

# Descent and Interior-point Methods

Convexity and Optimization – Part III

Lars-Åke Lindahl

$$\int_0^{\infty} e^{-x} \left(\frac{x}{s}\right)^a dx$$

Download free books at

LARS-ÅKE LINDAHL

---

# DESCENT AND INTERIOR-POINT METHODS

CONVEXITY AND  
OPTIMIZATION – PART III

Descent and Interior-point Methods: Convexity and Optimization – Part III

1<sup>st</sup> edition

© 2016 Lars-Åke Lindahl & [bookboon.com](http://bookboon.com)

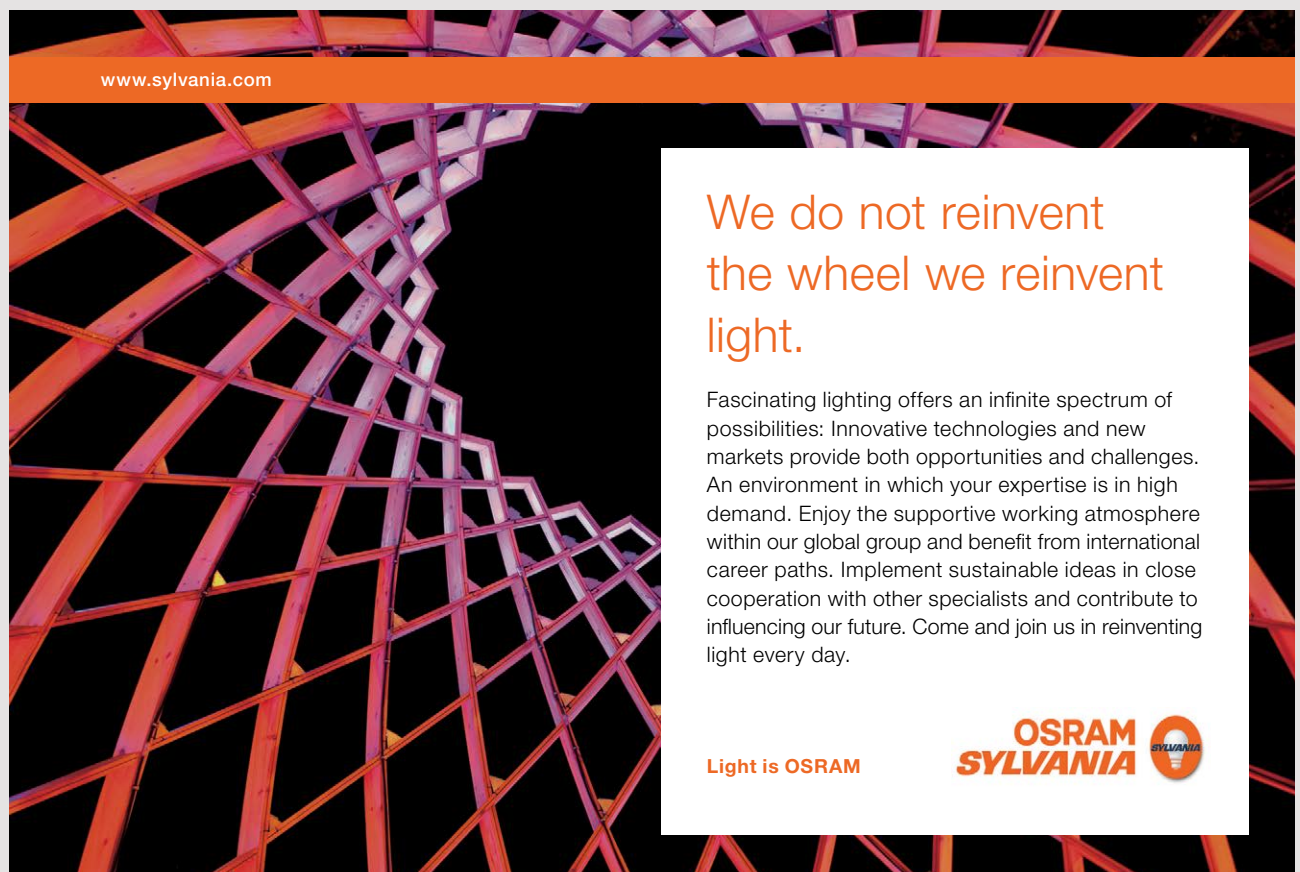
ISBN 978-87-403-1384-0

# CONTENTS

To see Part II, download: **Linear and Convex Optimization: Convexity and Optimization – Part II**

## Part I. Convexity

<b>1</b>	<b>Preliminaries</b>	<b>Part I</b>
<b>2</b>	<b>Convex sets</b>	<b>Part I</b>
2.1	Affine sets and affine maps	<b>Part I</b>
2.2	Convex sets	<b>Part I</b>
2.3	Convexity preserving operations	<b>Part I</b>
2.4	Convex hull	<b>Part I</b>
2.5	Topological properties	<b>Part I</b>
2.6	Cones	<b>Part I</b>
2.7	The recession cone	<b>Part I</b>
	Exercises	<b>Part I</b>




www.sylvania.com

We do not reinvent  
the wheel we reinvent  
light.

Fascinating lighting offers an infinite spectrum of possibilities: Innovative technologies and new markets provide both opportunities and challenges. An environment in which your expertise is in high demand. Enjoy the supportive working atmosphere within our global group and benefit from international career paths. Implement sustainable ideas in close cooperation with other specialists and contribute to influencing our future. Come and join us in reinventing light every day.

Light is OSRAM

**OSRAM SYLVANIA** 

<b>3</b>	<b>Separation</b>	<b>Part I</b>
3.1	Separating hyperplanes	<b>Part I</b>
3.2	The dual cone	<b>Part I</b>
3.3	Solvability of systems of linear inequalities	<b>Part I</b>
	Exercises	<b>Part I</b>
<b>4</b>	<b>More on convex sets</b>	<b>Part I</b>
4.1	Extreme points and faces	<b>Part I</b>
4.2	Structure theorems for convex sets	<b>Part I</b>
	Exercises	<b>Part I</b>
<b>5</b>	<b>Polyhedra</b>	<b>Part I</b>
5.1	Extreme points and extreme rays	<b>Part I</b>
5.2	Polyhedral cones	<b>Part I</b>
5.3	The internal structure of polyhedra	<b>Part I</b>
5.4	Polyhedron preserving operations	<b>Part I</b>
5.5	Separation	<b>Part I</b>
	Exercises	<b>Part I</b>
<b>6</b>	<b>Convex functions</b>	<b>Part I</b>
6.1	Basic definitions	<b>Part I</b>
6.2	Operations that preserve convexity	<b>Part I</b>
6.3	Maximum and minimum	<b>Part I</b>
6.4	Some important inequalities	<b>Part I</b>
6.5	Solvability of systems of convex inequalities	<b>Part I</b>
6.6	Continuity	<b>Part I</b>
6.7	The recessive subspace of convex functions	<b>Part I</b>
6.8	Closed convex functions	<b>Part I</b>
6.9	The support function	<b>Part I</b>
6.10	The Minkowski functional	<b>Part I</b>
	Exercises	<b>Part I</b>
<b>7</b>	<b>Smooth convex functions</b>	<b>Part I</b>
7.1	Convex functions on $\mathbb{R}$	<b>Part I</b>
7.2	Differentiable convex functions	<b>Part I</b>
7.3	Strong convexity	<b>Part I</b>
7.4	Convex functions with Lipschitz continuous derivatives	<b>Part I</b>
	Exercises	<b>Part I</b>

<b>8</b>	<b>The subdifferential</b>	<b>Part I</b>
8.1	The subdifferential	<b>Part I</b>
8.2	Closed convex functions	<b>Part I</b>
8.3	The conjugate function	<b>Part I</b>
8.4	The direction derivative	<b>Part I</b>
8.5	Subdifferentiation rules	<b>Part I</b>
	Exercises	<b>Part I</b>
	<b>Bibliographical and historical notices</b>	<b>Part I</b>
	<b>References</b>	<b>Part I</b>
	<b>Answers and solutions to the exercises</b>	<b>Part I</b>
	<b>Index</b>	<b>Part I</b>
	<b>Endnotes</b>	<b>Part I</b>
	<b>Part II. Linear and Convex Optimization</b>	
	<b>Preface</b>	<b>Part II</b>
	<b>List of symbols</b>	<b>Part II</b>
<b>9</b>	<b>Optimization</b>	<b>Part II</b>
9.1	Optimization problems	Part II
9.2	Classification of optimization problems	Part II
9.3	Equivalent problem formulations	Part II
9.4	Some model examples	Part II
	Exercises	Part II
<b>10</b>	<b>The Lagrange function</b>	<b>Part II</b>
10.1	The Lagrange function and the dual problem	Part II
10.2	John’s theorem	Part II
	Exercises	Part II
<b>11</b>	<b>Convex optimization</b>	<b>Part II</b>
11.1	Strong duality	Part II
11.2	The Karush-Kuhn-Tucker theorem	Part II
11.3	The Lagrange multipliers	Part II
	Exercises	Part II

<b>12</b>	<b>Linear programming</b>	<b>Part II</b>
12.1	Optimal solutions	Part II
12.2	Duality	Part II
	Exercises	Part II
<b>13</b>	<b>The simplex algorithm</b>	<b>Part II</b>
13.1	Standard form	Part II
13.2	Informal description of the simplex algorithm	Part II
13.3	Basic solutions	Part II
13.4	The simplex algorithm	Part II
13.5	Bland’s anti cycling rule	Part II
13.6	Phase 1 of the simplex algorithm	Part II
13.7	Sensitivity analysis	Part II
13.8	The dual simplex algorithm	Part II
13.9	Complexity	Part II
	Exercises	Part II
	<b>Bibliographical and historical notices</b>	<b>Part II</b>
	<b>References</b>	<b>Part II</b>
	<b>Answers and solutions to the exercises</b>	<b>Part II</b>
	<b>Index</b>	<b>Part II</b>
	<b>Part III. Descent and Interior-point Methods</b>	
	<b>Preface</b>	<b>ix</b>
	<b>List of symbols</b>	<b>x</b>
<b>14</b>	<b>Descent methods</b>	<b>1</b>
14.1	General principles	1
14.2	The gradient descent method	7
	Exercises	12
<b>15</b>	<b>Newton’s method</b>	<b>13</b>
15.1	Newton decrement and Newton direction	13
15.2	Newton’s method	22
15.3	Equality constraints	34
	Exercises	39

<b>16</b>	<b>Self-concordant functions</b>	<b>41</b>
16.1	Self-concordant functions	42
16.2	Closed self-concordant functions	47
16.3	Basic inequalities for the local seminorm	51
16.4	Minimization	56
16.5	Newton’s method for self-concordant functions	61
	Exercises	67
	Appendix	68
<b>17</b>	<b>The path-following method</b>	<b>73</b>
17.1	Barrier and central path	74
17.2	Path-following methods	78
<b>18</b>	<b>The path-following method with self-concordant barrier</b>	<b>83</b>
18.1	Self-concordant barriers	83
18.2	The path-following method	94
18.3	LP problems	108
18.4	Complexity	114
	Exercises	125
	<b>Bibliographical and historical notices</b>	<b>127</b>
	<b>References</b>	<b>128</b>
	<b>Answers and solution to the exercises</b>	<b>130</b>
	<b>Index</b>	<b>136</b>



# Preface

This third and final part of Convexity and Optimization discusses some optimization methods which when carefully implemented are efficient numerical optimization algorithms.

We begin with a very brief general description of descent methods and then proceed to a detailed study of Newton’s method. For a particular class of functions, the so-called self-concordant functions, discovered by Yurii Nesterov and Arkadi Nemirovski, it is possible to describe the convergence rate of Newton’s method with absolute constants, and we devote one chapter to this important class.

Interior-point methods are algorithms for solving constrained optimization problems. Contrary to the simplex algorithms, they reach the optimal solution by traversing the interior of the feasible region. Any convex optimization problem can be transformed into minimizing a linear function over a convex set by converting to the epigraph form and with a self-concordant function as barrier, and Nesterov and Nemirovski showed that the number of iterations of the path-following algorithm is bounded by a polynomial in the dimension of the problem and the accuracy of the solution. Their proof is described in this book’s final chapter.

Uppsala, April 2015  
*Lars-Åke Lindahl*

# List of symbols

bdry $X$	boundary of $X$ , see Part I
cl $X$	closure of $X$ , see Part I
dim $X$	dimension of $X$ , see Part I
dom $f$	the effective domain of $f$ : $\{x \mid -\infty < f(x) < \infty\}$ , see Part I
epi $f$	epigraph of $f$ , see Part I
ext $X$	set of extreme points of $X$ , see Part I
int $X$	interior of $X$ , see Part I
lin $X$	recessive subspace of $X$ , see Part I
recc $X$	recession cone of $X$ , see Part I
$e_i$	$i$ th standard basis vector $(0, \dots, 1, \dots, 0)$
$f'$	derivate or gradient of $f$ , see Part I
$f''$	second derivative or hessian of $f$ , see Part I
$v_{\max}, v_{\min}$	optimal values, see Part II
$B(a; r)$	open ball centered at $a$ with radius $r$
$\overline{B}(a; r)$	closed ball centered at $a$ with radius $r$
$Df(a)[v]$	differential of $f$ at $a$ , see Part I
$D^2f(a)[u, v]$	$\sum_{i,j=1}^n \frac{\partial^2 f}{\partial x_i \partial x_j}(a) u_i v_j$ , see Part I
$D^3f(a)[u, v, w]$	$\sum_{i,j,k=1}^n \frac{\partial^3 f}{\partial x_i \partial x_j \partial x_k}(a) u_i v_j w_k$ , see Part I
$\mathcal{E}(x; r)$	ellipsoid $\{y \mid \ y - x\ _x \leq r\}$ , p. 88
$L$	input length, p. 115
$L(x, \lambda)$	Lagrange function, see Part II
$\mathbf{R}_+, \mathbf{R}_{++}$	$\{x \in \mathbf{R} \mid x \geq 0\}, \{x \in \mathbf{R} \mid x > 0\}$
$\mathbf{R}_-$	$\{x \in \mathbf{R} \mid x \leq 0\}$
$\overline{\mathbf{R}}, \mathbf{R}, \underline{\mathbf{R}}$	$\mathbf{R} \cup \{\infty\}, \mathbf{R} \cup \{-\infty\}, \mathbf{R} \cup \{\infty, -\infty\}$
$S_{\mu,L}(X)$	class of $\mu$ -strongly convex functions on $X$ with $L$ -Lipschitz continuous derivative, see Part I
$\text{Var}_X(v)$	$\sup_{x \in X} \langle v, x \rangle - \inf_{x \in X} \langle v, x \rangle$ , p. 93
$X^+$	dual cone of $X$ , see Part I
$\mathbf{1}$	the vector $(1, 1, \dots, 1)$
$\lambda(f, x)$	Newton decrement of $f$ at $x$ , p. 16
$\pi_y$	translated Minkowski functional, p. 89
$\rho(t)$	$-t - \ln(1 - t)$ , p. 51
$\Delta x_{\text{nt}}$	Newton direction at $x$ , p. 15
$\nabla f$	gradient of $f$
$[x, y]$	line segment between $x$ and $y$
$]x, y[$	open line segment between $x$ and $y$
$\ \cdot\ _1, \ \cdot\ _2, \ \cdot\ _\infty$	$\ell^1$ -norm, Euclidean norm, maximum norm, see Part I
$\ \cdot\ _x$	the seminorm $\sqrt{\langle \cdot, f''(x) \cdot \rangle}$ , p. 18
$\ v\ _x^*$	dual local seminorm $\sup_{\ w\ _x \leq 1} \langle v, w \rangle$ , p. 92

# Chapter 14

## Descent methods

The most common numerical algorithms for minimization of differentiable functions of several variables are so-called *descent algorithms*. A descent algorithm is an iterative algorithm that from a given starting point generates a sequence of points with decreasing function values, and the process is stopped when one has obtained a function value that approximates the minimum value good enough according to some criterion. However, there is no algorithm that works for arbitrary functions; special assumptions about the function to be minimized are needed to ensure convergence towards the minimum point. Convexity is such an assumption, which makes it also possible in many cases to determine the speed of convergence.

This chapter describes descent methods in general terms, and we exemplify with the simplest descent method, the gradient descent method.

### 14.1 General principles

We shall study the optimization problem

$$(P) \quad \min f(x)$$

where  $f$  is a function which is defined and differentiable on an open subset  $\Omega$  of  $\mathbf{R}^n$ . We assume that the problem has a solution, i.e. that there is an optimal point  $\hat{x} \in \Omega$ , and we denote the optimal value  $f(\hat{x})$  as  $f_{\min}$ . A convenient assumption which, according to Corollary 8.1.7 in Part I, guarantees the existence of a (unique) optimal solution is that  $f$  is strongly convex and has some closed nonempty sublevel set.

Our aim is to generate a sequence  $x_1, x_2, x_3, \dots$  of points in  $\Omega$  from a given *starting point*  $x_0 \in \Omega$ , with decreasing function values and with the property that  $f(x_k) \rightarrow f_{\min}$  as  $k \rightarrow \infty$ . In the iteration leading from the

point  $x_k$  to the next point  $x_{k+1}$ , except when  $x_k$  is already optimal, one first selects a vector  $v_k$  such that the one-variable function  $\phi_k(t) = f(x_k + tv_k)$  is strictly decreasing at  $t = 0$ . Then, a *line search* is performed along the half-line  $x_k + tv_k$ ,  $t > 0$ , and a point  $x_{k+1} = x_k + h_k v_k$  satisfying  $f(x_{k+1}) < f(x_k)$  is selected according to specific rules.

The vector  $v_k$  is called the *search direction*, and the positive number  $h_k$  is called the *step size*. The algorithm is terminated when the difference  $f(x_k) - f_{\min}$  is less than a given tolerance.

Schematically, we can describe a typical descent algorithm as follows:

### Descent algorithm

**Given** a starting point  $x \in \Omega$ .

**Repeat**

1. Determine (if  $f'(x) \neq 0$ ) a search direction  $v$  and a step size  $h > 0$  such that  $f(x + hv) < f(x)$ .
2. *Update*:  $x := x + hv$ .

**until** stopping criterion is satisfied.

Different strategies for selecting the search direction, different ways to perform the line search, as well as different stop criteria, give rise to different algorithms, of course.

### Search direction

Permitted search directions in iteration  $k$  are vectors  $v_k$  which satisfy the inequality

$$\langle f'(x_k), v_k \rangle < 0,$$

because this ensures that the function  $\phi_k(t) = f(x_k + tv_k)$  is decreasing at the point  $t = 0$ , since  $\phi_k'(0) = \langle f'(x_k), v_k \rangle$ . We will study two ways to select the search direction.

The *gradient descent method* selects  $v_k = -f'(x_k)$ , which is a permissible choice since  $\langle f'(x_k), v_k \rangle = -\|f'(x_k)\|^2 < 0$ . Locally, this choice gives the fastest decrease in function value.

*Newton's method* assumes that the second derivative exists, and the search direction at points  $x_k$  where the second derivative is positive definite is

$$v_k = -f''(x_k)^{-1} f'(x_k).$$

This choice is permissible since  $\langle f'(x_k), v_k \rangle = -\langle f'(x_k), f''(x_k)^{-1} f'(x_k) \rangle < 0$ .

## Line search

Given the search direction  $v_k$  there are several possible strategies for selecting the step size  $h_k$ .

1. *Exact line search.* The step size  $h_k$  is determined by minimizing the one-variable function  $t \mapsto f(x_k + tv_k)$ . This method is used for theoretical studies of algorithms but almost never in practice due to the computational cost of performing the one-dimensional minimization.

2. The step size sequence  $(h_k)_{k=1}^{\infty}$  is given *a priori*, for example as  $h_k = h$  or as  $h_k = h/\sqrt{k+1}$  for some positive constant  $h$ . This is a simple rule that is often used in convex optimization.

3. The step size  $h_k$  at the point  $x_k$  is defined as  $h_k = \rho(x_k)$  for some given function  $\rho$ . This technique is used in the analysis of Newton's method for self-concordant functions.

4. *Armijo's rule.* The step size  $h_k$  at the point  $x_k$  depends on two parameters  $\alpha, \beta \in ]0, 1[$  and is defined as

$$h_k = \beta^m,$$

where  $m$  is the smallest nonnegative integer such that the point  $x_k + \beta^m v_k$



Discover the truth at [www.deloitte.ca/careers](http://www.deloitte.ca/careers)

**Deloitte.**

© Deloitte & Touche LLP and affiliated entities.

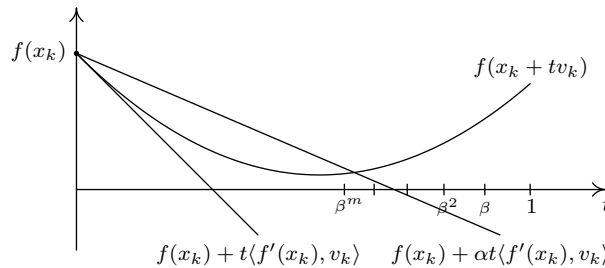
lies in the domain of  $f$  and satisfies the inequality

$$(14.1) \quad f(x_k + \beta^m v_k) \leq f(x_k) + \alpha \beta^m \langle f'(x_k), v_k \rangle.$$

Such an  $m$  certainly exists, since  $\beta^n \rightarrow 0$  as  $n \rightarrow \infty$  and

$$\lim_{t \rightarrow 0} \frac{f(x_k + tv_k) - f(x_k)}{t} = \langle f'(x_k), v_k \rangle < \alpha \langle f'(x_k), v_k \rangle.$$

The number  $m$  is determined by simple backtracking: Start with  $m = 0$  and examine whether  $x_k + \beta^m v_k$  belongs to the domain of  $f$  and inequality (14.1) holds. If not, increase  $m$  by 1 and repeat until the conditions are fulfilled. Figure 14.1 illustrates the process.



**Figure 14.1.** Armijo's rule: The step size is  $h_k = \beta^m$ , where  $m$  is the smallest nonnegative integer such that  $f(x_k + \beta^m v_k) \leq f(x_k) + \alpha \beta^m \langle f'(x_k), v_k \rangle$ .

The decrease in iteration  $k$  of function value per step size, i.e. the ratio  $(f(x_k) - f(x_{k+1})) / h_k$ , is for convex functions less than or equal to  $-\langle f'(x_k), v_k \rangle$  for any choice of step size  $h_k$ . With step size  $h_k$  selected according to Armijo's rule the same ratio is also  $\geq -\alpha \langle f'(x_k), v_k \rangle$ . With Armijo's rule, the decrease per step size is, in other words, at least  $\alpha$  of what the maximum might be. Typical values of  $\alpha$  in practical applications lie in the range between 0.01 and 0.3.

The parameter  $\beta$  determines how many backtracking steps are needed. The larger  $\beta$ , the more backtracking steps, i.e. the finer the line search. The parameter  $\beta$  is often chosen between 0.1 and 0.8.

Armijo's rule exists in different versions and is used in several practical algorithms.

### Stopping criteria

Since the optimum value is generally not known beforehand, it is not possible to formulate the stopping criterion directly in terms of the minimum.

Intuitively, it seems reasonable that  $x$  should be close to the minimum point if the derivative  $f'(x)$  is comparatively small, and the next theorem shows that this is indeed the case, under appropriate conditions on the objective function.

**Theorem 14.1.1.** *Suppose that the function  $f: \Omega \rightarrow \mathbf{R}$  is differentiable,  $\mu$ -strongly convex and has a minimum at  $\hat{x} \in \Omega$ . Then, for all  $x \in \Omega$*

$$(i) \quad f(x) - f(\hat{x}) \leq \frac{1}{2\mu} \|f'(x)\|^2 \quad \text{and}$$

$$(ii) \quad \|x - \hat{x}\| \leq \frac{1}{\mu} \|f'(x)\|.$$

*Proof.* Due to the convexity assumption,

$$(14.2) \quad f(y) \geq f(x) + \langle f'(x), y - x \rangle + \frac{1}{2}\mu \|y - x\|^2$$

for all  $x, y \in \Omega$ . The right-hand side of inequality (14.2) is a convex quadratic function in the variable  $y$ , which is minimized by  $y = x - \mu^{-1}f'(x)$ , and the minimum is equal to  $f(x) - \frac{1}{2}\mu^{-1}\|f'(x)\|^2$ . Hence,

$$f(y) \geq f(x) - \frac{1}{2}\mu^{-1}\|f'(x)\|^2$$

for all  $y \in \Omega$ , and we obtain the inequality (i) by choosing  $y$  as the minimum point  $\hat{x}$ .

Now, replace  $y$  with  $x$  and  $x$  with  $\hat{x}$  in inequality (14.2). Since  $f'(\hat{x}) = 0$ , the resulting inequality becomes

$$f(x) \geq f(\hat{x}) + \frac{1}{2}\mu \|x - \hat{x}\|^2,$$

which combined with inequality (i) gives us inequality (ii). □

We now return to the descent algorithm and our discussion of the the stopping criterion. Let

$$S = \{x \in \Omega \mid f(x) \leq f(x_0)\},$$

where  $x_0$  is the selected starting point, and assume that the sublevel set  $S$  is convex and that the objective function  $f$  is  $\mu$ -strongly convex on  $S$ . All the points  $x_1, x_2, x_3, \dots$  that are generated by the descent algorithm will of course lie in  $S$  since the function values are decreasing. Therefore, it follows from Theorem 14.1.1 that  $f(x_k) < f_{\min} + \epsilon$  if  $\|f'(x_k)\| < (2\mu\epsilon)^{1/2}$ .

As a stopping criterion, we can thus use the condition

$$\|f'(x_k)\| \leq \eta,$$

which guarantees that  $f(x_k) - f_{\min} \leq \eta^2/2\mu$  and that  $\|x_k - \hat{x}\| \leq \eta/\mu$ . A problem here is that the convexity constant  $\mu$  is known only in rare cases. So the stopping condition  $\|f'(x_k)\| \leq \eta$  can in general not be used to give precise bounds on  $f(x_k) - f_{\min}$ . But Theorem 14.1.1 verifies our intuitive feeling that the difference between  $f(x)$  and  $f_{\min}$  is small if the gradient of  $f$  at  $x$  is small enough.

### Convergence rate

Let us say that a convergent sequence  $x_0, x_1, x_2, \dots$  of points with limit  $\hat{x}$  converges *at least linearly* if there is a constant  $c < 1$  such that

$$(14.3) \quad \|x_{k+1} - \hat{x}\| \leq c\|x_k - \hat{x}\|$$


for all  $k$ , and that the convergence is *at least quadratic* if there is a constant  $C$  such that

$$(14.4) \quad \|x_{k+1} - \hat{x}\| \leq C\|x_k - \hat{x}\|^2$$


for all  $k$ .

SIMPLY CLEVER

ŠKODA



**We will turn your CV into an opportunity of a lifetime**



Do you like cars? Would you like to be a part of a successful brand? We will appreciate and reward both your enthusiasm and talent. Send us your CV. You will be surprised where it can take you.

Send us your CV on  
[www.employerforlife.com](http://www.employerforlife.com)





We also say that the convergence is *no better than linear* and *no better than quadratic* if

$$\liminf_{k \rightarrow \infty} \frac{\|x_{k+1} - \hat{x}\|}{\|x_k - \hat{x}\|^\alpha} > 0$$

for  $\alpha = 1$  and  $\alpha = 2$ , respectively.

Note that inequality (14.3) implies that the sequence  $(x_k)_0^\infty$  converges to  $\hat{x}$ , because it follows by induction that

$$\|x_k - \hat{x}\| \leq c^k \|x_0 - \hat{x}\|$$

for all  $k$ .

Similarly, inequality (14.4) implies that the sequence  $(x_k)_0^\infty$  converges to  $\hat{x}$  if the starting point  $x_0$  satisfies the condition  $\|x_0 - \hat{x}\| < C^{-1}$ , because we now have

$$\|x_k - \hat{x}\| \leq C^{-1} (C \|x_0 - \hat{x}\|)^{2^k}$$

for all  $k$ .

If an iterative method, when applied to functions in a given class of functions, always generates sequences that are at least linearly (quadratic) convergent and there is a sequence which does not converge better than linearly (quadratic), then we say that the method is *linearly (quadratic) convergent* for the function class in question.

## 14.2 The gradient descent method

In this section we analyze the gradient descent algorithm with constant step size. The iterative formulation of the variant of the algorithm that we have in mind looks like this:

### Gradient descent algorithm with constant step size

**Given** a starting point  $x$  and a step size  $h$ .

**Repeat**

1. Compute the search direction  $v = -f'(x)$ .
2. *Update:*  $x := x + hv$ .

**until** stopping criterion is satisfied.

The algorithm converges linearly to the minimum point for strongly convex functions with Lipschitz continuous derivatives provided that the step size is small enough and the starting point is chosen sufficiently close to the minimum point. This is the main content of the following theorem (and Example 14.2.1).

**Theorem 14.2.1.** *Let  $f$  be a function with a local minimum point  $\hat{x}$ , and suppose that there is an open neighborhood  $U$  of  $\hat{x}$  such that the restriction  $f|_U$  of  $f$  to  $U$  is  $\mu$ -strongly convex and differentiable with a Lipschitz continuous derivative and Lipschitz constant  $L$ . The gradient descent algorithm with constant step size  $h$  then converges at least linearly to  $\hat{x}$  provided that the step size is sufficiently small and the starting point  $x_0$  lies sufficiently close to  $\hat{x}$ .*

*More precisely: If the ball centered at  $\hat{x}$  and with radius equal to  $\|x_0 - \hat{x}\|$  lies in  $U$  and if  $h \leq \mu/L^2$ , and  $(x_k)_0^\infty$  is the sequence of points generated by the algorithm, then  $x_k$  lies in  $U$  and*

$$\|x_{k+1} - \hat{x}\| \leq c\|x_k - \hat{x}\|,$$

for all  $k$ , where  $c = \sqrt{1 - h\mu}$ .

*Proof.* Suppose inductively that the points  $x_0, x_1, \dots, x_k$  lie in  $U$  and that  $\|x_k - \hat{x}\| \leq \|x_0 - \hat{x}\|$ . Since the restriction  $f|_U$  is assumed to be  $\mu$ -strongly convex and since  $f'(\hat{x}) = 0$ ,

$$\langle f'(x_k), x_k - \hat{x} \rangle = \langle f'(x_k) - f'(\hat{x}), x_k - \hat{x} \rangle \geq \mu\|x_k - \hat{x}\|^2$$

according to Theorem 7.3.1 in Part I, and since the derivative is assumed to be Lipschitz continuous, we also have the inequality

$$\|f'(x_k)\| = \|f'(x_k) - f'(\hat{x})\| \leq L\|x_k - \hat{x}\|.$$

By combining these two inequalities, we obtain the inequality

$$\begin{aligned} \langle f'(x_k), x_k - \hat{x} \rangle &\geq \mu\|x_k - \hat{x}\|^2 = \frac{\mu}{2}\|x_k - \hat{x}\|^2 + \frac{\mu}{2}\|x_k - \hat{x}\|^2 \\ &\geq \frac{\mu}{2}\|x_k - \hat{x}\|^2 + \frac{\mu}{2L^2}\|f'(x_k)\|^2. \end{aligned}$$

Our next point  $x_{k+1} = x_k - hf'(x_k)$  therefore satisfies the inequality

$$\begin{aligned} \|x_{k+1} - \hat{x}\|^2 &= \|x_k - hf'(x_k) - \hat{x}\|^2 = \|(x_k - \hat{x}) - hf'(x_k)\|^2 \\ &= \|x_k - \hat{x}\|^2 - 2h\langle f'(x_k), x_k - \hat{x} \rangle + h^2\|f'(x_k)\|^2 \\ &\leq \|x_k - \hat{x}\|^2 - h\mu\|x_k - \hat{x}\|^2 - h\frac{\mu}{L^2}\|f'(x_k)\|^2 + h^2\|f'(x_k)\|^2 \\ &= (1 - h\mu)\|x_k - \hat{x}\|^2 + h\left(h - \frac{\mu}{L^2}\right)\|f'(x_k)\|^2. \end{aligned}$$

Hence,  $h \leq \mu/L^2$  implies that  $\|x_{k+1} - \hat{x}\|^2 \leq (1 - h\mu)\|x_k - \hat{x}\|^2$ , and this proves that the inequality of the theorem holds with  $c = \sqrt{1 - h\mu} < 1$ , and that the induction hypothesis is satisfied by the point  $x_{k+1}$ , too, since it lies closer to  $\hat{x}$  than the point  $x_k$  does. So the gradient descent algorithm converges at least linearly for  $f$  under the given conditions on  $h$  and  $x_0$ .  $\square$

We can obtain a slightly sharper result for  $\mu$ -strongly convex functions that are defined on the whole  $\mathbf{R}^n$  and have a Lipschitz continuous derivative.

**Theorem 14.2.2.** *Let  $f$  be a function in the class  $\mathcal{S}_{\mu,L}(\mathbf{R}^n)$ . The gradient descent method, with arbitrary starting point  $x_0$  and constant step size  $h$ , generates a sequence  $(x_k)_0^\infty$  of points that converges at least linearly to the function's minimum point  $\hat{x}$ , if*

$$0 < h \leq \frac{2}{\mu + L}.$$

More precisely,

$$(14.5) \quad \|x_k - \hat{x}\| \leq \left(1 - \frac{2h\mu L}{\mu + L}\right)^{k/2} \|x_0 - \hat{x}\|.$$

Moreover, if  $h = \frac{2}{\mu + L}$  then

$$(14.6) \quad \|x_k - \hat{x}\| \leq \left(\frac{Q-1}{Q+1}\right)^k \|x_0 - \hat{x}\| \quad \text{and}$$

$$(14.7) \quad f(x_k) - f_{\min} \leq \frac{L}{2} \left(\frac{Q-1}{Q+1}\right)^{2k} \|x_0 - \hat{x}\|^2,$$

where  $Q = L/\mu$  is the condition number of the function class  $\mathcal{S}_{\mu,L}(\mathbf{R}^n)$ .

I joined MITAS because  
I wanted **real responsibility**

The Graduate Programme  
for Engineers and Geoscientists  
[www.discovermitas.com](http://www.discovermitas.com)



**Month 16**

I was a construction  
supervisor in  
the North Sea  
advising and  
helping foremen  
solve problems

Real work  
International opportunities  
Three work placements



 **MAERSK**

*Proof.* The function  $f$  has a unique minimum point  $\hat{x}$ , according to Corollary 8.1.7 in Part I, and

$$\|x_{k+1} - \hat{x}\|^2 = \|x_k - \hat{x}\|^2 - 2h\langle f'(x_k), x_k - \hat{x} \rangle + h^2\|f'(x_k)\|^2,$$

just as in the proof of Theorem 14.2.1. Since  $f'(\hat{x}) = 0$ , it now follows from Theorem 7.4.4 in Part I (with  $x = \hat{x}$  and  $v = x_k - \hat{x}$ ) that

$$\langle f'(x_k), x_k - \hat{x} \rangle \geq \frac{\mu L}{\mu + L}\|x_k - \hat{x}\|^2 + \frac{1}{\mu + L}\|f'(x_k)\|^2,$$

which inserted in the above equation results in the inequality

$$\|x_{k+1} - \hat{x}\|^2 \leq \left(1 - \frac{2h\mu L}{\mu + L}\right)\|x_k - \hat{x}\|^2 + h\left(h - \frac{2}{\mu + L}\right)\|f'(x_k)\|^2.$$

So if  $h \leq 2/(\mu + L)$ , then

$$\|x_{k+1} - \hat{x}\| \leq \left(1 - \frac{2h\mu L}{\mu + L}\right)^{1/2}\|x_k - \hat{x}\|,$$

and inequality (14.5) now follows by iteration.

The particular choice of  $h = 2(\mu + L)^{-1}$  in inequality (14.5) gives us inequality (14.6), and the last inequality (14.7) follows from inequality (14.6) and Theorem 1.1.2 in Part I, since  $f'(\hat{x}) = 0$ .  $\square$

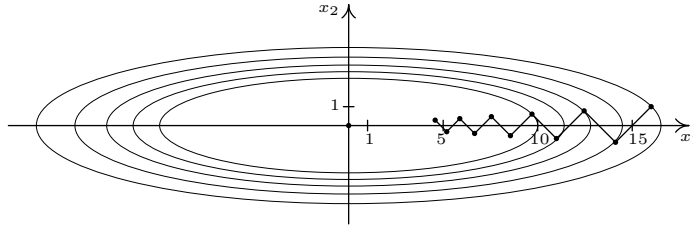
The rate of convergence in Theorems 14.2.1 and 14.2.2 depends on the condition number  $Q \geq 1$ . The smaller the  $Q$ , the faster the convergence. The constants  $\mu$  and  $L$ , and hence the condition number  $Q$ , are of course rarely known in practical examples, so the two theorems have a qualitative character and can rarely be used to predict the number of iterations required to achieve a certain precision.

Our next example shows that inequality (14.6) can not be sharpened.

**EXAMPLE 14.2.1.** Consider the function

$$f(x) = \frac{1}{2}(\mu x_1^2 + Lx_2^2),$$

where  $0 < \mu \leq L$ . This function belongs to the class  $\mathcal{S}_{\mu,L}(\mathbf{R}^2)$ ,  $f'(x) = (\mu x_1, Lx_2)$ , and  $\hat{x} = (0, 0)$  is the minimum point.



**Figure 14.2.** Some level curves for the function  $f(x) = \frac{1}{2}(x_1^2 + 16x_2^2)$  and the progression of the gradient descent algorithm with  $x^{(0)} = (16, 1)$  as starting point. The function's condition number  $Q$  is equal to 16, so the convergence to the minimum point  $(0, 0)$  is relatively slow. The distance from the generated point to the origin is improved by a factor of  $15/17$  in each iteration.

The gradient descent algorithm with constant step size  $h = 2(\mu + L)^{-1}$ , starting point  $x^{(0)} = (L, \mu)$ , and  $\alpha = \frac{Q-1}{Q+1}$  proceeds as follows

$$\begin{aligned} x^{(0)} &= (L, \mu) \\ f'(x^{(0)}) &= (\mu L, \mu L) \\ x^{(1)} &= x^{(0)} - hf'(x^{(0)}) = \alpha(L, -\mu) \\ f'(x^{(1)}) &= \alpha(\mu L, -\mu L) \\ x^{(2)} &= x^{(1)} - hf'(x^{(1)}) = \alpha^2(L, \mu) \\ &\vdots \\ x^{(k)} &= \alpha^k(L, (-1)^k \mu) \end{aligned}$$

Consequently,

$$\|x^{(k)} - \hat{x}\| = \alpha^k \sqrt{L^2 + \mu^2} = \alpha^k \|x^{(0)} - \hat{x}\|,$$

so inequality (14.6) holds with equality in this case. Cf. with figure 14.2.

Finally, it is worth noting that  $2(\mu + L)^{-1}$  coincides with the step size that we would obtain if we had used exact line search in each iteration step.  $\square$

The gradient descent algorithm is not invariant under affine coordinate changes. The speed of convergence can thus be improved by first making a coordinate change that reduces the condition number.

**EXAMPLE 14.2.2.** We continue with the function  $f(x) = \frac{1}{2}(\mu x_1^2 + Lx_2^2)$  in the previous example. Make the change of variables  $y_1 = \sqrt{\mu}x_1$ ,  $y_2 = \sqrt{L}x_2$ ,

and define the function  $g$  by

$$g(y) = f(x) = \frac{1}{2}(y_1^2 + y_2^2).$$

The condition number  $Q$  of the function  $g$  is equal to 1, so the gradient descent algorithm, started from an arbitrary point  $y^{(0)}$ , hits the minimum point  $(0, 0)$  after just one iteration.  $\square$

The gradient descent algorithm converges too slowly to be of practical use in realistic problems. In the next chapter we shall therefore study in detail a more efficient method for optimization, Newton's method.

## Exercises

**14.1** Perform three iterations of the gradient descent algorithm with  $(1, 1)$  as starting point on the minimization problem

$$\min x_1^2 + 2x_2^2.$$

**14.2** Let  $X = \{x \in \mathbf{R}^2 \mid x_1 > 1\}$ , let  $x^{(0)} = (2, 2)$ , and let  $f: X \rightarrow \mathbf{R}$  be the function defined by  $f(x) = \frac{1}{2}x_1^2 + \frac{1}{2}x_2^2$ .

a) Show that the sublevel set  $\{x \in X \mid f(x) \leq f(x^{(0)})\}$  is not closed.

b) Obviously,  $f_{\min} = \inf f(x) = \frac{1}{2}$ , but show that the gradient descent method, with  $x^{(0)}$  as starting point and with line search according to Armijo's rule with parameters  $\alpha \leq \frac{1}{2}$  and  $\beta < 1$ , generates a sequence  $x^{(k)} = (a_k, a_k)$ ,  $k = 0, 1, 2, \dots$ , of points that converges to the point  $(1, 1)$ . So the function values  $f(x^{(k)})$  converge to 1 and not to  $f_{\min}$ .

[Hint: Show that  $a_{k+1} - 1 \leq (1 - \beta)(a_k - 1)$  for all  $k$ .]

**14.3** Suppose that the gradient descent algorithm with constant step size converges to the point  $\hat{x}$  when applied to a continuously differentiable function  $f$ . Prove that  $\hat{x}$  is a stationary point of  $f$ , i.e. that  $f'(\hat{x}) = 0$ .

# Chapter 15

## Newton's method

In Newton's method for minimizing a function  $f$ , the search direction at a point  $x$  is determined by minimizing the function's Taylor polynomial of degree two, i.e. the polynomial

$$P(v) = f(x) + Df(x)[v] + \frac{1}{2}D^2f(x)[v, v] = f(x) + \langle f'(x), v \rangle + \frac{1}{2}\langle v, f''(x)v \rangle,$$

and since  $P'(v) = f'(x) + f''(x)v$ , we obtain the minimizing search vector as a solution to the equation

$$f''(x)v = -f'(x).$$

Each iteration is of course more laborious in Newton's method than in the gradient descent method, since we need to compute the second derivative and solve a quadratic equation to determine the search vector. However, as we shall see, this is more than compensated by a much faster convergence to the minimum value.

### 15.1 Newton decrement and Newton direction

Since the search directions in Newton's method are obtained by minimizing quadratic polynomials, we start by examining when such polynomials have minimum values, and since convexity is a necessary condition for quadratic polynomials to be bounded below, we can restrict ourself to the study of convex quadratic polynomials.

**Theorem 15.1.1.** *A quadratic polynomial*

$$P(v) = \frac{1}{2}\langle v, Av \rangle + \langle b, v \rangle + c$$

in  $n$  variables, where  $A$  is a positive semidefinite symmetric operator, is bounded below on  $\mathbf{R}^n$  if and only if the equation

$$(15.1) \quad Av = -b$$

has a solution.

The polynomial has a minimum if it is bounded below, and  $\hat{v}$  is a minimum point if and only if  $A\hat{v} = -b$ .

If  $\hat{v}$  is a minimum point of the polynomial  $P$ , then

$$(15.2) \quad P(v) - P(\hat{v}) = \frac{1}{2}\langle v - \hat{v}, A(v - \hat{v}) \rangle$$

for all  $v \in \mathbf{R}^n$ .

