



GROUPE D'ÉTUDES ET DE RECHERCHE
EN ANALYSE DES DÉCISIONS

Les Cahiers du GERAD

CITATION ORIGINALE / ORIGINAL CITATION

GERAD HEC Montréal
3000, ch. de la Côte-Sainte-Catherine
Montréal (Québec) Canada H3T 2A7

Tél. : 514 340-6053
Télec. : 514 340-5665
info@gerad.ca
www.gerad.ca

**Improved Primal Simplex:
A More General Theoretical Framework
and an Extended Experimental Analysis**

J. Omer, S. Rosat,
V. Raymond, F. Soumis

G-2014-13

March 2014

Improved Primal Simplex: A More General Theoretical Framework and an Extended Experimental Analysis

Jérémy Omer
Samuel Rosat
Vincent Raymond
François Soumis

*GERAD & Département de mathématiques et de génie industriel
École Polytechnique
Montréal (Québec) Canada, H3C 3A7*

jeremy.omer@gerad.ca
samuel.rosat@gerad.ca
vincent.raymond@polymtl.ca
francois.soumis@gerad.ca

March 2014

Les Cahiers du GERAD
G-2014-13

Copyright © 2014 GERAD

Abstract: In this article, we propose a general framework for an algorithm derived from the primal simplex that guarantees a strict improvement in the objective after each iteration. Our approach relies on the identification of *compatible* variables that ensure a nondegenerate iteration if pivoted into the basis. The problem of finding a strict improvement in the objective function is proved to be equivalent to two smaller problems respectively focusing on compatible and incompatible variables. We then show that the improved primal simplex (IPS) of Elhallaoui et al. is a particular implementation of this generic theoretical framework. The resulting new description of IPS naturally emphasizes what should be considered as necessary adaptations of the framework versus specific implementation choices. This provides original insight into IPS that allows for the identification of weaknesses and potential alternative choices that would extend the efficiency of the method to a wider set of problems. We perform experimental tests on an extended collection of data sets including instances of Mittelman's benchmark for linear programming. The results confirm the excellent potential of IPS and highlight some of its limits while showing a path toward an improved implementation of the generic algorithm.

Key Words: Linear programming, simplex, degeneracy, decomposition, primal algorithms.

1 Introduction

1.1 Degeneracy in the Primal Simplex

The primal simplex, described by Dantzig in 1947 (see Dantzig (1955)), was the first efficient algorithm to be developed for solving linear programs (LPs) and is still used for a large number of applications. It starts with a feasible solution obtained from a *basis* of the variable space by setting all the nonbasic variables at one of their bounds. It then iteratively improves the objective value until optimality is reached. At each iteration, the algorithm moves to an adjacent, improving basis by selecting a nonbasic entering variable and removing a variable from the basis through a simplex *pivot*. One of the major difficulties the algorithm may encounter is called *degeneracy*. It occurs when many the basic variables are at one of their bounds. In this case, there is a high probability that the variable selected to enter the basis cannot have its value modified without making the current solution infeasible. The resulting degenerate pivot leads to no change in the current solution and no improvement in the objective value. In a study of degeneracy in the simplex method, Perold (1980) states that typically if on average 20% of the basic variables are at one of their bounds, 50% of the iterations will be degenerate. Combined with the multiplicity of highly degenerate problems arising in industrial applications, this observation highlights the need for an efficient treatment of degeneracy in any good implementation of the primal simplex algorithm.

1.2 Dealing with Primal Degeneracy: A Short State of the Art

Several techniques have been developed to cope with degeneracy. One family of these methods focuses on the selection of the variable entering the basis. For instance, in Greenberg (1978), a vector is computed such that each of its positive entries corresponds to a nonbasic variable that will necessarily lead to a degenerate pivot if it is chosen to enter the basis. The number of operations is similar to the computation of a reduced cost, but the test is only heuristic since it does not identify all the nonbasic variables that would lead to a degenerate pivot. The perturbation (Benichou et al. 1977) and bound-shifting (Gill et al. 1989) techniques follow a different path. They apply a random modification of either the values of the degenerate variables or their bounds to put an end to a sequence of nonimproving pivots. The random modification avoids performing iterations without improvement, but it usually replaces them with small steps. Although not supported by a strong theory, the two methods, that mostly differ from an implementation point of view, perform well, and they appear in most efficient simplex codes. For further information, an excellent study of degeneracy emphasizing its computational aspects is given by Maros (2003).

Another approach aims to take advantage of degeneracy rather than just to minimize its negative effects. Geometrically, a degenerate vertex of the n -dimensional polyhedron described by the LP's constraints is the crossing point of more than $n - 1$ facets of the polyhedron. Degeneracy thus corresponds to a local excess of information, which suggests that a smaller problem may be considered locally to make progress from a degenerate solution. Perold (1980) introduces a particular degeneracy structure in the LU decomposition of the basis, which involves fewer calculations when performing degenerate pivots. Pan (2008) generalizes the definition of basis to include *deficient bases* containing p independent columns, with p lower than the number of rows. When degeneracy occurs, a deficient basis, smaller than the usual square basis matrix, may be used to perform the pivot calculations. Generalizing the dynamic constraint aggregation described by Elhallaoui et al. (2005) for set partitioning problems, the improved primal simplex (IPS) of Elhallaoui et al. (2011) makes the most of degeneracy by simultaneously reducing the number of rows and columns to perform as many cheap pivots as possible. Once all the interesting pivots are done, a *complementary* LP is solved to find a dual solution maximizing the minimum reduced cost. Either this reduced cost is null and optimality is proved, or it corresponds to a combination of primal variables that may enter the basis through a sequence of pivots ending with a strict improvement in the objective. Metrane et al. (2010) prove the equivalence of IPS to a particular column generation algorithm by respectively identifying the reduced and complementary problems of IPS with the master and pricing problems of column generation.

1.3 Contribution Statement

The present article is motivated by the promising experimental results for IPS that are reported by Elhallaoui et al. (2011) and Raymond et al. (2010b). However, Elhallaoui et al. (2011) conduct their tests on instances containing a majority of set-partitioning constraints, while the improvements proposed in Raymond et al. (2010b) are designed for only a part of this restricted benchmark. Based on this proof of concept, our intent is to investigate the performance of the algorithm on a much more diversified benchmark and, using a new, more general theoretical description of the algorithm, to highlight potential improvements.

With that in mind, our first contribution is to describe IPS for an LP with upper and lower bounds on the variables. In Elhallaoui et al. (2011), the algorithm is described for an LP in standard form, because every LP has an equivalent standard form, but this is not compatible with an efficient implementation. Including the bounds explicitly shows how degeneracy can be efficiently handled when a variable may be at either its lower or its upper bound, and it specifies the special treatment that unbounded variables should receive. The presentation of the algorithm is then generalized to a theoretical framework of which IPS is only one possible implementation. Our goal is to emphasize the steps in IPS that may be seen as implementation choices. By analyzing these choices independently of the core of the method, we can identify their main advantages and drawbacks and consider alternative implementations. Moreover, the new framework is purely primal, which simplifies the presentation and the proofs. We then extend the benchmark by including larger instances and problems that do not share the specific structure of the instances used by Elhallaoui et al. (2011) and Raymond et al. (2010b). A thorough analysis of the results provides a much more accurate description of the families of problems that may be solved efficiently using IPS while highlighting the implementation choices that would enlarge its domain of efficiency. The main contribution is thus to identify the limits of IPS while laying the theoretical foundation for future improvements.

The generic theoretical framework is developed in Section 2, and the implementation choices for IPS are described in Section 3. The results of the experimental tests are analyzed in Section 4, and in Section 5 we discuss directions for future research.

1.4 Notation

Lower-case bold symbols are used for column vectors and upper-case bold symbols denote matrices. For subsets $\mathcal{I} \subseteq \{1, \dots, m\}$ of row indices and $\mathcal{J} \subseteq \{1, \dots, n\}$ of column indices, the submatrix of \mathbf{A} with rows indexed by \mathcal{I} and columns indexed by \mathcal{J} is denoted $\mathbf{A}_{\mathcal{I}\mathcal{J}}$. Similarly, $\mathbf{A}_{\mathcal{I}}$ is the set of rows of \mathbf{A} indexed by \mathcal{I} , $\mathbf{A}_{\mathcal{J}}$ is the set of columns of \mathbf{A} indexed by \mathcal{J} , and $\mathbf{x}_{\mathcal{J}}$ is the subvector of variables corresponding to $\mathbf{A}_{\mathcal{J}}$. The vector of all zeros with dimension dictated by the context is denoted $\mathbf{0}$, and \mathbf{A}^T is the transpose of \mathbf{A} .

2 A Generic Decomposition Algorithm Taking Advantage of Degeneracy

IPS was designed to start from an extreme solution of the feasible domain, assuming nonnegative variables. We adapt it to any allow feasible starting solution, and we consider an LP with bounded variables:

$$\begin{cases} \min & \mathbf{c}^T \mathbf{x} \\ \text{s. t.} & \mathbf{A}\mathbf{x} = \mathbf{b} \\ & \mathbf{l} \leq \mathbf{x} \leq \mathbf{u}, \end{cases} \quad (\text{P})$$

where $\mathbf{x} \in \mathbb{R}^n$ is the vector of decision variables, $\mathbf{c} \in \mathbb{R}^n$ is the cost vector, $\mathbf{A} \in \mathbb{R}^{m \times n}$ is the constraint matrix, $\mathbf{b} \in \mathbb{R}^m$ is the right-hand side vector, and $\mathbf{l}, \mathbf{u} \in \mathbb{R} \cup \{\pm\infty\}$, $\mathbf{l} \leq \mathbf{u}$, are respectively the lower- and upper-bound vectors. We assume that \mathbf{A} is of full rank m with $m \leq n$ and that the feasible domain \mathcal{F}_P is nonempty.

Given a feasible solution $\mathbf{x} \in \mathcal{F}_P$, the indices of its components can be partitioned into four sets, \mathcal{P} , \mathcal{L} , \mathcal{U} , and \mathcal{M} . \mathcal{L} and \mathcal{U} respectively designate the set of indices of variables that are at their lower and

upper bounds. For the remaining variables—all of them being strictly within their bounds— \mathcal{P} designates a maximal subset such that $\mathbf{A}_{\mathcal{P}}$ is linearly independent, while \mathcal{M} gathers all the others. Note that $|\mathcal{P}| \leq m$, but nothing can be inferred for $|\mathcal{M}|$. For the sake of readability, we drop the index on \mathbf{x} for \mathcal{P} , \mathcal{L} , \mathcal{U} , and \mathcal{M} and the other sets introduced later (\mathcal{C} , \mathcal{I} , etc.). With this decomposition and a convenient reordering of the columns, $\mathbf{x} = (\mathbf{x}_{\mathcal{P}}, \mathbf{x}_{\mathcal{L}}, \mathbf{x}_{\mathcal{U}}, \mathbf{x}_{\mathcal{M}})$ with $\mathbf{l}_{\mathcal{P}} < \mathbf{x}_{\mathcal{P}} < \mathbf{u}_{\mathcal{P}}$, $\mathbf{x}_{\mathcal{L}} = \mathbf{l}_{\mathcal{L}}$, $\mathbf{x}_{\mathcal{U}} = \mathbf{u}_{\mathcal{U}}$, and $\mathbf{l}_{\mathcal{M}} < \mathbf{x}_{\mathcal{M}} < \mathbf{u}_{\mathcal{M}}$, and

$$\mathbf{A}_{\mathcal{P}}\mathbf{x}_{\mathcal{P}} + \mathbf{A}_{\mathcal{L}}\mathbf{l}_{\mathcal{L}} + \mathbf{A}_{\mathcal{U}}\mathbf{u}_{\mathcal{U}} + \mathbf{A}_{\mathcal{M}}\mathbf{x}_{\mathcal{M}} = \mathbf{b}. \quad (1)$$

$\mathbf{A}_{\mathcal{P}}$ is a basis of the linear span of $\{\mathbf{A}_{\mathcal{P}}, \mathbf{A}_{\mathcal{M}}\}$, denoted $\text{Span}(\mathbf{A}_{\mathcal{P}}, \mathbf{A}_{\mathcal{M}})$, and \mathcal{P} is therefore called the *reduced* or *working basis*. Its indices and the variables of $\mathbf{x}_{\mathcal{P}}$ are respectively referred to as *basic indices* and *basic variables*, and $p = |\mathcal{P}|$ denotes the cardinality of this working basis. Finally, the set of all nonbasic indices is denoted \mathcal{N} : $\mathcal{N} = \mathcal{L} \cup \mathcal{U} \cup \mathcal{M}$.

