

Notes for ECNU 2007 Short Course

October 2007

0. Introduction

- Purpose of the course: to offer the students a concise introduction to modern unconstrained and constrained optimization methods, algorithms, and computer software; functional optimization and applications will also be briefly addressed.
- Background mathematics includes first-year calculus and linear algebra, some knowledge of calculus of variations and PDE are helpful but not critical.
- Familiarity with MATLAB is important.
- Most of the distributed course material (Chapters 1, 2, 4, 5, 7, 9, 10, 12, 13, 14, 16, and Appendix A) are from the book

Andreas Antoniou and Wu-Sheng Lu, *Practical Optimization: Algorithms and Engineering Applications*, Springer, New York, March 2007.

The book has a web site: <http://www.ece.uvic.ca/~optimization/>

I. Numerical Optimization and Applications

1.1 Unconstrained Optimization (Chaps. 1, 2, 4, 5, 7, 9)

- Basic optimization problem:

$$\text{minimize } f(\mathbf{x}) \quad \text{for } \mathbf{x} \in R^n$$

- Assumption: $f(\mathbf{x})$ is twice continuously differentiable.
- An example: solve a system of algebraic equations $p_i(\mathbf{x}) = 0$ for $i = 1, 2, \dots, m$.
- Basic structure of a numerical algorithm for minimizing $f(\mathbf{x})$
 - Step 1: choose an initial point \mathbf{x}_0 , set a convergence tolerance ε , and set a counter $k = 0$.
 - Step 2: determine a search direction \mathbf{d}_k for reducing $f(\mathbf{x})$ from point \mathbf{x}_k .
 - Step 3: determine a step size α_k such that $f(\mathbf{x}_k + \alpha \mathbf{d}_k)$ is minimized for $\alpha \geq 0$, and construct $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$
 - Step 4: if $\|\alpha_k \mathbf{d}_k\| < \varepsilon$, stop and output a solution \mathbf{x}_{k+1} , else set $k := k + 1$ and repeat from Step 2.

- Comments:

- (a) Steps 2 and 3 are key steps of an optimization algorithm.
- (b) Different ways to accomplish Step 2 leads to different algorithms
- (c) Step 3 is a one-dimensional optimization problem and it is often called a line search step.

- Notation and basic quantities: \mathbf{x} : vector of n variables in column

$f(\mathbf{x})$: objective function

$\mathbf{g}(\mathbf{x})$: gradient (vector) of $f(\mathbf{x})$. A popular notation for $\mathbf{g}(\mathbf{x})$: $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x})$

$\mathbf{H}(\mathbf{x})$: Hessian (matrix) of $f(\mathbf{x})$. A popular notation for $\mathbf{H}(\mathbf{x})$: $\mathbf{H}(\mathbf{x}) = \nabla^2 f(\mathbf{x})$

Note: $\mathbf{H}(\mathbf{x})$ is a square, symmetric matrix and can be obtained as $\mathbf{H}(\mathbf{x}) = \nabla(\nabla^T f(\mathbf{x}))$

A point \mathbf{x} is called a stationary point if $\mathbf{g}(\mathbf{x}) = \mathbf{0}$.

- First-order necessary condition: If \mathbf{x}^* is a minimum point (minimizer), then $\nabla f(\mathbf{x}^*) = \mathbf{0}$. In other words, if \mathbf{x}^* is a minimum point, then it must be a stationary point.
- Second-order necessary condition: If \mathbf{x}^* is a minimum point (minimizer), then $\mathbf{H}(\mathbf{x}^*)$ is a positive semidefinite matrix, i.e., $\mathbf{H}(\mathbf{x}^*) \succeq \mathbf{0}$.
- Second-order sufficient condition: If $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\mathbf{H}(\mathbf{x}^*)$ is a positive definite matrix, i.e., $\mathbf{H}(\mathbf{x}^*) \succ \mathbf{0}$, then \mathbf{x}^* is a minimum point (minimizer).
- One-dimensional (1-D) optimization (line search)
 - ♦ Popular 1-D methods for unimodal functions include Dichotomous search, Fibonacci search, Golden-section search, quadratic interpolation method, and cubic interpolation method.

(we illustrate the basic ideas behind the Dichotomous search, Golden-section search, and quadratic interpolation methods.)

◆ 1-D optimization for general smooth functions: the inexact line search method by Roger Fletcher. Matlab codes implementing these search methods are available for download at

<http://www.ece.uvic.ca/~optimization/>

• Determination of a search direction \mathbf{d}_k :

◆ Basis of the methods: Taylor series of $f(\mathbf{x})$ in a small vicinity surrounding point \mathbf{x}_k :

$$f(\mathbf{x}_k + \boldsymbol{\delta}) = f(\mathbf{x}_k) + \nabla^T f(\mathbf{x}_k) \cdot \boldsymbol{\delta} + \frac{1}{2} \boldsymbol{\delta}^T \cdot \nabla^2 f(\mathbf{x}_k) \cdot \boldsymbol{\delta} + O(\|\boldsymbol{\delta}\|^3)$$

or

$$f(\mathbf{x}_k + \boldsymbol{\delta}) = f(\mathbf{x}_k) + \mathbf{g}^T(\mathbf{x}_k) \cdot \boldsymbol{\delta} + \frac{1}{2} \boldsymbol{\delta}^T \cdot \mathbf{H}(\mathbf{x}_k) \cdot \boldsymbol{\delta} + O(\|\boldsymbol{\delta}\|^3)$$

◆ Steepest descent method: $\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$, i.e., $\mathbf{d}_k = -\mathbf{g}(\mathbf{x}_k)$

□ Remarks:

(i) The concept of a descent algorithm $f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k)$ for $k = 0, 1, \dots$

(ii) The steepest descent algorithm is a descent algorithm.

(iii) The steepest descent algorithm is usually quite slow.

◆ Newton method: $\mathbf{d}_k = -\mathbf{H}^{-1}(\mathbf{x}_k) \mathbf{g}(\mathbf{x}_k)$

□ Remarks:

(i) if $\mathbf{H}(\mathbf{x}_k) = \mathbf{H}_k$ is positive definite for all \mathbf{x}_k , then the Newton algorithm is a descent algorithm.

(ii) if \mathbf{H}_k is not positive definite, then the search direction in Newton method is modified to

$$\mathbf{d}_k = -\hat{\mathbf{H}}_k^{-1} \mathbf{g}_k \text{ where } \hat{\mathbf{H}}_k = \frac{\mathbf{H}_k + \beta \mathbf{I}}{1 + \beta} \text{ with a } \beta \text{ such that } \mathbf{H}_k + \beta \mathbf{I} \succ \mathbf{0}$$

The Newton method with such modifications is a descent algorithm.

(iii) Typically Newton algorithm converges considerably faster than the steepest descent algorithm at a cost of increased computational intensity.

◆ Quasi-Newton methods: one can take a unifying look at the above two methods as

$$\mathbf{d}_k = -\mathbf{S}_k \mathbf{g}_k$$

where

$$\mathbf{S}_k = \begin{cases} \mathbf{I} & \text{for steepest descent} \\ \mathbf{H}_k^{-1} & \text{for Newton} \end{cases}$$

□ Quasi-Newton methods use gradient information to generate successive low-rank estimation of the inverse Hessian matrix. In this way, quasi-Newton algorithms provide convergent rates comparable with that of Newton algorithm with reduced complexity.

□ Notation: $\boldsymbol{\delta}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$, $\boldsymbol{\gamma}_k = \mathbf{g}_{k+1} - \mathbf{g}_k$

□ Requirements for \mathbf{S}_k : (i) Symmetric and positive definite, (ii) $\mathbf{S}_{k+1} \boldsymbol{\gamma}_k = \boldsymbol{\delta}_k$

□ Davidon-Fletcher-Powell (DFP) updating formula: $\mathbf{S}_0 = \mathbf{I}$, and

$$\mathbf{S}_{k+1} = \mathbf{S}_k + \frac{\boldsymbol{\delta}_k \boldsymbol{\delta}_k^T}{\boldsymbol{\delta}_k^T \boldsymbol{\gamma}_k} - \frac{\mathbf{S}_k \boldsymbol{\gamma}_k \boldsymbol{\gamma}_k^T \mathbf{S}_k}{\boldsymbol{\gamma}_k^T \mathbf{S}_k \boldsymbol{\gamma}_k}$$

□ Broyden-Fletcher-Goldfarb-Shanno (BFGS) updating formula: $\mathbf{S}_0 = \mathbf{I}$, and

$$\mathbf{S}_{k+1} = \mathbf{S}_k + \left(1 + \frac{\boldsymbol{\gamma}_k^T \mathbf{S}_k \boldsymbol{\gamma}_k}{\boldsymbol{\gamma}_k^T \boldsymbol{\delta}_k} \right) \frac{\boldsymbol{\delta}_k \boldsymbol{\delta}_k^T}{\boldsymbol{\gamma}_k^T \boldsymbol{\delta}_k} - \frac{\boldsymbol{\delta}_k \boldsymbol{\gamma}_k^T \mathbf{S}_k + \mathbf{S}_k \boldsymbol{\gamma}_k \boldsymbol{\delta}_k^T}{\boldsymbol{\gamma}_k^T \boldsymbol{\delta}_k}$$

1.2 Constrained Optimization: The Karush-Kuhn-Tucker (KKT) Conditions (Chap. 10)

- The constrained optimization problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to:} && a_i(\mathbf{x}) = 0 \quad \text{for } i = 1, 2, \dots, p \\ & && c_j(\mathbf{x}) \geq 0 \quad \text{for } j = 1, 2, \dots, q \end{aligned}$$

- Feasible region: $\{\mathbf{x} : a_i(\mathbf{x}) = 0 \text{ for } i = 1, 2, \dots, p, \text{ and } c_j(\mathbf{x}) \geq 0 \text{ for } j = 1, 2, \dots, q\}$
- A point \mathbf{x} is said to be feasible if it is in the feasible region (inside or on the boundary). We say an inequality constraint $c_j(\mathbf{x}) \geq 0$ is active at a feasible point \mathbf{x} , if $c_j(\mathbf{x}) = 0$. And inequality constraint $c_j(\mathbf{x}) \geq 0$ is said to be inactive at a feasible point \mathbf{x} if $c_j(\mathbf{x}) > 0$.
- KKT conditions (first-order necessary conditions for point \mathbf{x}^* to be a minimum point):
If \mathbf{x}^* is a local minimizer of the constrained problem, then
 - (a) $a_i(\mathbf{x}^*) = 0$ for $i = 1, 2, \dots, p$
 - (b) $c_j(\mathbf{x}^*) \geq 0$ for $j = 1, 2, \dots, q$
 - (c) there exist Lagrange multipliers λ_i^* for $1 \leq i \leq p$ and μ_j^* for $1 \leq j \leq q$ such that

$$\nabla f(\mathbf{x}^*) = \sum_{i=1}^p \lambda_i^* \nabla a_i(\mathbf{x}^*) + \sum_{j=1}^q \mu_j^* \nabla c_j(\mathbf{x}^*)$$

$$(d) \mu_j^* c_j(\mathbf{x}^*) = 0 \text{ for } j = 1, 2, \dots, q$$

$$(e) \mu_j^* \geq 0 \text{ for } j = 1, 2, \dots, q$$

Remarks:

- ♦ Conditions (a) and (b) are merely the constraints imposed in the problem.
- ♦ Condition (c) says at a local minimizer, the gradient of the objective function is a linear combination of the gradient of the constraint functions, and the coefficients of the combination are Lagrange multipliers. See Example 10.10 and Fig. 10.8 for illustration.
- ♦ Define Lagrangian $L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) - \sum_{i=1}^p \lambda_i a_i(\mathbf{x}) - \sum_{j=1}^q \mu_j c_j(\mathbf{x})$, the equation in (c) can be expressed as

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = \mathbf{0}$$

Also note that the equality and inequality constraints can be expressed in terms of Lagrangian as

$$\nabla_{\boldsymbol{\lambda}} L(\mathbf{x}^*, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{0}$$

and

$$\nabla_{\boldsymbol{\mu}} L(\mathbf{x}^*, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq \mathbf{0}$$

- ♦ Conditions in (d) are called *complementarity conditions*. It follows that if the k th inequality constraint is active at \mathbf{x}^* , then $\mu_k = 0$. In other words, the equation in (c) only contains those terms $\nabla c_j(\mathbf{x}^*)$ that are active at \mathbf{x}^* .
- ♦ The total number of equations p (from (a)) + n (from (c)) + q (from (d)) are equal to the total number of unknowns $\{\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*\}$.

1.3 Convex Sets, Convex Functions, and Convex Programming (CP) (Chaps. 2, 10, 12, 13, 14)

- ♦ A set \mathcal{R} is said to be convex if for every pair of points $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{R}$ and every real α with $0 \leq \alpha \leq 1$, point $\mathbf{x} = \alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2$ belongs to \mathcal{R} for all $0 \leq \alpha \leq 1$.

♦ A function $f(\mathbf{x})$ defined over a convex set \mathcal{R} is said to be convex if for every pair of points $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{R}$ and every real α with $0 \leq \alpha \leq 1$,

$$f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2) \leq \alpha f(\mathbf{x}_1) + (1 - \alpha) f(\mathbf{x}_2)$$

♦ Suppose $f(\mathbf{x}) \in C^1$. Function $f(\mathbf{x})$ is convex over a convex set \mathcal{R} if and only if

$$f(\mathbf{x}_1) \geq f(\mathbf{x}) + g(\mathbf{x})^T (\mathbf{x}_1 - \mathbf{x})$$

for all \mathbf{x} and $\mathbf{x}_1 \in \mathcal{R}$

♦ Using the above property, it can be shown that a function $f(\mathbf{x}) \in C^2$ is convex over a convex set \mathcal{R} if and only if the Hessian $\mathbf{H}(\mathbf{x})$ of $f(\mathbf{x})$ is positive semidefinite for $\mathbf{x} \in \mathcal{R}$.