If  $\hat{v}_1$  and  $\hat{v}_2$  are two minimum points, then  $\langle \hat{v}_1, A\hat{v}_1 \rangle = \langle \hat{v}_2, A\hat{v}_2 \rangle$ .

*Remark.* Another way to state that equation (15.1) has a solution is to say that the vector  $-b$ , and of course also the vector  $b$ , belongs to the range of the operator  $A$ . But the range of an operator on a finite dimensional space is equal to the orthogonal complement of the null space of the operator. Hence, equation (15.1) is solvable if and only if

$$Av = 0 \Rightarrow \langle b, v \rangle = 0.$$

**ie** business school

#1 EUROPEAN BUSINESS SCHOOL  
FINANCIAL TIMES 2013

#gobeyond

**MASTER IN MANAGEMENT**

**Because achieving your dreams is your greatest challenge.** IE Business School's Master in Management taught in English, Spanish or bilingually, trains young high performance professionals at the beginning of their career through an innovative and stimulating program that will help them reach their full potential.

- Choose your area of specialization.
- Customize your master through the different options offered.
- Global Immersion Weeks in locations such as London, Silicon Valley or Shanghai.

*Because you change, we change with you.*

www.ie.edu/master-management | mim.admissions@ie.edu |



*Proof.* First suppose that equation (15.1) has no solution. Then, by the remark above there exists a vector  $v$  such that  $Av = 0$  and  $\langle b, v \rangle \neq 0$ . It follows that

$$P(tv) = \frac{1}{2}\langle v, Av \rangle t^2 + \langle b, v \rangle t + c = \langle b, v \rangle t + c$$

for all  $t \in \mathbf{R}$ , and since the  $t$ -coefficient is nonzero, we conclude that the polynomial  $P(t)$  is unbounded below.

Next suppose that  $A\hat{v} = -b$ . Then

$$\begin{aligned} P(v) - P(\hat{v}) &= \frac{1}{2}(\langle v, Av \rangle - \langle \hat{v}, A\hat{v} \rangle) + \langle b, v \rangle - \langle b, \hat{v} \rangle \\ &= \frac{1}{2}(\langle v, Av \rangle - \langle \hat{v}, A\hat{v} \rangle) - \langle A\hat{v}, v \rangle + \langle A\hat{v}, \hat{v} \rangle \\ &= \frac{1}{2}(\langle v, Av \rangle + \langle \hat{v}, A\hat{v} \rangle - \langle A\hat{v}, v \rangle - \langle \hat{v}, Av \rangle) \\ &= \frac{1}{2}\langle v - \hat{v}, A(v - \hat{v}) \rangle \geq 0 \end{aligned}$$

for all  $v \in \mathbf{R}^n$ . This proves that the polynomial  $P(t)$  is bounded below, that  $\hat{v}$  is a minimum point, and that the equality (15.2) holds.

Since every positive semidefinite symmetric operator  $A$  has a unique positive semidefinite symmetric square root  $A^{1/2}$ , we can rewrite equality (15.2) as follows:

$$P(v) = P(\hat{v}) + \frac{1}{2}\langle A^{1/2}(v - \hat{v}), A^{1/2}(v - \hat{v}) \rangle = P(\hat{v}) + \frac{1}{2}\|A^{1/2}(v - \hat{v})\|^2.$$

If  $v$  is another minimum point of  $P$ , then  $P(v) = P(\hat{v})$ , and it follows that

$$A^{1/2}(v - \hat{v}) = 0.$$

Consequently,  $A(v - \hat{v}) = A^{1/2}(A^{1/2}(v - \hat{v})) = 0$ , i.e.  $Av = A\hat{v} = -b$ . Hence, every minimum point of  $P$  is obtained as a solution to equation (15.1).

Finally, if  $\hat{v}_1$  and  $\hat{v}_2$  are two minimum points of the polynomial, then  $A\hat{v}_1 = A\hat{v}_2 (= -b)$ , and it follows that  $\langle \hat{v}_1, A\hat{v}_1 \rangle = \langle \hat{v}_1, A\hat{v}_2 \rangle = \langle A\hat{v}_1, \hat{v}_2 \rangle = \langle A\hat{v}_2, \hat{v}_2 \rangle = \langle \hat{v}_2, A\hat{v}_2 \rangle$ .  $\square$

The problem to solve a convex quadratic optimization problem in  $\mathbf{R}^n$  is thus reduced to solving a quadratic system of linear equations in  $n$  variables (with a positive semidefinite coefficient matrix), which is a rather trivial numerical problem that can be performed with  $O(n^3)$  arithmetic operations.

We are now ready to define the main ingredients of Newton's method.

**Definition.** Let  $f: X \rightarrow \mathbf{R}$  be a twice differentiable function with an open subset  $X$  of  $\mathbf{R}^n$  as domain, and let  $x \in X$  be a point where the second derivative  $f''(x)$  is positive semidefinite.

By a *Newton direction*  $\Delta x_{\text{nt}}$  of the function  $f$  at the point  $x$  we mean a solution  $v$  to the equation

$$f''(x)v = -f'(x).$$

*Remark.* It follows from the remark after Theorem 15.1.1 that there exists a Newton direction at  $x$  if and only if

$$f''(x)v = 0 \Rightarrow \langle f'(x), v \rangle = 0.$$

The nonexistence of Newton directions at  $x$  is thus equivalent to the existence of a vector  $w$  such that  $f''(x)w = 0$  and  $\langle f'(x), w \rangle = 1$ .

The Newton direction  $\Delta x_{\text{nt}}$  is of course uniquely determined as

$$\Delta x_{\text{nt}} = -f''(x)^{-1}f'(x)$$

if the second derivative  $f''(x)$  is non-singular, i.e. positive definite.

A Newton direction  $\Delta x_{\text{nt}}$  is according to Theorem 15.1.1, whenever it exists, a minimizing vector for the Taylor polynomial

$$P(v) = f(x) + \langle f'(x), v \rangle + \frac{1}{2}\langle v, f''(x)v \rangle,$$

and the difference  $P(0) - P(\Delta x_{\text{nt}})$  is given by

$$P(0) - P(\Delta x_{\text{nt}}) = \frac{1}{2}\langle 0 - \Delta x_{\text{nt}}, f''(x)(0 - \Delta x_{\text{nt}}) \rangle = \frac{1}{2}\langle \Delta x_{\text{nt}}, f''(x)\Delta x_{\text{nt}} \rangle.$$

Using the Taylor approximation  $f(x+v) \approx P(v)$ , we conclude that

$$f(x) - f(x + \Delta x_{\text{nt}}) \approx P(0) - P(\Delta x_{\text{nt}}) = \frac{1}{2}\langle \Delta x_{\text{nt}}, f''(x)\Delta x_{\text{nt}} \rangle.$$

Hence,  $\frac{1}{2}\langle \Delta x_{\text{nt}}, f''(x)\Delta x_{\text{nt}} \rangle$  is (for small  $\Delta x_{\text{nt}}$ ) an approximation of the decrease in function value which is obtained by replacing  $f(x)$  with  $f(x + \Delta x_{\text{nt}})$ . This motivates our next definition.

**Definition.** The *Newton decrement*  $\lambda(f, x)$  of the function  $f$  at the point  $x$  is a quantity defined as

$$\lambda(f, x) = \sqrt{\langle \Delta x_{\text{nt}}, f''(x)\Delta x_{\text{nt}} \rangle}$$

if  $f$  has a Newton direction  $\Delta x_{\text{nt}}$  at  $x$ , and as

$$\lambda(f, x) = +\infty$$

if there is no Newton direction at  $x$ .

Note that the definition is independent of the choice of Newton direction at  $x$  in case of nonuniqueness of Newton direction. This follows immediately from the last statement in Theorem 15.1.1.

In terms of the Newton decrement, we thus have the following approximation

$$f(x) - f(x + \Delta x_{\text{nt}}) \approx \frac{1}{2}\lambda(f, x)^2$$

for small values of  $\Delta x_{\text{nt}}$ .

By definition  $f''(x)\Delta x_{\text{nt}} = -f'(x)$ , so it follows that the Newton decrement, whenever finite, can be computed using the formula

$$\lambda(f, x) = \sqrt{-\langle \Delta x_{\text{nt}}, f'(x) \rangle}.$$

In particular, if  $x$  is a point where the second derivative is positive definite, then

$$\lambda(f, x) = \sqrt{\langle f''(x)^{-1}f'(x), f'(x) \rangle}.$$

EXAMPLE 15.1.1. The convex one-variable function

$$f(x) = -\ln x, \quad x > 0$$

has Newton decrement

$$\lambda(f, x) = \sqrt{\langle x^2(-x^{-1}), -x^{-1} \rangle} = \sqrt{(-x) \cdot (-x^{-1})} = 1$$

at all points  $x > 0$ . □



**no.1**  
nine years  
in a row

Sweden  
Stockholm

## STUDY AT A TOP RANKED INTERNATIONAL BUSINESS SCHOOL

Reach your full potential at the Stockholm School of Economics, in one of the most innovative cities in the world. The School is ranked by the Financial Times as the number one business school in the Nordic and Baltic countries.

Visit us at [www.hhs.se](http://www.hhs.se)




At points  $x$  with a Newton direction it is also possible to express the Newton decrement in terms of the Euclidean norm  $\|\cdot\|$  as follows, by using the fact that  $f''(x)$  has a positive definite symmetric square root:

$$\lambda(f, x) = \sqrt{\langle f''(x)^{1/2} \Delta x_{\text{nt}}, f''(x)^{1/2} \Delta x_{\text{nt}} \rangle} = \|f''(x)^{1/2} \Delta x_{\text{nt}}\|.$$

The improvement in function value obtained by taking a step in the Newton direction  $\Delta x_{\text{nt}}$  is thus proportional to  $\|f''(x)^{1/2} \Delta x_{\text{nt}}\|^2$  and not to  $\|\Delta x_{\text{nt}}\|^2$ , a fact which motivates our introduction of the following seminorm.

**Definition.** Let  $f: X \rightarrow \mathbf{R}$  be a twice differentiable function with an open subset  $X$  of  $\mathbf{R}^n$  as domain, and let  $x \in X$  be a point where the second derivative  $f''(x)$  is positive semidefinite. The function  $\|\cdot\|_x: \mathbf{R}^n \rightarrow \mathbf{R}_+$ , defined by

$$\|v\|_x = \sqrt{\langle v, f''(x)v \rangle} = \|f''(x)^{1/2}v\|$$

for all  $v \in \mathbf{R}^n$ , is called the *local seminorm* at  $x$  of the function  $f$ .

It is easily verified that  $\|\cdot\|_x$  is indeed a seminorm on  $\mathbf{R}^n$ . Since

$$\{v \in \mathbf{R}^n \mid \|v\|_x = 0\} = \mathcal{N}(f''(x)),$$

where  $\mathcal{N}(f''(x))$  is the null space of  $f''(x)$ ,  $\|\cdot\|_x$  is a norm if and only if the positive definite second derivative  $f''(x)$  is nonsingular, i.e. positive definite.

At points  $x$  with a Newton direction, we now have the following simple relation between direction and decrement:

$$\lambda(f, x) = \|\Delta x_{\text{nt}}\|_x.$$

**EXAMPLE 15.1.2.** Let us study the Newton decrement  $\lambda(f, x)$  when  $f$  is a convex quadratic polynomial, i.e. a function of the form

$$f(x) = \frac{1}{2} \langle x, Ax \rangle + \langle b, x \rangle + c$$

with a positive semidefinite operator  $A$ . We have  $f'(x) = Ax + b$ ,  $f''(x) = A$  and  $\|v\|_x = \sqrt{\langle v, Av \rangle}$ , so the seminorms  $\|\cdot\|_x$  are the same for all  $x \in \mathbf{R}^n$ .

If  $\Delta x_{\text{nt}}$  is a Newton direction of  $f$  at  $x$ , then

$$A\Delta x_{\text{nt}} = -(Ax + b),$$

by definition, and it follows that  $A(x + \Delta x_{\text{nt}}) = -b$ . This implies that the function  $f$  is bounded below, according to Theorem 15.1.1.

So if  $f$  is not bounded below, then there are no Newton directions at any point  $x$ , which means that  $\lambda(f, x) = +\infty$  for all  $x$ .

Conversely, assume that  $f$  is bounded below. Then there exists a vector  $v_0$  such that  $Av_0 = -b$ , and it follows that

$$f''(x)(v_0 - x) = Av_0 - Ax = -b - Ax = -f'(x).$$

The vector  $v_0 - x$  is in other words a Newton direction of  $f$  at the point  $x$ , which means that the Newton decrement  $\lambda(f, x)$  is finite at all points  $x$  and is given by

$$\lambda(f, x) = \|v_0 - x\|_x.$$

If  $f$  is bounded below without being constant, then necessarily  $A \neq 0$  and we can choose a vector  $w$  such that  $\|w\|_x = \sqrt{\langle w, Aw \rangle} = 1$ . Let  $x_k = kw + v_0$ , where  $k$  is a positive number. Then

$$\lambda(f, x_k) = \|v_0 - x_k\|_{x_k} = k\|w\|_{x_k} = k,$$

and we conclude from this that  $\sup_{x \in \mathbf{R}^n} \lambda(f, x) = +\infty$ .

For constant functions  $f$ , the case  $A = 0$ ,  $b = 0$ , we have  $\|v\|_x = 0$  for all  $x$  and  $v$ , and consequently  $\lambda(f, x) = 0$  for all  $x$ .

In summary, we have obtained the following result:

The Newton decrement of downwards unbounded convex quadratic functions (which includes all non-constant affine functions) is infinite at all points. The Newton decrement of downwards bounded convex quadratic functions  $f$  is finite at all points, but  $\sup_x \lambda(f, x) = \infty$ , unless the function is constant.  $\square$

We shall give an alternative characterization of the Newton decrement, and for this purpose we need the following useful inequality.

**Theorem 15.1.2.** *Suppose  $\lambda(f, x) < \infty$ . Then*

$$|\langle f'(x), v \rangle| \leq \lambda(f, x) \|v\|_x$$

for all  $v \in \mathbf{R}^n$ .

*Proof.* Since  $\lambda(f, x)$  is assumed to be finite, there exists a Newton direction  $\Delta x_{\text{nt}}$  at  $x$ , and by definition,  $f''(x)\Delta x_{\text{nt}} = -f'(x)$ . Using the Cauchy-Schwarz inequality we now obtain:

$$\begin{aligned} |\langle f'(x), v \rangle| &= |\langle f''(x)\Delta x_{\text{nt}}, v \rangle| = |\langle f''(x)^{1/2}\Delta x_{\text{nt}}, f''(x)^{1/2}v \rangle| \\ &\leq \|f''(x)^{1/2}\Delta x_{\text{nt}}\| \|f''(x)^{1/2}v\| = \lambda(f, x) \|v\|_x. \end{aligned} \quad \square$$

**Theorem 15.1.3.** *Assume as before that  $x$  is a point where the second derivative  $f''(x)$  is positive semidefinite. Then*

$$\lambda(f, x) = \sup_{\|v\|_x \leq 1} \langle f'(x), v \rangle.$$

*Proof.* First assume that  $\lambda(f, x) < \infty$ . Then

$$\langle f'(x), v \rangle \leq \lambda(f, x)$$

for all vectors  $v$  such that  $\|v\|_x \leq 1$ , according to Theorem 15.1.2. In the case  $\lambda(f, x) = 0$  the above inequality holds with equality for  $v = 0$ , so assume that  $\lambda(f, x) > 0$ . For  $v = -\lambda(f, x)^{-1} \Delta x_{\text{nt}}$  we then have  $\|v\|_x = 1$  and

$$\langle f'(x), v \rangle = -\lambda(f, x)^{-1} \langle f'(x), \Delta x_{\text{nt}} \rangle = \lambda(f, x).$$

This proves that  $\lambda(f, x) = \sup_{\|v\|_x \leq 1} \langle f'(x), v \rangle$  for finite Newton decrements  $\lambda(f, x)$ .

Next assume that  $\lambda(f, x) = +\infty$ , i.e. that no Newton direction exists at  $x$ . By the remark after the definition of Newton direction, there exists a vector  $w$  such that  $f''(x)w = 0$  and  $\langle f'(x), w \rangle = 1$ . It follows that  $\|tw\|_x = t\|w\|_x = t\sqrt{\langle w, f''(x)w \rangle} = 0 \leq 1$  and  $\langle f'(x), tw \rangle = t$  for all positive numbers  $t$ , and this implies that  $\sup_{\|v\|_x \leq 1} \langle f'(x), v \rangle = +\infty = \lambda(f, x)$ .  $\square$

We sometimes need to compare  $\|\Delta x_{\text{nt}}\|$ ,  $\|f'(x)\|$  and  $\lambda(f, x)$ , and we can do so using the following theorem.

**#1**  
in eco-friendly  
attitude

**STUDY AT  
LINKÖPING UNIVERSITY, SWEDEN**  
RANKED AMONG TOP 50 UNIVERSITIES UNDER 50

Interested in Strategy and Management in International Organisations? Kick-start your career with a master's degree from Linköping University, Sweden.

→ **Click here!**

 **Linköping University**

**Theorem 15.1.4.** *Let  $\lambda_{\min}$  and  $\lambda_{\max}$  denote the smallest and the largest eigenvalue of the second derivative  $f''(x)$ , assumed to be positive semidefinite, and suppose that the Newton decrement  $\lambda(f, x)$  is finite. Then*

$$\lambda_{\min}^{1/2} \|\Delta x_{\text{nt}}\| \leq \lambda(f, x) \leq \lambda_{\max}^{1/2} \|\Delta x_{\text{nt}}\|$$

and

$$\lambda_{\min}^{1/2} \lambda(f, x) \leq \|f'(x)\| \leq \lambda_{\max}^{1/2} \lambda(f, x).$$

*Proof.* Let  $A$  be an arbitrary positive semidefinite operator on  $\mathbf{R}^n$  with smallest and largest eigenvalue  $\mu_{\min}$  and  $\mu_{\max}$  respectively. Then

$$\mu_{\min} \|v\| \leq \|Av\| \leq \mu_{\max} \|v\|$$

for all vectors  $v$ .

Since  $\lambda_{\min}^{1/2}$  and  $\lambda_{\max}^{1/2}$  are the smallest and the largest eigenvalues of the operator  $f''(x)^{1/2}$ , we obtain the two inequalities of our theorem by applying the general inequality to  $A = f''(x)^{1/2}$  and  $v = \Delta x_{\text{nt}}$ , and to  $A = f''(x)^{1/2}$  and  $v = f''(x)^{1/2} \Delta x_{\text{nt}}$ , noting that  $\|f''(x)^{1/2} \Delta x_{\text{nt}}\| = \lambda(f, x)$  and that

$$\|f''(x)^{1/2} (f''(x)^{1/2} \Delta x_{\text{nt}})\| = \|f''(x) \Delta x_{\text{nt}}\| = \|f'(x)\|. \quad \square$$

Theorem 15.1.4 is a local result, but if the function  $f$  is  $\mu$ -strongly convex, then  $\lambda_{\min} \geq \mu$ , and if the norm of the second derivative is bounded by some constant  $M$ , then  $\lambda_{\max} = \|f''(x)\| \leq M$  for all  $x$  in the domain of  $f$ . Therefore, we get the following corollary to Theorem 15.1.4.

**Corollary 15.1.5.** *If  $f: X \rightarrow \mathbf{R}$  is a twice differentiable  $\mu$ -strongly convex function, then*

$$\mu^{1/2} \|\Delta x_{\text{nt}}\| \leq \lambda(f, x) \leq \mu^{-1/2} \|f'(x)\|$$

for all  $x \in X$ . If moreover  $\|f''(x)\| \leq M$ , then

$$M^{-1/2} \|f'(x)\| \leq \lambda(f, x) \leq M^{1/2} \|\Delta x_{\text{nt}}\|.$$

The distance from an arbitrary point to the minimum point of a strongly convex function with bounded second derivative can be estimated using the Newton decrement, because we have the following result.

**Theorem 15.1.6.** *Let  $f: X \rightarrow \mathbf{R}$  be a  $\mu$ -strongly convex function, and suppose that  $f$  has a minimum at the point  $\hat{x}$  and that  $\|f''(x)\| \leq M$  for all  $x \in X$ . Then*

$$f(x) - f(\hat{x}) \leq \frac{M}{2\mu} \lambda(f, x)^2$$

and

$$\|x - \hat{x}\| \leq \frac{\sqrt{M}}{\mu} \lambda(f, x).$$

*Proof.* The theorem follows by combining Theorem 14.1.1 with the estimate  $\|f'(x)\| \leq M^{1/2} \lambda(f, x)$  from Corollary 15.1.5.  $\square$

The Newton decrement is invariant under surjective affine coordinate transformations. A slightly more general result is the following.

**Theorem 15.1.7.** *Let  $f$  be a twice differentiable function whose domain  $\Omega$  is a subset of  $\mathbf{R}^n$ , let  $A: \mathbf{R}^m \rightarrow \mathbf{R}^n$  be an affine map, and let  $g = f \circ A$ . Let furthermore  $x = Ay$  be a point in  $\Omega$ , and suppose that the second derivative  $f''(x)$  is positive semidefinite. The second derivative  $g''(y)$  is then positive semidefinite, and the Newton decrements of the two functions  $g$  and  $f$  satisfy the inequality*

$$\lambda(g, y) \leq \lambda(f, x).$$

*Equality holds if the affine map  $A$  is surjective.*

*Proof.* The affine map can be written as  $Ay = Cy + b$ , where  $C$  is a linear map and  $b$  is a vector, and the chain rule gives us the identities

$$\langle g'(y), w \rangle = \langle f'(x), Cw \rangle \quad \text{and} \quad \langle w, g''(y)w \rangle = \langle Cw, f''(x)Cw \rangle$$

for arbitrary vectors  $w$  in  $\mathbf{R}^m$ . It follows from the latter identity that the second derivative  $g''(y)$  is positive semidefinite if  $f''(x)$  is so, and that

$$\|w\|_y = \|Cw\|_x.$$

An application of Theorem 15.1.3 now gives

$$\lambda(g, y) = \sup_{\|w\|_y \leq 1} \langle g'(y), w \rangle = \sup_{\|Cw\|_x \leq 1} \langle f'(x), Cw \rangle \leq \sup_{\|v\|_x \leq 1} \langle f'(x), v \rangle = \lambda(f, x).$$

If the affine map  $A$  is surjective, then  $C$  is a surjective linear map, and hence  $v = Cw$  runs through all of  $\mathbf{R}^n$  as  $w$  runs through  $\mathbf{R}^m$ . In this case, the only inequality in the above chain of equalities and inequalities becomes an equality, which means that  $\lambda(g, y) = \lambda(f, x)$ .  $\square$

## 15.2 Newton's method

### The algorithm

Newton's method for minimizing a twice differentiable function  $f$  is a descent method, in which the search direction in each iteration is given by the Newton



direction  $\Delta x_{\text{nt}}$  at the current point. The stopping criterion is formulated in terms of the Newton decrement; the algorithm stops when the decrement is sufficiently small. In short, therefore, the algorithm looks like this:

### Newton's method

**Given** a starting point  $x \in \text{dom } f$  and a tolerance  $\epsilon > 0$ .

#### Repeat

1. Compute a Newton direction  $\Delta x_{\text{nt}}$  and the Newton decrement  $\lambda(f, x)$  at  $x$ .
2. *Stopping criterion:* **stop** if  $\lambda(f, x)^2 \leq 2\epsilon$ .
3. Determine a step size  $h > 0$ .
4. *Update:*  $x := x + h\Delta x_{\text{nt}}$ .

The step size  $h$  is set equal to 1 in each iteration in the so-called *pure* Newton method, while it is computed by line search with Armijo's rule or otherwise in *damped* Newton methods.

The stopping criterion is motivated by the fact that  $\frac{1}{2}\lambda(f, x)^2$  is an approximation to the decrease  $f(x) - f(x + \Delta x_{\text{nt}})$  in function value, and if this decrease is small, it is not worthwhile to continue.



"I studied English for 16 years but...  
...I finally learned to speak it in just six lessons"

Jane, Chinese architect

ENGLISH OUT THERE

Click to hear me talking before and after my unique course download

Newton's method generally works well for functions which are convex in a neighborhood of the optimal point, but it breaks down, of course, if it hits a point where the second derivative is singular and the Newton direction is lacking. We shall show that the pure method, under appropriate conditions on the objective function  $f$ , converges to the minimum point if the starting point is sufficiently close to the minimum point. To achieve convergence for arbitrary starting points, it is necessary to use methods with damping.

**EXAMPLE 15.2.1.** When applied to a downwards bounded convex quadratic polynomial

$$f(x) = \frac{1}{2}\langle x, Ax \rangle + \langle b, x \rangle + c,$$

Newton's pure method finds the optimal solution after just one iteration, regardless of the choice of starting point  $x$ , because  $f'(x) = Ax + b$ ,  $f''(x) = A$  and  $A\Delta x_{\text{nt}} = -(Ax + b)$ , so the update  $x^+ = x + \Delta x_{\text{nt}}$  satisfies the equation

$$f'(x^+) = Ax^+ + b = Ax + A\Delta x_{\text{nt}} + b = 0,$$

which means that  $x^+$  is the optimal point. □

## Invariance under change of coordinates

Unlike the gradient descent method, Newton's method is invariant under affine coordinate changes.

**Theorem 15.2.1.** *Let  $f: X \rightarrow \mathbf{R}$  be a twice differentiable function with a positive definite second derivative, and let  $(x_k)_0^\infty$  be the sequence generated by Newton's pure algorithm with  $x_0$  as starting point. Let further  $A: Y \rightarrow X$  be an affine coordinate transformation, i.e. the restriction to  $Y$  of a bijective affine map. Newton's pure algorithm applied to the function  $g = f \circ A$  with  $y_0 = A^{-1}x_0$  as the starting point then generates a sequence  $(y_k)_0^\infty$  with the property that  $Ay_k = x_k$  for each  $k$ .*

*The two sequences have identical Newton decrements in each iteration, and they therefore satisfy the stopping condition during the same iteration.*

*Proof.* The assertion about the Newton decrements follows from Theorem 15.1.7, and the relationship between the two sequences follows by induction if we show that  $Ay = x$  implies that  $A(y + \Delta y_{\text{nt}}) = x + \Delta x_{\text{nt}}$ , where  $\Delta x_{\text{nt}} = -f''(x)^{-1}f'(x)$  and  $\Delta y_{\text{nt}} = -g''(y)^{-1}g'(y)$  are the uniquely defined Newton directions at the points  $x$  and  $y$  of the respective functions.

The affine map  $A$  can be written as  $Ay = Cy + b$ , where  $C$  is an invertible linear map and  $b$  is a vector. If  $x = Ay$ , then  $g'(y) = C^T f'(x)$  and  $g''(y) =$

$C^T f''(x)C$ , by the chain rule. It follows that

$$\begin{aligned} C\Delta y_{\text{nt}} &= -Cg''(y)^{-1}g'(y) = -CC^{-1}f''(x)^{-1}(C^T)^{-1}C^T f'(x) \\ &= -f''(x)^{-1}f'(x) = \Delta x_{\text{nt}}, \end{aligned}$$

and hence

$$A(y + \Delta y_{\text{nt}}) = C(y + \Delta y_{\text{nt}}) + b = Cy + b + C\Delta y_{\text{nt}} = Ay + \Delta x_{\text{nt}} = x + \Delta x_{\text{nt}}. \quad \square$$

## Local convergence

We will now study convergence properties for the Newton method, starting with the pure method.

**Theorem 15.2.2.** *Let  $f: X \rightarrow \mathbf{R}$  be a twice differentiable,  $\mu$ -strongly convex function with minimum point  $\hat{x}$ , and suppose that the second derivative  $f''$  is Lipschitz continuous with Lipschitz constant  $L$ . Let  $x$  be a point in  $X$  and set*

$$x^+ = x + \Delta x_{\text{nt}},$$

where  $\Delta x_{\text{nt}}$  is the Newton direction at  $x$ . Then

$$\|x^+ - \hat{x}\| \leq \frac{L}{2\mu} \|x - \hat{x}\|^2.$$

Moreover, if the point  $x^+$  lies in  $X$  then

$$\|f'(x^+)\| \leq \frac{L}{2\mu^2} \|f'(x)\|^2.$$

*Proof.* The smallest eigenvalue of the second derivative  $f''(x)$  is greater than or equal to  $\mu$  by Theorem 7.3.2 in Part I. Hence,  $f''(x)$  is invertible and the largest eigenvalue of  $f''(x)^{-1}$  is less than or equal to  $\mu^{-1}$ , and it follows that

$$(15.3) \quad \|f''(x)^{-1}\| \leq \mu^{-1}.$$

To estimate the norm of  $x^+ - \hat{x}$ , we rewrite the difference as

$$(15.4) \quad \begin{aligned} x^+ - \hat{x} &= x + \Delta x_{\text{nt}} - \hat{x} = x - \hat{x} - f''(x)^{-1}f'(x) \\ &= f''(x)^{-1}(f''(x)(x - \hat{x}) - f'(x)) = -f''(x)^{-1}w \end{aligned}$$

with

$$w = f'(x) - f''(x)(x - \hat{x}).$$

For  $0 \leq t \leq 1$  we then define the vektor  $w(t)$  as

$$w(t) = f'(\hat{x} + t(x - \hat{x})) - tf''(x)(x - \hat{x}),$$

and note that  $w = w(1) - w(0)$ , since  $f'(\hat{x}) = 0$ . By the chain rule,

$$w'(t) = (f''(\hat{x} + t(x - \hat{x})) - f''(x))(x - \hat{x}),$$

and by using the Lipschitz continuity of the second derivative, we obtain the estimate

$$\begin{aligned} \|w'(t)\| &\leq \|f''(\hat{x} + t(x - \hat{x})) - f''(x)\| \|x - \hat{x}\| \\ &\leq L\|\hat{x} + t(x - \hat{x}) - x\| \|x - \hat{x}\| = L(1 - t)\|x - \hat{x}\|^2. \end{aligned}$$

Now integrate the above inequality over the interval  $[0, 1]$ ; this results in the inequality

$$\begin{aligned} (15.5) \quad \|w\| &= \left\| \int_0^1 w'(t) dt \right\| \leq \int_0^1 \|w'(t)\| dt \leq L\|x - \hat{x}\|^2 \int_0^1 (1 - t) dt \\ &= \frac{1}{2}L\|x - \hat{x}\|^2. \end{aligned}$$

By combining equality (15.4) with the inequalities (15.3) and (15.5) we obtain the estimate

$$\|x^+ - \hat{x}\| = \|f''(x)^{-1}w\| \leq \|f''(x)^{-1}\| \|w\| \leq \frac{L}{2\mu} \|x - \hat{x}\|^2,$$

which is the first claim of the theorem.

Excellent Economics and Business programmes at:



**university of  
 groningen**



**“The perfect start  
 of a successful,  
 international career.”**

**CLICK HERE**  
 to discover why both socially  
 and academically the University  
 of Groningen is one of the best  
 places for a student to be

[www.rug.nl/feb/education](http://www.rug.nl/feb/education)

To prove the second claim, we assume that  $x^+$  lies in  $X$  and consider for  $0 \leq t \leq 1$  the vectors

$$v(t) = f'(x + t\Delta x_{\text{nt}}) - t f''(x)\Delta x_{\text{nt}},$$

noting that

$$v(1) - v(0) = f'(x^+) - f''(x)\Delta x_{\text{nt}} - f'(x) = f'(x^+) + f'(x) - f'(x) = f'(x^+).$$

Since  $v'(t) = (f''(x + t\Delta x_{\text{nt}}) - f''(x))\Delta x_{\text{nt}}$ , it follows from the Lipschitz continuity that

$$\|v'(t)\| \leq \|f''(x + t\Delta x_{\text{nt}}) - f''(x)\| \|\Delta x_{\text{nt}}\| \leq L\|\Delta x_{\text{nt}}\|^2 t,$$

and by integrating this inequality, we obtain the desired estimate

$$\|f'(x^+)\| = \left\| \int_0^1 v'(t) dt \right\| \leq \int_0^1 \|v'(t)\| dt \leq \frac{L}{2} \|\Delta x_{\text{nt}}\|^2 \leq \frac{L}{2\mu^2} \|f'(x)\|^2,$$

where the last inequality follows from Corollary 15.1.5.  $\square$

One consequence of the previous theorem is that the pure Newton method converges quadratically when applied to functions with a positive definite second derivative that does not vary too rapidly in a neighborhood of the minimum point, provided that the starting point is chosen sufficiently close to the minimum point. More precisely, the following holds:

**Theorem 15.2.3.** *Let  $f: X \rightarrow \mathbf{R}$  be a twice differentiable,  $\mu$ -strongly convex function with minimum point  $\hat{x}$ , and suppose that the second derivative  $f''$  is Lipschitz continuous with Lipschitz constant  $L$ . Let  $0 < r \leq 2\mu/L$  and suppose that the open ball  $B(\hat{x}; r)$  is included in  $X$ .*

*Newton's pure method with starting point  $x_0 \in B(\hat{x}; r)$  will then generate a sequence  $(x_k)_0^\infty$  of points in  $\Omega$  such that*

$$\|x_k - \hat{x}\| \leq \frac{2\mu}{L} \left( \frac{L}{2\mu} \|x_0 - \hat{x}\| \right)^{2^k}$$

*for all  $k$ , and the sequence therefore converges to the minimum point  $\hat{x}$  as  $k \rightarrow \infty$ .*

The convergence is very rapid. For example,

$$\|x_k - \hat{x}\| \leq \frac{2\mu}{L} 2^{-2^k}$$

if the initial point is chosen such that  $\|x_0 - \hat{x}\| \leq \mu/L$ , and this implies that  $\|x_k - \hat{x}\| \leq 10^{-19} \mu/L$  already for  $k = 6$ .

*Proof.* We keep the notation of Theorem 15.2.2 and then have  $x_{k+1} = x_k^+$ , so if  $x_k$  lies in the ball  $B(\hat{x}; r)$ , then

$$(15.6) \quad \|x_{k+1} - \hat{x}\| \leq \frac{L}{2\mu} \|x_k - \hat{x}\|^2,$$

and this implies that  $\|x_{k+1} - \hat{x}\| < Lr^2/2\mu \leq r$ , i.e. the point  $x_{k+1}$  lies in the ball  $B(\hat{x}; r)$ . By induction, all points in the sequence  $(x_k)_0^\infty$  lie in  $B(\hat{x}; r)$ , and we obtain the inequality of the theorem by repeated application of inequality (15.6).  $\square$

## Global convergence

Newton's damped method converges, under appropriate conditions on the objective function, for arbitrary starting points. The damping is required only during an initial phase, because the step size becomes 1 once the algorithm has produced a point where the gradient is sufficiently small. The convergence is quadratic during this second stage.

The following theorem describes a convergence result for strongly convex functions with Lipschitz continuous second derivative.

**Theorem 15.2.4.** *Let  $f: X \rightarrow \mathbf{R}$  be a twice differentiable, strongly convex function with a Lipschitz continuous second derivative. Let  $x_0$  be a point in  $X$  and suppose that the sublevel set*

$$S = \{x \in X \mid f(x) \leq f(x_0)\}$$

*is closed.*

*Then,  $f$  has a unique minimum point  $\hat{x}$ , and Newton's damped algorithm, with  $x_0$  as initial point and with line search according to Armijo's rule with parameters  $0 < \alpha < \frac{1}{2}$  and  $0 < \beta < 1$ , generates a sequence  $(x_k)_0^\infty$  of points in  $S$  that converges towards the minimum point.*

*After an initial phase with damping, the algorithm passes into a quadratically convergent phase with step size 1.*

*Proof.* The existence of a unique minimum point is a consequence of Corollary 8.1.7 in Part I.

Suppose that  $f$  is  $\mu$ -strongly convex and let  $L$  be the Lipschitz constant of the second derivative. The sublevel set  $S$  is compact since it is bounded according to Theorem 8.1.6. It follows that the distance from the set  $S$  to the boundary of the open set  $X$  is positive. Fix a positive number  $r$  that is less than this distance and also satisfies the inequality

$$r \leq \mu/L.$$

Given  $x \in S$  we now define the point  $x^+$  by

$$x^+ = x + h\Delta x_{\text{nt}},$$

where  $h$  is the step size according to Armijo's rule. In particular,  $x_{k+1} = x_k^+$  for all  $k$ .


The core of the proof consists in showing that there are two positive constants  $\gamma$  and  $\eta \leq \mu r$  such that the following two implications hold for all  $x \in S$ :

- (i)  $\|f'(x)\| \geq \eta \Rightarrow f(x^+) - f(x) \leq -\gamma$ ;
- (ii)  $\|f'(x)\| < \eta \Rightarrow h = 1 \ \& \ \|f'(x^+)\| < \eta$ .

Suppose that we have managed to prove (i) and (ii). If  $\|f'(x_k)\| \geq \eta$  for  $0 \leq k < m$ , then

$$f_{\min} - f(x_0) \leq f(x_m) - f(x_0) = \sum_{k=0}^{m-1} (f(x_k^+) - f(x_k)) \leq -m\gamma,$$

because of property (i). This inequality can not hold for all  $m$ , and hence there is a smallest integer  $k_0$  such that  $\|f'(x_{k_0})\| < \eta$ , and this integer must



In the past four years we have drilled  
**89,000 km**  
That's more than **twice** around the world.

**Who are we?**  
We are the world's largest oilfield services company<sup>1</sup>. Working globally—often in remote and challenging locations—we invent, design, engineer, and apply technology to help our customers find and produce oil and gas safely.

**Who are we looking for?**  
Every year, we need thousands of graduates to begin dynamic careers in the following domains:

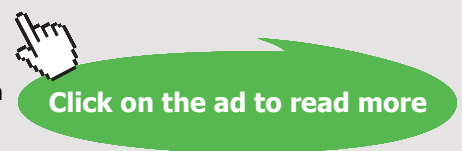
- **Engineering, Research and Operations**
- **Geoscience and Petrotechnical**
- **Commercial and Business**

**What will you be?**

**Schlumberger**

[careers.slb.com](http://careers.slb.com)

<sup>1</sup>Based on Fortune 500 ranking 2011. Copyright © 2015 Schlumberger. All rights reserved.



satisfy the inequality

$$k_0 \leq (f(x_0) - f_{\min})/\gamma.$$

It now follows by induction from (ii) that the step size  $h$  is equal to 1 for all  $k \geq k_0$ . The damped Newton algorithm is in other words a pure Newton algorithm from iteration  $k_0$  and onwards. Because of Theorem 14.1.1,

$$\|x_{k_0} - \hat{x}\| \leq \mu^{-1} \|f'(x_{k_0})\| < \mu^{-1} \eta \leq r \leq \mu L^{-1},$$

so it follows from Theorem 15.2.3 that the sequence  $(x_k)_0^\infty$  converges to  $\hat{x}$ , and more precisely, that the estimate

$$\|x_{k+k_0} - \hat{x}\| \leq \frac{2\mu}{L} \left( \frac{L}{2\mu} \|x_{k_0} - \hat{x}\| \right)^{2^k} \leq \frac{2\mu}{L} 2^{-2^k}$$

holds for  $k \geq 0$ .

It thus only remains to prove the existence of numbers  $\eta$  and  $\gamma$  with the properties (i) and (ii). To this end, let

$$S_r = S + \overline{B}(x; r);$$

the set  $S_r$  is a convex and compact subset of  $\Omega$ , and the two continuous functions  $f'$  and  $f''$  are therefore bounded on  $S_r$ , i.e. there are constants  $K$  and  $M$  such that

$$\|f'(x)\| \leq K \quad \text{and} \quad \|f''(x)\| \leq M$$

for all  $x \in S_r$ . It follows from Theorem 7.4.1 in Part I that the derivative  $f'$  is Lipschitz continuous on the set  $S_r$  with Lipschitz constant  $M$ , i.e.

$$\|f'(y) - f'(x)\| \leq M \|y - x\|$$

for  $x, y \in S_r$ .

We now define our numbers  $\eta$  and  $\gamma$  as

$$\eta = \min \left\{ \frac{3(1-2\alpha)\mu^2}{L}, \mu r \right\} \quad \text{and} \quad \gamma = \frac{\alpha\beta c\mu}{M} \eta^2, \quad \text{where } c = \min \left\{ \frac{1}{M}, \frac{r}{K} \right\}.$$

Let us first estimate the stepsize at a given point  $x \in S$ . Since

$$\|\Delta x_{\text{nt}}\| \leq \mu^{-1} \|f'(x)\| \leq \mu^{-1} K,$$

the point  $x + t\Delta x_{\text{nt}}$  lies in  $S_r$  and especially also in  $X$  if  $0 \leq t \leq r\mu K^{-1}$ . The function

$$g(t) = f(x + t\Delta x_{\text{nt}})$$



is therefore defined for these  $t$ -values, and since  $f$  is  $\mu$ -strongly convex and the derivative is Lipschitz continuous with constant  $M$  on  $S_r$ , it follows from Theorem 1.1.2 in Part I and Corollary 15.1.5 that

$$\begin{aligned} f(x + t\Delta x_{\text{nt}}) &\leq f(x) + t\langle f'(x), \Delta x_{\text{nt}} \rangle + \frac{1}{2}M\|\Delta x_{\text{nt}}\|^2 t^2 \\ &\leq f(x) + t\langle f'(x), \Delta x_{\text{nt}} \rangle + \frac{1}{2}M\mu^{-1}\lambda(f, x)^2 t^2 \\ &= f(x) + t\left(1 - \frac{1}{2}M\mu^{-1}t\right)\langle f'(x), \Delta x_{\text{nt}} \rangle. \end{aligned}$$

The number  $\hat{t} = c\mu$  lies in the interval  $[0, r\mu K^{-1}]$  and is less than or equal to  $\mu M^{-1}$ . Hence,  $1 - \frac{1}{2}M\mu^{-1}\hat{t} \geq \frac{1}{2} \geq \alpha$ , which inserted in the above inequality gives

$$f(x + \hat{t}\Delta x_{\text{nt}}) \leq f(x) + \alpha\hat{t}\langle f'(x), \Delta x_{\text{nt}} \rangle.$$

Now, let  $h = \beta^m$  be the step size given by Armijo's rule, which means that the Armijo algorithm terminates in iteration  $m$ . Since it does not terminate in iteration  $m - 1$ , we conclude that  $\beta^{m-1} > \hat{t}$ , i.e.

$$h \geq \beta\hat{t} = \beta c\mu,$$

and this gives us the following estimate for the point  $x^+ = x + h\Delta x_{\text{nt}}$ :

$$\begin{aligned} f(x^+) - f(x) &\leq \alpha h\langle f'(x), \Delta x_{\text{nt}} \rangle = -\alpha h\lambda(f, x)^2 \\ &\leq -\alpha\beta c\mu\lambda(f, x)^2 \leq -\alpha\beta c\mu M^{-1}\|f'(x)\|^2 = -\gamma\eta^{-2}\|f'(x)\|^2. \end{aligned}$$

So, if  $\|f'(x)\| \geq \eta$  then  $f(x^+) - f(x) \leq -\gamma$ , which is the content of implication (i).