2.1 A Primal Algorithm for Degenerate Problems

Like the primal simplex or gradient-descent methods, IPS is a primal algorithm: it starts from an initial feasible solution and iteratively improves it until optimality is reached. Primal algorithms mainly differ from each other in the improvement procedure, i.e., the way of finding a better solution. Some of them guarantee a strict improvement in the objective function, whereas others do not. Typically, Dantzig's primal simplex cannot guarantee this strict improvement when the instances are degenerate. As for IPS, it was designed to avoid nonimproving steps, although in practice degenerate pivots are still performed. The purpose of this section is to describe a general decomposition algorithm yielding a strict improvement at each step; IPS may be seen as one of its possible implementations. Moreover, we give a theoretical insight based on a purely primal approach. This approach both gives a rather geometric understanding of the algorithm and eases the proof of optimality.

In the rest of this paper, $\mathbf{x}^0 \in \mathcal{F}_{\mathcal{P}}$ denotes the current feasible solution of (P). To guarantee a strict improvement, the following problem must be addressed:

Improvement Problem (IMP): Given a feasible solution \mathbf{x}^0 , either supply a strictly improving vector $\mathbf{y} \in \mathbb{R}^n$, i.e., a vector such that $\mathbf{x}^0 + \mathbf{y} \in \mathcal{F}_{\mathcal{P}}$ and $\mathbf{c}^T \mathbf{y} < 0$, or assert that \mathbf{x}^0 is optimal for (P).

Note that in the case of a maximization problem, $\mathbf{c}^T \mathbf{y} < 0$ becomes $\mathbf{c}^T \mathbf{y} > 0$ since the aim is to increase the objective function. Although *Augmentation* is often found in the literature on integer programming, *Improvement* is more general than *Augmentation* or *Diminution* since it applies to both maximization and minimization problems. Given the current solution \mathbf{x}^0 , IMP comes down to finding an *improving feasible* direction as introduced below.

Definition 1 Vector $\mathbf{d} \in \mathbb{R}^n$ is a feasible direction at \mathbf{x}^0 if there exists a step $\rho > 0$ such that $\mathbf{x}^0 + \rho \cdot \mathbf{d} \in \mathcal{F}_{\mathcal{P}}$. Moreover, \mathbf{d} is an improving direction if $\mathbf{c}^T \mathbf{d} < 0$, i.e., if taking a positive step along \mathbf{d} yields an improvement in the objective value.

An improving vector \mathbf{y} is therefore characterized by a couple $(\mathbf{d}, \rho) \in \mathbb{R}^n \times \mathbb{R}_+$ where \mathbf{d} can be normalized at will as long as $\mathbf{x}^0 + \mathbf{y} = \mathbf{x}^0 + \rho \cdot \mathbf{d} \in \mathcal{F}_{\mathcal{P}}$. The set of all feasible directions at \mathbf{x}^0 can be described as the following cone section:

$$\Delta = \left\{ \mathbf{d} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{d} = 0, \mathbf{d}_{\mathcal{U}} \leq 0, \mathbf{d}_{\mathcal{L}} \geq 0, \sum_{i \in \mathcal{N}} w_i |d_i| = 1 \right\} \quad (2)$$

where $\mathbf{w} \in \mathbb{R}^n$, $\mathbf{w}_{\mathcal{P}} \geq 0$, $\mathbf{w}_{\mathcal{N}} > 0$. $\sum_{i \in \mathcal{N}} \mathbf{w}_i |d_i| = 1$ geometrically represents a normalization constraint. Without this constraint, Δ is a cone and if \mathbf{x}^0 is an extreme point of $\mathcal{F}_{\mathcal{P}}$, the extreme directions of this cone are the directions of the edges of $\mathcal{F}_{\mathcal{P}}$ at \mathbf{x}^0 .

The normalization constraint can be written in a way that makes it closer to a linear equality, which will prove useful later. With $\mathbf{w}_{\mathcal{L}} > 0$, $\mathbf{w}_{\mathcal{U}} < 0$ (and $\mathbf{w}_{\mathcal{P}} \geq 0$, $\mathbf{w}_{\mathcal{M}} > 0$ as previously), Equation (2) reads:

$$\Delta = \left\{ \mathbf{d} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{d} = 0, \mathbf{d}_{\mathcal{U}} \leq 0, \mathbf{d}_{\mathcal{L}} \geq 0, \sum_{i \in \mathcal{P} \cup \mathcal{M}} w_i |d_i| + \mathbf{w}_{\mathcal{U}}^T \mathbf{d}_{\mathcal{U}} + \mathbf{w}_{\mathcal{L}}^T \mathbf{d}_{\mathcal{L}} = 1 \right\}. \quad (3)$$

Moreover, given a feasible improving direction \mathbf{d} , the *maximal step* along \mathbf{d} at \mathbf{x}^0 , defined as $r(\mathbf{d}) = \max \{ \rho \mid \mathbf{x}^0 + \rho \mathbf{d} \in \mathcal{F}_{\mathcal{P}} \}$, is the most interesting step to take since it yields the greatest possible improvement along \mathbf{d} .

Remark 1 *In Dantzig's simplex algorithm, every pivot corresponds to taking a step in some extreme improving direction. However, what the simplex cannot guarantee is the feasibility of this direction. On the one hand, a degenerate pivot corresponds to a nonfeasible direction \mathbf{d} and, therefore, a step $r(\mathbf{d}) = 0$. On the other hand, nondegenerate simplex pivots exactly correspond to taking maximal steps in directions that have only one nonbasic nonzero entry.*

One way to tackle the improvement problem is to minimize the chosen objective over the domain Δ . Typically, that objective is related to the improvement in the objective function given by $\mathbf{d} \in \Delta$. IMP may thus come down to solving the following *greatest normalized improvement* program:

$$z_{\text{GNI}}^* = \min_{\mathbf{d} \in \mathbb{R}^n} \{ \mathbf{c}^T \mathbf{d} \mid \mathbf{d} \in \Delta \}. \quad (\text{GNI})$$

Proposition 1 \mathbf{x}^0 is optimal for (P) if and only if $z_{\text{GNI}}^* \geq 0$.

Proof. By construction of GNI, there exists a feasible improving direction if and only if $z_{\text{GNI}}^* \geq 0$. Therefore, the proposition holds. \square

Based on Proposition 1, the iterative procedure of Algorithm 1 is guaranteed to reach optimality. Moreover, each step $r(\mathbf{d}^k) \cdot \mathbf{d}^k$ taken during this procedure yields a strict improvement in the objective function since all the improving directions \mathbf{d}^k are feasible, thus $r(\mathbf{d}^k) > 0$ and each improvement is $z^{k+1} - z^k = r(\mathbf{d}^k) \cdot \mathbf{c}^T \mathbf{d}^k < 0$. Degeneracy is therefore avoided.

Algorithm 1 Nondecomposed Primal Improvement Procedure

1. Find an initial solution \mathbf{x}^0 ; $k \leftarrow 0$;
 2. If $z_{\text{GNI}}^* \geq 0$, return \mathbf{x}^k ;
 3. If $z_{\text{GNI}}^* < 0$: let \mathbf{d}^k be an optimal solution of GNI; $\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + r(\mathbf{d}^k) \cdot \mathbf{d}^k$; $k \leftarrow k + 1$; GOTO 1.
-

2.2 Of Compatibility and Decomposition: Turning Degeneracy into an Asset

Algorithm 1 not only avoids degeneracy but turns it into an asset by decomposing GNI into two smaller problems that, even if they are solved separately, assert that \mathbf{x}^0 is optimal or provide a feasible improving direction. This decomposition is the main topic of the present subsection. The main idea is to sequentially find nondegenerate improving simplex pivots if they exist, and otherwise to solve a more complicated problem similar to GNI. This decomposition relies on the notion of *compatibility*.

Definition 2 *Given a working basis \mathcal{P} , a vector $\mathbf{v} \in \mathbb{R}^m$ is compatible (with \mathcal{P}) if and only if $\mathbf{v} \in \text{Span}(\mathbf{A}_{\cdot \mathcal{P}})$.*

This notion is extended to the nonbasic columns of $\mathbf{A}_{\cdot \mathcal{N}}$, their indices, and the corresponding variables. The set of indices of the compatible columns of $\mathbf{A}_{\cdot \mathcal{N}}$ is denoted \mathcal{C} . Every other column of $\mathcal{I} = \mathcal{N} \setminus \mathcal{C}$ is said to be incompatible, and $(\mathcal{C}, \mathcal{I})$ forms a partition of \mathcal{N} .

Proposition 2 $\mathcal{M} \subset \mathcal{C}$.

Proof. By definition of \mathcal{P} , for every index $j \in \mathcal{M}$, the corresponding column $\mathbf{A}_{.j}$ is in $\text{Span}(\mathbf{A}_{.\mathcal{P}})$. \square

According to Proposition 2, $\mathcal{C} = \mathcal{C}_{\mathcal{L}} \cup \mathcal{C}_{\mathcal{U}} \cup \mathcal{M}$ and $\mathcal{I} = \mathcal{I}_{\mathcal{L}} \cup \mathcal{I}_{\mathcal{U}}$, where for any pair of sets $(\mathcal{X}, \mathcal{Y})$, $\mathcal{X}_{\mathcal{Y}} = \mathcal{X} \cap \mathcal{Y}$. All incompatible variables are thus at their lower or upper bounds. The following proposition indicates the value of compatible columns when solving IMP.

Proposition 3 *Compatible columns are exactly those that yield nondegenerate pivots if inserted into the working basis \mathcal{P} .*

Proof. Let $j \in \mathcal{N}$. Pivoting $\mathbf{A}_{.j}$ into the working basis leads to a solution $\mathbf{x}^1 \in \mathcal{F}_{\mathcal{P}}$ (where \mathbf{x}^1 may or may not be different from \mathbf{x}^0). Since only $\mathbf{A}_{.j}$ has been pivoted, for all $k \in \mathcal{N} \setminus \{j\}$, $\mathbf{x}_k^1 = \mathbf{x}_k^0$. The constraints of (P) written for both \mathbf{x}^0 and \mathbf{x}^1 are $\mathbf{A}_{.\mathcal{P}}\mathbf{x}_{\mathcal{P}}^0 + \mathbf{A}_{.j}\mathbf{x}_j^0 = \mathbf{A}_{.\mathcal{P}}\mathbf{x}_{\mathcal{P}}^1 + \mathbf{A}_{.j}\mathbf{x}_j^1$, so if $\mathbf{d} = \mathbf{x}^1 - \mathbf{x}^0$

$$\mathbf{A}_{.j}\mathbf{d}_j = \mathbf{A}_{.\mathcal{P}}\mathbf{d}_{\mathcal{P}}. \quad (4)$$

Since $\mathbf{A}_{.\mathcal{P}}$ is of full rank, $\mathbf{d}_{\mathcal{P}} = \mathbf{0}$ if and only if $\mathbf{d}_j = \mathbf{0}$. Therefore, $\mathbf{A}_{.j}$ yields a nondegenerate pivot if and only if $\mathbf{x}^1 \neq \mathbf{x}^0$, that is, if and only if $\mathbf{d}_j \neq \mathbf{0}$. Finally, $\mathbf{d}_j \neq \mathbf{0}$ and Equation (4) are simultaneously satisfied if and only if $\mathbf{A}_{.j} \in \text{Span}(\mathbf{A}_{.\mathcal{P}})$, so the result holds. \square

Note that bounds are not mentioned in the above proof. Since all the variables in \mathcal{P} are strictly within their bounds, they can increase or decrease without violating their bounds. Therefore, if $|d_j|$ is sufficiently small, for all $i \in \mathcal{P}$, $l_i \leq x_i^1 \leq u_i$.

