

Compact course notes

# COMBINATORICS AND OPTIMIZATION 466/666,

WINTER 2012

*Continuous optimization*

Lecturer: H. Wolkowicz, V. Cheung  
transcribed by: J. Lazovskis  
University of Waterloo  
April 6, 2012

---

---

## Contents

0.1	Introduction . . . . .	2
<b>1</b>	<b>Unconstrained optimization</b>	<b>2</b>
1.1	Definitions . . . . .	2
1.2	Optimality conditions . . . . .	3
1.3	Line search methods . . . . .	3
1.4	Trust region methods . . . . .	4
1.5	Conjugate gradient methods . . . . .	5
<b>2</b>	<b>Constrained optimization</b>	<b>6</b>
2.1	Feasibility and cones . . . . .	6
2.2	Convex analysis . . . . .	7
2.3	Duality . . . . .	8
2.4	Constraint qualifications . . . . .	9
2.5	Augmented Lagrangian method . . . . .	9
<b>3</b>	<b>Interior point methods</b>	<b>10</b>
3.1	Barrier functions . . . . .	10
3.2	Long-step IPM . . . . .	12
3.3	Extending IPM to SDP . . . . .	13

## 0.1 Introduction

**Definition 0.1.1.** An optimization problem is of the type

$$\begin{array}{ll} \min & f(x) \\ \text{s.t.} & c_i(x) = 0 \quad i \in \mathcal{E} \\ & c_j(x) \geq 0 \quad j \in \mathcal{I} \end{array}$$

where  $f(x)$  is the objective function, the  $c_i$  are the equality constraints, and the  $c_j$  are the inequality constraints.

- if  $\mathcal{E} = \mathcal{I} = \emptyset$ , then we have an unconstrained optimization problem
- otherwise the problem is one of constrained optimization

**Definition 0.1.2.** If not all the constraints are known at the time of formulation, a problem still can be created, based on how the model is expected to perform. In this case we call it a stochastic problem.

**Definition 0.1.3.** A set  $S \subset \mathbb{R}^n$  is termed convex if  $\lambda x + (1 - \lambda)y \in S$  for all  $x, y \in S$  and  $0 \leq \lambda \leq 1$ .

**Definition 0.1.4.** A function  $f : X \rightarrow Y$  is termed convex if  $X$  is convex and for all  $x, y \in X$  and  $\lambda \in [0, 1]$

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

Conversely, a function  $g$  is termed concave if  $-g$  is convex.

**Definition 0.1.5.** A convex optimization problem is one that has

- a convex objective function
- linear equality constraints
- concave inequality constraints

## 1 Unconstrained optimization

This is our main model that we will be using:

$$\begin{array}{ll} \min & f(x) \\ \text{s.t.} & c_i(x) = 0 \quad i \in \mathcal{E} \\ & c_j(x) \geq 0 \quad j \in \mathcal{I} \\ & x \in \Omega \subset \mathbb{R}^n \end{array} \quad \Longleftrightarrow \quad \begin{array}{ll} \min & f(x) \\ \text{s.t.} & g(x) \in K \\ & x \in \Omega \end{array}$$

Above right we have the objective function acting on the unknowns (also variables, parameters) subject to the equality, inequality, and simple constraints.

Above left,  $g(x) = \begin{pmatrix} (c_i(x))_{i \in \mathcal{E}} \\ (c_j(x))_{j \in \mathcal{I}} \end{pmatrix} \in \mathbb{R}^{m+p}$  for  $|\mathcal{E}| = m$  and  $|\mathcal{I}| = p$  where  $K$  is a cone (or is convex).

### 1.1 Definitions

**Definition 1.1.1.** A point  $x \in X$  is a global minimizer of a function  $f : X \rightarrow Y$  if  $f(x) \leq f(y) \forall y \in X$ .

**Definition 1.1.2.** A point  $x \in X$  is a local minimizer of a function  $f : X \rightarrow Y$  if there is some neighborhood  $N \ni x$  such that  $f(x) \leq f(y) \forall y \in N$ .

- To a local minimizer we may apply the adjectives weak, strict, and isolated.

**Definition 1.1.3.** Given a function  $f$ , the epigraph of  $f$  is “the region above  $f$ , i.e. the set

$$\text{epi}(f) := \{(r, x) \mid f(x) \leq r\}$$

**Remark 1.1.4.** Note that a function  $f$  is convex iff  $\text{epi}(f)$  is convex. Moreover,  $f$  being convex  $\implies f$  is locally Lipschitz  $\implies f$  is differentiable almost everywhere.

**Definition 1.1.5.** Suppose we have two sequences  $\{\eta_k\}$  and  $\{\nu_k\}$ . Then we say  
 $\{\eta_k\}$  is  $\mathcal{O}(\{\nu_k\}) \iff |\eta_k| \leq c|\nu_k|$  for all  $k$  for some constant  $c$   
 $\{\eta_k\}$  is  $o(\{\nu_k\}) \iff \frac{|\eta_k|}{|\nu_k|} \xrightarrow{k \rightarrow \infty} 0$

## 1.2 Optimality conditions

**Theorem 1.2.1.** [TAYLOR]

Taylor's theorem may be concisely stated, if  $x, p \in \mathbb{R}^n$  and  $t \in (0, 1)$ , as:

$$\begin{aligned} f(x+p) &= f(x) + \nabla f(x+tp)^T p \\ &= f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x+tp) p \\ &= f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x) p + o(\|p\|^2) \end{aligned}$$

**Definition 1.2.2.** If  $x^*$  in the domain of a function  $f$  is such that  $\nabla f(x^*) = 0$ , then  $x^*$  is termed a stationary point.

