving the trust-region subproblem

For brevity, we will temporarily drop the iteration subscript, and consider the problem of

$$\underset{s \in \mathbb{R}^n}{\text{(approximately) minimize}} \quad q(s) \equiv \langle s, g \rangle + \frac{1}{2} \langle s, Bs \rangle \text{ subject to } \|s\| \leq \Delta. \quad (3.3)$$

As we have already mentioned, our aim is to find s_* so that

$$q(s_*) \leq q(s^c) \text{ and } \|s_*\| \leq \Delta,$$

where s^c is the Cauchy point. We shall consider two approaches in this section. The first aims to solve (3.3) exactly, in which case our trust-region method will be akin to a Newton-like method. The second aims for an approximate solution using a conjugate-gradient like method. For simplicity, we shall only consider the ℓ_2 -trust region $\|s\| \leq \Delta$, mainly because there are very powerful methods in this case, but of course other norms are possible and are sometimes preferred in practice.

3.5.1 Solving the ℓ_2 -norm trust-region subproblem

There is a really powerful solution characterisation result for the ℓ_2 -norm trust-region subproblem.

Theorem 3.9. Any *global* minimizer s_* of $q(s)$ subject to $\|s\|_2 \leq \Delta$ satisfies the equation

$$(B + \lambda_* I)s_* = -g,$$

where $B + \lambda_* I$ is positive semi-definite, $\lambda_* \geq 0$ and $\lambda_*(\|s_*\|_2 - \Delta) = 0$. If $B + \lambda_* I$ is positive definite, s_* is unique.

This result is extraordinary as it is very unusual to be able to give necessary and sufficient *global* optimality conditions for a non-convex optimization problem (that is, a problem which might have a number of local minimizers). Even more extraordinary is the fact that the necessary and sufficient conditions are identical. But most crucially, these optimality conditions also suggest how we might solve the problem.

There are two cases to consider. If B is positive definite and the solution s to

$$Bs = -g \tag{3.4}$$

satisfies $\|s\|_2 \leq \Delta$, then it immediately follows that $s_* = s$ ($\lambda_* = 0$ in Theorem 3.9)—this potential solution may simply be checked by seeing if B has Cholesky factors and, if so, using these factors to solve (3.4) $Bs = -g$ and subsequently evaluate $\|s\|_2$. Otherwise, either B is positive definite but the solution to (3.4) satisfies $\|s\|_2 > \Delta$ or B is singular or indefinite. In these cases, Theorem 3.9 then says that s_* satisfies

$$(B + \lambda I)s = -g \text{ and } \langle s, s \rangle = \Delta^2, \tag{3.5}$$

which is a *nonlinear* (quadratic) system of algebraic equations in the $n+1$ unknowns s and λ . Thus, we now concentrate on methods for solving this system.

Suppose B has the spectral decomposition

$$B = U^T \Lambda U;$$

here U is a matrix of (orthonormal) eigenvectors while the diagonal matrix Λ is made up of eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$. Theorem 3.9 requires that $B + \lambda I$ be positive semi-definite, and so the solution (s, λ) to (3.5) that we seek necessarily satisfies $\lambda \geq -\lambda_1$. The first part of (3.5) enables us to write s explicitly in terms of λ , that is

$$s(\lambda) = -(B + \lambda I)^{-1}g;$$

we will temporarily disregard the possibility that the theorem permits a singular $B + \lambda I$. Notice that once we have found λ ,

$$(B + \lambda I)s = -g \tag{3.6}$$

is a linear system. In this case, we may substitute $s(\lambda)$ into the second part of (3.5) to reveal that

$$\psi(\lambda) \stackrel{\text{def}}{=} \|s(\lambda)\|_2^2 = \|U^T(\Lambda + \lambda I)^{-1}Ug\|_2^2 = \sum_{i=1}^n \frac{\gamma_i^2}{(\lambda_i + \lambda)^2} = \Delta^2, \tag{3.7}$$

where $\gamma_i = \langle e_i, Ug \rangle = \langle U^T e_i, g \rangle$. Thus to solve the trust-region subproblem, it appears that all we have to do is find a particular root of a univariate nonlinear equation.

We illustrate this in Figures 3.2–3.4.

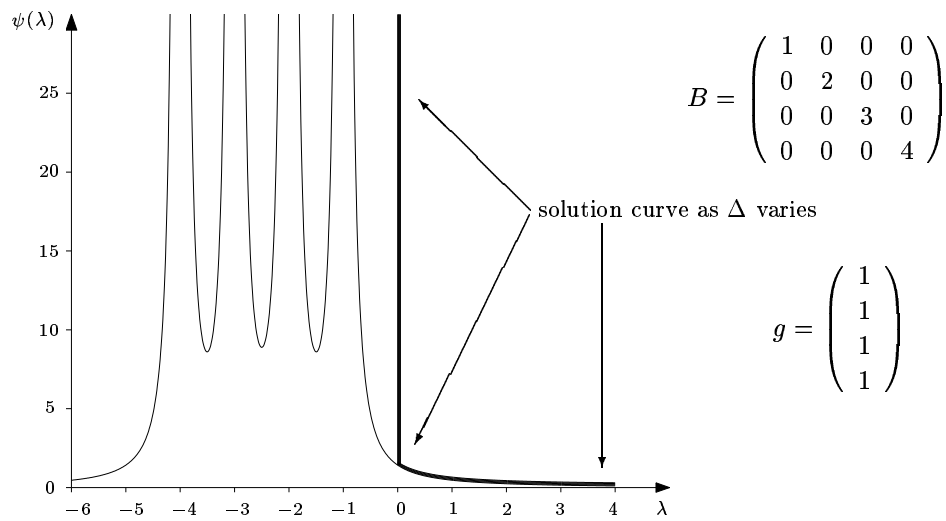


Figure 3.2: A plot of $\psi(\lambda)$ as λ varies from -6 to 4 . Note the poles at the negatives of the eigenvalues of H . The heavy curve plots λ against Δ ; the vertical component corresponds to interior solutions while the remaining segment indicates boundary solutions.

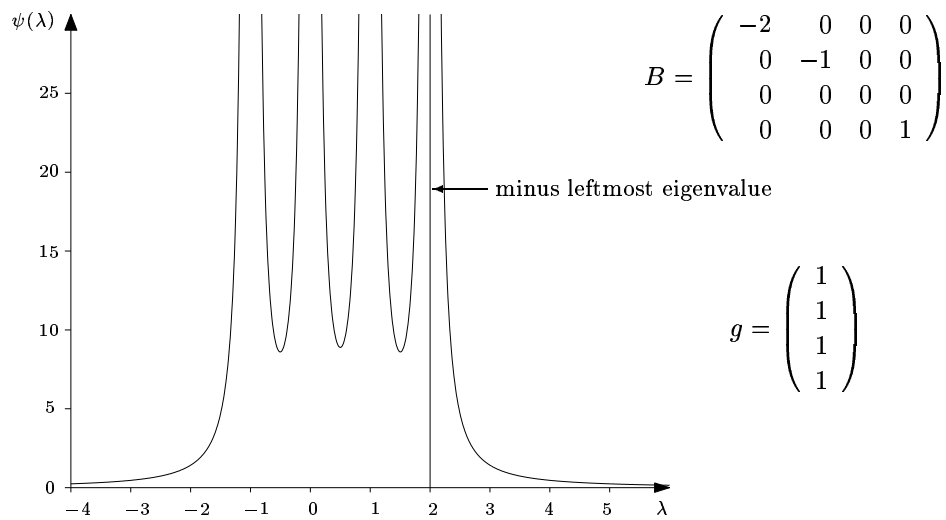


Figure 3.3: A plot of $\psi(\lambda)$ as λ varies from -4 to 5 . Again, note the poles at the negatives of the eigenvalues of H .

The first shows a convex example (B positive definite). For Δ^2 larger than roughly 1.5 , the solution to the problem lies in the interior of the trust region, and may be found directly from (3.4). When Δ is smaller than this, the solution lies on the boundary of the trust region, and can be found

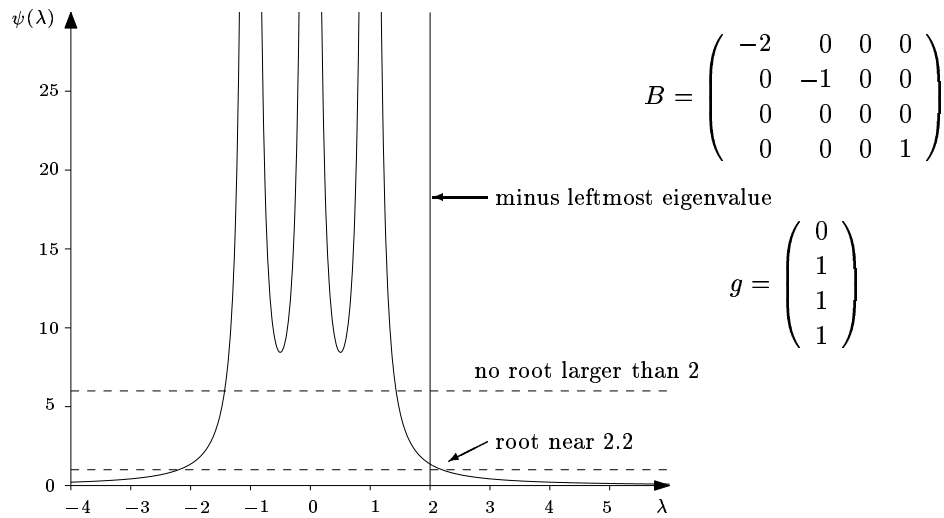


Figure 3.4: A plot of $\psi(\lambda)$ for the modified model as λ varies from -4 to 5 . Note that there is no solution with to the equation $\psi(\lambda) = \Delta^2$ with $\lambda \geq 2$ for Δ larger than roughly 1.2 .

as the right-most root of (3.7). The second example is non-convex (B indefinite). Now the solution must lie on the boundary of the trust region for all values of Δ , and again can be found as the right-most root of (3.7), to the right of $-\lambda_1$.

In both Figures 3.2 and 3.3 everything seems easy, and at least a semblance of an algorithm is obvious. But now consider the example in Figure 3.4. This example is especially chosen so that the coefficient γ_1 in (3.7) is zero, that is g is orthogonal to the eigenvector u_1 of B corresponding to the eigenvalue $\lambda_1 = -2$. Remember that Theorem 3.9 tells us that $\lambda \geq 2 = -\lambda_1$. But Figure 3.4 shows that there is no such root of (3.7) if Δ is larger than (roughly) 1.2 .

This is an example of what has become known as the *hard* case, which always arises when $\lambda_1 < 0$, $\langle u_1, g \rangle = 0$ and Δ is too big. What is happening? Quite simply, in the hard case $\lambda = -\lambda_1$ and (3.6) is a singular (but consistent) system—it is consistent precisely because $\langle u_1, g \rangle = 0$. But this system has other solutions $s + \alpha u_1$ for any α , because

$$(B + \lambda I)(s + \alpha u_1) = -g,$$

and u_1 is an eigenvector of $B + \lambda I$. The solution we require is that for which $\|s + \alpha u_1\|_2^2 = \Delta^2$, which is a quadratic equation for the unknown α , and either root suffices.

In the easy (that is not “hard”) case, it remains to see how best to solve $\|s(\lambda)\|_2 = \Delta$. The answer is blunt. Don’t! At least, not directly, since as the previous figures showed, $\psi(\lambda)$ is an unappealing function with many poles. It is far better to solve the equivalent *secular* equation

$$\phi(\lambda) \stackrel{\text{def}}{=} \frac{1}{\|s(\lambda)\|_2} - \frac{1}{\Delta} = 0,$$

as this has no poles, indeed it is an analytic function, and thus ideal for Newton’s method. We illustrate the secular equation in Figure 3.5.

Without giving details (for these, see the appendix, page 76), Newton’s method for the secular equation is as follows

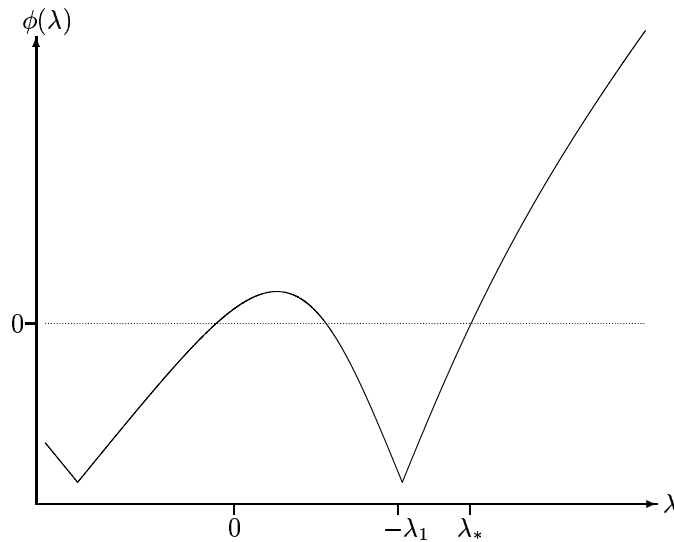


Figure 3.5: A plot of $\phi(\lambda)$ against λ for the problem of minimizing $-\frac{1}{4}s_1^2 + \frac{1}{4}s_2^2 + \frac{1}{2}s_1 + s_2$ subject to $\|s\|_2 \leq 0.4$.

Let $\lambda > -\lambda_1$ and $\Delta > 0$ be given.

Until “convergence” do:

Factorize $B + \lambda I = LL^T$.

Solve $LL^T s = -g$.

Solve $Lw = s$.

Replace λ by

$$\lambda + \left(\frac{\|s\|_2 - \Delta}{\Delta} \right) \left(\frac{\|s\|_2^2}{\|w\|_2^2} \right).$$

This is globally and ultimately quadratically convergent when started in the interval $[-\lambda_1, \lambda_*]$ except in the hard case, but needs to be safeguarded to make it robust for the hard and interior solution cases. Notice that the main computational cost per iteration is a Cholesky factorization of $B + \lambda I$, and while this may be reasonable for small problems, it may prove unacceptably expensive when the number of variables is large. We consider an alternative for this case next.

3.6 Solving the large-scale problem

Solving the large-scale trust-region subproblem using the above method is likely out of the question in all but very special cases. The obvious alternative is to use an iterative method to approximate its solution. The simplest approximation that is consistent with our fundamental requirement that we do as least as well as we would at the Cauchy point is to use the Cauchy point itself. Of course, this is simply the steepest descent method, and thus unlikely to be a practical method. The obvious generalization is the conjugate-gradient method, since the first step of CG is in the steepest-descent direction and, as subsequent CG steps further reduce the model, any step generated by the method is allowed by our theory. However, there are a number of other issues we need to address first. In

particular, what about the interaction between conjugate gradients and the trust region? And what if B is indefinite?

The conjugate-gradient method to find an approximation to a minimizer of $q(s)$ may be summarised as follows.

Given $s^0 = 0$, set $g^0 = g$, $d^0 = -g$ and $i = 0$.
 Until “breakdown” or g^i “small”, iterate:
 $\alpha^i = \|g^i\|_2^2 / \langle d^i, Bd^i \rangle$
 $s^{i+1} = s^i + \alpha^i d^i$
 $g^{i+1} = g^i + \alpha^i Bd^i$
 $\beta^i = \|g^{i+1}\|_2^2 / \|g^i\|_2^2$
 $d^{i+1} = -g^{i+1} + \beta^i d^i$
 and increase i by 1.

Notice that we have inserted a termination statement concerning “breakdown”. This is intended to cover the fatal case when $\langle d^i, Bd^i \rangle = 0$ (or, in practice, is close to zero), for which the iteration is undefined, and the non-fatal case when $\langle d^i, Bd^i \rangle < 0$ for which $q(s)$ is unbounded from below along the so-called *direction of negative curvature* d_i .

But what of the trust-region constraint? Here we have a crucial result.

Theorem 3.10. Suppose that the conjugate gradient method is applied to minimize $q(s)$ starting from $s^0 = 0$, and that $\langle d^i, Bd^i \rangle > 0$ for $0 \leq i \leq k$. Then the iterates s^j satisfy the inequalities

$$\|s^j\|_2 < \|s^{j+1}\|_2$$

for $0 \leq j \leq k - 1$.

Simply put, since the norm of the approximate solution generated by the conjugate gradients increases in norm at each iteration, if there is an iteration for which $\|s^j\|_2 > \Delta$, it must be that the solution to the trust-region subproblem lies on the trust-region boundary. That is $\|s_*\|_2 = \Delta$. This then suggests that we should apply the basic conjugate-gradient method above but terminate at iteration i if either (a) $\langle d^i, Bd^i \rangle \leq 0$, since this implies that $q(s)$ is unbounded along d^i , or (b) $\|s^i + \alpha^i d^i\|_2 > \Delta$, since this implies that the solution must lie on the trust-region boundary. In both cases, the simplest strategy is to stop on the boundary at $s = s^i + \alpha^B d^i$, where α^B chosen as positive root of the quadratic equation

$$\|s^i + \alpha^B d^i\|_2^2 = \Delta^2.$$

Crucially this s satisfies

$$q(s) \leq q(s^C) \text{ and } \|s\|_2 \leq \Delta$$

and thus Corollary 3.8 shows that the overall trust-region algorithm converges to a first-order critical point.

How good is this truncated conjugate-gradient strategy? In the convex case, it turns out to be very good. Indeed, no worse than half optimal!

Theorem 3.11. Suppose that the truncated conjugate gradient method is applied to approximately minimize $q(s)$ within $\|s\|_2 \leq \Delta$, and that B is positive definite. Then the computed and actual solutions to the problem, s and s_* , satisfy the bound $q(s) \leq \frac{1}{2}q(s_*)$.

In the non-convex (B_k indefinite) case, however, the strategy may be rather poor. For example, if $g = 0$ and B is indefinite, the above truncated conjugate-gradient method will terminate at $s = 0$, while the true solution lies on the trust-region boundary.

What can we do in the non-convex case? The answer is quite involved, but one possibility is to recall that conjugate-gradients is trying to solve the overall problem by successively solving the problem over a sequence of nested subspaces. As we saw, the CG method uses B -conjugate subspaces. But there is an equivalent method, the *Lanczos* method, that uses instead orthonormal bases. Essentially this may be achieved by applying the Gram-Schmidt procedure to the CG basis \mathcal{D}^i to build the equivalent basis $\mathcal{Q}^i = \{s \mid s = Q^i s_q \text{ for some } s_q \in \mathbb{R}^i\}$. It is easy to show that for this Q^i ,

$$Q^{iT}Q^i = I \text{ and } Q^{iT}BQ^i = T^i,$$

where T^i is tridiagonal, and $Q^{iT}g = \|g\|_2 e_1$, and it is trivial to generate Q^i from the CG \mathcal{D}^i . In this case the trust-region subproblem (3.3) may be rewritten as

$$s_q^i = \arg \min_{s_q \in \mathcal{R}^i} \|g\|_2 \langle e_1, s_q \rangle + \frac{1}{2} \langle s_q, T^i s_q \rangle \text{ subject to } \|s_q\|_2 \leq \Delta,$$

where $s^i = Q^i s_q^i$. Since T^i is tridiagonal, $T^i + \lambda I$ has very sparse Cholesky factors, and thus we can afford to solve this problem using the earlier secular equation approach. Moreover, since we will need to solve a sequence of related problems over nested subspaces, it is easy to imagine that one can use the solution for one problem to initialize the next. In practice, since the approach is equivalent to conjugate gradients, it is best to use CG until the trust-region boundary is reached and then to switch to the Lanczos method at that stage. Such a method has turned out to be most effective in practice.

4 INTERIOR-POINT METHODS FOR INEQUALITY CONSTRAINED OPTIMIZATION

Having given a break-neck description of methods for unconstrained minimization, we now turn our attention to the real problems of interest, namely those involving constraints. This section will focus on problems involving inequality constraints, while its successor will be concerned with equality constraints. But before we start, we need to discuss the conflicting nature of constrained optimization problems, and how we might deal with them.

Unconstrained minimization is “simple” because there is but one goal, namely to minimize the objective. This is not so for constrained minimization because there is now a conflict of requirements, the aforementioned objective minimization but at the same time a requirement of feasibility of the solution. While in some instances (such as for linear equality constraints and, to a certain extent,

all inequality constraints) it may be possible to generate feasible iterates, and thus to regain the advantages of having a single goal, this is not true for general constrained optimization.

4.1 Merit functions for constrained minimization

Most (but not all, see Section 5.4.3) constrained optimization techniques overcome this dichotomy by introducing a merit function to try to balance the two conflicting requirements of minimization and feasibility. Given parameters p , a composite function $\Phi(x, p)$ is a *merit function* if (some) minimizers of $\Phi(x, p)$ with respect to x approach those of $f(x)$ subject to the constraints as p approaches some set \mathcal{P} . Thus a merit function combines both optimality requirements into a single “artificial” objective function. In principal, it then only remains to use the best *unconstrained* minimization methods to solve the constrained problem. If only life were that simple!