From the perspective of the simplex algorithm, it will be efficient to give priority to the compatible columns since they are exactly those that yield nondegenerate pivots when inserted into the working basis \mathcal{P} . Based on the compatible-incompatible partition of the variables, GNI can be decomposed into the *Reduced*-GNI problem (R-GNI) and the *Complementary*-GNI problem (C-GNI) that respectively focus on the compatible and incompatible variables:

$$\begin{aligned} z_{\text{R-GNI}}^* &= \min_{(\mathbf{d}_{\mathcal{P}}, \mathbf{d}_{\mathcal{C}})} \mathbf{c}_{\mathcal{P}}^T \mathbf{d}_{\mathcal{P}} &+ \mathbf{c}_{\mathcal{C}}^T \mathbf{d}_{\mathcal{C}} \\ \text{s.t.} & \mathbf{A}_{.\mathcal{P}} \mathbf{d}_{\mathcal{P}} &+ \mathbf{A}_{.\mathcal{C}} \mathbf{d}_{\mathcal{C}} &= 0 \\ & \sum_{i \in \mathcal{P}} w_i |d_i| &+ \left(\sum_{i \in \mathcal{M}} w_i |d_i| + \mathbf{w}_{\mathcal{C}_{\mathcal{L}}}^T \mathbf{d}_{\mathcal{C}_{\mathcal{L}}} + \mathbf{w}_{\mathcal{C}_{\mathcal{U}}}^T \mathbf{d}_{\mathcal{C}_{\mathcal{U}}} \right) &= 1 \\ & \mathbf{d}_{\mathcal{C}_{\mathcal{L}}} \geq 0 &, \mathbf{d}_{\mathcal{C}_{\mathcal{U}}} \leq 0 \end{aligned} \quad (\text{R-GNI})$$

$$\begin{aligned} z_{\text{C-GNI}}^* &= \min_{(\mathbf{d}_{\mathcal{P}}, \mathbf{d}_{\mathcal{I}})} \mathbf{c}_{\mathcal{P}}^T \mathbf{d}_{\mathcal{P}} &+ \mathbf{c}_{\mathcal{I}}^T \mathbf{d}_{\mathcal{I}} \\ \text{s.t.} & \mathbf{A}_{.\mathcal{P}} \mathbf{d}_{\mathcal{P}} &+ \mathbf{A}_{.\mathcal{I}} \mathbf{d}_{\mathcal{I}} &= 0 \\ & \mathbf{w}_{\mathcal{P}}^T \mathbf{d}_{\mathcal{P}} &+ \mathbf{w}_{\mathcal{I}}^T \mathbf{d}_{\mathcal{I}} &= 1 \\ & \mathbf{d}_{\mathcal{I}_{\mathcal{L}}} \geq 0 &, \mathbf{d}_{\mathcal{I}_{\mathcal{U}}} \leq 0 \end{aligned} \quad (\text{C-GNI})$$

R-GNI and C-GNI are both smaller than GNI, so solving them separately must be much faster than solving GNI. The only remaining issue is to ensure that the search space for an improving direction Δ can be split into that of the compatible columns ($\mathcal{F}_{\text{R-GNI}}$) and that of the incompatible columns ($\mathcal{F}_{\text{C-GNI}}$), i.e., that solving GNI is equivalent to solving R-GNI and C-GNI.

Lemma 1 *A column $\mathbf{A}_{.j}$ is compatible if and only if there exists $\mathbf{v} \in \mathbb{R}^p$ such that $\mathbf{A}_{.j} = \mathbf{A}_{.\mathcal{P}}\mathbf{v}$. In this case, \mathbf{v} is unique.*

Proof. This lemma is a consequence of the definition of compatibility and of the fact that $\mathbf{A}_{.\mathcal{P}}$ is a full-rank matrix. \square

Theorem 1 *Assume $\mathbf{w}_{\mathcal{P}} = \mathbf{0}$, then $z_{\text{GNI}}^* = \min \{z_{\text{R-GNI}}^*, z_{\text{C-GNI}}^*\}$.*

Proof. Recall that $\mathbf{w}_{\mathcal{L}} > 0$, $\mathbf{w}_{\mathcal{U}} < 0$, $\mathcal{C} = \mathcal{C}_{\mathcal{L}} \cup \mathcal{C}_{\mathcal{U}} \cup \mathcal{M}$, and $\mathcal{I} = \mathcal{I}_{\mathcal{L}} \cup \mathcal{I}_{\mathcal{U}}$. Let $\mathbf{d} = (\mathbf{d}_{\mathcal{P}}, \mathbf{d}_{\mathcal{C}}, \mathbf{d}_{\mathcal{I}})$ be a solution of GNI. If $\mathbf{d}_{\mathcal{C}} = \mathbf{0}$ or $\mathbf{d}_{\mathcal{I}} = \mathbf{0}$, \mathbf{d} is respectively a solution of C-GNI or R-GNI and the result clearly holds.

Suppose now that $\mathbf{d}_{\mathcal{C}} \neq \mathbf{0}$ and $\mathbf{d}_{\mathcal{I}} \neq \mathbf{0}$. We first prove that \mathbf{d} can be written as a convex combination of \mathbf{u}' , the solution of R-GNI, and \mathbf{v}' , the solution of C-GNI. For every $j \in \mathcal{C}$, $\mathbf{A}_{\cdot j} \in \text{Span}(\mathbf{A}_{\cdot \mathcal{P}})$, so linear combinations of compatible columns $\mathbf{A}_{\cdot \mathcal{C}} \mathbf{d}_{\mathcal{C}}$ are also compatible. According to Lemma 1, there exists $\mathbf{u}_{\mathcal{P}} \in \mathbb{R}^{|\mathcal{P}|}$ such that $\mathbf{A}_{\cdot \mathcal{P}} \mathbf{u}_{\mathcal{P}} = \mathbf{A}_{\cdot \mathcal{C}} \mathbf{d}_{\mathcal{C}}$. Let $\mathbf{u} = (\mathbf{u}_{\mathcal{P}}, \mathbf{d}_{\mathcal{C}}, \mathbf{0})$ and $\mathbf{v} = \mathbf{d} - \mathbf{u} = (\mathbf{v}_{\mathcal{C}}, \mathbf{0}, \mathbf{d}_{\mathcal{I}})$. Let $\alpha_{\mathbf{u}} = \sum_{i \in \mathcal{M}} w_i |\mathbf{d}_i| + \mathbf{w}_{\mathcal{C}_{\mathcal{L}}}^T \mathbf{d}_{\mathcal{C}_{\mathcal{L}}} + \mathbf{w}_{\mathcal{C}_{\mathcal{U}}}^T \mathbf{d}_{\mathcal{C}_{\mathcal{U}}}$ and $\alpha_{\mathbf{v}} = \mathbf{w}^T \mathbf{v} = 1 - \alpha_{\mathbf{u}}$. Thus, $0 < \alpha_{\mathbf{u}}, \alpha_{\mathbf{v}} < 1$. Moreover, let $\mathbf{u}' = \mathbf{u}/\alpha_{\mathbf{u}}$ and $\mathbf{v}' = \mathbf{v}/\alpha_{\mathbf{v}}$ be the corresponding normalized directions. With these definitions, $\mathbf{d} = \alpha_{\mathbf{u}} \mathbf{u}' + \alpha_{\mathbf{v}} \mathbf{v}'$. By construction, \mathbf{u}' (resp. \mathbf{v}') is a solution of R-GNI (resp. C-GNI) and $\alpha_{\mathbf{u}} + \alpha_{\mathbf{v}} = \sum_{i \in \mathcal{M}} w_i |\mathbf{d}_i| + \mathbf{w}_{\mathcal{C}_{\mathcal{L}}}^T \mathbf{d}_{\mathcal{C}_{\mathcal{L}}} + \mathbf{w}_{\mathcal{C}_{\mathcal{U}}}^T \mathbf{d}_{\mathcal{C}_{\mathcal{U}}} = 1$, because $\mathbf{d} \in \Delta$. Therefore, $\mathbf{d} = \alpha_{\mathbf{u}} \mathbf{u}' + \alpha_{\mathbf{v}} \mathbf{v}'$ is a convex combination of a solution of R-GNI (\mathbf{u}') and a solution of C-GNI (\mathbf{v}').

Looking at the objective function, the convex combination reads $\mathbf{c}^T \mathbf{d} = \alpha_{\mathbf{u}} (\mathbf{c}^T \mathbf{u}') + \alpha_{\mathbf{v}} (\mathbf{c}^T \mathbf{v}')$. Either $\mathbf{c}^T \mathbf{d} \geq \mathbf{c}^T \mathbf{u}'$ or $\mathbf{c}^T \mathbf{d} \geq \mathbf{c}^T \mathbf{v}'$. However, since every solution of R-GNI and C-GNI is also a solution of GNI and \mathbf{d} is optimal for GNI, $\mathbf{c}^T \mathbf{d} \leq \mathbf{c}^T \mathbf{u}'$ and $\mathbf{c}^T \mathbf{d} \leq \mathbf{c}^T \mathbf{v}'$. One of the two inequalities is thus an equality and the second follows from $\mathbf{c}^T \mathbf{d} = \alpha_{\mathbf{u}} (\mathbf{c}^T \mathbf{u}') + \alpha_{\mathbf{v}} (\mathbf{c}^T \mathbf{v}')$. Therefore, $\mathbf{c}^T \mathbf{d} = \mathbf{c}^T \mathbf{u}' = \mathbf{c}^T \mathbf{v}'$, and since \mathbf{d} is an optimal solution of GNI, $z_{\text{R-GNI}}^* = z_{\text{C-GNI}}^* = z_{\text{GNI}}^*$. \square

Theorem 1 extends a theorem established by Rosat et al. (2013) for the set partitioning problem to general linear programming, and it has no equivalent in past presentations of IPS. It gives strong theoretical support to the compatible-incompatible partition and to the R-GNI/C-GNI decomposition of GNI: it states that one of the optimal solutions of GNI is an optimal solution of one of the decomposed problems. Solving them separately is thus equivalent to solving GNI, i.e., to determining whether or not \mathbf{x}^0 is optimal (see Proposition 1). Since the theorem requires that $\mathbf{w}_{\mathcal{P}} = 0$, we retain this assumption in the rest of this paper. This last condition is not overly restrictive. It is only natural that the emphasis is put on the nonbasic variables during the search for an improving direction. The same assumption is made by Elhallaoui et al. (2011) and by most pricing rules designed for the simplex algorithm.

Corollary 1 *If, for a current solution \mathbf{x} and the associated partition $(\mathcal{P}, \mathcal{C}, \mathcal{I})$, the optimal values of R-GNI and C-GNI are both nonnegative, then \mathbf{x} is optimal for (P).*

This corollary is a straightforward consequence of Theorem 1 and Proposition 1.

