

1 Linear Optimization Terminology

When it was first used in the context of operations research, the word *programming* conveyed the idea of *planning* or *scheduling*. For decades, *linear programming* and *mixed integer programming* have been synonymous with (very large) scheduling problems. The acronyms LP, QP, NLP, MIP, MINLP, SOCP, SDP have denoted increasingly complex problems in the field of *mathematical programming*.

It is hard to break away from such names, but in 2010 the Mathematical Programming Society officially changed its name to the Mathematical Optimization Society. We should now talk about LO problems (and QO, NLO, MILO, MINLO, SOCO, SDO problems) (!!).

2 The LO Problem

Mathematical Optimization systems such as CPLEX, Gurobi, Mosek, ... are designed to solve linear optimization problems in the form

LO1	$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && \phi(x) = c^T x \\ & \text{subject to} && \ell \leq \begin{pmatrix} x \\ Ax \end{pmatrix} \leq u, \end{aligned}$
-----	---

where A is a sparse matrix ($m \times n$). Typically $m < n$ or $m \ll n$, but some problems have $m \gg n$. In practice, many of the variables are likely to have “vanilla” bounds: $\ell_j = 0$, $u_j = \infty$. Sometimes there are *free variables* with both bounds infinite: $\ell_j = -\infty$, $u_j = \infty$. There may also be *fixed variables* for which $\ell_j = u_j = \beta$, say (especially with equality constraints $a_i^T x = \beta$). Other variables are simply called *bounded*.

The main question is *Which bounds are active at a solution?* Clearly, free variables should be the *best kind* because we know in advance that their bounds are not active. Good implementations (and text books!) take advantage of free variables.

Since bounds on variables are easier to work with than general inequalities, implementations introduce a full set of “slack variables” that are *defined* by the relation $s = Ax$. They then treat problem LO1 in the following equivalent form:

LO2	$\begin{aligned} & \underset{x \in \mathbb{R}^n, s \in \mathbb{R}^m}{\text{minimize}} && \phi(x) = c^T x \\ & \text{subject to} && (A \ -I) \begin{pmatrix} x \\ s \end{pmatrix} = 0, \quad \ell \leq \begin{pmatrix} x \\ s \end{pmatrix} \leq u. \end{aligned}$
-----	---

The matrix $-I$ is usually represented implicitly (sometimes it is $+I$), and its presence helps ensure that certain matrices have full rank. Otherwise, it is important to understand that solution algorithms may treat x and s equally. *They are the variables in the problem.*

As a result, there is no loss of generality if we describe algorithms in terms of the simpler problem

LO	$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && \phi(x) = c^T x \\ & \text{subject to} && Ax = b, \quad \ell \leq x \leq u, \end{aligned}$
----	---

where b could be zero or not. If we still define A to be $m \times n$, we now have $m \leq n$.

Commercial mathematical optimization systems are likely to contain the following main algorithms for LO problems:

Primal Simplex
Dual Simplex
Barrier (usually a primal-dual interior method).

The simplex algorithms maintain $Ax = b$ throughout, and a solution estimate x is “infeasible” if the bounds $\ell \leq x \leq u$ are not satisfied (to within some tolerance). In contrast, Barrier maintains $\ell < x < u$ strictly (for $l_j < u_j$), and x is “infeasible” if $Ax \neq b$.

3 Primal Simplex

All versions of the simplex method are *active-set methods* in which n constraints are active at the beginning and end of each iteration. The m constraints $Ax = b$ are always active, and an additional $n - m$ bounds (or temporary bounds) are active. This setup partitions x into $n - m$ *nonbasic variables* x_N that are temporarily frozen at their current values, and m *basic variables* that will be free to move at the next iteration after one of the nonbasic variables is chosen to move. For some permutation P , the columns of A are partitioned accordingly:

$$AP = (B \ N), \quad x = P \begin{pmatrix} x_B \\ x_N \end{pmatrix}, \quad Ax = Bx_B + Nx_N = b, \quad (1)$$

where the P must be chosen so that the $m \times m$ *basis matrix* B is nonsingular. We assume that linear systems of the form

$$Bx = v \quad \text{or} \quad B^T y = w \quad (2)$$

can be solved for any given vectors v and w . Often v and w are *sparse right-hand sides*, and for extremely large problems the sparsity must be taken into account. In particular, most right-hand sides v are columns of A itself. Typically there are only 1 to 10 nonzeros in each column of A *regardless of the problem dimensions*.

