

# A BILEVEL PROGRAMMING APPROACH TO OPTIMAL PRICE SETTING

*À Alain, en souvenir de randonnées bien arrosées sur la rivière l'Assomption  
et dans la réserve Mastigouche, aux temps héroïques du GERAD.*

Patrice Marcotte  
*DIRO and CRT*  
*Université de Montréal*  
marcotte@iro.umontreal.ca

Gilles Savard  
*MAGI and GERAD*  
*École Polytechnique de Montréal*  
gilles.savard@polymtl.ca

**Keywords:** Bilevel programming. Pricing. Revenue Management. Transportation.

**Abstract** In this paper, we survey applications and algorithms pertaining to an important class of price setting problems formulated as bilevel programs.

## **Introduction. Bilevel Programming.**

The present paper focuses on a specific class of bilevel problems that arise naturally when tariffs, tolls or devious taxes are imposed on a set of commodities. Not only does this class encompass several important optimization problems encountered in the transportation, telecommunication and airline industries, but its structure makes it amenable to efficient solution techniques. Our aim is to present several such models and briefly discuss algorithmic approaches, either exact or heuristic, that can be applied to large scale problems within this class.

Bilevel programming is the adequate framework for modelling asymmetric games where one player (the *leader*) calls the shots first, taking into account the optimal reaction of the second player (the *follower*).

If one denotes by  $x$  (respectively  $y$ ) the decision vector of the leader (respectively the follower), a bilevel program can be expressed as

$$\begin{aligned} & \min_{x,y} && f(x, y) \\ \text{subject to} &&& (x, y) \in X \\ &&& y \in \mathcal{S}(x), \end{aligned}$$

where  $\mathcal{S}(x)$  denotes the set of optimal solutions of a mathematical program parameterized in the leader's vector  $x$ , i.e.,

$$\begin{aligned} \mathcal{S}(x) &= \arg \min_y && g(x, y) \\ &\text{subject to} && (x, y) \in Y. \end{aligned}$$

The above formulation implicitly assumes that, if the lower level problem admits multiple solutions, ties are broken in favor of the leader. Alternative situations, where the follower reacts in an antagonistic fashion, have been analyzed by Loridan and Morgan [16].

Bilevel programs are closely related to mathematical programs with equilibrium constraints (MPECs), where the lower level corresponds to an equilibrium problem. Indeed, whenever the objective  $g$  of the lower level program is differentiable and convex in  $y$  and the set  $Y$  is convex,  $y \in \mathcal{S}(x)$  if and only if  $(x, y) \in Y$  and satisfies the variational inequality

$$\langle \nabla_y g(x, y), y - y' \rangle \leq 0$$

for all  $y'$  such that  $(x, y') \in Y$ . Letting  $Y(x) = \{y : (x, y) \in Y\}$  we obtain the one-level optimization formulation

$$\begin{aligned} & \min_{x,y} && f(x, y) \\ \text{subject to} &&& (x, y) \in X \\ &&& y \in Y(x) \\ &&& \langle \nabla_y g(x, y), y - y' \rangle \leq 0 \quad \forall y' \in Y(x) \end{aligned}$$

which subsumes the more general form of an MPEC:

$$\begin{aligned} & \min_{x,y} && f(x, y) \\ \text{subject to} &&& (x, y) \in X \\ &&& y \in Y(x) \\ &&& \langle G(x, y), y - y' \rangle \leq 0 \quad \forall y' \in Y(x), \end{aligned}$$

where the vector function  $G$  need not be a gradient mapping with respect to the variable  $y$ . Conversely, an MPEC can be reformulated as

a standard bilevel program by noting that a vector  $y$  is solution of the lower level variational inequality if and only if it globally minimizes, with respect to the argument  $y$ , the strongly convex function  $\text{gap}(x, y)$  defined as (see Fukushima [7]):

$$\text{gap}(x, y) = \max_{y' \in Y(x)} \langle G(x, y), y - y' \rangle - \frac{1}{2} \|y - y'\|^2.$$

The reader interested in the theory and applications of bilevel programming and MPEC is referred to the recent books by Shimizu, Ishizuka and Bard [20] and by Luo, Pang and Ralph [17].

Being generically nonconvex and nonsmooth, bilevel programs are difficult optimization problems. Even in the simple situation where both objectives are affine and the constraint sets are polyhedral, determining whether a solution is locally optimal is strongly NP-hard. This explains why global optimization techniques such as implicit enumeration, cutting planes or metaheuristics have been proposed for its solution (see e.g. [8] and [10]). These are most successful when the set  $\mathcal{S}(x)$  assumes a piecewise polyhedral structure, for instance when the lower level problem takes the form of a convex quadratic program. In the absence of such property, two main lines of attack have been pursued. The first, based on sensitivity analysis, adapts descent methods that are compatible with the optimality requirements of the follower, and relies on recent nonsmooth analysis results; the work of Kocvara, Outrata and Zowe [19] is typical of this trend. A drawback of this approach is that, even under strong regularity assumption, it may fail to uncover even a local optimum for the bilevel program.

A second approach applies standard optimization techniques to a one-level reformulation, smooth or not, of the bilevel program. A recent member of this family have been proposed by Scholtes and Stohr [21].

While interesting in their own right, the above approaches cannot solve the large scale applications that we have in mind, thereby justifying the development of specific algorithms.

## **A generic price setting problem**

Let  $x$  and  $y$  be real vectors that specify the respective levels of taxed and untaxed activities (commodities or services), and  $T$  denote the tax vector attached to  $x$ . For a given vector  $T$ , in control of the leader, the follower strives to minimize its operating costs, while the leader seeks to maximize the revenues raised from taxes. If one denote by  $F$  and  $f$  the leader's and follower's respective objective functions, the leader

maximizes his profit by solving the bilevel program

$$\begin{aligned} & \max_{T,x,y} && F(T, x, y) \\ \text{subject to} &&& (x, y) \in \arg \min_{(x', y') \in \Pi} f(T, x', y') \end{aligned} \quad (1.1)$$

where  $\Pi$  represents the constraint set of the second level player. From now on, we will record programs of the form (1.1) in the vertical format:

$$\begin{aligned} & \max_T && F(T, x, y) \\ & \min_{x,y} && f(T, x, y) \\ \text{subject to} &&& (x, y) \in \Pi. \end{aligned} \quad (1.2)$$

This seemingly simplistic model can cover a wide variety of situations. For instance the vector  $T$  may embody subsidies as well as taxes, while the vectors  $x$  and  $y$  may represent consumption or production levels. Alternatively, the lower level can represent the group behavior of economic agents competing for scarce resources; if the equilibrium states of the system are the solutions of a variational inequality parameterized in the leader's decision variables, we obtain an MPEC.