**Definition 1.2.3.** A matrix  $A \in M_{n \times n}$  is termed positive semi-definite if  $x^T A x \geq 0$  for all nonzero  $x \in \mathbb{R}^n$ . The matrix is termed positive definite if the inequality is strict.

**Theorem 1.2.4.** [FERMAT / FIRST ORDER NECESSARY OPTIMALITY]

Let  $f \in C^1$  in a neighborhood of a local minimum  $x^*$  of  $f$ . Then  $\nabla f(x^*) = 0$ .

**Theorem 1.2.5.** [SECOND ORDER NECESSARY OPTIMALITY]

Let  $f \in C^2$  in a neighborhood of a local minimum  $x^*$  of  $f$ . Then  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*) \geq 0$ .

**Theorem 1.2.6.** [SECOND ORDER SUFFICIENT OPTIMALITY]

Let  $f \in C^2$  in a neighborhood of a local minimum  $x^*$  of  $f$  with  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*) > 0$ . Then  $x^*$  is a strict local minimum.

**Theorem 1.2.7.** Consider  $x^* = \min q(x) = \frac{1}{2} x^T Q x - b^T x$ . Then the following are equivalent:

1.  $q(x)$  is bounded below
2.  $Q \geq 0, Qx = b$  is consistent
3.  $X^* = Q^{-1}b$  is a global optimum

**Remark 1.2.8.** If  $f$  is convex, then any local minimum is a global minimum. Moreover, if  $f \in C^1$ , then any stationary point is a global minimum.

## 1.3 Line search methods

**Definition 1.3.1.** Given an objective function  $f(x)$  and a starting point  $x_0$  and a search direction  $p$ , the line search method attempts to solve

$$\min_{\alpha > 0} f(x_0 + \alpha p)$$

And every next iteration is given by  $x_{k+1} = x_k + \alpha_k p_k$  for  $\alpha_k$  the step length.

**Proposition 1.3.2.** There are several descent directions that may be applied:

$$\begin{aligned} \text{steepest descent: } p_k &= -\nabla f(x_k) \\ \text{Newton's: } p_k &= -(\nabla^2 f(x_k))^{-1} \nabla f(x_k) \\ \text{quasi-Newton: } p_k &= -B_k^{-1} \nabla f(x_k) \\ \text{conjugate gradient: } p_k &= -\nabla f(x_k) + \beta_k p_{k-1} \end{aligned}$$

Note that Newton always has step length 1,  $B_k$  is some sort of approximation of  $\nabla^2 f(x_k)$ , and  $\beta_k$  ensures that  $p_k$  and  $p_{k-1}$  are conjugate.

**Definition 1.3.3.** The process of scaling is the making of the substitution  $Ay + a \rightarrow x$  in a problem.

**Theorem 1.3.4.** [WOLFE CONDITIONS]

Suppose the search direction at  $x_k$  is  $p_k$ , and  $\alpha_k \in \arg \min_{\alpha > 0} f(x_k + \alpha p_k)$  and the conditions:

- I.  $f(x_k + \alpha p_k) \leq f(x_k) + \alpha(c_1 \nabla f(x_k)^T p_k)$
- II.  $\nabla f(x_k + \alpha_k p_k)^T p_k \geq c_2 \nabla f(x_k)^T p_k$

are satisfied, where  $0 < c_1 < c_2 < 1$ . Then the search will go much faster. Sometimes we add

- III.  $|\nabla f(x_k + \alpha p_k)^T p_k| \leq c_2 |\nabla f(x_k)^T p_k|$

to replace II., which we then call the strong Wolfe conditions.

**Lemma 1.3.5.** Suppose that  $f$  is bounded below in the search direction  $p_k$  for  $f$  sufficiently smooth and  $0 < c_1 < c_2 < 1$ . Then there exist step lengths that satisfy the Wolfe conditions.

Proof: See page 35 in Nocedal & Wright.

**Theorem 1.3.6.** [ZOUTENDIJK]

Suppose that for  $\min f(x)$  with  $x_{k+1} = x_k + \alpha_k p_k$  the Wolfe conditions are satisfied, and

- $f$  is bounded below
- $f$  is  $C^1$  on a neighborhood  $N$  of  $x_0$
- $\nabla f$  is Lipschitz continuous on  $N$

Then, if  $\theta_k$  is the angle between  $p_k$  and  $-\nabla f(x_k)$ ,

$$\sum_{k \geq 0} \cos^2(\theta_k) \|\nabla f(x_k)\|^2 < \infty$$

**Remark 1.3.7.** The above, with some manipulation, implies that  $\lim_{k \rightarrow \infty} [\nabla f(x_k)] = 0$ .

**Definition 1.3.8.** For  $Q \in M_{n \times n}$ , define the weighted inner product  $\|\cdot\|_Q$  by  $\|x\|_Q^2 := x^T Q x$  for  $x \in \mathbb{R}^n$ .

**Lemma 1.3.9.** [KANTOROVICH]

Let  $Q \in M_{n \times n}$  with  $Q = Q^T > 0$  and  $x \in \mathbb{R}^n$ . Then

$$\frac{(x^T x)^2}{x^T Q x x^T Q^{-1} x} \geq \frac{4\lambda_{\min}(Q)\lambda_{\max}(Q)}{(\lambda_{\min}(Q) + \lambda_{\max}(Q))^2}$$

**Theorem 1.3.10.** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a  $C^2$  function. Suppose we apply exact line search to generate a sequence  $(x_k)$  with  $x_k \xrightarrow{k \rightarrow \infty} x^*$ . Moreover, suppose  $\nabla^2 f(x^*) > 0$  and  $\nabla f(x^*) = 0$ . Then

$$f(x_{k+1}) - f(x^*) \leq \left( \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2 (f(x_k) - f(x^*))$$

for  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  the eigenvalues of  $\nabla^2 f(x^*)$ .