Given a nonsingular B and values of x_N satisfying $\ell_N \leq x_N \leq u_N$, the Primal Simplex Method begins by solving

$$Bx_B = b - Nx_N$$

and taking x_B to be new values for the basic variables. From (1), this ensures that $Ax = b$. If the new x_B satisfies its bounds, the new x is “feasible” and Primal Simplex may proceed normally. Otherwise, an almost identical Phase 1 procedure is employed to push x_B toward feasibility.

The following steps are constructive. They illustrate how to generate a search direction Δx that is a *descent direction* (improving direction) for the objective $\phi(x)$. Primal Simplex does this by considering “Shall we move one of the nonbasic variables either up or down”.

If there is no such direction, the current x is an optimal solution. Otherwise it is good to move as far as possible along the search direction, because $\phi(x)$ is linear. Usually a basic variable reaches a bound and a basis exchange takes place. The process then repeats.

Similar steps are used by the *Reduced-Gradient Method* for problems in which $\phi(x)$ is a nonlinear objective. The main difference is that the RG method moves several nonbasic variables simultaneously (we call them *superbasic*) and it might be better not to move too far along the search direction.

The Primal Simplex Method without a_{ij} or B^{-1}

Search direction Δx Assume $\ell \leq x \leq u$ and $Ax = b$. To maintain $A(x + \Delta x) = b$ for some search direction Δx , we require $A\Delta x = 0$. Thus, if the nonbasic variables change in some direction Δx_N , we must solve

$$B\Delta x_B = -N\Delta x_N.$$

Effect on objective In order to choose a “good” Δx_N , suppose y and z_N are defined by

$$B^T y = c_B, \quad z_N = c_N - N^T y. \quad (3)$$

The objective changes from its current value $\phi = c^T x$ to $\bar{\phi}$ as follows:

$$\begin{aligned} \bar{\phi} &= c^T(x + \Delta x) \\ &= \phi + c_B^T \Delta x_B + c_N^T \Delta x_N \\ &= \phi + y^T B \Delta x_B + c_N^T \Delta x_N \\ &= \phi - y^T N \Delta x_N + c_N^T \Delta x_N \\ &= \phi + z_N^T \Delta x_N. \end{aligned}$$

The elements of z_N are the *reduced costs* or *reduced gradients*.

Choice of nonbasic to move Only one nonbasic variable is moved at a time. Suppose we try to move variable x_s . Let $\sigma = \pm 1$ because we’re not sure yet if x_s should increase or decrease. Then

$$\Delta x_N = \sigma e_s, \quad B \Delta x_B = -\sigma a_s, \quad \bar{\phi} = \phi + \sigma z_s. \quad (4)$$

Hence, the objective will improve if $\sigma = 1$, $z_s < 0$ or if $\sigma = -1$, $z_s > 0$. In general it is good if $|z_s|$ is large. The search for z_s is called *Pricing*.

Optimality No improvement is possible if one of the following conditions holds for every nonbasic variable x_j :

$$\begin{aligned} \ell_j < x_j < u_j & \quad \text{and} \quad z_j = 0, \\ x_j = \ell_j & \quad \text{and} \quad z_j \geq 0, \\ x_j = u_j & \quad \text{and} \quad z_j \leq 0. \end{aligned}$$

We know that $z_B \equiv c_B - B^T y = 0$ by definition of y . Thus we can say that a triple (x, y, z) solves problem LO if x is feasible ($Ax = b$ and $\ell \leq x \leq u$) and if

$$z = c - A^T y, \quad \min(x - \ell, z) \leq 0, \quad \min(u - x, -z) \leq 0.$$

Similar tests are suitable when the objective or constraints are nonlinear.

Step length If (x, y, z) is not optimal, the search direction is computed as in (4). The objective improves as much as possible in that direction if we choose a step length α to solve the one-dimensional problem

$$\max \alpha \quad \text{subject to} \quad \ell \leq x + \alpha \Delta x \leq u.$$

This is often called *the ratio test*. One of the variables x_r will reach its bound first. Note that the winner is sometimes x_s : it moves from its current value to one of its bounds. In this case, B is unaltered.