Let us first consider a basic model where the leader's revenues are proportional to tax and consumption levels, and where all constraints are linear. The resulting bilevel program, where both objectives are bilinear, takes the form

$$\begin{aligned} & \max_T && Tx \\ & \min_{x,y} && (c_1 + T)x + c_2y \\ \text{subject to} &&& A_1x + A_2y = b \\ &&& x, y \geq 0. \end{aligned} \quad (1.3)$$

Note that, from the leader's perspective, the objective function  $Tx$  is discontinuous at the points  $T$  that induce a change of optimal basis in the follower's linear program.

Assuming that the polyhedron  $\{(x, y) : A_1x + A_2y = b, x, y \geq 0\}$  is bounded and that the recourse polyhedron  $\{y : A_2y = b, y \geq 0\}$  is nonempty, the lower level always admits an optimal solution for any value of the tax vector  $T$ . Therefore, one can replace the lower level problem by its primal-dual optimality conditions to obtain the mathe-

mathematical program with linear and complementarity constraints

$$\begin{aligned}
 & \max_{T,x,y,\lambda} && Tx \\
 & \text{subject to} && A_1x + A_2y = b \\
 & && x, y \geq 0 \\
 & && \lambda A_1 \leq c_1 + T \\
 & && \lambda A_2 \leq c_2 \\
 & && (c_1 + T - \lambda A_1)x = 0 \\
 & && (c_2 - \lambda A_2)y = 0
 \end{aligned}$$

or the equivalent program

$$\begin{aligned}
 & \max_{T,x,y,\lambda} && \lambda b - (c_1x + c_2y) \\
 & \text{subject to} && A_1x + A_2y = b \\
 & && x, y \geq 0 \\
 & && \lambda A_1 \leq c_1 + T && (1.4) \\
 & && \lambda A_2 \leq c_2 \\
 & && (c_1 + T - \lambda A_1)x = 0 \\
 & && (c_2 - \lambda A_2)y = 0.
 \end{aligned}$$

It is not difficult to see that (1.4) admits an optimal solution of the form  $T = \lambda A_1 - c_1$ . Upon substitution, we obtain the simplified model

$$\begin{aligned}
 & \max_{x,y,\lambda} && \lambda b && - && (c_1x + c_2y) \\
 & \text{subject to} && \lambda A_2 \leq c_2 && && A_1x + A_2y = b && (1.5) \\
 & && && && x, y \geq 0 \\
 & && && && \boxed{(c_2 - \lambda A_2)y = 0}
 \end{aligned}$$

If we relax its complementarity constraint, (1.5) decomposes into two linear programs involving respectively the dual vector  $\lambda$  and the primal vectors  $x$  and  $y$ . The linear program associated with the primal variables corresponds to the lower level problem of (1.3) with taxes set at zero. The dual of the linear program associated with the dual variables is again the lower level linear program of (1.3), where the choice of activities is restricted to untaxed ones or, equivalently, where taxes are set to arbitrary high values.

Returning to the single-level formulation (1.5), one can penalize the complementarity constraint into the objective function, yielding a bilinear problem separable in the dual vector  $\lambda$  on the one hand, and in the

primal vectors  $x$  and  $y$  on the other hand:

$$\begin{aligned}
& \max_{x,y,\lambda} && \lambda b - (c_1x + c_2y) - M(c_2 - \lambda A_2)y \\
& \text{subject to} && \lambda A_2 \leq c_2 \\
& && A_1x + A_2y = b \\
& && x, y \geq 0.
\end{aligned} \tag{1.6}$$

Labbé, Marcotte and Savard [14] have established the existence of an exact penalty parameter  $M^*$  such that any optimal solution of (1.6) is also optimal for (1.5) whenever  $M$  exceeds  $M^*$ . Now, for fixed primal variables  $x$  and  $y$ , let us replace the objective function  $\lambda(b + MA_2y)$  by its dual objective; this yields the linear bilevel program:

$$\begin{aligned}
& \max_{x,y} && -c_1x - (M + 1)c_2y + c_2y' \\
& \text{subject to} && A_1x + A_2y = b \\
& && x, y \geq 0 \\
& && (1.7) \\
& \min_{y'} && c_2y' \\
& \text{subject to} && A_2y' = b + MA_2y \\
& && y' \geq 0.
\end{aligned}$$

An intuitive economic interpretation of the linear bilevel program (1.7) in terms of a *second best* alternative for the follower has been discussed in Labbé et al. [14].

We close this section with computational complexity considerations. In the absence of constraints on the tax vector  $T$ , we could not prove that (1.3) is NP-hard, although we suspect it is. However we proved that a variant of (1.3) where taxes are bounded from below is strongly NP-hard. The proof relies on a reduction from the ‘‘Hamiltonian Path’’ problem in a directed graph to a price setting problem with lower bound constraints.

## Toll setting

The problem of selecting optimal highway tolls clearly fits our price-setting framework. Let us consider a multicommodity network where each commodity  $k \in \mathcal{K}$  is associated with an origin-destination pair  $(o(k), d(k))$  of a transportation network  $G$  defined by a node set  $\mathcal{N}$  and an arc set  $\mathcal{A}$ , the latter being partitioned into the subset  $\mathcal{A}_1$  of toll arcs and the subset  $\mathcal{A}_2$  of toll-free arcs. With each arc  $a$  of  $\mathcal{A}_1$  is associated a generalized travel time composed of a fixed part  $c_{1a}$  representing the

travel time along arc  $a$  and an additional toll  $T_a$ , expressed in time units. Any toll-free arc  $a$  of  $\mathcal{A}_2$  bears a fixed unit travel cost  $c_{2a}$ . Tolls  $T_a$  cannot exceed a prescribed upper bound  $T_a^{\max}$ , which could be infinite.