To prove the remaining implication (ii), we return to the function  $g(t) = f(x + t\Delta x_{\text{nt}})$ , assuming that  $\|f'(x)\| < \eta$ . The function is well-defined for  $0 \leq t \leq 1$ , since

$$\|\Delta x_{\text{nt}}\| \leq \mu^{-1}\|f'(x)\| < \mu^{-1}\eta \leq r.$$

Moreover,

$$g'(t) = \langle f'(x + t\Delta x_{\text{nt}}), \Delta x_{\text{nt}} \rangle \text{ and } g''(t) = \langle \Delta x_{\text{nt}}, f''(x + t\Delta x_{\text{nt}})\Delta x_{\text{nt}} \rangle.$$

By Lipschitz continuity,

$$\begin{aligned} |g''(t) - g''(0)| &= |\langle \Delta x_{\text{nt}}, f''(x + t\Delta x_{\text{nt}})\Delta x_{\text{nt}} \rangle - \langle \Delta x_{\text{nt}}, f''(x)\Delta x_{\text{nt}} \rangle| \\ &\leq \|f''(x + t\Delta x_{\text{nt}}) - f''(x)\| \|\Delta x_{\text{nt}}\|^2 \leq tL\|\Delta x_{\text{nt}}\|^3, \end{aligned}$$

and it follows, since  $g''(0) = \lambda(f, x)^2$  and  $\|\Delta x_{\text{nt}}\| \leq \mu^{-1/2}\lambda(f, x)$ , that

$$g''(t) \leq \lambda(f, x)^2 + tL\|\Delta x_{\text{nt}}\|^3 \leq \lambda(f, x)^2 + tL\mu^{-3/2}\lambda(f, x)^3.$$

By integrating this inequality over the interval  $[0, t]$ , we obtain the inequality

$$g'(t) - g'(0) \leq t\lambda(f, x)^2 + \frac{1}{2}t^2L\mu^{-3/2}\lambda(f, x)^3.$$

But  $g'(0) = \langle f'(x), \Delta x_{\text{nt}} \rangle = -\lambda(f, x)^2$ , so it follows that

$$g'(t) \leq -\lambda(f, x)^2 + t\lambda(f, x)^2 + \frac{1}{2}t^2L\mu^{-3/2}\lambda(f, x)^3,$$

and further integration results in the inequality

$$g(t) - g(0) \leq -t\lambda(f, x)^2 + \frac{1}{2}t^2\lambda(f, x)^2 + \frac{1}{6}t^3L\mu^{-3/2}\lambda(f, x)^3.$$

Now, take  $t = 1$  to obtain

$$\begin{aligned} (15.7) \quad f(x + \Delta x_{\text{nt}}) &\leq f(x) - \frac{1}{2}\lambda(f, x)^2 + \frac{1}{6}L\mu^{-3/2}\lambda(f, x)^3 \\ &= f(x) - \lambda(f, x)^2\left(\frac{1}{2} - \frac{1}{6}L\mu^{-3/2}\lambda(f, x)\right) \\ &= f(x) + \langle f'(x), \Delta x_{\text{nt}} \rangle\left(\frac{1}{2} - \frac{1}{6}L\mu^{-3/2}\lambda(f, x)\right). \end{aligned}$$

Our assumption  $\|f'(x)\| < \eta$  implies that

$$\lambda(f, x) \leq \mu^{-1/2}\|f'(x)\| < \mu^{-1/2}\eta \leq \mu^{-1/2} \cdot 3(1-2\alpha)\mu^2L^{-1} = 3(1-2\alpha)\mu^{3/2}L^{-1}.$$

## American online LIGS University

is currently enrolling in the  
Interactive Online **BBA, MBA, MSc,**  
**DBA and PhD** programs:

- ▶ enroll **by September 30th, 2014** and
- ▶ **save up to 16%** on the tuition!
- ▶ pay in 10 installments / 2 years
- ▶ Interactive Online education
- ▶ visit [www.ligsuniversity.com](http://www.ligsuniversity.com) to  
find out more!

**Note:** LIGS University is not accredited by any  
nationally recognized accrediting agency listed  
by the US Secretary of Education.  
More info [here](#).



We conclude that

$$\frac{1}{2} - \frac{1}{6}L\mu^{-3/2}\lambda(f, x) > \alpha,$$

which inserted into inequality (15.7) gives us the inequality

$$f(x + \Delta x_{\text{nt}}) \leq f(x) + \alpha \langle f'(x), \Delta x_{\text{nt}} \rangle,$$

which tells us that the step size  $h$  is equal to 1.

The iteration leading from  $x$  to  $x^+ = x + h\Delta x_{\text{nt}}$  is therefore performed according to the pure Newton method. Due to the inequality

$$\|x - \hat{x}\| \leq \mu^{-1}\|f'(x)\| < \mu^{-1}\eta \leq r,$$

which holds by Theorem 14.1.1,  $x$  is a point in the ball  $B(\hat{x}; r)$ , so it follows from the local convergence Theorem 15.2.2 that

$$(15.8) \quad \|f'(x^+)\| \leq \frac{L}{2\mu^2}\|f'(x)\|^2.$$

Since  $\eta \leq \mu r \leq \mu^2/L$ ,

$$\|f'(x^+)\| < \frac{L}{2\mu^2}\eta^2 \leq \frac{\eta}{2} < \eta,$$

and the proof is now complete.  $\square$

By iterating inequality (15.8), one obtains in fact the estimate

$$\|f'(x_k)\| \leq \frac{2\mu^2}{L} \left( \frac{L}{2\mu^2} \|f'(x_{k_0})\| \right)^{2^{k-k_0}} < \frac{2\mu^2}{L} 2^{-2^{k-k_0}}$$

for  $k \geq k_0$ , and it now follows from Theorem 14.1.1 that

$$f(x_k) - f_{\min} < \frac{2\mu^3}{L^2} 2^{-2^{k-k_0+1}}$$

for  $k \geq k_0$ . Combining this estimate with the previously obtained bound on  $k_0$ , one obtains an upper bound on the number of iterations required to estimate the minimum value  $f_{\min}$  with a given accuracy. If

$$k > \frac{f(x_0) - f_{\min}}{\gamma} + \log_2 \log_2 \frac{2\mu^3}{L^2\epsilon},$$

then surely  $f(x_k) - f_{\min} < \epsilon$ . This estimate, however, is of no practical value, because the constants  $\gamma$ ,  $\mu$  and  $L$  are rarely known in concrete cases.

Another shortcoming of the classical convergence analysis of Newton's method is that the convergence constants, unlike the algorithm itself, depend on the coordinate system used. For self-concordant functions, it is however possible to carry out the convergence analysis without any unknown constants, as we shall do in Chapter 16.5.

### 15.3 Equality constraints

With only minor modifications, Newton's algorithm also works well when applied to convex optimization problems with constraints in the form of affine equalities.

Consider the convex optimization problem

$$(P) \quad \begin{array}{ll} \min & f(x) \\ \text{s.t.} & Ax = b \end{array}$$

where  $f: \Omega \rightarrow \mathbf{R}$  is a twice continuously differentiable convex function,  $\Omega$  is an open subset of  $\mathbf{R}^n$ , and  $A$  is an  $m \times n$ -matrix.

The problem's Lagrange function  $L: \Omega \times \mathbf{R}^m \rightarrow \mathbf{R}$  is given by

$$L(x, y) = f(x) + (Ax - b)^T y = f(x) + x^T A^T y - b^T y,$$

and according to the Karush–Kuhn–Tucker theorem (Theorem 11.2.1 in Part II), a point  $\hat{x}$  in  $\Omega$  is an optimal solution if and only if there is a vector  $\hat{y} \in \mathbf{R}^m$  such that

$$(15.9) \quad \begin{cases} f'(\hat{x}) + A^T \hat{y} = 0 \\ A\hat{x} = b. \end{cases}$$

Therefore, the minimization problem (P) is equivalent to the problem of solving the system (15.9) of linear equations.

EXAMPLE 15.3.1. When  $f$  is a convex quadratic function of the form

$$f(x) = \frac{1}{2} \langle x, Px \rangle + \langle q, x \rangle + r,$$

the linear system (15.9) becomes

$$\begin{cases} P\hat{x} + A^T \hat{y} = -q \\ A\hat{x} = b, \end{cases}$$

and this is a quadratic system of linear equations with a symmetric coefficient matrix of order  $m + n$ . The system has a unique solution if  $\text{rank } A = m$  and  $\mathcal{N}(A) \cap \mathcal{N}(P) = \{0\}$ . See exercise 15.4. In particular, there is a unique solution if the matrix  $P$  is positive definite and  $\text{rank } A = m$ .  $\square$

We now return to the general convex minimization problem (P). Let  $X$  denote the set of feasible points, so that

$$X = \{x \in \Omega \mid Ax = b\}.$$

In optimization problems without any constraints, the descent direction  $\Delta x_{\text{nt}}$  at the point  $x$  is a vector which minimizes the Taylor polynomial of degree

two of the function  $f(x+v)$ , and the minimization is over all vectors  $v$  in  $\mathbf{R}^n$ . As a new point  $x^+$  with function value less than  $f(x)$  we select  $x^+ = x + h\Delta x_{\text{nt}}$  with a suitable step size  $h$ . In constrained problems, the new point  $x^+$  has to be a feasible point, of course, and this requires that  $A\Delta x_{\text{nt}} = 0$ . The minimization of the Taylor polynomial is therefore restricted to vectors  $v$  that satisfy the condition  $Av = 0$ , and this means that we have to modify our previous definition of Newton direction and decrement as follows for constrained optimization problems.

**Definition.** In the equality constrained minimization problem (P), a vector  $\Delta x_{\text{nt}}$  is called a *Newton direction* at the point  $x \in X$  if there exists a vector  $w \in \mathbf{R}^m$  such that

$$(15.10) \quad \begin{cases} f''(x)\Delta x_{\text{nt}} + A^T w = -f'(x) \\ A\Delta x_{\text{nt}} = 0. \end{cases}$$

The quantity

$$\lambda(f, x) = \sqrt{\langle \Delta x_{\text{nt}}, f''(x)\Delta x_{\text{nt}} \rangle}$$

is called the *Newton decrement*.

.....Alcatel-Lucent 

[www.alcatel-lucent.com/careers](http://www.alcatel-lucent.com/careers)

What if you could build your future and create the future?

One generation's transformation is the next's status quo. In the near future, people may soon think it's strange that devices ever had to be "plugged in." To obtain that status, there needs to be "The Shift".

It follows from Example 15.3.1 that the Newton direction  $\Delta x_{\text{nt}}$  (if it exists) is an optimal solution to the minimization problem

$$\begin{aligned} \min \quad & f(x) + \langle f'(x), v \rangle + \frac{1}{2} \langle v, f''(x)v \rangle \\ \text{s.t.} \quad & Av = 0. \end{aligned}$$

And if  $(\Delta x_{\text{nt}}, w)$  is a solution to the system (15.10), then

$$\begin{aligned} -\langle f'(x), \Delta x_{\text{nt}} \rangle &= \langle f''(x)\Delta x_{\text{nt}} + A^T w, \Delta x_{\text{nt}} \rangle \\ &= \langle f''(x)\Delta x_{\text{nt}}, \Delta x_{\text{nt}} \rangle + \langle w, A\Delta x_{\text{nt}} \rangle \\ &= \langle f''(x)\Delta x_{\text{nt}}, \Delta x_{\text{nt}} \rangle + \langle w, 0 \rangle = \langle \Delta x_{\text{nt}}, f''(x)\Delta x_{\text{nt}} \rangle, \end{aligned}$$

so it follows that

$$\lambda(f, x) = \sqrt{-\langle f'(x), \Delta x_{\text{nt}} \rangle},$$

just as for unconstrained problems.

The objective function is decreasing in the Newton direction, because

$$\left. \frac{d}{dt} f(x + t\Delta x_{\text{nt}}) \right|_{t=0} = \langle f'(x), \Delta x_{\text{nt}} \rangle = -\lambda(f, x)^2 \leq 0,$$

so  $\Delta x_{\text{nt}}$  is indeed a descent direction.

Let  $P(v)$  denote the Taylor polynomial of degree two of the function  $f(x + v)$ . Then

$$\begin{aligned} f(x) - f(x + \Delta x_{\text{nt}}) &\approx P(0) - P(\Delta x_{\text{nt}}) \\ &= -\langle f'(x), \Delta x_{\text{nt}} \rangle - \frac{1}{2} \langle \Delta x_{\text{nt}}, f''(x)\Delta x_{\text{nt}} \rangle = \frac{1}{2} \lambda(f, x)^2, \end{aligned}$$

just as in the unconstrained case.

With our modified definition of the Newton direction, we can now copy Newton's method verbatim for convex minimization problem of the type

$$\begin{aligned} \min \quad & f(x) \\ \text{s.t.} \quad & Ax = b. \end{aligned}$$

The algorithm looks like this:

### Newton's method

**Given** a starting point  $x \in \Omega$  satisfying the constraint  $Ax = b$ , and a tolerance  $\epsilon > 0$ .

### Repeat

1. Compute the Newton direction  $\Delta x_{\text{nt}}$  at  $x$  by solving the system of equations (15.10), and compute the Newton decrement  $\lambda(f, x)$ .
2. *Stopping criterion:* **stop** if  $\lambda(f, x)^2 \leq 2\epsilon$ .
3. Compute a step size  $h > 0$ .
4. *Update:*  $x := x + h\Delta x_{\text{nt}}$ .

## Elimination of constraints

An alternative approach to the optimization problem

$$(P) \quad \begin{array}{ll} \min & f(x) \\ \text{s.t.} & Ax = b, \end{array}$$

with  $x \in \Omega$  as implicit condition and  $r = \text{rank } A$ , is to solve the system of equations  $Ax = b$  and to express  $r$  variables as linear combinations of the remaining  $n - r$  variables. The former variables can then be eliminated from the objective function, and we obtain in this way an optimization problem in  $n - r$  variables without explicit constraints, a problem that can be attacked with Newton's method. We will describe this approach in more detail and compare it with the method above.

Suppose that the set  $X$  of feasible points is nonempty, choose a point  $a \in X$ , and select an affine parametrization

$$x = \xi(z), \quad z \in \tilde{\Omega}$$

of  $X$  with  $\xi(0) = a$ . Since  $\{x \in \mathbf{R}^n \mid Ax = b\} = a + \mathcal{N}(A)$ , we can write the parametrization as

$$\xi(z) = a + Cz$$

where  $C: \mathbf{R}^p \rightarrow \mathbf{R}^n$  is an injective linear map, whose range  $\mathcal{V}(C)$  coincides with the null space  $\mathcal{N}(A)$  of the map  $A$ , and  $p = n - \text{rank } A$ . The domain  $\tilde{\Omega} = \{z \in \mathbf{R}^p \mid a + Cz \in \Omega\}$  is an open convex subset of  $\mathbf{R}^p$ .

A practical way to construct the parametrization is of course to solve the system  $Ax = b$  by Gaussian elimination.

Let us finally define the function  $\tilde{f}: \tilde{\Omega} \rightarrow \mathbf{R}$  by setting  $\tilde{f}(z) = f(\xi(z))$ . The problem (P) is then equivalent to the convex optimization problem

$$(\tilde{P}) \quad \min \tilde{f}(z)$$

which has no explicit constraints.

Let  $\Delta x_{\text{nt}}$  be a Newton direction of the function  $f$  at the point  $x$ , i.e. a vector that satisfies the system (15.10) for a suitably chosen vector  $w$ . We will show that the function  $\tilde{f}$  has a corresponding Newton direction  $\Delta z_{\text{nt}}$  at the point  $z = \xi^{-1}(x)$ , and that  $\Delta x_{\text{nt}} = C\Delta z_{\text{nt}}$ .

Since  $A\Delta x_{\text{nt}} = 0$  and  $\mathcal{N}(A) = \mathcal{V}(C)$ , there is a unique vector  $v$  such that  $\Delta x_{\text{nt}} = Cv$ . By the chain rule,  $\tilde{f}'(z) = C^T f'(x)$  and  $\tilde{f}''(z) = C^T f''(x)C$ , so it follows from the first equation in the system (15.10) that

$$\begin{aligned} \tilde{f}''(z)v &= C^T f''(x)Cv = C^T f''(x)\Delta x_{\text{nt}} = -C^T f'(x) - C^T A^T w \\ &= -\tilde{f}'(z) - C^T A^T w. \end{aligned}$$

A general result from linear algebra tells us that  $\mathcal{N}(S) = \mathcal{V}(S^T)^\perp$  for arbitrary linear maps  $S$ . Applying this result to the maps  $C^T$  and  $A$ , and using that  $\mathcal{V}(C) = \mathcal{N}(A)$ , we obtain the equality

$$\mathcal{N}(C^T) = \mathcal{V}(C)^\perp = \mathcal{N}(A)^\perp = \mathcal{V}(A^T)^{\perp\perp} = \mathcal{V}(A^T),$$

which implies that  $C^T A^T w = 0$ . Hence,

$$\tilde{f}''(z)v = -\tilde{f}'(z),$$

and  $v$  is thus a Newton direction of the function  $\tilde{f}$  at the point  $z$ . So,  $\Delta z_{\text{nt}} = v$  is the direction vector we are looking for.

The iteration step  $z \rightarrow z^+ = z + h\Delta z_{\text{nt}}$  in Newton's method for the unconstrained problem  $(\tilde{P})$  takes us from the point  $z = \xi^{-1}(x)$  in  $\tilde{\Omega}$  to the point  $z^+$  whose image in  $X$  is

$$\begin{aligned} \xi(z^+) &= \xi(z + h\Delta z_{\text{nt}}) = a + C(z + h\Delta z_{\text{nt}}) = a + Cz + hC(\Delta z_{\text{nt}}) \\ &= \xi(z) + h\Delta x_{\text{nt}} = x + h\Delta x_{\text{nt}}, \end{aligned}$$

and this is also the point we get by applying Newton's method to the point  $x$  in the constrained problem  $(P)$ .



**Empowering People. Improving Business.**

BI Norwegian Business School is one of Europe's largest business schools welcoming more than 20,000 students. Our programmes provide a stimulating and multi-cultural learning environment with an international outlook ultimately providing students with professional skills to meet the increasing needs of businesses.

BI offers four different two-year, full-time Master of Science (MSc) programmes that are taught entirely in English and have been designed to provide professional skills to meet the increasing need of businesses. The MSc programmes provide a stimulating and multi-cultural learning environment to give you the best platform to launch into your career.

- MSc in Business
- MSc in Financial Economics
- MSc in Strategic Marketing Management
- MSc in Leadership and Organisational Psychology

**BI NORWEGIAN BUSINESS SCHOOL**

EFMD **EQUIS ACCREDITED**

[www.bi.edu/master](http://www.bi.edu/master)



Also note that the Newton decrements are the same at corresponding points, because

$$\begin{aligned}\lambda(\tilde{f}, z)^2 &= -\langle \tilde{f}'(z), \Delta z_{\text{nt}} \rangle = -\langle C^T f'(x), \Delta z_{\text{nt}} \rangle = -\langle f'(x), C \Delta z_{\text{nt}} \rangle \\ &= -\langle f'(x), \Delta x_{\text{nt}} \rangle = \lambda(f, x)^2.\end{aligned}$$

In summary, we have arrived at the following result.

**Theorem 15.3.1.** *Let  $(x_k)_0^\infty$  be a sequence of points obtained by Newton's method applied to the constrained problem (P). Newton's method applied to the problem  $(\tilde{P})$ , obtained by elimination of the constraints and with  $\xi^{-1}(x_0)$  as initial point, will then generate a sequence  $(z_k)_0^\infty$  with the property that  $x_k = \xi(z_k)$  for all  $k$ .*

## Convergence analysis

No new convergence analysis is needed for the modified version of Newton's method, for we can, because of Theorem 15.3.1, apply the results of Theorem 15.2.4. If the restriction of the function  $f: \Omega \rightarrow \mathbf{R}$  to the set  $X$  of feasible points is strongly convex and the second derivative is Lipschitz continuous, then the same also holds for the function  $\tilde{f}: \tilde{\Omega} \rightarrow \mathbf{R}$ . (Cf. with exercise 15.5.) Assuming  $x_0$  to be a feasible starting point and the sublevel set  $\{x \in X \mid f(x) \leq f(x_0)\}$  to be closed, the damped Newton algorithm will therefore converge to the minimum point when applied to the constrained problem (P). Close enough to the minimum point, the step size  $h$  will also be equal to 1, and the convergence will be quadratic.

## Exercises

- 15.1** Determine the Newton direction, the Newton decrement and the local norm at an arbitrary point  $x > 0$  for the function  $f(x) = x \ln x - x$ .
- 15.2** Let  $f$  be the function  $f(x_1, x_2) = -\ln x_1 - \ln x_2 - \ln(4 - x_1 - x_2)$  with  $X = \{x \in \mathbf{R}^2 \mid x_1 > 0, x_2 > 0, x_1 + x_2 < 4\}$  as domain. Determine the Newton direction, the Newton decrement and the local norm at the point  $x$  when  
a)  $x = (1, 1)$       b)  $x = (1, 2)$ .
- 15.3** Determine a Newton direction, the Newton decrement and the local norm for the function  $f(x_1, x_2) = e^{x_1+x_2} + x_1 + x_2$  at an arbitrary point  $x \in \mathbf{R}^2$ .
- 15.4** Assume that  $P$  is a symmetric positive semidefinite  $n \times n$ -matrix and that

$A$  is an arbitrary  $m \times n$ -matrix. Prove that the matrix

$$M = \begin{bmatrix} P & A^T \\ A & 0 \end{bmatrix}$$

is invertible if and only if  $\text{rank } A = m$  and  $\mathcal{N}(A) \cap \mathcal{N}(P) = \{0\}$ .

**15.5** Assume that the function  $f: \Omega \rightarrow \mathbf{R}$  is twice differentiable and convex, let  $x = \xi(z) = a + Cz$  be an affine parametrization of the set

$$X = \{x \in \Omega \mid Ax = b\},$$

and define the function  $\tilde{f}$  by  $\tilde{f}(z) = f(\xi(z))$ , just as in Section 15.3. Let further  $\sigma$  denote the smallest eigenvalue of the symmetric matrix  $C^T C$ .

a) Prove that  $\tilde{f}$  is  $\mu\sigma$ -strongly convex if the restriction of  $f$  to  $X$  is  $\mu$ -strongly convex.

b) Assume that the matrix  $A$  has full rank and that there are constants  $K$  and  $M$  such that  $Ax = b$  implies

$$\left\| \begin{bmatrix} f''(x) & A^T \\ A & 0 \end{bmatrix}^{-1} \right\| \leq K \quad \text{and} \quad \|f''(x)\| \leq M.$$

Show that  $\tilde{f}$  is  $\mu$ -strongly convex with convexity constant  $\mu = \sigma K^{-2} M^{-1}$ .

## Need help with your dissertation?

Get in-depth feedback & advice from experts in your topic area. Find out what you can do to improve the quality of your dissertation!

Get Help Now



Go to [www.helpmyassignment.co.uk](http://www.helpmyassignment.co.uk) for more info



# Chapter 16

## Self-concordant functions

Self-concordant functions were introduced by Nesterov and Nemirovski in the late 1980s as a product of their analysis of the speed of convergence of Newton’s method. Classic convergence results for two times continuously differentiable functions assume that the second derivative is Lipschitz continuous, and the convergence rate depends on the Lipschitz constant. One obvious weakness of these results is that the value of the Lipschitz constant, unlike Newton’s method, is not invariant under affine coordinate transformations.

Suppose that a function  $f$ , which is defined on an open convex subset  $X$  of  $\mathbf{R}^n$ , has a Lipschitz continuous second derivative with Lipschitz constant  $L$ , i.e. that

$$\|f''(y) - f''(x)\| \leq L\|y - x\|$$

for all  $x, y \in X$ . For the restriction  $\phi_{x,v}(t) = f(x + tv)$  of  $f$  to a line through  $x$  with direction vector  $v$ , this means that

$$|\phi''_{x,v}(t) - \phi''_{x,v}(0)| = |\langle v, (f''(x + tv) - f''(x))v \rangle| \leq L\|x + tv - x\|\|v\|^2 = L|t|\|v\|^3.$$

So if the function  $f$  is three times differentiable, then consequently

$$|\phi'''_{x,v}(0)| = \lim_{t \rightarrow 0} \left| \frac{\phi''_{x,v}(t) - \phi''_{x,v}(0)}{t} \right| \leq L\|v\|^3.$$

But

$$\phi'''_{x,v}(0) = \sum_{i,j,k=1}^n \frac{\partial^3 f(x)}{\partial x_i \partial x_j \partial x_k} v_i v_j v_k = D^3 f(x)[v, v, v],$$

so a necessary condition for a three times differentiable function  $f$  to have a Lipschitz continuous second derivative with Lipschitz constant  $L$  is that

$$(16.1) \quad |D^3 f(x)[v, v, v]| \leq L\|v\|^3$$

for all  $x \in X$  and all  $v \in \mathbf{R}^n$ , and it is easy to show this is also a sufficient condition.

The reason why the value of the Lipschitz constant is not affinely invariant is that there is no natural connection between the Euclidean norm  $\|\cdot\|$  and the function  $f$ . The analysis of a function's behavior is simplified if we instead use a norm that is adapted to the form of the level surfaces, and for functions with a positive semidefinite second derivative  $f''(x)$ , such a (semi)norm is the local seminorm  $\|\cdot\|_x$ , introduced in the previous chapter and defined as  $\|v\|_x = \sqrt{\langle v, f''(x)v \rangle}$ . Nesterov–Nemirovski's stroke of genius consisted in replacing  $\|\cdot\|$  with the local seminorm  $\|\cdot\|_x$  in the inequality (16.1). For the function class obtained in this way, it is possible to describe the convergence rate of Newton's method in an affinely independent way and with absolute constants.

## 16.1 Self-concordant functions

We are now ready for Nesterov–Nemirovski's definition of self-concordance and for a study of the basic properties of self-concordant functions.

**Definition.** Let  $f: X \rightarrow \mathbf{R}$  be a three times continuously differentiable function with an open convex subset  $X$  of  $\mathbf{R}^n$  as domain. The function is called *self-concordant* if it is convex, and the inequality

$$(16.2) \quad |D^3 f(x)[v, v, v]| \leq 2(D^2 f(x)[v, v])^{3/2}$$

holds for all  $x \in X$  and all  $v \in \mathbf{R}^n$ .

Since  $D^2 f(x)[v, v] = \|v\|_x^2$ , where  $\|\cdot\|_x$  is the local seminorm defined by the function  $f$  at the point  $x$ , we can also write the defining inequality (16.2) as

$$|D^3 f(x)[v, v, v]| \leq 2\|v\|_x^3,$$

and it is this shorter version that we will prefer, when we work with a single function  $f$ .

*Remark 1.* There is nothing special about the constant 2 in inequality (16.2). If  $f$  satisfies the inequality  $|D^3 f(x)[v, v, v]| \leq K\|v\|_x^3$ , then the function  $F = \frac{1}{4}K^2 f$ , obtained from  $f$  by scaling, is self-concordant. The choice of 2 as the constant facilitates, however, the wording of a number of results.

*Remark 2.* For functions  $f$  defined on subsets of the real axis and  $v \in \mathbf{R}$ ,  $\|v\|_x^2 = f''(x)v^2$  and  $D^3 f(x)[v, v, v] = f'''(x)v^3$ . Hence, a convex function  $f: X \rightarrow \mathbf{R}$  is self-concordant if and only if

$$|f'''(x)| \leq 2f''(x)^{3/2}$$

for all  $x \in X$ .

*Remark 3.* In terms of the restriction  $\phi_{x,v}(t) = f(x + tv)$  of the function  $f$  to the line through  $x$  with direction  $v$ , we can equivalently write the inequality

$$|D^3 f(x + tv)[v, v, v]| \leq 2\|v\|_{x+tv}^3$$

as  $|\phi_{x,v}'''(t)| \leq 2\phi_{x,v}''(t)^{3/2}$ . A three times continuously differentiable convex function of several variables is therefore self-concordant if and only if all its restrictions to lines are self-concordant.

**EXAMPLE 16.1.1.** The convex function  $f(x) = -\ln x$  is self-concordant on its domain  $\mathbf{R}_{++}$ . Indeed, inequality (16.2) holds with equality for this function, since  $f''(x) = x^{-2}$  and  $f'''(x) = -2x^{-3}$ .  $\square$

**EXAMPLE 16.1.2.** Convex quadratic functions  $f(x) = \frac{1}{2}\langle x, Ax \rangle + \langle b, x \rangle + c$  are self-concordant since  $D^3 f(x)[v, v, v] = 0$  for all  $x$  and  $v$ .

Hence, affine functions are self-concordant, and the function  $x \mapsto \|x\|^2$ , where  $\|\cdot\|$  is the Euclidean norm, is self-concordant.  $\square$



**Brain power**

By 2020, wind could provide one-tenth of our planet's electricity needs. Already today, SKF's innovative know-how is crucial to running a large proportion of the world's wind turbines.

Up to 25 % of the generating costs relate to maintenance. These can be reduced dramatically thanks to our systems for on-line condition monitoring and automatic lubrication. We help make it more economical to create cleaner, cheaper energy out of thin air.

By sharing our experience, expertise, and creativity, industries can boost performance beyond expectations. Therefore we need the best employees who can meet this challenge!

The Power of Knowledge Engineering

Plug into The Power of Knowledge Engineering.  
Visit us at [www.skf.com/knowledge](http://www.skf.com/knowledge)

**SKF**

The expression

$$D^3 f(x)[u, v, w] = \sum_{i,k,k=1}^n \frac{\partial^3 f(x)}{\partial x_i \partial x_j \partial x_k} u_i v_j w_k$$

is a symmetric trilinear form in the variables  $u, v$ , and  $w$ , if the function  $f$  is three times continuously differentiable in a neighborhood of the point  $x$ . For self-concordant functions we have the following generalization of inequality (16.2) in the definition of self-concordance.

**Theorem 16.1.1.** *Suppose  $f: X \rightarrow \mathbf{R}$  is a self-concordant function. Then,*

$$|D^3 f(x)[u, v, w]| \leq 2\|u\|_x \|v\|_x \|w\|_x$$

for all  $x \in X$  and all vectors  $u, v, w$  in  $\mathbf{R}^n$ .

*Proof.* The proof is based on a general theorem on norms of symmetric trilinear forms, which is proven in an appendix to this chapter.

Assume first that  $x$  is a point where the second derivative  $f''(x)$  is positive definite. Then  $\|\cdot\|_x$  is a norm with  $\langle u, v \rangle_x = \langle u, f''(x)v \rangle$  as the corresponding scalar product. We can therefore apply Theorem 1 of the appendix to the symmetric trilinear form  $D^3 f(x)[u, v, w]$  with  $\|\cdot\|_x$  as the underlying norm, and it follows that

$$\sup_{u,v,w \neq 0} \frac{|D^3 f(x)[u, v, w]|}{\|u\|_x \|v\|_x \|w\|_x} = \sup_{v \neq 0} \frac{|D^3 f(x)[v, v, v]|}{\|v\|_x^3} \leq 2,$$

which is equivalent to the assertion of the theorem.

To cope with points where the second derivative is singular, we consider for  $\epsilon > 0$  the scalar product

$$\langle u, v \rangle_{x,\epsilon} = \langle u, f''(x)v \rangle + \epsilon \langle u, v \rangle,$$

where  $\langle \cdot, \cdot \rangle$  is the usual standard scalar product, and the corresponding norm

$$\|v\|_{x,\epsilon} = \sqrt{\langle v, v \rangle_{x,\epsilon}} = \sqrt{\|v\|_x^2 + \epsilon \|v\|^2}.$$

Obviously,  $\|v\|_x \leq \|v\|_{x,\epsilon}$  for all vectors  $v$ , and hence

$$|D^3 f(x)[v, v, v]| \leq 2\|v\|_{x,\epsilon}^3$$

for all  $v$ , since  $f$  is self-concordant. It now follows from Theorem 1 in the appendix that

$$\begin{aligned} |D^3 f(x)[u, v, w]| &\leq 2\|u\|_{x,\epsilon} \|v\|_{x,\epsilon} \|w\|_{x,\epsilon} \\ &= 2\sqrt{(\|u\|_x^2 + \epsilon \|u\|^2)(\|v\|_x^2 + \epsilon \|v\|^2)(\|w\|_x^2 + \epsilon \|w\|^2)}, \end{aligned}$$

and we get the sought-after inequality by letting  $\epsilon \rightarrow 0$ .  $\square$

**Theorem 16.1.2.** *The second derivative  $f''(x)$  of a self-concordant function  $f: X \rightarrow \mathbf{R}$  has the same null space  $\mathcal{N}(f''(x))$  at all points  $x \in X$ .*

*Proof.* We recall that  $\mathcal{N}(f''(x)) = \{v \mid \|v\|_x = 0\}$ .

Let  $x$  and  $y$  be two points in  $X$ . For reasons of symmetry, we only have to show the inclusion  $\mathcal{N}(f''(x)) \subseteq \mathcal{N}(f''(y))$ .

Assume therefore that  $v \in \mathcal{N}(f''(x))$  and let  $x^t = x + t(y - x)$ . Since  $X$  is an open convex set, there is certainly a number  $a > 1$  such that the points  $x^t$  lie in  $X$  for  $0 \leq t \leq a$ , and we now define a function  $g: [0, a] \rightarrow \mathbf{R}$  by setting

$$g(t) = D^2 f(x^t)[v, v] = \|v\|_{x^t}^2.$$

Then  $g(0) = \|v\|_x^2 = 0$  and  $g(t) \geq 0$  for  $0 \leq t \leq a$ , and since

$$g'(t) = D^3 f(x^t)[v, v, y - x],$$

it follows from Theorem 16.1.1 that

$$|g'(t)| \leq 2\|v\|_{x^t}^2 \|y - x\|_{x^t} = 2g(t)\|y - x\|_{x^t}.$$

But the seminorm

$$\|y - x\|_{x^t} = \sqrt{D^2 f(x^t)[y - x, y - x]}$$

depends continuously on  $t$ , and it is therefore bounded above by some constant  $C$  on the interval  $[0, a]$ . Hence,

$$|g'(t)| \leq 2Cg(t)$$

for  $0 \leq t \leq a$ . It now follows from Theorem 2 in the appendix to this chapter that  $g(t) = 0$  for all  $t$ , and in particular,  $g(1) = \|v\|_y^2 = 0$ , which proves that  $v \in \mathcal{N}(f''(y))$ . This proves the inclusion  $\mathcal{N}(f''(x)) \subseteq \mathcal{N}(f''(y))$ .  $\square$

Our next corollary is just a special case of Theorem 16.1.2, because  $f''(x)$  is non-singular if and only if  $\mathcal{N}(f''(x)) = \{0\}$ .

**Corollary 16.1.3.** *The second derivative of a self-concordant function is either non-singular at all points or singular at all points.*

A self-concordant function will be called *non-degenerate* if its second derivative is positive definite at all points, and by the above corollary, that is the case if the second derivative is positive definite at one single point.

A non-degenerate self-concordant function is in particular strictly convex.

## Operations that preserve self-concordance

**Theorem 16.1.4.** *If  $f$  is a self-concordant function and  $\alpha \geq 1$ , then  $\alpha f$  is self-concordant.*

*Proof.* If  $\alpha \geq 1$ , then  $\alpha \leq \alpha^{3/2}$ , and it follows that

$$\begin{aligned} |D^3(\alpha f)(x)[v, v, v]| &= \alpha |D^3 f(x)[v, v, v]| \leq 2\alpha (D^2 f(x)[v, v])^{3/2} \\ &\leq 2(\alpha D^2 f(x)[v, v])^{3/2} = 2(D^2(\alpha f)(x)[v, v])^{3/2}. \quad \square \end{aligned}$$

**Theorem 16.1.5.** *The sum  $f + g$  of two self-concordant functions  $f$  and  $g$  is self-concordant on its domain.*

*Proof.* We use the elementary inequality

$$a^{3/2} + b^{3/2} \leq (a + b)^{3/2},$$

which holds for all nonnegative numbers  $a, b$  (and is easily proven by squaring both sides) and the triangle inequality to obtain

What do you want to do?

No matter what you want out of your future career, an employer with a broad range of operations in a load of countries will always be the ticket. Working within the Volvo Group means more than 100,000 friends and colleagues in more than 185 countries all over the world. We offer graduates great career opportunities – check out the Career section at our web site [www.volvogroup.com](http://www.volvogroup.com). We look forward to getting to know you!

**VOLVO**  
AB Volvo (publ)  
[www.volvogroup.com](http://www.volvogroup.com)

VOLVO TRUCKS | RENAULT TRUCKS | MACK TRUCKS | VOLVO BUSES | VOLVO CONSTRUCTION EQUIPMENT | VOLVO PENTA | VOLVO AERO | VOLVO IT  
VOLVO FINANCIAL SERVICES | VOLVO 3P | VOLVO POWERTRAIN | VOLVO PARTS | VOLVO TECHNOLOGY | VOLVO LOGISTICS | BUSINESS AREA ASIA



$$\begin{aligned}
 |D^3(f+g)(x)[v, v, v]| &= |D^3f(x)[v, v, v] + D^3g(x)[v, v, v]| \\
 &\leq 2(D^2f(x)[v, v])^{3/2} + 2(D^2g(x)[v, v])^{3/2} \\
 &\leq 2(D^2f(x)[v, v] + D^2g(x)[v, v])^{3/2} \\
 &= 2(D^2(f+g)(x)[v, v])^{3/2}. \quad \square
 \end{aligned}$$

**Theorem 16.1.6.** *If the function  $f: X \rightarrow \mathbf{R}$  is self-concordant, where  $X$  is an open convex subset of  $\mathbf{R}^n$ , and  $A$  is an affine map from  $\mathbf{R}^m$  to  $\mathbf{R}^n$ , then the composition  $g = f \circ A$  is a self-concordant function on its domain  $A^{-1}(X)$ .*

*Proof.* The affine map  $A$  can be written as  $Ay = Cy + b$ , where  $C$  is a linear map and  $b$  is a vector. Let  $y$  be a point in  $A^{-1}(X)$  and let  $u$  be a vector in  $\mathbf{R}^m$ , and write  $x = Ay$  and  $v = Cu$ . According to the chain rule,

$$\begin{aligned}
 D^2g(y)[u, u] &= D^2f(Ay)[Cu, Cu] = D^2f(x)[v, v] \quad \text{and} \\
 D^3g(y)[u, u, u] &= D^3f(Ay)[Cu, Cu, Cu] = D^3f(x)[v, v, v],
 \end{aligned}$$

so it follows that

$$\begin{aligned}
 |D^3g(y)[u, u, u]| &= |D^3f(x)[v, v, v]| \leq 2(D^2f(x)[v, v])^{3/2} \\
 &= 2(D^2g(y)[u, u])^{3/2}. \quad \square
 \end{aligned}$$

**EXAMPLE 16.1.3.** It follows from Example 16.1.1 and Theorem 16.1.6 that the function  $f(x) = -\ln(b - \langle c, x \rangle)$  with domain  $\{x \in \mathbf{R}^n \mid \langle c, x \rangle < b\}$  is self-concordant.  $\square$

**EXAMPLE 16.1.4.** Suppose that the polyhedron

$$X = \bigcap_{j=1}^p \{x \in \mathbf{R}^n \mid \langle c_j, x \rangle \leq b_j\}$$

has nonempty interior. The function  $f(x) = -\sum_{j=1}^p \ln(b_j - \langle c_j, x \rangle)$ , with  $\text{int } X$  as domain, is self-concordant.  $\square$

## 16.2 Closed self-concordant functions

In Section 6.7 of Part I we studied the recessive subspace of arbitrary convex functions. The properties of the recessive subspace of a closed self-concordant function is given by the following theorem.

**Theorem 16.2.1.** *Suppose that  $f: X \rightarrow \mathbf{R}$  is a closed self-concordant function. The function's recessive subspace  $V_f$  is then equal to the null space  $\mathcal{N}(f''(x))$  of the second derivative  $f''(x)$  at an arbitrary point  $x \in X$ . Moreover,*

- (i)  $X = X + V_f$ .
- (ii)  $f(x + v) = f(x) + Df(x)[v]$  for all vectors  $v \in V_f$ .
- (iii) If  $\lambda(f, x) < \infty$ , then  $f(x + v) = f(x)$  for all  $v \in V_f$ .

*Proof.* Assertions (i) and (ii) are true for the recessive subspace of an arbitrary differentiable convex function according to Theorem 6.7.1, so we only have to prove the remaining assertions.

Let  $x$  be an arbitrary point in  $X$  and let  $v$  be an arbitrary vector in  $\mathbf{R}^n$ , and consider the restriction  $\phi_{x,v}(t) = f(x + tv)$  of  $f$  to the line through  $x$  with direction  $v$ . The domain of  $\phi_{x,v}$  is an open interval  $I = ]\alpha, \beta[$  around 0.

First suppose that  $v \in V_f$ . Then

$$\phi_{x,v}(t) = f(x) + tDf(x)[v]$$

for all  $t \in I$  because of property (ii), and it follows that

$$\|v\|_x^2 = D^2f(x)[v, v] = \phi_{x,v}''(0) = 0,$$

i.e. the vector  $v$  belongs to the null space of  $f''(x)$ . This proves the inclusion  $V_f \subseteq \mathcal{N}(f''(x))$ . Note that this inclusion holds for arbitrary twice differentiable convex functions without any assumptions concerning self-concordance and closedness.

To prove the converse inclusion  $\mathcal{N}(f''(x)) \subseteq V_f$ , we instead assume that  $v$  is a vector in  $\mathcal{N}(f''(x))$ . Since  $\mathcal{N}(f''(x + tv)) = \mathcal{N}(f''(x))$  for all  $t \in I$  due to Theorem 16.1.2, we now have

$$\phi_{x,v}''(t) = D^2f(x + tv)[v, v] = \|v\|_{x+tv}^2 = 0$$

for all  $t \in I$ , and it follows that

$$\phi_{x,v}(t) = \phi_{x,v}(0) + \phi_{x,v}'(0)t = f(x) + Df(x)[v]t.$$

If  $\beta < \infty$ , then  $x + \beta v$  is a boundary point of  $X$  and  $\lim_{t \rightarrow \beta} \phi_{x,v}(t) < \infty$ . However, according to Corollary 8.2.2 in Part I this is a contradiction to  $f$  being a closed function. Hence,  $\beta = \infty$ , and similarly,  $\alpha = -\infty$ . This means that  $I = ]-\infty, \infty[$ , and in particular,  $I$  contains the number 1. We conclude that the point  $x + v$  lies in  $X$  and that  $f(x + v) = \phi_{x,v}(1) = f(x) + Df(x)[v]$  for all  $x \in X$  and all  $v \in \mathcal{N}(f''(x))$ , and Theorem 6.7.1 now provides us with the inclusion  $\mathcal{N}(f''(x)) \subseteq V_f$ . Hence,  $V_f = \mathcal{N}(f''(x))$ .

