# Imputing a Variational Inequality Function or a Convex Objective Function: a Robust Approach

*Jerome Thai*

Electrical Engineering and Computer Sciences
University of California at Berkeley

January 17, 2017

# Imputing a Variational Inequality Function or a Convex Objective Function: a Robust Approach

by

Jérôme Thai

A technical report submitted in partial satisfaction of the
requirements for the degree of
Master of Science

in

Electrical Engineering and Computer Sciences

at the

University of California, Berkeley

Approved by:

Professor Alexandre Bayen, Research Advisor
Professor Laurent El Ghaoui

Fall 2016

The dissertation of Jérôme Thai, titled Imputing a Variational Inequality Function or a Convex Objective Function:
a Robust Approach, is approved:

 

Date _____

Date _____

University of California, Berkeley

# Imputing a Variational Inequality Function or a Convex Objective Function: a Robust Approach

# Abstract

Imputing a Variational Inequality Function or a Convex Objective Function:
a Robust Approach

by

Jérôme Thai

Master of Science in Electrical Engineering and Computer Sciences

University of California, Berkeley

Professor Alexandre Bayen, Chair

To impute the function of a variational inequality and the objective of a convex optimization problem from observations of (nearly) optimal decisions, previous approaches constructed inverse programming methods based on solving a convex optimization problem [17, 7]. However, we show that, in addition to requiring complete observations, these approaches are not robust to measurement errors, while in many applications, the outputs of decision processes are noisy and only partially observable from, *e.g.*, limitations in the sensing infrastructure. To deal with noisy and missing data, we formulate our inverse problem as the minimization of a weighted sum of two objectives: 1) a duality gap or Karush-Kuhn-Tucker (KKT) residual, and 2) a distance from the observations robust to measurement errors. In addition, we show that our method encompasses previous ones by generating a sequence of Pareto optimal points (with respect to the two objectives) converging to an optimal solution of previous formulations. To compare duality gaps and KKT residuals, we also derive new sub-optimality results defined by KKT residuals. Finally, an implementation framework is proposed with applications to delay function inference on the road network of Los Angeles, and consumer utility estimation in oligopolies.

# Contents

# Chapter 1

# Introduction

## 1.1 Motivation

Many decision processes are modeled as a Variational Inequality (VI) or Convex Optimization (CO) problem [15, 8]. However, the function that describes these processes are often difficult to estimate while their outputs (the decisions they describe) are often directly observable. For example, the traffic assignment problem considers a road network in which each road segment is associated to a delay that is a function of the volume of traffic on the arc [22]. The Wardrop's equilibrium principles [25] describe an equilibrium flow that is easily locally measurable by induction loop detectors or video cameras. While the delay functions are in general not observable, having accurate estimates of these functions is still crucial for urban planning. However, due to their cost of maintenance, traffic sensors are sparse, we thus present an approach robust to missing values and measurement errors. In consumer utility estimation, for example, the consumer is assumed to purchase various products from different companies in order to maximize a utility function minus the price paid, where the utility function measures the satisfaction the consumer receives from his purchases. In practice, the consumer's utility function is difficult to estimate but the consumer purchases, which is a function of the products' prices, are easily observable. We refer to [17, 7] for more examples, *e.g.*, value function estimation control.

## 1.2 Contributions and outline

Estimating the parameters of a process based on observations is related to various lines of work, *e.g.*, inverse reinforcement learning in robotics [20, 1], the inverse shortest path problem [10], recovering the parameters of the Lyapunov function given a linear control policy [9, §10.6]. The field of *structural estimation* in economics estimates the parameters of observed equilibrium models, *e.g.* imputing production and demand functions [23, 2, 4]. In general, *inverse problems* have been studied quite extensively and we refer to [17, 7] for more references on the subject. In [17] (resp. [7]), a program is proposed to impute a convex

objective (resp. a VI function) based on complete observations of nearly optimal decisions. The program is solved via CO.

After reviewing preliminary results in VI and CO in Section 2 and formally stating the problem in Section 3, our contributions in the remainder of the present article is as follows. In Section 4, we demonstrate that the methods presented in [17, 7] are in general not robust to noise and outliers in the data. In Section 5, we formulate our inverse problem as a weighted sum of a distance $r_{\text{obs}}$ from the observations and residual functions $r_{\text{eq}}$ in the form of duality gaps or Karush-Kuhn-Tucker (KKT) residuals, and show that our method is robust to noise and outliers while it avoids the disjunctive nature of the complementary condition. In Section 6, we show that the proposed weighted sum defines a set of Pareto efficient points whose closure contains a solution to the programs proposed in [17, 7]. Our method thus encompasses previous ones but performs better against noise and missing data. It also provides a conceptual way to recognize the implicit assumption of full noiseless observations made by previous inverse programming approaches. In Section 7, we compare the KKT residual and the duality gap and derive new sub-optimality results defined by the KKT residuals. In Section 8, an implementation framework is proposed. Finally, we apply our method to delay inference in the road network of Los Angeles, and consumer utility estimation and pricing in oligopolies in Sections 9 and 10.

# Chapter 2

# Preliminaries

## 2.1 Variational Inequality (VI) and Convex Optimization (CO)

VI is used to model a broad class of problems from economics, convex optimization, and game theory, see, *e.g.* [15], for a comprehensive treatment of the subject. Mathematically, a VI problem is defined as follows:

**Definition 2.1.** *Given a closed, convex set $\mathcal{K} \subseteq \mathbb{R}^n$ and a map $F : \mathcal{K} \to \mathbb{R}^n$, the VI problem, denoted VI($\mathcal{K}, F$), consists in finding a vector $\mathbf{x} \in \mathcal{K}$ such that*

$$F(\mathbf{x})^T(\mathbf{u} - \mathbf{x}) \geq 0, \, \forall \, \mathbf{u} \in \mathcal{K} \tag{2.1}$$

For the remainder of the article, we suppose that $\mathcal{K}$ is a polyhedron, written in standard form:[1]

$$\mathcal{K} = \{\mathbf{x} \in \mathbb{R}^n \,|\, \mathbf{Ax} = \mathbf{b}, \, \mathbf{x} \geq 0\} \tag{2.2}$$

This allows different characterizations of solutions to VI($\mathcal{K}, F$). We define the *primal-dual system* associated to the Linear Program (LP) $\min_{u \in \mathcal{K}} F(\mathbf{x})^T \mathbf{u}$:

**Definition 2.2.** *(See [3, Th. 1].) Given VI($\mathcal{K}, F$), and $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^n \times \mathbb{R}^n$, we define the associated primal-dual system as follows:*

$$\begin{aligned} F(\mathbf{x})^T\mathbf{x} &= \mathbf{b}^T\mathbf{y} \\ \mathbf{A}^T\mathbf{y} &\leq F(\mathbf{x}) \\ \mathbf{Ax} = \mathbf{b}, \, \mathbf{x} &\geq 0 \end{aligned} \tag{2.3}$$

In the above system, we say that $\mathbf{x}$ is *primal feasible* if $\mathbf{Ax} = \mathbf{b}$, $\mathbf{x} \geq 0$, and $(\mathbf{x}, \mathbf{y})$ is *dual feasible* if $\mathbf{A}^T\mathbf{y} \leq F(\mathbf{x})$. From LP strong duality, we have:

---

[1]The results in the present work can be generalized to conic-representable sets, but we choose to restrict $\mathcal{K}$ to a polyhedron for ease of notation. A discussion on the generalization is presented in Section 11.

**Theorem 2.1.** *(See [3, Th. 1].) Let $\mathcal{K}$ be a polyhedron given by (2.2).  Then $\mathbf{x} \in \mathbb{R}^n$ solves VI($\mathcal{K}, F$) if and only if there exists $\mathbf{y} \in \mathbb{R}^n$ such that the pair $(\mathbf{x}, \mathbf{y})$ satisfies the primal-dual system (2.3).*

We also define the *Karush-Kuhn-Tucker* (KKT) system of the VI($\mathcal{K}, F$):

**Definition 2.3.** *Let $\mathcal{K}$ be a polyhedron given by (2.2).  Given a map $F$ and $(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \in \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^n$, we define the associated KKT system as follows:*

$$
\begin{aligned}
F(\mathbf{x}) &= \mathbf{A}^T \mathbf{y} + \boldsymbol{\pi} \\
\mathbf{A}\mathbf{x} &= \mathbf{b} \\
\mathbf{x} &\geq 0,\ \boldsymbol{\pi} \geq 0,\ \mathbf{x}^T \boldsymbol{\pi} = 0
\end{aligned}
\tag{2.4}
$$

**Theorem 2.2.** *(See [8, §5.5.3].) Let $\mathcal{K}$ be a polyhedron given by (2.2).  Then a vector $\mathbf{x} \in \mathbb{R}^n$ solves VI($\mathcal{K}, F$) if and only if there exists $\mathbf{y}, \boldsymbol{\pi} \in \mathbb{R}^n$ such that the tuple $(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi})$ satisfies the KKT system (2.4).*

Convex Optimization (CO) is closely related to VI, see [8] for a comprehensive treatment on the subject.  A CO problem is defined as follows:

**Definition 2.4.** *Given a closed, convex set $\mathcal{K} \in \mathbb{R}^n$ and a convex potential $f : \mathcal{K} \to \mathbb{R}$, the CO problem, denoted CO($\mathcal{K}, f$), is a program of the form:*

$$
\min\ f(\mathbf{x}) \quad s.t. \quad \mathbf{x} \in \mathcal{K}
\tag{2.5}
$$

We have the following optimality condition to the CO($\mathcal{K}, f$):

**Theorem 2.3.** *(See [8, §4.2.3].) Given CO($\mathcal{K}, f$), suppose $f$ differentiable.  Then a vector $\mathbf{x} \in \mathcal{K}$ is an optimal solution to CO($\mathcal{K}, f$) if and only if:*

$$
\nabla f(\mathbf{x})^T (\mathbf{u} - \mathbf{x}) \geq 0,\ \forall\,\mathbf{u} \in \mathcal{K}
\tag{2.6}
$$

Hence VI($\mathcal{K}, F$) can be seen as a generalization of CO($\mathcal{K}, f$) where the gradient $\nabla f$ is substituted by a general map $F$.  Hence, when $f$ is differentiable, the primal-dual and KKT systems are both optimality conditions for the CO($\mathcal{K}, f$).

## 2.2  Approximate solutions

We now focus on the VI($\mathcal{K}, F$) since it encompasses CO($\mathcal{F}, f$).  The residual functions associated to the primal-dual and KKT systems are defined as

**Definition 2.5.** *(See [7].) Given VI($\mathcal{K}, F$), a residual function $r_{PD}$ of the primal-dual system (2.3) is a non-negative function which satisfies for all $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^n \times \mathbb{R}^n$ such that $\mathbf{A}\mathbf{x} = \mathbf{b}$, $\mathbf{x} \geq 0$, $\mathbf{A}^T \mathbf{y} \leq F(\mathbf{x})$:*

$$
r_{PD}(\mathbf{x}, \mathbf{y}) = 0 \quad \Longleftrightarrow \quad \text{(2.3) holds at } (\mathbf{x}, \mathbf{y})
\tag{2.7}
$$

**Definition 2.6.** *(See [17].) Given VI($\mathcal{K}, F$), a residual function $r_{KKT}$ of the primal-dual system* (2.4) *is a non-negative function which satisfies for all* $(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \in \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^n$ *such that* $\mathbf{A}\mathbf{x} = \mathbf{b}$, $\mathbf{x} \geq 0$, $\boldsymbol{\pi} \geq 0$, $\mathbf{A}^T\mathbf{y} \leq F(\mathbf{x})$

$$r_{KKT}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) = 0 \quad \Longleftrightarrow \quad (2.4) \text{ holds at } (\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \tag{2.8}$$

Residual functions are used as sub-optimality certificates in iterative methods for solving VI($\mathcal{K}, F$) and CO($\mathcal{K}, f$). As in [7] and [17], we specify $r_{\text{PD}}$ and $r_{\text{KKT}}$ as follows:

$$r_{\text{PD}}(\mathbf{x}) = F(\mathbf{x})^T\mathbf{x} - \mathbf{b}^T\mathbf{y} \tag{2.9}$$

$$r_{\text{KKT}}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) = \left\| \begin{bmatrix} \alpha(F(\mathbf{x}) - \mathbf{A}^T\mathbf{y} - \boldsymbol{\pi}) \\ \mathbf{x} \circ \boldsymbol{\pi} \end{bmatrix} \right\|_1 \tag{2.10}$$

where $\mathbf{x} \circ \boldsymbol{\pi} = [x_i \pi_i]_{i=1}^n$, $\|x\|_1 = \sum_{i=1}^n |x_i|$, and $\alpha > 0$ a weighting factor. The choice of $r_{\text{PD}}$ in (2.9) is natural since primal feasibility ($\mathbf{A}\mathbf{x} = \mathbf{b}$, $\mathbf{x} \geq 0$) and dual feasibility ($\mathbf{A}^T\mathbf{y} \leq F(\mathbf{x})$) imply that $r_{\text{PD}}$ is non-negative from weak LP duality [3, Cor. 1], and it is tied to the following optimality gaps for VI($\mathcal{K}, F$) and CO($\mathcal{K}, f$), taken from [15, §3.1.5] and [8, §9.3.1] respectively, for all $\mathbf{x} \in \mathcal{K}$:

$$r_{\text{VI}}(\mathbf{x}) = \max_{\mathbf{u} \in \mathcal{K}} F(\mathbf{x})^T(\mathbf{x} - \mathbf{u}) \tag{2.11}$$

$$r_{\text{CO}}(\mathbf{x}) = f(\mathbf{x}) - \min_{\mathbf{u} \in \mathcal{K}} f(\mathbf{u}) \tag{2.12}$$