A travel demand vector  $\{n^k\}_{k \in \mathcal{K}}$  induces the nodal demand vectors

$$b_i^k = \begin{cases} n^k & \text{if } i = o(k), \\ -n^k & \text{if } i = d(k), \\ 0 & \text{otherwise.} \end{cases}$$

The lower level variable  $x_a^k$  corresponds to the number of users of commodity  $k$  on arc  $a \in \mathcal{A}_1$  and the variable  $y_a^k$  to the number of users of commodity  $k$  on arc  $a \in \mathcal{A}_2$ .

Neglecting congestion effects, assuming that demand is fixed and that users minimize their individual generalized travel costs, the toll setting problem can be formulated as a bilevel program with bilinear objectives and linear constraints:

$$\begin{aligned} & \max_T \sum_{a \in \mathcal{A}_1} T_a \sum_{k \in \mathcal{K}} x_a^k \\ & \text{subject to } T_a \leq T_a^{\max} && \forall a \in \mathcal{A}_1 \\ & \min_{x,y} \sum_{k \in \mathcal{K}} \left( \sum_{a \in \mathcal{A}_1} (c_{1a} + T_a) x_a^k + \sum_{a \in \mathcal{A}_2} c_{2a} y_a^k \right) \\ & \text{subject to } \sum_{a \in i^+} (x_a^k + y_a^k) - \sum_{a \in i^-} (x_a^k + y_a^k) = b_i^k && \forall k \in \mathcal{K} \quad \forall i \in \mathcal{N} \\ & x_a^k \geq 0 && \forall k \in \mathcal{K} \quad \forall a \in \mathcal{A}_1 \\ & y_a^k \geq 0 && \forall k \in \mathcal{K} \quad \forall a \in \mathcal{A}_2, \end{aligned}$$

where  $i^+$  (respectively  $i^-$ ) denotes the set of arcs having  $i$  as their tail node (respectively head node).

One can introduce more realism into the previous model by considering congestion effects and/or a nonuniform distribution of the “time value of money” across the population. These features will be incorporated in the next applications.

## Price setting of telecommunication networks

Price and revenue optimization in the area of telecommunication is an active domain of research (see the book of Mitchell and Vogelsang [18]). Among various approaches to that important problem, one may mention those based on game theory (Cocchi et al. [4] and Bencheick [1]), opportunity costing (Wang et al. [23]) and network design (Girard

and Zidane [9]). In this section, we consider the optimal pricing of links in a packet-switched telecommunication network. We assume that users select a telecommunication provider according to two criteria, cost and quality of service, the latter being in direct relationship with the capacities of the links. The problem of the leader company is to price out the arcs of its subnetwork such as to maximize profit, while taking into account the user-optimized behavior of the customers. Whenever (i) the arrival rate of messages follows a Poisson process and is independent of service time, (ii) the length of a message is distributed according to an exponential random variable, then the average delay of a message through the network is given by the formula

$$T(f) = \frac{1}{\gamma} \sum_{a \in \mathcal{A}} \left( \frac{f_a}{C_a - f_a} + \mu f_a p_a \right), \quad (1.8)$$

where  $f_a$  denotes the total flow on arc  $a \in \mathcal{A}$ ,  $\gamma$  the total demand on the network,  $C_a$  the capacity of the arc  $a \in \mathcal{A}$ ,  $\mu$  the average message length and  $p_a$  the propagation delay along arc  $a \in \mathcal{A}$ . The first term in the summation, the node delay, reflects congestion at nodes. Under normal conditions, the second term (arc delay) is negligible and can safely be discarded.

Let  $d_a^k$  denote the average generalized cost along the arcs of the competing firms. Upon introduction of trade-off parameters  $\alpha_a$  that translate quality of service in terms of cost units, the price-setting problem takes the form of the bilevel program

$$\begin{aligned} & \max_T \quad \sum_{a \in \mathcal{A}_1} \sum_{k \in \mathcal{K}} T_a^k x_a^k \\ \text{subject to} \quad & T_a \leq T_a^{\max} \quad \forall a \in \mathcal{A}_1 \end{aligned}$$

$$\begin{aligned} & \min_{f, x, y} \quad \sum_{a \in \mathcal{A}} \frac{\alpha_a f_a}{\gamma(C_a - f_a)} + \sum_{k \in \mathcal{K}} \sum_{a \in \mathcal{A}_1} T_a^k x_a^k \\ & \quad + \sum_{k \in \mathcal{K}} \sum_{a \in \mathcal{A}_2} d_a^k y_a^k \\ \text{subject to} \quad & \sum_{a \in i^+} (x_a^k + y_a^k) - \sum_{a \in i^-} (x_a^k + y_a^k) = b_i^k \quad \forall k \in \mathcal{K} \quad \forall i \in \mathcal{N} \end{aligned}$$



$$\begin{aligned}
 f_a &= \sum_{k \in \mathcal{K}} x_a^k && \forall a \in \mathcal{A}_1 \\
 f_a &= \sum_{k \in \mathcal{K}} y_a^k && \forall a \in \mathcal{A}_2 \\
 x_a^k &\geq 0 && \forall k \in \mathcal{K} \quad \forall a \in \mathcal{A}_1 \\
 y_a^k &\geq 0 && \forall k \in \mathcal{K} \quad \forall a \in \mathcal{A}_2.
 \end{aligned}$$

In this formulation, arc prices are not required to be identical for all commodities. This opens room for discrimination of the users and would make the problem entirely separable by commodity, hence trivial to solve, were it not for the congestion effects.

### **Yield management in the airline industry**

Yield management in the airline industry has been an area of active research for the past few years (see [22]). It addresses four issues that deeply impact the industry revenues: forecasting, overbooking, seat allocation and pricing. Ideally, these four components should be part of an integrated profit-maximizing model. However, the complexity and the size of such a model would prevent the numerical resolution of problem instances of any realistic size. At the present time, the four issues are treated independently, even though the strong interaction between seat allocation and pricing is widely acknowledged. Our model, which addresses jointly the issues of seat allocation and pricing, is distinguished by three key features: its bilevel nature allows for an endogenous representation of the price-demand relationship, network interactions among the various carriers are explicitly considered, and the utility function of each user takes into account three criteria: fare, time and quality of service. Two parameters,  $\alpha$  and  $\beta$ , translate time and quality of service into cost units; these parameters are distributed across the population according to a density  $\phi^k(\alpha, \beta)$  which may vary by origin-destination pair  $k$ . This flexibility allows the model to assign a given user to a fare class according to his personal preference.