Finally, suppose that  $\lambda(f, x) < \infty$ . Then there exists, by definition, a Newton direction at  $x$ , and this implies, according to the remark after the definition of Newton direction, that the implication

$$f''(x)v = 0 \Rightarrow Df(x)[v] = 0$$

holds. Since  $V_f = \mathcal{N}(f''(x))$ , it now follows from assertion (ii) that  $f(x+v) = f(x)$  for all  $v \in V_f$ .  $\square$

The problem of minimizing a degenerate closed self-concordant function  $f: X \rightarrow \mathbf{R}$  with finite Newton decrement  $\lambda(f, x)$  at all points  $x \in X$  can be reduced to the problem of minimizing a non-degenerate closed self-concordant function as follows.

Assume that the domain  $X$  is a subset of  $\mathbf{R}^n$ , and let  $V_f$  denote the recessive subspace of  $f$ . Put  $m = \dim V_f^\perp$  and let  $A: \mathbf{R}^m \rightarrow \mathbf{R}^n$  be an arbitrary injective linear map onto  $V_f^\perp$ , and put  $X_0 = A^{-1}(X)$ . The set  $X_0$  is then an open subset of  $\mathbf{R}^m$ , and we obtain a function  $g: X_0 \rightarrow \mathbf{R}$  by defining  $g(y) = f(Ay)$  for  $y \in X_0$ .

The function  $g$  is self-concordant according to Theorem 16.1.6, and since  $(y, t)$  belongs to the epigraph of  $g$  if and only if  $(Ay, t)$  belongs to the epigraph of  $f$ , it follows that  $g$  is also a closed function.

**gaiteye**<sup>®</sup>  
Challenge the way we run

**EXPERIENCE THE POWER OF  
FULL ENGAGEMENT...**

.....

**RUN FASTER.  
RUN LONGER..  
RUN EASIER...**

**READ MORE & PRE-ORDER TODAY  
WWW.GAITEYE.COM**

Suppose  $v \in \mathcal{N}(g''(y))$ . Since  $g''(y) = A^T f''(Ay)A$ ,

$$\langle Av, f''(Ay)Av \rangle = \langle v, A^T f''(Ay)Av \rangle = \langle v, g''(y)v \rangle = 0,$$

which means that the vector  $Av$  belongs to  $\mathcal{N}(f''(Ay))$ , i.e. to the recessive subspace  $V_f$ . But  $Av$  also belongs to  $V_f^\perp$ , by definition, and  $V_f \cap V_f^\perp = \{0\}$ , so it follows that  $Av = 0$ . Hence  $v = 0$ , since  $A$  is an injective map. This proves that  $\mathcal{N}(g''(y)) = \{0\}$ , which means that  $g$  is a non-degenerate function.

Each vector  $x \in X$  has a unique decomposition  $x = x_1 + x_2$  with  $x_1 \in V_f^\perp$  and  $x_2 \in V_f$ , and  $x_1 (= x - x_2)$  lies in  $X$  according to Theorem 16.2.1. Consequently, there is a unique point  $y \in X_0$  such that  $Ay = x_1$ . Therefore,  $g(y) = f(Ay) = f(x_1) = f(x)$  by the same theorem.

The functions  $f$  and  $g$  thus have the same ranges, and  $\hat{y}$  is a minimum point of  $g$  if and only if  $A\hat{y}$  is a minimum point of  $f$ , and thereby also all points  $A\hat{y} + v$  with  $v \in V_f$  are minimum points of  $f$ .

We also note for future use that

$$\lambda(g, y) \leq \lambda(f, Ay) = \lambda(f, Ay + v)$$

for all  $y \in X_0$  and all  $v \in V_f$ , according to Theorem 15.1.7. (In the present case, the two Newton decrements are actually equal, which we leave as an exercise to show.)

**Corollary 16.2.2.** *A closed self-concordant function  $f: X \rightarrow \mathbf{R}$  is non-degenerate if its domain  $X$  does not contain any line.*

*Proof.* By Theorem 16.2.1,  $X = X + V_f$ . Hence, if  $f$  is degenerate, then  $X$  contains all lines through points in  $X$  with directions given by nonzero vectors in  $V_f$ . So the function must be non-degenerate if its domain does not contain any lines.  $\square$

**Corollary 16.2.3.** *A closed self-concordant function is non-degenerate if and only if it is strictly convex.*

*Proof.* The second derivative  $f''(x)$  of a non-degenerate self-concordant function  $f$  is positive definite for all  $x$  in its domain, and this implies that  $f$  is strictly convex.

The recessive subspace  $V_f$  of a degenerate function  $f$  is non-trivial, and the restriction  $\phi_{x,v}(t) = f(x + tv)$  of  $f$  to a line with a direction given by a nonzero vector  $v \in V_f$  is affine, according to Theorem 16.2.1. This prevents  $f$  from being strictly convex.  $\square$

## 16.3 Basic inequalities for the local seminorm

The graph of a convex function  $f$  lies above its tangent planes, and the vertical distance between the point  $(y, f(y))$  on the graph and the tangent plane through the point  $(x, f(x))$  is greater than or equal to  $\frac{1}{2}\mu\|y-x\|^2$  if  $f$  is  $\mu$ -strongly convex. The same distance is also bounded below if the function is self-concordant, but now by an expression that is a function of the local norm  $\|y-x\|_x$ . The actual function  $\rho$  is defined in the following lemma, which also describes all the properties of  $\rho$  that we will need.

**Lemma 16.3.1.** *Let  $\rho: ]-\infty, 1[ \rightarrow \mathbf{R}$  be the function*

$$\rho(t) = -t - \ln(1-t).$$

(i) *The function  $\rho$  is convex, strictly decreasing in the interval  $]-\infty, 0]$ , and strictly increasing in the interval  $[0, 1[$ , and  $\rho(0) = 0$ .*

(ii) *For  $0 \leq t < 1$ ,*

$$\rho(t) \leq \frac{t^2}{2(1-t)}.$$

*In particular,  $\rho(t) \leq t^2$  if  $0 \leq t \leq \frac{1}{2}$ .*

(iii) *If  $s < 1$  and  $t < 1$ , then  $\rho(s) + \rho(t) \geq -st$ .*

(iv) *If  $s \geq 0$ ,  $0 \leq t < 1$  and  $\rho(-s) \leq \rho(t)$ , then  $s \leq \frac{t}{1-t}$ .*

*Proof.* Assertion (i) follows easily by considering the sign of the derivative, and assertion (ii) follows from the Taylor series expansion, which gives

$$\rho(t) = \frac{1}{2}t^2 + \frac{1}{3}t^3 + \frac{1}{4}t^4 + \dots \leq \frac{1}{2}t^2(1+t+t^2+\dots) = \frac{1}{2}t^2(1-t)^{-1}$$

for  $0 \leq t < 1$ .

To prove (iii), we use the elementary inequality  $x - \ln(1+x) \geq 0$  and take  $x = st - s - t$ . This gives

$$\begin{aligned} st + \rho(s) + \rho(t) &= st - s - t - \ln(1-s) - \ln(1-t) \\ &= st - s - t - \ln(1+st-s-t) \geq 0. \end{aligned}$$

Since  $\rho$  is strictly decreasing in the interval  $]-\infty, 0]$ , assertion (iv) will follow once we show that  $\rho(-s) \geq \rho(t)$  when  $s = t/(1-t)$ . To show this inequality, let

$$g(t) = \rho\left(-\frac{t}{1-t}\right) - \rho(t)$$

for  $0 \leq t < 1$ . We simplify and obtain

$$g(t) = t - 1 + (1-t)^{-1} + 2\ln(1-t).$$

Since  $g(0) = 0$  and  $g'(t) = 1 + (1 - t)^{-2} - 2(1 - t)^{-1} = t^2(1 - t)^{-2} \geq 0$ , we conclude that  $g(t) \geq 0$  for all  $t \in [0, 1[$ , and this completes the proof of assertion (iv).  $\square$

The next theorem is used to estimate differences of the form  $\|w\|_y - \|w\|_x$ ,  $Df(y)[w] - Df(x)[w]$ , and  $f(y) - f(x) - Df(x)[y - x]$  in terms of  $\|w\|_x$ ,  $\|y - x\|_x$  and the function  $\rho$ .

**Theorem 16.3.2.** *Let  $f: X \rightarrow \mathbf{R}$  be a closed self-concordant function, and suppose that  $x$  is a point in  $X$  and that  $\|y - x\|_x < 1$ . Then,  $y$  is also a point in  $X$ , and the following inequalities hold for the vector  $v = y - x$  and arbitrary vectors  $w$ :*

$$(16.3) \quad \frac{\|v\|_x}{1 + \|v\|_x} \leq \|v\|_y \leq \frac{\|v\|_x}{1 - \|v\|_x}$$

$$(16.4) \quad \frac{\|v\|_x^2}{1 + \|v\|_x} \leq Df(y)[v] - Df(x)[v] \leq \frac{\|v\|_x^2}{1 - \|v\|_x}$$

$$(16.5) \quad \rho(-\|v\|_x) \leq f(y) - f(x) - Df(x)[v] \leq \rho(\|v\|_x)$$

$$(16.6) \quad (1 - \|v\|_x)\|w\|_x \leq \|w\|_y \leq \frac{\|w\|_x}{1 - \|v\|_x}$$

$$(16.7) \quad Df(y)[w] - Df(x)[w] \leq D^2f(x)[v, w] + \frac{\|v\|_x^2\|w\|_x}{1 - \|v\|_x} \leq \frac{\|v\|_x\|w\|_x}{1 - \|v\|_x}.$$

The left parts of the three inequalities (16.3), (16.4) and (16.5) are also satisfied with  $v = y - x$  for all  $y \in X$ .

*Proof.* We leave the proof that  $y$  belongs to  $X$  to the end and start by showing that the inequalities (16.3–16.7) hold under the additional assumption  $y \in X$ .

**I.** We begin with inequality (16.6). If  $\|w\|_x = 0$ , then  $\|w\|_z = 0$  for all  $z \in X$ , according to Theorem 16.1.2. Hence, the inequality holds in this case. Therefore, let  $w$  be an arbitrary vector with  $\|w\|_x \neq 0$ , let  $x^t = x + t(y - x)$ , and define the function  $\psi$  by

$$\psi(t) = \|w\|_{x^t}^{-1} = (D^2f(x^t)[w, w])^{-1/2}.$$

The function  $\psi$  is defined on an open interval that contains the interval  $[0, 1]$ ,  $\psi(0) = \|w\|_x^{-1}$  and  $\psi(1) = \|w\|_y^{-1}$ . It now follows, using Theorem 16.1.1, that

$$(16.8) \quad \begin{aligned} |\psi'(t)| &= \frac{1}{2} |(D^2f(x^t)[w, w])^{-3/2} D^3f(x^t)[w, w, v]| \\ &= \frac{1}{2} \|w\|_{x^t}^{-3} |D^3f(x^t)[w, w, v]| \leq \frac{1}{2} \|w\|_{x^t}^{-3} \cdot 2\|w\|_{x^t}^2\|v\|_{x^t} \\ &= \|w\|_{x^t}^{-1}\|v\|_{x^t} = \psi(t)\|v\|_{x^t}. \end{aligned}$$

If  $\|v\|_x = 0$ , then  $\|v\|_z = 0$  for all  $z \in X$ , and hence  $\psi'(t) = 0$  for  $0 \leq t \leq 1$ . This implies that  $\psi(1) = \psi(0)$ , i.e. that  $\|w\|_y = \|w\|_x$ . The inequalities (16.3) and (16.6) are thus satisfied in the case  $\|v\|_x = 0$ .

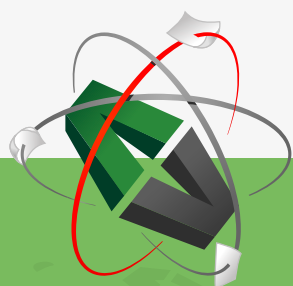
Assume henceforth that  $\|v\|_x \neq 0$ , and first take  $w = v$  in the definition of the function  $\psi$ . In this special case, inequality (16.8) simplifies to  $|\psi'(t)| \leq 1$  for  $t \in [0, 1]$ , and hence  $\psi(0) - 1 \leq \psi(1) \leq \psi(0) + 1$ , by the mean-value theorem. The right part of this inequality means that  $\|v\|_y^{-1} \leq \|v\|_x^{-1} + 1$ , which after rearrangement gives the left part of inequality (16.3). Note, that this is true even in the case  $\|v\|_x \geq 1$ .

Correspondingly, the left part of the same inequality gives rise to the right part of inequality (16.3), now under the assumption that  $\|v\|_x < 1$ .

To prove inequality (16.6), we return to the function  $\psi$  with a general  $w$ . Since  $\|tv\|_x = t\|v\|_x < 1$  for  $0 \leq t \leq 1$ , it follows from the already proven inequality (16.3) (with  $x^t = x + tv$  instead of  $y$ ) that

$$\|v\|_{x^t} = \frac{1}{t} \|tv\|_{x^t} \leq \frac{1}{t} \cdot \frac{\|tv\|_x}{1 - \|tv\|_x} = \frac{\|v\|_x}{1 - t\|v\|_x}.$$

This e-book  
is made with  
**SetaPDF**



PDF components for PHP developers

[www.setasign.com](http://www.setasign.com)

Insert this estimate into (16.8); this gives us the following inequality for the derivative of the function  $\ln \psi(t)$ :

$$|(\ln \psi(t))'| = \frac{|\psi'(t)|}{\psi(t)} = \|v\|_{x^t} \leq \frac{\|v\|_x}{1 - t\|v\|_x}.$$

Let us now integrate this inequality over the interval  $[0, 1]$ ; this results in the estimate

$$\begin{aligned} \left| \ln \frac{\|w\|_y}{\|w\|_x} \right| &= \left| \ln \frac{\psi(0)}{\psi(1)} \right| = |\ln \psi(1) - \ln \psi(0)| = \left| \int_0^1 (\ln \psi(t))' dt \right| \\ &\leq \int_0^1 \frac{\|v\|_x}{1 - t\|v\|_x} dt = -\ln(1 - \|v\|_x), \end{aligned}$$

which after exponentiation yields

$$1 - \|v\|_x \leq \frac{\|w\|_y}{\|w\|_x} \leq (1 - \|v\|_x)^{-1},$$

and this is inequality (16.6).

**II.** To prove the inequality (16.4), we define

$$\phi(t) = Df(x^t)[v],$$

where  $x^t = x + t(y - x)$ , as before. Then

$$\phi'(t) = D^2f(x^t)[v, v] = \|v\|_{x^t}^2 = t^{-2}\|tv\|_{x^t}^2,$$

so by using inequality (16.3), we obtain the inequality

$$\frac{\|v\|_x^2}{(1 + t\|v\|_x)^2} = \frac{1}{t^2} \frac{\|tv\|_x^2}{(1 + \|tv\|_x)^2} \leq \phi'(t) \leq \frac{1}{t^2} \frac{\|tv\|_x^2}{(1 - \|tv\|_x)^2} = \frac{\|v\|_x^2}{(1 - t\|v\|_x)^2}$$

for  $0 \leq t \leq 1$ . The left part of this inequality holds with  $v = y - x$  for all  $y \in X$ , and the right part holds if  $\|v\|_x < 1$ , and by integrating the inequality over the interval  $[0, 1]$ , we arrive at inequality (16.4).

**III.** To prove inequality (16.5), we start with the function

$$\Phi(t) = f(x^t) - Df(x)[v]t,$$

noting that

$$\Phi(1) - \Phi(0) = f(y) - f(x) - Df(x)[v]$$

and that

$$\Phi'(t) = Df(x^t)[v] - Df(x)[v].$$



By replacing  $y$  with  $x^t$  in inequality (16.4), we obtain the following inequality

$$\frac{t\|v\|_x^2}{1+t\|v\|_x} \leq \Phi'(t) \leq \frac{t\|v\|_x^2}{1-t\|v\|_x},$$

where the right part holds only if  $\|v\|_x < 1$ . By integrating the above inequality over the interval  $[0, 1]$ , we obtain

$$\rho(-\|v\|_x) = \int_0^1 \frac{t\|v\|_x^2}{1+t\|v\|_x} dt \leq \Phi(1) - \Phi(0) \leq \int_0^1 \frac{t\|v\|_x^2}{1-t\|v\|_x} dt = \rho(\|v\|_x),$$

i.e. inequality (16.5).

**IV.** The proof of inequality (16.7) is analogous to the proof of inequality (16.4), but this time our function  $\phi$  is defined as

$$\phi(t) = Df(x^t)[w].$$

Now,  $\phi'(t) = D^2f(x^t)[w, v]$  and  $\phi''(t) = D^3f(x^t)[w, v, v]$ , so it follows from Theorem 16.1.1 and inequality (16.6) that

$$|\phi''(t)| \leq 2\|w\|_{x^t}\|v\|_{x^t}^2 \leq 2\frac{\|w\|_x\|v\|_x^2}{(1-t\|v\|_x)^3}.$$

By integrating this inequality over the interval  $[0, s]$ , where  $s \leq 1$ , we get the estimate

$$\begin{aligned} \phi'(s) - \phi'(0) &\leq \int_0^s |\phi''(t)| dt \leq 2\|w\|_x \int_0^s \frac{\|v\|_x^2 dt}{(1-t\|v\|_x)^3} \\ &= \|w\|_x \left[ \frac{\|v\|_x}{(1-s\|v\|_x)^2} - \|v\|_x \right], \end{aligned}$$

and another integration over the interval  $[0, 1]$  results in the inequality

$$\phi(1) - \phi(0) - \phi'(0) \leq \int_0^1 (\phi'(s) - \phi'(0)) ds \leq \frac{\|w\|_x\|v\|_x^2}{1-\|v\|_x},$$

which is the left part of inequality (16.7).

By the Cauchy-Schwarz inequality,

$$\begin{aligned} D^2f(x)[v, w] &= \langle v, f''(x)w \rangle = \langle f''(x)^{1/2}v, f''(x)^{1/2}w \rangle \\ &\leq \|f''(x)^{1/2}v\| \|f''(x)^{1/2}w\| = \|v\|_x \|w\|_x, \end{aligned}$$

and we obtain the right part of inequality (16.7) by replacing  $D^2f(x)[v, w]$  with its majorant  $\|v\|_x\|w\|_x$ .

V. It now only remains to prove that the condition  $\|y - x\|_x < 1$  implies that the point  $y$  lies in  $X$ .

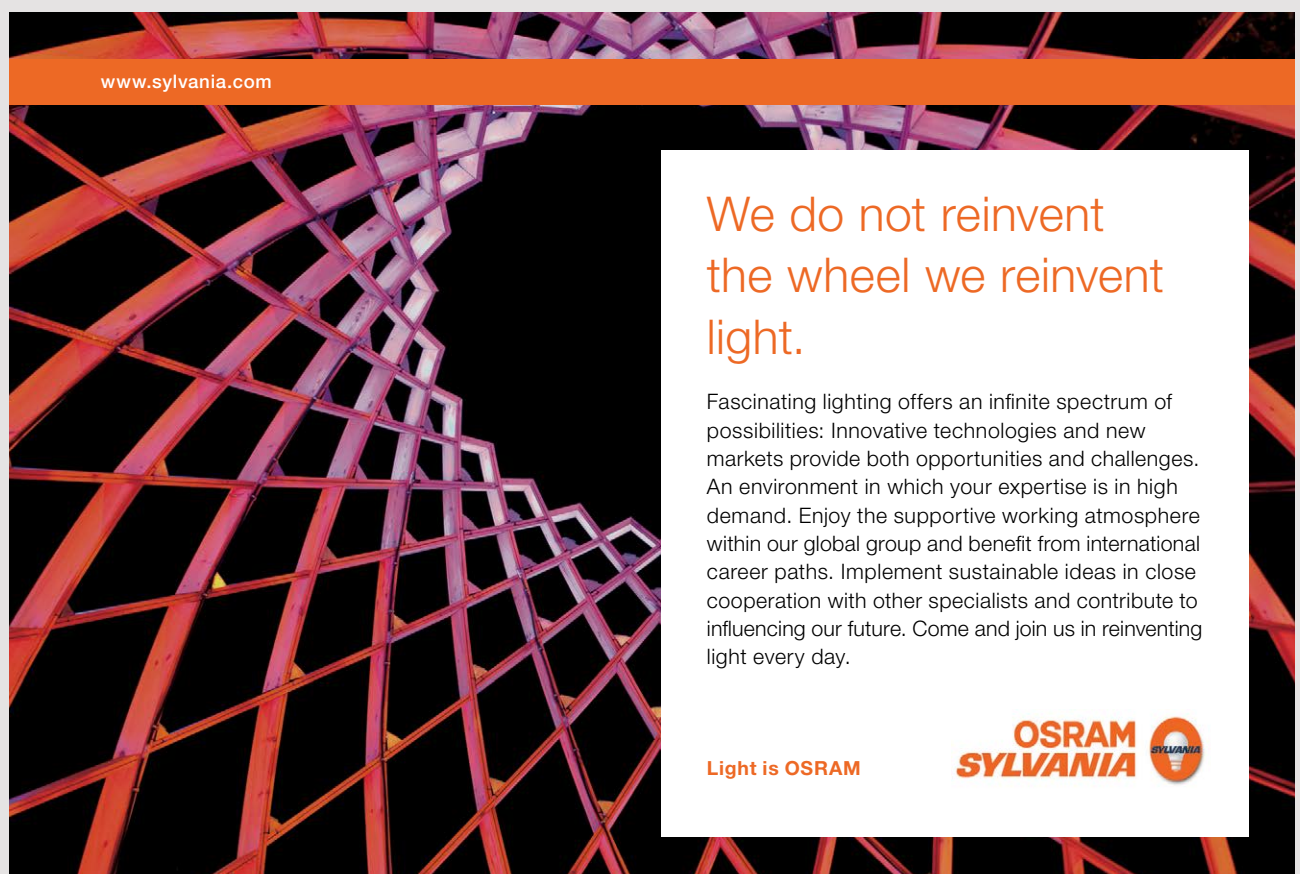
Assume the contrary. i.e. that there is a point  $y$  outside  $X$  such that  $\|y - x\|_x < 1$ . The line segment  $[x, y]$  then intersects the boundary of  $X$  in a point  $x + \bar{t}v$ , where  $\bar{t}$  is a number in the interval  $]0, 1[$ . The function  $\rho$  is increasing in the interval  $[0, 1[$ , and hence  $\rho(t\|v\|_x) \leq \rho(\|v\|_x)$  if  $0 \leq t < \bar{t}$ . It therefore follows from inequality (16.5) that

$$f(x + tv) \leq f(x) + tDf(x)[v] + \rho(t\|v\|_x) \leq f(x) + |Df(x)[v]| + \rho(\|v\|_x) < +\infty$$

for all  $t$  in the interval  $[0, \bar{t}[$ . However, this is a contradiction, because  $\lim_{t \rightarrow \bar{t}} f(x + tv) = +\infty$ , since  $f$  is a closed function and  $x + \bar{t}v$  is a boundary point. Thus,  $y$  is a point in  $X$ .  $\square$

## 16.4 Minimization

This section focuses on minimizing self-concordant functions, and the results are largely based on the following theorem, which also plays a significant role in our study of Newton's algorithm in the next section.



www.sylvania.com

We do not reinvent  
the wheel we reinvent  
light.

Fascinating lighting offers an infinite spectrum of possibilities: Innovative technologies and new markets provide both opportunities and challenges. An environment in which your expertise is in high demand. Enjoy the supportive working atmosphere within our global group and benefit from international career paths. Implement sustainable ideas in close cooperation with other specialists and contribute to influencing our future. Come and join us in reinventing light every day.

Light is OSRAM

OSRAM  
SYLVANIA

**Theorem 16.4.1.** *Let  $f: X \rightarrow \mathbf{R}$  be a closed self-concordant function, suppose that  $x \in X$  is a point with finite Newton decrement  $\lambda = \lambda(f, x)$ , let  $\Delta x_{\text{nt}}$  be a Newton direction at  $x$ , and define*

$$x^+ = x + (1 + \lambda)^{-1} \Delta x_{\text{nt}}.$$

*The point  $x^+$  is then a point in  $X$  and*

$$f(x^+) \leq f(x) - \rho(-\lambda).$$

*Remark.* So a minimum point  $\hat{x}$  of  $f$  must satisfy the inequality

$$f(\hat{x}) \leq f(x) - \rho(-\lambda)$$

for all  $x \in X$  with finite Newton decrement  $\lambda$ .

*Proof.* The vector  $v = (1 + \lambda)^{-1} \Delta x_{\text{nt}}$  has local seminorm

$$\|v\|_x = (1 + \lambda)^{-1} \|\Delta x_{\text{nt}}\|_x = \lambda(1 + \lambda)^{-1} < 1,$$

so it follows from Theorem 16.3.2 that the point  $x^+ = x + v$  lies in  $X$  and that

$$\begin{aligned} f(x^+) &\leq f(x) + Df(x)[v] + \rho(\|v\|_x) = f(x) + \frac{1}{1 + \lambda} \langle f'(x), \Delta x_{\text{nt}} \rangle + \rho\left(\frac{\lambda}{1 + \lambda}\right) \\ &= f(x) - \frac{\lambda^2}{1 + \lambda} - \frac{\lambda}{1 + \lambda} - \ln \frac{1}{1 + \lambda} = f(x) - \lambda + \ln(1 + \lambda) \\ &= f(x) - \rho(-\lambda). \end{aligned} \quad \square$$

**Theorem 16.4.2.** *The Newton decrement  $\lambda(f, x)$  of a downwards bounded closed self-concordant function  $f: X \rightarrow \mathbf{R}$  is finite at each point  $x \in X$  and  $\inf_{x \in X} \lambda(f, x) = 0$ .*

*Proof.* Let  $v$  be an arbitrary vector in the recessive subspace  $V_f = \mathcal{N}(f''(x))$ . Then

$$f(x + tv) = f(x) + t \langle f'(x), v \rangle$$

for all  $t \in \mathbf{R}$  according to Theorem 16.2.1, and since  $f$  is supposed to be bounded below, this implies that  $\langle f'(x), v \rangle = 0$ . This proves the implication

$$f''(x)v = 0 \Rightarrow \langle f'(x), v \rangle = 0,$$

which means that there exists a Newton direction at the point  $x$ . Hence,  $\lambda(f, x)$  is a finite number.

If there is a positive number  $\delta$  such that  $\lambda(f, x) \geq \delta$  for all  $x \in X$ , then repeated application of Theorem 16.4.1, with an arbitrary point  $x_0 \in X$  as starting point, results in a sequence  $(x_k)_0^\infty$  of points in  $X$ , defined as

$x_{k+1} = x_k^+$  and satisfying the inequality  $f(x_k) \leq f(x_0) - k\rho(-\delta)$  for all  $k$ . Since  $\rho(-\delta) > 0$ , this contradicts our assumption that  $f$  is bounded below. Thus,  $\inf_{x \in X} \lambda(f, x) = 0$ .  $\square$

**Theorem 16.4.3.** *All sublevel sets of a non-degenerate closed self-concordant function  $f: X \rightarrow \mathbf{R}$  are compact sets if  $\lambda(f, x_0) < 1$  for some point  $x_0 \in X$ .*

*Proof.* The sublevel sets are closed since the function is closed, and to prove that they are also bounded it is enough to prove that the particular sublevel set  $S = \{x \in X \mid f(x) \leq f(x_0)\}$  is bounded, because of Theorem 6.8.3 in Part I.

So, let  $x$  be an arbitrary point in  $S$ , and write  $r = \|x - x_0\|_{x_0}$  and  $\lambda_0 = \lambda(f, x_0)$  for short. Then

$$f(x) \geq f(x_0) + Df(x_0)[x - x_0] + \rho(-r),$$

according to Theorem 16.3.2, and

$$Df(x_0)[x - x_0] = \langle f'(x_0), x - x_0 \rangle \geq -\lambda(f, x_0)\|x - x_0\|_{x_0} = -\lambda_0 r,$$

by Theorem 15.1.2. Combining these two inequalities we obtain the inequality

$$f(x_0) \geq f(x) \geq f(x_0) - \lambda_0 r + \rho(-r),$$

which simplifies to

$$r - \ln(1 + r) = \rho(-r) \leq \lambda_0 r.$$

Hence,

$$(1 - \lambda_0)r \leq \ln(1 + r)$$

and it follows that  $r \leq r_0$ ,  $r_0$  being the unique positive root of the equation  $(1 - \lambda_0)r = \ln(1 + r)$ . The sublevel set  $S$  is thus included in the ellipsoid  $\{x \in \mathbf{R}^n \mid \|x - x_0\|_{x_0} \leq r_0\}$ , and it is therefore a bounded set.  $\square$

**Theorem 16.4.4.** *A closed self-concordant function  $f: X \rightarrow \mathbf{R}$  has a minimum point if  $\lambda(f, x_0) < 1$  for some point  $x_0 \in X$ .*

*Proof.* If in addition  $f$  is non-degenerate, then  $S = \{x \in X \mid f(x) \leq f(x_0)\}$  is a compact set according to the previous theorem, so the restriction of  $f$  to the sublevel set  $S$  attains a minimum, and this minimum is clearly a global minimum of  $f$ . The minimum point is furthermore unique, since non-degenerate self-concordant functions are strictly convex.

If  $f$  is degenerate, then there is a non-degenerate closed self-concordant function  $g: X_0 \rightarrow \mathbf{R}$  with the same range as  $f$ , according to the discussion following Theorem 16.2.1. The relationship between the two functions has the form  $g(y) = f(Ay + v)$ , where  $A$  is an injective linear map and  $v$  is

an arbitrary vector in the recessive subspace  $V_f$ . To the point  $x_0$  there corresponds a point  $y_0 \in X_0$  such that  $Ay_0 + v = x_0$  for some  $v \in V_f$ , and  $\lambda(g, y_0) \leq \lambda(f, x_0) < 1$ . By the already proven part of the theorem,  $g$  has a minimum point  $\hat{y}$ , and this implies that all points in the set  $A\hat{y} + V_f$  are minimum points of  $f$ .  $\square$

**Theorem 16.4.5.** *Every downwards bounded closed self-concordant function  $f: X \rightarrow \mathbf{R}$  has a minimum point.*

*Proof.* It follows from Theorem 16.4.2 that there is a point  $x_0 \in X$  such that  $\lambda(f, x_0) < 1$ , so the theorem is a corollary of Theorem 16.4.4.  $\square$

Our next theorem describes how well a given point approximates the minimum point of a closed self-concordant function.

**Theorem 16.4.6.** *Let  $f: X \rightarrow \mathbf{R}$  be a closed self-concordant function with a minimum point  $\hat{x}$ . If  $x \in X$  is an arbitrary point with Newton decrement  $\lambda = \lambda(f, x) < 1$ , then*

$$(16.9) \quad \rho(-\lambda) \leq f(x) - f(\hat{x}) \leq \rho(\lambda),$$

$$(16.10) \quad \frac{\lambda}{1 + \lambda} \leq \|x - \hat{x}\|_x \leq \frac{\lambda}{1 - \lambda},$$

$$(16.11) \quad \|x - \hat{x}\|_{\hat{x}} \leq \frac{\lambda}{1 - \lambda}.$$

*Remark.* Since  $\rho(t) \leq t^2$  if  $t \leq \frac{1}{2}$ , we conclude from inequality (16.9) that

$$f(x) - f_{\min} \leq \lambda(f, x)^2$$

as soon as  $\lambda(f, x) \leq \frac{1}{2}$ .

*Proof.* To simplify the notation, let  $v = x - \hat{x}$  and  $r = \|v\|_x$ .

The left part of inequality (16.9) follows directly from the remark after Theorem 16.4.1. To prove the right part of the same inequality, we recall the inequality

$$(16.12) \quad \langle f'(x), v \rangle \leq \lambda(f, x) \|v\|_x = \lambda r,$$

which we combine with the left part of inequality (16.5) in Theorem 16.3.2 and inequality (iii) in Lemma 16.3.1. This results in the following chain of inequalities:

$$\begin{aligned} f(\hat{x}) &= f(x - v) \geq f(x) + \langle f'(x), -v \rangle + \rho(-\|v\|_x) \\ &= f(x) - \langle f'(x), v \rangle + \rho(-r) \\ &\geq f(x) - \lambda r + \rho(-r) \geq f(x) - \rho(\lambda), \end{aligned}$$

and the proof of inequality (16.9) is now complete.

Since  $x - v = \hat{x}$  and  $f'(\hat{x}) = 0$ , it follows from inequality (16.12) and the left part of inequality (16.4) that

$$\lambda r \geq \langle f'(x), v \rangle = \langle f'(x - v), -v \rangle - \langle f'(x), -v \rangle \geq \frac{\| -v \|_x^2}{1 + \| -v \|_x} = \frac{r^2}{1 + r},$$

and by solving the inequality above with respect to  $r$ , we obtain the right part of inequality (16.10).

The left part of the same inequality obviously holds if  $r \geq 1$ . So assume that  $r < 1$ . Due to inequality (16.7),

$$\langle f'(x), w \rangle = \langle f'(x - v), -w \rangle - \langle f'(x), -w \rangle \leq \frac{\| -v \|_x \| -w \|_x}{1 - \| -v \|_x} = \frac{r}{1 - r} \|w\|_x,$$

and hence

$$\lambda = \sup_{\|w\|_x \leq 1} \langle f'(x), w \rangle \leq \frac{r}{1 - r},$$

which gives the left part of inequality (16.10).

To prove the remaining inequality (16.11), we use the left part of inequality (16.5) with  $y$  replaced by  $x$  and  $x$  replaced by  $\hat{x}$ , which results in the inequality

$$\rho(-\|x - \hat{x}\|_{\hat{x}}) \leq f(x) - f(\hat{x}).$$



Discover the truth at [www.deloitte.ca/careers](http://www.deloitte.ca/careers)

**Deloitte.**

© Deloitte & Touche LLP and affiliated entities.

According to the already proven inequality (16.9),  $f(x) - f(\hat{x}) \leq \rho(\lambda)$ , so it follows that  $\rho(-\|x - \hat{x}\|_{\hat{x}}) \leq \rho(\lambda)$ , and by Lemma 16.3.1, this means that  $\|x - \hat{x}\|_{\hat{x}} \leq \frac{\lambda}{1 - \lambda}$ .  $\square$

**Theorem 16.4.7.** *Let  $f$  be a closed self-concordant function whose domain  $X$  is a subset of  $\mathbf{R}^n$ , and suppose that*

$$\nu = \sup\{\lambda(f, x) \mid x \in X\} < 1.$$

*Then  $X$  is equal to the whole space  $\mathbf{R}^n$ , and  $f$  is a constant function.*

*Proof.* It follows from Theorem 16.4.4 that  $f$  has a minimum point  $\hat{x}$  and from inequality (16.9) in Theorem 16.4.6 that

$$\rho(-\nu) \leq f(x) - f(\hat{x}) \leq \rho(\nu)$$

for all  $x \in X$ . Thus,  $f$  is a bounded function, and since  $f$  is closed, this implies that  $X$  is a set without boundary points. Hence,  $X = \mathbf{R}^n$ .

Let  $v$  be an arbitrary vector in  $\mathbf{R}^n$ . By applying inequality (16.11) with  $x = \hat{x} + tv$ , we obtain the inequality

$$t\|v\|_{\hat{x}} = \|x - \hat{x}\|_{\hat{x}} \leq \frac{\lambda(f, x)}{1 - \lambda(f, x)} \leq \frac{\nu}{1 - \nu}$$

for all  $t > 0$ , and this implies that  $\|v\|_{\hat{x}} = 0$ . The recessive subspace  $V_f$  of  $f$  is in other words equal to  $\mathbf{R}^n$ , so  $f$  is a constant function according to Theorem 16.2.1.  $\square$

## 16.5 Newton's method for self-concordant functions

In this section we show that Newton's method converges when the objective function  $f: X \rightarrow \mathbf{R}$  is closed, self-concordant and bounded below. We shall also give an estimate of the number of iterations needed to obtain the minimum with a given accuracy  $\epsilon$  – an estimate that only depends on  $\epsilon$  and the difference between the minimum value and the function value at the starting point. The algorithm starts with a damped phase, which requires no line search as the step length at the point  $x$  can be chosen equal to  $1/(1 + \lambda(f, x))$ , and then enters into a pure phase with quadratic convergence, when the Newton decrement is sufficiently small.

## The damped phase

During the damped phase, the points  $x_k$  in Newton's algorithm are generated recursively by the equation

$$x_{k+1} = x_k + \frac{1}{1 + \lambda_k} v_k,$$

where  $\lambda_k = \lambda(f, x_k)$  is the Newton decrement at  $x_k$  and  $v_k$  is a Newton direction at the same point, i.e

$$f''(x_k)v_k = -f'(x_k).$$

According to Theorem 16.4.1, if the starting point  $x_0$  is a point in  $X$ , then all generated points  $x_k$  will lie in  $X$  and

$$f(x_{k+1}) - f(x_k) \leq \rho(-\lambda_k).$$

If  $\delta > 0$  and  $\lambda_k \geq \delta$ , then  $\rho(-\lambda_k) \geq \rho(-\delta)$ , because the function  $\rho(t)$  is decreasing for  $t < 0$ . So if  $x_N$  is the first point of the sequence that satisfies the inequality  $\lambda_N = \lambda(f, x_N) < \delta$ , then

$$\begin{aligned} f_{\min} - f(x_0) &\leq f(x_N) - f(x_0) = \sum_{k=0}^{N-1} (f(x_{k+1}) - f(x_k)) \\ &\leq -\sum_{k=0}^{N-1} \rho(-\lambda_k) \leq -\sum_{k=0}^{N-1} \rho(-\delta) = -N\rho(-\delta), \end{aligned}$$

which implies that  $N \leq (f(x_0) - f_{\min})/\rho(-\delta)$ . This proves the following theorem.

**Theorem 16.5.1.** *Let  $f: X \rightarrow \mathbf{R}$  be a closed, self-concordant and downwards bounded function. Using Newton's damped algorithm with step size as above, we need at most*

$$\left\lceil \frac{f(x_0) - f_{\min}}{\rho(-\delta)} \right\rceil$$

*iterations to generate a point  $x$  with Newton decrement  $\lambda(f, x) < \delta$  from an arbitrary starting point  $x_0$  in  $X$ .*

## Local convergence

We now turn to the study of Newton's pure method for starting points that are sufficiently close to the minimum point  $\hat{x}$ . For a corresponding analysis of Newton's damped method we refer to exercise 16.6.



**Theorem 16.5.2.** Let  $f: X \rightarrow \mathbf{R}$  be a closed self-concordant function, and suppose that  $x \in X$  is a point with Newton decrement  $\lambda(f, x) < 1$ . Let  $\Delta x_{\text{nt}}$  be a Newton direction at  $x$ , and let

$$x^+ = x + \Delta x_{\text{nt}}.$$

Then,  $x^+$  is a point in  $X$  and


$$\lambda(f, x^+) \leq \left( \frac{\lambda(f, x)}{1 - \lambda(f, x)} \right)^2.$$

*Proof.* The conclusion that  $x^+$  lies in  $X$  follows from Theorem 16.3.2, because  $\|\Delta x_{\text{nt}}\|_x = \lambda(f, x) < 1$ . To prove the inequality for  $\lambda(f, x^+)$ , we first use inequality (16.7) of the same theorem with  $v = x^+ - x = \Delta x_{\text{nt}}$  and obtain


$$\begin{aligned} \langle f'(x^+), w \rangle &\leq \langle f'(x), w \rangle + \langle f''(x)\Delta x_{\text{nt}}, w \rangle + \frac{\lambda(f, x)^2 \|w\|_x}{1 - \lambda(f, x)} \\ &= \langle f'(x), w \rangle + \langle -f'(x), w \rangle + \frac{\lambda(f, x)^2 \|w\|_x}{1 - \lambda(f, x)} = \frac{\lambda(f, x)^2 \|w\|_x}{1 - \lambda(f, x)}. \end{aligned}$$

But

$$\|w\|_x \leq \frac{\|w\|_{x^+}}{1 - \lambda(f, x)},$$

SIMPLY CLEVER


**We will turn your CV into an opportunity of a lifetime**



Do you like cars? Would you like to be a part of a successful brand? We will appreciate and reward both your enthusiasm and talent. Send us your CV. You will be surprised where it can take you.

Send us your CV on  
[www.employerforlife.com](http://www.employerforlife.com)

by inequality (16.6), so it follows that

$$\langle f'(x^+), w \rangle \leq \frac{\lambda(f, x)^2 \|w\|_{x^+}}{(1 - \lambda(f, x))^2},$$

and this implies that

$$\lambda(f, x^+) = \sup_{\|w\|_{x^+} \leq 1} \langle f'(x^+), w \rangle \leq \frac{\lambda(f, x)^2}{(1 - \lambda(f, x))^2}. \quad \square$$

We are now able to prove the following convergence result for Newton's pure method.

**Theorem 16.5.3.** *Suppose that  $f: X \rightarrow \mathbf{R}$  is a closed self-concordant function and that  $x_0$  is a point in  $X$  with Newton decrement*

$$\lambda(f, x_0) \leq \delta < \bar{\lambda} = \frac{1}{2}(3 - \sqrt{5}) = 0.381966 \dots$$

*Let the sequence  $(x_k)_0^\infty$  be recursively defined by*

$$x_{k+1} = x_k + v_k,$$

*where  $v_k$  is a Newton direction at the point  $x_k$ .*

*The sequence  $(f(x_k))_0^\infty$  converges to the minimum value  $f_{\min}$  of the function  $f$ , and if  $\epsilon > 0$  then*

$$f(x_k) - f_{\min} < \epsilon$$

*for  $k > A + \log_2(\log_2 B/\epsilon)$ , where  $A$  and  $B$  are constants that only depend on  $\delta$ .*

*Moreover, if  $f$  is a non-degenerate function, then  $(x_k)_0^\infty$  converges to the unique minimum point of  $f$ .*

*Proof.* The critical number  $\bar{\lambda}$  is a root of the equation  $(1 - \lambda)^2 = \lambda$ , and if  $0 \leq \lambda < \bar{\lambda}$  then  $\lambda < (1 - \lambda)^2$ .