**Theorem 2.4.** *(See [7, Th. 2].) Let $\mathcal{K}$ be a polyhedron given by* (2.2). *Then the following holds for any $\epsilon \geq 0$ and $\mathbf{x} \in \mathcal{K}$:*

$$r_{VI}(\mathbf{x}) \leq \epsilon \Longleftrightarrow \exists \mathbf{y} \in \mathbb{R}^n : \mathbf{A}^T\mathbf{y} \leq F(\mathbf{x}), r_{PD}(\mathbf{x}, \mathbf{y}) \leq \epsilon \tag{2.13}$$

*In addition, if $F$ is the gradient of a convex potential $f$, then, for all $\mathbf{x} \in \mathcal{K}$:*

$$r_{VI}(\mathbf{x}) \leq \epsilon \Longrightarrow r_{CO}(\mathbf{x}) \leq \epsilon \tag{2.14}$$

When primal and dual feasibilities hold, $r_{\text{PD}} \leq \epsilon$ is equivalent to $\epsilon$-suboptimality for VI($\mathcal{K}, F$) with respect to $r_{\text{VI}}$. When $f = \nabla F$, $r_{\text{PD}} \leq \epsilon$ is sufficient for $\epsilon$-suboptimality for CO($\mathcal{K}, f$) with respect to $r_{\text{CO}}$, but not necessary. To see this, consider a quadratic function $f : \mathbb{R} \to \mathbb{R}$ with minimum attained at $a > 0$:

$$\mathcal{K} = \mathbb{R}_+, \quad f(x) = (x - a)^2, \quad F(x) = \nabla f(x) = 2(x - a) \tag{2.15}$$

so $r_{\text{CO}}(a + \epsilon) = f(a + \epsilon) = \epsilon^2$ while $r_{\text{VI}}(a + \epsilon) = r_{\text{PD}}(a + \epsilon) = (a + \epsilon)F(a + \epsilon) = 2(a + \epsilon)\epsilon$ is arbitrarily large as $a$ goes to $+\infty$.

## 2.3 Distance from solutions

Assume VI$(\mathcal{K}, F)$ (resp. CO$(\mathcal{K}, f)$) has a unique solution $\mathbf{x}^\star$. An alternative sub-optimality condition is that $\|\mathbf{x} - \mathbf{x}^\star\| < \epsilon$ for $\mathbf{x} \in \mathcal{K}$. Main results rely on of strict and strong monotonicity of $F$ (resp. convexity of $f$):

**Definition 2.7.** *Given a convex set $\mathcal{K} \subseteq \mathbb{R}^n$ and a function $f : \mathcal{K} \to \mathbb{R}$, $f$ is said to be strictly convex on $\mathcal{K}$ if; $\forall \mathbf{x}, \mathbf{x}' \in \mathcal{K}$ and $\alpha \in (0, 1)$ such that $\mathbf{x} \neq \mathbf{x}'$*

$$f(\alpha \mathbf{x} + (1 - \alpha)\mathbf{x}') < \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{x}') \tag{2.16}$$

*is said to be strongly convex on $\mathcal{K}$ if; $\exists c > 0$ such that $\forall \alpha \in (0, 1)$, $\forall \mathbf{x}, \mathbf{x}' \in \mathcal{K}$:*

$$f(\alpha \mathbf{x} + (1 - \alpha)\mathbf{x}') \leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{x}') - \frac{c}{2}\alpha(1 - \alpha)\|\mathbf{x} - \mathbf{x}'\|^2 \tag{2.17}$$

**Definition 2.8.** *Given a convex set $\mathcal{K} \subseteq \mathbb{R}^n$ and a map $F : \mathcal{K} \to \mathbb{R}^n$, $F$ is said to be strictly monotone on $\mathcal{K}$ if*

$$(F(\mathbf{x}) - F(\mathbf{x}'))^T(\mathbf{x} - \mathbf{x}') \geq 0, \quad \forall \mathbf{x}, \mathbf{x}' \in \mathcal{K} \tag{2.18}$$

*strongly monotone on $\mathcal{K}$ if $\exists c > 0$ such that*

$$(F(\mathbf{x}) - F(\mathbf{x}'))^T(\mathbf{x} - \mathbf{x}') \geq c\|\mathbf{x} - \mathbf{x}'\|^2, \quad \forall \mathbf{x}, \mathbf{x}' \in \mathcal{K} \tag{2.19}$$

When $f$ is differentiable, $f$ is strictly (resp. strongly) convex is equivalent to $\nabla f$ is strictly (resp. strongly) monotone. Strong monotonicity allows us to bound $\|\mathbf{x} - \mathbf{x}^\star\|$ by the residual $r_{\text{VI}}(\mathbf{x})$ in (2.11):

**Theorem 2.5.** *(See [21, Th. 4.1].) If VI$(\mathcal{K}, F)$ is such that $\mathcal{K} \subseteq \mathbb{R}^n$ is closed convex and $F$ strongly monotone, VI$(\mathcal{K}, F)$ admits a unique solution $\mathbf{x}^\star$ and:*

$$\|\mathbf{x} - \mathbf{x}^\star\|_2 \leq \sqrt{r_{VI}(\mathbf{x})/c}, \quad \forall \mathbf{x} \in \mathcal{K} \tag{2.20}$$

*in addition, if $\exists f : F = \nabla f$, then $\mathbf{x}^\star$ is the unique solution to CO$(\mathcal{K}, f)$ and:*

$$\|\mathbf{x} - \mathbf{x}^\star\|_2 \leq \sqrt{2\,r_{CO}(\mathbf{x})/c}, \quad \forall \mathbf{x} \in \mathcal{K} \tag{2.21}$$

If $F$ is only strictly monotone, then VI$(\mathcal{K}, F)$ admits at most one solution [24]. If the solution $\mathbf{x}^\star$ exists, then strict monotonicity is not strong enough for a bound similar to (2.20).

# Chapter 3

# Problem statement

We present our problem statement in the most general case. We refer to Sections 9 and 10 for illustration of the problem in traffic assignment and consumer utility respectively. Let us consider a process in which decisions $\mathbf{x}$ are made by solving a parametric variational inequality $\mathrm{VI}(\mathcal{K}(\mathbf{p}), F(\cdot, \mathbf{p}))$, for a set of parameter values $\mathbf{p} \in \mathcal{P}$:

$$F(\mathbf{x}, \mathbf{p})^T(\mathbf{u} - \mathbf{x}) \geq 0, \quad \forall \mathbf{u} \in \mathcal{K}(\mathbf{p}) \tag{3.1}$$

$$\mathcal{K}(\mathbf{p}) := \{\mathbf{x} \in \mathbb{R}^n \ : \ \mathbf{A}(\mathbf{p})\mathbf{x} = \mathbf{b}(\mathbf{p}), \ \mathbf{x} \geq 0\} \tag{3.2}$$

where both the map $F(\cdot, \mathbf{p})$ and polyhedron $\mathcal{K}(\mathbf{p})$ depend on $\mathbf{p}$. The definitions and theorems in Section 2 apply for each $\mathbf{p} \in \mathcal{P}$, and the dependence of the residuals on $\mathbf{p}$ are made explicit with $r_{\mathrm{PD}}(\mathbf{x}, \mathbf{y}, \mathbf{p})$, $r_{\mathrm{KKT}}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}, \mathbf{p})$, etc.

**Inputs:** We are given $\mathbf{A}(\mathbf{p})$, $\mathbf{b}(\mathbf{p})$ for all $\mathbf{p}$, along with a parametric observation process $g(\cdot, \mathbf{p}) : \mathbb{R}^n \to \mathbb{R}^q$ and $N$ noisy observations

$$\mathbf{z}^{(j)} := g(\mathbf{x}^{(j)}, \mathbf{p}^{(j)}) + \mathbf{w}^{(j)}, \quad j = 1, \cdots, N \tag{3.3}$$

of (approximate) solutions $\mathbf{x}^{(j)}$ to $\mathrm{VI}(\mathcal{K}(\mathbf{p}^{(j)}), F(\cdot, \mathbf{p}^{(j)}))$ with random noise $\mathbf{w}^{(j)} \in \mathbb{R}^q$ and associated parameters $\mathbf{p}^{(j)}$. Unless $g(\cdot, \mathbf{p})$ is an injection from $\mathcal{K}(\mathbf{p})$ to $\mathbb{R}^q$ for all $\mathbf{p}$, the observation $\mathbf{z}^{(j)}$ contains in general less information than $\mathbf{x}^{(j)}$, thus (3.3) is our missing data model.

**Objective:** We want to impute the parametric map $F(\cdot, \mathbf{p})$ and the decision vectors $\mathbf{x}^{(j)}$ such that, for all $j$:

(a) $\mathbf{x}^{(j)}$ is an approximate solution to $\mathrm{VI}(\mathcal{K}(\mathbf{p}^{(j)}), F(\mathbf{p}^{(j)}))$.

(b) $\mathbf{x}^{(j)}$ agrees with the observations $\mathbf{z}^{(j)}$.

**Formalization:** Using Theorem 2.4, objective (a) consists in imputing a parametric map $F(\cdot, \mathbf{p})$ and a collection of decision vectors $\mathbf{x}^{(j)} \in \mathcal{K}(\mathbf{p}^{(j)})$, along with dual variables $\mathbf{y}^{(j)}$ with $\mathbf{A}(\mathbf{p}^{(j)})^T \mathbf{y}^{(j)} \leq F(\mathbf{x}^{(j)}, \mathbf{p}^{(j)})$, such that the following sum of residuals is minimized:

$$r_{\mathrm{eq}} := \sum_{j=1}^{N} r_{\mathrm{PD}}(\mathbf{x}^{(j)}, \mathbf{y}^{(j)}, \mathbf{p}^{(j)}) \tag{3.4}$$

Objective (b) consists in minimizing, with $\phi$ a non-negative convex function in $(\mathbf{x}, \mathbf{y})$ such that $\phi(\mathbf{x}, \mathbf{y}) = 0 \Leftrightarrow \mathbf{x} = \mathbf{y}$:

$$r_{\text{obs}} := \sum_{j=1}^{N} \phi\left(g(\mathbf{x}^{(j)}, \mathbf{p}^{(j)}), \, \mathbf{z}^{(j)}\right) \tag{3.5}$$

As discussed in [17, 7], the parametric map $F(\cdot, \mathbf{p})$ must be searched in a restricted space $\mathcal{F}$. Since the construction of $\mathcal{F}$ is not the focus of the present article, further details will be presented in Section 8.1.

# Chapter 4

# Previous methods

## 4.1 Inverse Variational Inequality

**Formulation:** Bertsimas et al. [7] impute $F(\cdot, \mathbf{p})$ given (perfect) observations $\mathbf{z}^{(j)} = \mathbf{x}^{(j)}$ (*i.e.* $g(\cdot, \mathbf{p}) =$ Id) of approximate solution to $\text{VI}(\mathcal{K}(\mathbf{p}^{(j)}), F(\cdot, \mathbf{p}^{(j)}))$ by setting objective $r_{\text{obs}}$ in (3.5) to zero and solving:

$$
\begin{aligned}
\min_{F, \mathbf{y}} \quad & r_{\text{eq}} = \sum_{j=1}^{N} r_{\text{PD}}(\mathbf{z}^{(j)}, \mathbf{y}^{(j)}, \mathbf{p}^{(j)}) \\
\text{s.t.} \quad & \mathbf{A}(\mathbf{p}^{(j)})^T \mathbf{y}^{(j)} \leq F(\mathbf{z}^{(j)}, \mathbf{p}^{(j)}), \quad \forall j
\end{aligned}
\tag{4.1}
$$

If $r_{\text{PD}}(\mathbf{x}, \mathbf{y}, \mathbf{p}) = F(\mathbf{x}, \mathbf{p})^T \mathbf{x} - \mathbf{b}(\mathbf{p})^T \mathbf{y}$ and $F(\cdot, \mathbf{p})$ is restricted to a finite dimensional affine parametrization $\sum_{i=1}^{K} a_i F_i(\cdot, \mathbf{p})$ with parameters $\mathbf{a} \in \mathbb{R}^K$ restricted to a convex set, then (4.1) is a convex program.