One other key aspect of the decomposition is that R-GNI may be further reduced by observing that the rank of its constraint matrix $\mathbf{A}_{\cdot \mathcal{P} \cup \mathcal{I}}$ is the same as the rank of $\mathbf{A}_{\cdot \mathcal{P}}$, i.e., p . Since R-GNI is feasible by construction, it contains $m - p$ redundant constraints that may be removed without harm. Assuming that the rows are permuted so that the first p constraints are independent, these rows are indexed by \mathcal{P} while the redundant ones have their index set denoted $\overline{\mathcal{P}} = \{p + 1, \dots, m\}$. Since we also have $\mathbf{w}_{\mathcal{P}} = 0$, R-GNI may be equivalently rewritten:

$$\begin{aligned} z_{\text{R-GNI}}^* = \min_{(\mathbf{d}_{\mathcal{P}}, \mathbf{d}_{\mathcal{C}})} \quad & \mathbf{c}_{\mathcal{P}}^T \mathbf{d}_{\mathcal{P}} \quad + \quad \mathbf{c}_{\mathcal{C}}^T \mathbf{d}_{\mathcal{C}} \\ \text{s.t.} \quad & \mathbf{A}_{\mathcal{P} \mathcal{P}} \mathbf{d}_{\mathcal{P}} \quad + \quad \mathbf{A}_{\mathcal{P} \mathcal{C}} \mathbf{d}_{\mathcal{C}} \quad = \quad 0 \\ & \sum_{i \in \mathcal{M}} w_i |\mathbf{d}_i| + \mathbf{w}_{\mathcal{C}_{\mathcal{L}}}^T \mathbf{d}_{\mathcal{C}_{\mathcal{L}}} + \mathbf{w}_{\mathcal{C}_{\mathcal{U}}}^T \mathbf{d}_{\mathcal{C}_{\mathcal{U}}} \quad = \quad 1 \\ & \mathbf{d}_{\mathcal{C}_{\mathcal{L}}} \geq 0 \quad , \quad \mathbf{d}_{\mathcal{C}_{\mathcal{U}}} \leq 0 \end{aligned}$$

This new form clearly emphasizes the potential advantage of solving R-GNI instead of GNI to find an improving feasible direction.

The previous results canonically lead to the procedure described in Algorithm 2. The main idea is to first look at compatible columns to perform nondegenerate pivots (R-GNI), and if no suitable columns are found, to look for a collection of incompatible columns that globally supplies an improvement (C-GNI).

Algorithm 2 Compatibility-based Decomposed Primal Improvement Procedure

1. Let \mathbf{x} be an initial feasible solution;
2. Build the partition $(\mathcal{P}, \mathcal{C}, \mathcal{I})$ corresponding to \mathbf{x} ;
3. Solve R-GNI: if $z_{\text{R-GNI}}^* \geq 0$, GOTO 6; else update \mathbf{x} by following the direction found by R-GNI;
4. If some criterion is satisfied, GOTO 2;
5. Build the partition $(\mathcal{P}, \mathcal{C}, \mathcal{I})$ corresponding to \mathbf{x} ;
6. Solve C-GNI: if $z_{\text{C-GNI}}^* \geq 0$, GOTO 8; else update \mathbf{x} by following the direction found by C-GNI;
7. If some criterion is satisfied, $z_{\text{R-GNI}}^* \leftarrow -1$, GOTO 5;
8. If $z_{\text{R-GNI}}^* \geq 0$ and $z_{\text{C-GNI}}^* \geq 0$ RETURN \mathbf{x} ; else GOTO 2.

At Step 7, $z_{\text{R-GNI}}^*$ is arbitrarily set to -1 to ensure that $z_{\text{R-GNI}}^* \geq 0$ and $z_{\text{C-GNI}}^* \geq 0$ (Step 8) can be satisfied only with R-GNI and C-GNI built from the same $(\mathcal{P}, \mathcal{C}, \mathcal{I})$ decomposition. With this precaution, Corollary 1 provides a proof of optimality for the algorithm.

Remark 2 *No theoretical ground requires R-GNI to be given priority over C-GNI. The order in which steps 3 and 6 are performed in Algorithm 2 is based on common sense: R-GNI is easier to solve than C-GNI because compatible columns guarantee nondegenerate pivots. Steps 3 and 6 could possibly be switched. Also, the criteria that appear in steps 4 and 7 may be set to allow for several instances of R-GNI or C-GNI to be solved consecutively. They may thus be used to emphasize the preference for one or the other problem.*

Algorithm 2 provides a strong basis for an efficient algorithm for degenerate LPs. Yet it is too generic to be implemented as such, and choices are necessary to make it a practical LP solver. Matters such as how to find an initial solution, when to update the $(\mathcal{P}, \mathcal{C}, \mathcal{I})$ partition, or when to switch from solving R-GNI to solving C-GNI are more than just implementation details. They are critical features of the algorithm that can dramatically change its computational behavior and the resulting running time. The next section thoroughly addresses the evolution from this theoretical scheme to a numerically efficient procedure.

3 Implementing the Framework: IPS

Practical implementations of Algorithm 2 already exist, even though the associated generic theory was never presented. One of them is, as already mentioned, the IPS algorithm of Elhallaoui et al. (2011). IPS includes practical implementation techniques, but it also makes implicit choices. In particular, all the weights in the normalization constraint are set to 1, and the complementary problem is transformed through a costly algebraic operation to reduce its size. These choices are discussed below, together with the most important implementation issues. Although some of them are discussed in Elhallaoui et al. (2011) and Raymond et al. (2010b), we adapt the presentation to emphasize the relationship between IPS and the aforementioned generic framework. New details are also needed because of the efficient treatment of upper-bounded and unbounded variables. Moreover, our intent to diversify the benchmark and the general view provided by the previous section give rise to several improvements.

3.1 Starting the Algorithm with an Initial Feasible Solution

One important improvement is that the theoretical presentation in the previous section does not assume that the initial feasible solution is basic. We use a simplex phase I to generate the initial solution in our tests, but IPS could start from a feasible solution generated, for instance, by an external heuristic procedure. The only additional work is to identify a maximal linearly independent set of columns within their bounds. This can be done, for instance, by performing Gaussian eliminations until reduction to row echelon form. The operation simultaneously provides the set \mathcal{P} and the $m - p$ redundant constraints of R-GNI.

3.2 Identifying Compatible Variables and Transforming the Complementary Problem

The efficiency of Algorithm 2 relies on the ability to quickly update the partition $(\mathcal{P}, \mathcal{C}, \mathcal{I})$. Once \mathcal{P} is known, identifying \mathcal{C} is equivalent to finding the set of nonbasic columns that belong to $\text{Span}(\mathbf{A}_{\cdot \mathcal{P}})$. A straightforward

way to do this is to complete the columns of $\mathbf{A}_{\mathcal{P}}$ with vectors of \mathbb{R}^m to form a basis \mathbf{B} of \mathbb{R}^m . Since the columns of $\mathbf{A}_{\mathcal{P}}$ are the first p columns, any vector $\mathbf{v} \in \mathbb{R}^m$ is written in the basis \mathbf{B} as $((\mathbf{B}^{-1}\mathbf{v})_{\mathcal{P}}, (\mathbf{B}^{-1}\mathbf{v})_{\overline{\mathcal{P}}})$, with $(\mathbf{B}^{-1}\mathbf{v})_{\mathcal{P}} \in \text{Span}(\mathbf{A}_{\mathcal{P}})$. This decomposition being unique, it follows that $\mathbf{v} \in \text{Span}(\mathbf{A}_{\mathcal{P}})$ if and only if $(\mathbf{B}^{-1}\mathbf{v})_{\overline{\mathcal{P}}} = \mathbf{0}$.

Elhallaoui et al. (2011) and Raymond et al. (2010b) always assume a complete basis is available because the initial solution comes from a simplex phase I. Here, the computations are minimized by choosing the simplest possible completion, i.e.,

$$\mathbf{B} = \begin{bmatrix} \mathbf{A}_{\mathcal{P}\mathcal{P}} & \mathbf{0} \\ \mathbf{A}_{\overline{\mathcal{P}}\mathcal{P}} & \mathbf{I}_{m-p} \end{bmatrix} \Leftrightarrow \mathbf{B}^{-1} = \begin{bmatrix} \mathbf{A}_{\mathcal{P}\mathcal{P}}^{-1} & \mathbf{0} \\ -\mathbf{A}_{\overline{\mathcal{P}}\mathcal{P}}\mathbf{A}_{\mathcal{P}\mathcal{P}}^{-1} & \mathbf{I}_{m-p} \end{bmatrix}.$$

As a consequence, the set of compatible columns is determined by forming the matrix $-\mathbf{A}_{\overline{\mathcal{P}}\mathcal{P}}\mathbf{A}_{\mathcal{P}\mathcal{P}}^{-1}\mathbf{A}_{\mathcal{P}\mathcal{N}} + \mathbf{A}_{\overline{\mathcal{P}}\mathcal{N}}$. Only for $\mathbf{x}_{\mathcal{M}}$ may this calculation be skipped, since those columns are known to be compatible.

In contrast to the identification of compatible columns, the transformation of C-GNI is not necessary for the execution of the algorithm. It is however an essential element of IPS. In C-GNI, the linear constraints $\mathbf{A}_{\mathcal{P}}\mathbf{d}_{\mathcal{P}} + \mathbf{A}_{\mathcal{I}}\mathbf{d}_{\mathcal{I}} = \mathbf{0}$ can be interpreted as follows: the weighted combination of the columns $\mathbf{A}_{\mathcal{I}}\mathbf{d}_{\mathcal{I}}$ that potentially enter the basis must be compatible. These constraints are equivalent to $\mathbf{A}_{\mathcal{I}}\mathbf{d}_{\mathcal{I}} \in \text{Span}(\mathbf{A}_{\mathcal{P}})$. As stated above, another way of putting this condition is

$$\left(-\mathbf{A}_{\overline{\mathcal{P}}\mathcal{P}}\mathbf{A}_{\mathcal{P}\mathcal{P}}^{-1}\mathbf{A}_{\mathcal{P}\mathcal{I}} + \mathbf{A}_{\overline{\mathcal{P}}\mathcal{I}}\right)\mathbf{d}_{\mathcal{I}} = \mathbf{0}. \quad (5)$$

Moreover, since $\mathbf{A}_{\mathcal{P}}$ is a full-rank matrix, for a given $\mathbf{d}_{\mathcal{I}}$ satisfying (5), there is a unique $\mathbf{d}_{\mathcal{P}}$ such that $\mathbf{A}_{\mathcal{P}}\mathbf{d}_{\mathcal{P}} + \mathbf{A}_{\mathcal{I}}\mathbf{d}_{\mathcal{I}} = \mathbf{0}$. This vector is given by $\mathbf{d}_{\mathcal{P}} = -\mathbf{A}_{\mathcal{P}\mathcal{P}}^{-1}\mathbf{A}_{\mathcal{P}\mathcal{I}}\mathbf{d}_{\mathcal{I}}$. Denoting $\bar{\mathbf{c}}^T = \mathbf{c}^T - \mathbf{c}_{\mathcal{P}}^T\mathbf{A}_{\mathcal{P}\mathcal{P}}^{-1}\mathbf{A}_{\mathcal{P}}$ to extend the usual notation of the reduced cost, the solution of C-GNI is thus equivalently found by solving a smaller LP that involves only the variables of $\mathbf{d}_{\mathcal{I}}$:

$$\begin{aligned} z_{\text{C-GNI}}^* &= \min_{\mathbf{d}_{\mathcal{I}}} \bar{\mathbf{c}}_{\mathcal{I}}^T \mathbf{d}_{\mathcal{I}} \\ \text{s.t.} \quad &\left(-\mathbf{A}_{\overline{\mathcal{P}}\mathcal{P}}\mathbf{A}_{\mathcal{P}\mathcal{P}}^{-1}\mathbf{A}_{\mathcal{P}\mathcal{I}} + \mathbf{A}_{\overline{\mathcal{P}}\mathcal{I}}\right)\mathbf{d}_{\mathcal{I}} = \mathbf{0} \\ &\mathbf{w}_{\mathcal{I}}^T \mathbf{d}_{\mathcal{I}} = 1 \\ &\mathbf{d}_{\mathcal{I}\mathcal{L}} \geq \mathbf{0}, \mathbf{d}_{\mathcal{I}\mathcal{U}} \leq \mathbf{0} \end{aligned} \quad (\text{CR-GNI})$$

This transformation of the complementary problem C-GNI into a smaller problem CR-GNI is similar to that performed in the reduced-gradient algorithms as presented by Kallio and Porteus (1978). That is, the reduction of a problem with $p + |\mathcal{I}|$ variables to $|\mathcal{I}|$ variables in such a way that $(\mathbf{d}_{\mathcal{P}}, \mathbf{d}_{\mathcal{I}})$ is easily inferred from $\mathbf{d}_{\mathcal{I}}$ resembles the reduced-gradient approach. However, the reduction in the number of constraints presented here is not mentioned in either the theory (see Kallio and Porteus (1978)) or the MINOS implementation of Wolfe's reduced-gradient (see Murtagh and Saunders (1978)), and neither work suggests adding a normalization constraint that could influence the objective.