Consider the case of equality constrained minimization, that is finding x_* to

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) \quad \text{subject to} \quad c(x) = 0. \quad (4.1)$$

A suitable merit function in this case is the *quadratic penalty function*

$$\Phi(x, \mu) = f(x) + \frac{1}{2\mu} \|c(x)\|_2^2, \quad (4.2)$$

where μ is a positive scalar parameter. It is easy to believe that if μ is small and we try to minimize $\Phi(x, \mu)$ much of the effort will be concentrated on making the second objective term $\frac{1}{2\mu} \|c(x)\|_2^2$ small, that is in forcing $c(x)$ to be small. But as f has a slight presence in the merit function, any remaining energy will be diverted to making $f(x)$ small amongst all of the values for which $c(x)$ is. Formally, it is easy to show that, under modest conditions, some minimizers of $\Phi(x, \mu)$ converge to solutions of (4.1) as μ approaches the set $\{0\}$ from above. Unfortunately, it is possible that $\Phi(x, \mu)$ may have other stationary points that are not solutions of (4.1)—indeed this must be the case if $c(x) = 0$ are inconsistent. The quadratic penalty function is but one of many merit functions for equality constrained minimization.

4.2 The logarithmic barrier function for inequality constraints

For the inequality constrained problem

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) \quad \text{subject to} \quad c(x) \geq 0 \quad (4.3)$$

the best known merit function is the *logarithmic barrier function*

$$\Phi(x, \mu) = f(x) - \mu \sum_{i=1}^m \log c_i(x),$$

where μ is again a positive scalar *barrier parameter*. Each logarithmic term $-\log c_i(x)$ becomes infinite as x approaches the boundary of the i -th inequality from the feasible side, and is undefined (effectively infinite) beyond there. The size of the logarithmic term is mitigated when μ is small, and it is then possible to get close to the boundary of the feasible region before its effect is felt, any minimization effort being directed towards reducing the objective. Once again, it is easy to show that, under modest conditions, some minimizers of $\Phi(x, \mu)$ converge to solutions of (4.3) as μ approaches the set $\{0\}$ from above. And once again a possible defect is that $\Phi(x, \mu)$ may have other, useless stationary points. The contours of a typical example are shown in Figure 4.1.

4.3 A basic barrier-function algorithm

The logarithmic barrier function is different in one vital aspect from the quadratic penalty function in that it requires that there is a *strictly* interior point. If we apply the obvious sequential minimization algorithm to $\Phi(x, \mu)$, a strictly interior starting point is required, and all subsequent iterates will be

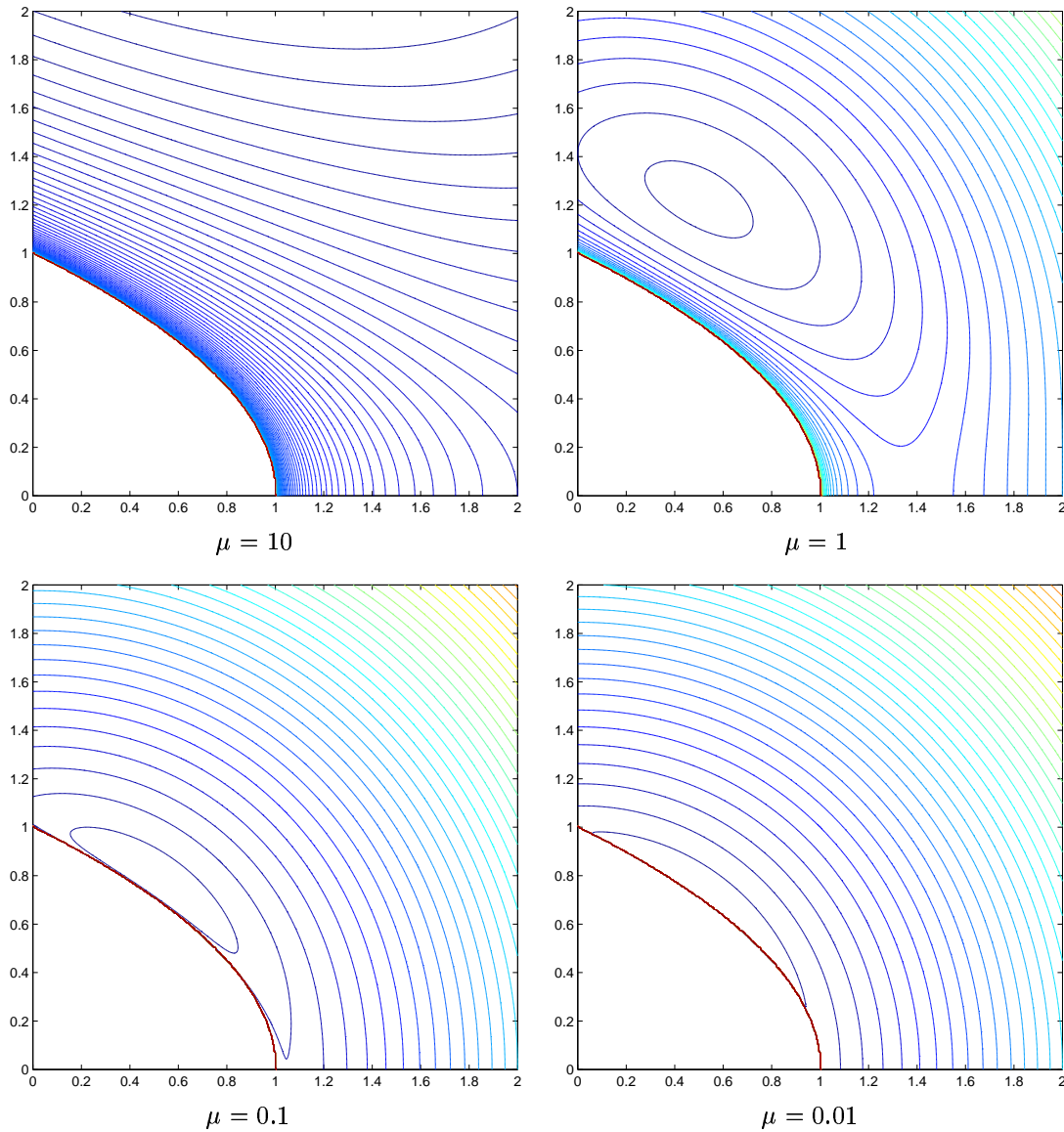


Figure 4.1: The logarithmic barrier function for $\min x_1^2 + x_2^2$ subject to $x_1 + x_2^2 \geq 1$. The contours for $\mu = 0.01$ are visually indistinguishable from $f(x)$ for feasible points.

strictly interior. The obvious “interior-point” algorithm is as follows.

Given $\mu_0 > 0$, set $k = 0$.
 Until “convergence”, iterate:
 Find x_k^s for which $c(x_k^s) > 0$.
 Starting from x_k^s , use an unconstrained
 minimization algorithm to find an
 “approximate” minimizer x_k of $\Phi(x, \mu_k)$.
 Compute $\mu_{k+1} > 0$ smaller than μ_k such
 that $\lim_{k \rightarrow \infty} \mu_{k+1} = 0$. and increase k by 1.

In practice it is common to choose $\mu_{k+1} = 0.1\mu_k$ or even $\mu_{k+1} = \mu_k^2$, while perhaps the obvious choice for a subsequent starting point is $x_{k+1}^s = x_k$.

Fortunately, as we have hinted, basic convergence for the algorithm is easily established. Recall that the *active set* $\mathcal{A}(x)$ at a point x is $\mathcal{A}(x) = \{i \mid c_i(x) = 0\}$. Then we have the following.

Theorem 4.1. Suppose that $f, c \in \mathcal{C}^2$, that $(y_k)_i \stackrel{\text{def}}{=} \mu_k / c_i(x_k)$ for $i = 1, \dots, m$, that

$$\|\nabla_x \Phi(x_k, \mu_k)\|_2 \leq \epsilon_k$$

where ϵ_k converges to zero as $k \rightarrow \infty$, and that x_k converges to x_* for which $\{a_i(x_*)\}_{i \in \mathcal{A}(x_*)}$ are linearly independent. Then x_* satisfies the first-order necessary optimality conditions for the problem

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) \quad \text{subject to} \quad c(x) \geq 0$$

and $\{y_k\}$ converge to the associated Lagrange multipliers y_* .

Notice here how the algorithm delivers something unexpected, namely estimates of the Lagrange multipliers. Also see the role played by the linear independence of the active constraint gradients, regrettably quite a strong constraint qualification.

4.4 Potential difficulties

As we now know that it suffices to (approximately) minimize $\Phi(x, \mu)$, how should we proceed? As $\Phi(x, \mu)$ is a smooth function, we can immediately appeal to the methods we discussed in Sections 2 and 3. But we need to be careful. Very, very careful.

We could use a linesearch method. Of note here is the fact that the barrier function has logarithmic singularities, indeed is undefined for infeasible points. Thus it makes sense to design a specialized linesearch to cope with the singularity of the log. Alternatively, we could use a trust-region method. Here we need to be able to instantly reject candidate steps for which $c(x_k + s_k) \not\geq 0$. More importantly, while all (consistent) trust-region norms are equivalent, (ideally) we should “shape” the trust region for any barrier-function model to cope with the contours of the singularity. This implies that

the trust-region shape may vary considerably from iteration to iteration, with its shape reflecting the eigenvalues arising from the singularity.

4.4.1 Potential difficulty I: ill-conditioning of the barrier Hessian

At the heart of both linesearch and trust-region methods is, of course, the Newton (second-order) model and related Newton direction. The computation of a Newton model/direction for the logarithmic barrier function is vital, and the resulting equations have a lot of (exploitable) structure. The gradient of the barrier function is

$$\nabla_x \Phi(x, \mu) = g(x) - \mu \sum_i a_i(x)/c_i(x) = g(x) - A^T(x)y(x) = g(x, y(x)),$$

where $y_i(x) \stackrel{\text{def}}{=} \mu/c_i(x)$ and $g(x, y)$ is the gradient of the Lagrangian function for (4.3). Likewise, the Hessian is

$$\nabla_{xx} \Phi(x, \mu) = H(x, y(x)) + \mu A^T(x)C^{-2}(x)A(x),$$

where $H(x, y(x)) = H(x) - \sum_{i=1}^m y_i(x)H_i(x)$ and $C(x) = \text{diag}(c_1(x), \dots, c_m(x))$, the diagonal matrix whose entries are the $c_i(x)$. Thus the Newton correction s^P from x for the barrier function satisfies

$$(H(x, y(x)) + \mu A^T(x)C^{-2}(x)A(x))s^P = -g(x, y(x)). \quad (4.4)$$

Since $y(x) = \mu C^{-1}(x)e$, (4.4) is sometimes written as

$$(H(x, y(x)) + A^T(x)C^{-1}(x)Y(x)A(x))s^P = -g(x, y(x)), \quad (4.5)$$

or

$$(H(x, y(x)) + A^T(x)Y^2(x)A(x)/\mu)s^P = -g(x, y(x)), \quad (4.6)$$

where $Y(x) = \text{diag}(y_1(x), \dots, y_m(x))$.

This is where we need to be careful. For we have the following estimates of the eigenvalues of the barrier function as we approach a solution.

Theorem 4.2. Suppose that the assumptions of Theorem 4.1 are satisfied, that $A_{\mathcal{A}}$ is the matrix whose rows are $\{a_i^T(x_*)\}_{i \in \mathcal{A}(x_*)}$, that $m_a = |\mathcal{A}(x_*)|$, and that x_* is *non-degenerate*, that is $(y_*)_i > 0$ for all $i \in \mathcal{A}(x_*)$. Then the Hessian matrix of the barrier function, $\nabla_{xx} \Phi(x_k, \mu_k)$, has m_a eigenvalues

$$\lambda_i(A_{\mathcal{A}}^T Y_{\mathcal{A}}^2 A_{\mathcal{A}})/\mu_k + O(1) \text{ for } i = 1, \dots, m_a$$

and the remaining $n - m_a$ eigenvalues

$$\lambda_i(N_{\mathcal{A}}^T H(x_*, y_*) N_{\mathcal{A}}) + O(\mu_k) \text{ for } i = 1, \dots, n - m_a$$

as $k \rightarrow \infty$, where $\lambda_i(\cdot)$ denotes the i -th eigenvalue of its matrix argument, $Y_{\mathcal{A}}$ is the diagonal matrix of active Lagrange multipliers at x_* and $N_{\mathcal{A}}$ is an orthogonal basis for the null-space of $A_{\mathcal{A}}$.

This demonstrates that the condition number of $\nabla_{xx} \Phi(x_k, \mu_k)$ is $O(1/\mu_k)$ as μ_k shrinks to zero, and suggests that it may not be straightforward to find the minimizer numerically. Look at how the contours around x_* in Figure 4.1 bunch together as μ approaches zero.

4.4.2 Potential difficulty II: poor starting points

As if this potential defect isn't serious enough, there is a second significant difficulty with the naive method we described earlier. This is that $x_{k+1}^s = x_k$ appears to be a very poor starting point for a Newton step just after the (small) barrier parameter is reduced. To see this suppose, as will be the case at the end of the minimization for the k -th barrier subproblem, that

$$0 \approx \nabla_x \Phi(x_k, \mu_k) = g(x_k) - \mu_k A^T(x_k) C^{-1}(x_k) e \approx g(x_k) - \mu_k A_{\mathcal{A}}^T(x_k) C_{\mathcal{A}}^{-1}(x_k) e,$$

the approximation being true because the neglected terms involve $y(x_k) = \mu_k / c_i(x_k)$ which converge to zero for inactive constraints. Then in the non-degenerate case, again roughly speaking, the Newton correction s^P for the new barrier parameter satisfies

$$\mu_{k+1} A_{\mathcal{A}}^T(x_k) C_{\mathcal{A}}^{-2}(x_k) A_{\mathcal{A}}(x_k) s^P \approx (\mu_{k+1} - \mu_k) A_{\mathcal{A}}^T(x_k) C_{\mathcal{A}}^{-1}(x_k) e \quad (4.7)$$

since

$$\nabla_x \Phi(x_k, \mu_{k+1}) \approx g(x_k) - \mu_{k+1} A_{\mathcal{A}}^T(x_k) C_{\mathcal{A}}^{-1}(x_k) e \approx (\mu_{k+1} - \mu_k) A_{\mathcal{A}}^T(x_k) C_{\mathcal{A}}^{-1}(x_k) e$$

and the $\mu_{k+1} A_{\mathcal{A}}^T(x_k) C_{\mathcal{A}}^{-2}(x_k) A_{\mathcal{A}}(x_k)$ term dominates $\nabla_{xx} \Phi(x_k, \mu_{k+1})$. If $A_{\mathcal{A}}(x_k)$ is full rank, then multiplying the approximation (4.7) from the left first by the generalized inverse, $(A_{\mathcal{A}} A_{\mathcal{A}}^T)^{-1} A_{\mathcal{A}}$ of $A_{\mathcal{A}}$ and then by $C_{\mathcal{A}}^2$ implies that

$$A_{\mathcal{A}}(x_k) s^P \approx \left(1 - \frac{\mu_k}{\mu_{k+1}}\right) c_{\mathcal{A}}(x_k)$$

from which a Taylor expansion of $c_{\mathcal{A}}(x_k + s^P)$ reveals that

$$c_{\mathcal{A}}(x_k + s^P) \approx c_{\mathcal{A}}(x_k) + A_{\mathcal{A}}(x_k) s^P \approx \left(2 - \frac{\mu_k}{\mu_{k+1}}\right) c_{\mathcal{A}}(x_k) < 0$$

whenever $\mu_{k+1} < \frac{1}{2}\mu_k$. Hence a Newton step will asymptotically be infeasible for anything but the most modest decrease in μ , and thus the method is unlikely to converge fast.

We will return to both of these issues shortly, but first we need to examine barrier methods in a seemingly different light.

4.5 A different perspective: perturbed optimality conditions

We now consider what, superficially, appears to be a completely different approach to inequality-constrained optimization. Recall from Theorem (1.8) that the first order optimality conditions for (4.3) are that there are Lagrange multipliers (or, as they are sometimes called, dual variables) y for which

$$\begin{aligned} g(x) - A^T(x)y &= 0 && \text{(dual feasibility)} \\ C(x)y &= 0 && \text{(complementary slackness) and} \\ c(x) &\geq 0 \text{ and } y \geq 0. \end{aligned}$$

Now consider the ‘‘perturbed’’ problem

$$\begin{aligned} g(x) - A^T(x)y &= 0 && \text{(dual feasibility)} \\ C(x)y &= \mu e && \text{(perturbed complementary slackness) and} \\ c(x) &> 0 \text{ and } y > 0, \end{aligned}$$

where $\mu > 0$.

Primal-dual path-following methods aim to track solutions to the system

$$g(x) - A^T(x)y = 0 \text{ and } C(x)y - \mu e = 0 \quad (4.8)$$

as μ shrinks to zero, while maintaining $c(x) > 0$ and $y > 0$. This approach has been amazingly successful when applied to linear programming problems, and has been extended to many other classes of convex optimization problems. Since (4.8) is simply a nonlinear system, an obvious (locally convergent) way to solve the system is, as always, to use Newton's method. It is easy to show that the Newton correction (s^{PD}, w) to (x, y) satisfies

$$\begin{pmatrix} H(x, y) & -A^T(x) \\ YA(x) & C(x) \end{pmatrix} \begin{pmatrix} s^{\text{PD}} \\ w \end{pmatrix} = - \begin{pmatrix} g(x) - A^T(x)y \\ C(x)y - \mu e \end{pmatrix}. \quad (4.9)$$

Using the second equation to eliminate w gives that

$$(H(x, y) + A^T(x)C^{-1}(x)YA(x)) s^{\text{PD}} = - (g(x) - \mu A^T(x)C^{-1}(x)e) = g(x, y(x)), \quad (4.10)$$

where, as before, $y(x) = \mu C^{-1}(x)e$. But now compare this with the Newton barrier system (4.5). Amazingly, the only difference is that the (left-hand-side) coefficient matrix in (4.5) mentions the specific $y(x)$ while that for (4.10) uses a generic y . And it is this difference that turns out to be crucial. The freedom to choose y in $H(x, y) + A^T(x)C^{-1}(x)YA(x)$ for the primal-dual approach proves to be vital. Making the primal choice $y(x) = \mu C^{-1}(x)e$ can be poor, while using a more flexible approach in which y is chosen by other means, such as through the primal-dual correction $y + w$ is often highly successful.

We now return to the potential difficulties with the primal approach we identified in Sections 4.4.1 and 4.4.2.

4.5.1 Potential difficulty II ... revisited

We first show that, despite our reservations in Section 4.4.2, the value $x_{k+1}^s = x_k$ can be a good starting point. The problem with the primal correction s^{P} is that the primal method has to choose $y = y(x_k^s) = \mu_{k+1}C^{-1}(x_k)e$, and this is a factor μ_{k+1}/μ_k too small to be a good Lagrange multiplier estimate—recall that Theorem 4.1 shows that $\mu_k C^{-1}(x_k)e$ converges to y_* .

But now suppose instead that we use the primal-dual correction s^{PD} and choose the “proper” $y = \mu_k C^{-1}(x_k)e$ rather than $y(x_k^s)$ —we know that this is a good choice insofar as this Newton step should decrease the dual infeasibility and complementary slackness since $(x_k, \mu_k C^{-1}(x_k)e)$ are already good estimates. In this case, arguing as before, in the non-degenerate case, the correction s^{PD} satisfies

$$\mu_k A_{\mathcal{A}}^T(x_k)C_{\mathcal{A}}^{-2}(x_k)A_{\mathcal{A}}(x_k)s^{\text{PD}} \approx (\mu_{k+1} - \mu_k)A_{\mathcal{A}}^T(x_k)C_{\mathcal{A}}^{-1}(x_k)e,$$

and thus if $A_{\mathcal{A}}(x_k)$ is full rank,

$$A_{\mathcal{A}}(x_k)s^{\text{PD}} \approx \left(\frac{\mu_{k+1}}{\mu_k} - 1 \right) c_{\mathcal{A}}(x_k).$$

Then using a Taylor expansion of $c_{\mathcal{A}}(x_k + s^{\text{PD}})$ reveals that

$$c_{\mathcal{A}}(x_k + s^{\text{PD}}) \approx c_{\mathcal{A}}(x_k) + A_{\mathcal{A}}(x_k)s^{\text{PD}} \approx \frac{\mu_{k+1}}{\mu_k} c_{\mathcal{A}}(x_k) > 0,$$

and thus $x_k + s^{\text{PD}}$ is feasible—the result is easy to show for inactive constraints. Hence, simply by using a different model Hessian we can compute a useful Newton correction from $x_{k+1}^s = x_k$ that both improves the violation of the optimality conditions (and ultimately leads to fast convergence) *and* stays feasible.