   **Limitations:** The above formulation, which we will refer to as *Inverse VI*, assumes that we have complete observations, which is not possible in many applications, such as in traffic assignment, see Section 9. In addition, (4.1) overlooks the measurement errors by tightly fitting an equilibrium model to the (complete) observations, thus attempting to explain random (irreducible) errors by a deterministic process. For example, consider the following process:

$$
\min_{\mathbf{x} \geq 0} (x - a)^2
\tag{4.2}
$$

where $a > 0$ needs to be imputed. The associated primal-dual system is:

$$
x(x - a) = 0, \; x \geq a, \; x \geq 0
\tag{4.3}
$$

Given $N$ observations $z^{(j)} \geq 0$, solving (4.1) applied to our particular case:

$$
\min_{\hat{a} \geq 0} \sum_{j=1}^{N} z^{(j)}(z^{(j)} - \hat{a}) \quad \text{s.t.} \quad \hat{a} \leq \min_{j} z^{(j)}
\tag{4.4}
$$

gives an imputed parameter $\hat{a} = \min_j z^{(j)}$. Independently of the data size, a single measurement error of $\delta$ in a set of perfect observations can induce a large mean residual error. In the above example, if $z^{(1)} = a - \delta$ and $z^{(j)} = a$ for $j = 2, \cdots, N$, then the imputed value is $\hat{a} = a - \delta$, with mean residual error:

$$\frac{1}{N} \sum_{j=1}^{N} z^{(j)}(z^{(j)} - \hat{a}) = \frac{(N-1)\, a\, \delta}{N} \longrightarrow a\delta \quad \text{as} \quad N \longrightarrow +\infty \tag{4.5}$$

## Inverse programming as a bilevel program

**Formulation:** An intuitive approach is via bilevel optimization in which the metric $r_{\text{obs}} = \sum_j \phi(\mathbf{x}^{(j)}, \mathbf{z}^{(j)})$ in (3.5) is minimized with $\mathbf{x}^{(j)}$ the decision vector predicted by the imputed process. We refer to, *e.g.*, [11] for the problem of OD matrix estimation given link cost functions and observed flows. Applying bilevel optimization to our function estimation problem:

$$\min_{F, \mathbf{x}, \mathbf{y}} \quad r_{\text{obs}} = \sum_{j=1}^{N} \phi\left(g(\mathbf{x}^{(j)}, p^{(j)}), \mathbf{z}^{(j)}\right)$$
$$\mathbf{x}^{(j)} \text{ is solution to } \text{VI}(\mathcal{K}(\mathbf{p}^{(j)}), F(\mathbf{p}^{(j)})), \quad \forall j \tag{4.6}$$

With a good choice of $\phi$, (4.6) can be robust to noise. For example, consider $N$ observations $z^{(j)}$ of the minimization process (4.2). Then, (4.6) becomes:

$$\min_{\hat{a} \geq 0,\, x} \sum_{j=1}^{N} \phi(x^{(j)},\, z^{(j)}) \quad \text{s.t.} \quad x^{(j)} \in \operatorname*{argmin}_{u \geq 0} (u - \hat{a})^2, \quad \forall j \tag{4.7}$$

We note that $\hat{a}$ is the sample mean when $\phi(x) = x^2$, while $\hat{a}$ is the sample median when $\phi(x) = |x|$. Hence, formulation (4.6) allows different choices of penalty functions $\phi$ on the observation residuals, thus a fitting more robust to noise. We will refer to (4.6) as the *Bilevel Program* (BP) in the context of inverse programming.

    **Limitations:** In general, the solution set of $\text{VI}(\mathcal{K}(\mathbf{p}^{(j)}), F(\mathbf{p}^{(j)}))$ does not have a closed-form expression, thus one approach replaces the constraint in (4.6) by the primal-dual system (2.3) or KKT system (2.4) to reduce (4.6) to a single-level program. However, the complementary condition $r_{\text{PD}}(\mathbf{z}^{(j)}, \mathbf{y}^{(j)}, \mathbf{p}^{(j)}) = 0$ in the constraints causes the standard Mangasarian-Fromovitz Constraint Qualification (MFCQ) to be violated at any feasible point [26], hence generating severe numerical difficulties, see [16, 18].

# Chapter 5

# Our method

## 5.1 A Weighted Sum Program

We minimize simultaneously objectives (3.4) and (3.5) subject to primal and dual feasibilities by considering the linear combination $w_{\text{eq}}r_{\text{eq}} + w_{\text{obs}}r_{\text{obs}}$:

$$
\begin{aligned}
\min_{F,\mathbf{x},\mathbf{y}} \quad & w_{\text{eq}}\,\Sigma_{j=1}^{N}r_{\text{PD}}(\mathbf{x}^{(j)},\mathbf{y}^{(j)},\mathbf{p}^{(j)}) + w_{\text{obs}}\,\Sigma_{j=1}^{N}\phi\left(g(\mathbf{x}^{(j)},p^{(j)}),\mathbf{z}^{(j)}\right) \\
\text{s.t.} \quad & \mathbf{x}^{(j)} \in \mathcal{K}(\mathbf{p}^{(j)}), \quad \forall\, j \\
& \mathbf{A}(\mathbf{p}^{(j)})^{T}\mathbf{y}^{(j)} \le F(\mathbf{x}^{(j)},\mathbf{p}^{(j)}), \quad \forall\, j
\end{aligned}
\tag{5.1}
$$

where $w_{\text{eq}}$ and $w_{\text{obs}}$ are positive scalars that articulate the preferences between the two objectives. This approach is known as the *Weighted Sum method* in Pareto Optimization (PO) theory, and is sufficient for Pareto optimality, *i.e.*, it is not possible to strictly decrease one objective among $r_{\text{eq}}$ and $r_{\text{obs}}$ without strictly increasing the other one, see, *e.g.*, [19, 13] for further details on PO.

One approach to explore the Pareto curve is shown in Algorithm 1. In step 1, we note that it is often desirable to scale the objective functions to have a consistent comparison between them. Varying the weights provides information about available trade-offs between the objectives. Specifically, for each of the different weights in step 2, if the solution is such that $r_{\text{eq}}$ and $r_{\text{obs}}$ are large, it means that either our model is not a good model to explain the observations, or that our observations are very noisy.

---

**Algorithm 1** `Weighted-sum(·)` Weighted sum method

---

1: Normalize objectives (3.4) and (3.5) for consistent comparisons.
2: Solve (5.1) with $w_{\text{obs}} + w_{\text{eq}} = 1$ and $w_{\text{obs}} \in \{0.001, 0.01, 0.1, 0.9, 0.99, 0.999\}$
3: Check the values of (3.4) and (3.5).

---

The proposed weighted Sum Program (WSP) is robust since it accommodates different penalty functions $\phi$ depending on the type of measurement errors, *e.g.* $\phi(\mathbf{x},\mathbf{z}) = \|\mathbf{x}-\mathbf{z}\|_1$ for robustness to outliers, and $\phi(\mathbf{x},\mathbf{z}) = \|\mathbf{x}-\mathbf{z}\|_2$ for robustness to Gaussian noise; see, *e.g.*, [8,

§6.1]. In addition, our WSP can be seen as a penalty method for constrained optimization that mitigates numerical difficulties by minimizing $r_{\mathrm{PD}}(\mathbf{z}, \mathbf{y}, \mathbf{p})$ instead of setting $r_{\mathrm{PD}}(\mathbf{z}, \mathbf{y}, \mathbf{p})$ to 0.

## 5.2 Example

Given $N$ observations $z^{(j)}$ of $\min\limits_{x \geq 0}(x - a)^2$, the WSP (5.1) is:

$$
\begin{aligned}
\min_{\hat{a}, x} \quad & w_{\mathrm{eq}} \textstyle\sum_{j=1}^{N} x^{(j)}(x^{(j)} - \hat{a}) + w_{\mathrm{obs}} \sum_{j=1}^{N} \phi(x^{(j)}, z^{(j)}) \\
\text{s.t.} \quad & x^{(j)} \geq 0, \quad \forall j \\
& 0 \leq \hat{a} \leq \min_{j} x^{(j)}
\end{aligned}
\tag{5.2}
$$

We now set $w_{\mathrm{obs}} = \alpha$, $w_{\mathrm{eq}} = 1 - \alpha$ for $\alpha \in (0, 1)$ and $\phi(x, y) = |x - y|$. Following the case study in Section 4.1, assume the observations are $z^{(1)} = a - \delta$ and $z^{(j)} = a$ for $j = 2, \cdots, N$, then the set of Pareto optimal points are

$$
\hat{a} = x^{(1)} \in [a - \delta, a], \quad x^{(j)} = a, \quad \text{for } j = 2, \cdots, N
\tag{5.3}
$$

Then, given estimate $\hat{a} \in [a - \delta, a]$, objectives $r_{\mathrm{eq}}$ in (3.4) and $r_{\mathrm{obs}}$ in (3.5) are:

$$
r_{\mathrm{eq}} = (N - 1)\, a\, (a - \hat{a})
\tag{5.4}
$$

$$
r_{\mathrm{obs}} = |a - \delta - \hat{a}| = \hat{a} + \delta - a
\tag{5.5}
$$

Solving $\min\limits_{\hat{a} \in [a-\delta,\, a]} w_{\mathrm{eq}} r_{\mathrm{eq}} + w_{\mathrm{obs}} r_{\mathrm{obs}} = (1 - \alpha)(N - 1)a(a - \hat{a}) + \alpha(\hat{a} + \delta - a)$:

$$
\begin{aligned}
w_{\mathrm{obs}} = \alpha < \tfrac{a(N-1)}{1+a(N-1)} \implies & \hat{a} = a, & r_{\mathrm{eq}} = 0, & & r_{\mathrm{obs}} = \delta \\
w_{\mathrm{obs}} = \alpha > \tfrac{a(N-1)}{1+a(N-1)} \implies & \hat{a} = a - \delta, & r_{\mathrm{eq}} = (N-1)a\delta, & & r_{\mathrm{obs}} = 0
\end{aligned}
\tag{5.6}
$$

In this case, if $w_{\mathrm{obs}}$ is close enough to 1, $r_{\mathrm{eq}}$ is large and equal to the one in the Inverse VI (see (4.5)), while with $w_{\mathrm{obs}}$ smaller, we have a small observation residual $r_{\mathrm{obs}}$ and $r_{\mathrm{eq}} = 0$. Thus, the estimation is good for $w_{\mathrm{obs}}$ close enough to 0 despite a fit to the data that is not perfect due to measurement errors.

In a second example, we randomly generate $N = 20$ independent and identically distributed (i.i.d.) samples $z^{(j)}$ from a Gaussian distribution with mean $a = 10$ and variance $\sigma = 5$. We apply our WSP (5.2) with $\phi(x, y) = (x - y)^2$. The estimates $\hat{a}$ are shown in Figure 5.1.a), and the values of the residuals $r_{\mathrm{obs}} = \sum_{j=1}^{N} x^{(j)}(x^{(j)} - \hat{a})$ and $r_{\mathrm{eq}} = \sum_{j=1}^{N}(x^{(j)} - z^{(j)})^2$ in Figure 5.1.b), for $w_{\mathrm{obs}} \in \{0.001, 0.01, 0.1, 0.9, 0.99, 0.999\}$ and $w_{\mathrm{eq}} = 1 - w_{\mathrm{obs}}$. In addition, we compare our method to the Inverse VI (4.4), which is tagged with label $w_{\mathrm{obs}} = 1$ in Figure 5.1. For $w_{\mathrm{obs}} < 0.1$, $r_{\mathrm{eq}} = 0$, $\hat{a} = 9.2$ is close to 10, and $r_{\mathrm{obs}}$ is small, while for $w_{\mathrm{obs}} > 0.99$ and for the Inverse VI, $r_{\mathrm{eq}}$ is large and $\hat{a}$ is largely under-estimating $a = 10$. In the presence of Gaussian noise, our WSP also performs well. Finally, we note that for large values of $w_{\mathrm{obs}}$, our method behaves similarly to the Inverse VI method.