Although the density of the constraint matrix changes as we transform C-GNI, in a way that is impossible to predict in the general case, it seems that an important reduction in the computational time could be achieved by considering an LP with fewer variables and constraints. In IPS, R-GNI has p constraints and $p + |\mathcal{C}|$ variables, and CR-GNI has $m - p$ constraints and $n - p - |\mathcal{C}|$ variables. If the number of simplex pivots is similar for the two problems, one would expect IPS to perform best on problems with $p \approx 0.5m$ and approximately as many compatible as incompatible variables.

Focusing now on the computational burden of the algebraic operations described above, we see that the method used to determine \mathcal{C} is closely linked to the choice of whether or not to transform C-GNI. If $-\mathbf{A}_{\overline{\mathcal{P}}\mathcal{P}}\mathbf{A}_{\mathcal{P}\mathcal{P}}^{-1}\mathbf{A}_{\mathcal{P}\mathcal{N}} + \mathbf{A}_{\overline{\mathcal{P}}\mathcal{N}}$ is computed to check the compatibility of the nonbasic variables, there is no other costly operation to execute in order to form CR-GNI. Therefore, solving CR-GNI should be faster than solving C-GNI. On the other hand, if a fast method for identifying the compatible columns is available, much of the time spent computing the constraints of CR-GNI would potentially be saved by keeping the original constraints in C-GNI. With this in mind, Raymond et al. (2010a) recently proposed a stochastic compatibility test that does approximately as many operations as the computation of a reduced cost and identifies all the compatible columns with an extremely small probability of error.

3.3 Setting the Normalization Constraint in the Complementary Problem

The second important implementation choice considered as a part of IPS by Elhallaoui et al. (2011) and Raymond et al. (2010b) is the normalization constraint of C-GNI. For simplicity, all the weights are set to ± 1 in IPS.

To better understand the effect of this choice, one may observe that CR-GNI could be equivalently solved by removing the normalization constraint and minimizing the normalized criterion $\bar{\mathbf{c}}_{\mathcal{I}}^T \mathbf{d}_{\mathcal{I}} / \mathbf{w}_{\mathcal{I}}^T \mathbf{d}_{\mathcal{I}}$. A normalized solution of CR-GNI would then be obtained as $\mathbf{d}_{\mathcal{I}}^* / \mathbf{w}_{\mathcal{I}}^T \mathbf{d}_{\mathcal{I}}^*$. This shows that the normalization constraint actually impacts the criterion guiding the search for an improving compatible combination of columns. For instance, even if it cannot lead to a compatible direction, if the search is restricted to solutions with one nonzero variable, $\bar{c}_i / w_i, i \in \mathcal{I}$ will be minimized. Setting all the weights to ± 1 thus leads to an extension of Dantzig's original pricing criterion (see Dantzig (1955)) to improving directions with more than one nonzero nonbasic variable. The main advantage of Dantzig's pricing is that the weights are easy to compute and do not introduce additional risks of numerical instability. On the other hand, there is no reason why they should lead to good directions.

3.4 Solving the Reduced Problem

The practical solution of R-GNI results from the immediate computation of its optimal solution, as indicated by the proposition below.

Proposition 4 R-GNI admits an optimal solution $(\mathbf{d}_{\mathcal{P}}^*, \mathbf{d}_{\mathcal{C}}^*)$ such that

$$\begin{cases} d_j^* = \frac{1}{w_j}, \text{ for some } j \in \arg \min \left(\left\{ \frac{\bar{c}_i}{w_i} : i \in \mathcal{C}_{\mathcal{L}} \cup \mathcal{C}_{\mathcal{U}} \right\} \cup \left\{ -\frac{|\bar{c}_i|}{w_i} : i \in \mathcal{M} \right\} \right) \\ d_i^* = 0, \text{ for all } i \in \mathcal{C} \setminus \{j\} \\ \mathbf{d}_{\mathcal{P}}^* = -\frac{1}{w_j} \mathbf{A}_{\mathcal{P}\mathcal{P}}^{-1} \mathbf{A}_{\cdot j} \end{cases} \quad (6)$$

Proof. Let $z_{\min} = \min \left(\left\{ \frac{\bar{c}_i}{w_i} : i \in \mathcal{C}_{\mathcal{L}} \cup \mathcal{C}_{\mathcal{U}} \right\} \cup \left\{ -\frac{|\bar{c}_i|}{w_i} : i \in \mathcal{M} \right\} \right)$. For any feasible solution \mathbf{d} of R-GNI, $\mathbf{d}_{\mathcal{P}} = -\mathbf{A}_{\mathcal{P}\mathcal{P}}^{-1} \mathbf{A}_{\cdot \mathcal{C}} \mathbf{d}_{\mathcal{C}}$, so the objective function $\mathbf{c}_{\mathcal{P}}^T \mathbf{d}_{\mathcal{P}} + \mathbf{c}_{\mathcal{C}}^T \mathbf{d}_{\mathcal{C}}$ is equivalently written

$$\begin{aligned} \sum_{i \in \mathcal{C}} \bar{c}_i d_i &= \sum_{i \in \mathcal{C}} \frac{\bar{c}_i}{w_i} w_i d_i \\ &\geq \sum_{i \in \mathcal{C}_{\mathcal{L}} \cup \mathcal{C}_{\mathcal{U}}} \frac{\bar{c}_i}{w_i} w_i d_i + \sum_{i \in \mathcal{M}} \frac{-|\bar{c}_i|}{w_i} w_i |d_i| \\ &\geq \min \left\{ \frac{\bar{c}_i}{w_i} : i \in \mathcal{C}_{\mathcal{L}} \cup \mathcal{C}_{\mathcal{U}} \right\} \sum_{i \in \mathcal{C}_{\mathcal{L}} \cup \mathcal{C}_{\mathcal{U}}} w_i d_i + \min \left\{ -\frac{|\bar{c}_i|}{w_i} : i \in \mathcal{M} \right\} \sum_{i \in \mathcal{M}} w_i |d_i| \\ &\geq z_{\min} \times \left(\mathbf{w}_{\mathcal{C}_{\mathcal{L}}}^T \mathbf{d}_{\mathcal{C}_{\mathcal{L}}} + \mathbf{w}_{\mathcal{C}_{\mathcal{U}}}^T \mathbf{d}_{\mathcal{C}_{\mathcal{U}}} + \sum_{i \in \mathcal{M}} w_i |d_i| \right) = z_{\min}. \end{aligned}$$

One can easily verify that a vector \mathbf{d}^* satisfying (6) is a feasible solution of R-GNI whose objective value is equal to z_{\min} , so it is optimal. \square

The consequence of Proposition 4 is that solving R-GNI and following the optimal solution until a bound is reached is equivalent to pivoting a compatible column into the working basis. As long as the variables of $\mathbf{x}_{\mathcal{P}}$ are strictly within their bounds, those pivots are guaranteed to be nondegenerate. In practice, this guarantee may be costly since it implies that variables at one of their bounds must be removed from \mathcal{P} and \mathcal{C} after each pivot. As a consequence, IPS solves R-GNI several times without updating the partition $(\mathcal{P}, \mathcal{C}, \mathcal{I})$

by running the primal simplex on the following *reduced problem* until a stopping criterion is reached:

$$\begin{cases} \min & \mathbf{c}_{\mathcal{P}}^T \mathbf{x}_{\mathcal{P}} + \mathbf{c}_{\mathcal{C}}^T \mathbf{x}_{\mathcal{C}} \\ \text{s. t.} & \mathbf{A}_{\mathcal{P}\mathcal{P}} \mathbf{x}_{\mathcal{P}} + \mathbf{A}_{\mathcal{C}\mathcal{P}} \mathbf{x}_{\mathcal{C}} = \tilde{\mathbf{b}}_{\mathcal{P}} \\ & \mathbf{l}_{\mathcal{P}} \leq \mathbf{x}_{\mathcal{P}} \leq \mathbf{u}_{\mathcal{P}}, \mathbf{l}_{\mathcal{C}} \leq \mathbf{x}_{\mathcal{C}} \leq \mathbf{u}_{\mathcal{C}} \end{cases} \quad (\text{RED})$$

Here $\tilde{\mathbf{b}}_{\mathcal{P}} = \mathbf{b}_{\mathcal{P}} - \mathbf{A}_{\mathcal{I}\mathcal{C}} \mathbf{l}_{\mathcal{I}\mathcal{C}} - \mathbf{A}_{\mathcal{I}\mathcal{U}} \mathbf{u}_{\mathcal{I}\mathcal{U}}$ is used instead of $\mathbf{b}_{\mathcal{P}}$ to ensure the feasibility of the current solution. The choice of the weights $\mathbf{w}_{\mathcal{C}}$ determines the pricing criterion used to solve RED. For instance, if all the weights are set to ± 1 , the primal simplex selects entering variables according to Dantzig's criterion.

The stopping criterion is that determined by Raymond et al. (2010b): the simplex is stopped after m iterations if optimality is not reached earlier. In contrast to Raymond et al. (2010b), because identifying the compatible variables can be costly, RED is built anew only if at least 10% of the variables in \mathcal{P} are equal to one of their bounds and if 30% of \mathcal{P} has changed since the last update of $(\mathcal{P}, \mathcal{C}, \mathcal{I})$. These modifications ensure that some progress is made between two consecutive manipulations of the constraints, and they keep the size of RED constant provided the degeneracy remains low enough.

3.5 Solving the Complementary Problem

Once an optimal solution of RED has been found, new improving directions are found or optimality is proved by solving CR-GNI. Theoretically, each time that CR-GNI is solved it provides a feasible descending direction. However, this property requires that the partition $(\mathcal{P}, \mathcal{C}, \mathcal{I})$ is updated each time. Furthermore, following the improving direction found by CR-GNI usually results in only a small modification of the current solution and the associated RED. For these reasons, the strict theoretical framework described in Algorithm 2 is usually inefficient in practice. As stated at the end of Subsection 3.4, the decomposition is thus updated only when the solution is degenerate enough and \mathcal{P} has changed significantly since the last update. Following the procedure of Raymond et al. (2010b), we then solve the complementary problem several times in a row. Each time we remove the variables that appear in the resulting direction from the problem and solve CR-GNI again to find a new direction involving different nonbasic variables. These directions are not directly followed; the nonzero variables are instead included in \mathcal{C} together with the other variables belonging to $\text{Span}(\mathcal{P} \cup \mathcal{C})$. This is done to grant the simplex more flexibility in choosing the order of the pivots when solving the resulting RED. In our tests, complementary problems were solved until at least 10% of the columns of \mathcal{I} were selected to be appended to the reduced problem.