## 1.4 Trust region methods

**Definition 1.4.1.** Given an objective function  $f(x)$ , a starting point  $x_0$  and a model  $m_k(x)$  of  $f$  around  $x_0$ , the trust region method attempts to solve

$$\min_p m_k(x_0 + p)$$

such that  $x_0 + p$  always lies inside some predefined trust region.

**Definition 1.4.2.** The trust region subproblem (TRS), for  $B_k \approx \nabla^2 f(x_k)$  is given by

$$\begin{aligned} \min \quad & f(x_k) + \nabla f(x_k)^T p_k + \frac{1}{2} p_k^T B_k p_k \\ \text{s.t.} \quad & \|p_k\| \leq \Delta_k \end{aligned}$$

This is quadratic minimization with one constraint, where we minimize over  $p_k$ .

**Definition 1.4.3.** Define the actual reduction and the predicted reduction in the ratio

$$\rho_k = \frac{f(x_k) - f(x_k + p_k)}{m_k(0) - m_k(p_k)}$$

Note that if  $\rho_k < 0$ , then the trust region is too large, and must be decreased. If  $\rho_k \approx 1$ , we may increase the trust region size safely.

**Theorem 1.4.4.** The point  $p^*$  solves the TRS iff there exists a  $\lambda$  such that

$$\left. \begin{aligned} (B_k + \lambda I)p^* &= -\nabla f(x_k)\lambda \\ B_k + \lambda I &\geq 0 \\ \|p^*\| &\leq \Delta_k \\ \lambda(\|\Delta^*\| - \Delta_k) &= 0 \end{aligned} \right\} \begin{array}{l} \text{dual feasibility} \\ \text{primal feasibility} \\ \text{complementary slackness} \end{array} \left. \vphantom{\begin{aligned} (B_k + \lambda I)p^* &= -\nabla f(x_k)\lambda \\ B_k + \lambda I &\geq 0 \\ \|p^*\| &\leq \Delta_k \\ \lambda(\|\Delta^*\| - \Delta_k) &= 0 \end{aligned}} \right\} \text{modern paradigm}$$

where  $\Delta_k \in (0, \Delta^*)$  for all  $k$ .

## 1.5 Conjugate gradient methods

**Definition 1.5.1.** A set of nonzero vectors  $\{v_0, \dots, v_n\}$  is termed conjugate wrt  $A$  if  $v_i^T A v_j = 0$  iff  $i \neq j$ .

**Lemma 1.5.2.** A conjugate set is linearly independent.

**Definition 1.5.3.** Suppose we begin with  $A \in M_{n \times n}$  and a problem

$$\min_x \frac{1}{2} x^T A x - b^T x = \varphi(x)$$

Then with a given set of conjugate vectors  $\{p_1, \dots, p_n\}$ , we solve

$$\min_{\alpha} \varphi(x_k + \alpha p_k)$$

This is termed the conjugate gradient method (CG).

**Definition 1.5.4.** With respect to the above, the expression  $\nabla \varphi(x) = r(x) = Ax - b$  is termed the residue.

**Theorem 1.5.5.** For any starting point, the conjugate gradient method converges in at most  $n$  steps.

**Proposition 1.5.6.** Recall, from above, that to find the set of conjugate direction vectors, we use the calculation  $p_k = -\nabla f(x_k) + \beta_k p_{k-1}$ . To make the CG method practical, we use

$$\alpha_k = \frac{r_k^T r_k}{p_k^T A p_k} \quad \beta_k = \frac{r_{k+1}^T r_k}{r_k^T r_k}$$

**Theorem 1.5.7.** If  $A$  has at most  $m$  distinct eigenvalues, then the CG method converges in at most  $m$  iterations.

## 2 Constrained optimization

**Remark 2.0.1.** Consider smooth functions  $f_1, \dots, f_n$  in an optimization problem

$$\min \quad \max\{f_1(x), \dots, f_n(x)\}$$

Such a problem may not have a smooth objective function. Then we may reformulate this equivalently as

$$\begin{aligned} \min \quad & t \\ \text{s.t.} \quad & t \geq f_1(x) \\ & \vdots \\ & t \geq f_n(x) \end{aligned}$$

which now has a smooth objective function and smooth constraints.

### 2.1 Feasibility and cones

**Definition 2.1.1.** Consider an optimization problem

$$\begin{aligned} \min \quad & f(x) \\ \text{s.t.} \quad & c_i(x) = 0 \quad \forall i \in \mathcal{E} \\ & c_k(x) \geq 0 \quad \forall k \in \mathcal{I} \\ & x \in \Omega \end{aligned}$$

Define the set of linearized feasible directions and the active set by

$$\begin{aligned} \mathcal{F} &:= \left\{ d \mid \begin{array}{l} \nabla c_i(x)^T d = 0 \quad \forall i \in \mathcal{E} \\ \nabla c_k(x)^T d \geq 0 \quad \forall i \in \mathcal{I} \cap \mathcal{A}(x) \end{array}, x \in \Omega \right\} \\ \mathcal{A}(x) &:= \mathcal{E} \cup \{i \in \mathcal{I} \mid c_i(x) = 0\} \end{aligned}$$

**Definition 2.1.2.** Let  $\Omega \subset \mathbb{R}^n$  and  $\bar{x} \in \Omega$ . Then a quantity  $d$  is termed a feasible direction to  $\Omega$  at  $\bar{x}$  iff there exists  $\bar{\alpha} > 0$  such that  $\bar{x} + \alpha d \in \Omega$  for all  $0 \leq \alpha < \bar{\alpha}$ .