Figure 5.1: **Imputation of the parametric program** $\min\limits_{x \geq 0}(x-a)^2$ **from** $N = 20$ **noisy observations with mean** $10$ **shown in Figure a). The estimates are shown by horizontal lines labelled by the value of** $w_{\mathrm{obs}}$ **used in the WSP** $(5.2)$**, at the exception of** $w_{\mathrm{obs}} = 1$ **for which the estimate is obtained via the Inverse VI** $(4.4)$**. The associated residuals** $r_{\mathrm{obs}}$ **and** $r_{\mathrm{eq}}$ **are shown in Figure b).**

# Chapter 6

# Relation to previous methods

## 6.1 Preliminary results

Intuitively, as $(w_{\mathrm{eq}}, w_{\mathrm{obs}})$ approaches $(0, 1)$, the WSP (5.1) mimics the Inverse VI (4.1), and as $(w_{\mathrm{eq}}, w_{\mathrm{obs}})$ approaches $(1, 0)$, it mimics the BP (4.6). Formally, given $f_1$, $f_2$ two non-negative continuous functions, a compact set $\mathcal{C} \subseteq \mathbb{R}^n$, and $w_1, w_2 > 0$, consider the general weighted sum program along with its solution set $\mathcal{S}(w_1, w_2)$ and the set $\mathcal{S}$ of all Pareto efficient points associated to it:

$$\min w_1 f_1(\mathbf{u}) + w_2 f_2(\mathbf{u}) \quad \text{s.t.} \quad \mathbf{u} \in \mathcal{C} \tag{6.1}$$

$$\mathcal{S}(w_1, w_2) := \arg\min_{\mathbf{u} \in \mathcal{C}} w_1 f_1(\mathbf{u}) + w_2 f_2(\mathbf{u}) \tag{6.2}$$

$$\mathcal{S} := \left\{ (w_1, w_2, \mathbf{u}^\star) : w_1 \in (0, 1), w_2 = 1 - w_1, \mathbf{u}^\star \in \mathcal{S}(w_1, w_2) \right\} \tag{6.3}$$

Since $\mathcal{C}$ is compact, $\mathcal{S}(w_1, w_2) \neq \emptyset$ for any $w_1, w_2$, hence $\mathcal{S}$ is well-defined. We also assume there exists $\mathbf{u} \in \mathcal{C}$ such that $f_1(\mathbf{u}) = 0$, and define the following constrained program and its approximate objective value $f_2^\star(\epsilon)$:

$$\min f_2(\mathbf{u}) \quad \text{s.t.} \quad f_1(\mathbf{u}) = 0, \mathbf{u} \in \mathcal{C} \tag{6.4}$$

$$f_2^\star(\epsilon) := \min_{\mathbf{u} \in \mathcal{C} : f_1(\mathbf{u}) \leq \epsilon} f_2(\mathbf{u}), \quad \forall \epsilon \geq 0 \tag{6.5}$$

**Lemma 6.1.** *Let $\mathcal{S}$ be a set described by (6.3). Then for any $(w_1, w_2, \mathbf{u}^\star) \in \mathcal{S}$:*

$$f_1(\mathbf{u}^\star) \leq (w_1^{-1} - 1) f_2^\star(0) \tag{6.6}$$

$$f_2(\mathbf{u}^\star) \leq f_2^\star(0) \tag{6.7}$$

**Proof.** Let $\mathbf{u} \in \mathcal{C}$ such that $f_1(\mathbf{u}) = 0$. For any $(w_1, w_2, \mathbf{u}^\star) \in \mathcal{S}$, we have $w_1 f_1(\mathbf{u}^\star) + w_2 f_2(\mathbf{u}^\star) \leq w_2 f_2(\mathbf{u})$, hence, from non-negativity of $f_1$, $f_2$ and positivity of $w_1$, $w_2$:

$$f_2(\mathbf{u}^\star) \leq f_2(\mathbf{u}) \tag{6.8}$$

$$f_1(\mathbf{u}^\star) \leq (w_2/w_1) f_2(\mathbf{u}) = ((1 - w_1)/w_1) f_2(\mathbf{u}) = (w_1^{-1} - 1) f_2(\mathbf{u}) \tag{6.9}$$

Since this is true for all $\mathbf{u} \in \mathcal{C}$ such that $f_1(\mathbf{u}) = 0$, minimizing $f_2$ for such $\mathbf{u}$ completes the proof. $\qquad\square$

**Lemma 6.2.** *(See [5]) Let $\mathcal{S}$ be a set described by (6.3). Then $\{f_2(\mathbf{u}^\star)\}_{\mathbf{u}^\star \in \mathcal{S}(w_1, w_2)}$ converges uniformly to $f_2^\star(0)$ as $w_1 \longrightarrow 1$. There also exists a solution $\bar{\mathbf{u}}$ to (6.4) and a sequence $(w_1^{(n)}, w_2^{(n)}, \mathbf{u}_n)_{n \in \mathbb{N}} \in \mathcal{S}^{\mathbb{N}}$ such that:*

$$(w_1^{(n)}, w_2^{(n)}, \mathbf{u}_n) \longrightarrow (1, 0, \bar{\mathbf{u}}) \quad as\ n \longrightarrow +\infty \tag{6.10}$$

*In addition, if (6.4) admits a unique solution $\bar{\mathbf{u}}$, any sequence $(w_1^{(n)}, w_2^{(n)}, \mathbf{u}_n)_{n \in \mathbb{N}} \in \mathcal{S}^{\mathbb{N}}$ such that $w_1^{(n)} \longrightarrow 1$ satisfies $\mathbf{u}^{(n)} \longrightarrow \bar{u}$.*

**Proof.** First, we want to prove that $f_2^\star(\cdot)$ is continuous at 0. We note that $f_2^\star(\cdot)$ is non-increasing on $\mathbb{R}_+$, hence it has a limit from the right at 0, which we denote $f_2^\star(0_+)$. Given any sequence $(\epsilon_n)_{n \in \mathbb{N}} \in \mathbb{R}_+^{\mathbb{N}}$ such that $\epsilon_n \to 0$, there exists a sequence $(\mathbf{u}_n)_{n \in \mathbb{N}}$ such that $\mathbf{u}_n \in \underset{\mathbf{u} \in \mathcal{C} : f_1(\mathbf{u}) \leq \epsilon_n}{\arg \min} \ f_2(\mathbf{u})$ for all $n \in \mathbb{N}$, since $\mathcal{C}$ is compact. Hence, $f_2^\star(\epsilon_n) = f_2(\mathbf{u}_n)$ for all $n \in \mathbb{N}$. From compactness, there exists a convergent subsequence $(\tilde{\epsilon}_n, \tilde{\mathbf{u}}_n)_{n \in \mathbb{N}}$ of $(\epsilon_n, \mathbf{u}_n)_{n \in \mathbb{N}}$, and its limit $(0, \bar{\mathbf{u}})$ is such that $\mathbf{u} \in \mathcal{C}$, $f_1(\bar{\mathbf{u}}) = 0$ and $f_2^\star(0_+) = f_2(\bar{\mathbf{u}}) \leq f_2^\star(0)$ from continuity of $f_1$ and $f_2$. By definition of $f_2^\star(0)$, we must have $f_2^\star(0_+) = f_2^\star(0)$. Hence $f_2^\star(\cdot)$ is continuous at 0.

To prove the first part of the lemma, we denote $g(w_1) := (w_1^{-1} - 1) f_2^\star(0)$. For any $(w_1, w_2, \mathbf{u}^\star) \in \mathcal{S}$, we have $f_1(\mathbf{u}^\star) \leq g(w_1)$ from lemma 6.1, hence $f_2^\star(g(w_1)) \leq f_2(\mathbf{u}^\star) \leq f_2^\star(0)$ by definition of $f_2^\star(\epsilon)$. Thus, by continuity of $f_2^\star(\cdot)$ at 0: $\forall \mathbf{u}^\star \in \mathcal{S}(w_1, w_2), \quad |f_2(\mathbf{u}^\star) - f_2^\star(0)| \leq |f_2^\star(g(w_1)) - f_2^\star(0)| \underset{w_1 \to 1}{\longrightarrow} 0$.

We prove the second part of the lemma. Given a sequence $(w_1^{(n)}, w_2^{(n)}, \mathbf{u}_n) \in \mathcal{S}^{\mathbb{N}}$ such that $w_1^{(n)} \longrightarrow 1$, consider a convergent subsequence of it from compactness of $\mathcal{C}$. Its limit $(1, 0, \bar{\mathbf{u}})$ is such that $\bar{\mathbf{u}} \in \mathcal{C}$, $f_1(\bar{\mathbf{u}}) = 0$, and $f_2(\bar{\mathbf{u}}) = f_2^\star(0)$ from continuity of $f_1$ and $f_2$. Hence $\bar{u}$ is a solution to (6.4), which gives the second result of the lemma.

For the third part of the lemma, we start from the proof of the second part and note that any convergent subsequence $(\tilde{w}_1^{(n)}, \tilde{w}_2^{(n)}, \tilde{\mathbf{u}}_n)$ of $(w_1^{(n)}, w_2^{(n)}, \mathbf{u}_n)$ is such that $\tilde{\mathbf{u}}_n$ converges to the unique solution $\bar{\mathbf{u}}$ to (6.4). Hence any convergent subsequence has the same limit $(1, 0, \bar{u})$, and $(w_1^{(n)}, w_2^{(n)}, \mathbf{u}_n)$ thus converges to $(1, 0, \bar{u})$. Since this is true for any sequence $(w_1^{(n)}, w_2^{(n)}, \mathbf{u}_n) \in \mathcal{S}^{\mathbb{N}}$ such that $w_1^{(n)} \longrightarrow 1$, we have the third result of the lemma. $\qquad\square$

## 6.2   Main results

To apply the results in Section 6.1 to our WSP, we substitute $\mathbf{u}$ with the tuple $(F(\cdot, \mathbf{p}), \{\mathbf{x}^{(j)}\}_j, \{\mathbf{y}^{(j)}\}_j)$ and the objectives $(f_1, f_2)$ with $(r_{\text{obs}}, r_{\text{eq}})$. Since the feasible set in (5.1) is closed, compactness is guaranteed with this assumption:

**Assumption 6.1.** *The variables $(F(\cdot, \mathbf{p}), \{\mathbf{x}^{(j)}\}_j, \{\mathbf{y}^{(j)}\}_j)$ of the WSP (5.1) are in a finite-dimensional bounded set.*

The finite dimension assumption is reasonable since restricting the map $F(\cdot, \mathbf{p})$ to a finite dimensional affine parametrization $\sum_{i=1}^{K} a_i F_i(\cdot, \mathbf{p})$ is intuitive (as in [17, 7]). The boundedness is essential since the primal variables $\mathbf{x}^{(j)}$, dual variables $\mathbf{y}^{(j)}$, and the parameters $\mathbf{a}$ have physical interpretations in terms of resource allocation, resource valuation, and variations of the map $F(\cdot, \mathbf{p})$ respectively, and thus are restricted to physically (or economically) reasonable ranges. Hence Assumption 6.1 is reasonable and guarantees compactness of the set of feasible variables of the WSP, and enables to apply the results in Section 6.1. From compactness, the minimal objective value $r_{\mathrm{eq}}^{\star}$ of the Inverse VI (4.1), and the minimal objective value $r_{\mathrm{obs}}^{\star}$ of the BP (4.6) are also attained.

**Theorem 6.1.** *Under Assumption 6.1, given $N$ approximate solutions $\mathbf{z}^{(j)} \in \mathcal{K}(\mathbf{p}^{(j)})$ to the problems $VI(\mathcal{K}(\mathbf{p}^{(j)}), F(\cdot, \mathbf{p}^{(j)}))$ for $j = 1, \cdots, N$, any optimal solution to the WSP (5.1) is such that $r_{obs} \leq r_{eq}^{\star}(w_{obs}^{-1} - 1)$ and $r_{eq} \leq r_{eq}^{\star}$. In addition, $r_{eq}$ converges uniformly to $r_{eq}^{\star}$ as $w_{obs} \longrightarrow 1$ and there exists a sequence of solutions to the WSP converging to a solution to the inverse VI (4.1).*

**Theorem 6.2.** *Under Assumption 6.1, given $N$ observations $\mathbf{z}^{(j)}$ in (3.3), any optimal solution to the WSP (5.1) is such that $r_{eq} \leq r_{obs}^{\star}(w_{eq}^{-1} - 1)$, $r_{obs} \leq r_{obs}^{\star}$, and In addition, $r_{obs}$ converges uniformly to $r_{obs}^{\star}$ as $w_{eq} \longrightarrow 1$ and there exsits a sequence of solutions to the WSP converging to a solution to the BP (4.6).*

Finally, the objective $r_{\mathrm{obs}}$ in our WSP (5.1) can be generalized, thus our WSP can be seen as a smoothing method for general bilevel programs where the complementary condition $r_{\mathrm{PD}} = 0$ is included in the objective in the form of a penalty function. Previous works have proposed smoothing methods via, *e.g.*, the perturbed Fischer-Burmeister function [12, §6.5] or a similar one [14], but our smoothing via residuals has a sub-optimality interpretation.