4.5.2 Primal-dual barrier methods

In order to globalize the primal-dual iteration, we simply need to build an appropriate model of the logarithmic barrier function within either a linesearch or trust-region framework for minimizing $\Phi(x, \mu_k)$. As we have already pointed out the disadvantages of only allowing the (primal) Hessian approximation $\nabla_{xx}\Phi(x_k, \mu_k)$, we instead prefer the more flexible search-direction model problem to (approximately)

$$\underset{s \in \mathbb{R}^n}{\text{minimize}} \quad \langle s, g(x, y(x)) \rangle + \frac{1}{2} \langle s, (H(x, y) + A^T(x)C^{-1}(x)YA(x))s \rangle, \quad (4.11)$$

possibly subject to a trust-region constraint. We have already noticed that the first-order term $g(x, y(x)) = \nabla_x \Phi(x, \mu)$ as $y(x) = \mu C^{-1}(x)e$, and thus the model gradient is that of the barrier function as required by our global convergence analyses of linesearch and trust-region methods. We have discounted always choosing $y = y(x)$ in (4.11), and have suggested that the choice $y = (\mu_{k-1}/\mu_k)y(x)$ when changing the barrier parameter results in good use of the starting point. Another possibility is to use $y = y^{\text{OLD}} + w^{\text{OLD}}$, where w^{OLD} is the primal-dual correction to the previous dual-variable estimates y^{OLD} . However, this needs to be used with care since there is no a priori assurance that $y^{\text{OLD}} + w^{\text{OLD}} > 0$, and indeed it is usual to prefer $y = \max(y^{\text{OLD}} + w^{\text{OLD}}, \epsilon(\mu_k)e)$ for some “small” $\epsilon(\mu_k) > 0$. The choice $\epsilon(\mu_k) = \mu_k^{1.5}$ leads to a realistic primal-dual method, although other precautions need sometimes to be taken.

4.5.3 Potential difficulty I ... revisited

We now return to the other perceived difficult with barrier or primal-dual path-following methods, namely that the inherent ill-conditioning in the barrier Hessian makes it hard to generate accurate Newton steps when the barrier parameter is small. Let \mathcal{I} be the set of inactive constraints at x_* , and denote the active and inactive components of c and y with suffices \mathcal{A} and \mathcal{I} respectively. Thus $c_{\mathcal{A}}(x_*) = 0$ and $c_{\mathcal{I}}(x_*) > 0$, while if the solution is non-degenerate, $y_{\mathcal{A}}(x_*) > 0$ and $y_{\mathcal{I}}(x_*) = 0$. As we have seen, the Newton correction s^{PD} satisfies (4.9), while the equivalent system (4.10) clearly has a condition number that approaches infinity as x and y reach their limits because $c_{\mathcal{A}}(x)$ approaches zero while $y_{\mathcal{A}}(x)$ approaches $y_{\mathcal{A}}(x_*) > 0$.

But now suppose that we separate (4.9) into

$$\begin{pmatrix} H(x, y) & -A_{\mathcal{A}}^T(x) & -A_{\mathcal{I}}^T(x) \\ Y_{\mathcal{A}}A_{\mathcal{A}}(x) & C_{\mathcal{A}}(x) & 0 \\ Y_{\mathcal{I}}A_{\mathcal{A}}(x) & 0 & C_{\mathcal{I}}(x) \end{pmatrix} \begin{pmatrix} s^{\text{PD}} \\ w_{\mathcal{A}} \\ w_{\mathcal{I}} \end{pmatrix} = - \begin{pmatrix} g(x) - A^T(x)y \\ C_{\mathcal{A}}(x)y_{\mathcal{A}} - \mu e \\ C_{\mathcal{I}}(x)y_{\mathcal{I}} - \mu e \end{pmatrix},$$

and then eliminate the variables $w_{\mathcal{I}}$, multiply the second equation by $Y_{\mathcal{A}}^{-1}$ and use $C_{\mathcal{I}}(x)y_{\mathcal{I}} = \mu e$, we obtain

$$\begin{pmatrix} H(x, y) + A_{\mathcal{I}}^T(x)C_{\mathcal{I}}(x)^{-1}Y_{\mathcal{I}}A_{\mathcal{I}}(x) & -A_{\mathcal{A}}^T(x) \\ A_{\mathcal{A}}(x) & C_{\mathcal{A}}(x)Y_{\mathcal{A}}^{-1} \end{pmatrix} \begin{pmatrix} s^{\text{PD}} \\ w_{\mathcal{A}} \end{pmatrix} = - \begin{pmatrix} g(x) - A_{\mathcal{A}}^T(x)y_{\mathcal{A}} - \mu A_{\mathcal{I}}^T(x)C_{\mathcal{I}}^{-1}(x)e \\ c_{\mathcal{A}}(x) - \mu Y_{\mathcal{A}}^{-1}e \end{pmatrix}. \quad (4.12)$$

But then we see that the terms involving inverses, $C_{\mathcal{I}}^{-1}(x)$ and $Y_{\mathcal{A}}^{-1}$, remain bounded, and indeed in the limit the system becomes

$$\begin{pmatrix} H(x, y) & -A_{\mathcal{A}}^T(x) \\ A_{\mathcal{A}}(x) & 0 \end{pmatrix} \begin{pmatrix} s^{\text{PD}} \\ w_{\mathcal{A}} \end{pmatrix} = - \begin{pmatrix} g(x) - A_{\mathcal{A}}^T(x)y_{\mathcal{A}} - \mu A_{\mathcal{I}}^T(x)C_{\mathcal{I}}^{-1}(x)e \\ 0 \end{pmatrix}$$

which is well behaved. Thus just because (4.10) is ill conditioned, this does not preclude us from finding s^{PD} from an equivalent, perfectly well-behaved system like (4.12).

4.6 A practical primal-dual method

Following on from the above, we now give the skeleton of a reasonable primal-dual method.

Given $\mu_0 > 0$ and feasible (x_0^s, y_0^s) , set $k = 0$.
 Until “convergence”, iterate:
 Inner minimization: starting from (x_k^s, y_k^s) , use an
 unconstrained minimization algorithm to find (x_k, y_k) for which
 $\|C(x_k)y_k - \mu_k e\| \leq \mu_k$ and $\|g(x_k) - A^T(x_k)y_k\| \leq \mu_k^{1.00005}$.
 Set $\mu_{k+1} = \min(0.1\mu_k, \mu_k^{1.9999})$.
 Find (x_{k+1}^s, y_{k+1}^s) using a primal-dual Newton step from (x_k, y_k) .
 If (x_{k+1}^s, y_{k+1}^s) is infeasible, reset (x_{k+1}^s, y_{k+1}^s) to (x_k, y_k) .
 Increase k by 1.

The inner minimization will be performed by either a linesearch or trust-region method for minimizing $\Phi(x, \mu_k)$, the stopping rules $\|C(x_k)y_k - \mu_k e\| \leq \mu_k$ and $\|g(x_k) - A^T(x_k)y_k\| \leq \mu_k^{1.00005}$ certainly being attainable as the first-order optimality condition for minimizing $\Phi(x, \mu_k)$ is that $g(x) - A^T(x)y = 0$, where $C(x)y = \mu_k e$. The extra step, in which the starting point is computed by performing a primal-dual Newton step from (x_k, y_k) , is simply included to generate a value that is already close to first order critical, and the stopping tolerances are specially chosen to encourage this. Indeed we have the following asymptotic convergence result.

Theorem 4.3. Suppose that $f, c \in \mathcal{C}^2$, that a subsequence $\{(x_k, y_k)\}$, $k \in \mathcal{K}$, of the practical primal-dual method converges to (x_*, y_*) satisfying second-order sufficiency conditions, that $A_A(x_*)$ is full-rank, and that $(y_*)_A > 0$. Then the starting point satisfies the inner-minimization termination test (i.e., $(x_k, y_k) = (x_k^s, y_k^s)$) for all k sufficiently large, and the whole sequence $\{(x_k, y_k)\}$ converges to (x_*, y_*) at a superlinear rate (with a Q-factor at least 1.9998).

This is a highly acceptable result, the convergence being essentially quadratic (which would correspond to a Q-factor of two—any sequence $\{\sigma_k\}$ is said to converge to σ_* with *Q-factor* at least q if $|\sigma_{k+1} - \sigma_*| \leq \gamma|\sigma_k - \sigma_*|^q$ for some $\gamma > 0$).

Primal-dual interior-point methods have the potential for both excellent theoretical and practical behaviour. There are polynomial interior-point algorithms for linear, (convex) quadratic and semi-definite programming. While it is unlikely that this is true for more general (nonconvex) problems, the barrier function globalization is most effective in practice, and the asymptotic behaviour is normally just as for the convex case. From a global perspective, it is very important that iterates are kept away from constraint boundary until near to convergence, as otherwise very slow progress

will be made—this is certainly born out in practice. Finally, while the methods we have discussed in this section have all required an interior starting point, it is possible to find one (if there is one!) by solving the “phase-one” problem to

$$\underset{(x,\gamma)}{\text{minimize}} \quad \gamma \quad \text{subject to} \quad c(x) + \gamma e \geq 0;$$

any feasible point (x, γ) for this auxiliary problem for which $\gamma < 0$ is suitable, for then $c(x) > 0$.

It is quite common in practice to replace the inequality $c_i(x) \geq 0$ by the equation $c_i(x) - s_i = 0$, and simple bound $s_i \geq 0$ on the *slack* variable s_i . This has the algebraic advantage that the inequality constraints are then all simple bounds and thus that barrier terms only appear on the diagonal of the Hessian model, but arguably the disadvantages that the dimensionality of the problem has been artificially increased, and that we now need to use some means of coping with equality constraints. We consider this latter point next.

5 SQP METHODS FOR EQUALITY CONSTRAINED OPTIMIZATION

In this final section, having already investigated very good methods for dealing with inequality constraints, we now turn our attention to the problem (4.1), in which there are only equality constraints on the variables. Of course in practice, there are frequently both equations and inequalities, and composite methods using the barrier/interior-point methods discussed in Section 4 and the SQP methods we shall consider here are often used. Alternatively, SQP methods themselves may easily be generalized to handle inequality constraints. For brevity we shall not consider such extensions further here.

5.1 Newton’s method for first-order optimality

Sequential Quadratic Programming (SQP) methods (sometimes called successive or recursive quadratic programming methods) are most naturally derived by considering the first-order necessary conditions for (4.1)—we will see where the names come from shortly. Recall at optimality we expect to have

$$g(x, y) \equiv g(x) - A^T(x)y = 0 \quad \text{and} \quad c(x) = 0. \quad (5.1)$$

This is a system of nonlinear equations in the variables x and the Lagrange multipliers y . Notice that the system is actually linear in y so that if x were known it would be straightforward to find y .

Suppose now that (x, y) is an approximation to a solution of (5.1). Then, as always, we might apply Newton’s method to try to improve (x, y) , and this leads us to construct a correction (s, w) for which

$$\begin{pmatrix} H(x, y) & -A^T(x) \\ A(x) & 0 \end{pmatrix} \begin{pmatrix} s \\ w \end{pmatrix} = - \begin{pmatrix} g(x, y) \\ c(x) \end{pmatrix}. \quad (5.2)$$

Newton’s method would then apply the same procedure to the “improved” estimate $(x_+, y_+) = (x + s, y + w)$.

There are a number of alternative formulations of (5.2). Firstly (5.2) may be written as the symmetric system of equations

$$\begin{pmatrix} H(x, y) & A^T(x) \\ A(x) & 0 \end{pmatrix} \begin{pmatrix} s \\ -w \end{pmatrix} = - \begin{pmatrix} g(x, y) \\ c(x) \end{pmatrix};$$

notice here that the coefficient matrix is indefinite because of its zero 2,2 block. Secondly, on writing $y_+ = y + w$, the equation becomes

$$\begin{pmatrix} H(x, y) & -A^T(x) \\ A(x) & 0 \end{pmatrix} \begin{pmatrix} s \\ y_+ \end{pmatrix} = - \begin{pmatrix} g(x) \\ c(x) \end{pmatrix},$$

or finally, in symmetric form,

$$\begin{pmatrix} H(x, y) & A^T(x) \\ A(x) & 0 \end{pmatrix} \begin{pmatrix} s \\ -y_+ \end{pmatrix} = - \begin{pmatrix} g(x) \\ c(x) \end{pmatrix}.$$

In practice we might prefer to approximate $H(x, y)$ by some symmetric B , and instead solve

$$\begin{pmatrix} B & A^T(x) \\ A(x) & 0 \end{pmatrix} \begin{pmatrix} s \\ -y_+ \end{pmatrix} = - \begin{pmatrix} g(x) \\ c(x) \end{pmatrix} = \begin{pmatrix} B & -A^T(x) \\ A(x) & 0 \end{pmatrix} \begin{pmatrix} s \\ y_+ \end{pmatrix}. \quad (5.3)$$

One could imagine solving these related systems by finding an LU factorization of the coefficient matrix in the unsymmetric case, or a symmetric-indefinite (a generalization of Cholesky) factorization in the symmetric case. Alternatively, if B is invertible, s and y_+ might be found successively by solving

$$A(x)B^{-1}A(x)^T y = -c + A(x)B^{-1}g \text{ and then } Bs = A(x)^T y - g$$

using symmetric factorizations of B and $A(x)B^{-1}A(x)^T$. For very large problems, iterative methods might be preferred, and here GMRES(k) or QMR, for the unsymmetric case, or MINRES or conjugate-gradients (restricted to the null-space of $A(x)$), for the symmetric case, have all been suggested. Thus there are many ways to solve the system(s) of linear equations that arise from SQP methods, and there is currently much interest in exploiting the structure in such systems to derive very efficient methods.

But where does the name “sequential quadratic programming” come from?

5.2 The Sequential Quadratic Programming iteration

A *quadratic program* is a problem involving the optimization of a quadratic function subject to a set of linear inequality and/or equality constraints. Consider the quadratic programming problem

$$\underset{s \in \mathbb{R}^n}{\text{minimize}} \quad \langle s, g(x) \rangle + \frac{1}{2} \langle s, Bs \rangle \quad \text{subject to} \quad A(x)s = -c(x). \quad (5.4)$$

Why this problem? Well, Theorem 1.2 indicates that $c(x) + A(x)s$ is a first-order (Taylor) approximation to the constraint function $c(x + s)$, while $\langle s, g(x) \rangle + \frac{1}{2} \langle s, Bs \rangle$ is potentially a second-order model of the decrease $f(x + s) - f(x)$. Thus one can argue that (5.4) gives a suitable (at least first-order) model of (4.1). An objection might be that really we should be aiming for true second-order approximations to all functions concerned, but this would lead to the significantly-harder minimization of a quadratic function subject to quadratic constraints—constraint curvature is a major obstacle.

The interesting feature of (5.4) is that it follows immediately from Theorem 1.6 that any first-order critical point of (5.4) is given by (5.3). Thus Newton-like methods for first-order optimality are equivalent to the solution of a sequence of related quadratic programs. Hence the name. Notice that if $B = H(x, y)$, solving (5.4) is actually Newton’s method for (5.1), and this suggests that B should be an approximation to the Hessian of the Lagrangian function, not the objective function. Clearly

the constraint curvature that we would have liked to have added to the linear approximations of the constraints has worked its way into the objective function!

To summarize, the basic SQP iteration is as follows.

Given (x_0, y_0) , set $k = 0$.

Until “convergence” iterate:

 Compute a suitable symmetric B_k using (x_k, y_k) .

 Find

$$s_k = \arg \min_{s \in \mathbb{R}^n} \langle s, g_k \rangle + \frac{1}{2} \langle s, B_k s \rangle \quad \text{subject to } A_k s = -c_k \quad (5.5)$$

 along with associated Lagrange multiplier estimates y_{k+1} .

 Set $x_{k+1} = x_k + s_k$ and increase k by 1.

The SQP method is both simple and fast. If $B_k = H(x_k, y_k)$, the method is Newton’s method for (5.1), and thus is quadratically convergent provided that (x_0, y_0) is sufficiently close to a first-order critical point (x_*, y_*) of (4.1) for which

$$\begin{pmatrix} H(x_*, y_*) & A^T(x_*) \\ A(x_*) & 0 \end{pmatrix}$$

is non-singular. Moreover, the method is superlinearly convergent when B_k is a “good” approximation to $H(x_k, y_k)$, and there is even no necessity that this be so for fast convergence. It should also be easy for the reader to believe that had we wanted to solve the problem (4.3) involving inequality constraints, the suitable SQP subproblem would be

$$\text{minimize } \langle s, g(x) \rangle + \frac{1}{2} \langle s, B s \rangle \quad \text{subject to } A(x)s \geq -c(x)$$

in which the nonlinear inequalities have been linearized.

But, as the reader will already have guessed, this basic iteration also has drawbacks, leading to a number of vital questions. For a start it is a Newton-like iteration, and thus may diverge from poor starting points. So how do we globalize this iteration? How should we pick B_k ? What should we do if (5.4) is unbounded from below? And precisely when is it unbounded?

The problem (5.4) only has a solution if the constraints $A(x)s = -c(x)$ are consistent. This is certainly the case if $A(x)$ is full rank, but may not be so if $A(x)$ is rank deficient—we shall consider alternatives that deal with this deficiency later. Applying Theorem 1.7 to (5.4), we deduce that any stationary point (s, y_+) satisfying (5.3) solves (5.4) only if B is positive semi-definite on the manifold $\{s : A(x)s = 0\}$ —if B is positive definite on the manifold (s, y_+) is the unique solution to the problem. If the m by n matrix $A(x)$ is full rank and the columns of $N(x)$ form a basis for the null-space of $A(x)$, it is easy to show that B being positive (semi-)definite on the manifold $\{s : A(x)s = 0\}$ is equivalent to $N(x)^T B N(x)$ being positive (semi-)definite which is in turn equivalent to the matrix

$$\begin{pmatrix} B & A^T(x) \\ A(x) & 0 \end{pmatrix}$$

(being non-singular and) having m negative eigenvalues. If B violates these assumptions, (5.4) is unbounded.

For the remainder of this section, we focus on methods to globalize the SQP iteration. And it should not surprise the reader that we shall do so by considering linesearch and trust-region schemes.

5.3 Linesearch SQP methods

The obvious way to embed the SQP step s_k within a linesearch framework is to pick $x_{k+1} = x_k + \alpha_k s_k$, where the step $\alpha_k > 0$ is chosen so that

$$\Phi(x_k + \alpha_k s_k, p_k) \text{ “<” } \Phi(x_k, p_k), \quad (5.6)$$

and where $\Phi(x, p)$ is a “suitable” merit function depending on parameters p_k . Of course it is then vital that s_k be a descent direction for $\Phi(x, p_k)$ at x_k , as otherwise there may be no α_k for which (5.6) is satisfied. As always with linesearch methods, this limits the choice of B_k , and it is usual to insist that B_k be positive definite—the reader may immediately object that this is imposing an unnatural requirement, since B_k is supposed to be approximating the (usually) indefinite matrix $H(x_k, y_k)$, and we can only sympathise with such a view!

What might a suitable merit function be? One possibility is to use the quadratic penalty function (4.2). In this case, we have the following result.

Theorem 5.1. Suppose that B_k is positive definite, and that (s_k, y_{k+1}) are the SQP search direction and its associated Lagrange multiplier estimates for the problem

$$\text{minimize } f(x) \text{ subject to } c(x) = 0 \\ x \in \mathbb{R}^n$$

at x_k . Then if x_k is not a first-order critical point, s_k is a descent direction for the quadratic penalty function $\Phi(x, \mu_k)$ at x_k whenever

$$\mu_k \leq \frac{\|c(x_k)\|_2}{\|y_{k+1}\|_2}.$$

We know that the parameter μ_k for the quadratic penalty function needs to approach zero for its minimizers to converge to those of (4.1), so Theorem 5.1 simply confirms this by suggesting how to adjust the parameter.

The quadratic penalty function has another role to play if the constraints are inconsistent. For consider the quadratic (Newton-like) model

$$\text{minimize } \langle s, g_k + A_k^T c_k / \mu_k \rangle + \frac{1}{2} \langle s, (B_k + 1/\mu_k A_k^T A_k) s \rangle \\ s \in \mathbb{R}^n$$

that might be used to compute a step s_k^Q from x_k . Stationary points of this model satisfy

$$(B_k + 1/\mu_k A_k^T A_k) s_k^Q = -(g_k + A_k^T c_k / \mu_k)$$

or, on defining $y_k^Q \stackrel{\text{def}}{=} -\mu_k^{-1}(c_k + A_k s_k^Q)$,

$$\begin{pmatrix} B_k & A_k^T \\ A_k & -\mu_k I \end{pmatrix} \begin{pmatrix} s_k^Q \\ -y_k^Q \end{pmatrix} = - \begin{pmatrix} g_k \\ c_k \end{pmatrix}. \quad (5.7)$$

But now compare this system with (5.3) that which defines the SQP step: the only difference is the vanishingly small 2,2 block $-\mu_k I$ in the coefficient matrix. While this indicates that Newton-like directions for the quadratic penalty function will become increasingly good approximations to SQP steps (and, incidentally, it can be shown that a Newton iteration for (4.2) with well chosen μ_k converges superlinearly under reasonable assumptions), the main point of the alternative (5.7) is that rank-deficiency in A_k is neutralised by the presence of 2,2 block term $-\mu_k I$. Nevertheless, the quadratic penalty function is rarely used, its place often being taken by non-differentiable exact penalty functions.