Let  $K(\lambda) = (1 - \lambda)^{-2}$ ; the function  $K$  is increasing in the interval  $[0, \bar{\lambda}[$  and  $K(\lambda)\lambda < 1$ . It therefore follows from Theorem 16.5.2 that the following inequality is true for all points  $x \in X$  with  $\lambda(f, x) \leq \delta < \bar{\lambda}$ :

$$\lambda(f, x^+) \leq K(\lambda(f, x)) \lambda(f, x)^2 \leq K(\delta) \lambda(f, x)^2 \leq K(\delta) \delta \lambda(f, x) \leq \lambda(f, x) \leq \delta.$$

Now, let  $\lambda_k = \lambda(f, x_k)$ . Due to the inequality above, it follows by induction that  $\lambda_k \leq \delta$  and that

$$\lambda_{k+1} \leq K(\delta) \lambda_k^2$$

for all  $k$ , and the latter inequality in turn implies that

$$\lambda_k \leq K(\delta)^{-1} (K(\delta) \lambda_0)^{2^k} \leq (1 - \delta)^2 (K(\delta) \delta)^{2^k}.$$

Hence,  $\lambda_k$  tends to 0 as  $k \rightarrow \infty$ , because  $K(\delta)\delta < 1$ . By the remark following Theorem 16.4.6,

$$f(x_k) - f_{\min} \leq \lambda_k^2,$$

if  $\lambda_k \leq \frac{1}{2}$ , so we conclude that

$$\lim_{k \rightarrow \infty} f(x_k) = f_{\min}.$$

To prove the remaining error estimate, we can without loss of generalization assume that  $\epsilon < \delta^2$ , because if  $\epsilon > \delta^2$  then already

$$f(x_0) - f_{\min} \leq \lambda(f, x_0)^2 \leq \delta^2 < \epsilon.$$

Let  $A$  and  $B$  be the constants defined by

$$A = -\log_2(-2\log_2(K(\delta)\delta)) \quad \text{and} \quad B = (1 - \delta)^4.$$

Then  $0 < B \leq 1$ , and  $\log_2(\log_2 B/\epsilon)$  is a well-defined number, since  $B/\epsilon \geq (1 - \delta)^4/\delta^2 = (K(\delta)\delta)^{-2} > 1$ . If  $k > A + \log_2(\log_2 B/\epsilon)$ , then

$$\lambda_k^2 \leq (1 - \delta)^4 (K(\delta)\delta)^{2^{k+1}} < \epsilon,$$

and consequently  $f(x_k) - f_{\min} \leq \lambda_k^2 < \epsilon$ .

If  $f$  is a non-degenerate function, then  $f$  has a unique minimum point  $\hat{x}$ , and it follows from inequality (16.11) in Theorem 16.4.6 that

$$\|x_k - \hat{x}\|_{\hat{x}} \leq \frac{\lambda_k}{1 - \lambda_k} \rightarrow 0, \quad \text{as } k \rightarrow \infty.$$

Since  $\|\cdot\|_{\hat{x}}$  is a proper norm, this means that  $x_k \rightarrow \hat{x}$ . □

When  $\delta = 1/3$ , the values of the constants in Theorem 16.5.3 are  $A = 0.268\dots$  and  $B = 16/81$ , and  $A + \log_2(\log_2 B/\epsilon) = 6.87$  for  $\epsilon = 10^{-30}$ . So with a starting point  $x_0$  satisfying  $\lambda(f, x_0) < 1/3$ , Newton's algorithm will produce a function value that approximates the minimum value with an error less than  $10^{-30}$  after at most 7 iterations.

## Newton's method for self-concordant functions

By combining Newton's damped method with  $1/(1 + \lambda(f, x))$  as damping factor and Newton's pure method, we arrive at the following variant of Newton's method.

**Newton's method**

**Given** a positive number  $\delta < \frac{1}{2}(3 - \sqrt{5})$ , a starting point  $x_0 \in X$ , and a tolerance  $\epsilon > 0$ .

1. *Initiate:*  $x := x_0$ .
2. Compute the Newton decrement  $\lambda = \lambda(f, x)$ .
3. Go to line 8 if  $\lambda < \delta$  else continue.
4. Compute a Newton direction  $\Delta x_{\text{nt}}$  at the point  $x$ .
5. *Update:*  $x := x + (1 + \lambda)^{-1} \Delta x_{\text{nt}}$ .
6. Go to line 2.
7. Compute the Newton decrement  $\lambda = \lambda(f, x)$ .
8. *Stopping criterion:* **stop** if  $\lambda < \sqrt{\epsilon}$ .  $x$  is an approximate optimal point.
9. Compute a Newton direction  $\Delta x_{\text{nt}}$  at the point  $x$ .
10. *Update:*  $x := x + \Delta x_{\text{nt}}$ .
11. Go to line 7.

Assuming that  $f$  is closed, self-concordant and downwards bounded, the damped phase of the algorithm, i.e. steps 2–6, continues during at most

$$\lfloor (f(x_0) - f_{\min}) / \rho(-\delta) \rfloor$$

I joined MITAS because  
I wanted **real responsibility**

The Graduate Programme  
for Engineers and Geoscientists  
[www.discovermitas.com](http://www.discovermitas.com)



**Month 16**  
I was a construction supervisor in the North Sea advising and helping foremen solve problems

Real work  
International opportunities  
Three work placements



**MAERSK**

iterations, and the pure phase 7–11 ends according to Theorem 16.5.3 after at most  $\lceil A + \log_2(\log_2 B/\epsilon) \rceil$  iterations. Therefore, we have the following result.

**Theorem 16.5.4.** *If the function  $f$  is closed, self-concordant and bounded below, then the above Newton method terminates at a point  $x$  satisfying  $f(x) < f_{\min} + \epsilon$  after at most*

$$\lfloor (f(x_0) - f_{\min})/\rho(-\delta) \rfloor + \lceil A + \log_2(\log_2 B/\epsilon) \rceil$$

*iterations, where  $A$  and  $B$  are the constants of Theorem 16.5.3.*

In particular,  $1/\rho(-\delta) = 21.905$  when  $\delta = 1/3$ , and the second term can be replaced by the number 7 when  $\epsilon \geq 10^{-30}$ . Thus, at most

$$\lfloor 22(f(x_0) - f_{\min}) \rfloor + 7$$

iterations are required to find an approximation to the minimum value that meets all practical requirements by a wide margin.

## Exercises

**16.1** Show that the function  $f(x) = x \ln x - \ln x$  is self-concordant on  $\mathbf{R}_{++}$ .

**16.2** Suppose  $f_i: X_i \rightarrow \mathbf{R}$  are self-concordant functions for  $i = 1, 2, \dots, m$ , and let  $X = X_1 \times X_2 \times \dots \times X_m$ . Prove that the function  $f: X \rightarrow \mathbf{R}$ , defined by

$$f(x_1, x_2, \dots, x_m) = f_1(x_1) + f_2(x_2) + \dots + f_m(x_m)$$

for  $x = (x_1, x_2, \dots, x_m) \in X$ , is self-concordant.

**16.3** Suppose that  $f: \mathbf{R}_{++} \rightarrow \mathbf{R}$  is a three times continuously differentiable, convex function, and that

$$|f'''(x)| \leq 3 \frac{f''(x)}{x} \quad \text{for all } x.$$

a) Prove that the function

$$g(x) = -\ln(-f(x)) - \ln x,$$

with  $\{x \in \mathbf{R}_{++} \mid f(x) < 0\}$  as domain, is self-concordant.

[Hint: Use that  $3a^2b + 3a^2c + 2b^3 + 2c^3 \leq 2(a^2 + b^2 + c^2)^{3/2}$  if  $a, b, c \geq 0$ .]

b) Prove that the function

$$F(x, y) = -\ln(y - f(x)) - \ln x$$

is self-concordant on the set  $\{(x, y) \in \mathbf{R}^2 \mid x > 0, y > f(x)\}$ .

**16.4** Show that the following functions  $f$  satisfy the conditions of the previous exercise:

a)  $f(x) = -\ln x$       b)  $f(x) = x \ln x$       c)  $f(x) = -x^p$ , where  $0 < p \leq 1$ .

**16.5** Let us write  $x'$  for  $(x_1, x_2, \dots, x_{n-1})$  when  $x = (x_1, x_2, \dots, x_n)$ , and let  $\|\cdot\|$  denote the Euclidean norm in  $\mathbf{R}^{n-1}$ . Let  $X = \{x \in \mathbf{R}^n \mid \|x'\| < x_n\}$ , and define the function  $f: X \rightarrow \mathbf{R}$  by  $f(x) = -\ln(x_n^2 - \|x'\|^2)$ . Prove that the following identity holds for all  $v \in \mathbf{R}^n$ :

$$D^2 f(x)[v, v] = \frac{1}{2} (Df(x)[v])^2 + 2 \frac{(x_n^2 - \|x'\|^2)(\|x'\|^2 \|v'\|^2 - \langle x', v' \rangle^2) + (v_n \|x'\|^2 - x_n \langle x', v' \rangle)^2}{(x_n^2 - \|x'\|^2)^2 \|x'\|^2},$$

and use it to conclude that  $f$  is a convex function and that  $\lambda(f, x) = 2$  for all  $x \in X$ .

**16.6** *Convergence for Newton's damped method.*

Suppose that the function  $f: X \rightarrow \mathbf{R}$  is closed and self-concordant, and define for points  $x \in X$  with finite Newton decrement the point  $x^+$  by

$$x^+ = x + \frac{1}{1 + \lambda(f, x)} \Delta x_{\text{nt}},$$

where  $\Delta x_{\text{nt}}$  is a Newton direction at  $x$ .

a) Then  $x^+$  is a point in  $X$ , according to Theorem 16.3.2. Show that

$$\lambda(f, x^+) \leq 2\lambda(f, x)^2,$$

and hence that  $\lambda(f, x^+) \leq \lambda(f, x)$  if  $\lambda(f, x) \leq \frac{1}{2}$ .

b) Suppose  $x_0$  is a point in  $X$  with Newton decrement  $\lambda(f, x_0) \leq \frac{1}{4}$ , and define the sequence  $(x_k)_0^\infty$  recursively by  $x_{k+1} = x_k^+$ . Show that

$$f(x_k) - f_{\min} \leq \frac{1}{4} \cdot \left(\frac{1}{2}\right)^{2^{k+1}},$$

and hence that  $f(x_k)$  converges quadratically to  $f_{\min}$ .

## Appendix

We begin with a result on tri-linear forms which was needed in the proof of the fundamental inequality  $|D^3 f(x)[u, v, w]| \leq 2\|u\|_x \|v\|_x \|w\|_x$  for self-concordant functions.

Fix an arbitrary scalar product  $\langle \cdot, \cdot \rangle$  on  $\mathbf{R}^n$  and let  $\|\cdot\|$  denote the corresponding norm, i.e.  $\|v\| = \langle v, v \rangle^{1/2}$ . If  $\phi(u, v, w)$  is a symmetric tri-linear form on  $\mathbf{R}^n \times \mathbf{R}^n \times \mathbf{R}^n$ , we define its norm  $\|\phi\|$  by

$$\|\phi\| = \sup_{u,v,w \neq 0} \frac{|\phi(u, v, w)|}{\|u\| \|v\| \|w\|}.$$

The numerator and the denominator in the expression for  $\|\phi\|$  are homogeneous of the same degree 3, hence

$$\|\phi\| = \sup_{(u,v,w) \in S^3} |\phi(u, v, w)|,$$

where  $S$  denotes the unit sphere in  $\mathbf{R}^n$  with respect to the norm  $\|\cdot\|$ , i.e.

$$S = \{u \in \mathbf{R}^n \mid \|u\| = 1\}.$$

It follows from the norm definition that

$$|\phi(u, v, w)| \leq \|\phi\| \|u\| \|v\| \|w\|$$

for all vectors  $u, v, w$  in  $\mathbf{R}^n$ .

**ie** business school

#1 EUROPEAN BUSINESS SCHOOL  
FINANCIAL TIMES 2013

**#gobeyond**

**MASTER IN MANAGEMENT**

**Because achieving your dreams is your greatest challenge.** IE Business School's Master in Management taught in English, Spanish or bilingually, trains young high performance professionals at the beginning of their career through an innovative and stimulating program that will help them reach their full potential.

- Choose your area of specialization.
- Customize your master through the different options offered.
- Global Immersion Weeks in locations such as London, Silicon Valley or Shanghai.

*Because you change, we change with you.*

www.ie.edu/master-management | mim.admissions@ie.edu

f t in YouTube

Since tri-linear forms are continuous and the unit sphere is compact, the least upper bound  $\|\phi\|$  is attained at some point  $(u, v, w) \in S^3$ , and we will show that the least upper bound is indeed attained at some point where  $u = v = w$ . This is the meaning of the following theorem.

**Theorem 1.** *Suppose that  $\phi(u, v, w)$  is a symmetric tri-linear form. Then*

$$\|\phi\| = \sup_{u, v, w \neq 0} \frac{|\phi(u, v, w)|}{\|u\| \|v\| \|w\|} = \sup_{v \neq 0} \frac{|\phi(v, v, v)|}{\|v\|^3}.$$

*Remark.* The theorem is a special case of the corresponding result for symmetric  $m$ -multilinear forms, but we only need the case  $m = 3$ . The general case is proved by induction.

*Proof.* Let

$$\|\phi\|' = \sup_{v \neq 0} \frac{|\phi(v, v, v)|}{\|v\|^3} = \sup_{\|v\|=1} |\phi(v, v, v)|.$$

We claim that  $\|\phi\| = \|\phi\|'$ . Obviously,  $\|\phi\|' \leq \|\phi\|$ , so we only have to prove the converse inequality  $\|\phi\| \leq \|\phi\|'$ .

To prove this inequality, we need the corresponding result for symmetric bilinear forms  $\psi(u, v)$ . To such a form there is associated a symmetric linear operator (matrix)  $A$  such that  $\psi(u, v) = \langle Au, v \rangle$ , and if  $e_1, e_2, \dots, e_n$  is an ON-basis of eigenvectors of  $A$  and  $\lambda_1, \lambda_2, \dots, \lambda_n$  denote the corresponding eigenvalues with  $\lambda_1$  as the one with the largest absolute value, and if  $u, v \in S$  are vectors with coordinates  $u_1, u_2, \dots, u_n$  and  $v_1, v_2, \dots, v_n$  with respect to the given ON-basis, then

$$\begin{aligned} |\psi(u, v)| &= \left| \sum_{i=1}^n \lambda_i u_i v_i \right| \leq \sum_{i=1}^n |\lambda_i| |u_i| |v_i| \leq |\lambda_1| \sum_{i=1}^n |u_i| |v_i| \\ &\leq |\lambda_1| \left( \sum_{i=1}^n u_i^2 \right)^{1/2} \left( \sum_{i=1}^n v_i^2 \right)^{1/2} = |\lambda_1| = |\psi(e_1, e_1)|, \end{aligned}$$

which proves that  $\sup_{(u, v) \in S^2} |\psi(u, v)| = \sup_{v \in S} |\psi(v, v)|$ .

We now return to the tri-linear form  $\phi(u, v, w)$ . Let  $(\hat{u}, \hat{v}, \hat{w})$  be a point in  $S^3$  where the least upper bound defining  $\|\phi\|$  is attained, i.e.

$$\|\phi\| = \phi(\hat{u}, \hat{v}, \hat{w}),$$

and consider the function

$$\psi(u, v) = \phi(u, v, \hat{w});$$

this is a symmetric bilinear form on  $\mathbf{R}^n \times \mathbf{R}^n$  and



$$\sup_{(u,v) \in S^2} |\psi(u, v)| = \|\phi\|.$$

But as already proven,

$$\sup_{(u,v) \in S^2} |\psi(u, v)| = \sup_{v \in S} |\psi(v, v)|.$$

Therefore, we conclude that we can without restriction assume that  $\hat{u} = \hat{v}$ .

We have in other words shown that the set

$$A = \{(v, w) \in S^2 \mid |\phi(v, v, w)| = \|\phi\|\}$$

is nonempty. The set  $A$  is a closed subset of  $S^2$ , and hence the number

$$\alpha = \max\{\langle v, w \rangle \mid (v, w) \in A\}$$

exists, and obviously  $0 \leq \alpha \leq 1$ .

Due to tri-linearity,

$$\phi(u + v, u + v, w) - \phi(u - v, u - v, w) = 4\phi(u, v, w).$$

So if  $u, v, w$  are arbitrary vectors in  $S$ , i.e. vectors with norm 1, then

$$\begin{aligned} 4|\phi(u, v, w)| &\leq |\phi(u + v, u + v, w)| + |\phi(u - v, u - v, w)| \\ &\leq |\phi(u + v, u + v, w)| + \|\phi\| \|u - v\|^2 \|w\| \\ &= |\phi(u + v, u + v, w)| - \|\phi\| \|u + v\|^2 + \|\phi\| (\|u + v\|^2 + \|u - v\|^2) \\ &= |\phi(u + v, u + v, w)| - \|\phi\| \|u + v\|^2 + \|\phi\| (2\|u\|^2 + 2\|v\|^2) \\ &= |\phi(u + v, u + v, w)| - \|\phi\| \|u + v\|^2 + 4\|\phi\|. \end{aligned}$$

Now choose  $(\bar{v}, \bar{w}) \in A$  such that  $\langle \bar{v}, \bar{w} \rangle = \alpha$ . By the above inequality, we then have

$$\begin{aligned} 4\|\phi\| &= 4|\phi(\bar{v}, \bar{v}, \bar{w})| = 4|\phi(\bar{v}, \bar{w}, \bar{v})| \\ &\leq |\phi(\bar{v} + \bar{w}, \bar{v} + \bar{w}, \bar{v})| - \|\phi\| \|\bar{v} + \bar{w}\|^2 + 4\|\phi\|, \end{aligned}$$

and it follows that

$$|\phi(\bar{v} + \bar{w}, \bar{v} + \bar{w}, \bar{v})| \geq \|\phi\| \|\bar{v} + \bar{w}\|^2.$$

Note that  $\|\bar{v} + \bar{w}\|^2 = \|\bar{v}\|^2 + \|\bar{w}\|^2 + 2\langle \bar{v}, \bar{w} \rangle = 2 + 2\alpha > 0$ . Therefore, we can form the vector  $\bar{z} = (\bar{v} + \bar{w}) / \|\bar{v} + \bar{w}\|$  and write the above inequality as

$$|\phi(\bar{z}, \bar{z}, \bar{v})| \geq \|\phi\|,$$

which implies that

$$(16.13) \quad |\phi(\bar{z}, \bar{z}, \bar{v})| = \|\phi\|$$

since  $\bar{z}$  and  $\bar{v}$  are vectors in  $S$ . We conclude that the pair  $(\bar{z}, \bar{v})$  is an element of the set  $A$ , and hence

$$\alpha \geq \langle \bar{z}, \bar{v} \rangle = \frac{\langle \bar{v}, \bar{v} \rangle + \langle \bar{w}, \bar{v} \rangle}{\|\bar{v} + \bar{w}\|} = \frac{1 + \alpha}{\sqrt{2 + 2\alpha}} = \sqrt{\frac{1 + \alpha}{2}}.$$

This inequality forces  $\alpha$  to be greater than or equal to 1. Hence  $\alpha = 1$  and

$$\langle \bar{z}, \bar{v} \rangle = 1 = \|\bar{z}\| \|\bar{v}\|.$$

So Cauchy–Schwarz’s inequality holds with equality in this case, and this implies that  $\bar{z} = \bar{v}$ . By inserting this in equality (16.13), we obtain the inequality

$$\|\phi\|' \geq \phi(\bar{v}, \bar{v}, \bar{v}) = \|\phi\|,$$

and the proof of the theorem is now complete.  $\square$

Our second result in this appendix is a uniqueness theorem for functions that satisfy a special differential inequality.

**Theorem 2.** *Suppose that the function  $y(t)$  is continuously differentiable in the interval  $I = [0, b[$ , that  $y(t) \geq 0$ ,  $y(0) = 0$  and  $y'(t) \leq Cy(t)^\alpha$  for some given constants  $C > 0$  and  $\alpha \geq 1$ . Then,  $y(t) = 0$  in the interval  $I$ .*

*Proof.* Let  $a = \sup\{x \in I \mid y(t) = 0 \text{ for } 0 \leq t \leq x\}$ . We will prove that  $a = b$  by showing that the assumption  $a < b$  gives rise to a contradiction.

By continuity,  $y(a) = 0$ . Choose a point  $c \in ]a, b[$  and let

$$M = \max\{y(t) \mid a \leq t \leq c\}.$$

Then choose a point  $d$  such that  $a < d < c$  and  $d - a \leq \frac{1}{2}C^{-1}M^{1-\alpha}$ . The maximum of the function  $y(t)$  on the interval  $[a, d]$  is attained at some point  $e$ , and by the least upper bound definition of the point  $a$ , we have  $y(e) > 0$ . Of course, we also have  $y(e) \leq M$ , so it follows that

$$\begin{aligned} y(e) &= y(e) - y(a) = \int_a^e y'(t) dt \leq C \int_a^e y(t)^\alpha dt \\ &\leq C(d - a)y(e)^\alpha \leq C(d - a)M^{\alpha-1}y(e) \leq \frac{1}{2}y(e), \end{aligned}$$

which is a contradiction.  $\square$

# Chapter 17

## The path-following method

In this chapter, we describe a method for solving the optimization problem

$$\begin{array}{ll} \min & f(x) \\ \text{s.t.} & x \in X \end{array}$$

when  $X$  is a closed subset of  $\mathbf{R}^n$  with nonempty interior and  $f$  is a continuous function which is differentiable in the interior of  $X$ . We assume throughout that  $X = \text{cl}(\text{int } X)$ . Pretty soon, we will restrict ourselves to convex problems, i.e. assume that  $X$  is a convex set and  $f$  is a convex function, in which case, of course, automatically  $X = \text{cl}(\text{int } X)$  for all sets with nonempty interior.

Descent methods require that the function  $f$  is differentiable in a neighborhood of the optimal point, and if the optimal point lies on the boundary of  $X$ , then we have a problem. One way to attack this problem is to choose a function  $F: \text{int } X \rightarrow \mathbf{R}$  with the property that  $F(x) \rightarrow +\infty$  as  $x$  goes to boundary of  $X$  and a parameter  $\mu > 0$ , and to minimize the function  $f(x) + \mu F(x)$  over  $\text{int } X$ . This function's minimum point  $\hat{x}(\mu)$  lies in the interior of  $X$ , and since  $f(x) + \mu F(x) \rightarrow f(x)$  as  $\mu \rightarrow 0$ , we can hope that the function value  $f(\hat{x}(\mu))$  should be close to the minimum value of  $f$ , if the parameter  $\mu$  is small enough. The function  $F$  acts as a barrier that prevents the approximating minimum point from lying on the boundary.

The function  $\mu^{-1}f(x) + F(x)$  has of course the same minimum point  $\hat{x}(\mu)$  as  $f(x) + \mu F(x)$ , and for technical reasons it works better to have the parameter in front of the objective function  $f$  than in front of the barrier function  $F$ . Henceforth, we will therefore instead, with new notation, examine what happens to the minimum point  $\hat{x}(t)$  of the function  $F_t(x) = tf(x) + F(x)$ , when the parameter  $t$  tends to  $+\infty$ .

## 17.1 Barrier and central path

### Barrier

We begin with the formal definition of a barrier.

**Definition.** Let  $X$  be a closed convex set with nonempty interior. A *barrier* to the set  $X$  is a differentiable function  $F: \text{int } X \rightarrow \mathbf{R}$  with the property that  $\lim_{k \rightarrow \infty} F(x_k) = +\infty$  for all sequences  $(x_k)_1^\infty$  that converge to a boundary point of  $X$ .

If a barrier function has a unique minimum point, then this point is called the *analytic center* of the set  $X$  (with respect to the barrier).

*Remark 1.* A convex function with an open domain goes to  $\infty$  at the boundary if and only if it is a closed function. Hence, if  $F: \text{int } X \rightarrow \mathbf{R}$  is convex and differentiable, then  $F$  is a barrier to  $X$  if and only if  $F$  is closed.

*Remark 2.* A strictly convex barrier function to a compact convex set has a unique minimum point in the interior of the set. So compact convex sets with nonempty interiors have analytic centers with respect to strictly convex barriers.



## STUDY AT A TOP RANKED INTERNATIONAL BUSINESS SCHOOL

Reach your full potential at the Stockholm School of Economics, in one of the most innovative cities in the world. The School is ranked by the Financial Times as the number one business school in the Nordic and Baltic countries.

Visit us at [www.hhs.se](http://www.hhs.se)





Now, let  $F$  be a barrier to the closed convex set  $X$ , and suppose that we want to minimize a given function  $f: X \rightarrow \mathbf{R}$ . For each real number  $t \geq 0$  we define the function  $F_t: \text{int } X \rightarrow \mathbf{R}$  by

$$F_t(x) = tf(x) + F(x).$$

In particular,  $F_0 = F$ . The following theorem is the basis for barrier-based interior-point methods for minimization.

**Theorem 17.1.1.** *Suppose that  $f: X \rightarrow \mathbf{R}$  is a continuous function, and let  $F$  be a downwards bounded barrier to the set  $X$ . Suppose that the functions  $F_t$  have minimum points  $\hat{x}(t)$  in the interior of  $X$  for each  $t > 0$ . Then,*

$$\lim_{t \rightarrow +\infty} f(\hat{x}(t)) = \inf_{x \in X} f(x).$$

*Proof.* Let  $v_{\min} = \inf_{x \in X} f(x)$  and  $M = \inf_{x \in \text{int } X} F(x)$ . (We do not exclude the possibility that  $v_{\min} = -\infty$ , but  $M$  is of course a finite number.)

Choose, given  $\eta > v_{\min}$ , a point  $x^* \in \text{int } X$  such that  $f(x^*) < \eta$ . Then

$$\begin{aligned} v_{\min} &\leq f(\hat{x}(t)) \leq f(\hat{x}(t)) + t^{-1}(F(\hat{x}(t)) - M) = t^{-1}(F_t(\hat{x}(t)) - M) \\ &\leq t^{-1}(F_t(x^*) - M) = f(x^*) + t^{-1}(F(x^*) - M). \end{aligned}$$

Since the right hand side of this inequality tends to  $f(x^*)$  as  $t \rightarrow +\infty$ , it follows that  $v_{\min} \leq f(\hat{x}(t)) < \eta$  for all sufficiently large numbers  $t$ , and this proves the theorem.  $\square$

In order to use the barrier method, one needs of course an appropriate barrier to the given set. For sets of the type

$$X = \{x \in \Omega \mid g_i(x) \leq 0, \quad i = 1, 2, \dots, m\}$$

we will use the *logarithmic barrier function*

$$(17.1) \quad F(x) = - \sum_{i=1}^m \ln(-g_i(x)).$$

Note that the barrier function  $F$  is convex if all functions  $g_i: \Omega \rightarrow \mathbf{R}$  are convex. In this case,  $X$  is a convex set, and the interior of  $X$  is nonempty if Slater's condition is satisfied, i.e. if there is a point  $\bar{x} \in \Omega$  such that  $g_i(\bar{x}) < 0$  for all  $i$ .

Other examples of barriers are the exponential barrier function

$$F(x) = \sum_{i=1}^m e^{-1/g_i(x)}$$

and the power function barriers

$$F(x) = \sum_{i=1}^m (-g_i(x))^{-p},$$

where  $p > 0$ .

## Central path

**Definition.** Let  $F$  be a barrier to the set  $X$  and suppose that the functions  $F_t$  have unique minimum points  $\hat{x}(t) \in \text{int } X$  for all  $t \geq 0$ . The curve  $\{\hat{x}(t) \mid t \geq 0\}$  is called the *central path* for the problem  $\min_{x \in X} f(x)$ .

Note that  $\hat{x}(0)$  is the analytic center of  $X$  with respect to the barrier  $F$ , so the central path starts at the analytic center.

Since the gradient is zero at an optimal point, we have

$$(17.2) \quad t f'(\hat{x}(t)) + F'(\hat{x}(t)) = 0$$

for all points on the central path. The converse is true if the objective function  $f$  and the barrier function  $F$  are convex, i.e.  $\hat{x}(t)$  is a point on the central path if and only if equation (17.2) is satisfied.

The logarithmic barrier  $F$  to  $X = \{x \in \Omega \mid g_i(x) \leq 0, i = 1, 2, \dots, m\}$  has derivative

$$F'(x) = - \sum_{i=1}^m \frac{1}{g_i(x)} g'_i(x),$$

so the central path equation (17.2) has in this case the following form for  $t > 0$ :

$$(17.3) \quad f'(\hat{x}(t)) - \frac{1}{t} \sum_{i=1}^m \frac{1}{g_i(\hat{x}(t))} g'_i(\hat{x}(t)) = 0.$$

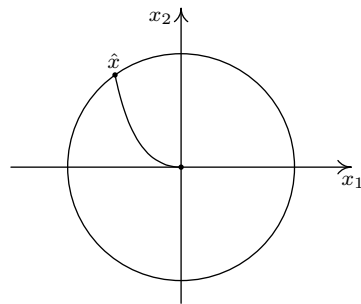
Let us now consider a convex optimization problem of the following type:

$$(P) \quad \begin{array}{ll} \min & f(x) \\ \text{s.t.} & g_i(x) \leq 0, \quad i = 1, 2, \dots, m \end{array}$$

We assume that Slater's condition is satisfied and that the problem has an optimal solution  $\hat{x}$ .

The corresponding Lagrange function  $L$  is given by

$$L(x, \lambda) = f(x) + \sum_{i=1}^m \lambda_i g_i(x),$$



**Figure 17.1.** The central path associated with the problem of minimizing the function  $f(x) = x_1 e^{x_1+x_2}$  over  $X = \{x \in \mathbf{R}^2 \mid x_1^2 + x_2^2 \leq 1\}$  with barrier function  $F(x) = (1 - x_1^2 - x_2^2)^{-1}$ . The minimum point is  $\hat{x} = (-0.5825, 0.8128)$ .

and it follows from equation (17.3) that  $L'_x(\hat{x}(t), \hat{\lambda}) = 0$ , if  $\hat{\lambda} \in \mathbf{R}_+^m$  is the vector defined by

$$\hat{\lambda}_i = -\frac{1}{tg_i(\hat{x}(t))}.$$

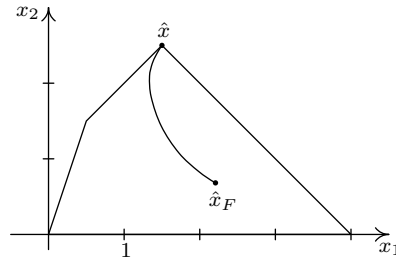
**#1**  
in eco-friendly attitude

**STUDY AT  
LINKÖPING UNIVERSITY, SWEDEN**  
RANKED AMONG TOP 50 UNIVERSITIES UNDER 50

Interested in Strategy and Management in International Organisations? Kick-start your career with a master's degree from Linköping University, Sweden.

→ **Click here!**

**Linköping University**



**Figure 17.2.** The central path for the LP problem  $\min_{x \in X} 2x_1 - 3x_2$  with  $X = \{x \in \mathbf{R}^2 \mid x_2 \geq 0, x_2 \leq 3x_1, x_2 \leq x_1 + 1, x_1 + x_2 \leq 4\}$  and logarithmic barrier. The point  $\hat{x}_F$  is the analytic center of  $X$ , and  $\hat{x} = (1.5, 2.5)$  is the optimal solution.

Since the Lagrange function is convex in the variable  $x$ , we conclude that  $\hat{x}(t)$  is a minimum point for the function  $L(\cdot, \hat{\lambda})$ . The value at  $\hat{\lambda}$  of the dual function  $\phi: \mathbf{R}_+^m \rightarrow \mathbf{R}$  to our minimization problem (P) is therefore by definition

$$\phi(\hat{\lambda}) = L(\hat{x}(t), \hat{\lambda}) = f(\hat{x}(t)) - m/t.$$

By weak duality,  $\phi(\hat{\lambda}) \leq f(\hat{x})$ , so it follows that

$$f(\hat{x}(t)) - m/t \leq f(\hat{x}).$$

We have thus arrived at the following approximation theorem, which for convex problems with logarithmic barrier provides more precise information than Theorem 17.1.1.

**Theorem 17.1.2.** *The points  $\hat{x}(t)$  on the central path for the convex minimization problem (P) with optimal solution  $\hat{x}$  and logarithmic barrier satisfy the inequality*

$$f(\hat{x}(t)) - f(\hat{x}) \leq \frac{m}{t}.$$

Note that the estimate of the theorem depends on the number of constraints but not on the dimension.

## 17.2 Path-following methods

A strategy for determining the optimal value of the convex optimization problem

$$\begin{aligned} \text{(P)} \quad & \min f(x) \\ & \text{s.t. } g_i(x) \leq 0, \quad i = 1, 2, \dots, m \end{aligned}$$



for twice continuously differentiable objective and constraint functions with an error that is less than or equal to  $\epsilon$ , would in light of Theorem 17.1.2 be to solve the optimization problem  $\min F_t(x)$  with logarithmic barrier  $F$  for  $t = m/\epsilon$ , using for example Newton's method. The strategy works for small problems and with moderate demands on accuracy, but better results are obtained by solving the problems  $\min F_t(x)$  for an increasing sequence of  $t$ -values until  $t \geq m/\epsilon$ .

A simple version of the barrier method or the *path-following method*, as it is also called, therefore looks like this:

### Path-following method

**Given** a starting point  $x = x_0 \in \text{int } X$ , a real number  $t = t_0 > 0$ , an update parameter  $\alpha > 1$  and a tolerance  $\epsilon > 0$ .

#### Repeat

1. Compute  $\hat{x}(t)$  by minimizing  $F_t = tf + F$  with  $x$  as starting point
2. *Update*:  $x := \hat{x}(t)$ .
3. *Stopping criterion*: **stop** if  $m/t \leq \epsilon$ .
4. *Increase  $t$* :  $t := \alpha t$ .

Step 1 is called an *outer iteration* or a *centering step* because it is about finding a point on the central path. To minimize the function  $F_t$ , Newton's method is used, and the iterations of Newton's method to compute  $\hat{x}(t)$  with  $x$  as the starting point are called *inner iterations*.

It is not necessary to compute  $\hat{x}(t)$  exactly in the outer iterations; the central path serves no other function than to lead to the optimal point  $\hat{x}$ , and good approximations to points on the central path will also give rise to a sequence of points which converges to  $\hat{x}$ .

The computational cost of the method obviously depends on the total number of outer iterations that have to be performed before the stopping criterion is met, and on the number of inner iterations in each outer iteration.

### The update parameter $\alpha$

The parameter  $\alpha$  (and the initial value  $t_0$ ) determines the number of outer iterations required to reach the stopping criterion  $t \geq m/\epsilon$ . If  $\alpha$  is small, i.e. close to 1, then many outer iterations are needed, but on the other hand, each outer iteration requires few inner iterations since the minimum point  $x = \hat{x}(t)$  of the function  $F_t$  is then a very good starting point in Newton's algorithm for the problem of minimizing the function  $F_{\alpha t}$ .

For large  $\alpha$  values the opposite is true; few outer iterations are needed, but each outer iteration now requires more Newton steps as the starting point  $\hat{x}(t)$  is farther from the minimum point  $\hat{x}(\alpha t)$ .

From experience, it turns out, however, that the above two effects tend to offset each other. The total number of Newton steps is roughly constant over a wide range of  $\alpha$ , and values of  $\alpha$  between 10 and 20 usually work well.

### The initial value $t_0$

The choice of the starting value  $t_0$  is also significant. A small value requires many outer iterations before the stopping criterion is met. A large value, on the other hand, requires many inner iterations in the first outer iteration before a sufficiently good approximation to the point  $\hat{x}(t_0)$  on the central path has been found. Since  $f(\hat{x}(t_0)) - f(\hat{x}) \approx m/t_0$ , it may be reasonable to choose  $t_0$  so that  $m/t_0$  is of the same magnitude as  $f(x_0) - f(\hat{x})$ . The problem, of course, is that the optimal value  $f(\hat{x})$  is not known a priori, so it is necessary to use a suitable estimate. If, for example, a feasible point  $\lambda$  for the dual problem is known and  $\phi$  is the dual function, then  $\phi(\lambda)$  can be used as an approximation of  $f(\hat{x})$ , and  $t_0 = m/(f(x_0) - \phi(\lambda))$  can be taken as initial  $t$ -value.

### The starting point $x_0$

The starting point  $x_0$  must lie in the interior of  $X$ , i.e. it has to satisfy all constraints with strict inequality. If such a point is not known in advance, then one can use the barrier method on an artificial problem to compute such



"I studied English for 16 years but...  
...I finally learned to speak it in just six lessons"

Jane, Chinese architect

ENGLISH OUT THERE

Click to hear me talking before and after my unique course download

a point, or to conclude that the original problem has no feasible points. The procedure is called *phase 1* of the path-following method and works as follows.

Consider the inequalities

$$(17.4) \quad g_i(x) \leq 0, \quad i = 1, 2, \dots, m$$

and suppose that the functions  $g_i: \Omega \rightarrow \mathbf{R}$  are convex and twice continuously differentiable. To determine a point that satisfies all inequalities strictly or to determine that there is no such point, we form the optimization problem

$$(17.5) \quad \begin{array}{ll} \min & s \\ \text{s.t.} & g_i(x) \leq s, \quad i = 1, 2, \dots, m \end{array}$$

in the variables  $x$  and  $s$ . This problem has strictly feasible points, because we can first choose  $x_0 \in \Omega$  arbitrarily and then choose  $s_0 > \max_i g_i(x_0)$ , and we obtain in this way a point  $(x_0, s_0) \in \Omega \times \mathbf{R}$  that satisfies the constraints with strict inequalities. The functions  $(x, s) \mapsto g_i(x) - s$  are obviously convex. We can therefore use the path-following method on the problem (17.5), and depending on the sign of the problem's optimal value  $v_{\min}$ , we get three cases.

$v_{\min} < 0$ : The system (17.4) has strictly feasible solutions. Indeed, if  $(x, s)$  is a feasible point for the problem (17.5) with  $s < 0$ , then  $g_i(x) < 0$  for all  $i$ . This means that it is not necessary to solve the optimization problem (17.5) with great accuracy. The algorithm can be stopped as soon as it has generated a point  $(x, s)$  with  $s < 0$ .

$v_{\min} > 0$ : The system (17.4) is infeasible. Also in this case, it is not necessary to solve the problem with great accuracy. We can stop as soon as we have found a feasible point for the dual problem with a positive value of the dual function, since this implies that  $v_{\min} > 0$ .

$v_{\min} = 0$ : If the greatest lower bound  $v_{\min} = 0$  is attained, i.e. if there is a point  $(\hat{x}, \hat{s})$  with  $\hat{s} = 0$ , then the system (17.4) is feasible but not strictly feasible. The system (17.4) is infeasible if  $v_{\min}$  is not attained. In practice, it is of course impossible to determine exactly that  $v_{\min} = 0$ ; the algorithm terminates with the conclusion that  $|v_{\min}| < \epsilon$  for some small positive number  $\epsilon$ , and we can only be sure that the system  $g_i(x) < -\epsilon$  is infeasible and that the system  $g_i(x) \leq \epsilon$  is feasible.

## Convergence analysis

At the beginning of outer iteration number  $k$ , we have  $t = \alpha^{k-1}t_0$ . The stopping criterion will be triggered as soon as  $m/(\alpha^{k-1}t_0) \leq \epsilon$ , i.e. when

$k - 1 \geq (\log(m/(\epsilon t_0)))/\log \alpha$ . The number of outer iterations is thus equal to

$$\left\lceil \frac{\log(m/(\epsilon t_0))}{\log \alpha} \right\rceil + 1$$

(for  $\epsilon \leq m/t_0$ ).

The path-following method therefore works, provided that the minimization problems

$$(17.6) \quad \begin{array}{ll} \min & tf(x) + F(x) \\ \text{s.t.} & x \in \text{int } X \end{array}$$

can be solved for  $t \geq t_0$ . Using Newton's method, this is true, for example, if the objective functions satisfy the conditions of Theorem 15.2.4, i.e. if  $F_t$  is strongly convex, has a Lipschitz continuous derivative and the sublevel set corresponding to the starting point is closed.

A question that remains to be resolved is whether the problem (17.6) gets harder and harder, that is requires more inner iterations, when  $t$  grows. Practical experience shows that this is not so – in most problems, the number of Newton steps seems to be roughly constant when  $t$  grows. For problems with self-concordant objective and barrier functions, it is possible to obtain exact estimates of the total number of iterations needed to solve the optimization problem (P) with a given accuracy, and this will be the theme in Chapter 18.

Excellent Economics and Business programmes at:



**university of  
 groningen**



**“The perfect start  
 of a successful,  
 international career.”**

**CLICK HERE**  
 to discover why both socially  
 and academically the University  
 of Groningen is one of the best  
 places for a student to be

[www.rug.nl/feb/education](http://www.rug.nl/feb/education)



# Chapter 18

## The path-following method with self-concordant barrier

### 18.1 Self-concordant barriers

**Definition.** Let  $X$  be a closed convex subset of  $\mathbf{R}^n$  with nonempty interior  $\text{int } X$ , and let  $\nu$  be a nonnegative number. A function  $f: \text{int } X \rightarrow \mathbf{R}$  is called a *self-concordant barrier to  $X$  with parameter  $\nu$* , or shorter a  *$\nu$ -self-concordant barrier*, if the function is closed, self-concordant and non-constant, and the Newton decrement satisfies the inequality

$$(18.1) \quad \lambda(f, x) \leq \nu^{1/2}$$

for all  $x \in \text{int } X$ .

It follows from Theorem 15.1.2 and Theorem 15.1.3 that inequality (18.1) holds if and only if

$$|\langle f'(x), v \rangle| \leq \nu^{1/2} \|v\|_x$$

for all vectors  $v \in \mathbf{R}^n$ , or equivalently, if and only if

$$(Df(x)[v])^2 \leq \nu D^2f(x)[v, v]$$

for all  $v \in \mathbf{R}^n$ .

A closed self-concordant function  $f: \Omega \rightarrow \mathbf{R}$  with the property that  $\sup_{x \in \Omega} \lambda(f, x) < 1$  is necessarily constant and the domain  $\Omega$  is equal to  $\mathbf{R}^n$ , according to Theorem 16.4.7. The parameter  $\nu$  of a self-concordant barrier must thus be greater than or equal to 1.

**EXAMPLE 18.1.1.** The function  $f(x) = -\ln x$  is a 1-self-concordant barrier to the interval  $[0, \infty[$ , because  $f$  is closed and self-concordant and  $\lambda(f, x) = 1$  for all  $x > 0$ .  $\square$

EXAMPLE 18.1.2. Convex quadratic functions

$$f(x) = \frac{1}{2}\langle x, Ax \rangle + \langle b, x \rangle + c$$

are self-concordant on  $\mathbf{R}^n$ , but they do not function as self-concordant barriers, because  $\sup \lambda(f, x) = \infty$  for all non-constant convex quadratic functions  $f$ , according to Example 15.1.2.  $\square$

We will show later that only subsets of halfspaces can have self-concordant barriers, so there is no self-concordant barrier to the whole  $\mathbf{R}^n$ .

EXAMPLE 18.1.3. Let  $g(x)$  be a non-constant convex, quadratic function. The function  $f$ , defined by

$$f(x) = -\ln(-g(x)),$$

is a 1-self-concordant barrier to the set  $X = \{x \in \mathbf{R}^n \mid g(x) \leq 0\}$ .