♦ The problem of minimizing a convex function $f(\mathbf{x})$ (with no constraints) is not very complicated, because the first-order necessary condition becomes a sufficient condition, and any local minimizer of a convex function is a global minimizer.

♦ A constrained minimization problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to:} && a_i(\mathbf{x}) = 0 \quad \text{for } i = 1, 2, \dots, p \\ & && c_j(\mathbf{x}) \geq 0 \quad \text{for } j = 1, 2, \dots, q \end{aligned}$$

is said to be a convex programming (CP) problem if the objective function is convex and the feasible region defined by the constraints, i.e., $\{\mathbf{x} : a_i(\mathbf{x}) = 0 \text{ for } i = 1, 2, \dots, p, \text{ and } c_j(\mathbf{x}) \geq 0 \text{ for } j = 1, 2, \dots, q\}$ is a convex set. In words, a CP problem is a problem of minimizing a convex function over a convex region. It can be shown that the feasible region $\{\mathbf{x} : a_i(\mathbf{x}) = 0 \text{ for } i = 1, 2, \dots, p, \text{ and } c_j(\mathbf{x}) \geq 0 \text{ for } j = 1, 2, \dots, q\}$ is convex if (i) all equality constraint functions $a_i(\mathbf{x})$ are linear, and (ii) all inequality functions $-c_j(\mathbf{x})$ are convex.

♦ An important property of the class of CP problems is that the KKT conditions become both necessary and sufficient. This implies that any $\{\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*\}$ satisfying the KKT conditions gives a local solution (minimizer) \mathbf{x}^* for the CP problem and any local minimizer of a CP problem is a global minimizer.

1.4 Wolfe Duality (Chap. 10)

Wolfe's theorem (1991) on duality is concerned with the general CP problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ (\mathcal{P}) \quad & \text{subject to:} && a_i(\mathbf{x}) = \mathbf{a}_i^T \mathbf{x} - b_i = 0 \quad \text{for } i = 1, 2, \dots, p \\ & && c_j(\mathbf{x}) \geq 0 \quad \text{for } j = 1, 2, \dots, q \end{aligned}$$

which is referred as the *primal problem*. Wolfe's theorem says that if $\{\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*\}$ solve the KKT system for the primal problem, then $\{\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*\}$ also solves the dual problem which is defined by

$$\begin{aligned} & \text{maximize} && L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \\ & && \mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu} \\ (\mathcal{D}) \quad & \text{subject to:} && \nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{0} \\ & && \boldsymbol{\mu} \geq \mathbf{0} \end{aligned}$$

In addition, $f(\mathbf{x}^*) = L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$.

♦ Proof: From the KKT conditions of the primal problem, we obtain $f(\mathbf{x}^*) = L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ and $\boldsymbol{\mu}^* \geq \mathbf{0}$. For a feasible set $\{\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}\}$, we have $\boldsymbol{\mu} \geq \mathbf{0}$ and $\nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{0}$, hence

$$L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = f(\mathbf{x}^*) \geq f(\mathbf{x}^*) - \sum_{i=1}^p \lambda_i^* a_i(\mathbf{x}^*) - \sum_{j=1}^q \mu_j^* c_j(\mathbf{x}^*) = L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$$

With $\mu \geq \mathbf{0}$ the Lagrangian $L(\mathbf{x}, \lambda, \mu)$ is convex, hence

$$L(\mathbf{x}^*, \lambda, \mu) \geq L(\mathbf{x}, \lambda, \mu) + (\mathbf{x}^* - \mathbf{x})^T \nabla_{\mathbf{x}} L(\mathbf{x}, \lambda, \mu) = L(\mathbf{x}, \lambda, \mu)$$

Therefore $L(\mathbf{x}^*, \lambda^*, \mu^*) \geq L(\mathbf{x}, \lambda, \mu)$, i.e., set $\{\mathbf{x}^*, \lambda^*, \mu^*\}$ solves the dual problem. ■

♦ If we define the duality gap $\delta(\mathbf{x}, \lambda, \mu) = f(\mathbf{x}) - L(\mathbf{x}, \lambda, \mu)$, then the KKT conditions imply that

$\delta(\mathbf{x}, \lambda, \mu) \geq 0$ for feasible $(\mathbf{x}, \lambda, \mu)$ and $\delta(\mathbf{x}^*, \lambda^*, \mu^*) = 0$. This is because for a feasible set $(\mathbf{x}, \lambda, \mu)$, we

have $\delta(\mathbf{x}, \lambda, \mu) = f(\mathbf{x}) - L(\mathbf{x}, \lambda, \mu) = \sum_{i=1}^p \lambda_i a_i(\mathbf{x}) + \sum_{j=1}^q \mu_j c_j(\mathbf{x}) = \sum_{j=1}^q \mu_j c_j(\mathbf{x}) \geq 0$, and

$$\delta(\mathbf{x}^*, \lambda^*, \mu^*) = \sum_{j=1}^q \mu_j^* c_j(\mathbf{x}^*) = 0$$

1.5 Linear Programming and Quadratic Programming (Chap. 12)

♦ Linear programming problems are of great importance in both theoretical and practical perspectives. The modern theory and algorithms of interior-point methods were initially developed for the LP problems in 1980's and the basic concepts and solution techniques for LP have since been extended to larger classes of CP and non-convex problems. On the other hand, there are a great many of practical problems that can be accurately or approximately modeled as LP problems.

♦ The standard-form linear programming (LP) problem is given by

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) = \mathbf{c}^T \mathbf{x} \\ (\mathcal{P}) \quad & \text{subject to:} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0} \end{aligned}$$

where matrix \mathbf{A} is of size $p \times n$ with $p < n$ and \mathbf{A} is assumed to have full row rank. It is important to realize that LP problems are convex programming problems. Hence the KKT conditions are necessary and sufficient conditions, and any local solution is a global solution.

♦ Moreover, the Wolfe theorem applies: the dual of the above problem is given by

$$\begin{aligned} & \text{maximize} && h(\lambda) = \mathbf{b}^T \lambda \\ (\mathcal{D}) \quad & \text{subject to:} && \mathbf{A}^T \lambda + \mu = \mathbf{c} \\ & && \mu \geq \mathbf{0} \end{aligned}$$

♦ The KKT conditions in this case are given by

$$\begin{aligned} & \mathbf{A}^T \lambda + \mu = \mathbf{c} \\ & \mathbf{A}\mathbf{x} = \mathbf{b} \\ & x_i \mu_i = 0 \quad \text{for } 1 \leq i \leq n \\ & \mathbf{x} \geq \mathbf{0}, \quad \mu \geq \mathbf{0} \end{aligned}$$

♦ A useful concept is the development of interior-point algorithms for LP is *central path*. To introduce it we first write the KKT conditions as

$$\begin{aligned} & \mathbf{A}\mathbf{x} = \mathbf{b} && \text{with } \mathbf{x} \geq \mathbf{0} \\ & \mathbf{A}^T \lambda + \mu = \mathbf{c} && \text{with } \mu \geq \mathbf{0} \\ & \mathbf{X}\mu = \mathbf{0} \end{aligned}$$

where $\mathbf{X} = \text{diag}\{x_1, x_2, \dots, x_n\}$. The central path for a standard-form LP is defined by a set of vectors $\{\mathbf{x}(\tau), \boldsymbol{\lambda}(\tau), \boldsymbol{\mu}(\tau)\}$ that satisfy

$$(C) \quad \begin{aligned} \mathbf{A}\mathbf{x} &= \mathbf{b} && \text{with } \mathbf{x} > \mathbf{0} \\ \mathbf{A}^T \boldsymbol{\lambda} + \boldsymbol{\mu} &= \mathbf{c} && \text{with } \boldsymbol{\mu} > \mathbf{0} \\ \mathbf{X}\boldsymbol{\mu} &= \tau \mathbf{e} && \text{with } \tau > 0, \mathbf{e} = [1 \ 1 \ \dots \ 1]^T \end{aligned}$$

The duality gap on the central path is found to be

$$\begin{aligned} \delta[\mathbf{x}(\tau), \boldsymbol{\lambda}(\tau)] &= \mathbf{c}^T \mathbf{x}(\tau) - \mathbf{b}^T \boldsymbol{\lambda}(\tau) \\ &= [\boldsymbol{\lambda}^T(\tau) \mathbf{A} + \boldsymbol{\mu}^T(\tau)] \mathbf{x}(\tau) - \mathbf{b}^T \boldsymbol{\lambda}(\tau) \\ &= \boldsymbol{\mu}^T(\tau) \mathbf{x}(\tau) = n\tau \end{aligned}$$

We see that duality gap approaches zero as parameter τ approaches zero. In other words, the central path defines a parameterized trajectory approaching to the solution of both primal and dual problems as τ approaches zero.

◆ Below we describe a primal-dual path-following algorithm for LP problems. The idea behind the algorithm is to start with a strictly feasible initial set of vectors $\mathbf{w}_0 = \{\mathbf{x}_0, \boldsymbol{\lambda}_0, \boldsymbol{\mu}_0\}$ and use first-order perturbations of the equations (C) to obtain a set of increment vectors to update point \mathbf{w}_k . So in the k th iteration we seek $\boldsymbol{\delta}_w = \{\boldsymbol{\delta}_x, \boldsymbol{\delta}_\lambda, \boldsymbol{\delta}_\mu\}$ such that the next iterate

$$\mathbf{w}_{k+1} = \{\mathbf{x}_{k+1}, \boldsymbol{\lambda}_{k+1}, \boldsymbol{\mu}_{k+1}\} = \{\mathbf{x}_k + \boldsymbol{\delta}_x, \boldsymbol{\lambda}_k + \boldsymbol{\delta}_\lambda, \boldsymbol{\mu}_k + \boldsymbol{\delta}_\mu\}$$

remains strictly feasible and approaches the central path defined by (C) with $\tau = \tau_{k+1}$. This leads (C) to a system of linear equations as follows

$$\begin{aligned} \mathbf{A}\boldsymbol{\delta}_x &= \mathbf{0} \\ \mathbf{A}^T \boldsymbol{\delta}_\lambda + \boldsymbol{\delta}_\mu &= \mathbf{0} \\ \mathbf{M}\boldsymbol{\delta}_x + \mathbf{X}\boldsymbol{\delta}_\mu &= \tau_{k+1} \mathbf{e} - \mathbf{X}\boldsymbol{\mu}_k \end{aligned}$$

where

$$\mathbf{M} = \text{diag}\{(\boldsymbol{\mu}_k)_1, (\boldsymbol{\mu}_k)_2, \dots, (\boldsymbol{\mu}_k)_n\}$$

◆ Additional details:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \alpha_k \boldsymbol{\delta}_w$$

where α_k is obtained by a line-search step to ensure the strict feasibility of point \mathbf{w}_{k+1} .

□ The value of parameter τ_{k+1} is determined by

$$\tau_{k+1} = \frac{\boldsymbol{\mu}_k^T \mathbf{x}_k}{n + \rho} \quad \text{with } \rho > \sqrt{n}$$

□ An interior-point algorithm with nonfeasible initialization can be developed in a similar way, see Sec. 12.5.2.

◆ LP problems may be formulated and solved in an alternative-form which are equivalent to LP problems in the standard form:

$$\begin{aligned} &\text{minimize } f(\mathbf{x}) = \mathbf{c}^T \mathbf{x} \\ &\text{subject to: } \mathbf{A}\mathbf{x} \geq \mathbf{b} \end{aligned}$$

Counterpart algorithms for LP problems in alternative-form can also be developed.

1.6 Semidefinite Programming (SDP) (Chap. 14)

♦ SDP is a class of CP problems that includes LP and convex quadratic programming problems as its subclasses and a lot more other CP problems of practical usefulness, 1990's have witnessed intensive research interest in SDP that has resulted in many efficient solution methods and software tools for SDP problems.

♦ Notation and definitions

We denote the set of $n \times n$ symmetric matrices by S^n . For $A, B \in S^n$, the inner product of A with B is defined by

$$A \cdot B = \text{trace}(AB) = \sum_{i=1}^n \sum_{j=1}^n a_{ij} b_{ij}$$

Matrix A being positive semidefinite is signified as $A \succeq \mathbf{0}$, and A being positive definite is denoted by $A \succ \mathbf{0}$.

♦ The primal SDP problem is defined as

$$\begin{aligned} (\mathcal{P}) \quad & \text{minimize} \quad C \cdot X \\ & \text{subject to: } A_i \cdot X = b_i \quad \text{for } i = 1, 2, \dots, p \\ & \quad \quad \quad X \succeq \mathbf{0} \end{aligned}$$

where C, X , and A_i are members of S^n . Note that if $C = \text{diag}\{c\}$, $A_i = \text{diag}\{a_i\}$, and $x = \text{diag}\{X\}$ for some vectors c and a_i , then the SDP problem becomes a standard LP problem:

$$\begin{aligned} & \text{minimize} \quad f(x) = c^T x \\ & \text{subject to: } Ax = b \\ & \quad \quad \quad x \geq \mathbf{0} \end{aligned}$$

where A is a matrix with a_i^T as its i -th row and $b = [b_1 \ b_2 \ \dots \ b_p]^T$.