Let  $T$  denote the fare vector and  $x$  (respectively  $y$ ) denote the passenger flow vector for the leader airline (respectively the competition), indexed by booking classes  $b$  and flights  $f$  connecting the origin-destination pair  $k$ . The revenue of the leader airline is

$$f(T, x, y) = \sum_k \sum_{f \in \mathcal{C}_1} \sum_{b \in \mathcal{B}(f)} T_{bf}^k \int_{\alpha} \int_{\beta} x_{bf}^k(\alpha, \beta) d\alpha d\beta,$$

where  $\mathcal{C}_1$  is the set of flights of the leader company and  $\mathcal{B}(f)$  is the set of booking classes available on flight  $f$ . A passenger whose utility parameters are  $\alpha$  and  $\beta$  will select a flight that minimizes his generalized cost

$$\text{fare} + (\alpha \times \text{flight length}) + (\beta \times \text{flight quality}).$$

Hence, at the lower level, the passengers are assigned to the flights and booking classes that minimize the objective

$$g(x, y) = \sum_k \left\{ \sum_{f \in \mathcal{C}_1} \sum_{b \in \mathcal{B}(f)} \int_{\alpha} \int_{\beta} (T_{bf}^k + \alpha l_f + \beta q_{bf}) x_{bf}^k(\alpha, \beta) d\alpha d\beta + \sum_{f \in \mathcal{C}_2} \sum_{b \in \mathcal{B}(f)} \int_{\alpha} \int_{\beta} (T_{bf}^k + \alpha l_f + \beta q_{bf}) y_{bf}^k(\alpha, \beta) d\alpha d\beta \right\},$$

where  $\mathcal{C}_2$  is the set of flights of the competing firms, and  $l_f$  and  $q_{bf}$  the respective length and quality parameters associated with flight  $f$  and booking class  $b$ . The functions  $x_{bf}^k(\alpha, \beta)$  and  $y_{bf}^k(\alpha, \beta)$  are the the “user group”  $(\alpha, \beta)$ . Their integrals represent the respective market shares of the leader airline and the competition.

Let  $F_1(k)$  (respectively  $F_2(k)$ ) denote the set of flights of the leader (respectively the competition) available for an origin-destination  $k$  for which the demand is  $d_k$ . Let  $A(f)$  denote the set of legs that make up flight  $f$ . Taking into account the flow conservation and capacity constraints, we can formulate the yield management model as bilevel program involving both finite-dimensional (fares) and infinite-dimensional (flow densities) decision variables:

$$\begin{aligned}
 & \max_{T,x,y} \sum_k \sum_{f \in \mathcal{C}_1} \sum_{b \in \mathcal{B}(f)} T_{bf}^k \int_{\alpha} \int_{\beta} x_{bf}^k(\alpha, \beta) d\alpha d\beta \\
 \text{subject to } & T_{bf}^{\min} \leq T_{bf} \leq T_{bf}^{\max} \quad \forall f \in \mathcal{C}_1 \quad \forall b \in \mathcal{B}(f) \\
 & \min_{x,y} \sum_k \left\{ \sum_{f \in \mathcal{C}_1} \sum_{b \in \mathcal{B}(f)} \int_{\alpha} \int_{\beta} (T_{bf}^k + \alpha l_f + \beta q_{bf}) x_{bf}^k(\alpha, \beta) d\alpha d\beta \right. \\
 & \quad \left. + \sum_{f \in \mathcal{C}_2} \sum_{b \in \mathcal{B}(f)} \int_{\alpha} \int_{\beta} (T_{bf}^k + \alpha l_f + \beta q_{bf}) y_{bf}^k(\alpha, \beta) d\alpha d\beta \right\} \\
 \text{subject to } & x_a = \sum_k \sum_{f: a \in A(f)} \sum_{b \in \mathcal{B}(f)} \int_{\alpha} \int_{\beta} x_{bf}^k(\alpha, \beta) d\alpha d\beta \quad \forall a \in \mathcal{A}_1 \\
 & \phi^k(\alpha, \beta) d^k = \sum_{b \in \mathcal{B}(f)} \left\{ \sum_{f \in F_1(k)} x_{bf}^k(\alpha, \beta) + \sum_{f \in F_2(k)} y_{bf}^k(\alpha, \beta) \right\} \\
 & \quad \forall k \in \mathcal{K} \\
 & x_a \leq x_a^{\max} \quad \forall a \in \mathcal{A}_1 \\
 & x_{bf}^k(\alpha, \beta) \geq 0 \quad \forall k \in \mathcal{K} \quad \forall f \in \mathcal{C}_1 \quad \forall b \in \mathcal{B}(f) \\
 & y_{bf}^k(\alpha, \beta) \geq 0 \quad \forall k \in \mathcal{K} \quad \forall f \in \mathcal{C}_2 \quad \forall b \in \mathcal{B}(f).
 \end{aligned}$$

Note that the above model is pessimistic in the sense that the competition’s capacity is assumed to be illimited; this assumption makes sense only if the leader company’s market share is small.

### Traffic management through link tolls

On a transportation network, Wardrop’s user principle states that, at equilibrium, flows should be assigned to shortest routes with respect to current travel delays. If the network is heavily congested, such an assignment may lead to an inefficient use of the network capacity. This situation can be remedied by adopting a marginal cost pricing scheme which would induce the network users to behave in a socially optimal fashion. However, this scheme is not alone to possess that property, and one might prefer to implement a regulating policy that minimizes the total amount of tolls raised from the users (see e.g. Hearn and Ramana [11] and Larsson and Patriksson [15]). Denoting by  $x^k$  the flow vector associated with origin-destination  $k$ ,  $x$  the vector of aggregate flows,  $x^*$  the vector of system-optimal flows,  $C_a(x)$  the congestion delay along arc  $a$  and  $\{X_k\}_{k \in \mathcal{K}}$  the sets of feasible flows per commodity, one looks for a

tariff schedule that solves the MPEC

$$\begin{aligned} \min_{T \geq 0} \quad & \sum_{k \in K} \sum_{a \in \mathcal{A}} T_a x_a^{k*} \\ \text{subject to} \quad & \sum_{a \in \mathcal{A}} \left( C_a \left( \sum_{l \in K} x_a^{l*} \right) + T_a \right) (x_a^{k*} - x_a^k) \leq 0 \quad \forall x^k \in X_k \quad \forall k \in K, \end{aligned}$$

where  $X_k = \{Bx^k = b^k, x^k \geq 0\}$ .