*Proof.* Let  $g(x) = \frac{1}{2}\langle x, Ax \rangle + \langle b, x \rangle + c$ , let  $v$  be an arbitrary vector in  $\mathbf{R}^n$ , and set

$$\alpha = -\frac{1}{g(x)}Dg(x)[v] \quad \text{and} \quad \beta = -\frac{1}{g(x)}D^2g(x)[v, v] = -\frac{1}{g(x)}\langle v, Av \rangle,$$

where  $x$  is an arbitrary point in the interior of  $X$ . Note that  $\beta \geq 0$  and that  $D^3g(x)[v, v, v] = 0$ . It therefore follows from the differentiation rules that

$$\begin{aligned} Df(x)[v] &= -\frac{1}{g(x)}Dg(x)[v] = \alpha, \\ D^2f(x)[v, v] &= \frac{1}{g(x)^2}(Dg(x)[v])^2 - \frac{1}{g(x)}D^2g(x)[v, v] = \alpha^2 + \beta \geq 0, \\ D^3f(x)[v, v, v] &= -\frac{2}{g(x)^3}(Dg(x)[v])^3 + \frac{3}{g(x)^2}D^2g(x)[v, v]Dg(x)[v] \\ &\quad - \frac{1}{g(x)}D^3g(x)[v, v, v] = 2\alpha^3 + 3\alpha\beta. \end{aligned}$$

The function  $f$  is convex since its second derivative is positive semidefinite, and it is closed since  $f(x) \rightarrow +\infty$  as  $g(x) \rightarrow 0$ . By squaring it is easy to show that the inequality  $|2\alpha^3 + 3\alpha\beta| \leq 2(\alpha^2 + \beta)^{3/2}$  holds for all  $\alpha \in \mathbf{R}$  and all  $\beta \in \mathbf{R}_+$ , and obviously  $\alpha^2 \leq \alpha^2 + \beta$ . This means that  $|D^3f(x)[v, v, v]| \leq 2(D^2f(x)[v, v])^{3/2}$  and that  $(Df(x)[v])^2 \leq D^2f(x)[v, v]$ . So  $f$  is 1-self-concordant.  $\square$

The following three theorems show how to build new self-concordant barriers from given ones.


**Theorem 18.1.1.** *If  $f$  is a  $\nu$ -self-concordant barrier to the set  $X$  and  $\alpha \geq 1$ , then  $\alpha f$  is an  $\alpha\nu$ -self-concordant barrier to  $X$ .*

*Proof.* The proof is left as a simple exercise. □

**Theorem 18.1.2.** *If  $f$  is a  $\mu$ -self-concordant barrier to the set  $X$  and  $g$  is a  $\nu$ -self-concordant barrier to the set  $Y$ , then the sum  $f + g$  is a self-concordant barrier with parameter  $\mu + \nu$  to the intersection  $X \cap Y$ . And  $f + c$  is a  $\mu$ -self-concordant barrier to  $X$  for each constant  $c$ .*

*Proof.* The sum  $f + g$  is a closed convex function, and it is self-concordant on the set  $\text{int}(X \cap Y)$  according to Theorem 16.1.5. To prove that the sum is a self-concordant barrier with parameter  $(\mu + \nu)$ , we assume that  $v$  is an arbitrary vector in  $\mathbf{R}^n$  and write  $a = D^2f(x)[v, v]$  and  $b = D^2g(x)[v, v]$ . We then have, by definition,

$$(Df(x)[v])^2 \leq \mu a \quad \text{and} \quad (Dg(x)[v])^2 \leq \nu b,$$



In the past four years we have drilled  
**89,000 km**  
That's more than **twice** around the world.

**Who are we?**  
We are the world's largest oilfield services company<sup>1</sup>. Working globally—often in remote and challenging locations—we invent, design, engineer, and apply technology to help our customers find and produce oil and gas safely.

**Who are we looking for?**  
Every year, we need thousands of graduates to begin dynamic careers in the following domains:

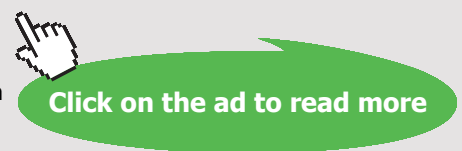
- **Engineering, Research and Operations**
- **Geoscience and Petrotechnical**
- **Commercial and Business**

**What will you be?**

**Schlumberger**

[careers.slb.com](http://careers.slb.com)

<sup>1</sup>Based on Fortune 500 ranking 2011. Copyright © 2015 Schlumberger. All rights reserved.



and using the inequality  $2\sqrt{\mu\nu ab} \leq \nu a + \mu b$  between the geometric and the arithmetic mean, we obtain the inequality

$$\begin{aligned} (D(f+g)(x)[v])^2 &= (Df(x)[v])^2 + (Dg(x)[v])^2 + 2Df(x)[v] \cdot Dg(x)[v] \\ &\leq \mu a + \nu b + 2\sqrt{\mu a \nu b} \leq \mu a + \nu b + \nu a + \mu b \\ &= (\mu + \nu)(a + b) = (\mu + \nu) D^2(f+g)(x)[v, v], \end{aligned}$$

which means that  $\lambda(f+g, x) \leq (\mu + \nu)^{1/2}$ .

The assertion about the sum  $f + c$  is trivial, since  $\lambda(f, x) = \lambda(f + c, x)$  for constants  $c$ .  $\square$

**Theorem 18.1.3.** *Suppose that  $A: \mathbf{R}^m \rightarrow \mathbf{R}^n$  is an affine map and that  $f$  is a  $\nu$ -self-concordant barrier to the subset  $X$  of  $\mathbf{R}^n$ . The composition  $g = f \circ A$  is then a  $\nu$ -self-concordant barrier to the inverse image  $A^{-1}(X)$ .*

*Proof.* The proof is left as an exercise.  $\square$

EXAMPLE 18.1.4. It follows from Example 18.1.1 and Theorems 18.1.2 and 18.1.3 that the function

$$f(x) = - \sum_{i=1}^m \ln(b_i - \langle a_i, x \rangle)$$

is an  $m$ -self-concordant barrier to the polyhedron

$$X = \{x \in \mathbf{R}^n \mid \langle a_i, x \rangle \leq b_i, \quad i = 1, 2, \dots, m\}. \quad \square$$

**Theorem 18.1.4.** *If  $f$  is a  $\nu$ -self-concordant barrier to the set  $X$ , then*

$$\langle f'(x), y - x \rangle \leq \nu$$

for all  $x \in \text{int } X$  and all  $y \in X$ .

*Remark.* It follows that a set with a self-concordant barrier must be a subset of some halfspace. Indeed, a set  $X$  with a  $\nu$ -self-concordant barrier is a subset of the closed halfspace  $\{y \in \mathbf{R}^n \mid \langle c, y \rangle \leq \nu + \langle c, x_0 \rangle\}$ , where  $x_0 \in \text{int } X$  is an arbitrary point with  $c = f'(x_0) \neq 0$ .

*Proof.* Fix  $x \in \text{int } X$  and  $y \in X$ , let  $x^t = x + t(y - x)$  and define the function  $\phi$  by setting  $\phi(t) = f(x^t)$ . Then  $\phi$  is certainly defined on the open interval  $]\alpha, 1[$  for some negative number  $\alpha$ , since  $x$  is an interior point. Moreover,

$$\phi'(t) = Df(x^t)[y - x],$$

and especially,  $\phi'(0) = Df(x)[y - x] = \langle f'(x), y - x \rangle$ . We will prove that  $\phi'(0) \leq \nu$ .



If  $\phi'(0) \leq 0$ , then we are done, so assume that  $\phi'(0) > 0$ . By  $\nu$ -self-concordance,

$$\phi''(t) = D^2 f(x^t)[y - x, y - x] \geq \nu^{-1} (Df(x^t)[y - x])^2 = \nu^{-1} \phi'(t)^2 \geq 0.$$

The derivative  $\phi'$  is thus increasing, and this implies that  $\phi'(t) \geq \phi'(0) > 0$  for  $t \geq 0$ . Furthermore,

$$\frac{d}{dt} \left( -\frac{1}{\phi'(t)} \right) = \frac{\phi''(t)}{\phi'(t)^2} \geq \frac{1}{\nu}$$

for all  $t$  in the interval  $[0, 1[$ , so by integrating the last mentioned inequality over the interval  $[0, \beta]$ , where  $\beta < 1$ , we obtain the inequality

$$\frac{1}{\phi'(0)} > \frac{1}{\phi'(0)} - \frac{1}{\phi'(\beta)} = \int_0^\beta \frac{d}{dt} \left( -\frac{1}{\phi'(t)} \right) dt \geq \frac{\beta}{\nu}.$$

Hence,  $\phi'(0) < \nu/\beta$  for all  $\beta < 1$ , which implies that  $\phi'(0) \leq \nu$ . □

**Theorem 18.1.5.** *Suppose that  $f$  is a  $\nu$ -self-concordant barrier to the set  $X$ . If  $x \in \text{int } X$ ,  $y \in X$  and  $\langle f'(x), y - x \rangle \geq 0$ , then*

$$\|y - x\|_x \leq \nu + 2\sqrt{\nu}.$$

*Remark.* If  $x \in \text{int } X$  is a minimum point, then  $\langle f'(x), y - x \rangle = 0$  for all points  $y \in X$ , since  $f'(x) = 0$ . Hence,  $\|y - x\|_x \leq \nu + 2\sqrt{\nu}$  for all  $y \in X$  if  $x$  is a minimum point.

*Proof.* Let  $r = \|y - x\|_x$ . If  $r \leq \sqrt{\nu}$ , then there is nothing to prove, so assume that  $r > \sqrt{\nu}$ , and consider for  $\alpha = \sqrt{\nu}/r$  the point  $z = x + \alpha(y - x)$ , which lies in the interior of  $X$  since  $\alpha < 1$ . By using Theorem 18.1.4 with  $z$  instead of  $x$ , the assumption  $\langle f'(x), y - x \rangle \geq 0$ , Theorem 16.3.2 and the equalities  $y - z = (1 - \alpha)(y - x)$  and  $z - x = \alpha(y - x)$ , we obtain the following chain of inequalities and equalities:

$$\begin{aligned} \nu &\geq \langle f'(z), y - z \rangle = (1 - \alpha) \langle f'(z), y - x \rangle \geq (1 - \alpha) \langle f'(z) - f'(x), y - x \rangle \\ &= \frac{1 - \alpha}{\alpha} \langle f'(z) - f'(x), z - x \rangle \geq \frac{1 - \alpha}{\alpha} \cdot \frac{\|z - x\|_x^2}{1 + \|z - x\|_x} \\ &= \frac{(1 - \alpha)\alpha \|y - x\|_x^2}{1 + \alpha \|y - x\|_x} = \frac{r\sqrt{\nu} - \nu}{1 + \sqrt{\nu}}. \end{aligned}$$

The inequality between the extreme ends simplifies to  $r \leq \nu + 2\sqrt{\nu}$ , which is the desired inequality. □

Given a self-concordant function  $f$  with the corresponding local seminorm  $\|\cdot\|_x$ , we set

$$\mathcal{E}(x; r) = \{y \in \mathbf{R}^n \mid \|y - x\|_x \leq r\}.$$

If  $f$  is non-degenerate, then  $\|\cdot\|_x$  is a norm at each point  $x \in \text{int } X$ , and the set  $\mathcal{E}(x; r)$  is a closed ellipsoid in  $\mathbf{R}^n$  with axis directions determined by the eigenvectors of the second derivative  $f''(x)$ .

For non-degenerate self-concordant barriers we now have the following corollary to Theorem 18.1.5.

**Theorem 18.1.6.** *Suppose that  $f$  is a non-degenerate  $\nu$ -self-concordant barrier to the closed convex set  $X$ . Then  $f$  attains a minimum if and only if  $X$  is a bounded set. The minimum point  $\hat{x}_f \in \text{int } X$  is unique in that case, and*

$$\mathcal{E}(\hat{x}_f; 1) \subseteq X \subseteq \mathcal{E}(\hat{x}_f; \nu + 2\sqrt{\nu}).$$

*Remark.* A closed self-concordant function whose domain does not contain any line, is automatically non-degenerate, so it is not necessary to state explicitly that a self-concordant barrier to a compact set should be non-degenerate.

## American online LIGS University

is currently enrolling in the  
Interactive Online **BBA, MBA, MSc,**  
**DBA and PhD** programs:

- ▶ enroll **by September 30th, 2014** and
- ▶ **save up to 16%** on the tuition!
- ▶ pay in 10 installments / 2 years
- ▶ Interactive Online education
- ▶ visit [www.ligsuniversity.com](http://www.ligsuniversity.com) to  
find out more!

**Note:** LIGS University is not accredited by any nationally recognized accrediting agency listed by the US Secretary of Education. More info [here](#).



*Proof.* The sublevel sets of a closed convex function are closed, so if  $X$  is a bounded set, then each sublevel set  $\{x \in \text{int } X \mid f(x) \leq \alpha\}$  is both closed and bounded, and this implies that  $f$  has a minimum, and the minimum point of a non-degenerate convex function is necessarily unique.

Conversely, assume that  $f$  has a minimum point  $\hat{x}_f$ . Then by the remark following Theorem 18.1.5,  $\|y - \hat{x}_f\|_{\hat{x}_f} \leq \nu + 2\sqrt{\nu}$  for all  $y \in X$ , and this amounts to the right inclusion in Theorem 18.1.6, which implies, of course, that  $X$  is a bounded set.

The remaining left inclusion follows from Theorem 16.3.2, which implies that the open ellipsoid  $\{y \in \mathbf{R}^n \mid \|y - x\|_x < 1\}$  is a subset of  $\text{int } X$  for each choice of  $x \in \text{int } X$ . The closure  $\mathcal{E}(x; 1)$  is therefore a subset of  $X$ , and we obtain the left inclusion by choosing  $x = \hat{x}_f$ .  $\square$

Given a self-concordant barrier to a set  $X$  we will need to compare the local seminorms  $\|v\|_x$  and  $\|v\|_y$  of a vector at different points  $x$  and  $y$ , and in order to achieve this we need a measure for the distance from  $y$  to  $x$  relative the distance from  $y$  to the boundary of  $X$  along the half-line from  $x$  through  $x$ . The following definition provides us with the relevant measure.

**Definition.** Let  $X$  be a closed convex subset of  $\mathbf{R}^n$  with nonempty interior. For each  $y \in \text{int } X$  we define a function  $\pi_y: \mathbf{R}^n \rightarrow \mathbf{R}_+$  by setting

$$\pi_y(x) = \inf\{t > 0 \mid y + t^{-1}(x - y) \in X\}.$$

Obviously,  $\pi_y(y) = 0$ . To determine  $\pi_y(x)$  if  $x \neq y$ , we consider the half-line from  $y$  through  $x$ ; if the half-line intersects the boundary of  $X$  in a point  $z$ , then  $\pi_y(x) = \|x - y\|/\|z - y\|$  (with respect to arbitrary norms), and if the entire half-line lies in  $X$ , then  $\pi_y(x) = 0$ . We note that  $\pi_y(x) < 1$  for interior points  $x$ , that  $\pi_y(x) = 1$  for boundary points  $x$ , and that  $\pi_y(x) > 1$  for points outside  $X$ .

We could also have defined the function  $\pi_y$  in terms of the Minkowski functional that was introduced in Section 6.10 of Part I, because

$$\pi_y(x) = \phi_{-y+X}(x - y),$$

where  $\phi_{-y+X}$  is the Minkowski functional of the set  $-y + X$ .

The following simple estimate of  $\pi_y(x)$  will be needed later on.

**Theorem 18.1.7.** *Let  $X$  be a compact convex set, let  $x$  and  $y$  be points in the interior of  $X$ , and suppose that*

$$B(x, r) \subseteq X \subseteq \overline{B}(0; R),$$

where the balls are given with respect to an arbitrary norm  $\|\cdot\|$ . Then

$$\pi_y(x) \leq \frac{2R}{2R + r}.$$

*Proof.* The inequality is trivially true if  $x = y$ , so suppose that  $x \neq y$ . The half-line from  $y$  through  $x$  intersects the boundary of  $X$  in a point  $z$  and  $\|z - y\| = \|z - x\| + \|x - y\|$ . Furthermore,  $\|z - x\| \geq r$  and  $\|x - y\| \leq 2R$ , so it follows that

$$\pi_y(x) = \frac{\|x - y\|}{\|z - y\|} = \left(1 + \frac{\|z - x\|}{\|x - y\|}\right)^{-1} \leq \left(1 + \frac{r}{2R}\right)^{-1} = \frac{2R}{2R + r}. \quad \square$$

The direction derivative  $\langle f'(x), v \rangle$  of a  $\nu$ -self-concordant barrier function  $f$  is bounded by  $\sqrt{\nu}\|v\|_x$ , by definition. Our next theorem shows that the same direction derivative is also bounded by a constant times  $\|v\|_y$ , if  $y$  is an arbitrary point in the domain of  $f$ . The two local norms  $\|v\|_x$  and  $\|v\|_y$  are also compared.

**Theorem 18.1.8.** *Let  $f$  be a  $\nu$ -self-concordant barrier to  $X$ , and let  $x$  and  $y$  be two points in the interior of  $X$ . Then, for all vectors  $v$*

$$(18.2) \quad |\langle f'(x), v \rangle| \leq \frac{\nu}{1 - \pi_y(x)} \|v\|_y$$

and

$$(18.3) \quad \|v\|_x \leq \frac{\nu + 2\sqrt{\nu}}{1 - \pi_y(x)} \|v\|_y.$$

*Proof.* The two inequalities hold if  $y = x$  since

$$|\langle f'(x), v \rangle| \leq \sqrt{\nu}\|v\|_x \leq \nu\|v\|_x$$

and  $\pi_x(x) = 0$ . They also hold if  $\|v\|_y = 0$ , i.e. if the vector  $v$  belongs to the recessive subspace of  $f$ , because then  $\|v\|_x = 0$  and  $\langle f'(x), v \rangle = 0$ . Assume henceforth that  $y \neq x$  and that  $\|v\|_y \neq 0$ .

First consider the case  $\|v\|_y = 1$ , and let  $s$  be an arbitrary number greater than  $\nu + 2\sqrt{\nu}$ . Then, by Theorems 16.3.2 and 18.1.5, we conclude that

- (i) The two points  $y \pm v$  lie in  $X$ .
- (ii) At least one of the two points  $x \pm \frac{s}{\|v\|_x}v$  lies outside  $X$ .

By the definition of  $\pi_y(x)$  there is a vector  $z \in X$  such that

$$x = y + \pi_y(x)(z - y),$$

and since

$$x \pm (1 - \pi_y(x))v = \pi_y(x)z + (1 - \pi_y(x))(y \pm v),$$

it follows from convexity that

- (iii) The two points  $x \pm (1 - \pi_y(x))v$  lie in  $X$ .

It now follows from (iii) and Theorem 18.1.4 that

$$\langle f'(x), \pm v \rangle = \frac{1}{1 - \pi_y(x)} \langle f'(x), x \pm (1 - \pi_y(x))v - x \rangle \leq \frac{\nu}{1 - \pi_y(x)},$$

which means that

$$|\langle f'(x), v \rangle| \leq \frac{\nu}{1 - \pi_y(x)}.$$

This proves inequality (18.2) for vectors  $v$  with  $\|v\|_y = 1$ , and if  $v$  is an arbitrary vector with  $\|v\|_y \neq 0$ , we obtain inequality (18.2) by replacing  $v$  in the inequality above with  $v/\|v\|_y$ .

By combining the two assertions (ii) and (iii) we conclude that

$$1 - \pi_y(x) < \frac{s}{\|v\|_x},$$

i.e. that

$$\|v\|_x < \frac{s}{1 - \pi_y(x)} = \frac{s}{1 - \pi_y(x)} \|v\|_y,$$

and since this holds for all  $s > \nu + 2\sqrt{\nu}$ , it follows that

$$\|v\|_x \leq \frac{\nu + 2\sqrt{\nu}}{1 - \pi_y(x)} \|v\|_y.$$

.....Alcatel-Lucent 

[www.alcatel-lucent.com/careers](http://www.alcatel-lucent.com/careers)



What if you could build your future and create the future?

One generation's transformation is the next's status quo. In the near future, people may soon think it's strange that devices ever had to be "plugged in." To obtain that status, there needs to be "The Shift".

This proves inequality (18.3) in the case  $\|v\|_y = 1$ , and since the inequality is homogeneous, it holds in general.  $\square$

**Definition.** Let  $\|\cdot\|_x$  be the local seminorm at  $x$  which is associated with the two times differentiable convex function  $f: X \rightarrow \mathbf{R}$ , where  $X$  is a subset of  $\mathbf{R}^n$ . The corresponding *dual local norm* is the function  $\|\cdot\|_x^*: \mathbf{R}^n \rightarrow \overline{\mathbf{R}}$ , which is defined by

$$\|v\|_x^* = \sup_{\|w\|_x \leq 1} \langle v, w \rangle$$

for all  $v \in \mathbf{R}^n$ .

The dual norm is easily verified to be subadditive and homogeneous, i.e.  $\|v + w\|_x^* \leq \|v\|_x^* + \|w\|_x^*$  and  $\|\lambda v\|_x^* = |\lambda| \|v\|_x^*$  for all  $v, w \in \mathbf{R}^n$  and all real numbers  $\lambda$ , but  $\|\cdot\|_x^*$  is a proper norm on the whole of  $\mathbf{R}^n$  only for points  $x$  where the second derivative  $f''(x)$  is positive definite, because  $\|v\|_x^* = \infty$  if  $v$  is a nonzero vector in the null space  $\mathcal{N}(f''(x))$  since  $\|tv\|_x = 0$  for all  $t \in \mathbf{R}$  and  $\langle v, tv \rangle = t\|v\|^2 \rightarrow \infty$  as  $t \rightarrow \infty$ . However,  $\|\cdot\|_x^*$  is always a proper norm when restricted to the subspace  $\mathcal{N}(f''(x))^\perp$ . See exercise 18.2.

By Theorem 15.1.3, we have the following expression for the Newton decrement  $\lambda(f, x)$  in terms of the dual local norm:

$$\lambda(f, x) = \|f'(x)\|_x^*.$$

The following variant of the Cauchy–Schwarz inequality holds for the local seminorm.

**Theorem 18.1.9.** *Assume that  $\|v\|_x^* < \infty$ . Then*

$$|\langle v, w \rangle| \leq \|v\|_x^* \|w\|_x$$

for all vectors  $w$ .

*Proof.* If  $\|w\|_x \neq 0$  then  $\pm w/\|w\|_x$  are two vectors with local seminorm equal to 1, so it follows from the definition of the dual norm that

$$\pm \frac{1}{\|w\|_x} \langle v, w \rangle = \langle v, \pm w/\|w\|_x \rangle \leq \|v\|_x^*,$$

and we obtain the sought inequality after multiplication by  $\|w\|_x$ .

If instead  $\|w\|_x = 0$ , then  $\|tw\|_x = 0$  for all real numbers  $t$ , and it follows from the supremum definition that  $t\langle v, w \rangle = \langle v, tw \rangle \leq \|v\|_x^* < \infty$  for all  $t$ . This being possible only if  $\langle v, w \rangle = 0$ , we conclude that the inequality applies in this case, too.  $\square$

Later we will need various estimates of  $\|v\|_x^*$ . Our first estimate is in terms of the width in different directions of the set  $X$ , and this motivates our next definition.

**Definition.** Given a nonempty subset  $X$  of  $\mathbf{R}^n$ , let  $\text{Var}_X: \mathbf{R}^n \rightarrow \overline{\mathbf{R}}$  be the function defined by

$$\text{Var}_X(v) = \sup_{x \in X} \langle v, x \rangle - \inf_{x \in X} \langle v, x \rangle.$$

$\text{Var}_X(v)$  is obviously a finite number for each  $v \in \mathbf{R}^n$  if the set  $X$  is bounded, and if  $v$  is a unit vector, then  $\text{Var}_X(v)$  measures the width of the set  $X$  in the direction of  $v$ .

Our next theorem shows how to estimate  $\|\cdot\|_x^*$  using  $\text{Var}_X$ .

**Theorem 18.1.10.** *Suppose that  $f: X \rightarrow \mathbf{R}$  is a closed self-concordant function with a bounded open convex subset  $X$  of  $\mathbf{R}^n$  as domain, and let  $\|\cdot\|_x^*$  be the dual local norm associated with the function  $f$  at the point  $x \in X$ . Then*

$$\|v\|_x^* \leq \text{Var}_X(v)$$

for all  $v \in \mathbf{R}^n$ .

*Proof.* It follows from the previous theorem that  $y$  is a point in  $\text{cl } X$  if  $x$  is a point in  $X$  and  $\|y - x\|_x \leq 1$ . Hence,

$$\begin{aligned} \|v\|_x^* &= \sup_{\|w\|_x \leq 1} \langle v, w \rangle = \sup_{\|y-x\|_x \leq 1} \langle v, y-x \rangle \leq \sup_{y \in \text{cl } X} \langle v, y-x \rangle = \sup_{y \in X} \langle v, y-x \rangle \\ &= \sup_{y \in X} \langle v, y \rangle - \langle v, x \rangle \leq \sup_{y \in X} \langle v, y \rangle - \inf_{y \in X} \langle v, y \rangle = \text{Var}_X(v). \quad \square \end{aligned}$$

We have previously defined the analytic center of a closed convex set  $X$  with respect to a given barrier as the unique minimum point of the barrier, provided that there is one. According to Theorem 18.1.6, every compact convex set with nonempty interior has an analytic center with respect to any given  $\nu$ -self-concordant barrier. We can now obtain an upper bound on the dual local norm  $\|v\|_x^*$  at an arbitrary point  $x$  in terms of the parameter  $\nu$  and the value of the dual norm at the analytic center.

**Theorem 18.1.11.** *Let  $X$  be a compact convex set, and let  $\hat{x}_f$  be the analytic center of the set with respect to a  $\nu$ -self-concordant barrier  $f$ . Then, for each vector  $v \in \mathbf{R}^n$  and each  $x \in \text{int } X$ ,*

$$\|v\|_x^* \leq (\nu + 2\sqrt{\nu}) \|v\|_{\hat{x}_f}^*.$$

*Proof.* Let  $B_1 = \mathcal{E}(x; 1)$  and  $B_2 = \mathcal{E}(\hat{x}_f; \nu + 2\sqrt{\nu})$ . Theorems 16.3.2 and 18.1.6 give us the inclusions  $B_1 \subseteq X \subseteq B_2$ , so it follows from the definition of the dual local norm that

$$\begin{aligned} \|v\|_x^* &= \sup_{\|w\|_x \leq 1} \langle v, w \rangle = \sup_{y \in B_1} \langle v, y - x \rangle \leq \sup_{y \in B_2} \langle v, y - x \rangle \\ &= \langle v, \hat{x}_f - x \rangle + \sup_{y \in B_2} \langle v, y - \hat{x}_f \rangle = \langle v, \hat{x}_f - x \rangle + \sup_{\|w\|_{\hat{x}_f} \leq \nu + 2\sqrt{\nu}} \langle v, w \rangle \\ &= \langle v, \hat{x}_f - x \rangle + (\nu + 2\sqrt{\nu}) \|v\|_{\hat{x}_f}^*. \end{aligned}$$

Since  $\|-v\|_x^* = \|v\|_x^*$ , we may now without loss of generality assume that  $\langle v, \hat{x}_f - x \rangle \leq 0$ , and this gives us the required inequality.  $\square$

## 18.2 The path-following method

### Standard form

Let us say that a convex optimization problem is in *standard form* if it is presented in the form

$$\begin{aligned} \min \quad & \langle c, x \rangle \\ \text{s.t.} \quad & x \in X \end{aligned}$$

where  $X$  is a compact convex set with nonempty interior and  $X$  is equipped with a  $\nu$ -self-concordant barrier function  $F$ .

*Remark.* One can show that every compact convex set  $X$  has a barrier function, but for a barrier function to be useful in a practical optimization problem, it has to be explicitly given so that it is possible to efficiently calculate its partial first and second derivatives.

The assumption that the set  $X$  is bounded is not particularly restrictive for problems with finite optimal values, for we can always modify such problems by adding artificial, very big bounds on the variables.

We also recall that an arbitrary convex problem can be transformed into an equivalent convex problem with a linear objective function by an epigraph formulation. (See Chapter 9.3 of Part II.)

**EXAMPLE 18.2.1.** Each LP problem with finite optimal value can be written in standard form after suitable transformations. By first identifying the affine hull of the polyhedron of feasible points with  $\mathbf{R}^n$  for an appropriate  $n$ , we can without restriction assume that the polyhedron has a nonempty interior, and by adding big bounds on the variables, if necessary, we can also assume



that our polyhedron  $X$  of feasible points is compact. And with  $X$  written in the form

$$(18.4) \quad X = \{x \in \mathbf{R}^n \mid \langle c_i, x \rangle \leq b_i, i = 1, 2, \dots, m\},$$

we get an  $m$ -self-concordant barrier  $F$  to  $X$ , by defining

$$F(x) = - \sum_{i=1}^m \ln(b_i - \langle c_i, x \rangle) \quad \square$$

EXAMPLE 18.2.2. Convex quadratic optimization problems, i.e. problems of the type

$$\begin{aligned} \min \quad & g(x) \\ \text{s.t.} \quad & x \in X \end{aligned}$$

where  $g$  is a convex quadratic function and  $X$  is a bounded polyhedron in  $\mathbf{R}^n$  with nonempty interior, can be transformed, using an epigraph formulation and an artificial bound  $M$  on the new variable  $s$ , to problems of the form

$$\begin{aligned} \min \quad & s \\ \text{s.t.} \quad & (x, s) \in Y \end{aligned}$$



**Maastricht University**

*Leading in Learning!*

**Join the best at  
the Maastricht University  
School of Business and  
Economics!**

**Top master's programmes**

- 33<sup>rd</sup> place Financial Times worldwide ranking: MSc International Business
- 1<sup>st</sup> place: MSc International Business
- 1<sup>st</sup> place: MSc Financial Economics
- 2<sup>nd</sup> place: MSc Management of Learning
- 2<sup>nd</sup> place: MSc Economics
- 2<sup>nd</sup> place: MSc Econometrics and Operations Research
- 2<sup>nd</sup> place: MSc Global Supply Chain Management and Change

Sources: Keuzegids Master ranking 2013; Elsevier 'Beste Studies' ranking 2012; Financial Times Global Masters in Management ranking 2012

**Maastricht University is the best specialist university in the Netherlands (Elsevier)**

**Visit us and find out why we are the best!  
Master's Open Day: 22 February 2014**

[www.mastersopenday.nl](http://www.mastersopenday.nl)



Click on the ad to read more

where  $Y = \{(x, s) \in \mathbf{R}^n \times \mathbf{R} \mid x \in X, g(x) \leq s \leq M\}$  is a compact convex set with nonempty interior. Now assume that the polyhedron  $X$  is given by equation (18.4) as an intersection of closed halfspaces. Then the function

$$F(x, s) = - \sum_{i=1}^m \ln(b_i - \langle c_i, x \rangle) - \ln(s - g(x)) - \ln(M - s)$$

is an  $(m + 2)$ -self-concordant barrier to  $Y$  according to Example 18.1.3.  $\square$

## Central path

We will now study the path-following method for the standard problem

$$\begin{aligned} \text{(SP)} \quad & \min \quad \langle c, x \rangle \\ & \text{s.t.} \quad x \in X \end{aligned}$$

where  $X$  is a compact convex subset of  $\mathbf{R}^n$  with nonempty interior, and  $F$  is a  $\nu$ -self-concordant barrier to  $X$ . The finite optimal value of the problem is denoted by  $v_{\min}$ .

For  $t \geq 0$  we define functions  $F_t: \text{int } X \rightarrow \mathbf{R}$  by

$$F_t(x) = t\langle c, x \rangle + F(x).$$

The functions  $F_t$  are closed and self-concordant, and since the set  $X$  is compact, each function  $F_t$  has a unique minimum point  $\hat{x}(t)$ . The central path  $\{\hat{x}(t) \mid t \geq 0\}$  is in other words well-defined, and its points satisfy the equation

$$(18.5) \quad tc + F'(\hat{x}(t)) = 0,$$

and the starting point  $\hat{x}(0)$  is by definition the analytic center  $\hat{x}_F$  of  $X$  with respect to the given barrier  $F$ .

We will use Newton's method to determine the minimum point  $\hat{x}(t)$ , and for that reason we need to calculate the Newton step and the Newton decrement with respect to the function  $F_t$  at points in the interior of  $X$ .

Since  $F_t''(x) = F''(x)$ , the local norm  $\|v\|_x$  of a vector  $v$  with respect to the function  $F_t$  is the same for all  $t \geq 0$ , namely

$$\|v\|_x = \sqrt{\langle v, F''(x)v \rangle}.$$

In contrast, Newton steps and Newton decrements depend on  $t$ ; the Newton step at the point  $x$  is equal to  $-F''(x)^{-1}F'_t(x)$  for the function  $F_t$ , and the decrement is given by

$$\lambda(F_t, x) = \sqrt{\langle F'_t(x), F''(x)^{-1}F'_t(x) \rangle} = \|F''(x)^{-1}F'_t(x)\|_x.$$

The following theorem is used to formulate the stopping criterion in the path-following method.

**Theorem 18.2.1.** (i) *The points  $\hat{x}(t)$  on the central path of the optimization problem (SP) satisfy the inequality*

$$\langle c, \hat{x}(t) \rangle - v_{\min} \leq \frac{\nu}{t}.$$

(ii) *Moreover, the inequality*

$$\langle c, x \rangle - v_{\min} \leq \frac{\nu + \kappa(1 - \kappa)^{-1}\sqrt{\nu}}{t}.$$

*holds for  $t > 0$  and all point  $x \in \text{int } X$  satisfying the condition*

$$\lambda(F_t, x) \leq \kappa < 1.$$

*Proof.* (i) Because of equation (18.5),  $c = -t^{-1}F'(\hat{x}(t))$ , and it therefore follows from Theorem 18.1.4 that

$$\langle c, \hat{x}(t) \rangle - \langle c, y \rangle = \frac{1}{t} \langle F'(\hat{x}(t)), y - \hat{x}(t) \rangle \leq \frac{\nu}{t}$$

for all  $y \in X$ . We obtain inequality (i) by choosing  $y$  as an optimal solution to the problem (SP).

(ii) Since  $\langle c, x \rangle - v_{\min} = (\langle c, x \rangle - \langle c, \hat{x}(t) \rangle) + (\langle c, \hat{x}(t) \rangle - v_{\min})$ , it suffices, due to the already proven inequality, to show that

$$(18.6) \quad \langle c, x \rangle - \langle c, \hat{x}(t) \rangle \leq \frac{\kappa}{1 - \kappa} \cdot \frac{\sqrt{\nu}}{t}$$

if  $x \in \text{int } X$  and  $\lambda(F_t, x) \leq \kappa < 1$ . But it follows from Theorem 16.4.6 that

$$\|x - \hat{x}(t)\|_{\hat{x}(t)} \leq \frac{\lambda(F_t, x)}{1 - \lambda(F_t, x)} \leq \frac{\kappa}{1 - \kappa},$$

so by using that  $tc = -F'(\hat{x}(t))$  and that  $F$  is  $\nu$ -self-concordant, we get the following chain of equalities and inequalities:

$$\begin{aligned} t(\langle c, x \rangle - \langle c, \hat{x}(t) \rangle) &= -\langle F'(\hat{x}(t)), x - \hat{x}(t) \rangle \leq \|F'(\hat{x}(t))\|_{\hat{x}(t)}^* \|x - \hat{x}(t)\|_{\hat{x}(t)} \\ &= \lambda(F, \hat{x}(t)) \|x - \hat{x}(t)\|_{\hat{x}(t)} \leq \sqrt{\nu} \frac{\kappa}{1 - \kappa}, \end{aligned}$$

which proves inequality (18.6). □

## Algorithm

The path-following algorithm for solving the standard problem

$$(SP) \quad \begin{array}{ll} \min & \langle c, x \rangle \\ \text{s.t.} & x \in X \end{array}$$

works in brief as follows.

We start with a parameter value  $t_0 > 0$  and a point  $x_0 \in \text{int} X$ , which is close enough to the point  $\hat{x}(t_0)$  on the central path. "Close enough" is expressed in terms of the Newton decrement  $\lambda(F_{t_0}, x_0)$ , which must be sufficiently small.

Then we update the parameter  $t$  by defining  $t_1 = \alpha t_0$  for a suitable  $\alpha > 1$  and minimize the function  $F_{t_1}$  using the damped Newton method with  $x_0$  as the starting point. Newton's method is terminated when it has reached a point  $x_1$ , which is sufficiently close to the minimum point  $\hat{x}(t_1)$  of  $F_{t_1}$ .

The procedure is then repeated with  $t_2 = \alpha t_1$  as new parameter and with  $x_1$  as starting point in Newton's method for minimization of the function  $F_{t_2}$ , etc. As a result we obtain a sequence  $t_0, x_0, t_1, x_1, t_2, x_2, \dots$  of parameter values and points, and the procedure is terminated when  $t_k$  has become sufficiently large with  $x_k$  as an approximate optimal point.



**> Apply now**

REDEFINE YOUR FUTURE  
**AXA GLOBAL GRADUATE  
PROGRAM 2015**

redefining / standards 

agence edg - © Photonostop

From this sketchy description of the algorithm it is clear that we need two parameters, one parameter  $\alpha$  to describe the update of  $t$ , and one parameter  $\kappa$  to define the stopping criterion in Newton's method. We shall estimate the total number of inner iterations, and the estimate will be the simplest and most obvious if one writes the update parameter  $\alpha$  in the form  $\alpha = 1 + \gamma/\sqrt{\nu}$ .

The following precise formulation of the path-following algorithm therefore contains the parameters  $\gamma$  and  $\kappa$ . The addition 'phase 2' is due to the need for an additional phase to generate feasible initial values  $x_0$  and  $t_0$ .

### Path-following algorithm, phase 2

**Given** an update parameter  $\gamma > 0$ , a neighborhood parameter  $0 < \kappa < 1$ , a tolerance  $\epsilon > 0$ , a starting point  $x_0 \in \text{int } X$ , and a starting value  $t_0 > 0$  such that  $\lambda(F_{t_0}, x_0) \leq \kappa$ .

1. *Initiate*:  $x := x_0$  and  $t := t_0$ .
2. *Stopping criterion*: **stop** if  $\epsilon t \geq \nu + \kappa(1 - \kappa)^{-1}\sqrt{\nu}$ .
3. *Increase  $t$* :  $t := (1 + \gamma/\sqrt{\nu})t$ .
4. *Update  $x$  by using Newton's damped method on the function  $F_t$  with the current  $x$  as starting point*:
  - (i) Compute the Newton decrement  $\lambda = \lambda(F_t, x)$ .
  - (ii) **quit** Newton's method if  $\lambda \leq \kappa$ , and go to line 2.
  - (iii) Compute the Newtonstep  $\Delta x_{\text{nt}} = -F''(x)^{-1}F'_t(x)$ .
  - (iv) *Uppdate*:  $x := x + (1 + \lambda)^{-1}\Delta x_{\text{nt}}$
  - (v) Go to (i).

We can now show the following convergence result.

**Theorem 18.2.2.** *Suppose that the above path-following algorithm is applied to the standard problem (SP) with a  $\nu$ -self-concordant barrier  $F$ . Then the algorithm stops with a point  $x \in \text{int } X$  which satisfies*

$$\langle c, x \rangle - v_{\min} \leq \epsilon.$$

*For each outer iteration, the number of inner iterations in Newton's algorithm is bounded by a constant  $K$ , and the total number of inner iterations in the path-following algorithm is bounded by*

$$C\sqrt{\nu} \ln\left(\frac{\nu}{t_0\epsilon} + 1\right),$$

*where the constants  $K$  and  $C$  only depend on  $\kappa$  and  $\gamma$ .*

*Proof.* Let us start by examining the inner loop 4 of the algorithm.

Each time the algorithm passes by line 2, it does so with a point  $x$  in  $\text{int } X$ , which belongs to a  $t$ -value with Newton decrement  $\lambda(F_t, x) \leq \kappa$ . In step 4, the function  $F_s$ , where  $s = (1 + \gamma/\sqrt{\nu})t$ , is then minimized

using Newton's damped method with  $y_0 = x$  as the starting point. The points  $y_k$ ,  $k = 1, 2, 3, \dots$ , generated by the method lie in  $\text{int } X$  according to Theorem 16.3.2, and the stopping condition  $\lambda(F_s, y_k) \leq \kappa$  implies, according to Theorem 16.5.1, that the algorithm terminates after at most  $\lfloor (F_s(x) - F_s(\hat{x}(s))) / \rho(-\kappa) \rfloor$  iterations, where  $\rho$  is the function

$$\rho(u) = -u - \ln(1 - u).$$

We shall show that there is a constant  $K$ , which only depends on the parameters  $\kappa$  and  $\gamma$ , so that

$$\left\lfloor \frac{F_s(x) - F_s(\hat{x}(s))}{\rho(-\kappa)} \right\rfloor \leq K,$$

and for that reason we need to estimate the difference  $F_s(x) - F_s(\hat{x}(s))$ , which we do in the next lemma.

**Lemma 18.2.3.** *Suppose that  $\lambda(F_t, x) \leq \kappa < 1$ . Then, for all  $s > 0$*

$$F_s(x) - F_s(\hat{x}(s)) \leq \rho(\kappa) + \frac{\kappa\sqrt{\nu}}{1 - \kappa} \cdot \left| \frac{s}{t} - 1 \right| + \nu \rho(1 - s/t).$$

*Proof of the lemma.* We start by writing

$$(18.7) \quad F_s(x) - F_s(\hat{x}(s)) = (F_s(x) - F_s(\hat{x}(t))) + (F_s(\hat{x}(t)) - F_s(\hat{x}(s))).$$

By using the equality  $tc = -F'(\hat{x}(t))$  and the inequality

$$|\langle F'(\hat{x}(t)), v \rangle| \leq \lambda(F, \hat{x}(t)) \|v\|_{\hat{x}(t)} \leq \sqrt{\nu} \|v\|_{\hat{x}(t)},$$

we obtain the following estimate of the first difference in the right-hand side of (18.7):

$$(18.8) \quad \begin{aligned} F_s(x) - F_s(\hat{x}(t)) &= F_t(x) - F_t(\hat{x}(t)) + (s - t) \langle c, x - \hat{x}(t) \rangle \\ &= F_t(x) - F_t(\hat{x}(t)) - (s/t - 1) \langle F'(\hat{x}(t)), x - \hat{x}(t) \rangle \\ &\leq F_t(x) - F_t(\hat{x}(t)) + |s/t - 1| \sqrt{\nu} \|x - \hat{x}(t)\|_{\hat{x}(t)}. \end{aligned}$$

By Theorem 16.4.6,

$$F_t(x) - F_t(\hat{x}(t)) \leq \rho(\lambda(F_t, x)) \leq \rho(\kappa)$$

and

$$\|x - \hat{x}(t)\|_{\hat{x}(t)} \leq \frac{\lambda(F_t, x)}{1 - \lambda(F_t, x)} \leq \frac{\kappa}{1 - \kappa}.$$

Therefore, it follows from inequality (18.8) that

$$(18.9) \quad F_s(x) - F_s(\hat{x}(t)) \leq \rho(\kappa) + \left| \frac{s}{t} - 1 \right| \cdot \frac{\kappa\sqrt{\nu}}{1 - \kappa}.$$

It remains to estimate the second difference

$$(18.10) \quad \begin{aligned} \phi(s) &= F_s(\hat{x}(t)) - F_s(\hat{x}(s)) \\ &= s\langle c, \hat{x}(t) \rangle - s\langle c, \hat{x}(s) \rangle + F(\hat{x}(t)) - F(\hat{x}(s)) \end{aligned}$$

in the right-hand side of (18.7).