**Definition 2.1.3.** Let  $\Omega \subset \mathbb{R}^n$  and  $\bar{x} \in \Omega$ . The tangent cone of  $\Omega$  at  $\bar{x}$  is defined to be

$$\begin{aligned} T(\Omega, \bar{x}) &:= \left\{ \alpha d \mid \exists (x_k)_{k=1}^\infty \subset \Omega \text{ s.t. } x_k \rightarrow \bar{x} \text{ and } d = \lim_{k \rightarrow \infty} \left[ \frac{x_k - \bar{x}}{\|x_k - \bar{x}\|} \right], \alpha \in \mathbb{R}_{\geq 0} \right\} \\ &= \overline{\text{cone}}(\Omega - \bar{x}) \end{aligned}$$

This is also termed the cone of limiting feasible directions.

· We note that  $T(\Omega, x)$  is always closed, and is convex if  $\Omega$  is convex.

**Definition 2.1.4.** Let  $\Omega \subset \mathbb{R}^n$  and  $x \in \Omega$ . The normal cone of  $\Omega$  at  $x$  is defined to be the set

$$N(\Omega, x) := \{v \mid \langle v, d \rangle \leq 0 \quad \forall d \in T(\Omega, x)\}$$

**Definition 2.1.5.** Let  $\Omega \subset \mathbb{R}^n$ . The polar cone of  $\Omega$  is defined to be the set

$$\Omega^+ := \{v \mid \langle v, d \rangle \geq 0 \quad \forall d \in \Omega\}$$

· In  $\mathbb{R}^n$ ,  $\langle v, d \rangle = v^T d$ , and we know  $\cos(\theta) = \frac{v^T d}{\|v\| \|d\|}$ .

**Definition 2.1.6.** Consider an optimization problem with local solution  $x^*$  and associated Lagrange multiplier  $\lambda^*$ . The critical cone of  $x^*$  with  $\lambda^*$  is defined to be the set

$$\mathcal{C}(x^*, \lambda^*) := \{w \in \mathcal{F}(x^*) \mid \nabla c_i(x^*)^T w = 0 \quad \forall i \in \mathcal{I} \cap \mathcal{A}(x^*) \text{ with } \lambda_i^* > 0\}$$

**Definition 2.1.7.** The set  $X \subset \mathbb{R}^n$  is termed an orthant iff it is the intersection of  $n$  pairwise orthogonal half-spaces of  $\mathbb{R}^n$ .

**Theorem 2.1.8.** Given a minimization problem as above,  $\bar{x} \in \arg \min_{x \in \mathcal{F}} f(x)$  implies  $\nabla f(\bar{x}) \in T(\mathcal{F}, \bar{x})$ .

*Proof:* Suppose the premise holds but not the conclusion.

Then there exists a  $d \in T(\mathcal{F}, \bar{x})$  such that  $\nabla f(\bar{x})^T d < 0$  for  $d = \lim_{k \rightarrow \infty} \left[ \frac{x_k - \bar{x}}{\|x_k - \bar{x}\|} \right]$ .

Then for some  $K \in \mathbb{N}$ , we have  $\nabla f(\bar{x})^T (x_k - \bar{x}) < 0$  for all  $k \geq K$ .

But then  $f(x_k) = f(\bar{x}) + \underbrace{\nabla f(\bar{x})^T (x_k - \bar{x})}_{<0} + o(\|x_k - \bar{x}\|)$ , implying  $f(x_k) < f(\bar{x})$ .

This is a contradiction. ■

**Corollary 2.1.9.** [FERMAT]

If  $\bar{x} \in \text{int}(\mathcal{F})$  and  $\bar{x} \in \arg \min_{x \in \mathcal{F}} f(x)$ , then  $\nabla f(\bar{x}) = 0$ .

**Corollary 2.1.10.** If  $f$  is a convex function and  $\mathcal{F}$  is a convex set, then

$$\bar{x} \in \arg \min_{x \in \mathcal{F}} f(x) \quad \text{iff} \quad \nabla f(\bar{x}) \in T(\mathcal{F}, \bar{x})^+ \quad \text{iff} \quad \nabla f(\bar{x}) \in (\mathcal{F} - \bar{x})^+$$

**Theorem 2.1.11.** [ROCKAFELLAR, PSHENICHNY]

If  $\bar{x} \in \text{int}(\mathcal{F})$  then  $\nabla f(\bar{x}) \in T(\mathcal{F}, \bar{x})^+$ .

## 2.2 Convex analysis

**Lemma 2.2.1.** Let  $K \neq \emptyset$  be a closed, bounded, convex set with  $\bar{x} \notin K$ . Then there exists a unique  $\bar{y} \in K$  with  $\bar{y} \in \arg \min_{y \in K} \{\|y - \bar{x}\|\}$ .

**Definition 2.2.2.** Let  $K \subset \mathbb{R}^n$  nontrivial. Then  $K$  is termed a cone if  $\alpha K \subset K$  for all  $\alpha \geq 0$ .  $K$  is termed a convex cone if it is a cone, and  $K + K \subset K$ .

**Definition 2.2.3.** A cone  $K$  is termed pointed iff  $K \cap -K = \{0\}$ .

**Definition 2.2.4.** A cone  $K$  is termed self-dual or self-polar if  $K^+ = K$ .

**Definition 2.2.5.** Let  $K$  be a cone. Then  $K = K^{++}$  iff  $K$  is a closed convex cone.

**Lemma 2.2.6.** [FARKAS]

Let  $A \in M_{n \times n}$ . Then equivalently

- I.  $Ax = b, x \geq 0$  is consistent
- II.  $A^T y \geq 0$  implies  $b^T y \geq 0$

**Theorem 2.2.7.** Suppose that for two convex sets  $C_1, C_2$  we have  $C_1 \cap \text{int}(C_2) = \emptyset$ . Then we can separate the two sets by a hyperplane.

