



A Simple New Algorithm for Quadratic Programming with Applications in Statistics

Mary C. Meyer

To cite this article: Mary C. Meyer (2013) A Simple New Algorithm for Quadratic Programming with Applications in Statistics, Communications in Statistics - Simulation and Computation, 42:5, 1126-1139, DOI: [10.1080/03610918.2012.659820](https://doi.org/10.1080/03610918.2012.659820)

To link to this article: <https://doi.org/10.1080/03610918.2012.659820>



Published online: 12 Dec 2012.



Submit your article to this journal [↗](#)



Article views: 503



View related articles [↗](#)



Citing articles: 6 View citing articles [↗](#)

A Simple New Algorithm for Quadratic Programming with Applications in Statistics

MARY C. MEYER

Department of Statistics, Colorado State University, Fort Collins, Colorado, USA

Problems involving estimation and inference under linear inequality constraints arise often in statistical modeling. In this article, we propose an algorithm to solve the quadratic programming problem of minimizing $\psi(\theta) = \theta'Q\theta - 2c'\theta$ for positive definite Q , where θ is constrained to be in a closed polyhedral convex cone $C = \{\theta : A\theta \geq d\}$, and the $m \times n$ matrix A is not necessarily full row rank. The three-step algorithm is intuitive and easy to code. Code is provided in the R programming language.

Keywords Cone projection; Constrained parameter estimation; Dual algorithm; Inequality constraints; Polar cone; Primal-dual methods; Restricted least squares; Shape-restricted regression

Mathematics Subject Classification Primary 65K05; Secondary 62G05

1. Motivation and Background

Let Q be a positive definite $n \times n$ matrix, let c be a vector in \mathbb{R}^n , let A_0 be an $m \times n$ irreducible matrix, and let d be a vector in the column space of A . The term “irreducible” means “nonredundant” and will be defined formally in Sec. 2. The quadratic programming problem

$$\text{find } \hat{\theta} \text{ to minimize } \theta'Q\theta - 2c'\theta \text{ subject to } A_0\theta \geq d \quad (1)$$

may be readily transformed into a “cone projection” problem by considering the Cholesky decomposition $U'U = Q$ and finding θ_0 such that $A_0\theta_0 = d$. Then define $\phi = U(\theta - \theta_0)$, $z = (U^{-1})'(c - Q\theta_0)$, and $A = A_0U^{-1}$ to get

$$\text{find } \hat{\phi} \text{ to minimize } \|z - \phi\|^2 \text{ subject to } A\phi \geq 0, \quad (2)$$

hence $\hat{\theta} = U^{-1}\hat{\phi} + \theta_0$. The set

$$C = \{\phi \in \mathbb{R}^n : A\phi \geq 0\} \quad (3)$$

is easily observed to be a polyhedral convex cone in \mathbb{R}^n , as each row of A defines a half-space, and C is the intersection of these half-spaces. Because the convex objective function is to be minimized over a convex set, a unique solution $\hat{\phi}$ exists.

Received December 3, 2009; Accepted January 10, 2012

Address correspondence to Mary C. Meyer, Department of Statistics, Colorado State University, 102 Statistics Building, Fort Collins, CO 80523-1877, USA; E-mail: meyer@stat.colostate.edu

An early algorithm for cone projection was proposed by Dykstra (1983). Using that \mathcal{C} is the intersection of half-spaces, it sequentially projects a residual vector onto the half-spaces, updating the residual vector at each step, and then repeats the process until convergence. The interior point algorithm for minimizing a quadratic function over a convex set is a gradient-based algorithm, first proposed by Karmarkar (1984) for linear programming. For more details, see Fang and Puthenpura (1993, chaps. 9 and 10). Fraser and Massam (1989) developed the mixed primal-dual bases algorithm for cone projections and applied it to concave nonparametric regression.

The interior point algorithm and the sequential projection algorithm are considered to converge in “infinitely many” steps, because the true solution is approached asymptotically and reached within a user-defined tolerance. In contrast, other algorithms that exploit the edges of the cone are guaranteed to produce the solution in a finite number of steps. These include the “nearest point” algorithm of Wilhelmsen (1976), the dual method of Goldfarb and Idnani (1983), the mixed primal-dual bases (MPDB) algorithm of Fraser and Massam (1989), and the nonnegative least squares (NNLS) algorithm of Lawson and Hanson (1974). The Goldfarb and Idnani algorithm is used in the R code `quadprog`; see <http://CRAN.R-project.org/package=quadprog>. The “hinge” algorithm proposed here can be summarized in three simple steps, and from a statistician’s point of view, it is more intuitive than other methods. The derivations and proofs use results from linear models rather than optimization theory.