Basis change In general, $r \neq s$ and one of the *basic* variables is first to reach a bound. Further improvement might be possible if we switch x_s and x_r in B and N . The basis is *updated*, $x \leftarrow x + \alpha \Delta x$, and all steps are repeated.

4 Simplex Facts

A few observations about typical simplex implementations can now be made.

- Variables may have any bounds, including ones where $\ell_j = u_j$.

- x includes a full set of slack variables. There are no “artificial variables”.
- Nonbasic variables need not be zero—they may have any value satisfying their bounds. “Nonbasic” means *temporarily frozen at current value*. Admittedly, most nonbasic variables are *equal* to one of their bounds (though not if both bounds are infinite!).
- Simplex methods start from a specified basis and solution (B, N, x_B, x_N) . If B is nonsingular, it lets us satisfy $Ax = b$.
- Free variables are usually in the basis. They are “best” because they always satisfy their (infinite) bounds and therefore won’t generate basis changes.
- (Partial Pricing) If $n \gg m$, the computation of the reduced-gradient vector $z_N = c_N - N^T y$ is expensive. We may compute only part of z_N each iteration and choose any element z_s that is reasonably large and has appropriate sign.
- (Steepest-Edge Pricing) Alternatively, the improvement per unit step is greater if $c^T \Delta x / \|\Delta x\|$ is large (and negative). This means $|z_s|/\omega_s$ should be large, where ω_s is defined as follows for each nonbasic variable x_s :

$$Bv_s = a_s, \quad \omega_s = \|\Delta x\| = \left\| \begin{pmatrix} -\sigma v_s \\ \sigma \end{pmatrix} \right\| = \left\| \begin{pmatrix} v_s \\ -1 \end{pmatrix} \right\|.$$

When the chosen a_s replaces the p th column of B , methods are known for updating all nonbasic $\omega_j^2 = (\|v_j\|_2^2 + 1)$ [4, 2], where $Bv_j = a_j$. It would be too expensive to solve for each v_j . Instead, the systems

$$B^T z = v_s \quad \text{and} \quad B^T w = e_p$$

are solved and then the steepest-edge weights are updated according to

$$\theta_j = \frac{w^T a_j}{v_{ps}}, \quad \bar{\omega}_j = \omega_j - 2\theta_j z^T a_j + \theta_j^2 \omega_s.$$

- We don’t invert B . We have to solve with B and B^T , using some factorization of B and updating the factors when one column of B is replaced by a column of N .
- After some number of updates (typically 50 to 100), the current B is factorized directly and the basic variables x_B are recomputed so that $Ax = b$. Any number of basic variables may prove to be infeasible. Phase 1 iterations are initiated or continued if necessary. (Thus, Phase 1 may be initiated more than once.)
- In general, x is regarded as feasible if $Ax = b$ to high accuracy and

$$\ell - \delta_P \leq x \leq u + \delta_P, \quad \text{or equivalently} \quad \min(x - \ell, u - x) \geq -\delta_P,$$

where δ_P is a small *feasibility tolerance*. Typically $\delta_P = 10^{-6}$, 10^{-7} , or 10^{-8} . It helps prevent Simplex from jumping in and out of Phase 1. It also allows for noisy data that might otherwise be “infeasible”.

- If x is feasible, (x, y, z) is regarded as optimal if nonbasic variables x_j satisfy

$$x_j > \ell_j \quad \text{and} \quad z_j \leq \delta_D \quad \text{or} \quad x_j < u_j \quad \text{and} \quad z_j \geq -\delta_D,$$

where δ_D is a small *optimality tolerance*. Typically $\delta_D = 10^{-6}$, 10^{-7} , or 10^{-8} . (Almost nobody needs 9 or more digits of precision in the optimal objective value.)

- Absolute tolerances of this nature assume that the problem is *well scaled*. This means that the largest elements in A , x , y and z should be of order 1. Automatic scaling of the rows and columns of A can usually achieve this (but wise modelers choose good units from the beginning).