♦ The dual SDP problem assumes the form

$$\begin{aligned} (\mathcal{D}) \quad & \text{maximize} \quad b^T y \\ & \text{subject to: } \sum_{i=1}^p y_i A_i + S = C \\ & \quad \quad \quad S \succeq \mathbf{0} \end{aligned}$$

By eliminating S and making simple variable changes, we can write the dual problem as

$$\begin{aligned} (\mathcal{D}') \quad & \text{minimize} \quad c^T x \\ & \text{subject to: } F(x) \succeq \mathbf{0} \end{aligned}$$

where

$$F(x) = F_0 + \sum_{i=1}^p x_i F_i$$

♦ The KKT conditions for SDP

$$\begin{aligned} & \sum_{i=1}^p y_i A_i + S = C \\ & A_i \cdot X = b_i \quad \text{for } i = 1, 2, \dots, p \\ & SX = \mathbf{0} \\ & X \succeq \mathbf{0}, \ S \succeq \mathbf{0} \end{aligned}$$

♦ An example: convex quadratic programming (QP) problems

The general convex QP problems can be expressed as

$$\begin{aligned} & \text{minimize} \quad \mathbf{x}^T \mathbf{H} \mathbf{x} + \mathbf{p}^T \mathbf{x} \quad \text{with } \mathbf{H} \succeq \mathbf{0} \\ & \text{subject to: } \mathbf{A} \mathbf{x} \geq \mathbf{b} \end{aligned}$$

which is equivalent to

$$\begin{aligned} & \text{minimize} \quad \delta \\ & \text{subject to: } \mathbf{x}^T \mathbf{H} \mathbf{x} + \mathbf{p}^T \mathbf{x} \leq \delta \\ & \quad \mathbf{A} \mathbf{x} \geq \mathbf{b} \end{aligned}$$

where, by using a decomposition of $\mathbf{H} = \hat{\mathbf{H}}^T \hat{\mathbf{H}}$, the first constraint can be expressed as

$$\delta - \mathbf{p}^T \mathbf{x} - (\hat{\mathbf{H}} \mathbf{x})^T (\hat{\mathbf{H}} \mathbf{x}) \geq 0$$

which is equivalent to

$$\mathbf{G}(\delta, \mathbf{x}) = \begin{bmatrix} \mathbf{I} & \hat{\mathbf{H}} \mathbf{x} \\ (\hat{\mathbf{H}} \mathbf{x})^T & \delta - \mathbf{p}^T \mathbf{x} \end{bmatrix} \succeq \mathbf{0}$$

In addition, the second constraint can be written as

$$\mathbf{F}(\mathbf{x}) = \mathbf{F}_0 + \sum_{i=1}^n x_i \mathbf{F}_i \succeq \mathbf{0}$$

where $\mathbf{F}_0 = -\text{diag}\{\mathbf{b}\}$, and $\mathbf{F}_i = \text{diag}\{\mathbf{a}_i\}$ with \mathbf{a}_i being the i -th column of \mathbf{A} .

By defining an augmented vector $\hat{\mathbf{x}} = \begin{bmatrix} \delta \\ \mathbf{x} \end{bmatrix}$, the convex QP problem is now formulated as the SDP problem

$$\begin{aligned} & \text{minimize} \quad \hat{\mathbf{c}}^T \hat{\mathbf{x}} \\ & \text{subject to: } \mathbf{E}(\hat{\mathbf{x}}) \succeq \mathbf{0} \end{aligned}$$

where $\hat{\mathbf{c}} = [1 \ 0 \ \dots \ 0]^T$, and $\mathbf{E}(\hat{\mathbf{x}}) = \text{diag}\{\mathbf{G}(\delta, \mathbf{x}), \mathbf{F}(\mathbf{x})\}$.

1.7 Second-Order Cone Programming (SOCP) (Chap. 14)

♦ The class of SOCP problem is a subset of the class of SDP problems but it is still large enough to include LP, QP, and many other CP problems of practical importance. Moreover, interior-point algorithms that are specifically designed for SOCP is considerably more efficient than what the best an SDP algorithm can offer.

♦ Notation and definitions

□ A second-order cone (also known as quadratic cone or Lorentz cone) of dimension n is defined as

$$\mathcal{K} = \left\{ \begin{bmatrix} t \\ \mathbf{u} \end{bmatrix} : t \in \mathbb{R}, \mathbf{u} \in \mathbb{R}^{n-1} \text{ for } \|\mathbf{u}\| \leq t \right\}$$

□ Examples: $n = 1$: a ray on the t -axis starting at $t = 0$, for $n = 2$ and 3 see Fig. 14.3.

□ Note that the second-order cone is convex in \mathbb{R}^n .

♦ The primal SOCP problem is given by

$$\begin{aligned} & \text{minimize} \quad \sum_{i=1}^q \hat{\mathbf{c}}_i^T \mathbf{x}_i \\ (\mathcal{P}) \quad & \text{subject to: } \sum_{i=1}^q \hat{\mathbf{A}}_i \mathbf{x}_i = \mathbf{b} \\ & \quad \mathbf{x}_i \in K_i \quad \text{for } i = 1, 2, \dots, q \end{aligned}$$

where $\hat{\mathbf{c}}_i \in R^{n_i \times 1}$, $\mathbf{x}_i \in R^{n_i \times 1}$, $\hat{\mathbf{A}}_i \in R^{m \times n_i}$, $\mathbf{b} \in R^{m \times 1}$, and K_i is a second-order cone of dimension n_i .

□ We note an analogy between SOCP and LP: both involve a linear objective function and a set of linear equality constraints. While the variable \mathbf{x} in an LP problem is constrained to the region $\{\mathbf{x} : \mathbf{x} \geq \mathbf{0}\}$ which is a convex cone, each variable vector \mathbf{x}_i in an SOCP problem is constrained to the second-order cone K_i .

♦ The dual SOCP problem assumes the form

$$\begin{aligned} & \text{maximize} && \mathbf{b}^T \mathbf{y} \\ (\mathcal{D}) \quad & \text{subject to:} && \hat{\mathbf{A}}_i^T \mathbf{y} + \mathbf{s}_i = \hat{\mathbf{c}}_i \quad \text{for } i = 1, 2, \dots, q \\ & && \mathbf{s}_i \in K_i \quad \text{for } i = 1, 2, \dots, q \end{aligned}$$

□ We also notice a similar analogy between the dual SOCP and dual LP problems.

♦ Now if we let

$$\mathbf{x} = -\mathbf{y}, \quad \hat{\mathbf{A}}_i^T = \begin{bmatrix} \mathbf{b}_i^T \\ \mathbf{A}_i^T \end{bmatrix} \quad \text{and} \quad \hat{\mathbf{c}}_i = \begin{bmatrix} d_i \\ \mathbf{c}_i \end{bmatrix}$$

then the dual SOCP problem can be expressed as

$$\begin{aligned} (\mathcal{D}') \quad & \text{minimize} && \mathbf{b}^T \mathbf{x} \\ & \text{subject to:} && \|\mathbf{A}_i^T \mathbf{x} + \mathbf{c}_i\| \leq \mathbf{b}_i^T \mathbf{x} + d_i \quad \text{for } i = 1, 2, \dots, q \end{aligned}$$

□ An example: Linear fractional problem

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^p \frac{1}{\mathbf{a}_i^T \mathbf{x} + c_i} \\ & \text{subject to:} && \mathbf{a}_i^T \mathbf{x} + c_i > 0 \quad \text{for } i = 1, 2, \dots, p \\ & && \mathbf{b}_i^T \mathbf{x} + d_i \geq 0 \quad \text{for } i = 1, 2, \dots, q \end{aligned}$$

By introducing auxiliary constraints

$$\frac{1}{\mathbf{a}_i^T \mathbf{x} + c_i} \leq \delta_i$$

i.e.,

$$\delta_i (\mathbf{a}_i^T \mathbf{x} + c_i) \geq 1$$

the linear fractional problem can be formulated as

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^p \delta_i \\ & \text{subject to:} && \delta_i (\mathbf{a}_i^T \mathbf{x} + c_i) \geq 1 \quad \text{for } i = 1, 2, \dots, p \\ & && \delta_i \geq 0 \\ & && \mathbf{b}_i^T \mathbf{x} + d_i \geq 0 \end{aligned}$$

We can show that $w^2 \leq uv$, $u \geq 0, v \geq 0$ if and only if

$$\left\| \begin{bmatrix} 2w \\ u - v \end{bmatrix} \right\| \leq u + v$$

Hence the first two constraints in the above problem are equivalent to

$$\left\| \begin{bmatrix} 2 \\ \mathbf{a}_i^T \mathbf{x} + c_i - \delta_i \end{bmatrix} \right\| \leq \mathbf{a}_i^T \mathbf{x} + c_i + \delta_i \quad \text{for } i = 1, 2, \dots, p$$

and the linear fractional problem can be converted into

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^p \delta_i \\ & \text{subject to:} && \left\| \begin{bmatrix} 2 \\ \mathbf{a}_i^T \mathbf{x} + c_i - \delta_i \end{bmatrix} \right\| \leq \mathbf{a}_i^T \mathbf{x} + c_i + \delta_i \quad \text{for } i = 1, 2, \dots, p \\ & && \mathbf{b}_i^T \mathbf{x} + d_i \geq 0 \end{aligned}$$

which is an SOCP problem.

Use SeDuMi to Solve LP, SDP and SCOP Problems: Remarks and Examples*

* This file was prepared by Wu-Sheng Lu, Dept. of Electrical and Computer Engineering, University of Victoria, revised on January 29, 2007.

The name of the toolbox, *SeDuMi*, stands for *self-dual minimization* as it implements a self-dual embedding technique for optimization over self-dual homogeneous cones. Below we give comments on the usage of this software for the solution of LP, SDP, and SOCP problems. Numerical examples are included for illustration purposes. Some of the examples are from the book:

A. Antoniou and W.-S. Lu, *Optimization: Algorithms and Applications*, Springer, 2007.

0. Presently the official web site for public-domain software SeDuMi is <http://sedumi.mcmaster.ca/>. As of January 2007, the latest version is SeDuMi version 1.1R3. However, this version clashes with MATLAB R2006b. So for users of MATLAB R2006b, the best version of SeDuMi right now appears to be SeDuMi 1.1R2 which can also be downloaded from the above site.

I. LP Problems

The solution of the primal standard-form LP problem

$$\text{minimize} \quad \mathbf{c}^T \mathbf{x} \quad (1a)$$

$$\text{subject to:} \quad \mathbf{Ax} = \mathbf{b} \quad (1b)$$

$$\mathbf{x} \geq \mathbf{0} \quad (1c)$$

can be obtained by using command $\mathbf{x} = \text{sedumi}(\mathbf{A}, \mathbf{b}, \mathbf{c})$; Alternatively, both the solution of the problem in (1) and the solution of the dual problem

$$\text{maximize} \quad \mathbf{b}^T \mathbf{y} \quad (2a)$$

$$\text{subject to:} \quad -\mathbf{A}^T \mathbf{y} \geq -\mathbf{c} \quad (2b)$$

can be found by using command $[\mathbf{x}, \mathbf{y}, \text{info}] = \text{sedumi}(\mathbf{A}, \mathbf{b}, \mathbf{c})$; where “info” contains information about validity of the solutions obtained:

- (1) $\text{pinf} = \text{dinf} = 0$: \mathbf{x} is an optimal solution and \mathbf{y} certifies optimality, viz. $\mathbf{b}^T \mathbf{y} = \mathbf{c}^T \mathbf{x}$ and $\mathbf{c} - \mathbf{A}^T \mathbf{y} \geq \mathbf{0}$. Stated otherwise, \mathbf{y} is an optimal solution to maximize $\mathbf{b}^T \mathbf{y}$ such that $\mathbf{c} - \mathbf{A}^T \mathbf{y} \geq \mathbf{0}$.

If $\text{size}(\mathbf{A}, 2) = \text{length}(\mathbf{b})$, then \mathbf{y} solves the linear program maximize $\mathbf{b}^T \mathbf{y}$ such that $\mathbf{c} - \mathbf{A}^T \mathbf{y} \geq \mathbf{0}$.

- (2) $\text{pinf} = 1$: there cannot be $\mathbf{x} \geq \mathbf{0}$ with $\mathbf{Ax} = \mathbf{b}$, and this is certified by \mathbf{y} , viz. $\mathbf{b}^T \mathbf{y} > 0$ and $\mathbf{A}^T \mathbf{y} \leq \mathbf{0}$. Thus \mathbf{y} is a Farkas solution.

- (3) $\text{dinf} = 1$: there cannot be \mathbf{y} such that $\mathbf{c} - \mathbf{A}^T \mathbf{y} \geq \mathbf{0}$, and this is certified by \mathbf{x} , viz. $\mathbf{c}^T \mathbf{x} < 0$, $\mathbf{Ax} = \mathbf{b}$, $\mathbf{x} \geq \mathbf{0}$. Thus \mathbf{x} is a Farkas solution.