An example of an application in statistics is the least-squares regression problem with linear inequality constraints on the coefficients. Let \mathbf{X} be an $n \times k$ fixed design matrix and $\boldsymbol{\beta}$ be a k -dimensional parameter vector, and consider the model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad (4)$$

where $E(\boldsymbol{\epsilon}) = \mathbf{0}$ and $\text{cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}$. Suppose that it is reasonable to impose constraints on the parameter vector in the form $\mathbf{A}_0 \boldsymbol{\beta} \geq \mathbf{d}$, where \mathbf{A}_0 is an $m \times k$ irreducible matrix. The problem of minimizing $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$ subject to the constraints is a quadratic programming problem with $\mathbf{Q} = \mathbf{X}'\mathbf{X}$ and $\mathbf{c} = \mathbf{X}'\mathbf{y}$. Methods for minimizing $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$ subject to various linear inequality and equality constraints were proposed by Judge and Takayama (1966); they used a simplex algorithm to obtain the optimal $\boldsymbol{\beta}$. Liew (1976) considered constraints $\mathbf{A}\boldsymbol{\beta} \geq \mathbf{d}$, and used the Dantzig–Cottle “principal pivoting” algorithm and provided an approximate variance-covariance matrix for $\hat{\boldsymbol{\beta}}$. Hawkins (1994) considered the problem of fitting monotonic polynomials to data using a primal-dual algorithm to adjust the coefficients iteratively. An exposition of estimation and testing under linear inequality constraints can be found in chapter 21 of Gourieroux and Monfort (1995), chapter 4 of Silvapulle and Sen (2005), and chapter 2 of Robertson et al. (1988).

Another traditional use of constrained estimation and inference is for nonparametric regression with shape restrictions. Suppose that interest is in fitting a scatterplot generated from

$$y_i = f(x_i) + \epsilon_i, \quad i = 1, \dots, n, \quad (5)$$

where a parametric form is unknown for f , but it may be assumed that f is nondecreasing. Defining $\theta_i = f(x_i)$, the problem is to find $\boldsymbol{\theta} \in \mathbb{R}^n$ to minimize $\|\mathbf{y} - \boldsymbol{\theta}\|^2$ given $\mathbf{A}\boldsymbol{\theta} \geq \mathbf{0}$, where the nonzero elements of the $(n-1) \times n$ matrix \mathbf{A} are $A_{i,i} = -1$ and $A_{i,i+1} = 1$, $i = 1, \dots, n-1$. This model was first considered by Brunk (1955) who proposed the pooled adjacent violators algorithm (PAVA) to obtain the solution. For estimation, assuming f is convex, the nonzero elements of the $(n-2) \times n$ constraint matrix are $A_{i,i} = t_{i+2} - t_{i+1}$,

$A_{i,i+1} = t_i - t_{i+2}$, and $A_{i,i+2} = t_{i+1} - t_i$. Fraser and Massam (1989) proposed the MPDB algorithm to solve the concave regression problem.

This article introduces a simple, three-step algorithm for cone projection that is fast and intuitive. The code, shown in the R programming language in Appendix A, is only a couple dozen lines. The speed is important for applications in iteratively reweighted cone projection algorithms, such as for generalized constrained regression, and in estimating the mixing distribution used for the tests of $H_0 : \mathbf{A}\boldsymbol{\theta} = \mathbf{b}$ versus $H_0 : \mathbf{A}\boldsymbol{\theta} > \mathbf{b}$ [see Gouriéroux et al. (1982), Raubertas et al. (1986), or Silvapulle and Sen (2005, chap. 4), or Robertson et al. (1988)]. The R package `ic.infer` supports inequality-constrained estimation and testing, and is described in Grömping (2010). The rest of the article is organized as follows. Some necessary theory about cones and projections onto cones is outlined in the next section, which is used to explain the steps of the algorithm and also used in the convergence results (Appendix B). This is meant to be accessible to people with some background in linear spaces, who are not necessarily familiar with convex optimization theory. For more details and proofs of the claims, see Silvapulle and Sen (2005, chap. 3), and van Eeden (2006) also provides theory for constrained estimation. The proposed cone projection (“hinge”) algorithm is described in Sec. 3, and some useful applications in statistics are given in Sec. 4.