# Chapter 7

# Comparison of the duality gap and the KKT residual

Given $N$ observations $\mathbf{z}^{(j)}$, for $j = 1, \cdots, N$, let $(F(\cdot, \mathbf{p}), \{\mathbf{x}^{(j)}\}_j, \{\mathbf{y}^{(j)}\}_j)$ be an optimal solution to the WSP (5.1). Then $r_{\text{obs}}$ in (3.5) measures how well $\mathbf{x}^{(j)}$ agree with the observations $\mathbf{z}^{(j)}$, while $r_{\text{eq}}$ in (3.4) measures how well the imputed process $\text{VI}(\mathcal{K}(\mathbf{p}), F(\cdot, \mathbf{p}))$ explains the imputed decision vectors $\mathbf{x}^{(j)}$. If the imputed map $F(\cdot, \mathbf{p})$ admits a unique solution $\hat{\mathbf{x}}(\mathbf{p})$ for all $\mathbf{p}$ (*e.g.*, from strict monotonicity), then an alternative metric to $r_{\text{eq}}$ is $\sum_{j=1}^{N} \|\mathbf{x}^{(j)} - \hat{\mathbf{x}}(\mathbf{p}^{(j)})\|$. If $F(\cdot, \mathbf{p})$ is strongly monotone with parameter $c$ for all $\mathbf{p}$, from (2.13) and (2.20):

$$\|\mathbf{x}^{(j)} - \hat{\mathbf{x}}(\mathbf{p}^{(j)})\|_2 \leq \sqrt{r_{\text{PD}}(\mathbf{x}^{(j)}, \mathbf{y}^{(j)}, \mathbf{p}^{(j)})/c} \quad \forall j \tag{7.1}$$

where $r_{\text{PD}}(\mathbf{x}^{(j)}, \mathbf{y}^{(j)}, \mathbf{p}^{(j)})$, $j = 1, \cdots, N$ are directly available from the WSP. Note that with only strict convexity of $F(\cdot, \mathbf{p})$, we can have $\|\mathbf{x}^{(j)} - \hat{\mathbf{x}}(\mathbf{p}^{(j)})\|_2 = \delta$ while $\sqrt{r_{\text{PD}}(\mathbf{x}^{(j)}, \mathbf{y}^{(j)}, \mathbf{p}^{(j)})}$ is infinitely small, as shown at the end of Section 2.3.

However, there is no result of the form $\|\mathbf{x} - \mathbf{x}^\star\| = \mathcal{O}(\sqrt{r_{\text{KKT}}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi})})$ to the best of our knowledge. We define the slack variables associated to the dual feasibility condition $\mathbf{A}^T \mathbf{y} \leq F(\mathbf{x})$:

$$\boldsymbol{\nu} := F(\mathbf{x}) - \mathbf{A}^T \mathbf{y} \tag{7.2}$$

which implies that dual feasibility is equivalent to $\boldsymbol{\nu} \geq 0$. We now derive a bound for the following generalized residuals:

$$r_{\text{PD}}^{\ell_p}(\mathbf{x}) = \|\boldsymbol{\nu} \circ \mathbf{x}\|_p = \left( \sum_{i=1}^{n} |\nu_i x_i|^p \right)^{1/p} \tag{7.3}$$

$$r_{\text{KKT}}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) = \left\| \begin{bmatrix} \alpha(\boldsymbol{\nu} - \boldsymbol{\pi}) \\ \mathbf{x} \circ \boldsymbol{\pi} \end{bmatrix} \right\|_p = \left( \sum_{i=1}^{n} \alpha^p |\nu_i - \pi_i|^p + |x_i \pi_i|^p \right)^{1/p} \tag{7.4}$$

where $\|x\|_p$ is the p-norm for $p \geq 1$, and $\mathbf{u} \circ \mathbf{v} = [u_i v_i]_{i=1}^n$ for all $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$. Since $\|\mathbf{x}\|_p \leq \|\mathbf{x}\|_1 \leq n^{1-1/p}\|\mathbf{x}\|_p$ for all $\mathbf{x} \in \mathbb{R}^n$, we have

$$r_{\mathrm{PD}}^{\ell_p}(\mathbf{x}, \mathbf{y}) \leq r_{\mathrm{PD}}^{\ell_1}(\mathbf{x}, \mathbf{y}) \leq n^{1-1/p} \cdot r_{\mathrm{PD}}^{\ell_p}(\mathbf{x}, \mathbf{y}) \tag{7.5}$$

$$r_{\mathrm{KKT}}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \leq r_{\mathrm{KKT}}^{\ell_1}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \leq n^{1-1/p} \cdot r_{\mathrm{KKT}}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \tag{7.6}$$

When primal and dual feasibilities hold, *i.e.* $\boldsymbol{\nu} \geq 0$, $\mathbf{Ax} = \mathbf{b}$, $\mathbf{x} \geq 0$, we note that $r_{\mathrm{PD}}^{\ell_1}$, $r_{\mathrm{KKT}}^{\ell_1}$ defined above correspond to $r_{\mathrm{PD}}$, $r_{\mathrm{KKT}}$ in (2.9), (2.8) since, for $r_{\mathrm{PD}}^{\ell_1}$:

$$r_{\mathrm{PD}}^{\ell_1}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n \nu_i x_i = \boldsymbol{\nu}^T \mathbf{x} = (F(\mathbf{x}) - \mathbf{A}^T \mathbf{y})^T \mathbf{x} = F(\mathbf{x})^T \mathbf{x} - \mathbf{b}^T \mathbf{y} \tag{7.7}$$

The results in Section 2 thus hold for $r_{\mathrm{PD}}^{\ell_p}$ and $r_{\mathrm{KKT}}^{\ell_p}$ with an additional $n^{1-1/p}$ factor, validating them as residuals for the primal-dual and KKT systems respectively. Before stating our main result of the section, we present a lemma:

**Lemma 7.1.** *Let $\mathcal{K}$ be a polyhedron given by (2.2). Then the following holds for any $\alpha > 0$, $p > 1$, $\mathbf{x} \in \mathcal{K}$, $\mathbf{y} \in \mathbb{R}^n$ such that $\mathbf{A}^T \mathbf{y} \leq F(\mathbf{x})$:*

$$\min_{\boldsymbol{\pi} \geq 0} r_{KKT}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) = \left( \sum_{i=1}^n \frac{(\nu_i x_i)^p}{\left(1 + (x_i/\alpha)^{\frac{p}{p-1}}\right)^{p-1}} \right)^{1/p} \tag{7.8}$$

*If $p = 1$, then for any $\alpha > 0$, $\mathbf{x} \in \mathcal{K}$, $\mathbf{y} \in \mathbb{R}^n$ such that $\mathbf{A}^T \mathbf{y} \leq F(\mathbf{x})$, we have:*

$$\min_{\boldsymbol{\pi} \geq 0} r_{KKT}^{\ell_1}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) = \sum_{i:x_i < \alpha} x_i \nu_i + \sum_{i:x_i > \alpha} \alpha \nu_i \tag{7.9}$$

**Proof.** For any $p \geq 1$, $\mathbf{x} \in \mathbb{R}^n$, and $\mathbf{y} \in \mathbb{R}^n$:

$$\min_{\boldsymbol{\pi} \geq 0} \left( r_{KKT}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \right)^p = \min_{\boldsymbol{\pi} \geq 0} \sum_{i=1}^n \alpha^p |\nu_i - \pi_i|^p + |x_i \pi_i|^p \tag{7.10}$$

$$= \sum_{i=1}^n \min_{\pi_i \geq 0} \{\alpha^p |\nu_i - \pi_i|^p + |x_i \pi_i|^p\} \tag{7.11}$$

When primal and dual feasibilities hold, $\mathbf{x} \geq 0$, $\boldsymbol{\nu} \geq 0$, which causes the map $\pi_i \geq 0 \mapsto \alpha^p |\nu_i - \pi_i|^p + |x_i \pi_i|^p$ to increase on $[\nu_i, +\infty)$ and thus to attain its minimum on $[0, \nu_i]$, on which it is also differentiable, for all $p \geq 1$, with gradient:

$$\pi_i \mapsto -p\alpha^p (\nu_i - \pi_i)^{p-1} + p x_i^p \pi_i^{p-1}, \quad i = 1, \cdots, n \tag{7.12}$$

When $p > 1$, the gradient vanishes at a unique point $\pi_i^\star$ in $[0, \nu_i]$:

$$\pi_i^\star = \frac{\nu_i}{1 + (x_i/\alpha)^{p/(p-1)}}, \quad i = 1, \cdots, n \tag{7.13}$$

Substituting in (7.11):

$$\min_{\pi_i \geq 0} \{\alpha^p |\nu_i - \pi_i|^p + |x_i \pi_i|^p\} = \frac{(\nu_i x_i)^p}{(1 + (x_i/\alpha)^{p/(p-1)})^{p-1}} \tag{7.14}$$

which gives the desired result for $p > 1$.

When $p = 1$, the map $\pi_i \geq 0 \mapsto \alpha|\nu_i - \pi_i| + |x_i \pi_i|$ is just affine on $[0, \nu_i]$, in which the minimum is, and thus attains it minimum at 0 if $x_i - \alpha \geq 0$, and $\nu_i$ if $x_i - \alpha < 0$. Hence:

$$\sum_i \min_{\pi_i \geq 0} \{\alpha|\nu_i - \pi_i| + |x_i \pi_i|\} = \sum_{i\,:\,x_i < \alpha} x_i \nu_i + \sum_{i\,:\,x_i > \alpha} \alpha \nu_i \tag{7.15}$$

which completes the proof. □

We are now present the main result of the section, where $\|\mathbf{x}\|_\infty = \max_i |x_i|$:

**Theorem 7.1.** *Let $\mathcal{K}$ be a polyhedron given by (2.2). Then the following holds for any $\alpha > 0$, $p \geq 1$, $\epsilon > 0$, $\mathbf{x} \in \mathcal{K}$, $\mathbf{y} \in \mathbb{R}^n$ such that $\mathbf{A}^T \mathbf{y} \leq F(\mathbf{x})$:*

$$r_{PD}^{\ell_p}(\mathbf{x}, \mathbf{y}) \leq \epsilon \quad \Longrightarrow \quad \exists\, \boldsymbol{\pi} \in \mathbb{R}^n : r_{KKT}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \leq \epsilon \tag{7.16}$$

*Reciprocally, for $p > 1$, we have; for all $\epsilon > 0$, $\mathbf{x} \in \mathcal{K}$, $\mathbf{y} \in \mathbb{R}^n : \mathbf{A}^T \mathbf{y} \leq F(\mathbf{x})$:*

$$\exists\, \boldsymbol{\pi} \in \mathbb{R}_+^n, r_{KKT}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \leq \epsilon \implies r_{PD}^{\ell_p}(\mathbf{x}, \mathbf{y}) \leq \epsilon \left(1 + (\|\mathbf{x}\|_\infty/\alpha)^{\frac{p}{p-1}}\right)^{\frac{p-1}{p}} \tag{7.17}$$

*When $p = 1$, we have; for all $\epsilon > 0$, $\mathbf{x} \in \mathcal{K}$, $\mathbf{y} \in \mathbb{R}^n$, and $\mathbf{A}^T \mathbf{y} \leq F(\mathbf{x})$:*

$$\exists\, \boldsymbol{\pi} \in \mathbb{R}_+^n, r_{KKT}^{\ell_1}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \leq \epsilon \implies r_{PD}^{\ell_1}(\mathbf{x}) \leq \epsilon \max\left(\|\mathbf{x}\|_\infty/\alpha, 1\right) \tag{7.18}$$

**Proof.** To prove (7.16) for $p > 1$, note that for all $\mathbf{x} \in \mathcal{K}$, $\mathbf{y} \in \mathbb{R}^n$ such that $\mathbf{A}^T \mathbf{y} \leq F(\mathbf{x})$, each term $(\nu_i x_i)^p / \left(1 + (x_i/\alpha)^{\frac{p}{p-1}}\right)^{p-1}$ in $\min_{\boldsymbol{\pi} \geq 0} r_{KKT}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi})$ given by (7.8) is less or equal than $|\nu_i x_i|^p$. Hence:

$$\min_{\boldsymbol{\pi} \geq 0} r_{KKT}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \leq r_{PD}^{\ell_p}(\mathbf{x}, \mathbf{y}) \tag{7.19}$$

which proves (7.16) for $p > 1$. For $p = 1$, (7.16) is true since, from $\mathbf{x} \geq 0$ and $\boldsymbol{\nu} \geq 0$:

$$\min_{\boldsymbol{\pi} \geq 0} r_{KKT}^{\ell_1}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) = \sum_{i\,:\,x_i < \alpha} x_i \nu_i + \sum_{i\,:\,x_i > \alpha} \alpha \nu_i \leq \sum_i x_i \nu_i = r_{PD}^{\ell_1}(\mathbf{x}, \mathbf{y}) \tag{7.20}$$