5 Phase 1

Primal Simplex performs a “Phase 1 iteration” whenever some of the current basic variables are infeasible. Phase 1 defines a *new problem every iteration*. The new problem is feasible (the violated bounds on x_B are relaxed to $\pm\infty$) and has a *new objective vector* \tilde{c} that seeks to reduce the *sum of infeasibilities* for the original problem. The elements of \tilde{c} are zero except for those corresponding to infeasible basic variables:

$$\tilde{c}_j = \begin{cases} +1 & \text{if } x_j > u_j + \delta_P, \\ -1 & \text{if } x_j < \ell_j - \delta_P. \end{cases}$$

A benefit is that *any number of infeasibilities* may be removed in a single iteration (not just 0 or 1). Special care is needed in defining the steplength α when the *last* infeasibility is removed.

In many implementations, the constraints $Ax = b$ include the objective vector as a constraint (with infinite bounds on the associated slack variable). For problem LO2 above, if the objective is in the last row of A , the constraint has the form $c^T x - s_m = 0$, where $-\infty \leq s_m \leq \infty$. The system “ $B^T y = c_B$ ” is therefore always of the form

$$\tilde{B}^T \tilde{y} = \begin{pmatrix} B^T & c_B \\ & -1 \end{pmatrix} \begin{pmatrix} y \\ y_m \end{pmatrix} = \begin{pmatrix} \tilde{c}_B \\ \gamma \end{pmatrix},$$

where $\gamma = 0$ in Phase 1 and $\tilde{c}_B = 0$ in Phase 2. Note the flexibility of this arrangement. The same factorization of \tilde{B} may be used for both phases, because only the right-hand side is being “juggled”. A single scan of the basic variables is used to define \tilde{c}_B , and then γ is set to 0 or 1 appropriately (or to -1 if x_B is feasible and the objective should be maximized).

6 Basis Factorization

Simplex and Reduced-Gradient implementations factorize B periodically using some kind of sparse LU procedure—a sparse form of Gaussian elimination: $B = LU$, where L and U are permuted sparse triangular matrices.

Certain updating methods allow the basis factorization to be a “black box”, but some methods update the LU factors directly and therefore affect the way in which L and U are computed and stored. For example, MA28 and MA48 are able to permute B to *block-triangular form* and then factorize each block separately, but this is not suitable for methods that update a single L and U .

The choice of permutations (the “pivot order”) is a delicate balance between sparsity and stability. A *stable* factorization ensures that at least one of the factors (typically L) is well-conditioned. The other factor then tends to reflect the condition of B .

If B is close to singular (especially the very first B), the factorization should be able to signal singularity and report which rows and columns of B appear to be the cause. Such *rank-revealing properties* cannot be guaranteed (there are pathological examples that defeat most practical strategies), but LUSOL (to be discussed later) has options that are increasingly successful at the expense of greater run-time and density of the factors.

If singularity is apparent, the offending basic variables (columns) are replaced by slack variables corresponding to the offending rows. Hence the need for a complete set of slack variables.

Ideally, the numerical values of the rejected basic variables are preserved in order to maintain $Ax = b$ as closely as possible. Hence the wish for nonbasic variables to take any value (within their bounds).

7 Basis Updating

Suppose a column \bar{a} replaces the p th column of B to give the next basis \bar{B} . Thus,

$$\bar{B} = B + (\bar{a} - a)e_p^T,$$

where e_p is the p th column of the identity. Given some kind of factorization of B , we need to obtain a factorization of \bar{B} . We study several methods that have varying degrees of stability, efficiency, and ease of implementation:

The Product-Form update
 The Bartels-Golub update
 The Forrest-Tomlin update
 The block-LU update.

8 Crash Procedures

This strange term refers to choosing an initial basis matrix B when none is provided by the user. Several approaches are described in [5, 8].

Triangular crash procedures make several passes through the columns of A , using the sparsity pattern to find a set of columns that form a permuted triangle. Each effective diagonal element should be reasonably large compared to other elements in the same column (say, no smaller than 0.1 times the largest nonzero in the column). The aim is to find an initial basis that is nonsingular and not too ill-conditioned.

Other procedures try to do more by using the simplex method itself. For example, a number of “cheap iterations” can be performed in which the Pricing step does very little searching to find a reduced gradient z_j that has suitable sign and is larger than a modest threshold.

9 Scaling

Unfortunately, some problems are formulated without concern about the size of the elements of A . Row and column scaling may be needed to maintain robustness of the simplex solver. Scaling often reduces the number of simplex iterations (although there are notorious cases that solve much more quickly without such treatment).