As in Raymond et al. (2010b), CR-GNI is solved with the dual simplex algorithm. An improvement is added to ensure that this process starts from a dual feasible solution. With that in mind, consider the dual formulation of CR-GNI with slack variables \mathbf{s} :

$$\begin{aligned} z_{\text{C-GNI}}^* &= \max_{(\boldsymbol{\pi}, y)} y \\ \text{s.t.} & \quad w_i y + \boldsymbol{\pi}^T (-\mathbf{A}_{\overline{\mathcal{P}}\mathcal{P}} \mathbf{A}_{\mathcal{P}\mathcal{P}}^{-1} \mathbf{A}_{\mathcal{P}i} + \mathbf{A}_{\overline{\mathcal{P}}i}) + s_i = \bar{c}_i, \quad \forall i \in \mathcal{I}_{\mathcal{C}} \\ & \quad w_i y + \boldsymbol{\pi}^T (-\mathbf{A}_{\overline{\mathcal{P}}\mathcal{P}} \mathbf{A}_{\mathcal{P}\mathcal{P}}^{-1} \mathbf{A}_{\mathcal{P}i} + \mathbf{A}_{\overline{\mathcal{P}}i}) - s_i = \bar{c}_i, \quad \forall i \in \mathcal{I}_{\mathcal{U}} \\ & \quad \mathbf{s} \geq \mathbf{0}, \boldsymbol{\pi} \in \mathbb{R}^{m-p}, y \in \mathbb{R} \end{aligned} \quad (7)$$

With $j \in \arg \min \{\bar{c}_i/w_i : i \in \mathcal{I}\}$, a feasible solution of (7) is obtained by setting $\boldsymbol{\pi} = \mathbf{0}$, $y = \bar{c}_j/w_j$, $s_j = 0$, and $s_i = \bar{c}_i/w_i - \bar{c}_j/w_j (\geq 0), \forall i \in \mathcal{I} \setminus \{j\}$. The associated basis is made up of y and $\{s_i\}_{i \in \mathcal{I} \setminus \{j\}}$, while $\boldsymbol{\pi}$ and s_j are nonbasic. Clearly, this solution is not optimal since the complementarity conditions would lead to only one nonzero variable in $\mathbf{d}_{\mathcal{I}}$, which contradicts the definition of an incompatible variable. On the other hand, giving this solution to a LP solver helps to avoid a potentially time-consuming phase I.

Once an optimal solution $\mathbf{d}_{\mathcal{I}}^*$ of CR-GNI has been found, all the variables such that $d_i^* \neq 0, i \in \mathcal{I}$, are removed from CR-GNI, which is equivalent to removing the associated constraints and slack variables in (7). As a consequence, the solution remains dual feasible and can be used to efficiently warm-start the solution process for this smaller CR-GNI.

3.6 Summary of the Algorithm

The choices and adaptations that are included in IPS to obtain an efficient practical implementation are summarized in Algorithm 2. We denote by \mathcal{P}^+ the subset of \mathcal{P} indexing the variables strictly within their bounds. Immediately after an update of $(\mathcal{P}, \mathcal{C}, \mathcal{I})$, $\mathcal{P} = \mathcal{P}^+$, but this may change after pivots are performed. The variables selected by CR-GNI to be included in RED are temporarily stored in \mathcal{J} .

Algorithm 3 Improved primal simplex with upper and lower bounds

1. Let \mathbf{x} be an initial feasible solution;
 2. Build the partition $(\mathcal{P}, \mathcal{C}, \mathcal{I})$ corresponding to \mathbf{x} ; $\mathcal{P}_0 \leftarrow \mathcal{P}$;
 3. Solve RED for m iterations or until optimality is reached; update the current solution \mathbf{x} ;
 4. If $|\mathcal{P}^+| < 0.9|\mathcal{P}|$ and $|\mathcal{P}_0 \cap \mathcal{P}^+| < 0.7|\mathcal{P}|$, GOTO 2;
else if RED is not solved to optimality, GOTO 3;
 6. Solve CR-GNI: let $\mathbf{d}_{\mathcal{I}}^*$ and $z_{\text{CR-GNI}}^*$ be the optimal solution and objective value;
 7. If $z_{\text{CR-GNI}}^* \geq 0$, GOTO 10;
else $\mathcal{I}^+ \leftarrow \{i \in \mathcal{I} : d_i^* \neq 0\}$, $\mathcal{J} \leftarrow \mathcal{J} \cup \mathcal{I}^+$;
 8. If $|\mathcal{J}| < 0.1|\mathcal{I}|$, remove $\mathbf{x}_{\mathcal{I}^+}$ from CR-GNI, GOTO 6;
 9. $\mathbf{A}_{\mathcal{J}} \leftarrow \mathbf{A}_{\mathcal{I}} \cap \text{Span}(\mathbf{A}_{\mathcal{P}}, \mathbf{A}_{\mathcal{J}})$; $\mathcal{C} \leftarrow \mathcal{C} \cap \mathcal{J}$, $\mathcal{I} \leftarrow \mathcal{I} \setminus \mathcal{J}$; $\mathcal{J} \leftarrow \emptyset$; GOTO 3;
 10. \mathbf{x} is optimal, RETURN \mathbf{x} .
-

4 Extended Computational Analysis of IPS

4.1 Solution of the LPs

The optimization library CPLEX 12.4¹ is called for the solution of every LP. The initial primal feasible solution is generated through a simplex phase I; the primal simplex is used to solve RED while the dual simplex is used for CR-GNI. IPS is then evaluated by comparing its performance with a direct solution of the complete problem P by the primal simplex of CPLEX, starting from the same initial primal feasible solution. It should be noted that, although both IPS and the direct solution process warm-start from feasible solutions, no presolve is done by CPLEX to take advantage of this warm start. For convenience, this direct primal simplex approach is referred to as DPS in the following discussion.

The default options of CPLEX are used for all the parameters except the pricing method. With default pricing, CPLEX automatically switches between the different criteria implemented for the primal simplex, depending on indicators that are unfortunately not documented or available to the user. There are several disadvantages to this automatic pricing. First and foremost, iteratively starting and stopping CPLEX has a strong undesirable impact on the choices it makes with regards to the selected criterion. As a consequence, the differences between the DPS and IPS performance could originate from the differing pricing methods, which would make interpretation difficult. Second, the choice of the pricing has important consequences on the number of pivots performed while solving an LP and the computational time for each pivot, as illustrated in Table 1. The automatic pricing would thus invalidate any comparison of the number of pivots. We therefore chose a specific pricing method for all the calls to CPLEX, both in DPS and IPS. To choose the method, we tested the three usually most efficient pricing methods implemented in CPLEX on the benchmark detailed in the next section. These methods are the automatic pricing, the devex method of Harris (1973), and an approximate version of the steepest-edge criterion first described by Goldfarb and Reid (1977). These methods respectively correspond to the values 0, 1, and 3 for the CPX_PARAM_PRIIND parameter. Table 1 shows the average runtime and the number of pivots for each problem family. For an easier comparison, we present the results for devex and steepest-edge as ratios relative to the default option. The results show that the approximate steepest-edge is neither systematically better nor worse than the automatic pricing, and so we used it in all the tests.

¹CPLEX is freely available for academic and research purposes under the IBM academic initiative: <http://www-03.ibm.com/ibm/university/academic>

Table 1: Comparison of three pricing criteria in CPLEX

Instance	default		devex (ratios)		steepest-edge (ratios)	
	pivots	cpu	pivots	cpu	pivots	cpu
ACP	17920	1.13 s	0.31	0.75	0.14	0.57
UBFA	128165	192.00 s	1.08	0.66	0.40	0.38
UFL	52926	17.56 s	1.66	1.00	0.77	0.77
VCS	169552	62.78 s	1.68	6.95	0.56	1.77
Mittelmann	672386	255.92 s	1.30	3.20	0.47	1.33

4.2 Description of the Benchmark

The performance of IPS was tested on five families of problems. The first four sets of instances correspond to specific types of problems: airline crew pairing problems (ACP), combined fleet assignment and aircraft routing problems (UBFA), uncapacitated warehouse location problems (UFL), and simultaneous vehicle and crew scheduling problems (VCS). Instances of these problems were used by Elhallaoui et al. (2011) because they exhibit significant degeneracy. We made two noteworthy modifications to this benchmark. The ten VCS instances of Elhallaoui et al. (2011) all had 2084 constraints and about 6000 to 10000 variables. We kept only the largest two, and we added three others involving two to fifty times more variables to observe the sensitivity of the algorithm to variations in the number of variables. Also, to evaluate how the algorithm performs when the variables are bounded, we derived the UBFA instances from the five largest fleet assignment instances of Elhallaoui et al. (2011) by explicitly adding upper bounds on the variables. The set-partitioning constraints implicitly constrain the variables to be less than or equal to one; we added the upper bound $x_i \leq 1$ to every variable x_i .

The instances used by Elhallaoui et al. (2011) and Raymond et al. (2010b) share common features that motivated their choice. They are all highly degenerate, their sizes are reasonable, the variables are all bounded below (by 0) but not above, and they include similar types of constraints. For instance, the ACP, UBFA, and VCS problems all contain a majority of set-partitioning constraints. Although we made some changes to the benchmark by considering larger instances and adding finite upper bounds to the fleet-assignment instances, it is difficult to make a thorough analysis of the strengths and weaknesses of IPS for such a homogeneous benchmark. We therefore added fifteen instances selected from Mittelmann's benchmark², which was built to enable a rigorous evaluation of the commercial and open implementations of the most efficient algorithms for LP. Testing IPS on these instances should help to confirm its strengths and identify its limits. Since the computational burden of the algebraic operations needed to identify compatible variables and transform C-GNI grows much more quickly with the number of constraints than with the number of variables, the dual form of the instances, containing significantly more constraints than variables, was solved by both IPS and DPS. This choice is also usually justified by efficiency concerns when solving an LP with the primal simplex.

Table 2 summarizes the characteristics of the thirty-six instances that we used. It focuses on three aspects of the instances: their sizes, the performance of DPS, and indicators of degeneracy. The header $\rho_{\mathbf{A}}$ stands for the density of the original constraint matrix \mathbf{A} . To characterize the degeneracy during the solution of the instances, we recorded both the average number of degenerate variables and the number of degenerate pivots. The instances were simply numbered from 1 to the number of instances for each type. The Mittelmann instances for which we solved the dual form are indicated with a ' symbol.

4.3 Experimental Results

The tests were all performed on an OpenSuse operating system with an Intel(R) Core(TM) i7-3770 CPU @ 3.40 GHz processor.

IPS as summarized in Algorithm 3 was called on the benchmark to obtain the results in Table 3. These results are relevant to a first analysis of the IPS performance. The headers are partitioned into three types of

²These instances are available online: <http://plato.asu.edu/ftp/lpcom.html>

Table 2: Characteristics of the benchmark

Instance	size of P			DPS solve		degeneracy	
	m: rows	n: cols	ρ_A	piv	cpu	var	piv
ACP1	823	8904	9.96E-03	3148	0.81	59.8%	66.4%
ACP2	531	5198	1.32E-02	1276	0.17	79.2%	76.8%
ACP3	825	8627	9.95E-03	3994	1.13	76.6%	86.6%
ACP4	426	7195	1.70E-02	1244	0.25	48.1%	43.2%
ACP5	801	8308	9.91E-03	3262	0.87	70.3%	84.0%
ACP6	646	7292	1.10E-02	2283	0.46	71.2%	73.8%
UBFA1	5182	23650	2.51E-03	11129	8.78	74.0%	65.2%
UBFA2	5182	23990	2.50E-03	15566	13.05	74.5%	66.9%
UBFA3	5182	24282	2.49E-03	67792	93.72	76.2%	68.5%
UBFA4	5182	24517	2.49E-03	142550	255.36	76.5%	73.2%
UBFA5	5182	24875	2.48E-03	53591	73.42	76.9%	85.8%
UFL1	7965	15330	2.41E-04	20327	5.31	57.3%	45.6%
UFL2	10440	20380	1.87E-04	36072	8.94	69.7%	60.2%
UFL3	15476	30452	1.27E-04	54069	17.57	75.1%	61.3%
UFL4	20534	40568	9.62E-05	60587	23.72	76.7%	71.4%
UFL5	25931	51462	7.65E-05	58382	23.21	80.2%	84.3%
VCS1	2084	10343	1.51E-02	21091	20.26	41.8%	70.3%
VCS2	2084	10150	1.57E-02	22980	21.65	44.5%	72.3%
VCS3	2085	26350	1.66E-02	42341	93.07	64.8%	47.9%
VCS4	1200	133572	1.70E-02	12479	68.43	50.2%	66.3%
VCS5	1600	570983	1.20E-02	25849	607.08	62.7%	66.9%
dano3	3202	15851	1.61E-03	13785	5.22	33.6%	60.1%
df001	6071	12230	4.80E-04	21895	8.77	48.8%	70.8%
fome12	24284	48920	1.20E-04	78964	66.61	49.6%	66.6%
l30	2701	16281	1.18E-03	123704	136.77	30.6%	98.0%
lp22	2958	16392	1.41E-03	36250	20.38	29.3%	80.0%
mod2	34774	66409	8.65E-05	38023	43.79	21.6%	12.1%
neos1'	1892	133473	1.86E-03	33580	111.98	85.8%	99.0%
neos2'	1560	134128	2.65E-03	39460	149.87	80.7%	99.2%
nug15	6330	22275	6.73E-04	92087	257.49	24.9%	55.6%
qap12	3192	8856	1.36E-03	20727	17.21	22.3%	32.8%
qap15	6330	22275	6.73E-04	108918	339.23	25.5%	60.6%
rail2586	2586	923269	3.36E-03	45855	1248.37	31.3%	36.0%
rail4284	4284	1096890	2.40E-03	88617	2520.75	29.2%	52.4%
rlfprim'	8052	74970	4.67E-04	3224	2.31	48.5%	100.0%
world	34506	67147	8.58E-05	44387	56.85	22.7%	13.7%