To prove (7.17), we note that for all $\mathbf{x} \in \mathcal{K}$, $\mathbf{y} \in \mathbb{R}^n$ such that $\mathbf{A}^T \mathbf{y} \leq F(\mathbf{x})$, each term $(\nu_i x_i)^p / \left(1 + (x_i/\alpha)^{\frac{p}{p-1}}\right)^{p-1}$ in $\min_{\boldsymbol{\pi} \geq 0} r_{KKT}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi})$ is greater or equal than $(\nu_i x_i)^p / \left(1 + (\|\mathbf{x}\|_\infty/\alpha)^{\frac{p}{p-1}}\right)^{p-1}$,

hence, for all $\boldsymbol{\pi} \in \mathbb{R}^n$ such that $\boldsymbol{\pi} \geq 0$:

$$r_{\text{KKT}}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \geq \min_{\boldsymbol{\pi} \geq 0} r_{\text{KKT}}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \tag{7.21}$$

$$\geq \left( \sum_i (\nu_i x_i)^p / \left( 1 + (\|\mathbf{x}\|_\infty/\alpha)^{\frac{p}{p-1}} \right)^{p-1} \right)^{1/p} \tag{7.22}$$

which proves (7.17) for $p > 1$. For $p = 1$, we have, with $(x)_+ = \max(x, 0)$:

$$r_{\text{PD}}^{\ell_1}(\mathbf{x}, \mathbf{y}) - \min_{\boldsymbol{\pi} \geq 0} r_{\text{KKT}}^{\ell_1}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) = \sum_{i\,:\,x_i > \alpha} (x_i - \alpha)\nu_i \tag{7.23}$$

$$\leq (\|\mathbf{x}\|_\infty - \alpha)_+ \sum_{i\,:\,x_i > \alpha} \nu_i \tag{7.24}$$

$$\leq \frac{(\|\mathbf{x}\|_\infty - \alpha)_+}{\alpha} \min_{\boldsymbol{\pi} \geq 0} r_{\text{KKT}}^{\ell_1}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \tag{7.25}$$

hence $r_{\text{PD}}^{\ell_1}(\mathbf{x}, \mathbf{y}) \leq (1 + (\|\mathbf{x}\|_\infty/\alpha - 1)_+) \min_{\boldsymbol{\pi} \geq 0} r_{\text{KKT}}^{\ell_1}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi})$. Finally, noting that $1 + (\|\mathbf{x}\|_\infty/\alpha - 1)_+ = \max(\|\mathbf{x}\|_\infty/\alpha, 1)$ completes the proof. $\qquad\square$

The first bound (7.16) in Theorem 5.1. is tight since, using Lemma 5.1, we have $\min_{\boldsymbol{\pi} \geq 0} r_{\text{KKT}}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \longrightarrow r_{\text{PD}}(\mathbf{x}, \mathbf{y})$ as $\alpha \longrightarrow +\infty$, for any $p \geq 1$. The bounds (7.17) and (7.18) are tight since we have equality in one dimension, *i.e.* $n = 1$. Combining (7.17), (7.5), (2.13) and (2.20) we have:

**Theorem 7.2.** *Let $\mathcal{K}$ be a polyhedron given by (2.2), and $F$ be a strongly monotone function with parameter $c > 0$. Then $VI(\mathcal{K}, F)$ admits a unique solution $\mathbf{x}^\star$ and; for any $\alpha > 0$, $p > 1$, $\epsilon > 0$, $\mathbf{x} \in \mathcal{K}$:*

$$\exists \mathbf{y}, \boldsymbol{\pi} \in \mathbb{R}^n \;:\; \mathbf{A}^T \mathbf{y} \leq F(\mathbf{x}),\; \boldsymbol{\pi} \geq 0,\; r_{KKT}^{\ell_p}(\mathbf{x}, \mathbf{y}, \boldsymbol{\pi}) \leq \epsilon$$
$$\implies \|\mathbf{x} - \mathbf{x}^\star\|_2 \leq \sqrt{n^{1-1/p} \cdot \epsilon \left( 1 + (\|\mathbf{x}\|_\infty/\alpha)^{\frac{p}{p-1}} \right)^{\frac{p-1}{p}} / c} \tag{7.26}$$

*For $p = 1$, the bound is $\|\mathbf{x} - \mathbf{x}^\star\|_2 \leq \sqrt{\epsilon \max (\|\mathbf{x}\|_\infty/\alpha,\, 1) / c}$.*

# Chapter 8

# Implementation

## 8.1 Affine parametrization

For tractability reasons, a classic approach consists in restricting the parametric map $F(\cdot, \mathbf{p})$ to be imputed to a finite dimensional affine parametric model

$$F(\cdot, \mathbf{p}) = F_0(\cdot, \mathbf{p}) + \sum_{i=1}^{K} a_i F_i(\cdot, \mathbf{p}), \quad \mathbf{a} \in \mathcal{A} \subseteq \mathbb{R}^K \tag{8.1}$$

where $F_i(\cdot, \mathbf{p})$, $i = 0, \cdots, K$ are pre-selected basis functions that typically contain prior knowledge on the candidate functions, and $\mathbf{a}$ is imputed in the set of allowable parameter vectors $\mathcal{A}$. For instance, if the true map $F^{\text{true}}(\cdot, \mathbf{p})$ is known to be increasing for all $\mathbf{p}$, then having a parameter space $\mathcal{A} \subseteq \mathbb{R}_+^K$ and increasing basis maps $F_i(\cdot, \mathbf{p})$ given any $(i, \mathbf{p})$ guarantees an increasing parametric map $F(\cdot, \mathbf{p})$ for all $\mathbf{a} \in \mathcal{A}$. In addition, the constant shift $F_0(\cdot, \mathbf{p})$ imposes a normalization on $F(\cdot, \mathbf{p})$ such that trivial solutions are excluded, *e.g.*, null maps where all of $\mathcal{K}$ is solution to the VI problem, and for which both non-negative objectives $r_{\text{eq}}$ (3.4) and $r_{\text{obs}}$ (3.5) can be minimized to zero.

A nonparametric estimation has also been considered in [7] using kernel methods and regularization methods from statistical learning. The methodology presented in the present article can also be extended to this approach.

## 8.2 Block-coordinate descent

Plugging in the affine parametrization (8.1) above, we solve the following WSP:

$$\begin{aligned} \min_{\mathbf{a}, \mathbf{x}, \mathbf{y}} \quad & w_{\text{eq}} \, \Sigma_{j=1}^{N} r_{\text{PD}}(\mathbf{x}^{(j)}, \mathbf{y}^{(j)}, \mathbf{p}^{(j)} \,|\, \mathbf{a}) + w_{\text{obs}} \, \Sigma_{j=1}^{N} \phi(g(\mathbf{x}^{(j)}, \mathbf{p}^{(j)}) - \mathbf{z}^{(j)}) \\ \text{s.t.} \quad & \mathbf{x}^{(j)} \in \mathcal{K}(\mathbf{p}^{(j)}), \quad \forall j \\ & \mathbf{A}(\mathbf{p}^{(j)})^T \mathbf{y}^{(j)} \leq F(\mathbf{x}^{(j)}, \mathbf{p}^{(j)} \,|\, \mathbf{a}), \quad \forall j \\ & \mathbf{a} \in \mathcal{A} \end{aligned}$$

where the dependencies in $\mathbf{a}$ are made explicit. Since the size of the inverse problem increases linearly with the number of observations $N$, but is separable into $N$ sub-problems with respect to the variables $\{\mathbf{x}^{(j)}, \mathbf{y}^{(j)}\}_{j=1,\cdots,N}$, we suggest to apply a Block-Coordinate Descent (BCD) algorithm to solve the WSP while avoiding the curse of dimensionality, see Algorithm 2. For the BCD, we cyclically update the $N$ vectors $\{\mathbf{x}^{(j)}\}_{j=1,\cdots,N}$, the $N$ vectors $\{\mathbf{y}^{(j)}\}_{j=1,\cdots,N}$, and the parameter vector $\mathbf{a}$. The sub-problems are:

$$\begin{aligned} \min_{\mathbf{x}^{(j)}} \quad & w_{\text{eq}} F(\mathbf{x}^{(j)}, \mathbf{p}^{(j)} \,|\, \mathbf{a})^T \mathbf{x}^{(j)} + w_{\text{obs}}\, \phi(g(\mathbf{x}^{(j)}, \mathbf{p}^{(j)}) - \mathbf{z}^{(j)}) \\ \text{s.t.} \quad & \mathbf{x}^{(j)} \in \mathcal{K}(\mathbf{p}^{(j)}) \\ & \mathbf{A}(\mathbf{p}^{(j)})^T \mathbf{y}^{(j)} \le F(\mathbf{x}^{(j)}, \mathbf{p}^{(j)} \,|\, \mathbf{a}) \end{aligned} \qquad (8.2)$$

$$\min_{\mathbf{y}^{(j)}} -\mathbf{b}^T \mathbf{y}^{(j)} \quad \text{s.t.} \quad \mathbf{A}(\mathbf{p}^{(j)})^T \mathbf{y}^{(j)} \le F(\mathbf{x}^{(j)}, \mathbf{p}^{(j)} \,|\, \mathbf{a}) \qquad (8.3)$$

$$\begin{aligned} \min_{\mathbf{a} \in \mathcal{A}} \quad & \sum_{i=1}^{N} F(\mathbf{x}^{(j)}, \mathbf{p}^{(j)} \,|\, \mathbf{a})^T \mathbf{x}^{(j)} \\ \text{s.t.} \quad & \mathbf{A}(\mathbf{p}^{(j)})^T \mathbf{y}^{(j)} \le F(\mathbf{x}^{(j)}, \mathbf{p}^{(j)} \,|\, \mathbf{a}), \quad \forall j \end{aligned} \qquad (8.4)$$

We note that steps 3 and 4 in Algorithm 2 can be done in parallel.

---

**Algorithm 2** BCD($\cdot$) Block descent algorithm for the inverse problem

---

1: **while** stopping criteria not met **do**
2:     $t := t + 1$
3:     $\mathbf{x}^{(j,t+1)} :=$ solution to (8.2) at $(\mathbf{y}^{(j)}, \mathbf{a}) = (\mathbf{y}^{(j,t)}, \mathbf{a}^{(t)})$ for $j = 1, \cdots, N$.
4:     $\mathbf{y}^{(j,t+1)} :=$ solution to (8.3) at $(\mathbf{x}^{(j)}, \mathbf{a}) = (\mathbf{x}^{(j,t+1)}, \mathbf{a}^{(t)})$ for $j = 1, \cdots, N$.
5:     $\mathbf{a}^{(t+1)} :=$ solution to (8.4) at $(\mathbf{x}^{(j)}, \mathbf{y}^{(j)}) = (\mathbf{x}^{(j,t+1)}, \mathbf{y}^{(j,t+1)})$ for all $j$.

---

# Chapter 9

# Application to Traffic Assignment

## 9.1 Model

A classic application of VI and CO is the traffic assignment problem, see, *e.g.* [22] for more details. Given road network modeled as a directed graph $(\mathcal{V}, \mathcal{E})$, with vertex set $\mathcal{V}$ and directed edge set $\mathcal{E}$, and a set of commodities $\mathcal{W} \subseteq \mathcal{V} \times \mathcal{V}$, a flow rate $d_k$ of a commodity $k$ must be routed from $s_k$ to $t_k$ for each $k = (s_k, t_k) \in \mathcal{C}$. The $k$-th *commodity flow vector* $\mathbf{x}^{(k)} = [x_e^{(k)}]_{e \in \mathcal{E}} \in \mathbb{R}_+^{\mathcal{E}}$ is feasible if it satisfies the flow equation at every vertex $i \in \mathcal{V}$:

$$\sum_{j \,:\, (j,i) \in \mathcal{E}} x_{(j,i)}^{(k)} - \sum_{j \,:\, (i,j) \in \mathcal{E}} x_{(i,j)}^{(k)} = \begin{cases} -d_k & \text{if} \quad i = s_k \\ d_k & \text{if} \quad i = t_k \\ 0 & \text{otherwise} \end{cases} \tag{9.1}$$

In matrix form, $\mathbf{x}^{(k)}$ is feasible if $\mathbf{N}\mathbf{x}^{(k)} = \mathbf{b}^{(k)}$, $\mathbf{x}^{(k)} \geq 0$, where $\mathbf{N}$ is the node-arc incidence matrix and $\mathbf{b}^{(k)} \in \mathbb{R}^{\mathcal{V}}$ the demand vector associated to commodity $k$ with entries such that $b_{s_k}^{(k)} = -d_k$, $b_{t_k}^{(k)} = d_k$, and $b_i^{(k)} = 0$, $\forall i \neq s_k, t_k$. Stacking everything together, we can simply rewrite the flow equations as $\mathbf{A}\mathbf{x} = \mathbf{b}$, $\mathbf{x} \geq 0$, where $\mathbf{x} = [x_e^{(k)}]_{e \in \mathcal{E}, k \in \mathcal{C}}$ is the overall flow vector. Following [6], the cost $c_e(x_e)$ of a road segment $e$ only depends on the flow $x_e$ of vehicles on this segment, where $x_e$ is expressed as $x_e = \sum_{k \in \mathcal{C}} x_e^{(k)}$, the sum of all the commodity flows. The cost functions $c_e(\cdot)$ are assumed to be continuous, positive, non-decreasing, and Beckmann et al. [6] proved that the *User Equilibrium* (UE), defined by [25], exists and is solution to the $\mathrm{CO}(\mathcal{K}, f)$ with potential:

$$f(\mathbf{x}) = \sum_{e \in \mathcal{E}} \int_0^{\sum_{k \in \mathcal{C}} x_e^{(k)}} c_e(u) du \tag{9.2}$$