The *non-differentiable exact penalty function* is given by

$$\Phi(x, \rho) = f(x) + \rho \|c(x)\| \quad (5.8)$$

for any norm $\|\cdot\|$ and scalar $\rho > 0$. Notice that the function is non-differentiable particularly when $c(x) = 0$, the very values we hope to attain! The following result helps explain why such a function is considered so valuable.

Theorem 5.2. Suppose that $f, c \in C^2$, and that x_* is an isolated local minimizer of $f(x)$ subject to $c(x) = 0$, with corresponding Lagrange multipliers y_* . Then x_* is also an isolated local minimizer of $\Phi(x, \rho)$ provided that

$$\rho > \|y_*\|_D,$$

where the *dual norm*

$$\|y\|_D = \sup_{x \neq 0} \frac{\langle y, x \rangle}{\|x\|}.$$

Notice that the fact that ρ merely needs to be larger than some critical value for $\Phi(x, \rho)$ to be usable to try to identify solutions to (4.1) is completely different to the quadratic penalty function, for which the parameter had to take on a limiting value.

More importantly, as we now see, $\Phi(x, \rho)$ may be used as a merit function for the SQP step.

Theorem 5.3. Suppose that B_k is positive definite, and that (s_k, y_{k+1}) are the SQP search direction and its associated Lagrange multiplier estimates for the problem

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) \quad \text{subject to} \quad c(x) = 0$$

at x_k . Then if x_k is not a first-order critical point, s_k is a descent direction for the non-differentiable penalty function $\Phi(x, \mu_k)$ at x_k whenever $\rho_k \geq \|y_{k+1}\|_D$.

Once again, this theorem indicates how ρ_k needs to be adjusted for use within a linesearch SQP framework.

Thus far, everything looks perfect. We have methods for globalizing the SQP iteration, an iteration that should ultimately converge very fast. But unfortunately, it is not as simple as that. For consider the example in Figure 5.1. Here the current iterate lies close to (actually on) the constraint, the SQP step moves tangentially from it, and thus moves away as the constraint is nonlinear, but unfortunately, at the same time, the value of the objective function rises. Thus any

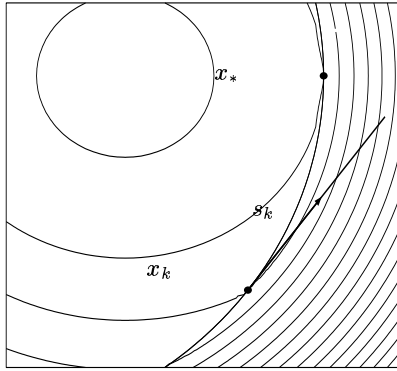


Figure 5.1: ℓ_1 non-differentiable exact penalty function ($\rho = 1$): $f(x) = 2(x_1^2 + x_2^2 - 1) - x_1$ and $c(x) = x_1^2 + x_2^2 - 1$. Solution: $x_* = (1, 0)$, $y_* = \frac{3}{2}$. The SQP direction using the optimal Hessian $H(x_*, y_*) = I$. Notice how the merit function increases at the point $x_k + s_k$.

merit function like (4.2) or (5.8) composed simply from positive combinations of the objective and (powers) of norms of constraint violations will increase after such an SQP step, and thus necessarily $\alpha_k \neq 1$ in (5.6)—worse still, this behaviour can happen arbitrarily close to the minimizer. This has the unfortunate side effect that it may happen that the expected fast convergence achievable by Newton-like methods will be thwarted by the merit function. That is, there is a serious mismatch between the global and local convergence needs of the SQP method. The fact that the merit function may prevent acceptance of the full SQP step is known as the *Maratos effect*.

The Maratos effect occurs because the curvature of the constraints is not adequately represented by linearization in the SQP model. In particular,

$$c(x_k + s_k) = O(\|s_k\|^2).$$

This suggests that we need to correct for this curvature. We may do this by computing a *second-order correction* from $x_k + s_k$, that is an extra step s_k^C for which

$$c(x_k + s_k + s_k^C) = o(\|s_k\|^2). \quad (5.9)$$

Since we do not want to destroy potential for fast convergence, we must also insist that the correction is small relative to the SQP step, and thus that

$$s_k^C = o(s_k). \quad (5.10)$$

There are a number of ways to compute a second-order correction. The first is simply to move back as quickly as possible towards the constraints. This suggests we compute a minimum (ℓ_2 -)norm solution to $c(x_k + s_k) + A(x_k + s_k)s_k^C = 0$. It is easy to check that the required solution satisfies

$$\begin{pmatrix} I & A^T(x_k + s_k) \\ A(x_k + s_k) & 0 \end{pmatrix} \begin{pmatrix} s_k^C \\ -y_{k+1}^C \end{pmatrix} = - \begin{pmatrix} 0 \\ c(x_k + s_k) \end{pmatrix}.$$

Since this requires that we re-evaluate the constraints *and* their Jacobian at $x_k + s_k$, we might hope instead to find a minimum norm solution to $c(x_k + s_k) + A(x_k)s_k^C = 0$, and thus that

$$\begin{pmatrix} I & A^T(x_k) \\ A(x_k) & 0 \end{pmatrix} \begin{pmatrix} s_k^C \\ -y_{k+1}^C \end{pmatrix} = - \begin{pmatrix} 0 \\ c(x_k + s_k) \end{pmatrix}.$$

A third amongst many other possibilities is to compute another SQP step from $x_k + s_k$, that is to compute s_k^C so that

$$\begin{pmatrix} B_k^C & A^T(x_k + s_k) \\ A(x_k + s_k) & 0 \end{pmatrix} \begin{pmatrix} s_k^C \\ -y_{k+1}^C \end{pmatrix} = - \begin{pmatrix} g(x_k + s_k) \\ c(x_k + s_k) \end{pmatrix},$$

where B_k^C is an approximation to $H(x_k + s_k, y_k^+)$. It can easily be shown that all of the above corrections satisfy (5.9)–(5.10). In Figure 5.2, we illustrate a second-order correction in action. It is

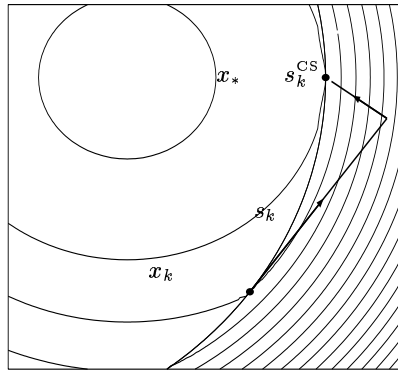


Figure 5.2: ℓ_1 non-differentiable exact penalty function ($\rho = 1$): $f(x) = 2(x_1^2 + x_2^2 - 1) - x_1$ and $c(x) = x_1^2 + x_2^2 - 1$ solution: $x_* = (1, 0)$, $y_* = \frac{3}{2}$. See that the second-order correction s_k^{CS} helps avoid the Maratos effect for the above problem with the ℓ_1 -penalty function. Notice how s_k^{CS} more than compensates for the increase in the merit function at the point $x_k + s_k$, and how much closer $x_k + s_k + s_k^{\text{CS}}$ is to x_* than is x_k .

possible to show that, under reasonable assumptions, any step $x_k + s_k + s_k^{\text{CS}}$ made up from the SQP step s_k and a second-order correction s_k^{CS} satisfying (5.9)–(5.10) will ultimately reduce (5.8). So now we can have both global and very fast asymptotic convergence at the expense of extra problem evaluations. Of course, we have stressed that a second SQP step gives a second-order correction, so another way of viewing this is to require that the merit function decreases at least every second iteration, and to tolerate non-monotonic behaviour in the interim.

5.4 Trust-region SQP methods

The main disadvantage of (at least naive) linesearch SQP methods is the unnatural requirement that B_k be positive definite. We saw the same restriction in the unconstrained case, although at least then there was some expectation that ultimately the true Hessian H_k would be positive (semi-)

definite. In the unconstrained case, indefinite model Hessians were better handled in a trust-region framework, and the same is true in the constrained case.

The obvious trust-region generalization of the basic SQP step-generation subproblem (5.4) is to find

$$s_k = \arg \min_{s \in \mathbb{R}^n} \langle s, g_k \rangle + \frac{1}{2} \langle s, B_k s \rangle \quad \text{subject to} \quad A_k s = -c_k \quad \text{and} \quad \|s\| \leq \Delta_k. \quad (5.11)$$

Since we do not require that B_k be positive definite, this allows us to use $B_k = H(x_k, y_k)$ if we so desire. However a few moments reflection should make it clear that such an approach has a serious flaw. Let Δ^{CRIT} be the least distance to the linearized constraints, i.e.

$$\Delta^{\text{CRIT}} \stackrel{\text{def}}{=} \min \|s\| \quad \text{subject to} \quad A_k s = -c_k.$$

The difficulty is that if $\Delta_k < \Delta^{\text{CRIT}}$, then there is *no solution* to the trust-region subproblem (5.11). This implies that unless $c_k = 0$, the subproblem is meaningless for all sufficiently small trust-region radius (see Figure 5.3). Thus we need to consider alternatives. In this section, we shall review the

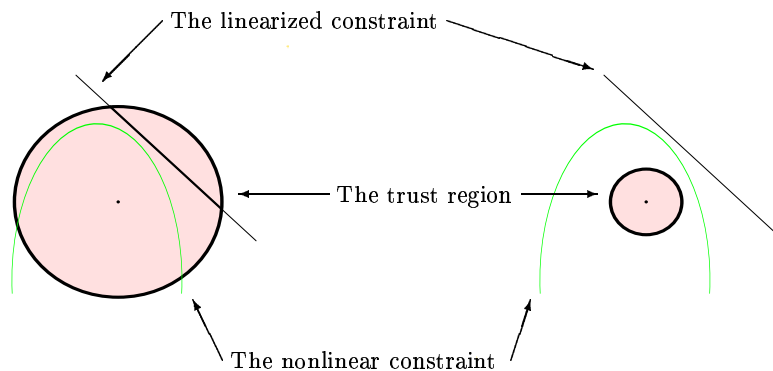


Figure 5.3: The intersection between the linearization of a nonlinear constraint and a spherical trust region. In the left figure, the trust-region radius is sufficiently large for the trust region and the linearized constraint to intersect. This is not so for the smaller trust region illustrated in the right figure.

$S\ell_p$ QP method of Fletcher, the composite step SQP methods due to Vardi, to Byrd and Omojokun, and to Celis, Dennis and Tapia, and the filter-SQP approach of Fletcher and Leyffer.

5.4.1 The $S\ell_p$ QP method

Our first trust-region approach is to try to minimize the ℓ_p -(exact) penalty function

$$\Phi(x, \rho) = f(x) + \rho \|c(x)\|_p \quad (5.12)$$

for sufficiently large $\rho > 0$ and some ℓ_p norm ($1 \leq p \leq \infty$). We saw in Section 5.3 that feasible minimizers of (5.12) may be solutions to (4.1) so long as $\rho > 0$ is large enough. Of course, as $\Phi(x, \rho)$ is non-differentiable, we cannot simply apply one of the unconstrained trust-region methods discussed in Section 3, but must instead build a specialized method.

Since we are discussing trust-region methods, a suitable model problem is the ℓ_p QP

$$\underset{s \in \mathbb{R}^n}{\text{minimize}} \quad f_k + \langle s, g_k \rangle + \frac{1}{2} \langle s, B_k s \rangle + \rho \|c_k + A_k s\|_p \quad \text{subject to} \quad \|s\| \leq \Delta_k.$$

This has the major advantage that the model problem is always consistent, since now the only constraint is the trust-region bound. In addition, when ρ and Δ_k are large enough, it can be shown that the model minimizer is the SQP direction so long as $A_k s = -c_k$ is consistent. Moreover, when the norms are polyhedral (e.g., the ℓ_1 or ℓ_∞ norms), ℓ_p QP is equivalent to a quadratic program.

To see this, consider for example the ℓ_1 QP model problem with an ℓ_∞ trust region

$$\underset{s \in \mathbb{R}^n}{\text{minimize}} \quad \langle s, g_k \rangle + \frac{1}{2} \langle s, B_k s \rangle + \rho \|c_k + A_k s\|_1 \quad \text{subject to} \quad \|s\|_\infty \leq \Delta_k.$$

But we can always write

$$c_k + A_k s = u - v, \quad \text{where} \quad (u, v) \geq 0.$$

Hence the ℓ_1 QP subproblem is equivalent to the quadratic program

$$\begin{aligned} & \underset{s \in \mathbb{R}^n, u, v \in \mathbb{R}^m}{\text{minimize}} && \langle s, g_k \rangle + \frac{1}{2} \langle s, B_k s \rangle + \rho \langle e, u + v \rangle \\ & \text{subject to} && A_k s - u + v = -c_k \\ & && u \geq 0, \quad v \geq 0 \\ & \text{and} && -\Delta_k e \leq s \leq \Delta_k e. \end{aligned}$$

Notice that the QP involves inequality constraints, but there are good methods (especially of the interior-point variety) for solving such problems. In particular, it is possible to exploit the structure of the u and v variables.

In order to develop a practical $S\ell_1$ QP method, it should not surprise the reader that we need to ensure that every step we generate achieves as much reduction in the model $f_k + \langle s, g_k \rangle + \frac{1}{2} \langle s, B_k s \rangle + \rho \|c_k + A_k s\|_p$ as would have been achieved at a Cauchy point. One such Cauchy point requires the solution to ℓ_1 LP model

$$\underset{s \in \mathbb{R}^n}{\text{minimize}} \quad \langle s, g_k \rangle + \rho \|c_k + A_k s\|_1 \quad \text{subject to} \quad \|s\|_\infty \leq \Delta_k,$$

which may be reformulated as a linear program. Fortunately approximate solutions to both ℓ_1 LP and ℓ_1 QP subproblems suffice. In practice it is also important to adjust ρ as the method progresses so as to ensure that ρ is larger than the (as yet unknown) $\|y_*\|_D$, and this may be achieved by using the available Lagrange multiplier estimates y_k . Such a scheme is globally convergent, but there is still a need for a second-order correction to prevent the Maratos effect and thus allow fast asymptotic convergence. If $c(x) = 0$ are inconsistent, the method converges to (locally) least value of the infeasibility $\|c(x)\|$ provided $\rho \rightarrow \infty$.

The alert reader will have noticed that in this section we have replaced the ℓ_2 trust-region of the unconstrained trust-region method by a box or ℓ_∞ trust-region. The reason for this apparent lack of consistency is that minimizing a quadratic subject to linear constraints *and* an additional quadratic trust-region is too hard. On the other hand, adding box-constraints does not increase the complexity of the resulting (quadratic programming) trust-region subproblem.

5.4.2 Composite-step methods

Another approach to avoid the difficulties caused by inconsistent QP subproblems is to separate the computation of the step into two stages. The aim of a *composite-step* method is to find

$$s_k = n_k + t_k,$$

where the *normal step* n_k moves towards feasibility of the linearized constraints (within the trust region), while the *tangential step* t_k reduces the model objective function (again within the trust-region) without sacrificing feasibility obtained from n_k . Of course since the normal step is solely concerned with feasibility, the model objective may get worse, and indeed it may not recover during the tangential step. The fact that the tangential step is required to maintain any gains in (linearized) feasibility achieved during the normal step implies that

$$A_k(n_k + t_k) = A_k n_k \quad \text{and hence that} \quad A_k t_k = 0.$$

We illustrate possible normal and tangential steps in Figure 5.4.

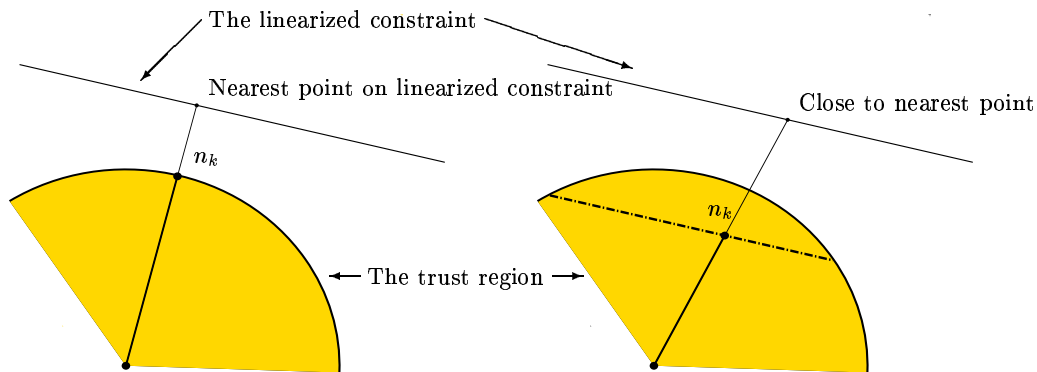


Figure 5.4: Computing the normal step. The left-hand figure shows the largest possible normal step. The right-hand figure illustrates a shorter normal step n , and the freedom this then allows for the tangential step—any point on the dotted line is a potential tangential step.

5.4.2.1 Constraint relaxation—Vardi’s method

Vardi’s approach is an early composite-step method. The normal step is found by relaxing the requirement

$$A_k s = -c_k \quad \text{and} \quad \|s\| \leq \Delta_k$$

to

$$A_k n = -\sigma_k c_k \quad \text{and} \quad \|n\| \leq \Delta_k,$$

where $\sigma_k \in [0, 1]$ is small enough so that there is a feasible n_k . Clearly $s = 0$ is feasible if $\sigma_k = 0$, and the largest possible σ_{\max} may be found by computing

$$\max_{\sigma \in (0,1]} \left[\min_{\|s\| \leq \Delta_k} \|A_k s + \sigma c_k\| = 0 \right].$$

In practice, some value between zero and σ_{\max} is chosen, since this gives some “elbow-room” in which to compute the tangential step. The main defect with the approach is that there may be no normal step if the linearized constraints are inconsistent.

Once a normal step has been determined, the tangential step is computed as the

$$\begin{aligned} & \text{(approximate) } \arg \min_{t \in \mathbb{R}^n} \quad \langle t, g_k + B_k n_k \rangle + \frac{1}{2} \langle t, B_k t \rangle \\ & \text{subject to} \quad A_k t = 0 \text{ and } \|n_k + t\| \leq \Delta_k. \end{aligned}$$

Although of historical interest, the method has been effectively superseded by the Byrd–Omojokun approach we describe next.

5.4.2.2 Constraint reduction—the Byrd–Omojokun method

The Byrd–Omojokun method aims to cope with the inconsistency issue that afflicts Vardi’s approach. Rather than relaxing the constraints, the normal step is now computed as

$$\text{approximately minimize } \|A_k n + c_k\| \text{ subject to } \|n\| \leq \Delta_k,$$

in order to achieve a reasonable improvement in linearized infeasibility that is consistent with the trust-region. The tangential step is then computed exactly as in Vardi’s method.

An important aspect is that it is possible to use the conjugate gradient method to solve both subproblems. This provides Cauchy points in both cases and allows the method to be used to solve large problems. The method has been shown to be globally convergent (under reasonable assumptions) using an ℓ_2 merit function, and is the basis of the successful KNITRO software package.

5.4.2.3 Constraint lumping—the Celis–Dennis–Tapia method

A third method which might be considered to be of the composite-step variety is that due to Celis, Dennis and Tapia. In this approach, the requirement that $A_k s = -c_k$ is replaced by requiring that

$$\|A_k s + c_k\| \leq \sigma_k$$

for some $\sigma_k \in [0, \|c_k\|]$. The value of σ_k is chosen so that the normal step n_k satisfies

$$\|A_k n + c_k\| \leq \sigma_k \text{ and } \|n\| \leq \Delta_k.$$

Having found a suitable normal step, the tangential step is found as an

$$\begin{aligned} & \text{(approximate) } \arg \min_{t \in \mathbb{R}^n} \quad \langle t, g_k + B_k n_k \rangle + \frac{1}{2} \langle t, B_k t \rangle \\ & \text{subject to} \quad \|A_k t + A_k n_k + c_k\| \leq \sigma_k \text{ and } \|t + n_k\| \leq \Delta_k. \end{aligned}$$

While finding a suitable σ_k is inconvenient, the real Achilles’ heel of this approach is that the tangential step subproblem is (much) harder than those we have considered so far. If the ℓ_2 -norm is used for the constraints, we need to find the minimizer of a quadratic objective within the intersection of two “spherical” regions. Unlike the case involving a single sphere (recall Section 3.5.1), it is not known if there is an efficient algorithm in the two-sphere case. Alternatively, if polyhedral (ℓ_1 or ℓ_∞) norms are used and B_k is indefinite, the subproblem becomes a non-convex quadratic program for which there is unlikely to be an efficient general-purpose algorithm—in the special case where B_k is positive semi-definite and the ℓ_∞ norm is used, the subproblem is a convex QP. For this reason, the Celis–Dennis–Tapia approach is rarely used in practice.