Example 1 Consider the standard-form LP problem in Example 11.9 [Antoniou and Lu]:

$$\begin{aligned}
&\text{minimize} && f(x) = 2x_1 + 9x_2 + 3x_3 \\
&\text{subject to:} && -2x_1 + 2x_2 + x_3 - x_4 = 1 \\
&&& x_1 + 4x_2 - x_3 - x_5 = 1 \\
&&& x_1 \geq 0, x_2 \geq 0, x_3 \geq 0, x_4 \geq 0, x_5 \geq 0
\end{aligned} \tag{3}$$

The MATLAB code listed below uses SeDuMi to solve the problem:

```

A = [-2 2 1 -1 0; 1 4 -1 0 -1];
b = [1 1]';
c = [2 9 3 0 0]';
x = sedumi(A,b,c);

```

which gives $x = [0 \ 0.3333 \ 0.3333 \ 0 \ 0]'$.

Example 2 Now let us consider the LP problem in Example 11.2 [Antoniou and Lu]:

$$\begin{aligned}
&\text{minimize} && -x_1 - 4x_2 \\
&\text{subject to:} && x_1 \geq 0 \\
&&& -x_1 \geq -2 \\
&&& x_2 \geq 0 \\
&&& -x_1 - x_2 + 3.5 \geq 0 \\
&&& -x_1 - 2x_2 + 6 \geq 0
\end{aligned} \tag{4}$$

By changing the variable x to y and defining

$$A = \begin{bmatrix} -1 & 1 & 0 & 1 & 1 \\ 0 & 0 & -1 & 1 & 2 \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ 4 \end{bmatrix}, \quad c = [0 \ 2 \ 0 \ 3.5 \ 6]^T$$

the problem in (4) can be expressed as the one in (2). Hence the LP problem in (4) can be solved by using the following MATLAB code:

```

b = [1 4]';
A = [-1 1 0 1 1; 0 0 -1 1 2];
c = [0 2 0 3.5 6]';
[x, y, info] = sedumi(A,b,c);

```

where y can be taken as the solution of the problem in (4). The result is given by $y = [0 \ 3]'$. In addition, output “info” provides the following information:

```

cpusec: 0.0938
iter: 4
feasratio: 1
pinf: 0
dinf: 0
numerr: 0

```

II. SOCP Problems

We now consider the SOCP formulation given by Eq. (14.104) (see [Antoniou and Lu]), i.e.,

$$\text{minimize} \quad \mathbf{b}^T \mathbf{x} \quad (5a)$$

$$\text{subject to:} \quad \|\mathbf{A}_i^T \mathbf{x} + \mathbf{c}_i\| \leq \mathbf{b}_i^T \mathbf{x} + d_i \quad \text{for } i=1, \dots, q \quad (5b)$$

where $\mathbf{b} \in R^{m \times 1}$, $\mathbf{x} \in R^{m \times 1}$, $\mathbf{A}_i \in R^{m \times (n_i-1)}$, $\mathbf{b}_i \in R^{m \times 1}$, $\mathbf{c}_i \in R^{(n_i-1) \times 1}$, and $d_i \in R$ for $1 \leq i \leq q$. In order to use the toolbox to solve the problem in (5), we define matrix \mathbf{A}_t , vectors \mathbf{b}_t and \mathbf{c}_t as follows:

$$\mathbf{A}_t = [\mathbf{A}_t^{(1)} \quad \mathbf{A}_t^{(2)} \quad \dots \quad \mathbf{A}_t^{(q)}]$$

$$\mathbf{A}_t^{(i)} = -[\mathbf{b}_i \quad \mathbf{A}_i]$$

$$\mathbf{b}_t = -\mathbf{b}$$

$$\mathbf{c}_t = [\mathbf{c}_t^{(1)}; \quad \mathbf{c}_t^{(2)}; \quad \dots \quad \mathbf{c}_t^{(q)}]$$

$$\mathbf{c}_t^{(i)} = [d_i; \quad \mathbf{c}_i]$$

where $\mathbf{A}_t \in R^{m \times n}$, $\mathbf{b}_t \in R^{m \times 1}$, and $\mathbf{c}_t \in R^{n \times 1}$ with $n = \sum_{i=1}^q n_i$.

In addition, define a q -dimensional vector $\mathbf{q} = [n_1 \quad n_2 \quad \dots \quad n_q]$ which describes the dimensions of the q conic constraints involved in (5b).

Once the data set $\{\mathbf{A}_t, \mathbf{b}_t, \mathbf{c}_t\}$ and vector \mathbf{q} have been prepared, the MATLAB commands for solving the SOCP problem in (5) are as follows:

```
K.q = q;
[xs,ys,info] = sedumi(At,bt,ct,K);
info
x = ys;
```

Example 3 Consider the problem described in Example 14.5 in [Antoniou and Lu]: Find the shortest distance between the two ellipses which are defined by

$$c_1(\mathbf{x}) = -[x_1 \quad x_2] \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + [x_1 \quad x_2] \begin{bmatrix} \frac{1}{2} \\ 0 \end{bmatrix} + \frac{3}{4} \geq 0$$

$$c_2(\mathbf{x}) = -\frac{1}{8}[x_3 \quad x_4] \begin{bmatrix} 5 & 3 \\ 3 & 5 \end{bmatrix} \begin{bmatrix} x_3 \\ x_4 \end{bmatrix} + [x_3 \quad x_4] \begin{bmatrix} \frac{11}{2} \\ 13 \end{bmatrix} - \frac{35}{2} \geq 0$$

The problem can be formulated as the constrained minimization problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) = (x_1 - x_3)^2 + (x_2 - x_4)^2 \\ & \text{subject to:} && c_1(\mathbf{x}) \geq 0 \quad \text{and} \quad c_2(\mathbf{x}) \geq 0 \end{aligned}$$

By introducing an upper bound δ for the square root of the objective function, the above problem can be expressed as

$$\begin{aligned} & \text{minimize} && \delta \\ & \text{subject to:} && \left[(x_1 - x_3)^2 + (x_2 - x_4)^2 \right]^{1/2} \leq \delta \\ & && \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 1/4 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 1/2 \\ 0 \end{bmatrix} \leq \frac{3}{4} \\ & && \begin{bmatrix} x_3 & x_4 \end{bmatrix} \begin{bmatrix} 5/8 & 3/8 \\ 3/8 & 5/8 \end{bmatrix} \begin{bmatrix} x_3 \\ x_4 \end{bmatrix} - \begin{bmatrix} x_3 & x_4 \end{bmatrix} \begin{bmatrix} 11/2 \\ 13/2 \end{bmatrix} \leq -\frac{35}{2} \end{aligned}$$

If we augment the decision vector by including the upper bound δ in it, i.e., $\mathbf{x} = [\delta \ x_1 \ x_2 \ x_3 \ x_4]^T$, then the above problem is equivalent to the problem in (5) with $q = 3$,

$$\begin{aligned} \mathbf{b} &= [1 \ 0 \ 0 \ 0 \ 0]^T \\ \mathbf{A}_1^T &= \begin{bmatrix} 0 & -1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & -1 \end{bmatrix}, \mathbf{A}_2^T = \begin{bmatrix} 0 & 0.5 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix} \\ \mathbf{A}_3^T &= \begin{bmatrix} 0 & 0 & 0 & -0.7071 & -0.7071 \\ 0 & 0 & 0 & -0.3536 & 0.3536 \end{bmatrix} \\ \mathbf{b}_1 &= \mathbf{b}, \mathbf{b}_2 = \mathbf{0}, \mathbf{b}_3 = \mathbf{0} \\ \mathbf{c}_1 &= \mathbf{0}, \mathbf{c}_2 = [-0.5 \ 0]^T, \mathbf{c}_3 = [4.2426 \ -0.7071]^T \\ d_1 &= 0, d_2 = 1, \text{ and } d_3 = 1. \end{aligned}$$

The MATLAB code listed below solves the SOCP problem using SeDuMi.

```

b = [1 0 0 0 0]';
A1 = [0 -1 0 1 0; 0 0 1 0 -1]';
A2 = [0 0.5 0 0 0; 0 0 1 0 0]';
A3 = [0 0 0 -0.7071 -0.7071; 0 0 0 -0.3536 0.3536]';
b1 = b;
b2 = zeros(5,1);
b3 = b2;
c1 = [0 0]';
c2 = [-0.5 0]';
c3 = [4.2426 -0.7071]';
d1 = 0;
d2 = 1;
d3 = 1;
At1 = -[b1 A1];
At2 = -[b2 A2];
At3 = -[b3 A3];

```

```

At = [At1 At2 At3];
bt = -b;
ct1 = [d1; c1];
ct2 = [d2; c2];
ct3 = [d3; c3];
ct = [ct1; ct2; ct3];
K.q = [size(At1,2) size(At2,2) size(At3,2)];
[xs, ys, info] = sedumi(At,bt,ct,K);
x = ys;
r_s = x(2:3);
s_s = x(4:5);
disp('solution points in regions R and S are:')
[r_s s_s]
disp('minimum distance:')
norm(r_s - s_s)

```

It was found that $r_s = [2.0447 \ 0.8527]'$ and $s_s = [2.5448 \ 2.4857]'$ which give the minimum distance 1.7078. Additional information provided by the toolbox is as follows:

```

info =
    cpusec: 0.0469
       iter: 9
  feastratio: 1.0000
       pinf: 0
       dinf: 0
      numerr: 0

```

III. SOCP Problems with Linear Constraints

SeDuMi can be used to solve SOCP problems with additional linear constraints. As an example, we consider the following SOCP problem

$$\text{minimize} \quad \mathbf{b}^T \mathbf{x} \quad (6a)$$

$$\text{subject to:} \quad \mathbf{D}^T \mathbf{x} + \mathbf{f} \geq \mathbf{0} \quad (6b)$$

$$\|\mathbf{A}_i^T \mathbf{x} + \mathbf{c}_i\| \leq \mathbf{b}_i^T \mathbf{x} + d_i \quad \text{for } i = 1, \dots, q \quad (6c)$$

where $\mathbf{D} \in \mathbb{R}^{m \times p}$ and $\mathbf{f} \in \mathbb{R}^{p \times 1}$ and other entries are defined in part II. In order to use the toolbox to solve the problem in (6), we define matrix \mathbf{A}_t , vectors \mathbf{b}_t and \mathbf{c}_t , which are obtained by augmenting the corresponding quantities used for the problem in (5), as follows:

$$\mathbf{A}_t = \begin{bmatrix} -\mathbf{D} & \mathbf{A}_t^{(1)} & \mathbf{A}_t^{(2)} & \dots & \mathbf{A}_t^{(q)} \end{bmatrix}$$

$$\mathbf{A}_t^{(i)} = -[\mathbf{b}_i \quad \mathbf{A}_i]$$

$$\mathbf{b}_t = -\mathbf{b}$$

$$\mathbf{c}_t = [\mathbf{f}; \quad \mathbf{c}_t^{(1)}; \quad \mathbf{c}_t^{(2)}; \quad \dots \quad \mathbf{c}_t^{(q)}]$$

$$\mathbf{c}_t^{(i)} = [d_i; \quad \mathbf{c}_i]$$

Now define $K.l = p$ where p is the number of linear constraints in (6b). And as before, define a q -dimensional vector $\mathbf{q} = [n_1 \ n_2 \ \dots \ n_q]$ to describe the dimensions of the q conic constraints in (6c). The MATLAB commands to solve the SOCP problem in (5) are as follows:

```
K.l = m;
K.q = q;
[xs,ys,info] = sedumi(At,bt,ct,K);
info
x = ys;
```

Example 4 Consider a problem similar to the one described in Example 14.5 of [Antoniou and Lu], to that we now add three additional linear constraints for variables x_3 and x_4 as follows:

$$\begin{aligned} -x_3 + 2.4 &\geq 0 \\ x_4 - 2.4 &\geq 0 \\ 1.5x_3 - x_4 + 2.4 &\geq 0 \end{aligned}$$

It can be readily verified that the problem at hand is a minimum distance problem for the ellipse \mathcal{R} and convex body \mathcal{S} shown in the next figure.

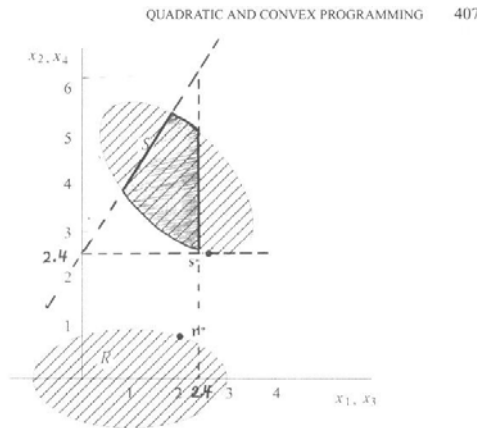


Figure 13.5. Distance between two ellipses (Example 13.5).