## Algorithms

### The price setting problem

The basic bilinear model can be solved as a linear bilevel program using reformulation (1.7). However, to preserve the structure of the problem, it is preferable to work directly on formulation (1.5), which can be solved by a Branch-and-Bound procedure where branching is performed with respect to the complementarity constraint, i.e., the alternative

$$(c_2 - \lambda A_2)_i = 0 \quad \text{or} \quad y_i = 0.$$

Whenever the index set  $P$  of positive variables  $y_j$  ( $j \in P$ ) in an optimal solution is specified at a given node of the enumeration tree, an optimal  $\lambda$ -vector can be recovered by solving the linear program:

$$\begin{aligned} \max_{\lambda} \quad & \lambda b \\ \text{subject to} \quad & (\lambda A_2 - c_2)_j \leq 0 \quad \forall j \notin P \\ & (\lambda A_2 - c_2)_j = 0 \quad \forall j \in P \\ & \lambda A_1 \leq c_1, \end{aligned}$$

whose dual

$$\begin{aligned} \min_{x,y} \quad & c_1 x + c_2 y \\ \text{subject to} \quad & A_1 x + A_2 y = b \\ & x \geq 0 \\ & y_j \text{ unrestricted } \forall j \in P \\ & y_j \geq 0 \quad \forall j \notin P \end{aligned}$$

has a structure similar to that of the original lower level LP. If this “inverse optimization” procedure is carried out at each node of the implicit enumeration tree, then the algorithm may heuristically be halted before termination and yet produce a solution of high quality. Actually,

the dual vector  $\lambda$  derived from an optimal solution of the lower level problem with  $T = 0$  often yields a very good initial solution.

In practice, it is frequently the case that the size of the taxed flow vector  $x$  is much smaller than the size of  $y$ . If this is the case, the knowledge of the vector  $x$  (not only the index set of its positive components, but also their values) allows to derive the values of  $y$  by solving the LP:

$$\begin{aligned} \min_y \quad & c_2 y \\ \text{subject to} \quad & A_2 y = b - A_1 x \\ & y \geq 0. \end{aligned}$$

As before, an optimal tax vector corresponding to the vector  $y$  can easily be recovered.

Another approach, which is better suited to the solution of large scale problems, consists in penalizing the complementarity constraint of (1.5) to obtain the bilinear program

$$\begin{aligned} \min_{x,y,\lambda} \quad & c_1 x + c_2 y - \lambda b + M(c_2 - \lambda A_2)y \\ \text{subject to} \quad & A_1 x + A_2 y = b \\ & x, y \geq 0 \\ & \lambda A_2 \leq c_2 \quad . \end{aligned} \tag{1.9}$$

For solving single-commodity instances of (1.9), variants of the Gauss-Seidel iterative procedure, coupled to a clever update of the penalty parameter  $M$ , performed very well.

### **Toll setting**

In contrast with the basic model, the multicommodity network model has three distinguished features: (i) its network structure, (ii) upper bounds on the tolls and (iii) tolls that apply to total flow, not individual commodity flows. The presence of a network structure and the fact that congestion is not taken into account (the flow of a given origin-destination pair will go through a unique path) allows for a mixed-integer formulation that involves a relatively small number of integer variables, i.e., one per toll arc and per commodity. We express the flow variables as proportions of the demand  $n_k$  associated with the origin-destination couple  $k$  and the nodal demands are set to  $e_i^k = \text{sgn}(b_i^k)$ ; this yields the equivalent formulation:

$$\begin{aligned}
& \max_{T,x,y,\lambda} \sum_{k \in \mathcal{K}} \sum_{a \in \mathcal{A}_1} n^k T_a^k \\
\text{subject to } & \sum_{a \in i^+} (x_a^k + y_a^k) - \sum_{a \in i^-} (x_a^k + y_a^k) = e_i^k \quad \forall i \in \mathcal{N} \quad \forall k \in \mathcal{K} \\
& \lambda_i^k - \lambda_j^k \leq c_{1a} + T_a \quad \forall a = (i, j) \in \mathcal{A}_1 \quad \forall k \in \mathcal{K} \\
& \lambda_i^k - \lambda_j^k \leq c_{2a} \quad \forall a \in \mathcal{A}_2 \quad \forall k \in \mathcal{K} \\
& \sum_{a \in \mathcal{A}_1} (c_{1a} x_a^k + T_a^k) + \sum_{a \in \mathcal{A}_2} c_{2a} y_a^k = \lambda_{o(k)}^k - \lambda_{d(k)}^k \quad \forall k \in \mathcal{K} \\
& -M x_a^k \leq T_a^k \leq M x_a^k \quad \forall k \in \mathcal{K} \quad \forall a \in \mathcal{A}_1 \\
& -M(1 - x_a^k) \leq T_a^k - T_a \leq M(1 - x_a^k) \quad \forall k \in \mathcal{K} \quad \forall a \in \mathcal{A}_1 \\
& x_a^k \in \{0, 1\} \quad \forall k \in \mathcal{K} \quad \forall a \in \mathcal{A}_1 \\
& y_a^k \geq 0 \quad \forall k \in \mathcal{K} \quad \forall a \in \mathcal{A}_2 \\
& T_a \leq T_a^{\max} \quad \forall a \in \mathcal{A}_1.
\end{aligned}$$

The reformulation has been obtained by first replacing the lower level problem of the original Toll Setting Problem by its optimality conditions. Next we have introduced variables

$$T_a^k = T_a x_a^k$$

and exploited the fact that the  $x_a^k$ 's are binary to linearize both the complementarity constraints and the objective function. The fourth constraint of the above program ensures that complementarity slackness of the lower level problem is satisfied, while the next two constraints ensure that the commodity toll  $T_a^k$  is equal to the actual common toll  $T_a$  whenever the commodity flow  $x_a^k$  is positive.