information. The first two columns record the number of times the partition $(\mathcal{P}, \mathcal{C}, \mathcal{I})$ was updated (“ \mathcal{PCT} ”) and the number of times the reduced problem was expanded by solving CR-GNI several times (“EXP”). The following columns give the proportion of the runtime spent in computing the partitions $(\mathcal{P}, \mathcal{C}, \mathcal{I})$, solving RED, and solving CR-GNI, as well as the total runtime of IPS. For an easier comparison, the last four columns contain the ratios of the performance measures of DPS relative to those of IPS, thus indicating the improvement factor of IPS. The headers “ piv_{RED} ” and “deg piv” respectively stand for the total number of primal simplex pivots and the proportion of degenerate pivots, and “cpu” stands for the total runtime. The quantity in the “ cpu_{RED} ” column is the ratio of the time spent solving RED relative to the DPS runtime. This column gives a measure of the improvement that could be expected for IPS if $(\mathcal{P}, \mathcal{C}, \mathcal{I})$ was updated rapidly and the solution of C-GNI was fast even if the problem is not transformed (since the transformation of C-GNI comes with costly updates of $(\mathcal{P}, \mathcal{C}, \mathcal{I})$). Average values are provided for the improvement factors of the first four types of problems, but given the heterogeneity of the Mittelmann instances it did not seem relevant to base an interpretation on average values. The analysis of the Mittelmann instances is instead done on a case-by-case basis.

The solution was on average 3.34 times faster with IPS than with DPS for the VCS instances. Similar results were observed by Elhallaoui et al. (2011) and Raymond et al. (2010b) for instances whose sizes are comparable with those of VCS1 and VCS2. The interesting point is that IPS remains efficient on VCS3,

Table 3: Evaluation of IPS against CPLEX

	phase count		cpu allocation				improvement factors			
	<i>PCI</i>	EXP	<i>PCI</i>	RED	CR-GNI	cpu	piv _{RED}	deg piv	cpu _{RED}	cpu
VCS1	3	13	4.2%	87.5%	6.5%	8.01 s	1.40	1.72	2.89	2.53
VCS2	4	14	5.7%	82.3%	9.8%	8.24 s	1.45	1.71	3.19	2.63
VCS3	5	14	10.4%	40.7%	45.4%	15.72 s	2.94	1.01	14.53	5.92
VCS4	2	1	8.0%	89.1%	2.4%	23.37 s	1.67	1.41	3.29	2.93
VCS5	3	6	5.8%	64.5%	27.9%	226.42 s	1.87	1.40	4.15	2.68
Average							1.87	1.45	5.61	3.34
ACP1	1	3	22.4%	70.9%	4.0%	0.48 s	1.50	1.25	2.40	1.70
ACP2	1	5	17.7%	60.8%	15.2%	0.16 s	1.10	0.97	1.79	1.09
ACP3	1	5	23.3%	55.5%	15.3%	0.34 s	2.78	1.42	6.03	3.35
ACP4	1	3	19.9%	71.1%	4.3%	0.21 s	1.16	1.36	1.69	1.20
ACP5	1	3	22.1%	66.9%	7.8%	0.34 s	2.00	1.29	3.89	2.60
ACP6	1	6	19.4%	56.8%	16.5%	0.28 s	1.64	1.16	2.92	1.66
Average							1.70	1.24	3.12	1.93
UBFA1	1	7	17.4%	45.3%	34.3%	6.08 s	1.12	1.90	3.19	1.44
UBFA2	1	8	16.5%	44.4%	35.6%	6.47 s	1.49	1.86	4.54	2.02
UBFA3	1	8	16.0%	48.5%	32.3%	6.72 s	5.94	1.56	28.74	13.94
UBFA4	1	8	14.6%	43.0%	39.4%	7.77 s	11.41	1.97	76.43	32.86
UBFA5	1	10	14.9%	49.3%	32.5%	7.26 s	4.26	1.89	20.53	10.12
Average							4.84	1.84	26.68	12.08
UFL1	3	15	18.4%	57.0%	21.3%	3.79 s	1.76	1.43	2.46	1.40
UFL2	3	13	25.9%	38.4%	31.8%	4.37 s	3.21	1.47	5.32	2.04
UFL3	4	12	32.6%	37.4%	26.8%	9.00 s	2.80	1.56	5.22	1.95
UFL4	3	11	32.2%	33.7%	31.4%	14.23 s	2.51	2.32	4.94	1.67
UFL5	4	10	43.8%	19.8%	33.9%	18.44 s	2.82	1.78	6.36	1.26
Average							2.62	1.71	4.86	1.66
dano3	3	6	10.6%	73.4%	13.7%	3.73 s	1.14	1.46	1.91	1.40
df001	3	13	11.1%	52.0%	35.0%	9.95 s	1.03	1.37	1.69	0.88
fome12	3	29	9.0%	29.9%	59.5%	187.39 s	0.70	1.03	1.19	0.36
l30	1	2	1.1%	55.3%	42.9%	28.56 s	4.66	1.03	8.66	4.79
lp22	3	10	4.1%	88.1%	7.3%	14.40 s	1.48	1.67	1.61	1.41
mod2	1	13	7.7%	82.7%	8.8%	64.42 s	0.73	1.20	0.82	0.68
neos1'	3	7	25.8%	2.3%	65.7%	8.10 s	20.33	1.20	602.02	13.82
neos2'	3	3	18.7%	4.9%	74.2%	12.77 s	11.62	1.14	237.88	11.74
nug15	4	7	7.6%	81.3%	10.9%	262.59 s	1.17	1.29	1.21	0.98
qap12	3	6	9.3%	79.3%	10.6%	23.49 s	0.84	0.93	0.92	0.73
qap15	4	8	7.0%	71.5%	21.3%	311.34 s	1.30	1.49	1.52	1.09
rail2586	6	13	13.8%	39.8%	44.5%	619.64 s	1.09	1.07	5.06	2.01
rail4284	5	12	7.6%	26.3%	65.2%	2048.01 s	1.14	0.93	4.67	1.23
rlfprim'	1	1	80.3%	0.6%	17.0%	2.49 s	1612.00	988.14	153.73	0.93
world	1	14	7.3%	82.5%	9.3%	73.27 s	0.78	1.23	0.94	0.78

VCS4, and VCS5, where the number of variables is two to fifty times larger. As expected, the total number of pivots and the proportion of degenerate pivots are lower for IPS than for DPS, and an important indicator of efficiency is that the time spent in transforming and solving CR-GNI remains a small part of the total runtime. The instance VCS4 seems to be an exception, but this is because that IPS performed many fewer pivots in RED, thus giving more importance to the partition updates and CR-GNI. For the ACP instances the global trends were similar to those for the VCS instances, which is consistent with the fact that both VCS and ACP mostly contain set-partitioning constraints.

Referring to Table 2, it appears that around 75% of the basic variables were degenerate for the UBFA instances. About one third of these variables are at their upper bounds ($x_i = 1$) while the rest are at their lower bounds ($x_i = 0$). Because of this high number of degenerate variables, 65% to 80% of the pivots of DPS are degenerate, and there is a great variability in the number of pivots although the instances all derive from the same problem and have similar sizes. The impressive average improvement by a factor of 12 for this class of problem reflects a stabilization of the solution process for these highly degenerate problems: as expected, the IPS runtimes are close for these instances, ranging from 6.1 s to 7.8 s, and they globally increase with the

instance sizes. These positive results validate the potential usefulness of IPS for degenerate problems with bounded variables.

The UFL instances clearly illustrate the impact of the partition updates on the overall execution of IPS. Although IPS achieves a considerable reduction in the number of pivots, leading to an average 2.6 improvement factor and a proportion of degenerate pivots that is 1.7 times smaller than that for DPS, the total improvement is not as good as that for the other families of problems. This is mostly due to the computational time spent in the update of $(\mathcal{P}, \mathcal{C}, \mathcal{I})$ (see column “ \mathcal{PCT} ”), which seems to grow much faster than the time spent solving RED as the number of constraints increases. As a consequence, the column “ cpu_{RED} ” shows a large improvement factor (4.9 on average) over DPS. This shows the potential for improvement in IPS if the partition can be updated quickly and the ensuing transformation of C-GNI avoided in such cases.

In contrast to the previously examined families of instances, the performance of IPS is much less homogeneous for the Mittelmann instances. A significant improvement (greater than 20%) was achieved for 7 of the 15 instances (dano3, l30, lp22, neos1’, neos2’, rail2586, rail4284), the computational time was significantly worse for 4 instances (fome12, mod2, gap12, world), and it was similar for the remaining 4 instances (df001, nug15, qap15, rlfprim’).

Starting with the largest improvements, we see that the success of IPS on the two similar instances neos1’ and neos2’ and on l30 is clearly due to the massive proportion of degenerate pivots (98% to 99%) when executing the primal simplex in DPS (see Table 2). IPS avoids most of the degenerate pivots by solving CR-GNI, and RED is extremely small. These results are consistent with the purpose of IPS.

The other four instances for which IPS performed significantly better than DPS (dano3, lp22, rail2586, rail4284) share several attributes: they have comparable constraint matrix densities and numbers of constraints, and the proportion of degenerate variables is close to 30%. Since the computational time spent updating $(\mathcal{P}, \mathcal{C}, \mathcal{I})$ depends on the number of constraints and the number of degenerate variables, it is logical that it is responsible for only a small part of the total computational time. On the other hand, although the overall cpu improvement factors are close (ranging from 1.2 to 2.0), the solution of CR-GNI represents a much greater proportion of the total runtime for the rail instances. As a consequence, a significant gain could be achieved if the solution of CR-GNI could be sped up for these instances; an improvement factor close to 5 was recorded for the solution of RED alone. Since rail2586 and rail4284 are instances of a set covering problem, it would be interesting to investigate the development of an algorithm specialized for such problems, comparable to the dynamic constraint-aggregation technique of Elhallaoui et al. (2005) for the set partitioning problem.

The remaining eight Mittelmann instances did not lead to a significant improvement. The possible identification of probable causes is supported by Table 4. This table is restricted to these eight instances. The bold font highlights the information most relevant to explaining the low improvement factors, and the instances with similar features are grouped in adjacent rows. The columns related to degeneracy and the improvement factors reproduce the corresponding columns of Tables 2 and 3. The second group of columns focuses on the reduction of P . To better emphasize the comparison with DPS, the computational time spent in the reduction is divided by the DPS runtime. The second column (“ $|\mathcal{C}| / |\mathcal{P}|$ ”) displays the maximum value of the ratio of the cardinality of \mathcal{C} and \mathcal{P} found after the partition updates. This value is an indicator of the potential progress made when solving RED: if it is large, there is a greater probability that many cheap pivots are performed. The third group of columns focuses on the solution of CR-GNI. The first column indicates the impact of the transformation of C-GNI on the density of its constraint matrix. The second column shows the CR-GNI computational time divided by the DPS runtime.