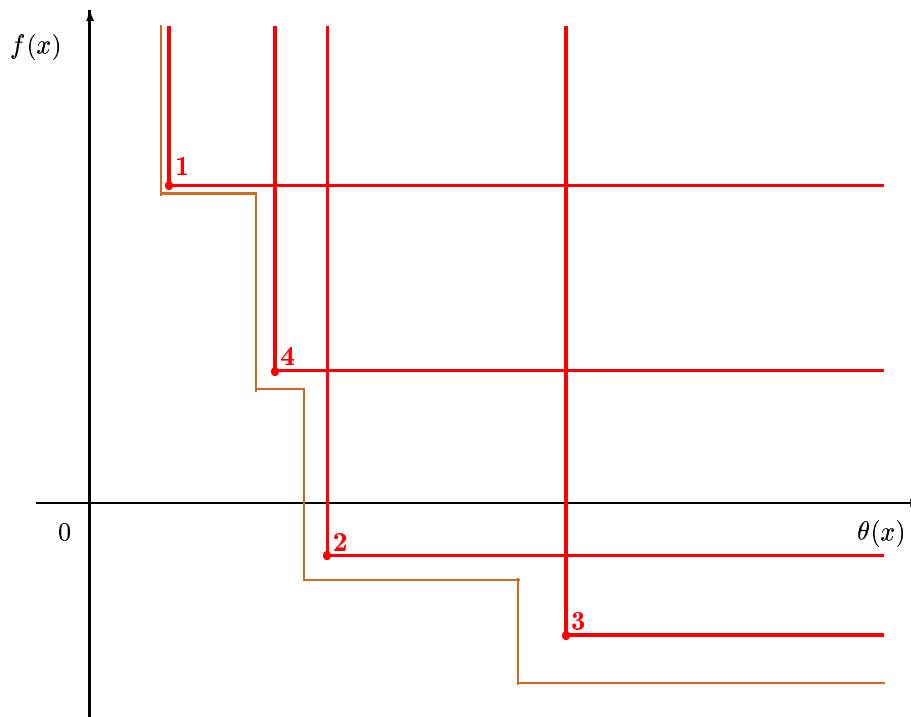


Figure 5.5: A filter with four entries.

5.4.3 Filter methods

The last SQP method we shall consider is the most recent. The approach taken is quite radical in that, unlike all of the methods we have considered so far, it makes no use of a merit function to force global convergence. The main objection to merit functions is that they depend, to a large degree, on arbitrary or a-priori unknown parameters. A secondary objection is that they tend to be overly conservative in accepting promising potential iterates. But if we wish to avoid merit functions, we need some other device to encourage convergence. The new idea is to use a “filter”

Let $\theta(x) = \|c(x)\|$ be some norm of the constraint violation at x . A *filter* is a set of pairs $\{(\theta_k, f_k)\}$ of violations and objective values such that no member dominates another, i.e., it does not happen that

$$\theta_i \ll \theta_j \text{ and } f_i \ll f_j$$

for any pair of filter points $i \neq j$ —the “ \ll ” here informally means “very slightly smaller than”. We illustrate a filter in Figure 5.5. A potential new entry to the “north-east” of any of the existing filter entries would not be permitted, and the forbidden region is the intersection of the solid horizontal and vertical lines emanating to the right and above each filter point. For theoretical reasons (akin to requiring sufficient decrease), we slightly enlarge the forbidden region by putting a small margin around each filter point, and this is illustrated in the figure by the dotted lines.

And now it is clear how to use a filter. Any potential SQP (or other) iterate $x_k + s_k$ will immediately be rejected if it lies in the forbidden filter region accumulated during the previous k iterations. This may be embedded in a trust-region framework, and a typical iteration might be as follows:

If possible find

$$s_k = \arg \min_{s \in \mathbb{R}^n} \langle s, g_k \rangle + \frac{1}{2} \langle s, B_k s \rangle \text{ subject to } A_k s = -c_k \text{ and } \|s\| \leq \Delta_k,$$

but otherwise, find s_k such that

$$\theta(x_k + s_k) < \theta_i \text{ for all } i \leq k.$$

If $x_k + s_k$ is “acceptable” for the filter, set $x_{k+1} = x_k + s_k$

and possibly add $(f(x_k + s_k), \theta(x_k + s_k))$ to the filter,

“prune” the filter, and increase Δ_k .

Otherwise reduce Δ_k and try again.

A few words of explanation are needed. The trust-region and linearized constraints will always be compatible if c_k is small enough so long as they are at $c(x) = 0$. Thus if the trust-region subproblem is incompatible, one remedy is simply to move closer to the constraints. This is known as a *restoration* step. By “pruning” the filter, we mean that a new point may completely dominate one or more existing filter points and, in this case, the dominated entry may be removed without altering the filter. For example, if a new entry were accepted to the “south-west” of point 4 in our figure, point 4 would be pruned.

While the basic filter idea is rather simple, in practice, it is significantly more complicated than this. In particular, there are theoretical reasons why some points that are acceptable to the filter should still be rejected if any decrease in the SQP model of the objective function is far from realized in practice.

CONCLUSION

We hope we have conveyed the impression that research into the design, convergence and implementation of algorithms for nonlinear optimization is an exciting and expanding area. We have only been able to outline the developments in the field, and have made no attempt to survey the vast literature that has built up over the last 50 years. Current algorithms for specialized problems like linear and quadratic programming and unconstrained optimization are well capable of solving problems involving millions of unknowns (and constraints), while those for generally constrained optimization routinely solve problems in the tens and, perhaps even, hundreds of thousands of unknowns and constraints. The next big goal is to be able to design algorithms that have some hope of finding global optima for large problems, the current state-of-the-art being for problems with tens or hundreds of unknowns. Clearly closing the gap between local and global optimization has some way to go!

APPENDIX A - SEMINAL BOOKS AND PAPERS

The following books and papers are classics in the field. Although many of them cover topics outside the material we have described, they are all worth reading. This section constitutes a personal view

of the most significant papers in the area. It is not meant to be a complete bibliography.

General text books

There are a large number of text books devoted to nonlinear (and even more for linear) programming. Those we find most useful and which emphasize practical methods are

J. Dennis and R. Schnabel, “Numerical Methods for Unconstrained Optimization and Non-linear Equations”, (republished by) SIAM (Classics in Applied Mathematics 16) (1996),

R. Fletcher, “Practical Methods of Optimization”, 2nd edition Wiley (1987), (republished in paperback 2000),

P. Gill, W. Murray and M. Wright, “Practical Optimization”, Academic Press (1981), and

J. Nocedal and S. Wright, “Numerical Optimization”, Springer Verlag (1999).

The first of these concentrates on unconstrained optimization, while the remainder cover (continuous) optimization in general.

Early quasi-Newton methods

These methods were introduced by

W. Davidon, “Variable metric method for minimization”, manuscript (1958), finally published *SIAM J. Optimization* **1** (1991) 1:17,

and championed by

R. Fletcher and M. J. D. Powell, “A rapidly convergent descent method for minimization”, *Computer J.* (1963) 163:168.

Although the so-called DFP method has been superseded by the more reliable BFGS method, it paved the way for a number of classes of important updates.

More modern quasi-Newton methods

Coincidentally, all of the papers

C. G. Broyden, “The convergence of a class of double-rank minimization algorithms”, *J. Inst. Math. Appls.*, **6** (1970) 76:90,

R. Fletcher, “A new approach to variable metric algorithms”, *Computer J.* (1970) **13** (1970) 317:322,

D. Goldfarb, “A family of variable metric methods derived by variational means”, *Math. Computation* **24** (1970) 23:26, and

D. F. Shanno, “Conditioning of quasi-Newton methods for function minimization”, *Math. Computation* **24** (1970) 647:657

appeared in the same year. The aptly-named BFGS method has stood the test of time well, and is still regarded as possibly the best secant updating formula.

Quasi-Newton methods for large problems

Limited memory methods are secant-updating methods that discard old information so as to reduce the amount of storage required when solving large problems. The methods first appeared in

J. Nocedal, “Updating quasi-Newton matrices with limited storage”, *Math. Computation* **35** (1980) 773:782, and

A. Buckley and A. Lenir, “QN-like variable storage conjugate gradients”, *Math. Programming* **27** (1983) 155:175.

Secant updating formulae proved to be less useful for large-scale computation, but a successful generalization, applicable to what are known as partially separable functions, was pioneered by

A. Griewank and Ph. Toint, “Partitioned variable metric updates for large structured optimization problems”, *Numerische Mathematik* **39** (1982) 119:137, see also 429:448, as well as

A. Griewank and Ph. Toint, “On the unconstrained optimization of partially separable functions”, in *Nonlinear Optimization 1981* (Powell, M., ed.) Academic Press (1982)

Conjugate gradient methods for large problems

Generalizations of Conjugate Gradient methods for non-quadratic minimization were originally proposed by

R. Fletcher and C. M. Reeves, “Function minimization by conjugate gradients”, *Computer J.* (1964) 149:154, and

E. Polak and G. Ribière, “Note sur la convergence de méthodes de directions conjuguées”, *Revue Française d’informatique et de recherche opérationnelle* **16** (1969) 35:43.

An alternative is to attempt to solve the (linear) Newton system by a conjugate-gradient like method. Suitable methods for terminating such a procedure while still maintaining fast convergence were proposed by

R. S. Dembo and T. Steihaug, “Truncated-Newton algorithms for large-scale unconstrained optimization”, *Math. Programming* **26** (1983) 190:212.

Non-monotone methods

While it is usual to think of requiring that the objective function decreases at every iteration, this is not actually necessary for convergence so long as there is some overall downward trend. The first method along these lines was by

L. Grippo, F. Lampariello and S. Lucidi, “A nonmonotone line search technique for Newton’s method”, *SIAM J. Num. Anal.*, **23** (1986) 707:716.

Trust-region methods

The earliest methods that might be regarded as trust-region methods are those by

K. Levenberg, “A method for the solution of certain problems in least squares”, *Quarterly J. Appl. Maths*, **2** (1944) 164:168, and

D. Marquardt, “An algorithm for least-squares estimation of nonlinear parameters” *SIAM J. Appl. Maths*, **11** (1963) 431:441

for the solution of nonlinear least-squares problems, although they are motivated from the perspective of modifying indefinite Hessians rather than restricting the step. Probably the first “modern” interpretation is by

S. Goldfeldt, R. Quandt and H. Trotter, “Maximization by quadratic hill-climbing”, *Econometrica*, **34** (1966) 541:551.

Certainly, the earliest proofs of convergence are given by

M. Powell, “A New Algorithm for Unconstrained Optimization”, in *Nonlinear Programming*, (Rosen, J., Mangasarian, O., and Ritter, K., eds.) Academic Press (1970),

while a good modern introduction is by

J. Moré, “Recent developments in algorithms and software for trust region methods”, in *Mathematical Programming: The State of the Art*, (Bachem, A., Grötschel, M., and Korte, B., eds.) Springer Verlag (1983).

You might want to see our book

A. Conn, N. Gould and Ph. Toint, “Trust-region methods”, SIAM (2000)

for a comprehensive history and review of the large variety of articles on trust-region methods.

Trust-region subproblems

Almost all you need to know about solving small-scale trust-region subproblems is contained in the paper

J. Moré and D. Sorensen, “Computing a trust region step”, *SIAM J. Sci. Stat. Comp.* **4** (1983) 533:572.

Likewise

T. Steihaug, “The conjugate gradient method and trust regions in large scale optimization”, *SIAM J. Num. Anal.* **20** (1983) 626:637

provides the basic truncated conjugate-gradient approach used so successfully for large-scale problems. More recently¹

N. Gould, S. Lucidi, M. Roma and Ph. Toint, “Solving the trust-region subproblem using the Lanczos method”, *SIAM J. Optimization* **9** (1999) 504:525

show how to improve on Steihaug’s approach by moving around the trust-region boundary. A particularly nice new paper by

Y. Yuan, “On the truncated conjugate-gradient method”, *Math. Programming*, **87** (2000) 561:573

proves that Steihaug’s approximation gives at least 50% of the optimal function decrease when applied to convex problems.

The Symmetric Rank-One quasi-Newton approximation

Since trust-region methods allow non-convex models, perhaps the simplest of all Hessian approximation methods, the Symmetric Rank-One update, is back in fashion. Although it is unclear who first suggested the method,

¹We would hate to claim “seminal” status for one of our own papers!

C. Broyden, “Quasi-Newton methods and their application to function minimization”, *Math. Computation* **21** (1967) 577:593

is the earliest reference that we know of. Its revival in fortune is due² to

A. Conn, N. Gould and Ph. Toint, “Convergence of quasi-Newton matrices generated by the Symmetric Rank One update” *Math. Programming*, **50** (1991) 177:196 (see also *Math. Comp.* **50** (1988) 399:430), and

R. Byrd, H. Khalfan and R. Schnabel “Analysis of a symmetric rank-one trust region method” *SIAM J. Optimization* **6** (1996) 1025:1039,

and it has now taken its place alongside the BFGS method as the pre-eminent updating formula.

More non-monotone methods

Non-monotone methods have also been proposed in the trust-region case. The basic reference here is the paper by

Ph. Toint, “A non-monotone trust-region algorithm for nonlinear optimization subject to convex constraints”, *Math. Programming*, **77** (1997) 69:94.

Barrier function methods

Although they appear to have originated in a pair of unpublished University of Oslo technical reports by K. Frisch in the mid 1950s, (logarithmic) barrier function were popularized by

A. Fiacco and G. McCormick, “The sequential unconstrained minimization technique for nonlinear programming: a primal-dual method”, *Management Science* **10** (1964) 360:366; see also *ibid* (1964) 601:617.

A full early history is given in the book

A. Fiacco and G. McCormick, “Nonlinear programming: sequential unconstrained minimization techniques” (1968), republished as *Classics in Applied Mathematics 4*, SIAM (1990).

The worsening conditioning of the Hessian was first highlighted by

F. Lootsma, “Hessian matrices of penalty functions for solving constrained optimization problems”, *Philips Research Reports*, **24** (1969) 322:331, and

W. Murray, “Analytical expressions for eigenvalues and eigenvectors of the Hessian matrices of barrier and penalty functions”, *J. Optimization Theory and Applications*, **7** (1971) 189:196,

although recent work by

M. Wright, “Ill-conditioning and computational error in interior methods for nonlinear programming”, *SIAM J. Optimization* **9** (1999) 84:111, and

S. Wright, “Effects of finite-precision arithmetic on interior-point methods for nonlinear programming”, *SIAM J. Optimization* **12** (2001) 36:78

²See previous footnote ...

demonstrates that this “defect” is far from fatal.

Interior-point methods

The interior-point revolution was started by

N. Karmarkar, “A new polynomial-time algorithm for linear programming”, *Combinatorica* **4** (1984) 373:395.

It did not take long for

P. Gill, W. Murray, M. Saunders, J. Tomlin and M. Wright, “On projected Newton barrier methods for linear programming and an equivalence to Karmarkar’s projective method”, *Math. Programming*, **36** (1986) 183:209

to realize that this radical “new” approach was actually something that nonlinear programmers had tried (but, most unfortunately, discarded) in the past.

SQP methods

The first SQP method was proposed in the overlooked 1963 Harvard Master’s thesis of R. Wilson. The generic linesearch SQP method is that of

B. Pschenichny, “Algorithms for general problems of mathematical programming”, *Kibernetika*, **6** (1970) 120:125,

while there is a much larger variety of trust-region SQP methods, principally because of the constraint incompatibility issue.

Merit functions for SQP

The first use of an exact penalty function to globalize the SQP method was by

S. Han, “A globally convergent method for nonlinear programming”, *J. Optimization Theory and Applies*, **22** (1977) 297:309, and

M. Powell, “A fast algorithm for nonlinearly constrained optimization calculations”, in *Numerical Analysis, Dundee 1977* (G. Watson, ed) Springer Verlag (1978) 144:157.

The fact that such a merit function may prevent full SQP steps was observed N. Maratos in his 1978 U. of London Ph. D. thesis, while methods for combating the Maratos effect were subsequently proposed by

R. Fletcher, “Second-order corrections for non-differentiable optimization”, in *Numerical Analysis, Dundee 1981* (G. Watson, ed) Springer Verlag (1982) 85:114, and

R. Chamberlain, M. Powell, C. Lemaréchal, and H. Pedersen, “The watchdog technique for forcing convergence in algorithms for constrained optimization”, *Math. Programming Studies*, **16** (1982) 1:17.

An SQP method that avoids the need for a merit function altogether by staying feasible is given by

J. Bonnans, E. Panier, A. Tits, and J. Zhou, “Avoiding the Maratos effect by means of a nonmonotone linesearch II. Inequality constrained problems—feasible iterates”, *SIAM J. Num. Anal.*, **29** (1992) 1187:1202.

Hessian approximations

There is a vast literature on suitable Hessian approximations for use in SQP methods. Rather than point at individual papers, a good place to start is

P. Boggs and J. Tolle, “Sequential quadratic programming”, *Acta Numerica* **4** (1995) 1:51,

but see also our paper

N. Gould and Ph. Toint, “SQP methods for large-scale nonlinear programming”, in *System modelling and optimization, methods, theory and applications* (M. Powell and S. Scholtes, eds.) Kluwer (2000) 149:178.

Trust-region SQP methods

Since the trust-region and the linearized constraints may be incompatible, almost all trust-region SQP methods modify the basic SQP method in some way. The $S\ell_1$ QP method is due to

R. Fletcher, “A model algorithm for composite non-differentiable optimization problems”, *Math. Programming Studies*, **17** (1982) 67:76.

Methods that relax the constraints include those proposed by

A. Vardi, “A trust region algorithm for equality constrained minimization: convergence properties and implementation”, *SIAM J. Num. Anal.*, **22** (1985) 575:591, and

M. Celis, J. Dennis and R. Tapia, “A trust region strategy for nonlinear equality constrained optimization”, in *Numerical Optimization 1984* (P. Boggs, R. Byrd and R. Schnabel, eds), SIAM (1985) 71:82,

as well as a method that appeared in the 1989 U. of Colorado at Boulder Ph. D. thesis of E. Omojokun, supervised by R. Byrd. The Filter-SQP approach may be found in³

R. Fletcher and S. Leyffer, “Nonlinear programming without a penalty function”, *Math. Programming*, **91** (2002) 239:269.

Modern methods for nonlinear programming

Many modern methods for nonlinearly constrained optimization tend to be SQP-interior-point hybrids. A good example is due to

R. Byrd, J. Gilbert and J. Nocedal, “A trust region method based on interior point techniques for nonlinear programming”, *Math. Programming A* **89** (2000) 149:185,

and forms the basis for the excellent KNITRO package.

³Once again, see previous footnote ...

APPENDIX B - OPTIMIZATION RESOURCES ON THE WORLD-WIDE-WEB

B.1 Answering questions on the web

A good starting point for finding out more about optimization are the two lists of Frequently Asked Questions (FAQs) on optimization. The Linear Programming FAQ,

www-unix.mcs.anl.gov/otc/Guide/faq/linear-programming-faq.html ,

is dedicated to question on linear optimization problems as well as certain aspects of mixed integer linear programming. The Nonlinear Programming FAQ,

www-unix.mcs.anl.gov/otc/Guide/faq/nonlinear-programming-faq.html ,

offers a concise introduction to nonlinear optimization. The NEOS guide,

www-fp.mcs.anl.gov/otc/Guide ,

provides an overview of optimization and the solvers available. It contains the optimization tree,

www-fp.mcs.anl.gov/otc/Guide/OptWeb ,

a dichotomy of optimization problems. Both sites are maintained by the Optimization Technology Center

www.ece.nwu.edu/OTC ,

a loose collaboration between Argonne National Laboratory and Northwestern University in the USA.

Hans Mittelmann of Arizona State University maintains a decision tree for optimization software,

plato.la.asu.edu/guide.html ,

and he also provides a useful set of benchmarks for optimization software,

plato.la.asu.edu/bench.html .

Harvey Greenberg's Mathematical Programming Glossary,

www.cudenver.edu/hgreenbe/glossary/glossary.html

contains brief definitions of commonly used expressions in optimization and operations research. The usenet newsgroup

`sci.op-research`

is dedicated to answering questions on optimization and operations research. Brian Borchers edits a weekly digest of postings to it. You can receive the digest by sending an email to

listserv@listserv.okstate.edu

with the message

SUBSCRIBE ORCS-L Your Name .

B.2 Solving optimization problems on the web

B.2.1 The NEOS server

Probably the most important and useful optimization site on the web is the NEOS server⁴ at

`www-neos.mcs.anl.gov/neos`

which allows you to solve optimization problems over the internet. NEOS handles several thousand (!) submissions per week. The server provides a wide choice of state-of-the-art optimization software which can be used remotely without the need to install or maintain any software.

The problems should preferably be formulated in a modelling language such as AMPL⁵ or GAMS⁶ (see Section B.2.3). However, some solvers also accept problem descriptions in other formats such as C or fortran source code or the verbose linear programming MPS format.