This formulation allows to solve to optimality small instances of the multicommodity toll problem by any integer programming code (see Brotcorne [3] and Brotcorne et al [2]).

For large scale instances we adopt, again, a Gauss-Seidel approach. However, while the treatment of upper bounds in the single-commodity model constitutes a mere technicality, the presence of multiple commodities makes the resolution of (1.9) by a Gauss-Seidel strategy significantly more complex. Specifically, the inverse optimization procedure described earlier must now rely on the solution of a multicommodity flow problem with upper bounds on arc flows. To avoid this difficulty, we propose simply to replace the global toll vector  $T$  by commodity tolls  $T^k$  that must satisfy the compatibility constraint

$$T_k = T_1 \quad \forall k.$$

By penalizing into the objective this constraint and the complementarity constraints (expressed in the primal-dual equality form), one obtains the bilinear-quadratic program

$$\begin{aligned}
 \max_{\lambda, x, y} \quad & \sum_{k \in K} \left[ \lambda^k b^k - (c_1 x^k + c_2 y^k) - M_1((c_1 + T^k)x^k + c_2 y^k - \lambda^k b^k) \right. \\
 & \left. - M_2 \|T^k - T^1\|^2 \right] \\
 \text{subject to} \quad & A_1 x^k + A_2 y^k = b^k \quad \forall k \in K \\
 & x^k, y^k \geq 0 \quad \forall k \in K \\
 & \lambda^k A_1 \leq c_1 + T^k \quad \forall k \in K \\
 & \lambda^k A_2 \leq c_2 \quad \forall k \in K.
 \end{aligned}$$

Once the quadratic penalty term is linearized (à la Frank-Wolfe), the Gauss-Seidel strategy may be applied to the resulting bilinear program. Provided that the penalty parameters  $M_1$  and  $M_2$  are calibrated in a suitable way, the coupling of the Gauss-Seidel and inverse optimization strategies could uncover solutions of large-scale problems whose objective values were typically within one percent of the optimal values, when these were known (see [2]).

### Price setting of telecommunication networks

The algorithmic approaches presented for the multicommodity network model can be extended to the nonlinear telecommunication model. To see this, let us replace the lower level by its optimality conditions:

$$\begin{aligned}
 \max_{T, x, y, \lambda} \quad & \sum_{a \in \mathcal{A}_1} \sum_{k \in \mathcal{K}} T_a^k x_a^k \\
 \text{subject to} \quad & T_a \leq T_a^{\max} \quad \forall a \in \mathcal{A}_1 \\
 & T_a^k + \frac{\alpha_a C_a}{\gamma(C_a - f_a)^2} - \lambda_j^k + \lambda_i^k - \mu_a^k = 0 \quad \forall a \in \mathcal{A}_1 \quad \forall k \in \mathcal{K}
 \end{aligned}$$

$$\begin{aligned}
d_a^k + \frac{\alpha_a C_a}{\gamma(C_a - f_a)^2} - \lambda_j^k + \lambda_i^k - \mu_a^k &= 0 & \forall a \in \mathcal{A}_2 \quad \forall k \in \mathcal{K} \\
\sum_{a \in i^+} (x_a^k + y_a^k) - \sum_{a \in i^-} (x_a^k + y_a^k) &= b_i^k & \forall i \in \mathcal{N} \quad \forall k \in \mathcal{K} \\
f_a &= \sum_{k \in \mathcal{K}} x_a^k & \forall a \in \mathcal{A}_1 \\
f_a &= \sum_{k \in \mathcal{K}} y_a^k & \forall a \in \mathcal{A}_2 \\
\mu_a^k x_a^k = 0, \quad \mu_a^k \geq 0, \quad x_a^k \geq 0 & & \forall k \in \mathcal{K} \quad \forall a \in \mathcal{A}_1 \\
\mu_a^k y_a^k = 0, \quad \mu_a^k \geq 0, \quad y_a^k \geq 0 & & \forall k \in \mathcal{K} \quad \forall a \in \mathcal{A}_2.
\end{aligned}$$

Using the same argument as the one developed for the general model, we obtain, after straightforward albeit tedious calculations:

$$\begin{aligned}
\sum_{a \in \mathcal{A}_1} \sum_{k \in \mathcal{K}} T_a^k x_a^k &= \sum_{i \in \mathcal{N}} \sum_{k \in \mathcal{K}} \lambda_i^k b_i^k - \sum_{a \in \mathcal{A}_1} \frac{\alpha_a C_a f_a}{\gamma(C_a - f_a)^2} \\
&\quad - \sum_{a \in \mathcal{A}_2} \frac{\alpha_a C_a f_a}{\gamma(C_a - f_a)^2} - \sum_{a \in \mathcal{A}_2} \sum_{k \in \mathcal{K}} d_a^k f_a^k.
\end{aligned}$$

Now we linearize both the congestion function (by a piecewise linear curve) and the complementarity constraints (as previously and with the help of the binary variables  $z_a^k$ ). If  $m_{ap}$  corresponds to a breakpoint of the piecewise linear curve,  $w_{ap}$  to the binary variable associated with the  $p$ th segment of the curve, and  $t_{ap}$  the variable associated with the convex combination of the  $p^{\text{th}}$  segment, one obtains the mixed integer linear program

$$\begin{aligned}
\max_{T, x, y, w, z} \quad & \sum_{i \in \mathcal{N}} \sum_{k \in \mathcal{K}} \lambda_i^k b_i^k - \sum_{a \in \mathcal{A}_1} \sum_{p=1}^{n+1} \frac{\alpha_a C_a}{\gamma} t_{ap} h(m_{ap}) \\
& - \sum_{a \in \mathcal{A}_2} \sum_{p=1}^{n+1} \frac{\alpha_a C_a}{\gamma} t_{ap} h(m_{ap}) \\
& - \sum_{a \in \mathcal{A}_2} \sum_{k \in \mathcal{K}} d_a^k f_a^k
\end{aligned}$$