**Definition 2.2.8.** Given a problem with

$$\begin{aligned} Ax &= b \\ bx &\leq d \\ g(x) &\leq 0 \end{aligned}$$

linear constraints and  $g$  convex, the generalized Slater CQ is

$$\text{there exists } \hat{x} \text{ such that } A\hat{x} = b, \quad B\hat{x} \leq d, \quad g(\hat{x}) = 0$$

**Remark 2.2.9.** The GSCQ implies the weakest CQ.

## 2.3 Duality

**Definition 2.3.1.** Given convex cones  $K, L$  define the primal and dual problems to be

$$\begin{aligned} p^* &= \min_{\substack{\langle c, x \rangle \\ \text{s.t. } Ax \geq_K b \\ x \geq_L 0}} & \alpha^* &= \max_{\substack{\langle b, y \rangle \\ \text{s.t. } A^*y \leq_{L^+} c \\ y \geq_{K^+} 0}} \end{aligned}$$

where  $A^*$  is found through the adjoint linear transformation, i.e.  $\langle A^*y, x \rangle = \langle y, Ax \rangle$  for all  $x, y$ .

**Definition 2.3.2.** Suppose that  $f$  is convex,  $g$  is  $K$ -convex for  $K$  a closed convex cone and  $C$  is convex. Then for the problem

$$\begin{aligned} \min & f(x) \\ \text{s.t. } & g(x) \leq_K 0 \\ & x \in C \end{aligned}$$

define the perturbation function  $w(\varepsilon) = \min_{x \in C} \{f(x) \mid g(x) \leq_K \varepsilon\}$  for  $g : X \rightarrow Y$ .

**Proposition 2.3.3.** Let  $\varepsilon \in \mathbb{R}^n$  and  $K = \mathbb{R}_+^m$  with  $\Gamma = \{\varepsilon \in Y \mid \text{there exists } x \in C \text{ with } g(x) \leq_K \varepsilon\}$ . Then

- i.  $\Gamma$  is a convex set
- ii.  $w(\varepsilon)$  is a convex function on its domain (where it is finite)
- iii.  $w$  is non-increasing in  $\varepsilon$

**Theorem 2.3.4.** Suppose there exists  $\hat{x} \in C$  with  $g(\hat{x}) > 0$  (that is,  $g(\hat{x}) \in -\text{int}(K)$ ), so SCQ is satisfied. If  $w(0)$  is finite, then there exists equivalently

- an optimal Lagrange multiplier
- $\lambda^* \geq_K 0$  such that  $w(0) = \min_{x \in C} \{f(x) + \langle \lambda^*, g(x) \rangle\}$

Moreover, if the minimum is attained at  $x^* \in C$  and  $x^*$  is feasible (that is,  $g(x^*) \leq_K 0$ ), then  $x^*$  solves the convex program, and  $\langle \lambda^*, g(x^*) \rangle = 0$ .

**Remark 2.3.5.** Suppose that we have a *nonlinear problem*

$$\begin{aligned} \min_x & f(x) \\ \text{s.t. } & c_i(x) = 0 \quad \forall i \in \mathcal{E} \\ & c_k(x) \geq 0 \quad \forall k \in \mathcal{I} \\ & x \in \mathbb{R}^n \end{aligned}$$

To provide a lower bound on the optimal solution, we form the *Lagrangian*

$$\mathcal{L}(x, \lambda) = f(x) - \sum_{i \in \mathcal{E}, \mathcal{I}} \lambda_i c_i(x) \quad \text{for } \lambda_i \geq 0 \text{ if } i \in \mathcal{I}$$

Then we define the *dual functional*

$$g(\lambda) = \inf_x \{\mathcal{L}(x, \lambda)\}$$

And finally we have the *dual problem*

$$\begin{aligned} \max_{\lambda} & g(\lambda) \\ \text{s.t. } & \lambda \geq 0 \quad \forall i \in \mathcal{I} \end{aligned}$$

**Proposition 2.3.6.** If the weakest constraint qualifications (WCQ) hold at  $x^*$  (that is,  $T(\Omega, x^*) = \mathcal{F}(x^*)$ ), then the KKT conditions hold at  $x^*$ , that is, there exists a  $\lambda^*$  such that

$$\begin{aligned} \nabla_x \mathcal{L}(x^*, \lambda^*) &= 0 & \lambda_i &= 0 \quad \forall i \in \mathcal{I} & \text{(dual feasibility)} \\ c_i(x^*) &= 0 & \forall i &\in \mathcal{E} & \text{(primal feasibility)} \\ c_j(x^*) &\geq 0 & \forall j &\in \mathcal{I} \\ \lambda_i^* c_i(x^*) &= 0 & \forall i &\in \mathcal{E}, \mathcal{I} & \text{(complementary slackness)} \\ \lambda_j^* &\geq 0 & \forall j &\in \mathcal{I} \end{aligned}$$



**Definition 2.3.7.** The strict complementarity conditions hold at  $x^*$  if the KKT conditions hold with for some Lagrangian multiplier  $\lambda^*$  such that  $\lambda_j^* > 0$  for all  $j \in \mathcal{I}$ .

**Proposition 2.3.8.** Let  $\Omega \subset \mathbb{R}^n$  and  $x^* \in \Omega$ . If  $N(\Omega, x^*) = -\mathcal{F}(x^*)^+$ , then  $\Omega$  is convex and WCQ holds.