However, cost functions $c_e$ are in general unknown, other as through empirical modeling such as the BPR function, while total flows $x_e = \sum_{k \in \mathcal{C}} x_e^{(k)}$ are measurable, but only on a small subset of arcs in the network, due to the cost of deploying and maintaining a sensing

infrastructure in a large urban area. With $g(\cdot)$ our fixed observation function (due to a fixed sensing infrastructure), we want to estimate delay functions from partial and noisy observations $\mathbf{z}^{(j)} = g(\mathbf{x}^{(j)}) + \mathbf{w}^{(j)}$ of flows $\mathbf{x}^{(j)}$ associated to different traffic demands $\mathbf{b}(\mathbf{p}^{(j)})$ and with noise $\mathbf{w}^{(j)}$, where each superscript $j$ refers to different demand levels, *e.g.*, morning or evening commutes. The imputed delay functions can be used to control or re-design the road network. See Figure 9.1 for an example.
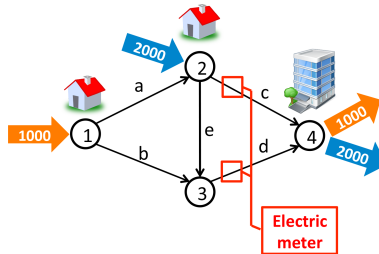


Figure 9.1: **Example of a morning commute on a simple road network with arcs $\{a, b, c, d, e\}$, and two commodities $1$ and $2$ with commodity flows $x_a^{(k)}, x_b^{(k)}, x_c^{(k)}, x_d^{(k)}, x_e^{(k)}$, $k \in \{1, 2\}$. A flow of 1000 veh/hour in $c_1$ is known to be routed along the shortest paths from nodes 1 to 4 and a flow of 2000 veh/hour in commodity $c_2$ is routed from 2 to 4, resulting in a UE flow on the network. Given only measurement of $z_1 = (x_c^{(1)} + x_c^{(2)}) + w_1$ and $z_2 = (x_d^{(1)} + x_d^{(2)}) + w_2$ with noise $w_1$, $w_2$, how can we impute the delay functions on each arc?**

## 9.2   Parametrization

We want to fit polynomial edge cost functions that are positive and non-decreasing. Hence we use the following parametrization, for all $e \in \mathcal{E}$:

$$c_e(x_e \,|\, \mathbf{a}) = d_e + d_e \sum_{i=1}^{K} a_i (x_e/m_e)^i, \quad \mathbf{a} = [a_i]_{i=1}^{K} \in \mathbb{R}_+^K \tag{9.3}$$

where $m_e$ is the capacity of road segment $e$ (typically proportional to the numer of lanes), and $d_e$ is the known free-flow travel time. Here, $d_e$ is the shift discussed in Section 8.1 to restrict the parameters $a_i$. The potential function $f$ (which does not depend on the parameter $\mathbf{p}$) is then, using the expression in (9.2):

$$f(\mathbf{x} \,|\, \mathbf{a}) = f_0(\mathbf{x}) + \sum_{i=1}^{K} a_i f_i(\mathbf{x})$$

$$f_i(\mathbf{x}) = \sum_{e \in \mathcal{E}} \frac{d_e}{m_e^i} \int_0^{\sum_{k \in \mathcal{C}} x_e^{(k)}} u^i du = \sum_{e \in \mathcal{E}} \frac{d_e}{m_e^i} \frac{\left( \sum_{k \in \mathcal{C}} x_e^{(k)} \right)^{i+1}}{i+1} \quad i = 0, 1, \cdots, K$$

We are now in position to use our method with the basis map functions:

$$F_i(\mathbf{x}) = \nabla f_i(\mathbf{x}) = [\partial f_i(\mathbf{x})/\partial x_e^{(k)}]_{e \in \mathcal{E}, k \in \mathcal{C}} = \left[ d_e \frac{\left( \sum_{k \in \mathcal{C}} x_e^{(k)} \right)^i}{m_e^i} \right]_{e \in \mathcal{E}, k \in \mathcal{C}} \tag{9.4}$$
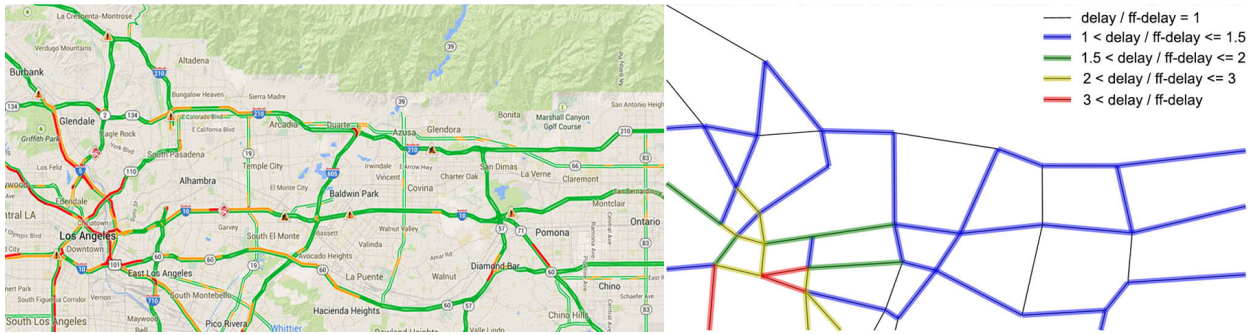
## 9.3 Numerical experiments



Figure 9.2: **Left: Highway network of L.A. in morning rush hour on 2014-06-12 at 9:14 AM from Google Maps; right: The network in UE with the resulting delays under demand 1.2\*b. The congested area is near central L.A.**

We consider the highway network near Los Angeles with 44 nodes and 122 arcs; see Figure 9.2. The roads characteristics (geometry, capacity, free flow delay) are obtained from OpenStreetMaps. The OD demands $\mathbf{b}$ are based on data from the Census Bureau and calibrated to represent a static morning rush hour model. We consider $N = 4$ equilibria $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \mathbf{x}^{(4)}$ associated to four demand vectors $\mathbf{b}(\mathbf{p}^{(j)}) \in \mathbb{R}^{|\mathcal{C}||\mathcal{V}|}$, $j \in \{1, 2, 3, 4\}$ obtained by scaling $\mathbf{b}$ with respective factors .5, 0.8, 1, 1.2. The measurements are obtained by solving the traffic assignment problem:

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{b}(\mathbf{p}), \ \mathbf{x} \geq 0 \tag{9.5}$$

with potential function $f$ given by (9.2), constraints $\mathbf{A}$ given by (9.1), demand vectors $\mathbf{b}(\mathbf{p}^{(j)})$, $j \in \{1, 2, 3, 4\}$, and for two types of delay functions:

$$c^{\text{poly}}(x_e) = d_e(1 + 0.15(x_e/m_e)^4) \tag{9.6}$$

$$c^{\text{hyper}}(x_e) = 1 - 3.5/3 + 3.5/(3 - x_e/m_e) \tag{9.7}$$

where (9.6) is estimated by the Bureau of Public Roads (BPR), and (9.7) is hyperbolic delay similar to the BPR one. These two functions are considered the ground truth delay functions

and we want to recover them from the observations $\mathbf{z}^{(j)} = g(\mathbf{x}^{(j)}) = [\sum_{k \in \mathcal{C}} x_e^{(k)}]_{e \in \mathcal{E}^{\text{obs}}}$, where $\mathcal{E}^{\text{obs}} \subseteq \mathcal{E}$ is the set of observed edge flows. We normalize $r_{\text{eq}}$ and $r_{\text{obs}}$ and solve the WSP (5.1) using the BCD algorithm discussed in Section 8.2. For $w_{\text{obs}} = 0.001, 0.01, 0.1, 0.5, 0.9, 0.99, 0.999$ and $w_{\text{eq}} = 1 - w_{\text{obs}}$, Figure 9.3 provides the error $\sum_{j=1}^{N} \|\mathbf{x}^{(j)} - \hat{\mathbf{x}}(\mathbf{p}^{(j)})\|$, where $\mathbf{x}^{(j)}$ are the ground-truth equilibrium flows and $\hat{\mathbf{x}}(\mathbf{p}^{(j)})$ the estimated ones.
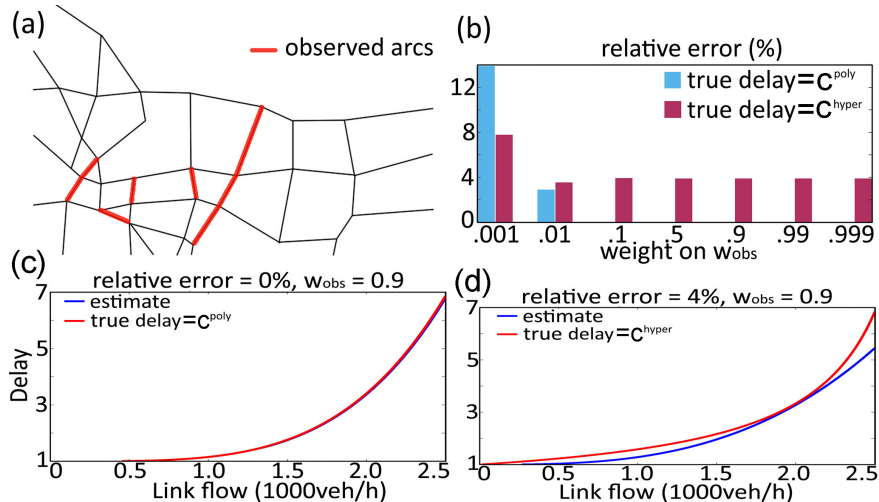


Figure 9.3: **Imputation of the delay maps $c^{\text{poly}}$, $c^{\text{hyper}}$ with parametric map given by (9.3). The relative error on the flow predicted by the imputed map is small for $w_{\text{obs}}$ large enough as shown in (b). With accurate measurements, we suggest to solve the WSP with $w_{\text{obs}} = 0.9$, which gives the estimated cost function for the BPR cost function in (c) and hyperbolic cost function in (d).**

In a second experiment, we study the sensitivity of our estimation algorithm to four sets of observed links, see Figure 9.4. The parameters $\mathbf{a}$ imputed by our latency inference methodology give a delay function $1 + \sum_{i=1}^{6} a_i x^i$ for each of the four sensor configurations. In case 1, we have a very good match between the estimated delay function and the true one because we observe the entire network, while in case 4, the measurements do not provide additional information because they are already contained in the given OD demands, see Figure 9.4.
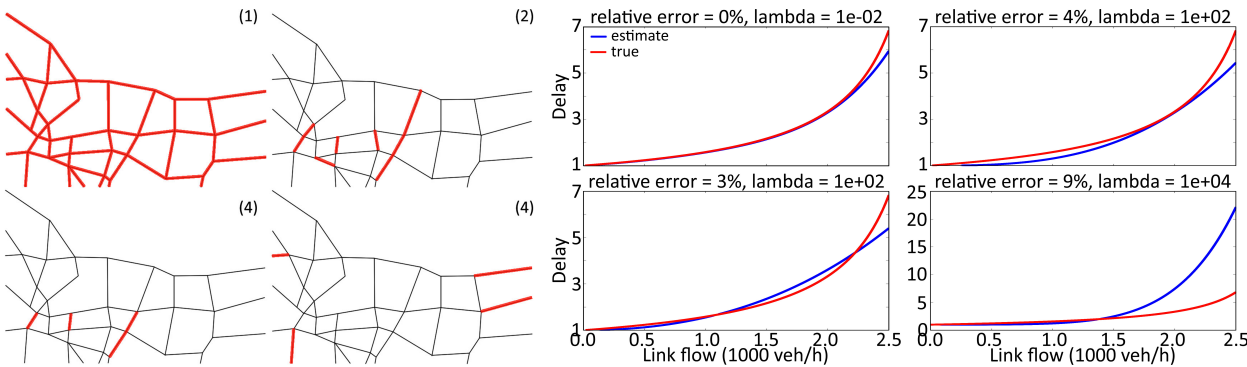
Figure 9.4: **Left: the 4 sensor configurations: (1) all arcs are observed; (2) 10 arcs are observed in the congested area; (3) 4 arcs are observed in the congested area; (4) 4 arcs are observed at the boundaries of the region, where the inflows are already known from the OD demands.**

# Chapter 10

# Application to Consumer Utility

## 10.1 Model

We also consider an oligopoly in which $n$ firms produce each one a product indexed by $i = 1, \cdots, n$ with prices $\mathbf{p} = [p_i]_{i=1}^n$. We suppose that the consumer purchases a quantity $x_i$ of product $i$ in order to maximize a non-decreasing and concave utility function $U(\mathbf{x})$ minus the price paid $\mathbf{p}^T \mathbf{x}$, where $\mathbf{x} = [x_i]_{i=1}^n$ is the overall demand, hence the optimization problem and parametric map:

$$\min_{\mathbf{x} \geq 0} f(\mathbf{x}) = \mathbf{p}^T \mathbf{x} - U(\mathbf{x}) \quad \implies \quad F(\mathbf{x}, \mathbf{p}) = \mathbf{p} - \nabla U(\mathbf{x}) \tag{10.1}$$

However, the utility $U : \mathbb{R}^n \to \mathbb{R}$ is not known in practice, and the inverse problem consists in imputing $U$ based on $N$ observations of pairs $(\mathbf{p}^{(j)}, \mathbf{x}^{(j)})$, $j = 1, \cdots, N$ of prices and associated demands. The imputed utility $U$ is then used by company producing $i$ to set a price $p_i$ in order to achieve a target consumer demand $x_i^{\text{des}}$ in its product. In oligopolies, the price of each product is publicly available and each firm in general knows its own demand $x_i$, however it may only have partial information on other demands. For example, if there are $n = 5$ firms and consumer demand in product 1 produced by firm 1 is not known, then we only observe the vector $\mathbf{z} = g(\mathbf{x}) = [x_2, x_3, x_4, x_5]^T$.