subject to

$$\begin{aligned}
 T_a &\leq T_a^{\max} && \forall a \in \mathcal{A}_1 \\
 \lambda_j^k - \lambda_i^k - \sum_{p=1}^{n+1} \frac{\alpha_a C_a}{\gamma} t_{ap} g(m_{ap}) &\leq d_a^k && \forall k \in \mathcal{K} \quad \forall a \in \mathcal{A}_2 \\
 f_a &= \sum_{p=1}^{n+1} t_{ap} m_{ap} && \forall a \in \mathcal{A} \\
 f_a &= \sum_{k \in \mathcal{K}} x_a^k && \forall a \in \mathcal{A}_1 \\
 f_a &= \sum_{k \in \mathcal{K}} y_a^k && \forall a \in \mathcal{A}_2 \\
 \sum_{p=1}^{n+1} t_{ap} &= 1 && \forall a \in \mathcal{A} \\
 t_{a1} &\leq w_{a1} \quad t_{a(n+1)} \leq w_{an} && \forall a \in \mathcal{A} \\
 t_{ap} &\leq w_{a(p-1)} + w_{ap} && p = 2, \dots, n \quad \forall a \in \mathcal{A} \\
 \sum_{p=1}^n w_{ap} &= 1 && \forall a \in \mathcal{A} \\
 \sum_{a \in i^+} (x_a^k + y_a^k) - \sum_{a \in i^-} (x_a^k + y_a^k) &= e_i^k && \forall i \in \mathcal{N} \quad \forall k \in \mathcal{K} \\
 d_a^k + \sum_{p=1}^{n+1} \frac{\alpha_a C_a}{\gamma} t_{ap} g(m_{ap}) - \lambda_j^k + \lambda_i^k - M z_a^k &\leq 0 && \forall a \in \mathcal{A}_2 \quad \forall k \in \mathcal{K} \\
 f_a^k &\leq M(1 - z_a^k) && \forall a \in \mathcal{A}_2 \quad \forall k \in \mathcal{K} \\
 f_a^k &\geq 0 && \forall a \in \mathcal{A} \quad \forall k \in \mathcal{K} \\
 w_{ap} &\in \{0, 1\} && p \in \{1, 2, \dots, n\} \quad \forall a \in \mathcal{A} \\
 t_{ap} &\geq 0 && p \in \{1, 2, \dots, n+1\} \quad \forall a \in \mathcal{A} \\
 z_a^k &\in \{0, 1\} && \forall a \in \mathcal{A}_2 \quad \forall k \in \mathcal{K}.
 \end{aligned}$$

The above model can only be solved to optimality for small instances. For larger instances, Julsain [13] implemented the following procedure: for fixed price levels, an equilibrium flow solution of the lower level problem traffic assignment problem is obtained; next, the optimal price schedule corresponding to this flow is determined by solving a linear

inverse optimization problem. The process is iterated until no change is observed, at which point the best solution obtained by the inverse optimization procedure is retained.

### Yield management in the airline industry

The yield management problem is by far the most ambitious model presented in this survey, involving the solution of an infinite-dimensional lower level problem. Working in infinite dimension could, surprisingly, prove an asset from the computational point of view, since the continuous distribution of the trade-off values  $\alpha$  and  $\beta$  across the population smoothes out the lower level reaction to the upper level fares and ensures, under weak conditions, that the lower level solution is unique. It is then conceivable to address the bilevel problem as a single-level differentiable optimization problem, whose objective can be computed by tedious but rather straightforward implicit derivation rules. Alternatively, one could discretize the distribution functions and apply the algorithmic ideas presented for the other applications. A mixed continuous-discrete approach based on both ideas is currently investigated.

### Traffic management through link tolls

In contrast with the previous applications, the traffic management problem is not a *bona fide* bilevel program but actually an inverse optimization problem that can be solved easily. Let us first analyze the case of a single origin (single commodity), which can be expressed, using obvious vector and matrix notation, as

$$\begin{aligned} & \min_{T \geq 0} \quad Tx^* \\ \text{subject to} \quad & \langle C(x^*) + T, x^* - x \rangle \leq 0 \quad \forall x \in X = \{Bx = b, x \geq 0\}. \end{aligned}$$

Since the solution of the lower level variational inequality is known to be  $x^*$ , this vector is a solution of the primal-dual system

$$\begin{aligned} Bx^* &= b \\ x^* &\geq 0 \\ \lambda B &\leq C(x^*) + T \\ \langle C(x^*) + T - \lambda B, x^* \rangle &= 0. \end{aligned} \tag{1.10}$$

If  $x_a^*$  is positive, the corresponding toll is set to  $T_a = (\lambda B - C(x^*))_a$ ; otherwise, the toll  $T_a$  can be set to any value that exceeds the maximum equilibrium delay, which is known. To simplify the presentation, we

assume that all components of the vector  $x^*$  are positive. The traffic management problem can then be reduced to the linear program

$$\begin{aligned} \min_{\lambda} \quad & \langle \lambda B - C(x^*), x^* \rangle \\ \text{subject to} \quad & \lambda B - C(x^*) \geq 0, \end{aligned}$$

whose dual

$$\max_{x \in X} C(x^*)x$$

consists in finding a longest path tree in an acyclic network, and is solvable by a greedy algorithm in polynomial time (Dial [5]).

In the multicommodity case, let  $\bar{x}_a$  denote the total flow on arc  $a$ . The equivalent of (1.10) is

$$\begin{aligned} Bx^{k*} &= b^k & \forall k \in \mathcal{K} \\ x^{k*} &\geq 0 & \forall k \in \mathcal{K} \\ \lambda^k B &\leq C(\bar{x}^*) + T & \forall k \in \mathcal{K} \\ \sum_{k \in \mathcal{K}} \langle C(\bar{x}^*) + T - \lambda^k B, x^{k*} \rangle &= 0. \end{aligned}$$

The dual of the problem of minimizing  $T\bar{x}^*$  subject to the last two constraints is

$$\begin{aligned} \max_{x, \alpha} \quad & \langle C(\bar{x}^*), \alpha \bar{x}^* - \bar{z} \rangle \\ \text{subject to} \quad & \bar{z} = \sum_{k \in \mathcal{K}} z^k \\ & \bar{z} \leq (1 + \alpha) \bar{x}^* \\ & Bz^k = \alpha b^k, z^k \geq 0 \quad \forall k \in \mathcal{K} \end{aligned}$$

which, after performing the scaling  $z = \alpha x'$ , becomes

$$\begin{aligned} \max_{\alpha} \quad \phi(\alpha) \quad & := \quad \alpha \max_{x' \in \prod_k X_k} \langle C(\bar{x}^*), \bar{x}^* - \bar{x}' \rangle \\ \text{subject to} \quad & \bar{x}' \leq (1 + 1/\alpha) \bar{x}^*. \end{aligned} \quad (1.11)$$

The function  $\phi$  is concave, increasing, piecewise linear, and its maximum is achieved for any sufficiently large value of the variable  $\alpha$ . For such values of  $\alpha$ , an optimal tax vector will be recovered from the dual vector  $\lambda$  associated with the capacity constraint (1.11) of a linear multicommodity flow problem. A slightly different multicommodity formulation involving a nonnegativity constraint on the sum of commodity flows was derived by Dial [6].