2. The Constraint Cone and Its Polar Cone

The convex polyhedral cone (3) is considered for an $m \times n$ irreducible constraint matrix \mathbf{A} , where “irreducible” means that no row of \mathbf{A} is a positive linear combination of other rows, and there is not a positive linear combination of rows of \mathbf{A} that equals the zero vector. Note that if \mathbf{A} is full row rank, then it is irreducible. If a row is a positive linear combination of other rows, it can be removed without affecting the problem, and if the origin can be written as a positive linear combination of rows, then there is an implicit equality constraint in the matrix \mathbf{A} , which can be dealt with separately. Silvapulle and Sen (2005) use the term “tight” to describe this nonredundancy of \mathbf{A} .

Let m_1 be the dimension of the space spanned by the rows of \mathbf{A} . If m_1 is less than n , then the cone contains a linear space V of dimension $n - m_1$; this is the null space of \mathbf{A} . Let $\Omega = \mathcal{C} \cap V^\perp$, where V^\perp is the linear space orthogonal to V . Then Ω is a polyhedral convex cone that does not contain a linear space of dimension one or greater, and sits in an m_1 -dimensional subspace of \mathbb{R}^n . The “edges” or “generators” of Ω are vectors $\boldsymbol{\delta}^1, \dots, \boldsymbol{\delta}^{m_2}$ in Ω such that

$$\Omega = \left\{ \boldsymbol{\phi} \in \mathbb{R}^n : \boldsymbol{\phi} = \sum_{j=1}^{m_2} b_j \boldsymbol{\delta}^j, \text{ where } b_j \geq 0, \ j = 1, \dots, m_2 \right\},$$

and hence

$$\mathcal{C} = \left\{ \boldsymbol{\phi} \in \mathbb{R}^n : \boldsymbol{\phi} = \mathbf{v} + \sum_{j=1}^{m_2} b_j \boldsymbol{\delta}^j, \text{ where } b_j \geq 0, \ j = 1, \dots, m_2 \text{ and } \mathbf{v} \in V \right\}. \quad (6)$$

If $m_1 = m$, then $m_2 = m$ and the edges of Ω are the columns of the matrix $\Delta = \mathbf{A}'(\mathbf{A}\mathbf{A}')^{-1}$. For the case of more constraints than dimensions ($m > m_1$), the edges of the cone are obtained as follows (see Meyer, 1999, for proof). Define $\boldsymbol{\gamma}^1, \dots, \boldsymbol{\gamma}^m$ to be the rows of $-\mathbf{A}$. Suppose $J \subseteq \{1, \dots, m\}$ and let $S = \mathcal{L}(\{\boldsymbol{\gamma}^j, j \in J\})$, where \mathcal{L} denotes “the space spanned by.” If $\dim(S) = m_1 - 1$, then $S^\perp \cap V^\perp$ is a line through the origin containing the vectors $\boldsymbol{\delta}$ and $-\boldsymbol{\delta}$, say, where $\boldsymbol{\delta} \perp \boldsymbol{\gamma}^j$ for all $j \in J$. If $\langle \boldsymbol{\delta}, \boldsymbol{\gamma}^i \rangle \leq 0$ for all $i \notin J$, then $\boldsymbol{\delta}$ is an edge of

Ω . Conversely, all edges are of this form. In the case of more constraints than dimensions, the number m_2 of edges may be considerably larger than m .

Because V is orthogonal to Ω , $\hat{\phi}$ is the sum of the projections of y onto V and Ω so that $\hat{\phi}$ can be written as $v + \Delta b$, where $v \in V$, the columns of Δ are $\delta_1, \dots, \delta_m$, and $b \geq 0$. Note that this representation is not necessarily unique if there are more constraints than dimensions, but $\hat{\phi}$ is unique. Necessary and sufficient conditions for $\hat{\phi} \in \mathcal{C}$ to minimize $\|y - \phi\|^2$ are $\langle y - \hat{\phi}, \hat{\phi} \rangle = 0$ and $\langle y - \hat{\phi}, \phi \rangle \leq 0$, for all $\phi \in \mathcal{C}$; because the δ vectors are generators of the cone, these conditions may be written as

$$\langle y - \hat{\phi}, \hat{\phi} \rangle = 0; \quad \langle y - \hat{\phi}, v \rangle \leq 0 \quad \text{for all } v \in V; \quad \text{and} \quad \langle y - \hat{\phi}, \delta^j \rangle \leq 0 \quad \text{for } j = 1, \dots, m_2. \quad (7)$$