The MATLAB code listed below solves this SOCP problem using SeDuMi.

```
D = [0 0 0; 0 0 0; 0 0 0; -1 0 1.5; 0 1 -1];
f = [2.4 -2.4 2.4]';
b = [1 0 0 0 0]';
A1 = [0 -1 0 1 0; 0 0 1 0 -1]';
A2 = [0 0.5 0 0 0; 0 0 1 0 0]';
A3 = [0 0 0 -0.7071 -0.7071; 0 0 0 -0.3536 0.3536]';
b1 = b;
b2 = zeros(5,1);
b3 = b2;
c1 = [0 0]';
c2 = [-0.5 0]';
c3 = [4.2426 -0.7071]';
d1 = 0;
```



```

d2 = 1;
d3 = 1;
At1 = -[b1 A1];
At2 = -[b2 A2];
At3 = -[b3 A3];
At = [-D At1 At2 At3];
bt = -b;
ct1 = [d1; c1];
ct2 = [d2; c2];
ct3 = [d3; c3];
ct = [f; ct1; ct2; ct3];
K.l = size(D,2);
K.q = [size(At1,2) size(At2,2) size(At3,2)];
[xs, ys, info] = sedumi(At,bt,ct,K);
x = ys;
r_s = x(2:3);
s_s = x(4:5);
disp('solution points in regions R and S are:')
[r_s s_s]
disp('minimum distance:')
norm(r_s - s_s)

```

It was found that $r_s = [1.9518 \ 0.8795]'$ and $s_s = [2.4000 \ 2.5362]'$ which give the minimum distance 1.7163. Additional information provided by the toolbox are as follows:

```

info =
    cpusec: 0.6406
    iter: 13
    feastratio: 1.0000
    pinf: 0
    dinf: 0
    numerr: 0

```

IV. SDP Problems

Consider the “standard” dual SDP problem

$$\text{minimize } \mathbf{c}^T \mathbf{y} \quad (7a)$$

$$\text{subject to: } \mathbf{F}_0 + \mathbf{y}_1 \mathbf{F}_1 + \cdots + \mathbf{y}_p \mathbf{F}_p \succeq \mathbf{0} \quad (7b)$$

where vector $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_p]^T$, and \mathbf{F}_i are symmetric matrices. In order to use SeDuMi to solve the SDF problem, we define

$$\mathbf{b}_t = -\mathbf{c};$$

$$\mathbf{c}_t = \text{vec}(\mathbf{F}_0);$$

$$\mathbf{A}_t(:, i) = -\text{vec}(\mathbf{F}_i) \text{ for } i = 1, 2, \dots, p$$

where $\text{vec}(\mathbf{F})$ converts matrix \mathbf{F} into a column vector by stacking all columns of \mathbf{F} . The MATLAB code using commands from SeDuMi is as follows.

```

p = length(c);
bt = -c;
ct = vec(F0);
for i = 1:p, At(:,i) = -vec(Fi); end
K.s = size(F0,1);
[x, y, info] = sedumi(At,bt,ct,K);
info

```

where y gives the solution of the problem in (7) and $info$ provides information about the validity of the solution.

Example 5 Solve the problem in Example 14.1 of [Antoniou and Lu]: Find scalars y_1 , y_2 , and y_3 such that the maximum eigenvalue of $F = A_0 + y_1 A_1 + y_2 A_2 + y_3 A_3$ with

$$A_0 = \begin{bmatrix} 2 & -0.5 & -0.6 \\ -0.5 & 2 & 0.4 \\ -0.6 & 0.4 & 3 \end{bmatrix}, A_1 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, A_2 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}, A_3 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

is minimized. This problem can be formulated as the SDP problem

$$\begin{aligned} & \text{minimize } t & (8a) \\ & \text{subject to: } tI - (A_0 + y_1 A_1 + y_2 A_2 + y_3 A_3) \succeq 0 & (8b) \end{aligned}$$

With $y = [y_1 \ y_2 \ y_3 \ t]^T$, $F_0 = -A_0$, $F_1 = -A_1$, $F_2 = -A_2$, and $c = [0 \ 0 \ 0 \ 1]^T$, problem (8) becomes the standard SDP problem in (7). The MATLAB code listed below uses SeDuMi to solve the problem in (8).

```

A0 = [2 -0.5 -0.6; -0.5 2 0.4; -0.6 0.4 3];
A1 = [0 1 0; 1 0 0; 0 0 0];
A2 = [0 0 1; 0 0 0; 1 0 0];
A3 = [0 0 0; 0 0 1; 0 1 0];
F0 = -A0;
F1 = -A1;
F2 = -A2;
F3 = -A3;
F4 = eye(3);
At = [-vec(F1) vec(F2) vec(F3) vec(F4)];
bt = -[0 0 0 1]';
ct = vec(F0);
K.s = size(F0,1);
[x,y,info] = sedumi(At,bt,ct,K);
info

```

The result obtained is given by $y = [0.5 \ 0.6 \ -0.4 \ 3]^T$; which says the minimum of the largest eigenvalue is 3 that can be achieved by the choice $y_1 = 0.5$, $y_2 = 0.6$, and $y_3 = -0.4$. The 'info' gives the following:

```

cpusec: 0.5000
iter: 5
feasratio: 1.0000
pinf: 0
dinf: 0
numerr: 0

```

V. SDP Problems with Linear Constraints

In principal linear constraints can be formulated as SDP constraints (see Chap. 14 of [Antoniou and Lu]), but one can take the advantage of SeDuMi that allows both linear and SDP constraints simultaneously to make the MATLAB code more efficient. To this end, we consider the problem

$$\text{minimize } \mathbf{c}^T \mathbf{y} \quad (9a)$$

$$\text{subject to: } \mathbf{A}\mathbf{y} \geq \mathbf{b} \quad (9b)$$

$$\mathbf{F}_0 + y_1 \mathbf{F}_1 + \cdots + y_p \mathbf{F}_p \succeq \mathbf{0} \quad (9c)$$

Similar to that in part (IV), we define

$$\begin{aligned} \mathbf{b}_t &= -\mathbf{c}; \\ \mathbf{c}_t &= \text{vec}(\mathbf{F}_0); \\ \mathbf{A}_t(:, i) &= -\text{vec}(\mathbf{F}_i) \text{ for } i = 1, 2, \dots, p \end{aligned}$$

The MATLAB code using commands from SeDuMi is as follows.

```
p = length(c);
btt = -c;
ctt = [-b; vec(F0)];
for i = 1:p, At(:,i) = -vec(Fi); end
Att = [-A; At];
K.l = size(A,1);
K.s = size(F0,1);
[x, y, info] = sedumi(At,bt,ct,K);
info
```

where \mathbf{y} gives the solution of the problem in (9) and info provides information about the validity of the solution.

Example 6 We now consider a problem similar to the one in Example 5, but we add additional constraints that parameters y_1, y_2, y_3 satisfy $0.7 \leq y_1 \leq 1$, $0 \leq y_2 \leq 0.3$, and $y_3 \geq 0$. These additional linear constraints can be put into a matrix form $\mathbf{A}\mathbf{y} \geq \mathbf{b}$ with

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}, \mathbf{b} = \begin{bmatrix} 0.7 \\ -1 \\ 0 \\ -0.3 \\ 0 \end{bmatrix}$$

where vector \mathbf{y} is defined as in Example 5, i.e., $\mathbf{y} = [y_1 \ y_2 \ y_3 \ t]^T$. The MATLAB code solving this problem using SeDuMi is as follows:

```
A = [1 0 0 0; -1 0 0 0; 0 1 0 0; 0 -1 0 0; 0 0 1 0];
b = [0.7 -1 0 -0.3 0]';
A0 = [2 -0.5 -0.6; -0.5 2 0.4; -0.6 0.4 3];
```

```

A1 = [0 1 0; 1 0 0; 0 0 0];
A2 = [0 0 1; 0 0 0; 1 0 0];
A3 = [0 0 0; 0 0 1; 0 1 0];
F0 = -A0;
F1 = -A1;
F2 = -A2;
F3 = -A3;
F4 = eye(3);
At = -[vec(F1) vec(F2) vec(F3) vec(F4)];
Att = [-A; At];
btt = -[0 0 0 1]';
ct = vec(F0);
ctt = [-b; ct];
K.l = size(A,1);
K.s = size(F0,1);
[x,y,info] = sedumi(Att,btt,ctt,K);
info

```

The result obtained is given by $y = [1.0 \ 0.3 \ -0.0 \ 3.1556]'$; which says the minimum of the largest eigenvalue is 3.1556 that can be achieved by the choice $y_1 = 1.0$, $y_2 = 0.3$, and $y_3 = 0$. The info gives the following:

```

iter: 9
feasratio: 0.9967
pinf: 0
dinf: 0
numerr: 0
timing: [0.0313 0.2969 0.0156]
cpusec: 0.3438

```

VI. Other Problems

6.1 $[x, y, \text{info}] = \text{sedumi}(A, b, 0)$ solves the feasibility problem: Find x such that $Ax = b$ and $x \geq 0$.

Example 7 With A and b given in Example 1, $[x, y, \text{info}] = \text{sedumi}(A, b, 0)$ yields $x = [1.3130 \ 1.2353 \ 2.8405 \ 1.6851 \ 2.4136]^T$ and

```

info =
cpusec: 0
iter: 1
feasratio: 1
pinf: 0
dinf: 0
numerr: 0

```

6.2 $[x, y, \text{info}] = \text{sedumi}(A, 0, c)$ solves the feasibility problem: Find y such that $A^T y \leq c$.

Example 8 Let

$$A = \begin{bmatrix} -1 & 1 & 0 & 1 & 1 \\ 0 & 0 & -1 & 1 & 2 \end{bmatrix}, \quad c = [0 \ 2 \ 0 \ 3.5 \ 6]^T$$

Applying $[x, y, \text{info}] = \text{sedumi}(A, 0, c)$ yields

$y = [0.9810 \ 1.2670]^T$ and

info =

cpusec: 0.0625

iter: 3

feasratio: 1

pinf: 0

dinf: 0

numerr: 0

Variational Methods in Optimization

References

1. I. M. Gelfand and S. V. Formin, *Calculus of Variations*, Dover Publications Inc., 2000.
2. L. Rudin, S. Osher, and E. Fatemi, “Nonlinear total variation based noise removal algorithms,” *Physica D.*, vol. 60, pp. 259-268, 1992.
3. L. Rudin and S. Osher, “Total variation based image restoration with free local constraints,” *Proc. ICIP*, pp. 31-35, Austin, TX., 1994.
4. S. Osher, M. Burger, D. Goldfarb, J. Xu, and W. Yin, “An iterative regularization method for total variation-based image restoration,” *SIAM Journal on Multiscale Modeling and Simulation*, vol. 4, no. 2, pp. 460-489, 2005.
5. Y. Meyer, *Oscillating Patterns in Image Processing and Nonlinear Evolution Equations*, University Lecture Series, vol. 22, American Mathematical Society, Providence, RI., 2001.
6. UCLA Reports: <http://www.math.ucla.edu/applied/cam/index.html>

1. Elements of Calculus of Variations

1.1 Functionals

Roughly speaking, a functional is a “function of function”. We all know the definition of function, say $f(x)$ – which is taken to mean a correspondence that assigns a definite (real) number for each variable value x in a certain domain. Similarly, a *functional* is a correspondence that assigns a real number to each *function* (such as a curve, or a surface, etc.) that belongs to some class of functions.

For example, for each curve that connects two fixed points \mathbf{a} and \mathbf{b} , the length of the curve is a functional. As another example, the time required for a mass point traveling down along a smooth and monotonically decreasing plane curve with two fixed ends is a functional. From these examples, we see that a functional is defined over a certain class of functions.

1.2 Variational Problems: Examples

Example 1 A smooth plane curve that joins points $\mathbf{p}_a = (a, A)$ and $\mathbf{p}_b = (b, B)$ can be described by a differentiable function $y = y(x)$ with $A = y(a)$ and $B = y(b)$, and its length is known to be

$$J[y] = \int_a^b \sqrt{1 + y'^2} \, dx \quad (1)$$

Here $J[y]$ is obviously a functional for the class of smooth functions $y(x)$ that satisfy the boundary conditions $A = y(a)$ and $B = y(b)$. A simple variational problem associated with the functional in (1) is to find the shortest curve joining points \mathbf{p}_a and \mathbf{p}_b .

Typically, for the one-variable case a functional of interest assumes the form

$$J[y] = \int_a^b F(x, y, y') dx \quad (2)$$

and for the two-variable case

$$J[u] = \iint_R F(x, y, u, u_x, u_y) dx dy \quad (3)$$

where $u = u(x, y)$, $u_x = \partial u / \partial x$ and $u_y = \partial u / \partial y$.

Example 2 Let $u(x, y)$ denote a two-variable smooth function that models the light intensity of an (analog) image over region R . The total variation (TV) of the image is defined as

$$J[u] = \iint_R \sqrt{u_x^2 + u_y^2} dx dy \quad (4)$$

which is obviously a functional of type (3) for the class of smooth functions defined over region R . It has been argued that the noise removal problem for images can be cast as minimizing the TV subject to some statistical constraints.

1.3 Several Lemmas [1]

The key analytical tool for investigating various image processing problems in a variational optimization setting is the Euler-Lagrange (EL) equation. The lemmas presented below are instrumental in deriving the EL equation.

In what follows the set of continuous functions over interval $[a, b]$ is denoted by $C(a, b)$ and the set of functions that are continuous and have up to p -th order continuous derivatives on $[a, b]$ is denoted by $D_p(a, b)$.

Lemma 1 If $\alpha(x)$ is continuous on $[a, b]$, and if

$$\int_a^b \alpha(x) h(x) dx = 0$$

for every function $h(x) \in C(a, b)$ such that $h(a) = h(b) = 0$, then $\alpha(x) = 0$ for all x in $[a, b]$.

Proof. The lemma will be proved by the method of contradiction. Suppose the function $\alpha(x)$ is nonzero, say positive, at some point in $[a, b]$. Then $\alpha(x)$ is also positive in some interval $[x_1, x_2]$ contained in $[a, b]$. If we set

$$h(x) = (x - x_1)(x_2 - x)$$

for x in $[x_1, x_2]$ and $h(x) = 0$ otherwise, then $h(x)$ obviously satisfies the condition of the lemma. However,

$$\int_a^b \alpha(x)h(x)dx = \int_{x_1}^{x_2} \alpha(x)(x-x_1)(x_2-x)dx > 0$$

because the integrand is positive (except at x_1 and x_2). This contradiction proves the lemma. ■

Lemma 2 If $\alpha(x)$ is continuous in $[a, b]$, and if

$$\int_a^b \alpha(x)h'(x)dx = 0$$

for every $h(x) \in D_1(a, b)$ such that $h(a) = h(b) = 0$, then $\alpha(x) = c$ for all x in $[a, b]$, where c is a constant.