The first group of instances, mod2, qap12, and world, exhibit low proportions of degenerate variables and pivots when solved by DPS. Since the main benefits of IPS are the opportunity to solve a reduced problem and to perform fewer degenerate pivots, these two instances leave little room for improvement. The IPS runtime is larger than that for DPS, and so is the time spent solving RED. It appears that IPS is not appropriate for such instances. A simple adaptation would be to wait for an important proportion of degenerate pivots to

Table 4: Focus on the unimproved sets CPLEX

	degeneracy		reduction		solve CR-GNI		improvement factors	
	var	piv	IPS/DPS	$ \mathcal{C} / \mathcal{P} $	$\rho_{\text{CR-GNI}}/\rho_{\mathbf{A}}$	IPS/DPS	cpu_{RED}	cpu
mod2	21.6%	12.1%	11.3%	76.9%	2.23	12.9%	0.82	0.68
qap12	22.3%	32.8%	4.1%	185.6%	85.7	14.5%	0.92	0.73
world	22.7%	13.7%	9.3%	78.6%	2.1	12.0%	0.94	0.78
nug15	24.9%	55.6%	2.9%	248.0%	442.5	11.1%	1.21	0.98
qap15	25.5%	60.6%	2.0%	237.5%	95.1	19.5%	1.52	1.09
df001	48.8%	70.8%	4.8%	35.9%	2.5	39.7%	1.69	0.88
fome12	49.6%	66.6%	9.7%	21.9%	1.8	167.3%	1.19	0.36
rlfprim'	48.5%	100.0%	85.9%	2.5%	3.4	18.3%	153.73	0.93

be performed before starting the algorithm. With such a modification, IPS would be equivalent to DPS on such instances.

The following two instances, nug15 and qap15, also have low proportions of degenerate variables, but their proportions of degenerate pivots respectively reach 55.6% and 60.6%. As expected, the improvement column “ cpu_{RED} ” thus displays a potential improvement for IPS. The second prominent point relates to the increase in the constraint matrix density after the transformation. It appears that the matrix of CR-GNI is respectively about 95 and 440 times denser than \mathbf{A} for nug15 and qap15. This certainly affects the solution of CR-GNI and is another motivation for IPS not to resort to the transformation of \mathbf{A} .

Remark 3 *Overall, we observed that the transformation of C-GNI leads to a less dense matrix only for the ACP, VCS, and UBFA instances, which concurs with the global deduction that IPS has a special affinity with certain classes of problems, such as set partitioning problems.*

Regarding the third group (df001 and fome12), the most probable cause of the mediocre performance of IPS is the number of compatible variables after each reduction. The direct consequence is that after the reduction, a large majority of the variables in RED belong to the working basis, leaving little hope for important progress toward optimality. This also has a direct impact on the time spent solving the complementary problem since many variables have to be returned to RED before it includes a sufficient number of variables. A simple adaptation would be to estimate the number of compatible variables before starting IPS. However, that is currently costly since it requires the computation of $-\mathbf{A}_{\overline{\mathcal{P}}}\mathbf{A}_{\overline{\mathcal{P}}}^{-1}\mathbf{A}_{\overline{\mathcal{P}}\mathcal{N}} + \mathbf{A}_{\overline{\mathcal{P}}\mathcal{N}}$, so a faster method would need to be used.

Finally, rlfprim' is a special case because the initial feasible solution is optimal. With DPS, 3224 degenerate pivots are performed before proving optimality, whereas IPS solves CR-GNI once before declaring that the solution is optimal. This extreme efficiency is sadly counterbalanced by the fact that the sole partition update takes more than 85% of the DPS runtime. Once again, the only possible fix for such situations would be to run IPS with fast partition updates and no transformation of C-GNI.

5 Conclusion

We have developed a new decomposition-based primal algorithm for degenerate LPs. The decomposition relies on the partition of the nonbasic variables into compatible (\mathcal{C}) and incompatible variables (\mathcal{I}), the first group containing the variables that can be pivoted into the working basis \mathcal{P} with a strict improvement in the objective value. This separates the overall improvement problem, GNI, into two smaller problems, R-GNI and C-GNI, respectively focusing on \mathcal{C} and \mathcal{I} . We have shown that solving R-GNI and C-GNI is equivalent to solving GNI, and we have proposed a decomposition procedure to find the optimum of the LP.

Among the multiple possible implementations of this generic framework is the IPS algorithm, first described by Elhallaoui et al. (2011). The following developments focus on the choices to be made and the issues to be addressed when moving from the theoretical procedure to IPS. Some new insight into the solution

of R-GNI and C-GNI results from this approach that considers IPS as a particular implementation of a much more general algorithm. We make a link between the normalization weights and the pricing criterion; we explain how the algorithm can start from any feasible solution, whether or not it is basic; and we highlight the importance of warm-starting the solution of C-GNI. It appears that the most important choices relate to the identification of the compatible variables, the transformation of the constraint matrix in C-GNI, and the choice of the weights appearing in its normalization constraint. For instance, the first two choices involve costly algebraic operations, while the normalization weights are simply set to ± 1 , thus extending the Dantzig pricing criterion to improving directions with multiple nonbasic nonzero variables.

Finally, we conducted experiments on a diversified version of the instances of Elhallaoui et al. (2011) extended with fifteen data sets from Mittelmann’s benchmark. IPS outperformed direct solution with the primal simplex on several classes of problems including set partitioning, set covering, and warehouse location problems. It also achieved impressive gains for some very degenerate instances from Mittelmann’s benchmark. On the other hand, it gave no improvement on half of Mittelmann’s benchmark. A thorough analysis of these results identified criteria that should be checked before switching from the standard simplex algorithm to IPS. It is necessary to observe a sufficiently large number of degenerate variables and pivots, otherwise there is no reason why IPS should perform better than the primal simplex. Moreover, the number of compatible variables must be large enough, or solving RED will not make significant progress toward optimality. As for the pricing criteria, perturbations, bound shifting, and other techniques included in efficient implementations of the simplex algorithm, IPS should not be used blindly but only when certain criteria are met.

We also demonstrated the limits of the current update of the partition $(\mathcal{P}, \mathcal{C}, \mathcal{I})$, the most important one being that it may be time consuming. This indicates the necessity for an implementation of IPS that would update $(\mathcal{P}, \mathcal{C}, \mathcal{I})$ quickly and solve C-GNI efficiently even without transformation of its constraint matrix. The positive-edge criterion described by Raymond et al. (2010a) could be considered since it is a fast method for the identification of compatible variables. Moreover, we noticed that although the transformation leads to a much smaller form of C-GNI, it may also increase the density of the constraint matrix. Further experimentation will be needed to get more insight into this issue.

Our theoretical and experimental analyses suggest numerous interesting directions for further research. For instance, since solving C-GNI or R-GNI does not always provide feasible improving directions in IPS, the important nondegeneracy property of the generic theoretical scheme is lost for the sake of practical efficiency. As a consequence, it is also interesting to interpret IPS as an analytical—as opposed to systematic—partial pricing strategy for the primal simplex algorithm. Let $\mathcal{J} \subset \mathcal{N}$ be the set of nonbasic variables that are considered during the pricing. This set is initialized with the set of compatible variables \mathcal{C} . Later on, \mathcal{J} goes through several consecutive phases of expansions and reductions respectively triggered by insufficient progress or problematic degeneracy. The reduction is done by keeping in \mathcal{P} only the variables strictly within their bounds and updating \mathcal{C} accordingly. The expansion is performed by solving a sequence of complementary problems, as described in Subsection 3.5. This expansion is an unusually complex procedure for a partial pricing, which is why we qualified it as “analytical.” This interpretation could certainly be fruitful for further improvements of IPS.

McCormick and Shioura (2000) proved that a particular choice of the weights appearing in the normalization constraints makes the primal improving procedure polynomial if the improvement problem is considered as an oracle. This result underlines the special attention that should be given to this constraint in future developments. With additional calculations, the weights could be modified to emulate the steepest-edge rule of Goldfarb and Reid (1977), the devex of Harris (1973), or any other normalized pricing rule. One may also devise another normalization resulting in a faster solution of C-GNI.

From a more practical point of view, using an open optimization library would be valuable. As in Towhidi et al. (2012), the integration of IPS as an internal procedure of the COIN-OR CLP solver would have several benefits. The extended control over the simplex algorithm would allow us to adapt the choices made during the solution process to the particularities of IPS. Moreover, with complete information on the simplex algorithm available, it would be possible to fill some gaps in the interpretation of the performance of IPS.

Finally, we have focused on the primal version of the simplex algorithm. However, better results are often reported for the dual simplex, and it is also a powerful reoptimization tool used for instance in mixed integer programming to explore the branch-and-bound tree. For these reasons, it would be interesting to develop a dual version of IPS to take advantage of dual degeneracy when solving LPs with the dual simplex.

References

- Benichou, M., J.M. Gauthier, G. Hentges, G. Ribiere. 1977. The efficient solution of large-scale linear programming problems: Some algorithmic techniques and computational results. *Mathematical Programming* **13** 280–322. doi:10.1007/BF01584344.
- Dantzig, G.B. 1955. The general simplex method for minimizing a linear form under inequality constraints. *Pacific Journal of Mathematics* **5** 183–195.
- Elhallaoui, I., A. Metrane, G. Desaulniers, F. Soumis. 2011. An improved primal simplex algorithm for degenerate linear programs. *INFORMS Journal on Computing* **23** 569–577.
- Elhallaoui, I., D. Villeneuve, F. Soumis, G. Desaulniers. 2005. Dynamic aggregation of set-partitioning constraints in column generation. *Operations Research* **53** 632–645.
- Gill, P.E., W. Murray, M.A. Saunders, M.H. Wright. 1989. A practical anti-cycling procedure for linearly constrained optimization. *Mathematical Programming* **45** 437–474. doi:10.1007/BF01589114.
- Goldfarb, D., J. Reid. 1977. A practicable steepest-edge simplex algorithm. *Mathematical Programming* **12** 361–371.
- Greenberg, H.J. 1978. *Design and Implementation of Optimization Software*, chap. Pivot selection tactics. Sijthoff & Noordhoff, 109–142.
- Harris, P.M.J. 1973. Pivot selection method of the devex LP code. *Mathematical Programming* **5** 1–28.
- Kallio, M., E.L. Porteus. 1978. A class of methods for linear programming. *Mathematical Programming* **14** 161–169. doi:10.1007/BF01588963.
- Maros, I. 2003. *Computational Techniques of the Simplex Method*. International Series in Operations Research and Management Science, Kluwer Academic Publishers.
- McCormick, S.T., A. Shioura. 2000. Minimum ratio canceling is oracle polynomial for linear programming but not strongly polynomial even for networks. *Operations Research Letters* **27** 199–207.
- Metrane, A., F. Soumis, I. Elhallaoui. 2010. Column generation decomposition with the degenerate constraints in the subproblem. *European Journal of Operational Research* **207** 37–44.
- Murtagh, B.A., M.A. Saunders. 1978. Large-scale linearly constrained optimization. *Mathematical Programming* **14** 41–72. doi:10.1007/BF01588950.
- Pan, P.-Q. 2008. A primal deficient-basis simplex algorithm for linear programming. *Applied Mathematics and Computation* **196** 898–912.
- Perold, A.F. 1980. A degeneracy exploiting LU factorization for the simplex method. *Mathematical Programming* **19** 239–254.
- Raymond, V., F. Soumis, A. Metrane, J. Desrosiers. 2010a. Positive edge: A pricing criterion for the identification of non-degenerate simplex pivots. *Les Cahiers du GERAD* G-2010-61, HEC Montréal, Canada.
- Raymond, V., F. Soumis, D. Orban. 2010b. A new version of the improved primal simplex for degenerate linear programs. *Computers & Operations Research* **37** 91–98.
- Rosat, S., I. Elhallaoui, F. Soumis, A. Lodi. 2013. Integral simplex using decomposition with primal cuts. *Les Cahiers du GERAD* G-2013-79, HEC Montréal, Canada.
- Towhidi, M., J. Desrosiers, F. Soumis. 2012. The positive edge pivot rule within COIN-OR's CLP. *Les Cahiers du GERAD* G-2012-77, HEC Montréal, Canada.