Proposition 2.1. *Let $\hat{\phi}$ be the unique minimizer of $\|y - \phi\|^2$ over $\phi \in \mathcal{C}$, and write $\hat{\phi} = v + \Delta b$ for $b \geq 0$ and $v \in V$. Let $J \subseteq \{1, \dots, m\}$ index the nonzero elements of b ; that is, $j \in J$ if $b_j > 0$. Then $\hat{\phi}$ is the projection of y onto the linear space spanned by δ^j , $j \in J$, and the basis vectors for V .*

Hence, the projection of a vector onto a polyhedral convex cone is the ordinary least-squares projection onto the linear space spanned by a basis for V and a subset of cone edges indexed by a subset J of $\{1, \dots, m\}$. The projection lands on a *face* of the cone; the faces can be defined as

$$\mathcal{F}_J = \left\{ \phi \in \mathbb{R}^n : \phi = v + \sum_{j \in J} b_j \delta_j, \quad b_j > 0, \quad v \in V \right\}.$$

Note that $\mathcal{F}_\emptyset = V$ (where \emptyset denotes the empty set), and the interior of the constraint cone is the face where $J = \{1, \dots, m\}$.

The set of points whose projection onto \mathcal{C} lands in V (or whose projection onto Ω is the origin) is also called the *polar cone*. The edges of the polar cone are $\gamma_1, \dots, \gamma_m$, that is, the rows of $-A$, and the polar cone may be seen to be the collection of vectors in V^\perp that make obtuse angles with all vectors in \mathcal{C} . The next proposition implies that the polar cone is analogous to the orthogonal space in the linear model.

Proposition 2.2. *Let $\hat{\rho}$ be the projection of y onto the polar cone, and let $\hat{\phi}$ be the projection of y onto the constraint cone. Then $\hat{\phi} + \hat{\rho} = y$.*

Therefore, the projection onto the constraint cone may alternatively be found via projection onto the polar cone, and that the polar cone is analogous to the orthogonal space in the linear model.

3. The Hinge Algorithm

The hinge algorithm for cone projection uses either the constraint cone edges $\delta_1, \dots, \delta_{m_2}$ or the polar cone edges $\gamma_1, \dots, \gamma_m$, so the cone can be specified either as (3) or as (6). In the following, the δ vectors are used, but the algorithm is identical if the γ vectors are substituted, except that the residual of the projection onto the constraint cone is obtained. This is useful for problems with more constraints than dimensions, as the number of constraint cone edges can be large and finding them may be computationally intensive.

Proposition 1 tells us that the minimizer of (2) subject to $A\phi \geq \mathbf{0}$ with irreducible $m \times n$ A can be solved by finding $J \subseteq \{1, \dots, m\}$ where $\hat{\phi}$ lands on \mathcal{F}_J . The hinge algorithm arrives at the appropriate set J through a series of guesses J_k . At a typical iteration, the current estimate ϕ^k is the least-squares regression of z on the space spanned by δ^j , for $j \in J_k$. (The δ^j , $j \in J$, were originally called “hinges” since in the convex regression problem, for which the algorithm was initially devised, the $(t_j, \hat{\phi}_j)$, $j \in J$, are the points at which the line segments in the piecewise linear fit change slope, and the bends are allowed to go only one way.) The initial guess J_0 can be any subset of $\{1, \dots, m_2\}$ for which the corresponding δ^j , $j \in J$, form a linearly independent set. The hinge algorithm can be summarized in three steps. At the k th iteration:

1. Project z onto the linear space spanned by $\{\delta^j, j \in J_k\}$, to get $\phi^k = \sum_{j \in J_k} b_j^{(k)} \delta_j$.
2. Check to see if ϕ^k satisfies the constraints, that is, if all $b_j^{(k)}$ are nonnegative:
 - If yes, go to step 3.
 - If no, choose j for which $b_j^{(k)}$ is smallest, and remove it from the set; go to step 1.
3. Compute $\langle y - \phi^k, \delta^j \rangle$ for each $j \notin J_k$. If these are all nonpositive, then stop. If not, choose j for which this inner product is largest, add it to the set, and go to step 1.