Proof. Let c be a constant given by

$$c = \frac{1}{b-a} \int_a^b \alpha(x)dx$$

which implies that

$$\int_a^b [\alpha(x) - c]dx = 0,$$

and define function $h(x)$ as

$$h(x) = \int_a^x [\alpha(\xi) - c]d\xi.$$

Evidently, $h(x)$ belongs to $D_1(a, b)$ and satisfies $h(a) = 0$ and $h(b) = 0$. Then on the one hand,

$$\int_a^b [\alpha(x) - c]h'(x)dx = \int_a^b \alpha(x)h'(x)dx - c[h(b) - h(a)] = 0.$$

while on the other hand,

$$\int_a^b [\alpha(x) - c]h'(x)dx = \int_a^b [\alpha(x) - c]^2 dx.$$

Therefore, $\alpha(x) - c = 0$ i.e. $\alpha(x) = c$ for all x in $[a, b]$. ■

Lemma 3 If $\alpha(x)$ is continuous in $[a, b]$, and if

$$\int_a^b \alpha(x)h''(x)dx = 0$$

for every $h(x) \in D_2(a, b)$ such that $h(a) = h(b) = 0$ and $h'(a) = h'(b) = 0$, then $\alpha(x) = c_0 + c_1x$ for all x in $[a, b]$, where c_0 and c_1 are constants.

Proof. Let c_0 and c_1 be constants defined by the conditions

$$\int_a^b [\alpha(x) - c_0 - c_1x] dx = 0$$

$$\int_a^b dx \int_a^x [\alpha(\xi) - c_0 - c_1\xi] d\xi = 0$$

and let

$$h(x) = \int_a^x d\xi \int_a^\xi [\alpha(t) - c_0 - c_1t] dt$$

It is easy to verify that $h(x) \in D_2(a, b)$ and satisfies $h(a) = h(b) = 0$ and $h'(a) = h'(b) = 0$. Then on the one end,

$$\begin{aligned} & \int_a^b [\alpha(x) - c_0 - c_1x] h''(x) dx \\ &= \int_a^b \alpha(x) h''(x) dx - c_0 [h'(b) - h'(a)] - c_1 \int_a^b x h''(x) dx \\ &= -c_1 [bh'(b) - ah'(a)] - c_1 [h(b) - h(a)] = 0 \end{aligned}$$

while on the other hand,

$$\int_a^b [\alpha(x) - c_0 - c_1x] h''(x) dx = \int_a^b [\alpha(x) - c_0 - c_1x]^2 dx = 0$$

It follows that $\alpha(x) - c_0 - c_1x = 0$, hence $\alpha(x) = c_0 + c_1x$ for all x in $[a, b]$. ■

Lemma 4 If $\alpha(x)$ and $\beta(x)$ are continuous in $[a, b]$, and if

$$\int_a^b [\alpha(x)h(x) + \beta(x)h'(x)] dx = 0 \quad (5)$$

for every function $h(x) \in D_1(a, b)$ such that $h(a) = h(b) = 0$, then $\beta(x)$ is differentiable, and $\beta'(x) = \alpha(x)$ for all x in $[a, b]$.

Proof. Let $A(x)$ be the function defined by

$$A(x) = \int_a^x \alpha(\xi) d\xi$$

Using integration by parts, we can verify that

$$\int_a^b \alpha(x)h(x)dx = -\int_a^b A(x)h'(x)dx$$

hence (5) can be written as

$$\int_a^b [-A(x) + \beta(x)]h'(x)dx = 0 \quad (6)$$

By Lemma 2, (6) implies that $\beta(x) - A(x) = \text{const}$, i.e., $\beta(x) = A(x) + \text{const}$, thus function $\beta(x)$ is differentiable and by definition of $A(x)$, $\beta'(x) = \alpha(x)$. ■

1.4 Variation of a Functional

Let $J[y]$ be a functional defined on some normed linear space and let

$$\Delta J[h] = J[y + h] - J[y]$$

be its increment corresponding to the increment $h = h(x)$ of the “independent variable” $y = y(x)$. If y is fixed, then $\Delta J[h]$ is a functional of h , in general a nonlinear functional. Suppose that

$$\Delta J[h] = \varphi[h] + \varepsilon \|h\|$$

where $\varphi[h]$ is a linear functional and $\varepsilon \rightarrow 0$ as $\|h\| \rightarrow 0$. Then the functional $J[y]$ is said to be *differentiable*, and the principal linear part of the increment $\Delta J[h]$, i.e., the linear functional $\varphi[h]$ which differs from $\Delta J[h]$ by an infinitesimal of order higher than 1 relative to $\|h\|$, is called the variation of $J[y]$ and is denoted by $\delta J[h]$. It can be shown that the variation of a differentiable functional is unique [1].

We now use the concept of the variation of a functional to establish a necessary condition for a functional to have an extremum.

Similar to the concept of the extremum of a function from analysis, we say that the functional $J[y]$ has a relative extremum for $y = \hat{y}$ if $J[y] - J[\hat{y}]$ does not change its sign in some neighborhood of the curve $y = \hat{y}(x)$.

At this point, it is important that the meaning of the term “neighborhood” is explicitly clarified. To this end we need to define two norms, one for space $C(a, b)$ and the other for $D_1(a, b)$:

♦ For a function $y(x)$ in space $C(a, b)$, we define

$$\|y\|_0 = \max_{a \leq x \leq b} |y(x)|$$

♦ For a function $y(x)$ in space $D_1(a, b)$, we define

$$\|y\|_1 = \max_{a \leq x \leq b} |y(x)| + \max_{a \leq x \leq b} |y'(x)|$$

Depending on which of these two norms (topology) is employed, a *neighborhood* of the curve $y = \hat{y}(x)$ may be referred to those y such that $\|y - \hat{y}\|_0 < \varepsilon$, or to those y such that $\|y - \hat{y}\|_1 < \varepsilon$. Accordingly, we shall say that the functional $J[y]$ has a weak extremum for $y = \hat{y}$, if there exist an $\varepsilon > 0$, such that $J[y] - J[\hat{y}]$ has the same sign for all y in the neighborhood of $\hat{y}(x)$ characterized by $\|y - \hat{y}\|_1 < \varepsilon$. On the other hand, we shall say that the functional $J[y]$ has a strong extremum for $y = \hat{y}$, if there exist an $\varepsilon > 0$, such that $J[y] - J[\hat{y}]$ has the same sign for all y in the neighborhood of $\hat{y}(x)$ characterized by $\|y - \hat{y}\|_0 < \varepsilon$.

Theorem 1 [1] A necessary condition for the differentiable functional $J[y]$ to have an extremum for $y = \hat{y}$ is that its variation vanish for $y = \hat{y}$, namely,

$$\delta J[h] = 0$$

for $y = \hat{y}$ and all admissible h .

Proof. To be explicit, suppose $J[y]$ has a minimum for $y = \hat{y}$. By the definition of the variation $\delta J[h]$, we have

$$\Delta J[h] = \varphi[h] + \varepsilon \|h\| \tag{7}$$

where $\varepsilon \rightarrow 0$ as $\|h\| \rightarrow 0$. Hence, for sufficiently small $\|h\|$, the sign of $\Delta J[h]$ is the same as the sign of $\delta J[h]$. Now, suppose that $\delta J[h_0] \neq 0$ for some admissible h_0 . Then for any $\alpha > 0$, no matter how small, we have

$$\delta J[-\alpha h_0] = -\delta J[\alpha h_0]$$

because $\delta J[h]$ is a linear functional of h . Hence (7) can be made to have either sign for arbitrary small $\|h\|$, which is not possible according to the hypothesis $J[y]$ has a minimum for $y = \hat{y}$. This contradiction proves the theorem. ■

1.5 Euler-Lagrange Equation: The One-Variable Case

We now study the simplest variational problem with one variable: *Let $F(x, y, z)$ be a function with continuous first and second (partial) derivatives with respect to all its variables. Then, among all functions in $D_1(a, b)$ that satisfy the boundary conditions*

$$y(a) = A, \quad y(b) = B, \quad (8)$$

find the function for which the functional

$$J[y] = \int_a^b F(x, y, y') dx \quad (9)$$

has a weak extremum.

Suppose we give $y(x)$ an increment $h(x)$. In order for the function $y(x) + h(x)$ to continue to satisfy the boundary condition (8), we must impose

$$h(a) = h(b) = 0 \quad (10)$$

The corresponding increment of the functional (9) is given by

$$\begin{aligned} \Delta J &= J[y + h] - J[y] = \int_a^b F(x, y + h, y' + h') dx - \int_a^b F(x, y, y') dx \\ &= \int_a^b [F(x, y + h, y' + h') - F(x, y, y')] dx \end{aligned}$$

Using Taylor's theorem, it follows that

$$\Delta J = \int_a^b [F_y(x, y, y')h + F_{y'}(x, y, y')h'] dx + \dots \quad (11)$$

where the subscripts denote partial derivatives with respect to the corresponding arguments, and the dots denote terms of order higher than 1 relative to h and h' . The integral in the right-hand side of (11) represents the principal linear part of the increment ΔJ , hence the variation of $J[y]$ is given by

$$\delta J = \int_a^b [F_y(x, y, y')h + F_{y'}(x, y, y')h'] dx$$

By Theorem 1, a necessary condition for $J[y]$ to have an extreme for $y = y(x)$ is that

$$\delta J = \int_a^b (F_y h + F_{y'} h') dx = 0 \quad (12)$$

for all admissible h . According to Lemma 2, (12) implies that

$$F_y - \frac{d}{dx} F_{y'} = 0 \quad (13)$$

The equation in (13) is known as the *Euler-Lagrange equation*. Thus we have proved

Theorem 2 Let $J[y]$ be a functional of the form

$$J[y] = \int_a^b F(x, y, y') dx$$

defined on the set of functions which have continuous first derivatives in $[a, b]$ and satisfy the boundary conditions $y(a) = A$, $y(b) = B$. Then a necessary condition for $J[y]$ to have an extremum for a given $y(x)$ is that $y(x)$ satisfies the Euler-Lagrange equation

$$F_y - \frac{d}{dx} F_{y'} = 0$$

1.6 Euler-Lagrange Equation: The Two-Variable Case

In this section we consider functionals depending on functions of *two variables* that assume the form of

$$J[u] = \iint_R F(x, y, u, u_x, u_y) dx dy \quad (14)$$

where R is some closed region and u_x, u_y are partial derivatives of $u = u(x, y)$.

Similar to Lemma 1 (for the one-variable case), we have

Lemma 5 If $\alpha(x, y)$ is a fixed function which is continuous in a closed region R , and if the integral

$$\iint_R \alpha(x, y) h(x, y) dx dy$$

vanishes for every function $h(x, y)$ which has continuous first and second derivatives in R and equals zero on the boundary Γ of region R , then $\alpha(x, y) = 0$ every in R .

The proof of the above lemma is similar in spirit to that of Lemma 1 and the interested reader is referred to [1].

In order to apply the necessary condition for an extremum of the functional (14), i.e., $\delta J = 0$, we must first calculate the variation δJ . Let $h(x, y)$ be an arbitrary function which has continuous first and second derivative and vanishes in the boundary Γ of R . Then if $u(x, y)$ belongs to the domain of definition, i.e., $u(x, y)$ and its first and second derivatives are continuous in R and $u(x, y)$ takes given values on boundary Γ , so does $u(x, y) + h(x, y)$. It follows that

$$\begin{aligned}
\Delta J &= J[u+h] - J[u] \\
&= \iint_R \left[F(x, y, u+h, u_x+h_x, u_y+h_y) - F(x, y, u, u_x, u_y) \right] dx dy \\
&= \iint_R \left[F_u h + F_{u_x} h_x + F_{u_y} h_y \right] dx dy + \dots\dots\dots
\end{aligned}$$

The integral in the last equality represents the principal linear part of the increment ΔJ , hence the variation of $J[u]$ is given by

$$\delta J = \iint_R \left[F_u h + F_{u_x} h_x + F_{u_y} h_y \right] dx dy.$$

By using Green's theorem

$$\iint_R \left(\frac{\partial Q}{\partial x} - \frac{\partial P}{\partial y} \right) dx dy = \int_{\Gamma} (P dx + Q dy)$$

we can write

$$\begin{aligned}
&\iint_R (F_{u_x} h_x + F_{u_y} h_y) dx dy \\
&= \iint_R \left[\frac{\partial}{\partial x} (F_{u_x} h) + \frac{\partial}{\partial y} (F_{u_y} h) \right] dx dy - \iint_R \left(\frac{\partial}{\partial x} F_{u_x} + \frac{\partial}{\partial y} F_{u_y} \right) h dx dy \\
&= \int_{\Gamma} (F_{u_x} h dy - F_{u_y} h dx) - \iint_R \left(\frac{\partial}{\partial x} F_{u_x} + \frac{\partial}{\partial y} F_{u_y} \right) h dx dy \\
&= -\iint_R \left(\frac{\partial}{\partial x} F_{u_x} + \frac{\partial}{\partial y} F_{u_y} \right) h dx dy
\end{aligned}$$

where the last equality is obtained by the fact that $h(x, y)$ assumes zero value on boundary Γ . Consequently, the variation of the functional $J[u]$ becomes

$$\delta J = \iint_R \left(F_u - \frac{\partial}{\partial x} F_{u_x} - \frac{\partial}{\partial y} F_{u_y} \right) h(x, y) dx dy \quad (15)$$

Thus the condition $\delta J = 0$ implies that the integral in (15) vanishes for any $h(x, y)$ satisfying the zero boundary condition. By Lemma 5, this leads to the second-order partial differential equation, again known as *the Euler-Lagrange equation*:

$$F_u - \frac{\partial}{\partial x} F_{u_x} - \frac{\partial}{\partial y} F_{u_y} = 0 \quad (16)$$

2. Functions of Bounded Variation

2.1 Bounded Variation Functions of One Variable

Let $u(x)$ be a function defined over $[a, b]$ on which it always assumes finite values. For a set grid points $a = x_0 < x_1 < \dots < x_n = b$, we construct the sum

$$V_u(x_0, \dots, x_n) = \sum_{i=1}^n |u(x_i) - u(x_{i-1})|$$

which is called *the variation of $u(x)$ with respect to grid set (x_0, \dots, x_n)* . If there a constant M such that

$$\sup_{x_0, \dots, x_n} V_u(x_0, \dots, x_n) \leq M < \infty$$

regardless of the number of grid points, then $u(x)$ is called a function of *bounded variation* (BV), and

$$\vee_a^b(u) = \sup_{x_0, \dots, x_n} V_u(x_0, \dots, x_n)$$

is called the total variation (TV) of $u(x)$ on $[a, b]$. In what follows, “ $u(x)$ is BV” is taken to mean that $u(x)$ is a function of BV.