There are a number of solvers implementing algorithms for nonlinearly constrained optimization problems. Most are hybrids, and thus capable of handling both equality and inequality constraints. There are at least three interior point solvers (see Section 4).

KNITRO (with a silent “K”), is a primal-dual interior-point method which uses trust regions.

LOQO is based on an infeasible primal-dual interior-point method. It uses a linesearch and a version of a filter to enforce global convergence.

MOSEK can only be used to solve *convex* large-scale smooth nonlinear optimization problems. It does not work for *nonconvex* problems.

There are at least three solvers implementing SQP algorithms (see Section 5).

DONLP2 implements a linesearch SQP algorithm with an exact non-differentiable ℓ_1 -penalty function as a merit function. It uses dense linear algebra.

FILTER implements a trust-region SQP algorithm which is suitable for solving large nonlinearly constrained problems with small degrees of freedom. It uses a filter (see Section 5.4.3) to promote global convergence.

SNOPT implements a linesearch SQP algorithm which uses an augmented Lagrangian as a merit function. It maintains a positive definite limited memory approximation of the Hessian of the Lagrangian.

There is also a range of other solvers not covered in this article.

CONOPT is a feasible path method based on the generalized reduced gradient algorithm.

LANCELOT implements an augmented Lagrangian algorithm. It uses a trust-region to promote global convergence.

MINOS implements a sequential linearly constrained algorithm. Steplength control is heuristic (for want of a suitable merit function), but superlinear convergence is often achieved.

⁴J. Czyzyk, M. Mesnier and J. Moré. The NEOS server. *IEEE Journal on Computational Science and Engineering*, 5:68–75, 1998.

⁵R. Fourer, D. Gay and B. Kernighan. *AMPL: A modelling Language for Mathematical Programming*. Boyd & Fraser Publishing Company, Massachusetts, 1993.

⁶A. Brooke, D. Kendrick, A. Meeraus and R. Raman. *GAMS A user's guide*. GAMS Developments Corporation, 1217 Potomac Street, N.W., Washington DC 20007, USA, December 1998.

PATHNLP finds stationary points for the nonlinear problem by solving the Karush-Kuhn-Tucker conditions (see Theorems 1.6 and 1.8), written as a mixed complementarity problem, using the PATH solver.

Consult the NEOS guide (see Section B.1) for appropriate contacts.

A wide range of other optimization problems can also be solved such as semi-infinite optimization, mixed integer linear and nonlinear optimization, semidefinite optimization, complementarity problems, non-differentiable optimization, and unconstrained and stochastic optimization problems. The fact that the server maintains state-of-the-art optimization software makes it suitable for medium to large scale applications.

Users with their own copy of the modelling systems AMPL or GAMS can even invoke the NEOS solvers out of their local AMPL or GAMS session using KESTREL,

`www-neos.mcs.anl.gov/neos/kestrel.html` .

This is very convenient as it makes it possible to post- or pre-process the models using a local copy of the modelling tool.

B.2.2 Other online solvers

The system `www-Nimbus`, from

`nimbus.mit.jyu.fi` ,

is designed to solve (small) multi-objective optimization problems. It consists of a sequence of menus to input the multi-objective problem as well as some facilities for displaying the solution. It requires the user to interactively guide the optimization and requires some familiarity with multi-objective terminology. An online tutorial guides the user through the process. Certain topology optimization problems can be solved at

`www.topopt.dtu.dk`

The input is via a GUI and the solution is also display graphically. The system `Baron`,

`archimedes.scs.uiuc.edu/baron/availability.html` ,

allows the solution of small global optimization problems online.

B.2.3 Useful sites for modelling problems prior to online solution

AMPL (A Mathematical Programming Language)

`www.ampl.com`

is a modelling language for optimization problems. The site lists extensions to the book, allows the solution of example models and contains a list of available solvers. Further AMPL models can be found at the following sites:

NLP models by Bob Vanderbei:

`www.sor.princeton.edu/~rvdb/ampl/nlmodels` .

MINLP and MPEC models by Sven Leyffer:

`www.maths.dundee.ac.uk/~sleyffer/MacMINLP` and

`www.maths.dundee.ac.uk/~sleyffer/MacMPEC` .

The COPS collection of Jorge Moré:

`www-unix.mcs.anl.gov/~more/cops` .

These sites are especially useful to help with your own modelling exercises.

GAMS (the General Algebraic Modelling System)

www.gams.com

is another modelling language. The site contains documentation on GAMS and some example models. More GAMS models can be found on the GAMS-world pages. These are sites, dedicated to important modelling areas, see

www.gamsworld.org.

It also offers a translation service from one modelling language to another.

Recently, optimization solvers have also been interfaced to matlab at

tomlab.biz/.

B.2.4 Free optimization software

An extension of MPS to nonlinear optimization, SIF (standard input format), can be used to model optimization problems. The reference document can be found at

www.numerical.rl.ac.uk/lancelot/sif/sifhtml.html.

A collection of optimization problems in SIF is available at CUTEr can be found via

www.cse.clrc.ac.uk/Activity/CUTEr.

Two solvers, LANCELOT

www.cse.clrc.ac.uk/Activity/LANCELOT

and GALAHAD

www.cse.clrc.ac.uk/Activity/GALAHAD

are available freely for non-commercial users.

AMPL and some solvers are also available freely in limited size student versions, which allow the solution of problems of up to 300 variables and constraints, see

netlib.bell-labs.com/netlib/ampl/student/.

B.3 Optimization reports on the web

Optimization online,

www.optimization-online.org,

is an e-print site for papers on optimization. It is sponsored by the Mathematical Programming Society. It allows you to search for preprints on certain subjects. A monthly digest summarizes all monthly submissions.

The two main optimization journals, Mathematical Programming and SIAM Journal on Optimization maintain free sites with access to titles and abstracts, see

link.springer.de/link/service/journals/10107/,

and

www.siam.org/journals/siopt/siopt.htm.

APPENDIX C - SKETCHES OF PROOFS

Theorems 1.1 and 1.2 can be found in any good book on analysis. Theorem 1.1 follows directly by considering the remainders of truncated Taylor expansions of the univariate function $f(x + \alpha s)$ with $\alpha \in [0, 1]$, while Theorem 1.2 uses the Newton formula

$$F(x + s) = F(x) + \int_0^1 \nabla_x F(x + \alpha s) s d\alpha.$$

Proof of Theorem 1.3

Suppose otherwise, that $g(x_*) \neq 0$. A Taylor expansion in the direction $-g(x_*)$ gives

$$f(x_* - \alpha g(x_*)) = f(x_*) - \alpha \|g(x_*)\|^2 + O(\alpha^2).$$

For sufficiently small α , $\frac{1}{2}\alpha \|g(x_*)\|^2 \geq O(\alpha^2)$, and thus

$$f(x_* - \alpha g(x_*)) \leq f(x_*) - \frac{1}{2}\alpha \|g(x_*)\|^2 < f(x_*).$$

This contradicts the hypothesis that x_* is a local minimizer.

Proof of Theorem 1.4

Again, suppose otherwise that $\langle s, H(x_*)s \rangle < 0$. A Taylor expansion in the direction s gives

$$f(x_* + \alpha s) = f(x_*) + \frac{1}{2}\alpha^2 \langle s, H(x_*)s \rangle + O(\alpha^3),$$

since $g(x_*) = 0$. For sufficiently small α , $-\frac{1}{4}\alpha^2 \langle s, H(x_*)s \rangle \geq O(\alpha^3)$, and thus

$$f(x_* + \alpha s) \leq f(x_*) + \frac{1}{4}\alpha^2 \langle s, H(x_*)s \rangle < f(x_*).$$

Once again, this contradicts the hypothesis that x_* is a local minimizer.

Proof of Theorem 1.5

By continuity $H(x)$ is positive definite for all x in a open ball \mathcal{N} around x_* . The generalized mean value theorem then says that if $x_* + s \in \mathcal{N}$, there is a value z between the points x_* and $x_* + s$ for which

$$f(x_* + s) = f(x_*) + \langle s, g(x_*) \rangle + \frac{1}{2} \langle s, H(z)s \rangle = f(x_*) + \frac{1}{2} \langle s, H(z)s \rangle > f(x_*)$$

for all nonzero s , and thus x_* is an isolated local minimizer.

Proof of Theorem 1.6

We consider feasible perturbations about x_* . Consider a vector valued C^2 (C^3 for Theorem 1.7) function $x(\alpha)$ of the scalar α for which $x(0) = x_*$ and $c(x(\alpha)) = 0$. (The constraint qualification is that all such feasible perturbations are of this form). We may then write

$$x(\alpha) = x_* + \alpha s + \frac{1}{2}\alpha^2 p + O(\alpha^3) \tag{C.1}$$

and we require that

$$\begin{aligned} 0 &= c_i(x(\alpha)) = c_i(x_* + \alpha s + \frac{1}{2}\alpha^2 p + O(\alpha^3)) \\ &= c_i(x_*) + \langle a_i(x_*), \alpha s + \frac{1}{2}\alpha^2 p \rangle + \frac{1}{2}\alpha^2 \langle s, H_i(x_*)s \rangle + O(\alpha^3) \\ &= \alpha \langle a_i(x_*), s \rangle + \frac{1}{2}\alpha^2 (\langle a_i(x_*), p \rangle + \langle s, H_i(x_*)s \rangle) + O(\alpha^3) \end{aligned}$$

using Taylor's theorem. Matching similar asymptotic terms, this implies that for such a feasible perturbation

$$A(x_*)s = 0 \tag{C.2}$$

and

$$\langle a_i(x_*), p \rangle + \langle s, H_i(x_*)s \rangle = 0 \tag{C.3}$$

for all $i = 1, \dots, m$. Now consider the objective function

$$\begin{aligned} f(x(\alpha)) &= f(x_* + \alpha s + \frac{1}{2}\alpha^2 p + O(\alpha^3)) \\ &= f(x_*) + \langle g(x_*), \alpha s + \frac{1}{2}\alpha^2 p \rangle + \frac{1}{2}\alpha^2 \langle s, H(x_*)s \rangle + O(\alpha^3) \\ &= f(x_*) + \alpha \langle g(x_*), s \rangle + \frac{1}{2}\alpha^2 (\langle g(x_*), p \rangle + \langle s, H(x_*)s \rangle) + O(\alpha^3). \end{aligned} \tag{C.4}$$

This function is unconstrained along $x(\alpha)$, so we may deduce, as in Theorem 1.3, that

$$\langle g(x_*), s \rangle = 0 \text{ for all } s \text{ such that } A(x_*)s = 0. \tag{C.5}$$

If we let S be a basis for the null-space of $A(x_*)$, we may write

$$g(x_*) = A^T(x_*)y_* + Sz_* \tag{C.6}$$

for some y_* and z_* . Since, by definition, $A(x_*)S = 0$, and as it then follows from (C.5) that $g^T(x_*)S = 0$, we have that

$$0 = S^T g(x_*) = S^T A^T(x_*)y_* + S^T S z_* = S^T S z_*.$$

Hence $S^T S z_* = 0$ and thus $z_* = 0$ since S is of full rank. Thus (C.6) gives

$$g(x_*) - A^T(x_*)y_* = 0. \tag{C.7}$$

Proof of Theorem 1.7

We have shown that

$$f(x(\alpha)) = f(x_*) + \frac{1}{2}\alpha^2 (\langle p, g(x_*) \rangle + \langle s, H(x_*)s \rangle) + O(\alpha^3) \tag{C.8}$$

for all s satisfying $A(x_*)s = 0$, and that (C.7) holds. Hence, necessarily,

$$\langle p, g(x_*) \rangle + \langle s, H(x_*)s \rangle \geq 0 \tag{C.9}$$

for all s and p satisfying (C.2) and (C.3). But (C.7) and (C.3) combine to give

$$\langle p, g(x_*) \rangle = \sum_{i=1}^m (y_*)_i \langle p, a_i(x_*) \rangle = - \sum_{i=1}^m (y_*)_i \langle s, H_i(x_*)s \rangle$$

and thus (C.9) is equivalent to

$$\left\langle s, \left(H(x_*) - \sum_{i=1}^m (y_*)_i H_i(x_*) \right) s \right\rangle \equiv \langle s, H(x_*, y_*)s \rangle \geq 0$$

for all s satisfying (C.2).

Proof of Theorem 1.8

As in the proof of Theorem 1.6, we consider feasible perturbations about x_* . Since any constraint that is inactive at x_* (i.e., $c_i(x_*) > 0$) will remain inactive for small perturbations, we need only consider perturbations that are constrained by the constraints active at x_* , (i.e., $c_i(x_*) = 0$). Let \mathcal{A} denote the indices of the active constraints. We then consider a vector valued C^2 (C^3 for Theorem 1.9) function $x(\alpha)$ of the scalar α for which $x(0) = x_*$ and $c_i(x(\alpha)) \geq 0$ for $i \in \mathcal{A}$. In this case, assuming that $x(\alpha)$ may be expressed as (C.1), we require that

$$\begin{aligned} 0 &\leq c_i(x(\alpha)) = c(x_* + \alpha s + \frac{1}{2}\alpha^2 p + O(\alpha^3)) \\ &= c_i(x_*) + \langle a_i(x_*), \alpha s + \frac{1}{2}\alpha^2 p \rangle + \frac{1}{2}\alpha^2 \langle s, H_i(x_*)s \rangle + O(\alpha^3) \\ &= \alpha \langle a_i(x_*), s \rangle + \frac{1}{2}\alpha^2 (\langle a_i(x_*), p \rangle + \langle s, H_i(x_*)s \rangle) + O(\alpha^3) \end{aligned}$$

for all $i \in \mathcal{A}$. Thus

$$\langle s, a_i(x_*) \rangle \geq 0 \tag{C.10}$$

and

$$\langle p, a_i(x_*) \rangle + \langle s, H_i(x_*)s \rangle \geq 0 \text{ when } \langle s, a_i(x_*) \rangle = 0 \tag{C.11}$$

for all $i \in \mathcal{A}$. The expansion of $f(x(\alpha))$ (C.4) then implies that x_* can only be a local minimizer if

$$\mathcal{S} = \{s \mid \langle s, g(x_*) \rangle < 0 \text{ and } \langle s, a_i(x_*) \rangle \geq 0 \text{ for } i \in \mathcal{A}\} = \emptyset.$$

But then the result follows directly from Farkas' Lemma—a proof of this famous result is given, for example, as Lemma 9.2.4 in R. Fletcher “Practical Methods of Optimization”, Wiley (1987, 2nd edition).

Farkas' Lemma. Given any vectors g and a_i , $i \in \mathcal{A}$, the set

$$\mathcal{S} = \{s \mid \langle s, g \rangle < 0 \text{ and } \langle s, a_i \rangle \geq 0 \text{ for } i \in \mathcal{A}\}$$

is empty if and only if

$$g = \sum_{i \in \mathcal{A}} y_i a_i$$

for some $y_i \geq 0$, $i \in \mathcal{A}$

Proof of Theorem 1.9

The expansion (C.4) for the change in the objective function will be dominated by the first-order term $\alpha \langle s, g(x_*) \rangle$ for feasible perturbations unless $\langle s, g(x_*) \rangle = 0$, in which case the expansion (C.8) is relevant. Thus we must have that (C.9) holds for all feasible s for which $\langle s, g(x_*) \rangle = 0$. The latter requirement gives that

$$0 = \langle s, g(x_*) \rangle = \sum_{i \in \mathcal{A}} y_i \langle s, a_i(x_*) \rangle,$$

and hence that either $y_i = 0$ or $\langle s, a_i(x_*) \rangle = 0$ (or both).

We now focus on the *subset* of all feasible arcs that ensure $c_i(x(\alpha)) = 0$ if $y_i > 0$ and $c_i(x(\alpha)) \geq 0$ if $y_i = 0$ for $i \in \mathcal{A}$. For those constraints for which $c_i(x(\alpha)) = 0$, we have that (C.2) and (C.3) hold, and thus for such perturbations $s \in \mathcal{N}_+$. In this case

$$\langle p, g(x_*) \rangle = \sum_{i \in \mathcal{A}} y_i \langle p, a_i(x_*) \rangle = \sum_{\substack{i \in \mathcal{A} \\ y_i > 0}} y_i \langle p, a_i(x_*) \rangle = - \sum_{\substack{i \in \mathcal{A} \\ y_i > 0}} y_i \langle s, H_i(x_*) s \rangle = - \sum_{i \in \mathcal{A}} y_i \langle s, H_i(x_*) s \rangle$$

This combines with (C.9) to give that

$$\langle s, H(x_*, y_*) s \rangle \equiv \left\langle s, \left(H(x_*) - \sum_{i=1}^m (y_*)_i H_i(x_*) \right) s \right\rangle = \langle p, g(x_*) \rangle + \langle s, H(x_*) s \rangle \geq 0.$$

for all $s \in \mathcal{N}_+$, which is the required result.

Proof of Theorem 1.10

Consider any feasible arc $x(\alpha)$. We have seen that (C.10) and (C.11) hold, and that first-order feasible perturbations are characterized by \mathcal{N}_+ . It then follows from (C.11) that

$$\langle p, g(x_*) \rangle = \sum_{i \in \mathcal{A}} y_i \langle p, a_i(x_*) \rangle = \sum_{\substack{i \in \mathcal{A} \\ \langle s, a_i(x_*) \rangle = 0}} y_i \langle p, a_i(x_*) \rangle \geq - \sum_{\substack{i \in \mathcal{A} \\ \langle s, a_i(x_*) \rangle = 0}} y_i \langle s, H_i(x_*) s \rangle = - \sum_{i \in \mathcal{A}} y_i \langle s, H_i(x_*) s \rangle,$$

and hence by assumption that

$$\langle p, g(x_*) \rangle + \langle s, H(x_*) s \rangle \geq \left\langle s, \left(H(x_*) - \sum_{i=1}^m (y_*)_i H_i(x_*) \right) s \right\rangle \equiv \langle s, H(x_*, y_*) s \rangle > 0$$

for all $s \in \mathcal{N}_+$. But this then combines with (C.4) and (C.10) to show that $f(x(\alpha)) > f(x_*)$ for all sufficiently small α .

Proof of Theorem 2.1

From Taylor's theorem (Theorem 1.1), and using the bound

$$\alpha \leq \frac{2(\beta - 1) \langle p, g(x) \rangle}{\gamma(x) \|p\|_2^2},$$

we have that

$$\begin{aligned} f(x + \alpha p) &\leq f(x) + \alpha \langle p, g(x) \rangle + \frac{1}{2} \gamma(x) \alpha^2 \|p\|^2 \\ &\leq f(x) + \alpha \langle p, g(x) \rangle + \alpha(\beta - 1) \langle p, g(x) \rangle \\ &= f(x) + \alpha \beta \langle p, g(x) \rangle. \end{aligned}$$

Proof of Corollary 2.2

Theorem 2.1 shows that the linesearch will terminate as soon as $\alpha^{(l)} \leq \alpha_{\max}$. There are two cases to consider. Firstly, it may be that α_{init} satisfies the Armijo condition, in which case $\alpha_k = \alpha_{\text{init}}$. If not, there must be a last linesearch iteration, say the l th, for which $\alpha^{(l)} > \alpha_{\max}$ (if the linesearch has not already terminated). Then $\alpha_k \geq \alpha^{(l+1)} = \tau \alpha^{(l)} > \tau \alpha_{\max}$. Combining these two cases gives the required result.

Proof of Theorem 2.3

We shall suppose that $g_k \neq 0$ for all k and that

$$\lim_{k \rightarrow \infty} f_k > -\infty.$$

From the Armijo condition, we have that

$$f_{k+1} - f_k \leq \alpha_k \beta \langle p_k, g_k \rangle$$

for all k , and hence summing over the first j iterations

$$f_{j+1} - f_0 \leq \sum_{k=0}^j \alpha_k \beta \langle p_k, g_k \rangle.$$

Since the left-hand side of this inequality is, by assumption, bounded below, so is the sum on right-hand-side. As this sum is composed of negative terms, we deduce that

$$\lim_{k \rightarrow \infty} \alpha_k \langle p_k, g_k \rangle = 0.$$

Now define the two sets

$$\mathcal{K}_1 = \left\{ k \mid \alpha_{\text{init}} > \frac{2\tau(\beta - 1) \langle p_k, g_k \rangle}{\gamma \|p_k\|_2^2} \right\}$$

and

$$\mathcal{K}_2 = \left\{ k \mid \alpha_{\text{init}} \leq \frac{2\tau(\beta - 1) \langle p_k, g_k \rangle}{\gamma \|p_k\|_2^2} \right\},$$

where γ is the assumed uniform Lipschitz constant. For $k \in \mathcal{K}_1$,

$$\alpha_k \geq \frac{2\tau(\beta - 1) \langle p_k, g_k \rangle}{\gamma \|p_k\|_2^2}$$

in which case

$$\alpha_k \langle p_k, g_k \rangle \leq \frac{2\tau(\beta - 1)}{\gamma} \left(\frac{\langle p_k, g_k \rangle}{\|p_k\|} \right)^2 < 0.$$

Thus

$$\lim_{k \in \mathcal{K}_1 \rightarrow \infty} \frac{|\langle p_k, g_k \rangle|}{\|p_k\|_2} = 0. \quad (\text{C.12})$$

For $k \in \mathcal{K}_2$,

$$\alpha_k \geq \alpha_{\text{init}}$$

in which case

$$\lim_{k \in \mathcal{K}_2 \rightarrow \infty} |\langle p_k, g_k \rangle| = 0. \quad (\text{C.13})$$

Combining (C.12) and (C.13) gives the required result.