Nominally, the nonzero elements of A should be of order 1. We mean that there is no reason to have very large entries, but it is hard to know if *small* entries are small through poor scaling or because they intentionally have little effect on the associated constraints. As suggested by Robert Fourer (private communication, 1979), a good general approach is to make several passes through the columns and rows of A , computing the geometric mean of the nonzeros in each (scaled) column or row and using that as the new scale factor. (The geometric mean is approximated by the square root of the product of the largest and smallest nonzero in each column.)

On problem `pilotja` [9, 7], the scaling procedure in MINOS proceeds as follows:

	Min elem	Max elem	Max col ratio
After 0	2.00E-06	5.85E+06	189805175.80
After 1	1.27E-03	1.00E+02	78436.79
After 2	1.59E-02	6.30E+01	3962.95
After 3	2.02E-02	4.94E+01	2441.93
After 4	2.11E-02	4.74E+01	2247.57

The `Min` and `Max` columns refer to the extreme nonzero $|a_{ij}|$ in all of A . The last column `Max col ratio` means $\max_j \{\max_i |a_{ij}| / \min_i |a_{ij}|\}$ among nonzero $|a_{ij}|$. The MATLAB function `gmscale.m` implements the same scaling procedure.

10 Anti-degeneracy

The steplength procedure may return an embarrassing value—namely $\alpha = 0$. The objective function then does not improve. If a sequence of zero steplengths repeats indefinitely, the simplex method is said to be *cycling* (and convergence does not occur).

Indeed it is common for *many* basic variables to be essentially equal to one of their bounds. The current solution is then called *primal degenerate*. The search direction Δx might want to move several basic variables outside their bounds, and care is needed to ensure that $|\Delta x_r|$ is sufficiently large compared to other possible choices.

A naive but effective approach is to make small perturbations to the problem data (typically to the vector b) to reduce the probability that any basic variable would be on its bound.

A more elaborate approach was introduced as the EXPAND procedure in [3]. The key idea is to relax all of the bounds $\ell \leq x \leq u$ by a tiny amount every simplex iteration. In particular, when a variable leaves the basis, it may become nonbasic at a *slightly infeasible value* (no greater than δ_p). Over a period of K iterations, where $K = 10,000$ say, there may be many tiny steplengths α (at which points the simplex method is effectively *stalling*), but it is unlikely that any sequence of iterations will recur.

Unfortunately, the bound relaxation cannot continue indefinitely. After K iterations, infeasible nonbasic variables must be put on their bounds, the basic variables recomputed, and the EXPAND process restarted. There is a danger of the same K iterations repeating, and indeed, Hall and McKinnon [6] have constructed examples that cause EXPAND to cycle. Nevertheless, EXPAND continues to be used successfully within MINOS and SQOPT.

Many things contribute to efficient implementations of the simplex method. We have touched on some. Bixby [1] describes the remarkable progress made during 1987–2001. A significant item is taking advantage of *sparse* right-hand sides v and w in (2).

References

- [1] R. E. Bixby. Solving real-world linear programs: a decade and more of progress. *Operations Research*, 50(1):3–15, 2002.
- [2] J. J. Forrest and D. Goldfarb. Steepest-edge simplex algorithms for linear programming. *Math. Program.*, 57:341–374, 1992.
- [3] P. E. Gill, W. Murray, M. A. Saunders, and M. H. Wright. A practical anti-cycling procedure for linear and nonlinear programming. *Math. Program.*, 45:437–474, 1989.
- [4] D. Goldfarb and J. K. Reid. A practicable steepest-edge simplex algorithm. *Math. Program.*, 12:361–371, 1977 (where the title should have read *practical!*).
- [5] N. I. M. Gould and J. K. Reid. New crash procedures for large-scale systems of linear constraints. *Math. Program.*, 45:475–501, 1989.
- [6] J. A. J. Hall and K. I. M. McKinnon. The simplest examples where the simplex method cycles and conditions where EXPAND fails to prevent cycling. *Math. Program.*, 100(1):133–150, 2004.
- [7] LPnetlib collection of LP problems in MATLAB format. <http://www.cise.ufl.edu/research/sparse/mat/LPnetlib>.
- [8] I. Maros and G. Mitra. Strategies for creating advanced bases for large-scale linear programming problems. *INFORMS J. on Computing*, 10:248–260, 1998.
- [9] Netlib collection of LP problems in MPS format. <http://www.netlib.org/lp/data>.