## 2.4 Constraint qualifications

There are several main constraint qualifications:

LICQ	holds at $x \in \Omega$ if	$\{\nabla c_i(x) \mid i \in A(x)\}$ is linearly independent
MFCQ	holds at $x \in \Omega$ if	<ol style="list-style-type: none"> <li>1. there exists <math>w</math> such that <math>\nabla c_i(x)^T w = 0</math> for all <math>i \in \mathcal{E}</math></li> <li>2. <math>\{\nabla c_i(x) \mid i \in \mathcal{E}\}</math> is linearly independent</li> </ol>
WCQ	holds at $x \in \Omega$ if	all constraints are linear

We note that  $\text{LICQ} \implies \text{MFCQ} \implies \text{WCQ}$ .

**Corollary 2.4.1.** If LICQ holds, then none of the active constraint gradients can be zero.

**Definition 2.4.2.** Given a feasible point  $x^*$  in an optimization problem, a sequence  $(z_k)_{k=1}^\infty$  is termed a feasible sequence approaching  $x^*$  iff  $z_k \xrightarrow{k \rightarrow \infty} x^*$  and  $z_k \in \Omega$  for all  $k$ .

· We note that for any  $x^*$  feasible, the inclusion  $T(\Omega, x^*) \subset \mathcal{F}(x^*)$  always holds.

**Lemma 2.4.3.** Let  $x^* \in \Omega$  and LICQ holds at  $x^*$ . Let  $d \in \mathcal{F}(x^*)$ . Then for all  $t_k > 0$  with  $t_k \xrightarrow{k \rightarrow \infty} 0$ , there exists  $(z_k)_{k=1}^\infty$  such that

1.  $z_k \in \Omega$  for all  $k$
2.  $z_k \xrightarrow{k \rightarrow \infty} x^*$
3.  $d = \lim_{k \rightarrow \infty} \left[ \frac{z_k - x^*}{t_k} \right]$
4.  $c_i(z_k) = t_k \nabla c_i(x^*)^T d$  for all  $i \in A(x^*)$

**Corollary 2.4.4.**  $\text{LICQ} \implies \text{WCQ}$ , i.e.  $T(\Omega, x^*) = \mathcal{F}(x^*)$ .

## 2.5 Augmented Lagrangian method

**Definition 2.5.1.** For an equality-constrained non-linear problem, define the augmented Lagrangian to be the equation

$$\begin{aligned} \mathcal{L}_A(x, \lambda, \mu) &= f(x) - \sum_{i \in \mathcal{E}} \lambda_i c_i(x) + \frac{\mu}{2} \sum_{i \in \mathcal{E}} (c_i(x))^2 \\ &= \mathcal{L}(x, \lambda) + \frac{\mu}{2} \sum_{i \in \mathcal{E}} (c_i(x))^2 \end{aligned}$$

Augmenting the Lagrangian proves useful, as we may adjust  $\mu$  as desired. Moreover, first order conditions are unchanged from the original problem, as

$$\begin{aligned} \nabla_x \mathcal{L}_A(x, \lambda, \mu) &= \nabla_x \mathcal{L}(x, \lambda) + \mu \sum_{i \in \mathcal{E}} c_i(x) \nabla c_i(x) \\ &= \nabla_x f(x) + \sum_{i \in \mathcal{E}} (\mu c_i(x) - \lambda_i) \nabla c_i(x) \end{aligned}$$

**Theorem 2.5.2.** Suppose that  $x^*$  is a local solution to (ECNLP) that satisfies the KKT conditions with Lagrangian multiplier  $\lambda^*$ , as well as the second order sufficiency conditions. Then there exists a  $\mu_0 \in \mathbb{R}_{\geq 0}$  such that for all  $\mu \geq \mu_0$ ,  $x^*$  solves the problem

$$\min_x \mathcal{L}_A(x, \lambda^*, \mu)$$

**Theorem 2.5.3.** For fixed  $\mu$ , if  $\frac{\|\lambda - \lambda^*\|}{\mu}$  is small, then the following method works for solving (ECNLP).

1.  $x \leftarrow \arg \min \mathcal{L}_A(\cdot, \lambda, \mu)$
2.  $\lambda \leftarrow \lambda - \mu c(x)$

In this case a solution  $x^*$  to (ECNLP) will be the limit of the iterates under the above instructions.

We may update  $\mu$  in the following fashion. Choose  $\beta > 1$  not too large and not too small such that convergence is not too slow and the problem does not degenerate. Then set

$$\mu_{k+1} \leftarrow \begin{cases} \beta \mu_k & \text{if } \|c(x_k)\| > \gamma \|c(x_{k-1})\| \text{ for some fixed } \gamma \in (0, 1) \\ \mu_k & \text{else} \end{cases}$$

### 3 Interior point methods

To the classical optimization problem so far we have seen several approaches:

1. merit functions
2. quadratic penalty methods
3.  $\ell^2$  penalty method
4. augmented Lagrangian method

All these deal with so-called “exterior point methods,” which involve approaching the feasible region from the outside.

#### 3.1 Barrier functions

**Definition 3.1.1.** In an optimization problem, a barrier function is added to the objective function to prevent it from going near the boundary of the feasible region. We may define such a function in several ways:

$$\begin{aligned} \cdot \text{ inverse barrier function: } \hat{B}(x) &:= \begin{cases} \sum_{i \in \mathcal{I}} \frac{1}{c_i(x)} & \text{if } c_i(x) > 0 \ \forall i \in \mathcal{I} \\ \infty & \text{else} \end{cases} \\ \cdot \text{ log barrier function: } B(x) &:= \begin{cases} \sum_{i \in \mathcal{I}} \log(c_i(x)) & \text{if } c_i(x) > 0 \ \forall i \in \mathcal{I} \\ \infty & \text{else} \end{cases} \end{aligned}$$

Note that the log barrier function, while extended-real valued, is continuous.