The function  $\hat{x}(s)$  is continuously differentiable. This follows from the implicit function theorem, because  $\hat{x}(s)$  satisfies the equation

$$sc + F'(\hat{x}(s)) = 0,$$

and the second derivative  $F''(x)$  is continuous and non-singular everywhere. By implicit differentiation,

$$c + F''(\hat{x}(s))\hat{x}'(s) = 0,$$



**Empowering People. Improving Business.**

BI Norwegian Business School is one of Europe's largest business schools welcoming more than 20,000 students. Our programmes provide a stimulating and multi-cultural learning environment with an international outlook ultimately providing students with professional skills to meet the increasing needs of businesses.

BI offers four different two-year, full-time Master of Science (MSc) programmes that are taught entirely in English and have been designed to provide professional skills to meet the increasing need of businesses. The MSc programmes provide a stimulating and multi-cultural learning environment to give you the best platform to launch into your career.

- MSc in Business
- MSc in Financial Economics
- MSc in Strategic Marketing Management
- MSc in Leadership and Organisational Psychology

**BI NORWEGIAN BUSINESS SCHOOL**

EFMD **EQUIS** ACCREDITED

[www.bi.edu/master](http://www.bi.edu/master)

which means that

$$\hat{x}'(s) = -F''(\hat{x}(s))^{-1}c.$$

It now follows from equation (18.10) that the difference  $\phi(s)$  is continuously differentiable with derivative

$$\begin{aligned}\phi'(s) &= \langle c, \hat{x}(t) \rangle - \langle c, \hat{x}(s) \rangle - s \langle c, \hat{x}'(s) \rangle - \langle F'(\hat{x}(s)), \hat{x}'(s) \rangle \\ &= \langle c, \hat{x}(t) - \hat{x}(s) \rangle - s \langle c, \hat{x}'(s) \rangle + s \langle c, \hat{x}'(s) \rangle \\ &= \langle c, \hat{x}(t) - \hat{x}(s) \rangle,\end{aligned}$$

and a further differentiation gives

$$\begin{aligned}\phi''(s) &= -\langle c, \hat{x}'(s) \rangle = \langle c, F''(\hat{x}(s))^{-1}c \rangle \\ &= \langle s^{-1}F'(\hat{x}(s)), s^{-1}F''(\hat{x}(s))^{-1}F'(\hat{x}(s)) \rangle \\ &= s^{-2} \langle F'(\hat{x}(s)), F''(\hat{x}(s))^{-1}F'(\hat{x}(s)) \rangle = s^{-2} \lambda(F, \hat{x}(s))^2 \leq \nu s^{-2}.\end{aligned}$$

Now note that  $\phi(t) = \phi'(t) = 0$ . By integrating the inequality for  $\phi''(s)$  over the interval  $[t, u]$ , we therefore obtain the following estimate for  $u \geq t$ :

$$\phi'(u) = \phi'(u) - \phi'(t) \leq \int_t^u \nu s^{-2} ds = \nu(t^{-1} - u^{-1}).$$

Integrating once more over the interval  $[t, s]$  results in the inequality

$$\begin{aligned}(18.11) \quad F_s(\hat{x}(t)) - F_s(\hat{x}(s)) &= \phi(s) = \int_t^s \phi'(u) du \leq \nu \int_t^s (t^{-1} - u^{-1}) du \\ &= \nu \left( \frac{s}{t} - 1 - \ln \frac{s}{t} \right) = \nu \rho(1 - s/t)\end{aligned}$$

for  $s \geq t$ . The same conclusion is also reached for  $s < t$  by first integrating the inequality for  $\phi''(s)$  over the interval  $[u, t]$ , and then the resulting inequality for  $\phi'(u)$  over the interval  $[s, t]$ .

The inequality in the lemma is now finally a consequence of equation (18.7) and the estimates (18.9) and (18.11).  $\square$

*Continuation of the proof of Theorem 18.2.2.* By using the lemma's estimate of the difference  $F_s(x) - F_s(\hat{x}(s))$  when  $s = (1 + \gamma/\sqrt{\nu})t$ , we obtain the inequality

$$\left| \frac{F_s(x) - F_s(\hat{x}(s))}{\rho(-\kappa)} \right| \leq \left| \frac{\rho(\kappa) + \gamma\kappa(1 - \kappa)^{-1} + \nu \rho(-\gamma\nu^{-1/2})}{\rho(-\kappa)} \right|,$$

and  $\nu \rho(-\gamma\nu^{-1/2}) \leq \frac{1}{2}\gamma^2$ , because  $\rho(u) = -u - \ln(1 - u) \leq \frac{1}{2}u^2$  for  $u < 0$ . The number of inner iterations in each outer iteration is therefore bounded by the constant

$$K = \left\lceil \frac{\rho(\kappa) + \gamma\kappa(1 - \kappa)^{-1} + \frac{1}{2}\gamma^2}{\rho(-\kappa)} \right\rceil,$$



which only depends on the parameters  $\kappa$  and  $\gamma$ . For example,  $K = 5$  if  $\kappa = 0.4$  and  $\gamma = 0.32$ .

We now turn to the number of outer iterations. Set

$$\beta(\kappa) = \nu + \kappa(1 - \kappa)^{-1}\sqrt{\nu}.$$

Suppose that the stopping condition  $\epsilon t \geq \beta(\kappa)$  is triggered during iteration number  $k$  when  $t = (1 + \gamma/\sqrt{\nu})^k t_0$ . Because of Theorem 18.2.1, the current point  $x$  then satisfies the condition

$$\langle c, x \rangle - v_{\min} \leq \epsilon,$$

which shows that  $x$  approximates the minimum point with prescribed accuracy.

Since  $k$  is the least integer satisfying the inequality  $(1 + \gamma/\sqrt{\nu})^k \geq \beta(\kappa)/t_0\epsilon$ , we have

$$k = \left\lceil \frac{\ln(\beta(\kappa)/t_0\epsilon)}{\ln(1 + \gamma/\sqrt{\nu})} \right\rceil.$$

To simplify the denominator, we use the fact that  $\ln(1 + \gamma x)$  is a concave function. This implies that  $\ln(1 + \gamma x) \geq x \ln(1 + \gamma)$  if  $0 \leq x \leq 1$ , and hence

$$\ln(1 + \gamma/\sqrt{\nu}) \geq \frac{\ln(1 + \gamma)}{\sqrt{\nu}}.$$

Furthermore,  $\beta(\kappa) = \nu + \kappa(1 - \kappa)^{-1}\sqrt{\nu} \leq \nu + \kappa(1 - \kappa)^{-1}\nu = (1 - \kappa)^{-1}\nu$ . This gives us the estimate

$$k \leq \left\lceil \frac{\sqrt{\nu} \ln((1 - \kappa)^{-1}\nu/t_0\epsilon)}{\ln(1 + \gamma)} \right\rceil \leq K' \sqrt{\nu} \ln\left(\frac{\nu}{t_0\epsilon} + 1\right)$$

for the number of outer iterations with a constant  $K'$  that only depends on  $\kappa$  and  $\gamma$ , and by multiplying this with the constant  $K$  we obtain the corresponding estimate for the total number of inner iterations.  $\square$

## Phase 1

In order to use the path-following algorithm, we need a  $t_0 > 0$  and a point  $x_0 \in \text{int } X$  with Newton decrement  $\lambda(F_{t_0}, x_0) \leq \kappa$  to start from. Since the central path begins in the analytic center  $\hat{x}_F$  of  $X$  and  $\lambda(F, \hat{x}_F) = 0$ , it can be expected that  $(x_0, t_0)$  is good enough as a starting pair if only  $x_0$  is close enough to  $\hat{x}_F$  and  $t_0 > 0$  is sufficiently small. Indeed, this is true, and we shall show that one can generate such a pair by solving an artificial problem, given that one knows a point  $\bar{x} \in \text{int } X$ .

Therefore, let  $G_t: \text{int } X \rightarrow \mathbf{R}$ , where  $0 \leq t \leq 1$ , be the functions defined by

$$G_t(x) = -t\langle F'(\bar{x}), x \rangle + F(x).$$

The functions  $G_t$  are closed and self-concordant, and they have unique minimum points  $\bar{x}(t)$ .

Note that  $G_0 = F$ , and hence  $\bar{x}(0) = \hat{x}_F$ . Since  $G'_t(x) = -tF'(\bar{x}) + F'(x)$ ,  $G'_1(\bar{x}) = 0$ , and this means that  $\bar{x}$  is the minimum point of the function  $G_1$ . Hence,  $\bar{x}(1) = \bar{x}$ . The curve  $\{\bar{x}(t) \mid 0 \leq t \leq 1\}$  thus starts in the analytic center  $\hat{x}_F$  and ends in the given point  $\bar{x}$ . By using the path-following method, now following the curve *backwards*, we will therefore obtain a suitable starting point for phase 2 of the algorithm.

We use Newton's damped method to minimize  $G_t$  and note that  $G''_t = F''$  for all  $t$ , so the local norm with respect to the function  $G_t$  coincides with the local norm with respect to the function  $F$ , and we can thus unambiguously use the symbol  $\|\cdot\|_x$  for the local norm at the point  $x$ .

The algorithm for determining a starting pair  $(x_0, t_0)$  now looks like this.

## Need help with your dissertation?

Get in-depth feedback & advice from experts in your topic area. Find out what you can do to improve the quality of your dissertation!

Get Help Now



Go to [www.helpmyassignment.co.uk](http://www.helpmyassignment.co.uk) for more info



### Path-following algorithm, phase 1

**Given**  $\bar{x} \in \text{int } X$ , and parameters  $0 < \gamma < \frac{1}{2}\sqrt{\nu}$  and  $0 < \kappa < 1$ .

1. *Initiate*:  $x := \bar{x}$  and  $t := 1$ .
2. *Stopping criterion*: **stop** if  $\lambda(F, x) < \frac{3}{4}\kappa$  and set  $x_0 = x$ .
3. *Decrease t*:  $t := (1 - \gamma/\sqrt{\nu})t$ .
4. *Update x by using Newton's damped method on the function  $G_t$  with the current x as starting point*:
  - (i) Compute  $\lambda = \lambda(G_t, x)$ .
  - (ii) **quit** Newton's method if  $\lambda \leq \kappa/2$ , and go to line 2.
  - (iii) Compute the Newton step  $\Delta x_{\text{nt}} = -F''(x)^{-1}G'_t(x)$ .
  - (iv) *Update*:  $x := x + (1 + \lambda)^{-1}\Delta x_{\text{nt}}$ .
  - (v) Go to (i).

When the algorithm has stopped with a point  $x_0$ , we define  $t_0$  by setting

$$t_0 = \max\{t \mid \lambda(F_t, x_0) \leq \kappa\}.$$

The number of iterations in phase 1 is given by the following theorem.

**Theorem 18.2.4.** *Phase 1 of the path-following algorithm stops with a point  $x_0 \in \text{int } X$  after at most*

$$C\sqrt{\nu} \ln\left(\frac{\nu}{1 - \pi_{\hat{x}_F}(\bar{x})} + 1\right)$$

*inner iterations, where the constant  $C$  only depends on  $\kappa$  and  $\gamma$ , the number  $t_0$  satisfies the conditions  $\lambda(F_{t_0}, x_0) \leq \kappa$  and  $t_0 \geq \kappa/4 \text{Var}_X(c)$ .*

*Proof.* We start by estimating the number of inner iterations in each outer iteration; this number is bounded by the quotient

$$\frac{G_s(x) - G_s(\bar{x}(s))}{\rho(-\kappa/2)},$$

where  $s = (1 - \gamma/\sqrt{\nu})t$ , and Lemma 18.2.3 gives us the majorant

$$\rho(\kappa/2) + \frac{\kappa\sqrt{\nu}}{2 - \kappa} \cdot \frac{\gamma}{\sqrt{\nu}} + \nu \rho(\gamma/\sqrt{\nu})$$

for the numerator of the quotient. By Lemma 16.3.1,  $\nu\rho(\gamma/\sqrt{\nu}) \leq \gamma^2$ , so the number of inner iterations in each outer iteration is bounded by the constant

$$\frac{\rho(\kappa/2) + \kappa(2 - \kappa)^{-1}\gamma + \gamma^2}{\rho(-\kappa/2)}.$$

We now consider the outer iterations. Since  $F' = G'_t + tF'(\bar{x})$ ,

$$(18.12) \quad \begin{aligned} \lambda(F, x) &= \|F'(x)\|_x^* = \|G'_t(x) + tF'(\bar{x})\|_x^* \leq \|G'_t(x)\|_x^* + t\|F'(\bar{x})\|_x^* \\ &= \lambda(G_t, x) + t\|F'(\bar{x})\|_x^*. \end{aligned}$$

It follows from Theorem 18.1.11 that

$$\|F'(\bar{x})\|_x^* \leq (\nu + 2\sqrt{\nu})\|F'(\bar{x})\|_{\hat{x}_F}^* \leq 3\nu\|F'(\bar{x})\|_{\hat{x}_F}^*,$$

and from Theorem 18.1.8 that

$$\|F'(\bar{x})\|_{\hat{x}_F}^* = \sup_{\|v\|_{\hat{x}_F} \leq 1} \langle F'(\bar{x}), v \rangle \leq \frac{\nu}{1 - \pi_{\hat{x}_F}(\bar{x})}.$$

Hence

$$(18.13) \quad \|F'(\bar{x})\|_x^* \leq \frac{3\nu^2}{1 - \pi_{\hat{x}_F}(\bar{x})}.$$

During outer iteration number  $k$ , we have  $t = (1 - \gamma/\sqrt{\nu})^k$  and the point  $x$  satisfies the condition  $\lambda(G_t, x) \leq \kappa/2$  when Newton's method stops. So it follows from inequality (18.12) and the estimate (18.13) that the stopping condition  $\lambda(F, x) < \frac{3}{4}\kappa$  in line 2 of the algorithm is fulfilled if

$$\frac{1}{2}\kappa + \frac{3\nu^2}{1 - \pi_{\hat{x}_F}(\bar{x})}(1 - \gamma/\sqrt{\nu})^k \leq \frac{3}{4}\kappa,$$

i.e. if

$$k \ln(1 - \gamma/\sqrt{\nu}) < -\ln\left(\frac{12\kappa^{-1}\nu^2}{1 - \pi_{\hat{x}_F}(\bar{x})}\right).$$

By using the inequality  $\ln(1 - x) \leq -x$ , which holds for  $0 < x < 1$ , we see that the stopping condition is fulfilled for

$$k > \frac{\sqrt{\nu}}{\gamma} \ln\left(\frac{12\kappa^{-1}\nu^2}{1 - \pi_{\hat{x}_F}(\bar{x})}\right).$$

So the number of outer iterations is less than

$$K\sqrt{\nu} \ln\left(\frac{\nu}{1 - \pi_{\hat{x}_F}(\bar{x})} + 1\right),$$

where the constant  $K$  only depends on  $\kappa$  and  $\gamma$ , and this proves the estimate of the theorem, since the number of inner iterations in each outer iteration is bounded by a constant, which only depends on  $\kappa$  and  $\gamma$ .

The definition of  $t_0$  implies that  $\kappa = \lambda(F_{t_0}, x_0)$ , so we get the following inequalities with the aid of Theorem 18.1.10:

$$\begin{aligned}\kappa &= \lambda(F_{t_0}, x_0) = \|F'_{t_0}(x_0)\|_{x_0}^* = \|t_0 c + F'(x_0)\|_{x_0}^* \leq t_0 \|c\|_{x_0}^* + \|F'(x_0)\|_{x_0}^* \\ &= t_0 \|c\|_{x_0}^* + \lambda(F, x_0) \leq t_0 \text{Var}_X(c) + \frac{3}{4}\kappa.\end{aligned}$$

It follows that

$$t_0 \geq \frac{\kappa}{4 \text{Var}_X c}. \quad \square$$

The following complexity result is now obtained by combining the two phases of the path-following algorithm.

**Theorem 18.2.5.** *A standard problem (SP) with  $\nu$ -self-concordant barrier, tolerance level  $\epsilon > 0$  and starting point  $\bar{x} \in \text{int } X$  can be solved with at most*

$$C\sqrt{\nu} \ln(\nu\Phi/\epsilon + 1)$$

*Newton steps, where*

$$\Phi = \frac{\text{Var}_X(c)}{1 - \pi_{\hat{x}_F}(\bar{x})}$$

*and the constant  $C$  only depends on  $\gamma$  and  $\kappa$ .*

**Brain power**

By 2020, wind could provide one-tenth of our planet's electricity needs. Already today, SKF's innovative know-how is crucial to running a large proportion of the world's wind turbines.

Up to 25 % of the generating costs relate to maintenance. These can be reduced dramatically thanks to our systems for on-line condition monitoring and automatic lubrication. We help make it more economical to create cleaner, cheaper energy out of thin air.

By sharing our experience, expertise, and creativity, industries can boost performance beyond expectations. Therefore we need the best employees who can meet this challenge!

The Power of Knowledge Engineering

Plug into The Power of Knowledge Engineering.  
Visit us at [www.skf.com/knowledge](http://www.skf.com/knowledge)

**SKF**

*Proof.* Phase 1 provides a starting point  $x_0$  and an initial value  $t_0$  for phase 2, satisfying the condition  $t_0 \geq \kappa/(4 \text{Var}_X(c))$ . The number of inner iterations in phase 2 is therefore bounded by

$$O(1)\sqrt{\nu} \ln\left(\frac{4\nu \text{Var}_X(c)}{\kappa\epsilon} + 1\right) = O(1)\sqrt{\nu} \ln\left(\frac{\nu \text{Var}_X(c)}{\epsilon} + 1\right).$$

So the total number of inner iterations in the two phases is

$$\begin{aligned} O(1)\sqrt{\nu} \ln\left(\frac{\nu}{1 - \pi_{\hat{x}_F}(\bar{x})} + 1\right) + O(1)\sqrt{\nu} \ln\left(\frac{\nu \text{Var}_X(c)}{\epsilon} + 1\right) \\ = O(1)\sqrt{\nu} \ln(\nu\Phi/\epsilon + 1). \quad \square \end{aligned}$$

*Remark.* The algorithms in this section provide nice theoretical complexity results, but they are not suitable for practical use. The main limitation is the low updating factor  $(1 + O(1)\nu^{-1/2})$  of the penalty parameter  $t$ , which implies that the total number of Newton steps will be proportional to  $\sqrt{\nu}$ . For an LP problem with  $n = 1000$  variables and  $m = 10000$  inequalities, one would need to solve hundreds of linear equations with 1000 variables, which requires far more time than what is needed by the simplex algorithm. In the majority of outer iterations, one can, however, in practice increase the penalty parameter much faster than what is needed for the theoretical worst case analysis, without necessarily having to increase the number of Newton steps to maintain proximity to the central path. There are good practical implementations of the algorithm that use various dynamic strategies to control the penalty parameter  $t$ , and as a result only a moderate total number of Newton steps is needed, regardless of the size of the problem.

### 18.3 LP problems

We now apply the algorithm in the previous section on LP problems. Consider a problem of the type

$$(18.14) \quad \begin{aligned} \min \quad & \langle c, x \rangle \\ \text{s.t.} \quad & Ax \leq b \end{aligned}$$

where  $A = [a_{ij}]$  is an  $m \times n$ -matrix. We assume that the polyhedron

$$X = \{x \in \mathbf{R}^n \mid Ax \leq b\}$$

of feasible points is bounded and has a nonempty interior. The boundedness assumption implies that  $m > n$ .

The  $i$ th row of the matrix  $A$  is denoted by  $a_i$ , that is  $a_i = [a_{i1} \ a_{i2} \ \dots \ a_{in}]$ . The matrix product  $a_i x$  is thus well-defined.

As a barrier to the set  $X$  we use the  $m$ -self-concordant function

$$F(x) = - \sum_{i=1}^m \ln(b_i - a_i x).$$

The path-following algorithm started from an arbitrary point  $\bar{x} \in \text{int } X$  results in an  $\epsilon$ -solution, i.e. a point with a value of the objective function that approximates the optimal value with an error less than  $\epsilon$ , after at most

$$O(1)\sqrt{m} \ln(m\Phi/\epsilon + 1)$$

inner iterations, where

$$\Phi = \frac{\text{Var}_X(c)}{1 - \pi_{\hat{x}_F}(\bar{x})}.$$

We now estimate the number of arithmetic operations (additions, subtractions, multiplications and divisions) that are required during phase 2 of the algorithm to obtain this  $\epsilon$ -solution.

For each inner iteration of the Newton algorithm, we first have to compute the gradient and the hessian of the barrier function at the current point  $x$ , i.e.

$$F'(x) = \sum_{i=1}^m \frac{a_i^T}{b_i - a_i x} \quad \text{och} \quad F''(x) = \sum_{i=1}^m \frac{a_i^T a_i}{(b_i - a_i x)^2}.$$

This can be done with  $O(mn^2)$  arithmetic operations. The Newton direction  $\Delta x_{\text{nt}}$  at  $x$  is obtained as solution to the quadratic system

$$F''(x)\Delta x_{\text{nt}} = -(tc + F'(x))$$

of linear equations, and using Gaussian elimination, we find the solution after  $O(n^3)$  arithmetic operations. Finally,  $O(n)$  additional arithmetic operations, including one square root extraction, are needed to compute the Newton decrement  $\lambda = \lambda(F_t, x)$  and the new point  $x^+ = x + (1 + \lambda)^{-1}\Delta x_{\text{nt}}$ .

The corresponding estimate of the number of operations is also true for phase 1 of the algorithm.

The gradient and hessian computation is the most costly of the above computations since  $m > n$ . The total number of arithmetic operations in each iteration is therefore  $O(mn^2)$ , and by multiplying with the number of inner iterations, the overall arithmetic cost of the path-following algorithm is estimated to be no more than  $O(m^{3/2}n^2) \ln(m\Phi/\epsilon + 1)$  operations.

The resulting approximate minimum point  $\hat{x}(\epsilon)$  is an interior point of the polyhedron  $X$ , but the minimum is of course attained at an extreme point

on the the boundary of  $X$ . However, there is a simple procedure, called *purification* and described below, which starting from  $\hat{x}(\epsilon)$  finds an extreme point  $\hat{x}$  of  $X$  after no more than  $O(mn^2)$  arithmetic operations and with an objective function value that does not exceed the value at  $\hat{x}(\epsilon)$ . This means that we have the following result.

**Theorem 18.3.1.** *For the LP problem (18.14) at most*

$$O(m^{3/2}n^2) \ln(m\Phi/\epsilon + 1)$$

*arithmetic operations are needed to find an extreme point  $\hat{x}$  of the polyhedron of feasible points that approximates the minimum value with an error less than  $\epsilon$ .*

## Purification

The proof of the following theorem describes an algorithm for how to generate an extreme point with a value of the objective function that does not exceed the value at a given interior point of the polyhedron of feasible points.

What do you want to do?

No matter what you want out of your future career, an employer with a broad range of operations in a load of countries will always be the ticket. Working within the Volvo Group means more than 100,000 friends and colleagues in more than 185 countries all over the world. We offer graduates great career opportunities – check out the Career section at our web site [www.volvogroup.com](http://www.volvogroup.com). We look forward to getting to know you!

**VOLVO**  
AB Volvo (publ)  
[www.volvogroup.com](http://www.volvogroup.com)

VOLVO TRUCKS | RENAULT TRUCKS | MACK TRUCKS | VOLVO BUSES | VOLVO CONSTRUCTION EQUIPMENT | VOLVO PENTA | VOLVO AERO | VOLVO IT  
VOLVO FINANCIAL SERVICES | VOLVO 3P | VOLVO POWERTRAIN | VOLVO PARTS | VOLVO TECHNOLOGY | VOLVO LOGISTICS | BUSINESS AREA ASIA



**Theorem 18.3.2.** *Let*

$$\begin{aligned} \min \quad & \langle c, x \rangle \\ \text{s.t.} \quad & Ax \leq b \end{aligned}$$

*be an LP problem with  $n$  variables and  $m$  constraints, and suppose that the polyhedron  $X$  of feasible points is line-free and that the objective function is bounded below on  $X$ . For each point of  $X$  we can generate an extreme point of  $X$  with a value of the objective function that does not exceed the value at the given point with an algorithm using at most  $O(mn^2)$  arithmetic operations.*

*Proof.* The idea is very simple: Follow a half-line from the given point  $x^{(0)}$  with non-increasing function values until hitting upon a point  $x^{(1)}$  in a face  $F_1$  of the polyhedron  $X$ . Then follow a half-line in the face  $F_1$  with non-increasing function values until hitting upon a point  $x^{(2)}$  in the intersection  $F_1 \cap F_2$  of two faces, etc. After  $n$  steps, one has reached a point  $x^{(n)}$  in the intersection of  $n$  (independent) faces, i.e. an extreme point, with a function value that is less than or equal to the value at the starting point.

To estimate the number of arithmetic operation we have to study the above procedure in a little more detail.

We start by defining  $v^{(1)} = \mathbf{e}_1$  if  $c_1 < 0$ ,  $v^{(1)} = -\mathbf{e}_1$  if  $c_1 > 0$ , and  $v^{(1)} = \pm \mathbf{e}_1$  if  $c_1 = 0$ , where the sign in the latter case should be chosen so that the half-line  $x^{(0)} + tv^{(1)}$ ,  $t \geq 0$ , intersects the boundary of the polyhedron; this is possible since the polyhedron is assumed to be line-free. In the first two cases, the half-line also intersects the boundary of the polyhedron, because  $\langle c, x^{(0)} + tv^{(1)} \rangle = \langle c, x^{(0)} \rangle - t|c_1|$  tends to  $-\infty$  as  $t$  tends to  $\infty$  and the objective function is assumed to be bounded below on  $X$ . The intersection point  $x^{(1)} = x^{(0)} + t_1v^{(1)}$  between the half-line and the boundary of  $X$  can be computed with  $O(mn)$  arithmetic operations, since we only have to compute the vectors  $b - Ax^{(0)}$  and  $Av^{(1)}$ , and quotients between their coordinates in order to find the nonnegative parameter value  $t_1$ .

After renumbering the equations, we may assume that the point  $x^{(1)}$  lies in the hyperplane  $a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1$ . We now eliminate the variable  $x_1$  from the constraints and the objective function, which results in a system of the form

$$(18.15) \quad \left\{ \begin{array}{l} x_1 + a'_{12}x_2 + \cdots + a'_{1n}x_n = b'_1 \\ A' \begin{bmatrix} x_2 \\ \vdots \\ x_n \end{bmatrix} \leq b' \end{array} \right.$$

where  $A'$  is an  $(m - 1) \times (n - 1)$ -matrix, and in a new objective function

$$c'_2x_2 + \cdots + c'_nx_n + d',$$

which is the restriction of the original objective function to the current face. The number of operations required to perform the eliminations is  $O(mn)$ .

After  $O(mn)$  operations we have thus managed to find a point  $x^{(1)}$  in a face  $F_1$  of  $X$  with an objective function value  $\langle c, x^{(1)} \rangle = \langle c, x^{(0)} \rangle - t_1|c_1|$  not exceeding  $\langle c, x^{(0)} \rangle$ , and determined the equation of the face and the restriction of the objective function to the face. We now have a problem of lower dimension  $n - 1$  and with  $m - 1$  constraints.

We continue by choosing a descent vector  $v^{(2)}$  for the objective function that is parallel to the face  $F_1$ , and we achieve this by defining  $v^{(2)}$  so that  $v_2^{(2)} = \pm 1$ ,  $v_3^{(2)} = \cdots = v_n^{(2)} = 0$  (and  $v_1^{(2)} = -a'_{12}v_2^{(2)}$ ), where the sign of  $v_2^{(2)}$  should be chosen so that the objective function is non-decreasing along the half-line  $x^{(1)} + tv^{(2)}$ ,  $t \geq 0$ , and the half-line intersects the relative boundary of  $F_1$ . This means that  $v_2^{(2)} = 1$  if  $c'_2 < 0$  and  $v_2^{(2)} = -1$  if  $c'_2 > 0$ , while the sign of  $v_2^{(2)}$  is determined by the requirement that the half-line should intersect the boundary in the case  $c'_2 = 0$ .

We then determine the intersection between the half-line  $x^{(1)} + tv^{(2)}$ ,  $t \geq 0$ , and the relative boundary of  $F_1$ , which occurs in one of the remaining hyperplanes. If this hyperplane is the hyperplane  $a'_{21}x_2 + \cdots + a'_{2n}x_n = b'_2$ , say, we eliminate the variable  $x_2$  from the remaining constraints and the objective function. All this can be done with at most  $O(mn)$  operations and results in a point  $x^{(2)}$  in the intersection of two faces, and the new value of the objective function is  $\langle c, x^{(2)} \rangle = \langle c, x^{(1)} \rangle - t_2|c'_2| \leq \langle c, x^{(1)} \rangle$ .

After  $n$  iterations, which together require at most  $nO(mn) = O(mn^2)$  arithmetic operations, we have reached an extreme point  $\hat{x} = x^{(n)}$  with a function value that does not exceed the value at the starting point  $x^{(0)}$ . The coordinates of the extreme point are obtained by solving a triangular system of equations, which only requires  $O(n^2)$  operations. The total number of operations is thus  $O(mn^2)$ .  $\square$

EXAMPLE 18.3.1. We exemplify the purification algorithm with the LP problem

$$\begin{array}{ll} \min & -2x_1 + x_2 + 3x_3 \\ \text{s.t.} & \begin{cases} -x_1 + 2x_2 + x_3 \leq 4 \\ -x_1 + x_2 + x_3 \leq 2 \\ x_1 - 2x_2 \leq 1 \\ x_1 - x_2 - 2x_3 \leq 1 \end{cases} \end{array}$$

Starting from the interior point  $x^{(0)} = (1, 1, 1)$  with objective function value  $c^T x^{(0)} = 2$ , we shall find an extreme point with a lower value.

Since  $c_1 = -2 < 0$ , we begin by choosing  $v^{(1)} = (1, 0, 0)$  and by determining the point of intersection between the half-line  $x = x^{(0)} + tv^{(1)} = (1+t, 1, 1)$ ,  $t \geq 0$ , and the boundary of the polyhedron of feasible points. We find that the point  $x^{(1)} = (3, 1, 1)$ , corresponding to  $t = 2$ , satisfies all constraints and the third constraint with equality. So  $x^{(1)}$  lies in the face obtained by intersecting the polyhedron  $X$  with the supporting hyperplane  $x_1 - 2x_2 = 1$ . We eliminate  $x_1$  from the objective function and from the remaining constraints using the equation of this hyperplane, and consider the restriction of the objective function to the corresponding face, i.e. the function  $f(x) = -3x_2 + 3x_3 - 2$  restricted to the polyhedron given by the system

$$\begin{cases} x_1 - 2x_2 & = 1 \\ & x_3 \leq 5 \\ -x_2 + x_3 & \leq 3 \\ x_2 - 2x_3 & \leq 0 \end{cases}$$

The  $x_2$ -coefficient of our new objective function is negative, so we follow the half-line  $x_2 = 1 + t$ ,  $x_3 = 1$ ,  $t \geq 0$ , in the hyperplane  $x_1 - 2x_2 = 1$  until it hits a new supporting hyperplane, which occurs for  $t = 1$ , when it intersects

**gaiteye**<sup>®</sup>  
*Challenge the way we run*

**EXPERIENCE THE POWER OF  
FULL ENGAGEMENT...**

**RUN FASTER.  
RUN LONGER..  
RUN EASIER...**

**READ MORE & PRE-ORDER TODAY  
WWW.GAITEYE.COM**

the hyperplane  $x_2 - 2x_3 = 0$  in the point  $x^{(2)} = (5, 2, 1)$ . Elimination of  $x_2$  results in the objective function  $f(x) = -3x_3 - 2$  and the system

$$\begin{cases} x_1 - 2x_2 & = 1 \\ x_2 - 2x_3 & = 0 \\ x_3 & \leq 5 \\ -x_3 & \leq 3 \end{cases}$$

Our new half-line in the face  $F_1 \cap F_2$  is given by the equation  $x_3 = 1 + t$ ,  $t \geq 0$ , and the halfline intersects the third hyperplane  $x_3 = 5$  when  $t = 4$ , i.e. in a point with  $x_3$ -coordinate equal to 5. Back substitution gives  $x^{(3)} = (21, 10, 5)$ , which is an extreme point with objective function value equal to  $-17$ .  $\square$

## 18.4 Complexity

By the *complexity* of a problem we here mean the number of arithmetic operations needed to solve it, and in this section we will study the complexity of LP problems with rational coefficients. The *solution* of an LP problem consists by definition of the problem's optimal value and, provided the value is finite, of an optimal point. All known estimates of the complexity depend not only on the number of variables and constraints, but also on the size of the coefficients, and an appropriate measure of the size of a problem is given by the number of binary bits needed to represent all its coefficients.

**Definition.** The *input length* of a vector  $x = (x_1, x_2, \dots, x_n)$  in  $\mathbf{R}^n$  is the integer  $\ell(x)$  defined as

$$\ell(x) = \sum_{j=1}^n \lceil \log_2(|x_j| + 1) \rceil.$$

The number of digits in the binary expansion of a positive integer  $z$  is equal to  $\lceil \log_2(|z| + 1) \rceil$ . The binary representation of a negative integer  $z$  requires one bit more in order to take care of the sign, and so does the representation of  $z = 0$ . The number of bits to represent an arbitrary vector  $x$  in  $\mathbf{R}^n$  with integer coordinates is therefore at most  $\ell(x) + n$ .

The norm of a vector can be estimated using the input length, and we shall need the following simple estimate in the two cases  $p = 1$  and  $p = 2$ .

**Lemma 18.4.1.**  $\|x\|_p \leq 2^{\ell(x)}$  for all  $x \in \mathbf{R}^n$  and all  $p \geq 1$ .

*Proof.* The inequality is a consequence of the following trivial inequalities  $\sum_{j=1}^n a_j \leq \prod_{j=1}^n (a_j + 1)$ ,  $a^p + 1 \leq (a + 1)^p$  and  $\log_2(a + 1) \leq \lceil \log_2(a + 1) \rceil$ ,

which hold for nonnegative numbers  $a, a_j$ , and imply that

$$\|x\|_p^p = \sum_{j=1}^n |x_j|^p \leq \prod_{j=1}^n (|x_j|^p + 1) \leq \prod_{j=1}^n (|x_j| + 1)^p \leq 2^{p\ell(x)}. \quad \square$$

We will now study LP problems of the type

$$\begin{aligned} \text{(LP)} \quad & \min \langle c, x \rangle \\ & \text{s.t. } Ax \leq b \end{aligned}$$

where all coefficients of the  $m \times n$ -matrix  $A = [a_{ij}]$  and of the vectors  $b$  and  $c$  are integers. Every LP problem with rational coefficients can obviously be replaced by an equivalent problem of this type after multiplication with a suitable least common denominator. The polyhedron of feasible points will be denoted by  $X$  so that

$$X = \{x \in \mathbf{R}^n \mid Ax \leq b\}.$$

**Definition.** The two integers

$$\ell(X) = \ell(A) + \ell(b) \quad \text{and} \quad L = \ell(X) + \ell(c) + m + n,$$

where  $\ell(A)$  denotes the input length of the matrix  $A$ , considered as a vector in  $\mathbf{R}^{mn}$ , are called the *input length of the polyhedron  $X$*  and the *input length of the given LP problem (LP)*, respectively.

The main result of this section is the following theorem, which implies that there is a solution algorithm that is polynomial in the input length of the LP problem.

**Theorem 18.4.2.** *There is an algorithm which solves the LP problem (LP) with at most  $O((m+n)^{7/2}L)$  arithmetic operations.*

*Proof. I.* We begin by noting that we can without restriction assume that the polyhedron  $X$  of feasible points is line-free. Indeed, we can, if necessary replace the problem (LP) with the equivalent and line-free LP problem

$$\begin{aligned} \min \quad & \langle c, x^+ \rangle - \langle c, x^- \rangle \\ \text{s.t.} \quad & \begin{cases} Ax^+ - Ax^- \leq b \\ -x^+ \leq 0 \\ -x^- \leq 0. \end{cases} \end{aligned}$$

This LP problem in  $n' = 2n$  variables and with  $m' = m + 2n$  constraints has input length

$$\begin{aligned} L' &= 2\ell(A) + 2n + \ell(b) + 2\ell(c) + m' + n' \\ &\leq 2(\ell(A) + \ell(b) + \ell(c) + m + n) + 4n = 2L + 4n \leq 6L, \end{aligned}$$

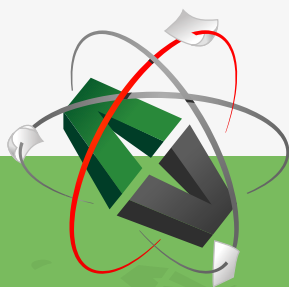
so any algorithm that solves this problem with  $O((m' + n')^{7/2}L')$  operations also solves problem (LP) with  $O((m + n)^{7/2}L)$  operations since  $m' + n' \leq 4(m + n)$  and  $L' \leq 6L$ .

From now on, we therefore assume that  $X$  is a *line-free polyhedron*, and for nonempty polyhedra  $X$  this implies that  $m \geq n$  and that  $X$  has at least one extreme point.

The assertion of the theorem is also trivially true for LP problems with only one variable, so we assume that  $m \geq n \geq 2$ . Finally, we can naturally assume that all the rows of the matrix  $A$  are nonzero, for if the  $k$ th row is identically zero, then the corresponding constraint can be deleted if  $b_k \geq 0$ , while the polyhedron  $X$  of feasible point is empty if  $b_k < 0$ . In the future, we can thus make use of the inequalities

$$\ell(X) \geq \ell(A) \geq m \geq n \geq 2 \text{ and } L \geq \ell(X) + m + n \geq \ell(X) + 4.$$

This e-book  
is made with  
**SetaPDF**



PDF components for PHP developers

[www.setasign.com](http://www.setasign.com)



**II.** Under the above assumptions, we will prove the theorem by showing:

1. With  $O(m^{7/2}L)$  operations, one can determine whether the optimal value of the problem is  $+\infty$ ,  $-\infty$  or finite, i.e. whether there are any feasible points or not, and if there are feasible points whether the objective functions is bounded below or not.
2. Given that the optimal value is finite, one can then determine an optimal solution with  $O(m^{3/2}n^2L)$  operations.

Since the proof of statement 1 uses the solution of an appropriate auxiliary LP problem with finite value, we begin by showing statement 2.

**III.** As a first building block we need a lemma that provides information about the extreme points of the polyhedron  $X$  in terms of its input length.

**Lemma 18.4.3.** (i) Let  $\hat{x}$  be an extreme point of the polyhedron  $X$ . Then, the following inequality holds for all nonzero coordinates  $\hat{x}_j$ :

$$2^{-\ell(X)} \leq |\hat{x}_j| \leq 2^{\ell(X)}.$$

Thus, all extreme points of  $X$  lie in the cube  $\{x \in \mathbf{R}^n \mid \|x\|_\infty \leq 2^{\ell(X)}\}$ .

(ii) If  $\hat{x}$  and  $\tilde{x}$  are two extreme points of  $X$  and  $\langle c, \hat{x} \rangle \neq \langle c, \tilde{x} \rangle$ , then

$$|\langle c, \hat{x} \rangle - \langle c, \tilde{x} \rangle| \geq 4^{-\ell(X)}.$$

*Proof.* To prove the lemma, we begin by recalling *Hadamard's inequality* for  $k \times k$ -matrices  $C = [c_{ij}]$  with columns  $C_{*1}, C_{*2}, \dots, C_{*k}$ , and which reads as follows:

$$|\det C| \leq \prod_{j=1}^k \|C_{*j}\|_2 = \prod_{j=1}^k \left( \sum_{i=1}^k c_{ij}^2 \right)^{1/2}.$$

The inequality is geometrically obvious – the left-hand side  $|\det C|$  is the volume of a (hyper)parallelepiped, spanned by the matrix columns, while the right-hand side is the volume of a (hyper)cuboid whose edges are of the same length as the edges of the parallelepiped.

By combining Hadamard's inequality with Lemma 18.4.1, we obtain the inequality

$$|\det C| \leq \prod_{j=1}^k 2^{\ell(C_{*j})} = 2^{\ell(C)}.$$

If  $C$  is a quadratic submatrix of the matrix  $[A \ b]$ , then obviously  $\ell(C) \leq \ell(A) + \ell(b) = \ell(X)$ , and it follows from the above inequality that

$$(18.16) \quad |\det C| \leq 2^{\ell(X)}.$$

Now let  $\hat{x}$  be an extreme point of the polyhedron  $X$ . According to Theorem 5.1.1 in Part I, there is a set  $\{i_1, i_2, \dots, i_n\}$  of row indices such that the extreme point  $\hat{x}$  is obtained as the unique solution to the equation system

$$\sum_{j=1}^n a_{ij}x_j = b_i, \quad i = i_1, i_2, \dots, i_n.$$

By Cramer's rule, we can write the solution in the form

$$\hat{x}_j = \frac{\Delta_j}{\Delta},$$

where  $\Delta$  is the determinant of the coefficient matrix and  $\Delta_j$  is the determinant obtained by replacing column number  $j$  in  $\Delta$  with the right-hand side of the equation system. The determinants  $\Delta$  and  $\Delta_j$  are integers, and their absolute values are at most equal to  $2^{\ell(X)}$ , because of inequality (18.16). This leads to the following estimates for all nonzero coordinates  $\hat{x}_j$ , i.e. for all  $j$  with  $\Delta_j \neq 0$ :

$$|\hat{x}_j| = |\Delta_j|/|\Delta| \leq 2^{\ell(X)}/1 = 2^{\ell(X)} \quad \text{and} \quad |\hat{x}_j| = |\Delta_j|/|\Delta| \geq 1/2^{\ell(X)} = 2^{-\ell(X)},$$

which is assertion (i) of the lemma.

(ii) The value of the objective function at the extreme point  $\hat{x}$  is

$$\langle c, \hat{x} \rangle = \left( \sum_{j=1}^n c_j \Delta_j \right) / \Delta = T / \Delta,$$

where the numerator  $T$  is an integer. If  $\tilde{x}$  is another extreme point, then of course we also have  $\langle c, \tilde{x} \rangle = T' / \Delta'$  for some integer  $T'$  and determinant  $\Delta'$  with  $|\Delta'| \leq 2^{\ell(X)}$ . It follows that the difference

$$\langle c, \tilde{x} \rangle - \langle c, \hat{x} \rangle = (T\Delta' - T'\Delta) / \Delta\Delta'$$

is either equal to zero or, if the numerator is nonzero, an integer with absolute value  $\geq 1/|\Delta\Delta'| \geq 4^{-\ell(X)}$ .  $\square$

**IV.** We shall use the path-following method, but this assumes that the polyhedron of feasible points is bounded and that there is an inner point from which to start phase 1. To get around this difficulty, we consider the following auxiliary problems in  $n + 1$  variables and  $m + 2$  linear constraints:

$$\begin{aligned} (\text{LP}_M) \quad & \min \quad \langle c, x \rangle + Mx_{n+1} \\ & \text{s.t.} \quad \begin{cases} Ax + (b - \mathbf{1})x_{n+1} \leq b \\ x_{n+1} \leq 2 \\ -x_{n+1} \leq 0. \end{cases} \end{aligned}$$



Here,  $M$  is a positive integer,  $\mathbf{1}$  denotes the vector  $(1, 1, \dots, 1)$  in  $\mathbf{R}^m$ , and  $x$  is as before the  $n$ -tuple  $(x_1, x_2, \dots, x_n)$ .