## 10.2 Parametrization

Similarly to [17], we consider a quadratic parametrization for the utility $U$, *i.e.* $U(\mathbf{x} \,|\, Q, \mathbf{r}) = \mathbf{x}^T Q \mathbf{x} + 2\mathbf{r}^T \mathbf{x}$, hence the parametric potential is

$$f(\mathbf{x}, \mathbf{p} \,|\, Q, \mathbf{r}) = \mathbf{p}^T \mathbf{x} - \left( \mathbf{x}^T Q \mathbf{x} + 2\mathbf{r}^T \mathbf{x} \right), \quad (Q, \mathbf{r}) \in \mathcal{A} \tag{10.2}$$

$$\mathcal{A} = \{ (Q, \mathbf{r}) : Q\mathbf{x}_{\max} + \mathbf{r} \geq 0, \ \mathbf{r} \geq 0, \ Q \preceq 0 \} \tag{10.3}$$

where $\mathcal{A}$ is chosen such that $U(\cdot \mid Q, \mathbf{r})$ is concave and non-decreasing on the demand range $[0, \mathbf{x}_{\max}]$. The parametric map $F(\cdot, \mathbf{p} \mid Q, \mathbf{r})$ is then:

$$F(\mathbf{x}, \mathbf{p} \mid Q, \mathbf{r}) = \mathbf{p} - 2Q\mathbf{x} - 2\mathbf{r} \qquad (10.4)$$
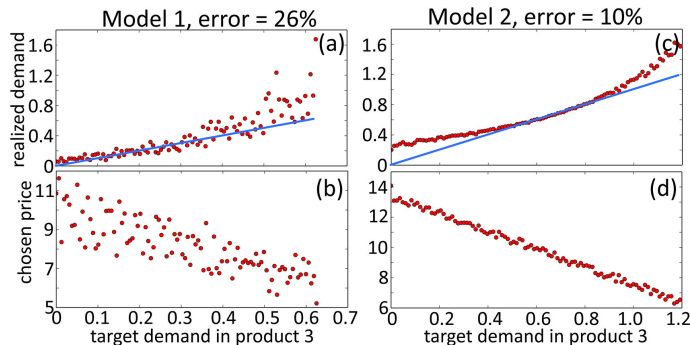
## 10.3  Numerical experiments



Figure 10.1: **Use of the imputed utility to price product 3 for different target demands** $x_3^{\mathbf{des}}$. **In (b), the prices are scattered due to correlations with other prices in model 1, while in (d), the prices vary linearly with** $x_3^{\mathbf{des}}$ **since the prices in model 2 are more uncorrelated. In (a), (c) the blue line is the** $x = y$ **line. For both models, the imputed utility performs well with relative errors of 26% and 10% on the training data and target demands** $\mathbf{x}_3^{\mathbf{des}}$ **close to realized demands** $\mathbf{x}_3^{\mathbf{real}}$.

We consider the case of $n = 5$ firms competing for the same market. At the time period $j$, let $\mathbf{x}^{(j)} \in \mathbb{R}_+^5$ be the consumer demand in response to the prices $\mathbf{p}^{(j)} \in \mathbb{R}_+^5$ set by each firm, sampled uniformly as i.i.d. random vectors in $[8, 12]^5$. We assume that the third firm only observes the demand $\mathbf{z}^{(j)} = [x_2^{(j)}, \cdots, x_5^{(j)}]^T$ in products from firms 2, 3, 4, 5 over $N = 200$ time periods along with the prices $\mathbf{p}^{(j)}$. The demand $\mathbf{x}^{(j)}$ incurred by prices $\mathbf{p}^{(j)}$ are assumed to be solution of the convex optimization model (10.1) with underlying consumer utility function $U^{\mathrm{real}}(\mathbf{x}) = \mathbf{1}^T \sqrt{\mathbf{A}\mathbf{x} + \mathbf{b}}$. Firm 3 wants to impute $U^{\mathrm{real}}$ using the parametric utility given by (10.2). The numerical results are shown in Figure 10.1 with two models for $\mathbf{A} = 50(\mathbf{I} + \mathbf{B})$ in $U^{\mathrm{real}}$: *model 1* where $\mathbf{B}_{ij}$ is sampled uniformly in $[0, 0.3]$ for $i \neq j$, and *model 2* where $\mathbf{B}_{ij}$ is sampled from $0.5 \cdot$Bernoulli$(0.3)$.

# Chapter 11

# Remark: possible extensions to convex cones

The results in the present article can be generalized to conic representable sets, such as the feasible set for the primal variables, which can be restricted to the generalized polyhedron, with $\mathcal{C} \subseteq \mathbb{R}^n$ a convex cone:

$$\mathcal{K} = \{\mathbf{x} \in \mathbb{R}^n \,|\, \mathbf{A}\mathbf{x} = \mathbf{b},\, \mathbf{x} \in \mathcal{C}\} \tag{11.1}$$

In the article, $\mathcal{C}$ is the non-negative orthant $\mathbb{R}^n_+$, but more complex cones include the second-order cone or the cone of positive semi-definite matrices. Using the standard notation: $\mathbf{x} \geq_{\mathcal{C}} \mathbf{y} \iff \mathbf{x} - \mathbf{y} \geq_{\mathcal{C}} 0 \iff \mathbf{x} - \mathbf{y} \in \mathcal{C}$, the primal-dual and KKT systems in Section 2 can be generalized to:

$$
\begin{array}{llll}
\text{Primal-dual} & F(\mathbf{x})^T \mathbf{x} = \mathbf{b}^T \mathbf{y} & \text{KKT} & F(\mathbf{x}) = \mathbf{A}^T \mathbf{y} + \boldsymbol{\pi} \\
\text{system:} & F(\mathbf{x}) \geq_{\mathcal{C}} \mathbf{A}^T \mathbf{y} & \text{system:} & \mathbf{A}\mathbf{x} = \mathbf{b} \\
& \mathbf{A}\mathbf{x} = \mathbf{b},\, \mathbf{x} \in \mathcal{C} & & \mathbf{x} \in \mathcal{C},\, \boldsymbol{\pi} \in \mathcal{C},\, \mathbf{x}^T \boldsymbol{\pi} = 0
\end{array}
\tag{11.2}
$$

Hence, for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, primal feasibility becomes $\mathbf{A}\mathbf{x} = \mathbf{b},\, \mathbf{x} \in \mathcal{C}$ and dual feasibility becomes $F(\mathbf{x}) \geq_{\mathcal{C}} \mathbf{A}^T \mathbf{y}$. All the results in the article hold at the exception of Lemma 7.1, and Theorems 7.1 and 7.2, which require that $\mathcal{C} = \mathbb{R}^n_+$.

# Bibliography

[1] Pieter Abbeel and Andrew Ng. "Apprenticeship Learning via Inverse Reinforcement Learning". In: *21th International Conference on Machine Learning* 69 (2004), pp. 1–8.

[2] Daniel Ackerberg et al. "Econometric Tools for Analyzing Market Outcomes". In: *Handbook of Econometrics* 6 (2007).

[3] M. Aghassi, D. Bertsimas, and G. Perakis. "Solving asymmetric variational inequalities via convex optimization". In: *Operations Research Letters 34* 5 (2006), pp. 481–490.

[4] P. Bajari, C. Benkard, and J. Levin. "Estimating dynamic models of imperfect competition". In: *Econometrica* 75 (2007), pp. 1331–1370.

[5] B. Bank et al. *Non-linear parametric optimization.* Basel - Boston - Stuttgart: Birkhauser Verlag, 1983.

[6] M. Beckmann, C. B. McGuire, and C. B. Winsten. *Studies in the Economics of Transportation.* Ed. by Yale University Press. 1956.

[7] Dimitris Bertsimas, Vishal Gupta, and Ioannis Ch. Paschalidis. "Data-Driven Estimation in Equilibrium Using Inverse Optimization". In: *Mathematical Programming* 153 (2015), pp. 595–633.

[8] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization.* Cambridge University Press, Mar. 2004.

[9] S. Boyd et al. *Linear Matrix Inequalities in Systems and Control Theory.* Philadelphia: SIAM books, 1994.

[10] D. Burton, W. R. Pulleyblank, and Ph. L. Toint. "The inverse shortest path problem with upper bounds on shortest path costs". In: *Network Optimization* 450 (1997), pp. 156–171.

[11] Y. Chen and M. Florian. *Congested O-D trip demand adjustment problem: bilevel programming formulation and optimality conditions.* Kluwer Academic Publishers, Dordrecht, 1998, pp. 1–22.

[12] S. Dempe. *Foundations of Bilevel Programming.* Springer, 2002.

[13] M. Ehrgott. *Multicriteria optimization.* Springer Verlag, Berlin, 2005.

[14] F. Facchinei, H. Jiang, and L. Qi. "A smoothing method for mathematical programs with equilibrium constraints". In: *Mathematical Programming* 85 (1999), pp. 107–134.

[15] Francisco Facchinei and Jong-Shi Pang. *Finite-Dimensional Variational Inequalities and Complementarity Problems*. Springer, New York, 2003.

[16] Houyuan Jiang and Daniel Ralph. "QPECgen, a MATLAB Generator for Mathematical Programs with Quadratic Objectives and Affine Variational Inequality Constraints". In: *Computational Optimization and Applications* 13 (1999), pp. 25–59.

[17] Arezou Keshavarz, Yang Wang, and Stephen Boyd. "Imputing a Convex Objective Function". In: *IEEE International Symposium on Intelligent Control* (2011), pp. 613–619.

[18] Zhi-Quan Luo, Jong-Shi Pang, and Daniel Ralph. *Mathematical Programs with Equilibrium Constraints*. Cambridge University Press, Nov. 1996.

[19] R. T. Marler and J. S. Arora. "Survey of multi-objective optimization methods for engineering". In: *Structural and Multidisciplinary Optimization* 26.6 (2004), pp. 369–395.

[20] Andrew Ng and Stuart Russell. "Algorithms for inverse reinforcement learning". In: *17th International Conference on Machine Learning* (2000), pp. 663–670.

[21] Jong-Shi Pang. "A Posteriori Error Bounds for the Linearly-Constrained Variational Inequality Problem". English. In: *Mathematics of Operations Research* 12.3 (1987), pp. 474–484. ISSN: 0364765X.

[22] Michael Patriksson. *The Traffic Assignment Problem - Models and Methods*. Dover Publications, 2015.

[23] J. Rust. "Structural Estimation of Markov decision processes". In: *Handbook of Econometrics* 4 (1994), pp. 3081–3143.

[24] Gesualdo Scutari et al. "Convex Optimization, Game Theory, and Variational Inequality Theory". In: *IEEE Signal Processing Magazine* 27 (May 2010), pp. 35–49.

[25] J. G. Wardrop and J. I. Whitehead. "Correspondence. Some Theoretical Aspects of Road Traffic Research". In: *ICE Proceedings: Engineering Divisions 1* 1 (1952), pp. 325–378.

[26] J. J. Ye and D. L. Zhu. "Optimality conditions for bilevel programming problems". In: *Optimization* 33 (1995), pp. 9–27.