Intuitively, at each stage, the new edge is added where it is “most needed,” and other edges are removed if the new fit does not satisfy the constraints. Although the δ^j , $j = 1, \dots, m_2$, might not form a linearly independent set, at each step, the δ^j , $j \in J_k$ are linearly independent. Suppose that at step 2, we have that J_k defines a linearly independent set of δ^j vectors. For all vectors δ^j such that the indices $J_k \cup \{j\}$ do not define a linearly independent set, we have $\langle y - \phi^k, \delta^j \rangle = 0$, so these j are not added. Since the stopping criteria are defined by (7), it is clear that if the algorithm ends, it gives the correct solution. The only thing that requires proof is that the algorithm does end, that is, it does not produce the same set of edges twice, which would result in an infinite loop. The proofs are deferred to Appendix B.

4. Applications in Statistics

4.1. Constrained Least-Squares Regression

Consider the model (4) where it is known that $A\beta \geq \mathbf{0}$ for irreducible A . Let $\hat{\beta}$ be the constrained least-squares estimator for β , and let $\tilde{\beta}$ be the unconstrained estimator $(X'X)^{-1}X'y$. The conditions (7) lead to $(y - X\hat{\beta})'X\hat{\beta} = 0$ and $(y - X\hat{\beta})'X\beta \leq 0$ for all β such that $A\beta \geq \mathbf{0}$. Now

$$\|X\tilde{\beta} - X\beta\|^2 = \|X\tilde{\beta} - X\hat{\beta}\|^2 + \|X\hat{\beta} - X\beta\|^2 + 2(X\tilde{\beta} - X\hat{\beta})'(X\hat{\beta} - X\beta).$$

The last term is

$$(X\tilde{\beta} - X\hat{\beta})'(X\hat{\beta} - X\beta) = (y - X\hat{\beta})'(X\hat{\beta} - X\beta) - (y - X\tilde{\beta})'(X\hat{\beta} - X\beta),$$

where the second term on the right is zero by principles of ordinary least-squares regression, and the first is positive. Hence, if the true β satisfies the constraints, the constrained estimate of $E(y)$ is closer to the truth than the unconstrained estimate, with equality only if the two estimates coincide.

A well-known example is the simple linear regression model $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$, where the slope β_1 must be nonnegative. The constrained least-squares estimate is easy to obtain: it is the ordinary least-squares estimate unless that does not provide a positive slope,

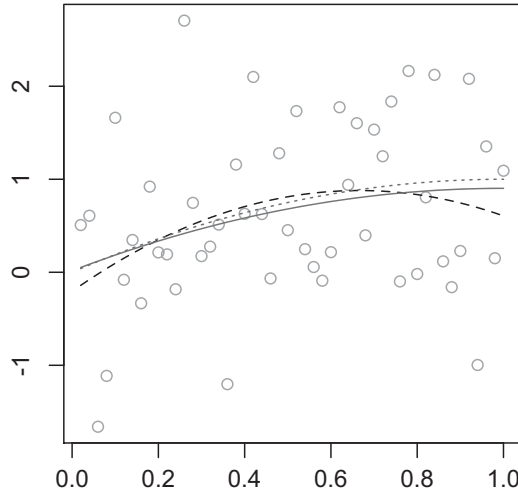


Figure 1. Quadratic fits to a scatterplot simulated using $f(x) = 1 - (1 - x)^2$ (shown as the dotted curve) and $n = 50$ observations. The dashed curve is the unconstrained fit, and the solid curve is constrained to be nondecreasing.

in which case $\hat{\beta}_1 = 0$ and $\hat{\beta}_0 = \bar{y}$. Now suppose that the researcher wishes to allow some curvature in the regression function, so the quadratic model $y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \epsilon_i$ is considered, retaining the assumption that the expected value of y is nondecreasing in x . The constraint is $\beta_1 + 2\beta_2 x \geq 0$ over the range of the data, so if $x \in [0, 1]$, these can be written as $A\beta \geq \mathbf{0}$ with

$$A = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 1 & 2 \end{pmatrix},$$

and $\hat{\beta}$ is obtained by using the hinge algorithm. To illustrate, the scatterplot of $n = 50$ points shown in Fig. 1 was generated using equally spaced x_i , the regression function $f(x_i) = 1 - (1 - x_i)^2$, and independent standard normal errors. Least-squares quadratic fits are shown, where the dashed curve represents the unconstrained fit and the solid curve is constrained to be monotone. The dotted curve is the true regression function.