Let  $X'$  denote the polyhedron of feasible points for the problem  $(LP_M)$ . Since  $(x, x_{n+1}) = (0, 1)$  satisfies all constraints with strict inequality,  $(0, 1)$  is an inner point in  $X'$ .

We obtain the following estimates for the input length  $\ell(X')$  of the polyhedron  $X'$  and the input length  $L(M)$  of problem  $(LP_M)$ :

$$\begin{aligned}
 (18.17) \quad \ell(X') &= \ell(A) + \sum_{i=1}^m \lceil \log_2(|b_i - 1| + 1) \rceil + 1 + 1 + \ell(b) + 2 \\
 &\leq \ell(X) + 4 + \sum_{i=1}^m (1 + \lceil \log_2(1 + |b_i|) \rceil) \\
 &= \ell(X) + 4 + m + \ell(b) \leq 2\ell(X) + 4 \leq 2L - 4,
 \end{aligned}$$

$$\begin{aligned}
 (18.18) \quad L(M) &= \ell(X') + \ell(c) + \lceil \log_2(M + 1) \rceil + m + n + 3 \\
 &\leq 2\ell(X) + 2\ell(c) + \lceil \log_2 M \rceil + m + n + 8 \\
 &= 2L + \lceil \log_2 M \rceil - (m + n) + 8 \leq 2L + \lceil \log_2 M \rceil + 4.
 \end{aligned}$$

www.sylvania.com

**We do not reinvent the wheel we reinvent light.**

Fascinating lighting offers an infinite spectrum of possibilities: Innovative technologies and new markets provide both opportunities and challenges. An environment in which your expertise is in high demand. Enjoy the supportive working atmosphere within our global group and benefit from international career paths. Implement sustainable ideas in close cooperation with other specialists and contribute to influencing our future. Come and join us in reinventing light every day.

Light is OSRAM **OSRAM SYLVANIA**

The reason for studying our auxiliary problem  $(LP_M)$  is given in the following lemma.

**Lemma 18.4.4.** *Assume that problem (LP) has a finite value. Then:*

- (i) *Problem  $(LP_M)$  has a finite value for each integer  $M > 0$ .*
- (ii) *If  $(\hat{x}, 0)$  is an optimal solution to problem  $(LP_M)$ , then  $\hat{x}$  is an optimal solution to the original problem (LP).*
- (iii) *Assume that  $M \geq 2^{4L}$  and that the extreme point  $(\hat{x}, \hat{x}_{n+1})$  of  $X'$  is an optimal solution to problem  $(LP_M)$ . Then,  $\hat{x}_{n+1} = 0$ , so  $\hat{x}$  is an optimal solution to problem (LP).*

*Proof.* (i) The assumption of finite value means that the polyhedron  $X$  is nonempty and that the objective function  $\langle c, x \rangle$  is bounded below on  $X$ , and by Theorem 12.1.1 in Part II, this implies that the vector  $c$  lies in the dual cone of the recession cone  $\text{recc } X$ . Since

$$\begin{aligned} \text{recc } X' &= \{(x, x_{n+1}) \mid Ax + (b - \mathbf{1})x_{n+1} \leq 0, x_{n+1} = 0\} \\ &= \text{recc } X \times \{0\}, \end{aligned}$$

the dual cone of  $\text{recc } X'$  is equal to  $(\text{recc } X)^+ \times \mathbf{R}$ . We conclude that the vector  $(c, M)$  lies in the dual cone  $(\text{recc } X')^+$ , which means that the objective function of problem  $(LP_M)$  is bounded below on the nonempty set  $X'$ . Hence, our auxiliary problem has a finite value.

The polyhedron  $X'$  is line-free, since

$$\begin{aligned} \text{lin } X' &= \{(x, x_{n+1}) \mid Ax + (b - \mathbf{1})x_{n+1} = 0, x_{n+1} = 0\} \\ &= \text{lin } X \times \{0\} = \{(0, 0)\}. \end{aligned}$$

(ii) The point  $(x, 0)$  is feasible for problem  $(LP_M)$  if and only if  $x$  belongs to  $X$ , i.e. is feasible for our original problem (LP). So if  $(\hat{x}, 0)$  is an optimal solution to the auxiliary problem, then in particular

$$\langle c, \hat{x} \rangle = \langle c, \hat{x} \rangle + M \cdot 0 \leq \langle c, x \rangle + M \cdot 0 = \langle c, x \rangle$$

for all  $x \in X$ , which shows that  $\hat{x}$  is an optimal solution to problem (LP).

(iii) Assume that  $(\hat{x}, \hat{x}_{n+1})$  is an extreme point of the polyhedron  $X'$  and an optimal solution to problem  $(LP_M)$ . By Lemma 18.4.3, applied to the polyhedron  $X'$ , and the estimate (18.17), we then have the inequality

$$(18.19) \quad \|\hat{x}\|_\infty \leq 2^{\ell(X')} \leq 2^{2\ell(X)+4} \leq 2^{2L-4},$$

so it follows by using Lemma 18.4.1 that

$$\begin{aligned} |\langle c, \hat{x} \rangle| &\leq \sum_{j=1}^n |c_j| |\hat{x}_j| \leq \|c\|_1 \|\hat{x}\|_\infty \leq 2^{\ell(c)} \cdot 2^{2\ell(X)+4} \leq 2^{2\ell(X)+2\ell(c)+4} \\ &\leq 2^{2L-2m-2n+4} \leq 2^{2L-4}. \end{aligned}$$

Assume that  $\hat{x}_{n+1} \neq 0$ . Then  $\hat{x}_{n+1} \geq 2^{-\ell(X')} \geq 2^{-2L}$ , according to Lemma 18.4.3. The optimal value  $\hat{v}_M$  of the auxiliary problem (LP<sub>M</sub>) therefore satisfies the inequality

$$\hat{v}_M = \langle c, \hat{x} \rangle + M\hat{x}_{n+1} \geq M\hat{x}_{n+1} - |\langle c, \hat{x} \rangle| \geq M \cdot 2^{-2L} - 2^{2L-4}.$$

Let now  $x$  be an arbitrary extreme point of  $X$ . Since  $(x, 0)$  is a feasible point for problem (LP<sub>M</sub>) and since  $\|x\|_\infty \leq 2^{\ell(X)}$  by lemma 18.4.3, the optimal value  $\hat{v}_M$  must also satisfy the inequality

$$\hat{v}_M \leq \langle c, x \rangle + M \cdot 0 \leq |\langle c, x \rangle| \leq \|c\|_1 \cdot \|x\|_\infty \leq 2^{\ell(c)+\ell(X)} = 2^{L-m-n} \leq 2^{L-4}.$$

By combining the two inequalities for  $\hat{v}_M$ , we obtain the inequality

$$2^{L-4} \geq M \cdot 2^{-2L} - 2^{2L-4},$$

which implies that

$$M \leq 2^{3L-4} + 2^{4L-4} < 2^{4L}.$$

So if  $M \geq 2^{4L}$ , then  $\hat{x}_{n+1} = 0$ . □

**V.** We are now ready for the main step in the proof of Theorem 18.4.2.

**Lemma 18.4.5.** *Suppose that problem (LP) has a finite value. The path-following algorithm, applied to the problem (LP<sub>M</sub>) with  $\|x\|_\infty \leq 2^{2L}$  as an additional constraint,  $M = 2^{4L}$ ,  $\epsilon = 2^{-4L}$ , and  $(0, 1)$  as starting point for phase 1, and complemented with a subsequent purification operation, generates an optimal solution to problem (LP) after at most  $O(m^{3/2}n^2L)$  arithmetic operations.*

*Proof.* It follows from the previous lemma and the estimate (18.19) that the LP problem (LP<sub>M</sub>) has an optimal solution  $(\hat{x}, 0)$  which satisfies the additional constraint  $\|\hat{x}\|_\infty \leq 2^{2L}$  if  $M = 2^{4L}$ . The LP problem obtained from (LP<sub>M</sub>) by adding the  $2n$  constraints

$$x_j \leq 2^{2L} \quad \text{and} \quad -x_j \leq 2^{2L}, \quad j = 1, 2, \dots, n,$$

therefore has the same optimal value as (LP<sub>M</sub>).

The extended problem has  $m + 2n + 2 = O(m)$  linear constraints, and the point  $\bar{z} = (\bar{x}, \bar{x}_{n+1}) = (0, 1)$  is an interior point of the compact polyhedron of feasible points, which we denote by  $Z$ . By Theorem 18.3.1, the path-following algorithm with  $\epsilon = 2^{-4L}$  and  $\bar{z}$  as the starting point therefore stops after  $O((m + 2n + 2)^{3/2}n^2) \ln((m + 2n + 2)\Phi/\epsilon + 1) = O(m^{3/2}n^2) \ln(m2^{4L}\Phi + 1)$  arithmetic operations at a point in the polyhedron  $X'$  and with a value of the objective function that approximates the optimal value  $\hat{v}_M$  with an error less than  $2^{-4L}$ .

Purification according to the method in Theorem 18.3.2 leads to an extreme point  $(\hat{x}, \hat{x}_{n+1})$  of  $X'$  with a value of the objective function less than  $\hat{v}_M + 2^{-4L}$ , and since  $2^{-4L} = 4^{-2L} < 4^{-\ell(X')}$ , it follows from Lemma 18.4.3 that  $(\hat{x}, \hat{x}_{n+1})$  is an optimal solution to  $(LP_M)$ . By Lemma 18.4.4, this implies that  $\hat{x}$  is an optimal solution to the original problem (LP).

The purification process requires  $O(mn^2)$  arithmetic operations, so the total arithmetic cost is

$$O(mn^2) + O(m^{3/2}n^2) \ln(m2^{4L}\Phi + 1) = O(m^{3/2}n^2) \ln(m2^{4L}\Phi + 1)$$

operations. It thus only remains to prove that  $\ln(m2^{4L}\Phi + 1) = O(L)$ , and since  $m \leq L$ , this will follow if we show that  $\ln \Phi = O(L)$ .



Discover the truth at [www.deloitte.ca/careers](http://www.deloitte.ca/careers)

**Deloitte.**

© Deloitte & Touche LLP and affiliated entities.



By definition,

$$\Phi = \text{Var}_Z(c, M) \cdot \frac{1}{1 - \pi_{\hat{z}_F}(\bar{z})},$$

where  $\hat{z}_F$  is the analytic center of  $Z$  with respect to the relevant logarithmic barrier  $F$ . The absolute value of the objective function at an arbitrary point  $(x, x_{n+1}) \in Z$  can be estimated by

$$|\langle c, x \rangle + Mx_{n+1}| \leq \|c\|_1 \|x\|_\infty + 2M \leq 2^{\ell(c)+2L} + 2 \cdot 2^{4L} \leq 2^{4L+2},$$

and the maximal variation of the function is at most twice this value. Hence,

$$\text{Var}_Z(c, M) \leq 2^{4L+3}.$$

The second component of  $\Phi$  is estimated using Theorem 18.1.7. Let  $\bar{B}_\infty(a, a_{n+1}; r)$  denote the closed ball of radius  $r$  in  $\mathbf{R}^{n+1} = \mathbf{R}^n \times \mathbf{R}$  with center at the point  $(a, a_{n+1})$  and with distance given by the maximum norm, i.e.

$$\bar{B}_\infty(a, a_{n+1}; r) = \{(x, x_{n+1}) \in \mathbf{R}^n \times \mathbf{R} \mid \|x - a\|_\infty \leq r, |x_{n+1} - a_{n+1}| \leq r\}.$$

The polyhedron  $Z$  is by definition included in the ball  $\bar{B}_\infty(0, 0; 2^{2L})$ . On the other hand, the tiny ball  $\bar{B}_\infty(\bar{z}; 2^{-L})$  is included in  $Z$ , for if  $\|x\|_\infty \leq 2^{-L}$  and  $|x_{n+1} - 1| \leq 2^{-L}$ , then

$$\begin{aligned} \sum_{j=1}^n a_{ij}x_j + (b_i - 1)x_{n+1} - b_i &= \sum_{j=1}^n a_{ij}x_j + b_i(x_{n+1} - 1) - x_{n+1} \\ &\leq \sum_{j=1}^n |a_{ij}||x_j| + |b_i||x_{n+1} - 1| - x_{n+1} \leq 2^{-L} \left( \sum_{j=1}^n |a_{ij}| + |b_i| \right) - (1 - 2^{-L}) \\ &\leq 2^{-L+\ell(X)} + 2^{-L} - 1 \leq 2^{-4} + 2^{-L} - 1 < 0, \end{aligned}$$

which proves that the  $i$ th inequality of the system  $Ax + (b - 1)x_{n+1} \leq b$  holds with strict inequality for  $i = 1, 2, \dots, m$ , and the remaining inequalities that define the polyhedron  $Z$  are obviously strictly satisfied.

It therefore follows from Theorem 18.1.7 that

$$\pi_{\hat{z}_F}(\bar{z}) \leq \frac{2 \cdot 2^{2L}}{2 \cdot 2^{2L} + 2^{-L}},$$

and that consequently

$$\frac{1}{1 - \pi_{\hat{z}_F}(\bar{z})} \leq 2 \cdot 2^{3L} + 1 < 2^{3L+2}.$$

This implies that  $\Phi \leq 2^{4L+3} \cdot 2^{3L+2} = 2^{7L+5}$ . Hence,  $\ln \Phi = O(L)$ , which completes the proof of the lemma.  $\square$

**VI.** It remains to show that  $O(m^{7/2}L)$  operations are sufficient to decide whether the optimal value of the original problem (LP) is  $+\infty$ ,  $-\infty$  or finite.

To decide whether the value is  $+\infty$  or not, i.e. whether the polyhedron  $X$  is empty or not, we consider the artificial LP problem

$$\begin{aligned} \min \quad & x_{n+1} \\ \text{s.t.} \quad & \begin{cases} Ax - \mathbf{1}x_{n+1} \leq b \\ -x_{n+1} \leq 0 \end{cases} \end{aligned}$$

This problem has feasible points since  $(0, t)$  satisfies all constraints for sufficiently large positive numbers  $t$ . The optimal value of the problem is apparently greater than or equal to zero, and it is equal to zero if and only if  $X \neq \emptyset$ .

So we can decide whether the polyhedron  $X$  is empty or not by determining an optimal solution to the artificial problem. The input length of this problem is  $\ell(X) + 2m + n + 4$ , and since this number is  $\leq 2L$ , it follows from Lemma 18.4.5 that we can decide whether  $X$  is empty or not with  $O(m^{3/2}n^2L)$  arithmetic operations.

Note that we do not need to solve the artificial problem exactly. If the value is greater than zero, then, because of Lemma 18.4.3, it is namely greater than or equal to  $2^{-2L}$ . It is therefore sufficient to determine a point that approximates the value with an error of less than  $2^{-2L}$  to know if the value is zero or not.

**VII.** If the polyhedron  $X$  is nonempty, we have as the next step to decide whether the objective function is bounded below. This is the case if and only if the dual problem to problem (LP) has feasible points, and this dual maximization problem is equivalent to the minimization problem

$$\begin{aligned} \min \quad & \langle -b, y \rangle \\ \text{s.t.} \quad & \begin{cases} A^T y \leq c \\ -A^T y \leq -c \\ -y \leq 0, \end{cases} \end{aligned}$$

which is a problem with  $m$  variables,  $2n + m (= O(m))$  constraints and input length

$$2\ell(A) + m + 2\ell(c) + \ell(b) + m + (2n + m) \leq 2L + m \leq 3L.$$

So it follows from step VI that we can decide whether the dual problem has any feasible points with  $O(m^{7/2}L)$  operations.

The proof of Theorem 18.4.2 is now complete. □

## Exercises

- 18.1** Show that if the functions  $f_i$  are  $\nu_i$ -self-concordant barriers to the subsets  $X_i$  of  $\mathbf{R}^{n_i}$ , then  $f(x_1, \dots, x_m) = f_1(x_1) + \dots + f_m(x_m)$  is a  $(\nu_1 + \dots + \nu_m)$ -self-concordant barrier to the product set  $X_1 \times \dots \times X_m$ .
- 18.2** Prove that the dual local norm  $\|v\|_x^*$  that is associated with the function  $f$  is finite if and only if  $v$  belongs to  $\mathcal{N}(f''(x))^\perp$ , and that the restriction of  $\|\cdot\|_x^*$  to  $\mathcal{N}(f''(x))^\perp$  is a proper norm.
- 18.3** Let  $X$  be a closed proper convex cone with nonempty interior, let  $\nu \geq 1$  be a real number, and suppose that the function  $f: \text{int } X \rightarrow \mathbf{R}$  is closed and self-concordant and that  $f(tx) = f(x) - \nu \ln t$  for all  $x \in \text{int } X$  and all  $t > 0$ . Prove that
- a)  $f'(tx) = t^{-1}f'(x)$       b)  $f'(x) = -f''(x)x$       c)  $\lambda(f, x) = \sqrt{\nu}$ .
- The function  $f$  is in other words a  $\nu$ -self-concordant barrier to  $X$ .
- 18.4** Show that the nonnegative orthant  $X = \mathbf{R}_+^n$ ,  $\nu = n$  and the logarithmic barrier  $f(x) = -\sum_{i=1}^n \ln x_i$  fulfill the assumptions of the previous exercise.
- 18.5** Let  $X = \{(x, x_{n+1}) \in \mathbf{R}^n \times \mathbf{R} \mid x_{n+1} \geq \|x\|_2\}$ .
- a) Show that the function  $f(x) = -\ln(x_{n+1}^2 - (x_1^2 + \dots + x_n^2))$  is self-concordant on  $\text{int } X$ .

SIMPLY CLEVER

ŠKODA



We will turn your CV into  
an opportunity of a lifetime



Do you like cars? Would you like to be a part of a successful brand?  
We will appreciate and reward both your enthusiasm and talent.  
Send us your CV. You will be surprised where it can take you.

Send us your CV on  
[www.employerforlife.com](http://www.employerforlife.com)



b) Show that  $X$ ,  $\nu = 2$  and  $f$  fulfill the assumptions of exercise 18.3. The function  $f$  is thus a 2-self-concordant barrier to  $X$ .

**18.6** Suppose that the function  $f: \mathbf{R}_{++} \rightarrow \mathbf{R}$  is convex, three times continuously differentiable and that

$$|f'''(x)| \leq 3 \frac{f''(x)}{x}$$

for all  $x > 0$ . The function

$$F(x, y) = -\ln(y - f(x)) - \ln x$$

with  $X = \{(x, y) \in \mathbf{R}^2 \mid x > 0, y > f(x)\}$  as domain is self-concordant according to exercise 16.3. Show that  $F$  is a 2-self-concordant barrier to the closure  $\text{cl } X$ .

**18.7** Prove that the function

$$F(x, y) = -\ln(y - x \ln x) - \ln x$$

is a 2-self-concordant barrier to the epigraph

$$\{(x, y) \in \mathbf{R}^2 \mid y \geq x \ln x, x \geq 0\}.$$

**18.8** Prove that the function

$$G(x, y) = -\ln(\ln y - x) - \ln y$$

is a 2-self-concordant barrier to the epigraph  $\{(x, y) \in \mathbf{R}^2 \mid y \geq e^x\}$ .

I joined MITAS because  
I wanted **real responsibility**

The Graduate Programme  
for Engineers and Geoscientists  
[www.discovermitas.com](http://www.discovermitas.com)



**Month 16**

I was a construction  
supervisor in  
the North Sea  
advising and  
helping foremen  
solve problems

Real work  
International opportunities  
Three work placements



 **MAERSK**



# Bibliographical and historical notices

Newton's method is a classic iterative algorithm for finding critical points of differentiable functions, and it was proven by Kantorovich [1] that the algorithm converges quadratically when the function has a Lipschitz continuous, positive definite second derivatives in a neighborhood of the critical point, provided the starting point is selected close enough.

Barrier methods for solving nonlinear optimization problems were first used during the 1950s. The central path with logarithmic barriers was studied by Fiacco and McCormick, and their book on sequential minimization techniques – Fiacco–McCormick [1], first published in 1968 – is the standard work in the field. The methods worked well in practice, for the most part, but there were no theoretical complexity results. They lost in popularity in the 1970s and then experienced a renaissance in the wake of Karmarkar's discovery.

Karmarkar's [1] polynomial algorithm for linear programming proceeds by mapping the polyhedron of feasible points and the current approximate solution  $x_k$  onto a new polyhedron and a new point  $x'_k$  which is located near the center of the new polyhedron, using a projective scaling transformation. Thereafter, a step is taken in the transformed space which results in a point  $x_{k+1}$  with a lower objective function value. The progress is measured by means of a logarithmic potential function.

It was soon noted that Karmarkar's potential-reducing algorithm was akin to previously studied path-following methods, and Renegar [1] and Gonzaga [1] managed to show that the path-following method with logarithmic barrier is polynomial for LP problems.

A general introduction to linear programming and the algorithm development in the area until the late 1980s (the ellipsoid method, Karmarkar's algorithm, etc.) is given by Goldfarb–Todd [1]. An overview of potential-reducing algorithms is given by Todd [1], while Gonzaga [2] describes the evolution of path-following algorithms until 1992.

A breakthrough in convex optimization occurred in the late 1980s, when Yurii Nesterov discovered that Gonzaga’s and Renegar’s proof only used two properties of the logarithmic barrier function, namely, that it satisfies the two differential inequalities, which with Nesterov’s terminology means that the barrier is self-concordant with finite parameter  $\nu$ . Since explicit computable self-concordant barriers exist for a number of important types of convex sets, the theoretical complexity results for linear programming could now be extended to a large class of convex optimization problems, and Nemirovskii together with Nesterov developed algorithms for convex optimization based on self-concordant barriers. See Nesterov–Nesterovski [1].

A modern textbook on convex optimization, which in addition to theory and algorithms also contains lots of interesting applications from a variety of fields, is the book by Boyd–Vandenberghe [1].

## References

Boyd, S. & Vandenberghe, L.

- [1] *Convex Optimization*, Cambridge Univ. Press, Cambridge, UK, 2004.

Fiacco, A.V. & McCormick, G.P.

- [1] *Nonlinear Programming: Sequential Unconstrained Minimization Techniques*. Society for Industrial and Applied Mathematics, 1990. (First published in 1968 by Research Analysis Corporation.)

Goldfarb, D.G. & Todd, M.J.

- [1] Linear programming. Chapter 2 in Nemhauser, G.L. et al. (eds.), *Handbooks in Operations Research and Management Science, vol. 1: Optimization*, North-Holland, 1989.

Gonzaga, C.C.

- [1] An algorithm for solving linear programming problems in  $O(n^3L)$  operations. Pages 1–28 in Megiddo, N. (ed.), *Progress in Mathematical Programming: Interior-Point and Related Methods*, Springer-Verlag, 1988.  
[2] Path-Following Methods for Linear Programming, *SIAM Rev.* 34 (1992), 167–224.

Kantorovich, L.V.

- [1] *Functional Analysis and Applied Mathematics*. National Bureau of Standards, 1952. (First published in Russian in 1948.)

Karmarkar, N.

- [1] A new polynomial-time algorithm for linear programming, *Combinatorica* 4 (1984), 373–395.

Nesterov, Y. & Nemirovskii, A.

- [1] *Interior-Point Polynomial Algorithms in Convex Programming*. Society for Industrial and Applied Mathematics, 1994.

Renegar, J.

- [1] A polynomial-time algorithm based on Newton's method for linear programming, *Math. Programm.* 40 (1988), 59–94.

Todd, M.

- [1] Potential-reduction methods in mathematical programming, *Math. Program.* 76 (1997), 3–45.

**ie** business school

#1 EUROPEAN BUSINESS SCHOOL  
FINANCIAL TIMES  
2013

#gobeyond

**MASTER IN MANAGEMENT**

**Because achieving your dreams is your greatest challenge.** IE Business School's Master in Management taught in English, Spanish or bilingually, trains young high performance professionals at the beginning of their career through an innovative and stimulating program that will help them reach their full potential.

- Choose your area of specialization.
- Customize your master through the different options offered.
- Global Immersion Weeks in locations such as London, Silicon Valley or Shanghai.

*Because you change, we change with you.*

www.ie.edu/master-management | mim.admissions@ie.edu |

# Answers and solution to the exercises

## Chapter 14

14.1  $x_1 = (\frac{4}{9}, -\frac{1}{9}), x_2 = (\frac{2}{27}, \frac{2}{27}), x_3 = (\frac{8}{243}, -\frac{2}{243})$ .

14.3  $hf'(x_k) = f(x_k) - f(x_{k+1}) \rightarrow f(\hat{x}) - f(\hat{x}) = 0$  and  $hf'(x_k) \rightarrow hf'(\hat{x})$ .  
 Hence,  $f'(\hat{x}) = 0$ .

## Chapter 15

15.1  $\Delta x_{nt} = -x \ln x, \lambda(f, x) = \sqrt{x} \ln x, \|v\|_x = |v|/\sqrt{x}$ .

15.2 a)  $\Delta x_{nt} = (\frac{1}{3}, \frac{1}{3}), \lambda(f, x) = \sqrt{\frac{1}{3}}, \|v\|_x = \frac{1}{2}\sqrt{5v_1^2 + 2v_1v_2 + 5v_2^2}$

b)  $\Delta x_{nt} = (\frac{1}{3}, -\frac{2}{3}), \lambda(f, x) = \sqrt{\frac{1}{3}}, \|v\|_x = \frac{1}{2}\sqrt{8v_1^2 + 8v_1v_2 + 5v_2^2}$ .

15.3  $\Delta x_{nt} = (v_1, v_2)$ , where  $v_1 + v_2 = -1 - e^{-(x_1+x_2)}$ ,

$\lambda(f, x) = e^{(x_1+x_2)/2} + e^{-(x_1+x_2)/2}, \|v\|_x = e^{(x_1+x_2)/2}|v_1 + v_2|$ .

15.4 If  $\text{rank } A < m$ , then  $\text{rank } M < m + n$ , and if  $\mathcal{N}(A) \cap \mathcal{N}(P)$  contains a nonzero vector  $x$ , then  $M \begin{bmatrix} x \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ . Hence, the matrix  $M$  has no inverse in these cases.

Conversely, suppose that  $\text{rank } A = m$ , i.e. that  $\mathcal{N}(A^T) = \{0\}$ , and that  $\mathcal{N}(A) \cap \mathcal{N}(P) = \{0\}$ . We show that the coefficient matrix  $M$  is invertible by showing that the homogeneous system

$$\begin{cases} Px + A^T y = 0 \\ Ax = 0 \end{cases}$$

has no other solutions than the trivial one,  $x = 0$  and  $y = 0$ .

By multiplying the first equation from the left by  $x^T$  we obtain

$$0 = x^T Px + x^T A^T y = x^T Px + (Ax)^T y = x^T Px,$$

and since  $P$  is positive semidefinite, it follows that  $Px = 0$ . The first equation now gives  $A^T y = 0$ . Hence,  $x \in \mathcal{N}(A) \cap \mathcal{N}(P)$  and  $y \in \mathcal{N}(A^T)$ , which means that  $x = 0$  and  $y = 0$ .

15.5 a) By assumption,  $\langle v, f''(x)v \rangle \geq \mu\|v\|^2$  if  $Av = 0$ . Since  $AC = 0$ , we conclude that

$$\begin{aligned}\langle w, \tilde{f}''(z)w \rangle &= \langle w, C^T f''(x)Cw \rangle = \langle Cw, f''(x)Cw \rangle \geq \mu\|Cw\|^2 \\ &= \mu\langle w, C^T Cw \rangle \geq \mu\sigma\|w\|^2\end{aligned}$$

for all  $w \in \mathbf{R}^p$ , which shows that the function  $\tilde{f}$  is  $\mu\sigma$ -strongly convex.

b) The assertion follows from a) if we show that the restriction of  $f$  to  $X$  is a  $K^{-2}M^{-1}$ -strongly convex function. So assume that  $x \in X$  and that  $Av = 0$ . Then

$$\begin{bmatrix} f''(x) & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} v \\ 0 \end{bmatrix} = \begin{bmatrix} f''(x)v \\ 0 \end{bmatrix}$$

and due to the bound on the norm of the inverse matrix, we conclude that

$$\|v\| \leq K\|f''(x)v\|.$$

The positive semidefinite second derivative  $f''(x)$  has a positive semidefinite square root  $f''(x)^{1/2}$  and  $\|f''(x)^{1/2}\| = \|f''(x)\|^{1/2} \leq M^{1/2}$ . It follows that

$$\begin{aligned}\|f''(x)v\|^2 &= \|f''(x)^{1/2}f''(x)^{1/2}v\|^2 \leq \|f''(x)^{1/2}\|^2\|f''(x)^{1/2}v\|^2 \\ &\leq M\|f''(x)^{1/2}v\|^2 = M\langle v, f''(x)v \rangle,\end{aligned}$$

which inserted in the above inequality results in the inequality

$$\langle v, f''(x)v \rangle \geq K^{-2}M^{-1}\|v\|^2.$$

## Chapter 16

16.2 Let  $P_i$  denote the projection of  $\mathbf{R}^{n_1} \times \dots \times \mathbf{R}^{n_m}$  onto then  $i$ th factor  $\mathbf{R}^{n_i}$ . Then  $f(x) = \sum_{i=1}^m f_i(P_i x)$ , so it follows from Theorems 16.1.5 and 16.1.6 that  $f$  is self-concordant.

16.3 a) The function  $g$  is convex, since  $g''(x) = \frac{f'(x)^2}{f(x)^2} - \frac{f''(x)}{f(x)} + \frac{1}{x^2} \geq 0$ .

$$g'''(x) = -\frac{f'''(x)}{f(x)} + 3\frac{f'(x)f''(x)}{f(x)^2} - 2\frac{f'(x)^3}{f(x)^3} - \frac{2}{x^3}$$

implies that

$$|g'''(x)| \leq 3\frac{f'''(x)}{x|f(x)|} + 3\frac{|f'(x)|f''(x)}{f(x)^2} + 2\frac{|f'(x)|^3}{|f(x)|^3} + 2\frac{1}{x^3}.$$

The inequality  $|g'''(x)| \leq 2g''(x)^{3/2}$ , which proves that the function  $g$  is self-concordant, is now obtained by choosing  $a = \sqrt{f''(x)/|f(x)|}$ ,  $b = |f'(x)|/|f(x)|$  and  $c = 1/x$  in the equality

$$3a^2b + 3a^2c + 2b^3 + 2c^3 \leq 2(a^2 + b^2 + c^2)^{3/2}.$$

To prove this inequality, we can due to homogeneity assume that

$$a^2 + b^2 + c^2 = 1.$$

Inserting  $a^2 = 1 - b^2 - c^2$  into the inequality, we can rewrite it as  $(b + c)(3 - (b + c)^2) \leq 2$ , which holds since  $x(3 - x^2) \leq 2$  for  $x \geq 0$ .

16.3 b) Let  $\phi(t) = F(x_0 + \alpha t, y_0 + \beta t)$  be the restriction of  $F$  to an arbitrary line through the point  $(x_0, y_0)$  in  $\text{dom } F$ . We will prove that  $\phi$  is self-concordant, and we have to treat the cases  $\alpha = 0$  and  $\alpha \neq 0$  separately. If  $\alpha = 0$ , then  $\phi(t) = -\ln(\beta t + a) + b$ , where  $a = y_0 - f(x_0)$  and  $b = -\ln x_0$ , so  $\phi$  is self-concordant in this case.

To prove the case  $\alpha \neq 0$ , we note that  $f(x) - Ax - B$  satisfies the assumptions of the exercise for each choice of the constants  $A$  and  $B$ , and hence  $h(x) = -\ln(Ax + B - f(x)) - \ln x$  is self-concordant according to the result in a). But  $\phi(t) = h(\alpha t + x_0)$ , where  $A = \beta/\alpha$  and  $B = y_0 - \beta x_0/\alpha$ . Thus,  $\phi$  is self-concordant.



## STUDY AT A TOP RANKED INTERNATIONAL BUSINESS SCHOOL

Reach your full potential at the Stockholm School of Economics, in one of the most innovative cities in the world. The School is ranked by the Financial Times as the number one business school in the Nordic and Baltic countries.

Visit us at [www.hhs.se](http://www.hhs.se)



16.6 a) Set  $\lambda = \lambda(f, x)$  and use the inequalities (16.7) and (16.6) in Theorem 16.3.2 with  $y = x^+$  and  $v = x^+ - x = (1 + \lambda)^{-1}\Delta x_{\text{nt}}$ . This results in the inequality

$$\begin{aligned} \langle f'(x^+), w \rangle &\leq \langle f'(x), w \rangle + \frac{1}{1 + \lambda} \langle f''(x)\Delta x_{\text{nt}}, w \rangle + \frac{\lambda^2 \|w\|_x}{(1 + \lambda)^2(1 - \lambda/(1 + \lambda))} \\ &= \langle f'(x), w \rangle - \frac{1}{1 + \lambda} \langle f'(x), w \rangle + \frac{\lambda^2}{1 + \lambda} \|w\|_x \\ &= \frac{\lambda}{1 + \lambda} \langle f'(x), w \rangle + \frac{\lambda^2}{1 + \lambda} \|w\|_x \\ &\leq \frac{\lambda}{1 + \lambda} \lambda \|w\|_x + \frac{\lambda^2}{1 + \lambda} \|w\|_x = \frac{2\lambda^2}{1 + \lambda} \|w\|_x \\ &\leq \frac{2\lambda^2 \|w\|_{x^+}}{(1 + \lambda)(1 - \lambda/(1 + \lambda))} = 2\lambda^2 \|w\|_{x^+} \end{aligned}$$

with  $\lambda(f, x^+) \leq 2\lambda^2$  as conclusion.

## Chapter 18

18.1 Follows from Theorems 18.1.3 and 18.1.2.

18.2 To prove the implication  $\|v\|_x^* < \infty \Rightarrow v \in \mathcal{N}(f''(x))^\perp$  we write  $v$  as  $v = v_1 + v_2$  with  $v_1 \in \mathcal{N}(f''(x))$  and  $v_2 \in \mathcal{N}(f''(c))^\perp$ , noting that  $\|v_1\|_x = 0$ . Hence  $\|v\|_1^2 = \langle v_1, v_1 \rangle = \langle v, v_1 \rangle \leq \|v\|_x^* \|v_1\|_x = 0$ , and we conclude that  $v_1 = 0$ . This proves that  $v$  belongs to  $\mathcal{N}(f''(x))^\perp$ .

Given  $v \in \mathcal{N}(f''(x))^\perp$  there exists a vector  $u$  such that  $v = f''(x)u$ . We shall prove that  $\|v\|_x^* = \|u\|_x$ . From this follows that  $\|v\|_x^* < \infty$  and that  $\|\cdot\|_x^*$  is a norm on the subspace  $\mathcal{N}(f''(x))^\perp$  of  $\mathbf{R}^n$ .

Let  $w \in \mathbf{R}^n$  be arbitrary. By Cauchy–Schwarz’s inequality,

$$\begin{aligned} \langle v, w \rangle &= \langle f''(x)u, w \rangle = \langle f''(x)^{1/2}u, f''(x)^{1/2}w \rangle \\ &\leq \|f''(x)^{1/2}u\| \|f''(x)^{1/2}w\| = \|u\|_x \|v\|_x, \end{aligned}$$

and this implies that  $\|v\|_x^* \leq \|u\|_x$ . Suppose  $v \neq 0$ . Then  $u$  does not belong to  $\mathcal{N}(f''(x))$ , which means that  $\|u\|_x \neq 0$ , and for  $w = u/\|u\|_x$  we get the identity

$$\langle v, w \rangle = \|u\|_x^{-1} \langle f''(x)^{1/2}u, f''(x)^{1/2}u \rangle = \|u\|_x^{-1} \|f''(x)^{1/2}u\|^2 = \|u\|_x,$$

which proves that  $\|v\|_x^* = \|u\|_x$ . If on the other hand  $v = 0$ , then  $u$  is a vector in  $\mathcal{N}(f''(x))$  so we have  $\|v\|_x^* = \|u\|_x$  in this case, too.

- 18.3 a) Differentiate the equality  $f(tx) = f(x) - \nu \ln t$  with respect to  $x$ .  
 b) Differentiate the equality obtained in a) with respect to  $t$  and then take  $t = 1$ .  
 c) Since  $X$  does not contain any line,  $f$  is a non-degenerate self-concordant function, and it follows from the result in b) that  $x$  is the unique Newton direction of  $f$  at the point  $x$ . By differentiating the equality  $f(tx) = f(x) - \nu \ln t$  with respect to  $t$  and then putting  $t = 1$ , we obtain  $\langle f'(x), x \rangle = -\nu$ . Hence

$$\nu = -\langle f'(x), x \rangle = -\langle f'(x), \Delta x_{\text{nt}} \rangle = \lambda(f, x)^2.$$

- 18.5 Define  $g(x, x_{n+1}) = (x_1^2 + \dots + x_n^2) - x_{n+1}^2 = \|x\|^2 - x_{n+1}^2$ , so that

$$f(x) = -\ln(-g(x, x_{n+1})),$$

and let  $w = (v, v_{n+1})$ . Then

$$\begin{aligned} Dg &= Dg(x, x_{n+1})[w] = 2(\langle v, x \rangle - x_{n+1}v_{n+1}), \\ D^2g &= D^2g(x, x_{n+1})[w, w] = 2(\|v\|^2 - v_{n+1}^2), \\ D^3g &= D^3g(x, x_{n+1})[w, w, w] = 0, \\ Df &= Df(x, x_{n+1})[w] = -\frac{1}{g}Dg \\ D^2f &= D^2f(x, x_{n+1})[w, w] = \frac{1}{g^2}((Dg)^2 - gD^2g), \\ D^3f &= D^3f(x, x_{n+1})[w, w, w] = \frac{1}{g^3}(-2(Dg)^3 + 3gDgD^2g). \end{aligned}$$

Consider the difference

$$\Delta = (Dg)^2 - gD^2g = 4(\langle x, v \rangle - x_{n+1}v_{n+1})^2 + 2(x_{n+1}^2 - \|x\|^2)(\|v\|^2 - v_{n+1}^2).$$

Since  $x_{n+1} > \|x\|$ , we have  $\Delta \geq 0$  if  $|v_{n+1}| \leq \|v\|$ . So suppose that  $|v_{n+1}| > \|v\|$ . Then

$$\begin{aligned} |x_{n+1}v_{n+1} - \langle x, v \rangle| &\geq x_{n+1}|v_{n+1}| - |\langle x, v \rangle| \\ &\geq x_{n+1}|v_{n+1}| - \|x\|\|v\| \geq 0, \end{aligned}$$

and it follows that

$$\begin{aligned} \Delta &\geq 4(x_{n+1}|v_{n+1}| - \|x\|\|v\|)^2 + 2(x_{n+1}^2 - \|x\|^2)(\|v\|^2 - v_{n+1}^2) \\ &= 2(x_{n+1}|v_{n+1}| - \|x\|\|v\|)^2 + 2(x_{n+1}\|v\| - \|x\||v_{n+1}|)^2 \geq 0. \end{aligned}$$

This shows that  $D^2f = \Delta/g^2 \geq 0$ , so  $f$  is a convex function.

To prove that the function is self-concordant, we shall show that

$$4(D^2f)^3 - (D^3f)^2 \geq 0.$$



After simplification we obtain

$$4(D^2f)^3 - (D^3f)^2 = g^{-4}(D^2g)^2(3(Dg)^2 - 4gD^2g),$$

and the problem has now been reduced to showing that the difference

$$\begin{aligned}\Delta' &= 3(Dg)^2 - 4gD^2g \\ &= 12(\langle x, v \rangle - x_{n+1}v_{n+1})^2 + 8(x_{n+1}^2 - \|x\|^2)(\|v\|^2 - v_{n+1}^2)\end{aligned}$$

is nonnegative. This is obvious if  $|v_{n+1}| \leq \|v\|$ , and if  $|v_{n+1}| > \|v\|$  then we get in a similar way as above

$$\begin{aligned}\Delta' &\geq 12(x_{n+1}|v_{n+1}| - \|x\|\|v\|)^2 + 8(x_{n+1}^2 - \|x\|^2)(\|v\|^2 - v_{n+1}^2) \\ &= 4(x_{n+1}|v_{n+1}| - \|x\|\|v\|)^2 + 8(x_{n+1}\|v\| - \|x\|v_{n+1})^2 \geq 0.\end{aligned}$$

18.6 Let  $w = (u, v)$  be an arbitrary vector in  $\mathbf{R}^2$ . Writing  $a = 1/(y - f(x))$ ,  $b = -1/x$ ,  $A = f'(x)$  and  $B = f''(x)$  for short, where  $a > 0$  and  $B \geq 0$ , we obtain

$$\begin{aligned}DF(x, y)[w] &= (aA + b)u - av \\ D^2F(x, y)[w, w] &= (aB + a^2A^2 + b^2)u^2 - 2a^2Auv + a^2v^2,\end{aligned}$$

and

$$\begin{aligned}2D^2F(x, y)[w, w] - (DF(x, y)[w])^2 &= a^2A^2u^2 + b^2u^2 + a^2v^2 + 2abuv - 2a^2Auv - 2abAu^2 + 2aBu^2 \\ &= (aAu - bu - av)^2 + 2aBu^2 \geq 0.\end{aligned}$$

So  $F$  is a 2-self-concordant function.

18.7 Use the previous exercise with  $f(x) = x \ln x$ .

18.8 Taking  $f(x) = -\ln x$  in exercise 18.5, we see that

$$F(x, y) = -\ln(\ln x + y) - \ln x$$

is a 2-self-concordant barrier to the closure of the region  $-y < \ln x$ . Since  $G(x, y) = F(y, -x)$ , it then follows from Theorem 18.1.3 that  $G$  is a 2-self-concordant barrier to the region  $y \geq e^x$ .

# Index

- analytic center, 74
- Armijo's rule, 3
- barrier, 74
- central path, 76
- convergence
  - linear, 6, 7
  - quadratic, 6, 7
- damped Newton method, 23
- descent algorithm, 1
- dual local norm, 92
- gradient descent method, 2, 7
- inner iteration, 79
- input length, 114
- line search, 2
- linear convergence, 6, 7
- local seminorm, 18
- logarithmic barrier, 75
- $\nu$ -self-concordant barrier, 83
- Newton
  - decrement, 16, 35
  - direction, 15, 35
  - method, 2, 23, 66
- non-degenerate, 45
- outer iteration, 79
- path-following method, 79
- phase 1, 81
- pure Newton method, 23
- purification, 110
- quadratic convergence, 6, 7
- search direction, 2
- self-concordant, 42
- standard form, 94
- step size, 2