The general barrier method algorithm works as follows:

- start with  $\mu_0 > 0$
- for  $k = 0, 1, \dots$ 
  - find  $x_k \in \arg \min_x \{f(x) + \mu_k B(x) \mid c_i(x) = 0 \ \forall i \in \mathcal{E}\}$
  - choose  $\mu_{k+1} \in (0, \mu_k)$
- end

**Proposition 3.1.2.** Let  $\Omega = \text{cl}\left\{x \mid \begin{matrix} c_i(x) = 0 & \forall i \in \mathcal{E} \\ c_j(x) > 0 & \forall j \in \mathcal{I} \end{matrix}\right\} = \text{cl}(\hat{\Omega})$ . Then every limit point  $\bar{x}$  of  $(x_k)_{k \in \mathbb{N}}$  generated by the general barrier method with  $(\mu_k)_{k \in \mathbb{N}}$  and  $\mu_k \xrightarrow{k \rightarrow \infty} 0$  is a global solution of (NLP).

Proof: Let  $y \in \Omega$ ,

Then there exists a sequence  $(y_\ell)_{\ell=1}^\infty \subset \hat{\Omega}$  such that  $y_\ell \xrightarrow{\ell \rightarrow \infty} y$ ,  
This implies that for all  $\ell$  and for all  $k$ ,

$$f(x_k) + \mu_k \beta(x_k) \leq f(y_\ell) + \mu_k \beta(y_\ell)$$

Taking the limit as  $k \rightarrow \infty$ ,

$$f(\bar{x}) \leq \lim_{k \rightarrow \infty} [f(x_k) + \mu_k \beta(x_k)] \leq f(y_\ell)$$

Then taking the limit as  $\ell \rightarrow \infty$ , we find that  $f(\bar{x}) \leq f(y)$ . ■

**Theorem 3.1.3.** [FUNDAMENTAL THEOREM OF LPs]

A linear program (LP) is exactly one of the following:

1. infeasible
2. unbounded
3. solvable (implying strong duality)

The (LP) may be solved by the interior point method with

$$(p_\mu) \quad \min_x \quad f_p(x) = c^T x - \mu \sum_{i=1}^n \log(x_i) \\ \text{s.t.} \quad Ax = b$$

**Remark 3.1.4.** For such a problem,  $(p_\mu)$  has a unique solution  $x(\mu)$  for each  $\mu > 0$  if  $Ax = b$  is consistent.

**Definition 3.1.5.** The set  $\{x(\mu) \mid \mu > 0\}$  is termed the central path. The analytic center of the set of optimal solutions is defined by

$$x_\infty := \arg \min_x \left\{ -\sum_{i=1}^n \log(x_i) \mid Ax = b, x \geq 0 \right\}$$

The general primal/dual interior point method works as follows:

· Initialize:

$$\begin{aligned} x^0 &> 0 \\ y^0 & \\ s^0 &> 0 \\ 0 &\leq \sigma_{\min} < \sigma_{\max} \leq 1 \\ \text{tol} &> 0 \end{aligned}$$

· for  $k = 0, 1, \dots$ :

- $\mu_k \leftarrow (x_k^T s_k)/n$
- if  $\mu_k < \text{tol}$  and  $\|r_d\| < \text{tol}$  and  $\|r_p\| < \text{tol}$ :  
break
- else:
- pick  $\sigma_k \in [\sigma_{\min}, \sigma_{\max}]$
- solve the system

$$\begin{bmatrix} 0 & A^T & I \\ A & 0 & 0 \\ S & 0 & X \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta y \\ \Delta s \end{bmatrix} = \begin{bmatrix} -r_d \\ -r_p \\ -x_k \circ s_k + \sigma_k \mu_k e \end{bmatrix}$$

where  $S = \text{diag}(s)$  and  $X = \text{diag}(x)$   
 · pick  $\alpha_k \in (0, 1)$  such that

$$(x_{k+1}, y_{k+1}, s_{k+1}) = (k_k + \alpha \Delta x, y_k + \alpha \Delta y, s_k + \alpha \Delta s)$$

satisfying  $x_{k+1} > 0$  and  $s_{k+1} > 0$  and the centrality criterion

· end

The centrality criterion is a restriction on  $\theta \in [0, 1)$  and  $\gamma \in (0, 1]$  so that every iterate is not be too far from the central path  $C$  in terms of some neighborhood of  $C$ . We define the centrality measure to be  $\mu = x^T s / n$ . We speak in terms of the neighborhoods

$$\begin{aligned} N_2(\theta) &= \{(x, y, s) \mid A^T y + s = c, Ax = b, x > 0, s > 0, \|x \circ s - \mu e\|_2 \leq \theta \mu\} \\ N_{-\infty}(\gamma) &= \{(x, y, s) \mid A^T y + s = c, Ax = b, x > 0, s > 0, x_i s_i \geq \gamma \mu \text{ for all } i\} \end{aligned}$$

Then if  $\theta \in (0, 1)$  for all  $(x, y, s) \in N_2(\theta)$ , we will have that for all  $i$ ,

$$|x_i s_i - \mu| \leq \|x \circ s - \mu e\|_2 \leq \theta \mu \implies x_i s_i \geq \mu - \theta \mu = (1 - \theta) \mu \quad \text{and} \quad N_2(\theta) \subset N_{-\infty}(1 - \theta)$$

### 3.2 Long-step IPM

**Remark 3.2.1.** For all  $\varepsilon > 0$  with initial duality measure  $\mu_0$ , the long-step interior point method takes  $k = \mathcal{O}(n |\log(\varepsilon)|)$  steps to reduce the duality measure by a factor of  $\varepsilon$ , i.e. to find  $(x_k, y_k, s_k)$  such that

$$\mu_k = \frac{x_k^T s_k}{n} \leq \varepsilon \mu_0$$

**Lemma 3.2.2.** For all  $u, v \in \mathbb{R}^n$  with  $u^T v \geq 0$ ,

$$\|u \circ v\|_2 \leq 2^{-3/2} \|u + v\|_2^2$$

where  $\circ$  is the Hadamard product, for which  $(u \circ v)_i = u_i v_i$ .