Proof of Theorem 2.4

Follows immediately from Theorem 2.3, since for $p_k = -g_k$,

$$\min(|\langle p_k, g_k \rangle|, |\langle p_k, g_k \rangle| / \|p_k\|_2) = \|g_k\|_2 \min(1, \|g_k\|_2)$$

and thus

$$\lim_{k \rightarrow \infty} \min(|\langle p_k, g_k \rangle|, |\langle p_k, g_k \rangle| / \|p_k\|_2) = 0$$

implies that $\lim_{k \rightarrow \infty} g_k = 0$.

Proof of Theorem 2.5

Let $\lambda^{\min}(B_k)$ and $\lambda^{\max}(B_k)$ be the smallest and largest eigenvalues of B_k . By assumption, there are bounds $\lambda^{\min} > 0$ and λ^{\max} such that

$$\lambda^{\min} \leq \lambda^{\min}(B_k) \leq \frac{\langle s, B_k s \rangle}{\|s\|^2} \leq \lambda^{\max}(B_k) \leq \lambda^{\max}$$

for any nonzero vector s . Thus

$$|\langle p_k, g_k \rangle| = |\langle g_k, B_k^{-1} g_k \rangle| \geq \lambda_{\min}(B_k^{-1}) \|g_k\|_2^2 = \frac{1}{\lambda_{\max}(B_k)} \|g_k\|_2^2 \geq \lambda_{\max}^{-1} \|g_k\|_2^2.$$

In addition

$$\|p_k\|_2^2 = \langle g_k, B_k^{-2} g_k \rangle \leq \lambda_{\max}(B_k^{-2}) \|g_k\|_2^2 = \frac{1}{\lambda_{\min}(B_k^2)} \|g_k\|_2^2 \leq \lambda_{\min}^{-2} \|g_k\|_2^2,$$

and hence

$$\|p_k\|_2 \leq \lambda_{\max}^{-1} \|g_k\|_2,$$

which leads to

$$\frac{|\langle p_k, g_k \rangle|}{\|p_k\|_2} \geq \frac{\lambda_{\min}}{\lambda_{\max}} \|g_k\|_2.$$

Thus

$$\min(|\langle p_k, g_k \rangle|, |\langle p_k, g_k \rangle| / \|p_k\|_2) \geq \lambda_{\max}^{-1} \|g_k\|_2 \min(\|g_k\|_2, \lambda_{\min}).$$

and hence

$$\lim_{k \rightarrow \infty} \min(|\langle p_k, g_k \rangle|, |\langle p_k, g_k \rangle| / \|p_k\|_2) = 0$$

implies, as before, that $\lim_{k \rightarrow \infty} g_k = 0$.

Proof of Theorem 2.6

Consider the sequence of iterates x_k , $k \in \mathcal{K}$, whose limit is x_* . By continuity, H_k is positive definite for all such k sufficiently large. In particular, we have that there is a $k_0 \geq 0$ such that

$$\langle p_k, H_k p_k \rangle \geq \frac{1}{2} \lambda_{\min}(H_*) \|p_k\|_2^2$$

for all $k \in \mathcal{K} \geq k_0$, where $\lambda_{\min}(H_*)$ is the smallest eigenvalue of $H(x_*)$. We may then deduce that

$$|\langle p_k, g_k \rangle| = -\langle p_k, g_k \rangle = \langle p_k, H_k p_k \rangle \geq \frac{1}{2} \lambda_{\min}(H_*) \|p_k\|_2^2. \quad (\text{C.14})$$

for all such k , and also that

$$\lim_{k \in \mathcal{K} \rightarrow \infty} p_k = 0$$

since Theorem 2.5 implies that at least one of the left-hand sides of (C.14) and

$$\frac{|\langle p_k, g_k \rangle|}{\|p_k\|_2} = -\frac{\langle p_k, g_k \rangle}{\|p_k\|_2} \geq \frac{1}{2} \lambda_{\min}(H_*) \|p_k\|_2$$

converges to zero for all such k .

From Taylor's theorem, there is a z_k between x_k and $x_k + p_k$ such that

$$f(x_k + p_k) = f_k + \langle p_k, g_k \rangle + \frac{1}{2} \langle p_k, H(z_k) p_k \rangle.$$

Thus, the Lipschitz continuity of H gives that

$$\begin{aligned} f(x_k + p_k) - f_k - \frac{1}{2}\langle p_k, g_k \rangle &= \frac{1}{2}(\langle p_k, g_k \rangle + \langle p_k, H(z_k)p_k \rangle) \\ &= \frac{1}{2}(\langle p_k, g_k \rangle + \langle p_k, H_k p_k \rangle) + \frac{1}{2}\langle p_k, (H(z_k) - H_k)p_k \rangle \\ &\leq \frac{1}{2}\gamma\|z_k - x_k\|_2\|p_k\|_2^2 \leq \frac{1}{2}\gamma\|p_k\|_2^3 \end{aligned} \quad (\text{C.15})$$

since $H_k p_k + g_k = 0$. Now pick k sufficiently large so that

$$\gamma\|p_k\|_2 \leq \lambda_{\min}(H_*)(1 - 2\beta).$$

In this case, (C.14) and (C.15) give that

$$f(x_k + p_k) - f_k \leq \frac{1}{2}\langle p_k, g_k \rangle + \frac{1}{2}\lambda_{\min}(H_*)(1 - 2\beta)\|p_k\|_2^2 \leq \frac{1}{2}(1 - (1 - 2\beta))\langle p_k, g_k \rangle = \beta\langle p_k, g_k \rangle,$$

and thus that a unit stepsize satisfies the Armijo condition, which proves (i).

To obtain the remaining results, note that $\|H_k^{-1}\|_2 \leq 2/\lambda_{\min}(H_*)$ for all sufficiently large $k \in \mathcal{K}$. The iteration gives

$$x_{k+1} - x_* = x_k - x_* - H_k^{-1}g_k = x_k - x_* - H_k^{-1}(g_k - g(x_*)) = H_k^{-1}(g(x_*) - g_k - H_k(x_* - x_k)).$$

But Theorem 1.2 gives that

$$\|g(x_*) - g_k - H_k(x_* - x_k)\|_2 \leq \gamma\|x_* - x_k\|_2^2.$$

Hence

$$\|x_{k+1} - x_*\|_2 \leq \gamma\|H_k^{-1}\|_2\|x_* - x_k\|_2^2$$

which is (iii) when $\kappa = 2\gamma/\lambda_{\min}(H_*)$. Result (ii) follows since once an iterate becomes sufficiently close to x_* , (iii) implies that the next is even closer.

Conjugate Gradient methods (Section 2)

All of the results given here are easy to verify, and may be found in any of the books of suggested background reading material. The result that any $p_k = p^i$ is a descent direction follows immediately since the fact that p^i minimizes $q(p)$ in \mathcal{D}^i implies that

$$p^i = p^{i-1} - \frac{\langle g_k, d^{i-1} \rangle}{\langle d^{i-1}, B_k d^{i-1} \rangle} d^{i-1}.$$

Thus

$$\langle g_k, p^i \rangle = \langle g_k, p^{i-1} \rangle - \frac{(\langle g_k, d^{i-1} \rangle)^2}{\langle d^{i-1}, B_k d^{i-1} \rangle},$$

from which it follows that $\langle g_k, p^i \rangle < \langle g_k, p^{i-1} \rangle$. The result then follows by induction, since

$$\langle g_k, p^1 \rangle = -\frac{\|g_k\|_2^2}{\langle g_k, B_k g_k \rangle} < 0.$$

Proof of Theorem 3.1

Firstly note that, for all $\alpha \geq 0$,

$$m_k(-\alpha g_k) = f_k - \alpha\|g_k\|_2^2 + \frac{1}{2}\alpha^2\langle g_k, B_k g_k \rangle. \quad (\text{C.16})$$

If g_k is zero, the result is immediate. So suppose otherwise. In this case, there are three possibilities:

- (i) the curvature $\langle g_k, B_k g_k \rangle$ is not strictly positive; in this case $m_k(-\alpha g_k)$ is unbounded from below as α increases, and hence the Cauchy point occurs on the trust-region boundary.
- (ii) the curvature $\langle g_k, B_k g_k \rangle > 0$ and the minimizer of $m_k(-\alpha g_k)$ occurs at or beyond the trust-region boundary; once again, the the Cauchy point occurs on the trust-region boundary.
- (iii) the curvature $\langle g_k, B_k g_k \rangle > 0$ and the minimizer of $m_k(-\alpha g_k)$, and hence the Cauchy point, occurs before the trust-region is reached.

We consider each case in turn;

Case (i). In this case, since $\langle g_k, B_k g_k \rangle \leq 0$, (C.16) gives

$$m_k(-\alpha g_k) = f_k - \alpha \|g_k\|_2^2 + \frac{1}{2} \alpha^2 \langle g_k, B_k g_k \rangle \leq f_k - \alpha \|g_k\|_2^2 \quad (\text{C.17})$$

for all $\alpha \geq 0$. Since the Cauchy point lies on the boundary of the trust region

$$\alpha_k^c = \frac{\Delta_k}{\|g_k\|}. \quad (\text{C.18})$$

Substituting this value into (C.17) gives

$$f_k - m_k(s_k^c) \geq \|g_k\|_2^2 \frac{\Delta_k}{\|g_k\|} \geq \kappa_s \|g_k\|_2 \Delta_k \geq \frac{1}{2} \kappa_s \|g_k\|_2 \Delta_k \quad (\text{C.19})$$

since $\|g_k\|_2 \geq \kappa_s \|g_k\|$.

Case (ii). In this case, let α_k^* be the unique minimizer of (C.16); elementary calculus reveals that

$$\alpha_k^* = \frac{\|g_k\|_2^2}{\langle g_k, B_k g_k \rangle}. \quad (\text{C.20})$$

Since this minimizer lies on or beyond the trust-region boundary (C.18) and (C.20) together imply that

$$\alpha_k^c \langle g_k, B_k g_k \rangle \leq \|g_k\|_2^2.$$

Substituting this last inequality in (C.16), and using (C.18) and $\|g_k\|_2 \geq \kappa_s \|g_k\|$, it follows that

$$f_k - m_k(s_k^c) = \alpha_k^c \|g_k\|_2^2 - \frac{1}{2} [\alpha_k^c]^2 \langle g_k, B_k g_k \rangle \geq \frac{1}{2} \alpha_k^c \|g_k\|_2^2 = \frac{1}{2} \|g_k\|_2^2 \frac{\Delta_k}{\|g_k\|} \geq \frac{1}{2} \kappa_s \|g_k\|_2 \Delta_k.$$

Case (iii). In this case, $\alpha_k^c = \alpha_k^*$, and (C.16) becomes

$$f_k - m_k(s_k^c) = \frac{\|g_k\|_2^4}{\langle g_k, B_k g_k \rangle} - \frac{1}{2} \frac{\|g_k\|_2^4}{\langle g_k, B_k g_k \rangle} = \frac{1}{2} \frac{\|g_k\|_2^4}{\langle g_k, B_k g_k \rangle} \geq \frac{1}{2} \frac{\|g_k\|_2^2}{1 + \|B_k\|},$$

where

$$|\langle g_k, B_k g_k \rangle| \leq \|g_k\|_2^2 \|B_k\|_2 \leq \|g_k\|_2^2 (1 + \|B_k\|_2)$$

because of the Cauchy-Schwarz inequality.

The result follows since it is true in each of the above three possible cases. Note that the “1+” is only needed to cover case where $B_k = 0$, and that in this case, the “min” in the theorem might actually be replaced by $\kappa_s \Delta_k$.

Proof of Corollary 3.2

Immediate from Theorem 3.1 and the requirement that $m_k(s_k) \leq m_k(s_k^c)$.

Proof of Lemma 3.3

The mean value theorem gives that

$$f(x_k + s_k) = f(x_k) + \langle s_k, \nabla_x f(x_k) \rangle + \frac{1}{2} \langle s_k, \nabla_{xx} f(\xi_k) s_k \rangle$$

for some ξ_k in the segment $[x_k, x_k + s_k]$. Thus

$$\begin{aligned} |f(x_k + s_k) - m_k(s_k)| &= \frac{1}{2} |\langle s_k, H(\xi_k) s_k \rangle - \langle s_k, B_k s_k \rangle| \leq \frac{1}{2} |\langle s_k, H(\xi_k) s_k \rangle| + \frac{1}{2} |\langle s_k, B_k s_k \rangle| \\ &\leq \frac{1}{2} (\kappa_h + \kappa_b) \|s_k\|_2^2 \leq \frac{1}{2} \kappa_l^2 (\kappa_h + \kappa_b) \|s_k\|^2 \leq \kappa_d \Delta_k^2 \end{aligned}$$

using the triangle and Cauchy-Schwarz inequalities.

Proof of Lemma 3.4

By definition,

$$1 + \|B_k\|_2 \leq \kappa_h + \kappa_b,$$

and hence for any radius satisfying the given (first) bound,

$$\kappa_s \Delta_k \leq \frac{\|g_k\|_2}{\kappa_h + \kappa_b} \leq \frac{\|g_k\|_2}{1 + \|B_k\|_2}.$$

As a consequence, Corollary 3.2 gives that

$$f_k - m_k(s_k) \geq \frac{1}{2} \|g_k\|_2 \min \left[\frac{\|g_k\|_2}{1 + \|B_k\|_2}, \kappa_s \Delta_k \right] = \frac{1}{2} \kappa_s \|g_k\|_2 \Delta_k. \quad (\text{C.21})$$

But then Lemma 3.3 and the assumed (second) bound on the radius gives that

$$|\rho_k - 1| = \left| \frac{f(x_k + s_k) - m_k(s_k)}{f_k - m_k(s_k)} \right| \leq 2 \frac{\kappa_d \Delta_k^2}{\kappa_s \|g_k\|_2 \Delta_k} = \frac{2\kappa_d}{\kappa_s} \frac{\Delta_k}{\|g_k\|_2} \leq 1 - \eta_v. \quad (\text{C.22})$$

Therefore, $\rho_k \geq \eta_v$ and the iteration is very successful.

Proof of Lemma 3.5

Suppose otherwise that Δ_k can become arbitrarily small. In particular, assume that iteration k is the first such that

$$\Delta_{k+1} \leq \kappa_\epsilon. \quad (\text{C.23})$$

Then since the radius for the previous iteration must have been larger, the iteration was unsuccessful, and thus $\gamma_d \Delta_k \leq \Delta_{k+1}$. Hence

$$\Delta_k \leq \epsilon \min \left(\frac{1}{\kappa_s (\kappa_h + \kappa_b)}, \frac{\kappa_s (1 - \eta_v)}{2\kappa_d} \right) \leq \|g_k\| \min \left(\frac{1}{\kappa_s (\kappa_h + \kappa_b)}, \frac{\kappa_s (1 - \eta_v)}{2\kappa_d} \right).$$

But this contradicts the assertion of Lemma 3.4 that the k -th iteration must be very successful.

Proof of Lemma 3.6

The mechanism of the algorithm ensures that $x_* = x_{k_0+1} = x_{k_0+j}$ for all $j > 0$, where k_0 is the index of the last successful iterate. Moreover, since all iterations are unsuccessful for sufficiently large k , the sequence $\{\Delta_k\}$ converges to zero. If $\|g_{k_0+1}\| > 0$, Lemma 3.4 then implies that there must be a successful iteration of index larger than k_0 , which is impossible. Hence $\|g_{k_0+1}\| = 0$.

Proof of Theorem 3.7

Lemma 3.6 shows that the result is true when there are only a finite number of successful iterations. So it remains to consider the case where there are an infinite number of successful iterations. Let \mathcal{S} be the index set of successful iterations. Now suppose that

$$\|g_k\| \geq \epsilon \quad (\text{C.24})$$

for some $\epsilon > 0$ and all k , and consider a successful iteration of index k . The fact that k is successful, Corollary 3.2, Lemma 3.5, and the assumption (C.24) give that

$$f_k - f_{k+1} \geq \eta_s [f_k - m_k(s_k)] \geq \delta_\epsilon \stackrel{\text{def}}{=} \frac{1}{2} \eta_s \epsilon \min \left[\frac{\epsilon}{1 + \kappa_b}, \kappa_s \kappa_\epsilon \right]. \quad (\text{C.25})$$

Summing now over all successful iterations from 0 to k , it follows that

$$f_0 - f_{k+1} = \sum_{\substack{j=0 \\ j \in \mathcal{S}}}^k [f_j - f_{j+1}] \geq \sigma_k \delta_\epsilon,$$

where σ_k is the number of successful iterations up to iteration k . But since there are infinitely many such iterations, it must be that

$$\lim_{k \rightarrow \infty} \sigma_k = +\infty.$$

Thus (C.24) can only be true if f_{k+1} is unbounded from below, and conversely, if f_{k+1} is bounded from below, (C.24) must be false, and there is a subsequence of the $\|g_k\|$ converging to zero.

Proof of Corollary 3.8

Suppose otherwise that f_k is bounded from below, and that there is a subsequence of successful iterates, indexed by $\{t_i\} \subseteq \mathcal{S}$, such that

$$\|g_{t_i}\| \geq 2\epsilon > 0 \quad (\text{C.26})$$

for some $\epsilon > 0$ and for all i . Theorem 3.7 ensures the existence, for each t_i , of a first successful iteration $\ell_i > t_i$ such that $\|g_{\ell_i}\| < \epsilon$. That is to say that there is another subsequence of \mathcal{S} indexed by $\{\ell_i\}$ such that

$$\|g_k\| \geq \epsilon \text{ for } t_i \leq k < \ell_i \text{ and } \|g_{\ell_i}\| < \epsilon. \quad (\text{C.27})$$

We now restrict our attention to the subsequence of successful iterations whose indices are in the set

$$\mathcal{K} \stackrel{\text{def}}{=} \{k \in \mathcal{S} \mid t_i \leq k < \ell_i\},$$

where t_i and ℓ_i belong to the two subsequences defined above.

The subsequences $\{t_i\}$, $\{\ell_i\}$ and \mathcal{K} are all illustrated in Figure C.1, where, for simplicity, it is assumed that all iterations are successful. In this figure, we have marked position j in each of the subsequences represented in abscissa when j belongs to that subsequence. Note in this example that $\ell_0 = \ell_1 = \ell_2 = \ell_3 = \ell_4 = \ell_5 = 8$, which we indicated by arrows from $t_0 = 0$, $t_1 = 1$, $t_2 = 2$, $t_3 = 3$, $t_4 = 4$ and $t_5 = 7$ to $k = 9$, and so on.

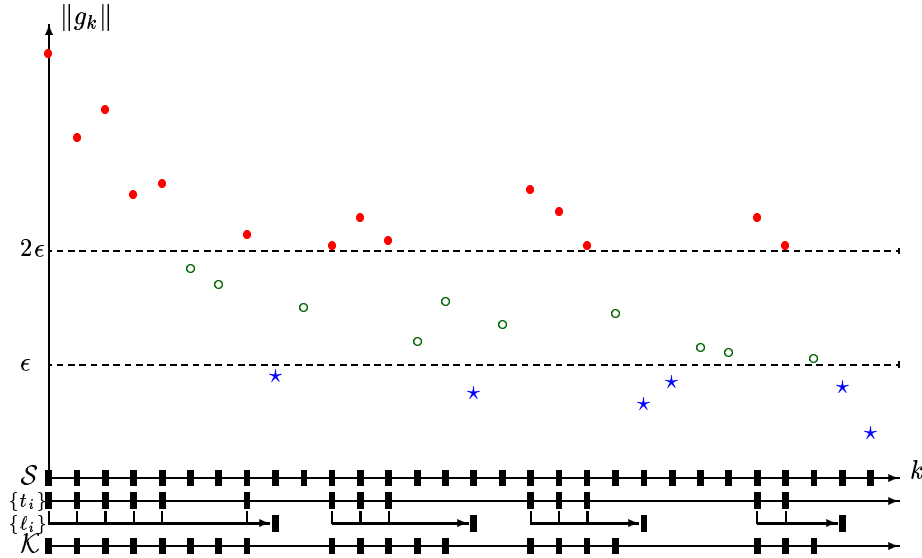


Figure C.1: The subsequences of the proof of Corollary 3.8

As in the previous proof, it immediately follows that

$$f_k - f_{k+1} \geq \eta_s [f_k - m_k(s_k)] \geq \frac{1}{2} \eta_s \epsilon \min \left[\frac{\epsilon}{1 + \kappa_b}, \kappa_s \Delta_k \right] \quad (\text{C.28})$$

holds for all $k \in \mathcal{K}$ because of (C.27). Hence, since $\{f_k\}$ is, by assumption, bounded from below, the left-hand side of (C.28) must tend to zero when k tends to infinity, and thus that

$$\lim_{\substack{k \rightarrow \infty \\ k \in \mathcal{K}}} \Delta_k = 0.$$

As a consequence, the second term dominates in the minimum of (C.28) and it follows that, for $k \in \mathcal{K}$ sufficiently large,

$$\Delta_k \leq \frac{2}{\epsilon \eta_s \kappa_s} [f_k - f_{k+1}].$$

We then deduce from this bound that, for i sufficiently large,

$$\|x_{t_i} - x_{\ell_i}\| \leq \sum_{\substack{j=t_i \\ j \in \mathcal{K}}}^{\ell_i-1} \|x_j - x_{j+1}\| \leq \sum_{\substack{j=t_i \\ j \in \mathcal{K}}}^{\ell_i-1} \Delta_j \leq \frac{2}{\epsilon \eta_s \kappa_s} [f_{t_i} - f_{\ell_i}]. \quad (\text{C.29})$$

But, because $\{f_k\}$ is monotonic and, by assumption, bounded from below, the right-hand side of (C.29) must converge to zero. Thus $\|x_{t_i} - x_{\ell_i}\|$ tends to zero as i tends to infinity, and hence, by continuity, $\|g_{t_i} - g_{\ell_i}\|$ also tends to zero. However this is impossible because of the definitions of $\{t_i\}$ and $\{\ell_i\}$, which imply that $\|g_{t_i} - g_{\ell_i}\| \geq \epsilon$. Hence, no subsequence satisfying (C.26) can exist.