Example 1 Let $u(x)$ be a function defined on $[a, b]$ that satisfies the Lipschitz condition, namely, there exists a constant M such that, for any x and x' on $[a, b]$, $|u(x) - u(x')| \leq M |x - x'|$ holds. It follows that

$$V_u(x_0, \dots, x_n) = \sum_{i=1}^n |u(x_i) - u(x_{i-1})| \leq \sum_{i=1}^n M |x_i - x_{i-1}| = M(b-a)$$

which implies that $\vee_a^b(u) \leq M(b-a)$, hence $u(x)$ is BV. In addition, if we use Lagrange's formula $u(x_1) - u(x_2) = u'(\xi) \cdot (x_1 - x_2)$ for some $\xi \in [x_1, x_2]$, then we conclude that if $u(x)$ has bounded first-order derivative on $[a, b]$, then $u(x)$ is BV. ■

Example 2 A continuous function is not necessarily BV. For example, function $u(x)$ defined by

$$u(x) = \begin{cases} x \sin \frac{1}{x} & \text{for } 0 < x \leq 1 \\ 0 & \text{for } x = 0 \end{cases}$$

is continuous on $[0, 1]$. If we take

$$x_n = \frac{1}{n\pi + \frac{\pi}{2}}, \quad n = 0, 1, \dots, m$$

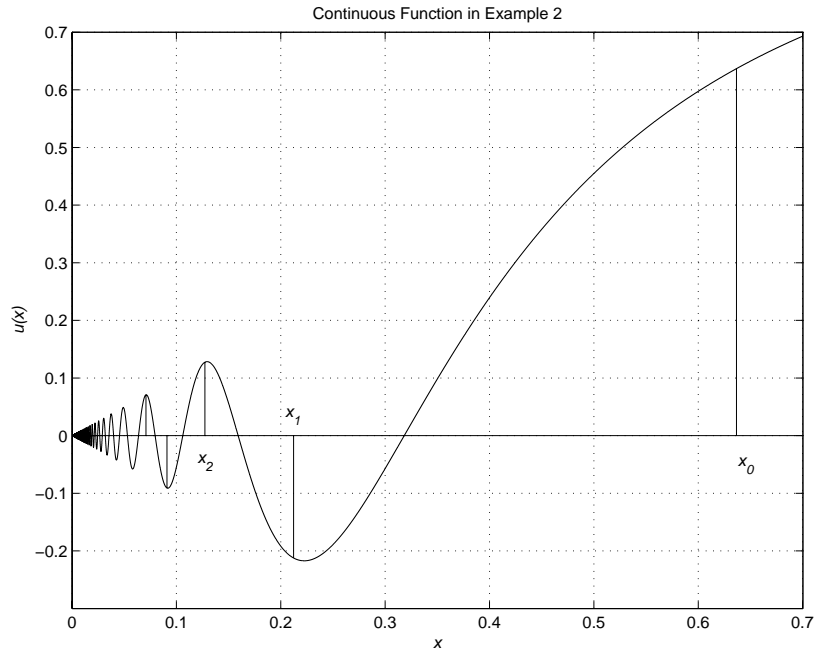
together with points 0 and 1 as the grid points, then the associated variation (see the figure below) can be estimated as

$$V_u(0, x_m, \dots, x_0, 1) \geq \sum_{i=1}^m \left| x_i \sin \frac{1}{x_i} - x_{i-1} \sin \frac{1}{x_{i-1}} \right| \geq \frac{1}{\pi} \sum_{i=1}^m \frac{1}{i}$$

Hence

$$\sup_m V_u(0, x_m, \dots, x_0, 1) \geq \frac{1}{\pi} \sum_{i=1}^m \frac{1}{i} \rightarrow \infty$$

and $u(x)$ is not BV.



■

Some properties and characterization of BV functions are listed below without proof.

- ♦ If $u(x)$ is monotonic on $[a, b]$, then $u(x)$ is BV and $\bigvee_a^b(u) = |u(a) - u(b)|$.
- ♦ If $u(x)$ is BV, then it is a bounded function, i.e., there exists a constant M such that $|u(x)| \leq M$ for $x \in [a, b]$.
- ♦ If $u(x)$ is BV and $\bigvee_a^b(u) = 0$, then $u(x)$ on $[a, b]$ is a constant.
- ♦ If $u(x)$ and $v(x)$ are BV on $[a, b]$ and α and β are constants, then $\alpha u(x) + \beta v(x)$ is BV.
- ♦ If $u(x)$ and $v(x)$ are BV on $[a, b]$, then $u(x) \cdot v(x)$ is BV.
- ♦ If $u(x)$ is BV on $[a, b]$ and c is any real value with $a < c < b$, then

$$\bigvee_a^b(u) = \bigvee_a^c(u) + \bigvee_c^b(u)$$

Conversely, if $u(x)$ is BV on $[a, c]$ and $[c, b]$, then $u(x)$ is also BV on $[a, b]$.

- ♦ If $u(x)$ is BV on $[a, b]$, then $\bigvee_a^x(u)$ is a monotonically increasing function of x on $[a, b]$.

♦ **Jordan's Decomposition Theorem** If $u(x)$ is defined on $[a, b]$ and is BV, then $u(x)$ can be expressed as the difference of two monotonically increasing functions, e.g.,

$$u(x) = \varphi(x) - \psi(x) \text{ with } \varphi(x) = \frac{1}{2} \left[\bigvee_a^x(u) + u(x) \right], \psi(x) = \frac{1}{2} \left[\bigvee_a^x(u) - u(x) \right]$$

♦ If $u(x)$ is BV on $[a, b]$, then it has at most a countable infinite number of points on $[a, b]$, where $u(x)$ is discontinuous and all discontinuous points are of first kind; $u(x)$ is Riemann-integrable; $u(x)$ has finite derivative almost everywhere on $[a, b]$; and its derivative $u'(x)$ is Lebesgue-integrable on $[a, b]$.

- ♦ If $u(x)$ is BV on $[a, b]$, then

$$\frac{d}{dx} \bigvee_a^x(u) = |u'(x)| \tag{17}$$

holds almost everywhere on $[a, b]$.

2.2 Bounded Variation Functions of Two Variables

From (17), we can write $d \int_a^x (u)'(x) dx$, hence the total variation (TV) of $u(x)$ can be computed as

$$\int_a^b (u)'(x) dx = \int_a^b |u'(x)| dx \quad (18)$$

The expression of TV in (18) provides a natural way to extend the concept of BV from the one-variable case to the two-variable case.

Definition A function $u(x, y)$ defined in a compact region R is said to be a function of bounded variation (BV) if it has a finite total variation (TV) which is defined by

$$\int_R \sqrt{u_x^2 + u_y^2} dx dy \quad (19)$$

where $u_x = \partial u(x, y) / \partial x$, $u_y = \partial u(x, y) / \partial y$. Note that the integrand in (19), i.e. $\sqrt{u_x^2 + u_y^2}$ is just the Euclid norm of the gradient of $u(x, y)$, the TV can also be written as

$$\int_R \|\nabla u(x, y)\|_2 dx dy \quad (20)$$

3. Image De-Noising by TV Minimization

3.1 Osher's Model for Images

Given a noisy image which is represented by its light intensity $u_0(x, y)$ over a bounded open set R in the two-dimensional vector space, we want to obtain a decomposition

$$u_0(x, y) = u(x, y) + n(x, y) \quad (21)$$

where $u(x, y)$ represents the true image and $n(x, y)$ is the noise with zero mean and standard deviation σ^2 . The de-noising problem at hand is to reconstruct $u(x, y)$ from $u_0(x, y)$. The model in (21) has been known as the *Osher model* [5] since the variational optimization method for image de-noising [2].

In the case of image de-blurring, model (21) is modified to

$$u_0(x, y) = (\mathcal{A}u)(x, y) + n(x, y) \quad (22)$$

where \mathcal{A} is a compact operator on $L^2(R)$ to model the de-blurring process [3].

The main difference between the methods proposed in [2][3] with many well known earlier attempts is that here images are regarded as two-variable functions of bounded variation and the

de-noising and de-blurring problems are treated by minimizing the total variation of the image subject to some variational constraints derived from the model used and the statistics of the noise component. It is this BV framework that allows analysis and processing of piecewise smooth functions with jumps as seen in virtually all natural and synthetic images and is proven proper for many basic image processing tasks.

3.2 Problem Formulation

The TV minimization problem for image de-noising can be formulated as

$$\begin{aligned} \text{minimize } \quad & \mathcal{V}_R(u) = \iint_R \sqrt{u_x^2 + u_y^2} \, dx \, dy \\ \text{subject to: } \quad & \text{constraints involving } u_0(x, y) \end{aligned} \quad (23)$$

It follows from model (21) that

$$\iint_R u(x, y) \, dx \, dy = \iint_R u_0(x, y) \, dx \, dy \quad (24)$$

and

$$\frac{1}{2} \iint_R [u(x, y) - u_0(x, y)]^2 \, dx \, dy = \sigma^2 \quad (25)$$

3.3 A Solution Procedure

By defining the *Lagrangian* of the problem (23), where the constraints are given in (24) and (25), as the functional

$$L(u, \lambda_1, \lambda_2) = \iint_R \left\{ \sqrt{u_x^2 + u_y^2} + \lambda_1(u - u_0) + \lambda_2 \left[\frac{1}{2}(u - u_0)^2 - \sigma^2 \right] \right\} dx \, dy \quad (26)$$

and applying the Euler-Lagrange equation (16) to functional (26), it can be concluded that the solution of problem (23) – (26) must satisfy the second-order partial differential equation

$$\frac{\partial}{\partial x} \left(\frac{u_x}{\sqrt{u_x^2 + u_y^2}} \right) + \frac{\partial}{\partial y} \left(\frac{u_y}{\sqrt{u_x^2 + u_y^2}} \right) - \lambda_1 - \lambda_2(u - u_0) = 0 \quad (27a)$$

subject to boundary condition

$$\frac{\partial u(x, y)}{\partial n} = 0 \quad \text{on the boundary } \Gamma = \partial R \quad (27b)$$

where n denotes the normal along the boundary of region R , which is a reasonable requirement for images defined on a rectangular region and undergone a symmetric extension along the boundary of R . Another advantage of imposing (27b) is that the Lagrange multiplier λ_1 can now be dropped from the equation because a solution of

$$\frac{\partial}{\partial x} \left(\frac{u_x}{\sqrt{u_x^2 + u_y^2}} \right) + \frac{\partial}{\partial y} \left(\frac{u_y}{\sqrt{u_x^2 + u_y^2}} \right) - \lambda(u - u_0) = 0 \quad (28a)$$

$$\frac{\partial u(x, y)}{\partial n} = 0 \quad \text{on } \Gamma = \partial R \quad (28b)$$

automatically satisfies $\iint_R (u - u_0) dx dy = 0$, i.e., constraint (24).

In [2], the solution procedure for problem (28) uses a nonlinear *parabolic* equation with time as an evolution parameter. Numerically this is equivalent to the gradient descent method. In doing so, we need to solve

$$\frac{\partial u}{\partial t} = \frac{\partial}{\partial x} \left(\frac{u_x}{\sqrt{u_x^2 + u_y^2}} \right) + \frac{\partial}{\partial y} \left(\frac{u_y}{\sqrt{u_x^2 + u_y^2}} \right) - \lambda(u - u_0) \quad \text{for } t > 0, (x, y) \in R \quad (29a)$$

$$u(x, y, 0) \text{ given} \quad (29b)$$

$$\frac{\partial u(x, y)}{\partial n} = 0 \quad \text{on } \Gamma = \partial R \quad (29c)$$

$$\frac{1}{2} \iint_R [u(x, y, t) - u_0(x, y)]^2 dx dy \equiv \sigma^2 \quad (29d)$$

where the initial value in (29b) is typically taken to be the original noisy image $u_0(x, y)$. As time t increases, the solution of the (29) converges to its “steady state” i.e. becomes time invariant thus the limiting solution of (29) satisfies (28), a de-noised version of the image.