While the multicommodity problem is considerably more difficult to solve than finding a longest path tree in an acyclic network, it is nonetheless manageable, even for large scale networks. Actually, by penalizing the upper bound constraints, an approximate solution can be computed by solving a convex flow problem which is at least as easy to solve as the traffic assignment problem which had to be solved in order to determine the system-optimal flow vector  $x^*$  in the first place.

If only a subset of arcs is subject to tolls, then it might be impossible to induce a system-optimal flow pattern. If one adopts as maximization criterion the social surplus, then the problem becomes an authentic bilevel program that can be addressed by the techniques presented earlier.

## Conclusion

In this paper, we have presented a short survey of pricing situations modeled as bilevel programs, as well as several avenues for their numerical resolution. We firmly believe that this approach will gain in popularity both in the economic and mathematical programming communities, and that the day is not far where these models will be routinely solved for near-optimal solutions.

## References

- [1] Bencheick, A., “Conception et tarification d’un réseau privé multidébit à large bande dans un environnement concurrentiel”, Ph.D. thesis, Institut National de Recherche Scientifique, INRS-Télécommunication (1997).
- [2] Brotcorne, L., Labbé, M., Marcotte, P. and Savard, G., “A Bilevel Model for Toll Optimization on a Multicommodity Transportation Network”, forthcoming in *Transportation Science*.
- [3] Brotcorne, L., “Approches opérationnelles et stratégiques des problèmes de trafic routier”, Ph.D. thesis, Université Libre de Bruxelles, February 1998.
- [4] Cocchi, R., Shenker, S., Estrin, D. and Zhang, L., “Pricing in computer networks: motivation, formulation, and example”, *IEEE/ACM Transactions on Networking* **1** (1993) 614–627.
- [5] Dial, R. B., “Minimal-revenue congestion pricing Part II: An efficient algorithm for the general case”, *Transportation Research B* **34** (2000) 645–665.

- [6] Dial, R. B., “Minimal-revenue congestion pricing part I: A fast algorithm for the single-origin case”, *Transportation Research B* **33** (1999) 189–202.
- [7] Fukushima, M., “Equivalent differentiable optimization problems and descent methods for asymmetric variational inequality problems”, *Mathematical Programming* **53** (1992) 99–110.
- [8] Gendreau, M., Marcotte, P. and Savard, G., “A Hybrid tabu-ascent algorithm for the linear bilevel programming problem”, *Journal of Global Optimization* **8** (1996) 217–232.
- [9] Girard, A. and Zidane, R., “Revenue Optimization of B-ISDN Networks”, *IEEE Transactions on Communications* **43** (1995) 15–25.
- [10] Hansen, P., Jaumard, B. and Savard, G., “New branch-and-bound rules for linear bilevel programming”, *JOTA* **22** (1992) 1194–1217.
- [11] Hearn, D.W. and Ramana, M.V., “Solving congestion toll pricing models”, in: *Equilibrium and Advanced Transportation Modelling*, Patrice Marcotte and Sang Nguyen (eds.), Kluwer, Dordrecht, pp. 109-123, 1998.
- [12] Kočvara, M. and Outrata, J.V., “A nonsmooth approach to optimization problems with equilibrium constraints”, in *Complementarity and variational problems. State of the art*, M.C. Ferris and J.S. Pang (eds.), SIAM, Philadelphia (1997).
- [13] Julsain, H., “Tarification dans les réseaux de télécommunications: une approche par programmation mathématique à deux niveaux”, Mémoire de maîtrise, École Polytechnique de Montréal, 1999.
- [14] Labbé, M., Marcotte, P. and Savard, G., “A bilevel model of taxation and its application to optimal highway pricing”, *Management Science* **44** (1998) 1595–1607 .
- [15] Larsson, T. and Patriksson, M., “Traffic management through link tolls - an approach utilizing side constrained traffic equilibrium models”, in: *Equilibrium and Advanced Transportation Modelling*, Patrice Marcotte and Sang Nguyen (eds.), Kluwer, Dordrecht, pp. 125-151, 1998.
- [16] Loridan, P. and Morgan, J., “ $\epsilon$ -regularized two-level optimization problems: approximation and existence results”, *Fifth French-German Conference on Optimization*, Lecture Notes in Mathematics 1405, Springer-Verlag (1989) 99–113.
- [17] Luo, Z.Q., Pang, J.S. and Ralph, D., *Mathematical Programs with Equilibrium Constraints*, Cambridge University Press, Cambridge (1996).

- [18] Mitchell, B.M. and Vogelsang, I., *Telecommunications pricing: Theory and practice*, Cambridge University Press, (1992).
- [19] Outrata, J.V., Kocvara, M. and Zowe, J., *Nonsmooth Approach to Optimization Problems with Equilibrium Constraints*. Kluwer Academic Publishers, Dordrecht, (1998).
- [20] Shimizu, K., Ishizuka, Y. and Bard, J.F., *Nondifferentiable and Two-Level Mathematical Programming*, Kluwer, Boston (1997).
- [21] Scholtes, S. and Stohr, M., “Exact penalization of mathematical programs with equilibrium constraints” *SIAM J. Control and Optimization* **37**, pp.617-652, (1999).
- [22] Special issue on Yield Management, *Transportation Science* **33** (1999).
- [23] Wang, Q., Peha, J.M. and Sirbu, M., “Optimal Pricing for Integrated-Services Networks”, in *Internet Economics*, Joseph Bailey and Lee McKnight editors, MIT Press, (1997) 353–376.