**Lemma 3.2.3.** If  $(x, y, s) \in N_{-\infty}(\gamma)$  for  $\gamma \in (0, 1]$  fixed, and  $(\Delta x, \Delta y, \Delta s)$  solves

$$\begin{bmatrix} 0 & A^T & I \\ A & 0 & 0 \\ S & 0 & X \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta y \\ \Delta s \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ -x \circ s + \sigma \mu e \end{bmatrix}$$

then we have the following three results:

1.  $\|\Delta x \circ \Delta s\|_2 \leq 2^{-3/2} \left(1 + \frac{1}{\gamma}\right) n \mu$
2.  $\Delta x^T \Delta s = 0$
3.  $(x(\alpha), y(\alpha), s(\alpha)) \in N_{-\infty}(\gamma)$  where

$$\begin{aligned} x(\alpha) &= x + \alpha \Delta x \\ y(\alpha) &= y + \alpha \Delta y \\ s(\alpha) &= s + \alpha \Delta s \end{aligned} \quad \text{for any} \quad \alpha \in \left[0, 2^{3/2} \cdot \frac{\gamma(1-\gamma)}{1+\gamma} \cdot \frac{\sigma}{n}\right]$$

**Theorem 3.2.4.** Given  $\gamma, \sigma_{\min}, \sigma_{\max}$  in the long-step IPM path for each  $k$ , setting

$$\alpha_k = 2^{3/2} \cdot \frac{\gamma(1-\gamma)}{1+\gamma} \cdot \frac{\sigma_k}{n}$$

there exists  $\delta > 0$ , independent of  $n$ , such that

$$\mu_{k+1} \leq \left(1 - \frac{\delta}{n}\right) \mu_k$$

Proof: Performing routine calculations we get the result.

$$\begin{aligned}
\mu_{k+1} &= \mu_k \alpha_k \\
&= \left( 1 - 2^{3/2} \frac{\gamma(1-\gamma)}{1+\gamma} \cdot \frac{1}{n} \cdot \sigma(1-\sigma) \right) \mu_k \\
&\leq \left( 1 - \frac{1}{n} \cdot \underbrace{2^{3/2} \cdot \frac{\gamma(1-\gamma)}{1+\gamma} \cdot \min \left\{ \frac{\sigma_{\min}(1-\sigma_{\min})}{\sigma_{\max}(1-\sigma_{\max})} \right\}}_{\sim \delta} \right) \mu_k
\end{aligned}$$

■

**Theorem 3.2.5.** Fix  $\varepsilon \in [0, 1]$ ,  $\gamma \in (0, 1)$ ,  $0 \leq \delta_{\min} \leq \delta_{\max} \leq 1$ , an initial point  $(x_0, y_0, s_0) \in N_{-\infty}(\gamma)$ . Then for  $\delta$  as in the above proof,

$$\mu_k \leq \varepsilon \mu_0 \quad \text{for all} \quad k \geq \frac{\delta}{n} |\log(\varepsilon)|$$

### 3.3 Extending IPM to SDP

**Definition 3.3.1.** In graph theory, a common problem is the max-cut problem. Given a graph  $G = (V, E)$  with weighted edges  $e \in E$ , what is the cut of maximum size?

- a cut of  $G$  is a partition  $\{U_1, U_2\}$  of  $V$  such that  $U_1 \cup U_2 = V$
- the size of a cut  $\{U_1, U_2\}$  is the sum of edge weights of edges that are not completely within  $U_1$  or  $U_2$

Here we will consider the problem with unweighted edges, that is, where all edges have an equal weight of 1.

**Definition 3.3.2.** For  $A$  the adjacency matrix of  $G$ , define the Laplacian matrix of  $G$  to be

$$L = \text{diag}(Ae) - A$$

This matrix is positive semi-definite (PSD) and singular.

Now we may formulate the max-cut problem in an optimization manner. Here the vector  $x$  is basically the set of vertices  $V$  of  $G$  arranged in a vector.

$$\begin{aligned}
&\max_x \quad \frac{1}{2} x^T L x \\
&\text{s.t.} \quad x_i^2 = 1 \quad \text{for all } i
\end{aligned}$$

The SDP relaxation of this problem is given by

$$\begin{aligned}
&\max_X \quad \langle L, X \rangle \\
&\text{s.t.} \quad \text{diag}(X) = e \\
&\quad \quad X \geq 0
\end{aligned}$$

The dual to the original problem is given by

$$\begin{aligned}
&\min_{\lambda} \quad \lambda^T e \\
&\text{s.t.} \quad \text{diag}(\lambda) - \frac{1}{4} L \geq 0
\end{aligned}$$

**Remark 3.3.3.** In a more general setting, an SDP program and its dual are given by

$$\begin{aligned}
(P) \quad &\inf_X \quad \langle C, X \rangle \\
&\text{s.t.} \quad \langle A_i, X \rangle = b_i \text{ for all } i \\
&\quad \quad X \geq 0 \\
(D) \quad &\sup_y \quad b^T y \\
&\text{s.t.} \quad C - \sum_i y_i A_i \geq 0
\end{aligned}$$