Notice that the Lagrange multiplier λ in (29a) is *time-dependent* and must be evaluated in order to keep the time-evolution of the solution $u(x, y, t)$ going forward. By multiplying (29a) by $(u - u_0)$, using (29d), integrating by parts over region R and using (29c), and assuming that the steady state has been reached, we obtain a formula for updating the Lagrange multiplier $\lambda(t)$ as

$$\lambda(t) = -\frac{1}{2\sigma^2} \iint_R \left[\sqrt{u_x^2 + u_y^2} - \left(\frac{(u_0)_x \cdot u_x}{\sqrt{u_x^2 + u_y^2}} + \frac{(u_0)_y \cdot u_y}{\sqrt{u_x^2 + u_y^2}} \right) \right] dx dy \quad (30)$$

3.4 A Difference Scheme for Solving (29), (30)

Given below is a difference scheme proposed in [2] for computing a numerical solution of problem (28) by solving (30) on an evolving time grids. For the sake of notation simplicity, in the difference scheme a *square* region $R = [0, 1] \times [0, 1]$ has been assumed. The scheme can be applied to non-square images after some straightforward modifications.

$$x_i = ih, \quad y_j = jh, \quad \text{for } i, j = 0, 1, \dots, N, \quad \text{with } Nh = 1$$

$$t_n = n\Delta t, \quad \text{for } n = 0, 1, \dots$$

$$u_{i,j}^n = u(x_i, y_j, t_n)$$

$$u_{i,j}^0 = u_0(x_i, y_j)$$

For $i, j = 0, 1, \dots, N-1$, do

$$u_{i,j}^{n+1} = u_{i,j}^n$$

$$+ \frac{\Delta t}{h} \left[\Delta_-^x \left(\frac{\Delta_+^x u_{i,j}^n}{\left[(\Delta_+^x u_{i,j}^n)^2 + (m(\Delta_+^y u_{i,j}^n, \Delta_-^y u_{i,j}^n))^2 \right]^{1/2}} \right) \right]$$

$$+ \frac{\Delta t}{h} \left[\Delta_-^y \left(\frac{\Delta_+^y u_{i,j}^n}{\left[(\Delta_+^y u_{i,j}^n)^2 + (m(\Delta_+^x u_{i,j}^n, \Delta_-^x u_{i,j}^n))^2 \right]^{1/2}} \right) \right]$$

$$- \Delta t \cdot \lambda^n (u_{i,j}^n - u_0(ih, jh))$$

with

$$u_{0,j}^n = u_{1,j}^n, u_{N,j}^n = u_{N-1,j}^n, u_{i,0}^n = u_{i,1}^n, \text{ and } u_{i,N}^n = u_{i,N-1}^n.$$

$$\Delta_{\mp}^x u_{i,j} = \mp(u_{i\mp 1,j} - u_{i,j})$$

$$\Delta_{\mp}^y u_{i,j} = \mp(u_{i,j\mp 1} - u_{i,j})$$

$$m(a, b) = \min \text{mod}(a, b) = \left(\frac{\text{sgn } a + \text{sgn } b}{2} \right) \min(|a|, |b|)$$

where λ^n is computed based on (30) as

$$\lambda^n = -\frac{h}{2\sigma^2} \left[\sum_i \sum_j \sqrt{(\Delta_+^x u_{i,j}^n)^2 + (\Delta_+^y u_{i,j}^n)^2} - \frac{(\Delta_+^x u_{i,j}^0)(\Delta_+^x u_{i,j}^n)}{\sqrt{(\Delta_+^x u_{i,j}^n)^2 + (\Delta_+^y u_{i,j}^n)^2}} - \frac{(\Delta_+^y u_{i,j}^0)(\Delta_+^y u_{i,j}^n)}{\sqrt{(\Delta_+^x u_{i,j}^n)^2 + (\Delta_+^y u_{i,j}^n)^2}} \right]$$

A step size restriction is imposed for numerical stability of the scheme:

$$\frac{\Delta t}{h^2} \leq c \quad \text{for some constant } c \quad (31)$$

3.5 MATLAB Code and Results

Listed below are several MATLAB functions written by Guy Gilboa of UCLA that can be used to demonstrate the Rudin-Osher-Fatemi algorithm [2]. The main routine named `demo_tv.m` is given in the first m-file which calls for the second and third MATLAB routines.

1. File Name: `demo_tv.m`

```
%%% TV demo %%%
%%% By Guy Gilboa
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
% Based on: [ROF92] L. Rudin, S. Osher, E. Fatemi,
% "Nonlinear Total Variation based noise removal algorithms",
% Physica D 60 259-268, 1992.
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
%%% Simple TV denoising %%%
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
I = imread('camera256.tif'); % load image
I = double(I(11:138,11:138)); % cut a piece, convert to double

% params
iter = 80;
% dt=0.2; eps=1;

%%% Add noise
std_n = 20; % Gaussian noise standard deviation
In = randn(size(I))*std_n; % White Gaussian noise
I0 = I + In; % noisy input image
% show original and noisy images
close all
figure(1); imshow(uint8(I)); title('Original')
figure(2); imshow(uint8(I0)); title('Noisy image')
% denoise image by using tv for some iterations
J = tv(I0,iter);
figure(3); imshow(uint8(J)); title('Denoised image')

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
%%% TV denoising with (scalar) data fidelity term %%%
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
%% Here we assume the noise variance is known.
%% Therefore we can compute lambda (instead of fixing
%% iter). TV therefore can run automatically quite well.
%% The process is slower.
J = I0;
% params
ep_J = 0.01; % minimum mean change in image J
lam = 0; J_old = 0;
iter = 10; dt = 0.2; eps = 1;
```

```

var_n = std_n^2; % noise variance
i = 0;
while (mean(mean(abs(J - J_old))) > ep_J), % iterate until convergence
    J_old = J;
    J = tv(J,iter,dt,eps,lam,I0); % scalar lam
    lam = calc_lam(J,I0,var_n,eps); % update lambda (fidelity term)
end % for i
figure(4); imshow(uint8(J)); title('Denoised image with lambda')

```

2. File Name: tv.m

```

function J = tv(I,iter,dt,ep,lam,I0,C)

%% Private function: tv (by Guy Gilboa).
%% Total Variation denoising.
%% Example: J = tv(I,iter,dt,ep,lam,I0)
%% Input: I - image (double array gray level 1-256),
%%        iter - num of iterations,
%%        dt - time step [0.2],
%%        ep - epsilon (of gradient regularization) [1],
%%        lam - fidelity term lambda [0],
%%        I0 - input (noisy) image [I0=I]
%%        (default values are in [])
%% Output: evolved image

if ~exist('ep')
    ep = 1;
end
if ~exist('dt')
    dt = ep/5; % dt below the CFL bound
end
if ~exist('lam')
    lam = 0;
end
if ~exist('I0')
    I0 = I;
end
if ~exist('C')
    C = 0;
end
[ny,nx] = size(I);
ep2 = ep^2;

for i=1:iter, %% do iterations
    % estimate derivatives
    I_x = (I(:,[2:nx nx]) - I(:,[1 1:nx-1]))/2;
    I_y = (I([2:ny ny],:) - I([1 1:ny-1],:))/2;
    I_xx = I(:,[2:nx nx]) + I(:,[1 1:nx-1]) - 2*I;
    I_yy = I([2:ny ny],:) + I([1 1:ny-1],:) - 2*I;
    Dp = I([2:ny ny],[2:nx nx]) + I([1 1:ny-1],[1 1:nx-1]);
    Dm = I([1 1:ny-1],[2:nx nx]) + I([2:ny ny],[1 1:nx-1]);

```

```

    I_xy = (Dp-Dm)/4;
% compute flow
    Num = I_xx.*(ep2+I_y.^2)-2*I_x.*I_y.*I_xy+I_yy.*(ep2+I_x.^2);
    Den = (ep2+I_x.^2+I_y.^2).^(3/2);
    I_t = Num./Den + lam.*(I0-I+C);
    I = I+dt*I_t; %% evolve image by dt
end % for i
%% return image
%J = I*Imean/mean(mean(I)); % normalize to original mean
J = I;

```

3. File Name: calc_lam.m

```

function lam = calc_lam(I,I0,var_n,ep)

%% private function: calc_lam (by Guy Gilboa).
%% calculate scalar fidelity term of the Total-Variation
%% input: current image, original noisy image, noise variance, epsilon
%% output: lambda
%% example: lam=calc_lam(I,I0,var_n,ep)

if ~exist('ep') % good for 256 gray-level
    ep = 1;
end

[ny,nx] = size(I); ep2 = ep^2;

% estimate derivatives
I_x = (I(:,[2:nx nx])-I(:,[1 1:nx-1]))/2;
I_y = (I([2:ny ny],:)-I([1 1:ny-1],:))/2;
I_xx = I(:,[2:nx nx])+I(:,[1 1:nx-1])-2*I;
I_yy = I([2:ny ny],:)+I([1 1:ny-1],:)-2*I;
    Dp = I([2:ny ny],[2:nx nx])+I([1 1:ny-1],[1 1:nx-1]);
    Dm = I([1 1:ny-1],[2:nx nx])+I([2:ny ny],[1 1:nx-1]);
    I_xy = (Dp-Dm)/4;
% compute numerator and denominator
    Num = I_xx.*(ep2+I_y.^2)-2*I_x.*I_y.*I_xy+I_yy.*(ep2+I_x.^2);
    Den = (ep2+I_x.^2+I_y.^2).^(3/2);
    Div = Num./Den;
    %%%% fidelity term
lam = mean(mean(Div.*(I-I0)))./var_n;

```

The above MATLAB code was tested by applying it to a small part of a test image named cameraman. The image has a size of 128 by 128 and was taken from the tif file camera256(11:138,11:138) and was contaminated by a Gaussian white noise with zero mean and standard deviation = 20.

The results in term of the original, noisy, and de-noised images are depicted in the following plots. The peak-signal-to-noise ratio (PSNR) was improved from 22.0647 dB to 29.5268 dB.

Original



Noisy image



Denoised image with lambda



For comparison purposes, Donoho's well-known wavelet de-noising algorithm using Daubechies wavelets of length 2, 4, and 6 were also applied to the same test image. The best result was achieved by the D2 (Haar) wavelet that yielded PSNR = 27.0670dB.

The plot below displays the de-noised image by wavelet (on the left) and de-noised image by TV (on the right) minimization. It is observed that the TV-minimization based algorithm yields improved de-noising results in terms of PSNR as well as visual examination.



4. Image De-Blurring by TV Minimization

In [3], the image de-blurring problem is investigated in a similar variational optimization framework. Here the image model is given by (22), i.e.,

$$u_0(x, y) = (\mathcal{A}u)(x, y) + n(x, y) \quad \text{for } (x, y) \text{ in } R \quad (32)$$

where \mathcal{A} is typically a convolutional type integral operator that can be used for modeling several common blurring processes such as averaging, Gaussian low-pass, Laplacian of Gaussian, and motion blurs.

The variational optimization problem in this case can be formulated as

$$\begin{aligned} \text{minimize} \quad & \mathcal{V}_R(u) = \iint_R \sqrt{u_x^2 + u_y^2} \, dx \, dy \\ \text{subject to:} \quad & \iint_R u(x, y) \, dx \, dy = \iint_R u_0(x, y) \, dx \, dy \\ & \iint_R [\mathcal{A}u(x, y) - u_0(x, y)]^2 \, dx \, dy = \sigma^2 \end{aligned} \quad (33)$$

The Euler-Lagrange equation corresponding to problem (33) is given by

$$\frac{\partial}{\partial x} \left(\frac{u_x}{\sqrt{u_x^2 + u_y^2}} \right) + \frac{\partial}{\partial y} \left(\frac{u_y}{\sqrt{u_x^2 + u_y^2}} \right) - \lambda A^* (Au - u_0) = 0 \quad (34a)$$

$$\frac{\partial u(x, y)}{\partial n} = 0 \quad \text{on } \Gamma = \partial R \quad (34b)$$

where \mathcal{A}^* denotes the *adjoint operator* for operator \mathcal{A} . In [3], it is proposed that (34) be solved by calculating the steady-state solution of the parabolic partial differential equation

$$\frac{\partial u}{\partial t} = \frac{\partial}{\partial x} \left(\frac{u_x}{\sqrt{u_x^2 + u_y^2}} \right) + \frac{\partial}{\partial y} \left(\frac{u_y}{\sqrt{u_x^2 + u_y^2}} \right) - \lambda A^* (Au - u_0) \quad \text{for } t > 0, (x, y) \in R \quad (35a)$$

$$\frac{\partial u(x, y)}{\partial n} = 0 \quad \text{on } \Gamma = \partial R \quad (35b)$$

where, using a technique similar to that for the de-noising case, the time-varying Lagrange multiplier $\lambda = \lambda(t)$ is found to be

$$\lambda(t) = \frac{\iint_R \left[\frac{\partial}{\partial x} \left(\frac{u_x}{\sqrt{u_x^2 + u_y^2}} \right) + \frac{\partial}{\partial y} \left(\frac{u_y}{\sqrt{u_x^2 + u_y^2}} \right) \right] \cdot A^* (Au - u_0) dx dy}{\iint_R [A^* (Au - u_0)]^2 dx dy} \quad (36)$$