Proof of Theorem 3.9

The constraint $\|s\|_2 \leq \Delta$ is equivalent to

$$\frac{1}{2} \Delta^2 - \frac{1}{2} \langle s, s \rangle \geq 0. \quad (\text{C.30})$$

Applying Theorem 1.7 to the problem of minimizing $q(s)$ subject to (C.30) gives

$$g + Bs_* = -\lambda_* s_* \quad (\text{C.31})$$

for some Lagrange multiplier $\lambda_* \geq 0$ for which either $\lambda_* = 0$ or $\|s_*\|_2 = \Delta$ (or both). It remains to show that $B + \lambda_* I$ is positive semi-definite.

If s_* lies in the interior of the trust-region, necessarily $\lambda_* = 0$, and Theorem 1.8 implies that $B + \lambda_* I = B$ must be positive semi-definite. Likewise if $\|s_*\|_2 = \Delta$ and $\lambda_* = 0$, it follows from Theorem 1.8 that necessarily $\langle v, Bv \rangle \geq 0$ for all $v \in \mathcal{N}_+ = \{v \mid \langle s_*, v \rangle \geq 0\}$. If $v \notin \mathcal{N}_+$, then $-v \in \mathcal{N}_+$, and thus $\langle v, Bv \rangle \geq 0$ for all v . Thus the only outstanding case is where $\|s_*\|_2 = \Delta$ and $\lambda_* > 0$. In this case, Theorem 1.8 shows that $\langle v, (B + \lambda_* I)v \rangle \geq 0$ for all $v \in \mathcal{N}_+ = \{v \mid \langle s_*, v \rangle = 0\}$, so it remains to consider $\langle v, Bv \rangle$ when $\langle s_*, v \rangle \neq 0$.

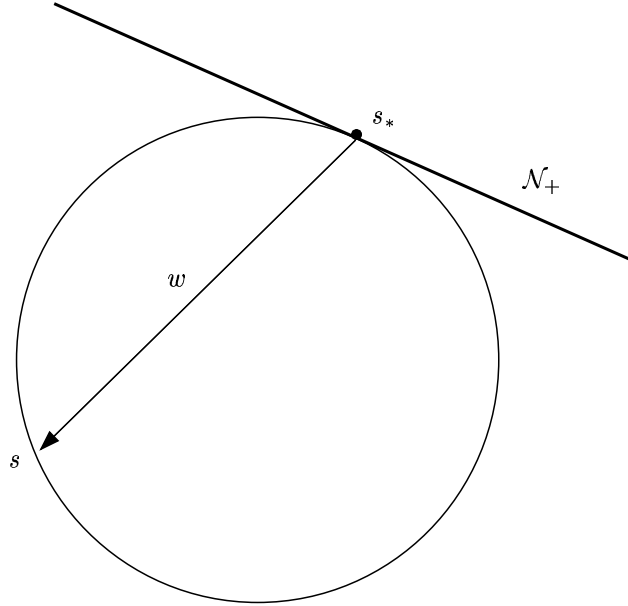


Figure C.2: Construction of “missing” directions of positive curvature.

Let s be any point on the boundary of the trust-region, and let $w = s - s_*$, as in Figure C.2. Then

$$-\langle w, s_* \rangle = \langle s_* - s, s_* \rangle = \frac{1}{2} \langle s_* - s, s_* - s \rangle = \frac{1}{2} \langle w, w \rangle \quad (\text{C.32})$$

since $\|s\|_2 = \Delta = \|s_*\|_2$. Combining this with (C.31) gives

$$q(s) - q(s_*) = \langle w, g + Bs_* \rangle + \frac{1}{2} \langle w, Bw \rangle = -\lambda_* \langle w, s_* \rangle + \frac{1}{2} \langle w, Bw \rangle = \frac{1}{2} \langle w, (B + \lambda_* I)w \rangle, \quad (\text{C.33})$$

and thus necessarily $\langle w, (B + \lambda_* I)w \rangle \geq 0$ since s_* is a global minimizer. It is easy to show that

$$s = s_* - 2 \frac{\langle s_*, v \rangle}{\langle v, v \rangle} v$$

lies on the trust-region boundary, and thus for this s , w is parallel to v from which it follows that $\langle v, (B + \lambda_* I)v \rangle \geq 0$.

When $B + \lambda_* I$ is positive definite, $s_* = -(B + \lambda_* I)^{-1}g$. If this point is on the trust-region boundary, while s is any value in the trust-region, (C.32) and (C.33) become $-\langle w, s_* \rangle \geq \frac{1}{2}\langle w, w \rangle$ and $q(s) \geq q(s_*) + \frac{1}{2}\langle w, (B + \lambda_* I)w \rangle$ respectively. Hence, $q(s) > q(s_*)$ for any $s \neq s_*$. If s_* is interior, $\lambda_* = 0$, B is positive definite, and thus s_* is the unique unconstrained minimizer of $q(s)$.

Newton's method for the secular equation (Section 3)

Recall that the Newton correction at λ is $-\phi(\lambda)/\phi'(\lambda)$. Since

$$\phi(\lambda) = \frac{1}{\|s(\lambda)\|_2} - \frac{1}{\Delta} = \frac{1}{(\langle s(\lambda), s(\lambda) \rangle)^{\frac{1}{2}}} - \frac{1}{\Delta},$$

it follows, on differentiating, that

$$\phi'(\lambda) = -\frac{\langle s(\lambda), \nabla_\lambda s(\lambda) \rangle}{(\langle s(\lambda), s(\lambda) \rangle)^{\frac{3}{2}}} = -\frac{\langle s(\lambda), \nabla_\lambda s(\lambda) \rangle}{\|s(\lambda)\|_2^3}.$$

In addition, on differentiating the defining equation

$$(B + \lambda I)s(\lambda) = -g,$$

it must be that

$$(B + \lambda I)\nabla_\lambda s(\lambda) + s(\lambda) = 0.$$

Notice that, rather than the value of $\nabla_\lambda s(\lambda)$, merely the numerator

$$\langle s(\lambda), \nabla_\lambda s(\lambda) \rangle = -\langle s(\lambda), (B + \lambda I)(\lambda)^{-1}s(\lambda) \rangle$$

is required in the expression for $\phi'(\lambda)$. Given the factorization $B + \lambda I = L(\lambda)L^T(\lambda)$, the simple relationship

$$\langle s(\lambda), (B + \lambda I)^{-1}s(\lambda) \rangle = \langle s(\lambda), L^{-T}(\lambda)L^{-1}(\lambda)s(\lambda) \rangle = \langle L^{-1}(\lambda)s(\lambda), L^{-1}(\lambda)s(\lambda) \rangle = \|w(\lambda)\|_2^2$$

where $L(\lambda)w(\lambda) = s(\lambda)$ then justifies the Newton step.

Proof of Theorem 3.10

We first show that

$$\langle d^i, d^j \rangle = \frac{\|g^i\|_2^2}{\|g^j\|_2^2} \|d^j\|_2^2 > 0 \tag{C.34}$$

for all $0 \leq j \leq i \leq k$. For any i , (C.34) is trivially true for $j = i$. Suppose it is also true for all $i \leq l$. Then, the update for d^{l+1} gives

$$d^{l+1} = -g^{l+1} + \frac{\|g^{l+1}\|_2^2}{\|g^l\|_2^2} d^l.$$

Forming the inner product with d^j , and using the fact that $\langle d^j, g^{l+1} \rangle = 0$ for all $j = 0, \dots, l$, and (C.34) when $j = l$, reveals

$$\langle d^{l+1}, d^j \rangle = -\langle g^{l+1}, d^j \rangle + \frac{\|g^{l+1}\|_2^2}{\|g^l\|_2^2} \langle d^l, d^j \rangle = \frac{\|g^{l+1}\|_2^2}{\|g^l\|_2^2} \frac{\|g^l\|_2^2}{\|g^j\|_2^2} \|d^j\|_2^2 = \frac{\|g^{l+1}\|_2^2}{\|g^j\|_2^2} \|d^j\|_2^2 > 0.$$

Thus (C.34) is true for $i \leq l + 1$, and hence for all $0 \leq j \leq i \leq k$.

We now have from the algorithm that

$$s^i = s^0 + \sum_{j=0}^{i-1} \alpha^j d^j = \sum_{j=0}^{i-1} \alpha^j d^j$$

as, by assumption, $s^0 = 0$. Hence

$$\langle s^i, d^i \rangle = \left\langle \sum_{j=0}^{i-1} \alpha^j d^j, d^i \right\rangle = \sum_{j=0}^{i-1} \alpha^j \langle d^j, d^i \rangle > 0 \quad (\text{C.35})$$

as each $\alpha^j > 0$, which follows from the definition of α^j , since $\langle d^j, H d^j \rangle > 0$, and from relationship (C.34). Hence

$$\begin{aligned} \|s^{i+1}\|_2^2 &= \langle s^{i+1}, s^{i+1} \rangle = \langle s^i + \alpha^i d^i, s^i + \alpha^i d^i \rangle \\ &= \langle s^i, s^i \rangle + 2\alpha^i \langle s^i, d^i \rangle + \alpha^{i2} \langle d^i, d^i \rangle > \langle s^i, s^i \rangle = \|s^i\|_2^2 \end{aligned}$$

follows directly from (C.35) and $\alpha^i > 0$ which is the required result.

Proof of Theorem 3.11

The proof is elementary but rather complicated. See

Y. Yuan, “On the truncated conjugate-gradient method”, *Mathematical Programming*, **87** (2000) 561:573

for full details.

Proof of Theorem 4.1

Let $\mathcal{A} = \mathcal{A}(x_*)$, and $\mathcal{I} = \{1, \dots, m\} \setminus \mathcal{A}$ be the indices of constraints that are active and inactive at x_* . Furthermore let subscripts \mathcal{A} and \mathcal{I} denote the rows of matrices/vectors whose indices are indexed by these sets. Denote the left generalized inverse of $A_{\mathcal{A}}^T(x)$ by

$$A_{\mathcal{A}}^+(x) = (A_{\mathcal{A}}(x)A_{\mathcal{A}}^T(x))^{-1} A_{\mathcal{A}}(x)$$

at any point for which $A_{\mathcal{A}}(x)$ is full rank. Since, by assumption, $A_{\mathcal{A}}(x_*)$ is full rank, these generalized inverses exists, and are bounded and continuous in some open neighbourhood of x_* .

Now let

$$(y_k)_i = \frac{\mu_k}{c_i(x_k)}$$

for $i = 1, \dots, m$, as well as

$$(y_*)_{\mathcal{A}} = A_{\mathcal{A}}^+(x_*)g(x_*)$$

and $(y_*)_{\mathcal{I}} = 0$. If $\mathcal{I} \neq \emptyset$, then

$$\|(y_k)_{\mathcal{I}}\|_2 \leq 2\mu_k \sqrt{|\mathcal{I}|} / \min_{i \in \mathcal{I}} |c_i(x_*)| \quad (\text{C.36})$$

for all sufficiently large k . It then follows from the inner-iteration termination test that

$$\begin{aligned} \|g(x_k) - A_{\mathcal{A}}^T(x_k)(y_k)_{\mathcal{A}}\|_2 &\leq \|g(x_k) - A^T(x_k)y_k\|_2 + \|A_{\mathcal{I}}^T(x_k)(y_k)_{\mathcal{I}}\|_2 \\ &\leq \bar{\epsilon}_k \stackrel{\text{def}}{=} \epsilon_k + \mu_k \frac{2\sqrt{|\mathcal{I}|}\|A_{\mathcal{I}}\|_2}{\min_{i \in \mathcal{I}} |c_i(x_*)|}. \end{aligned} \quad (\text{C.37})$$

Hence

$$\|A_{\mathcal{A}}^+(x_k)g(x_k) - (y_k)_{\mathcal{A}}\|_2 = \|A_{\mathcal{A}}^+(x_k)(g(x_k) - A_{\mathcal{A}}^T(x_k)(y_k)_{\mathcal{A}})\|_2 \leq 2\|A_{\mathcal{A}}^+(x_k)\|_2 \bar{\epsilon}_k.$$

Then

$$\|(y_k)_{\mathcal{A}} - (y_*)_{\mathcal{A}}\|_2 \leq \|A_{\mathcal{A}}^+(x_*)g(x_*) - A_{\mathcal{A}}^+(x_k)g(x_k)\|_2 + \|A_{\mathcal{A}}^+(x_k)g(x_k) - (y_k)_{\mathcal{A}}\|_2$$

which, in combination with (C.36) and convergence of x_k , implies that $\{y_k\}$ converges to y_* . In addition, continuity of the gradients and (C.37) implies that

$$g(x_*) - A^T(x_*)y_* = 0$$

while the fact that $c(x_k) > 0$ for all k , the definition of y_k and y_* (and the implication that $c_i(x_k)(y_k)_i = \mu_k$) shows that $c(x_*) \geq 0$, $y_* \geq 0$ and $c_i(x_*)(y_*)_i = 0$. Hence (x_*, y_*) satisfies the first-order optimality conditions.

Proof of Theorem 4.2

A formal proof is given by

W. Murray, ‘‘Analytical expressions for eigenvalues and eigenvectors of the Hessian matrices of barrier and penalty functions’’, *J. Optimization Theory and Applics*, **7** (1971) 189:196.

By way of a sketch, let $Q(x)$ and $N(x)$ be orthonormal bases for the range- and null-spaces of $A_{\mathcal{A}(x_*)}(x)$, and let $A_{\mathcal{I}}(x)$ be the matrix whose rows are $\{a_i^T(x)\}_{i \notin \mathcal{A}(x_*)}$. As we have shown, the required Hessian may be expressed (in decreasing terms of asymptotic dominance) as

$$\nabla_{xx}\Phi(x, \mu) = A_{\mathcal{A}}^T(x)Y_{\mathcal{A}}^2(x)A_{\mathcal{A}}(x)/\mu + H(x, y(x)) + \mu A_{\mathcal{I}}^T(x)C_{\mathcal{I}}^{-2}(x)A_{\mathcal{I}}(x).$$

Since the eigenvalues of $\nabla_{xx}\Phi(x, \mu)$ are not affected by orthonormal transformations, on pre- and post-multiplying $\nabla_{xx}\Phi(x, \mu)$ by $(Q(x) \ N(x))$ and its transpose, we see that the required eigenvalues are those of

$$\begin{pmatrix} Q(x)^T A_{\mathcal{A}}^T(x) Y_{\mathcal{A}}^2(x) A_{\mathcal{A}}(x) Q(x) / \mu + Q(x)^T H(x, y(x)) Q(x) & Q(x)^T H(x, y(x)) N(x) \\ N(x)^T H(x, y(x)) Q(x) & N(x)^T H(x, y(x)) N(x) \end{pmatrix} + O(\mu), \quad (\text{C.38})$$

where we have used the relationship $A(x)N(x) = 0$. The dominant eigenvalues are those arising from the 1,1 block of (C.38), and these are those of $A_{\mathcal{A}}^T(x)Y_{\mathcal{A}}^2(x)A_{\mathcal{A}}(x)/\mu$ with an $O(1)$ error. Since the remaining eigenvalues must occur for eigenvectors orthogonal to those giving the 1,1 block, they will asymptotically be those of the 2,2 block, and thus those of $N(x)^T H(x, y(x))N(x)$ with an $O(\mu)$ term.

Proof of Theorem 4.3

The proof of this result is elementary, but rather long and involved. See

N. Gould, D. Orban, A. Sartenaer and Ph. L. Toint, ‘‘Superlinear convergence of primal-dual interior point algorithms for nonlinear programming’’, *SIAM J. Optimization*, **11**(4) (2001) 974:1002

for full details.

Proof of Theorem 5.1

The SQP search direction s_k and its associated Lagrange multiplier estimates y_{k+1} satisfy

$$B_k s_k - A_k^T y_{k+1} = -g_k \quad (\text{C.39})$$

and

$$A_k s_k = -c_k. \quad (\text{C.40})$$

Pre-multiplying (C.39) by s_k and using (C.40) gives that

$$\langle s_k, g_k \rangle = -\langle s_k, B_k s_k \rangle + \langle s_k, A_k^T y_{k+1} \rangle = -\langle s_k, B_k s_k \rangle - \langle c_k, y_{k+1} \rangle. \quad (\text{C.41})$$

Likewise (C.40) gives

$$\frac{1}{\mu_k} \langle s_k, A_k^T c_k \rangle = -\frac{\|c_k\|_2^2}{\mu_k}. \quad (\text{C.42})$$

Combining (C.41) and (C.42), and using the positive definiteness of B_k , the Cauchy-Schwarz inequality and the fact that $s_k \neq 0$ if x_k is not critical, yields

$$\begin{aligned} \langle s_k, \nabla_x \Phi(x_k) \rangle &= \left\langle s_k, g_k + \frac{1}{\mu_k} A_k^T c_k \right\rangle = -\langle s_k, B_k s_k \rangle - \langle c_k, y_{k+1} \rangle - \frac{\|c_k\|_2^2}{\mu_k} \\ &< -\|c_k\|_2 \left(\frac{\|c_k\|_2}{\mu_k} - \|y_{k+1}\|_2 \right) \leq 0 \end{aligned}$$

because of the required bound on μ_k .

Proof of Theorem 5.2

The proof is slightly complicated as it uses the calculus of non-differentiable functions. See Theorem 14.3.1 in

R. Fletcher, "Practical Methods of Optimization", Wiley (1987, 2nd edition),

where the converse result, that if x_* is an isolated local minimizer of $\Phi(x, \rho)$ for which $c(x_*) = 0$ then x_* solves the given nonlinear program so long as ρ is sufficiently large, is also given. Moreover, Fletcher shows (Theorem 14.3.2) that x_* cannot be a local minimizer of $\Phi(x, \rho)$ when $\rho < \|y_*\|_D$.

Proof of Theorem 5.3

For small steps α , Taylor's theorem applied separately to f and c , along with (C.40), gives that

$$\begin{aligned} \Phi(x_k + \alpha s_k, \rho_k) - \Phi(x_k, \rho_k) &= \alpha \langle s_k, g_k \rangle + \rho_k (\|c_k + \alpha A_k s_k\| - \|c_k\|) + O(\alpha^2) \\ &= \alpha \langle s_k, g_k \rangle + \rho_k (\|(1 - \alpha)c_k\| - \|c_k\|) + O(\alpha^2) \\ &= \alpha (\langle s_k, g_k \rangle - \rho_k \|c_k\|) + O(\alpha^2). \end{aligned}$$

Combining this with (C.41), and once again using the positive definiteness of B_k , the Hölder inequality (that is that $\langle u, v \rangle \leq \|u\| \|v\|_D$ for any u, v) and the fact that $s_k \neq 0$ if x_k is not critical, yields

$$\begin{aligned} \Phi(x_k + \alpha s_k, \rho_k) - \Phi(x_k, \rho_k) &= -\alpha (\langle s_k, B_k s_k \rangle + \langle c_k, y_{k+1} \rangle + \rho_k \|c_k\|) + O(\alpha^2) \\ &< -\alpha (-\|c_k\| \|y_{k+1}\|_D + \rho_k \|c_k\|) + O(\alpha^2) \\ &= -\alpha \|c_k\| (\rho_k - \|y_{k+1}\|_D) + O(\alpha^2) < 0 \end{aligned}$$

because of the required bound on ρ_k , for sufficiently small α . Hence sufficiently small steps along s_k from non-critical x_k reduce $\Phi(x, \rho_k)$.