Another important example is the “warped plane” model $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{1i} x_{2i} + \epsilon_i$, where $E(y)$ is constrained to be increasing in both variables. The constraints are $\beta_1 + \beta_3 x_{2i} \geq 0$, for all x_{2i} , and $\beta_2 + \beta_3 x_{1i} \geq 0$, for all x_{1i} . If both predictors are confined to the unit interval, this provides four constraints on three parameters. Constrained and unconstrained fits are shown in Fig. 2 for data simulated from the surface $f(x_1, x_2) = x_1 x_2$, unit error variance, and $n = 100$, with the predictor values forming a grid in the unit intervals.

4.2. Monotone Regression

There is an elegant closed-form solution for the monotone regression problem provided by Brunk (1955):

$$\hat{\theta}_i = \min_{v \geq i} \max_{u \leq i} \frac{1}{v - u + 1} \sum_{j=u}^v y_j.$$

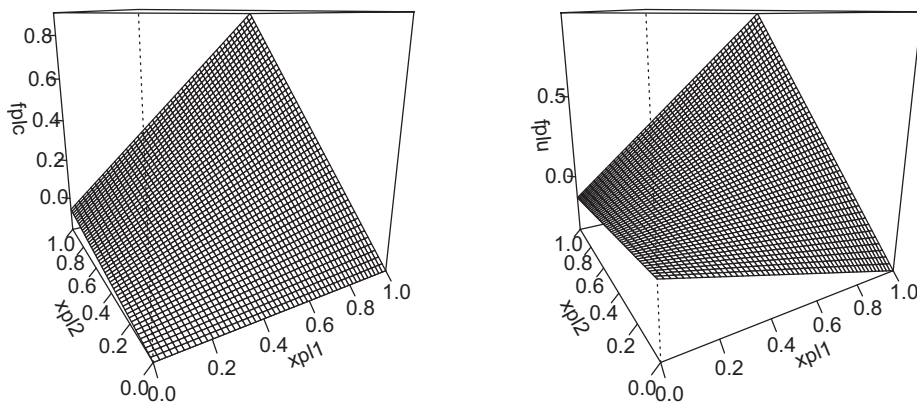


Figure 2. (a) Warped plane fit to dataset $y_i = x_1 x_2 + \epsilon_i$, where the expected value of y is constrained to be increasing in both predictors. (b) Unconstrained fit.

PAVA is a well-known method for finding this solution; see Silvapulle and Sen (2005, p. 47) for more details. The PAVA algorithm begins by setting $\theta = y$, $l = 1$, and $u = 2$. Then the following steps are performed in a loop: if $\theta_l \leq \theta_u$, then $l = u$, $u = u + 1$. Else, the θ_i , $i = l, \dots, u$, are set to the average of the y_i , for $i = l, \dots, u$, and l and u are updated as follows. If $l > 1$, then $l = l - 1$, but if $l = 1$ then $u = u + 1$. The loop ends when $u > n$.

The monotone regression problem is an interesting application of the hinge algorithm when the edges of the constraint cone are used. During the implementation, no hinge indices are removed from the sets J_k at any iteration, so the number of iterations is the number of jumps in the monotone regression estimator. The δ vectors can be written as $\delta_i^j = j - n$, for $i = 1, \dots, j$, and $\delta_i^j = j$, for $i = j + 1, \dots, n$, for $j = 1, \dots, m$, and V is spanned by the one vector. Suppose at some iteration, we have $J_k = \{j_1, \dots, j_p\}$, where the elements of J_k are ordered so that $j_1 < \dots < j_p$. Then it is easily observed that for $1 \leq l \leq p$,

$$\bar{y} + b_{j_1} + \dots + b_{j_l} = \frac{1}{j_{l+1} - j_k} \sum_{i=j_l}^{j_{l+1}-1} y_i.$$

It can be shown that the first hinge index, say l , is chosen so that for any $j_1 < l$ and $j_2 > l$,

$$\frac{1}{l - j_1} \sum_{i=j_1}^{l-1} y_i \leq \frac{1}{j_2 - l} \sum_{i=l}^{j_2-1} y_i.$$

This means that the coefficient on δ^l remains positive after any subsequent additions of edges, and a similar argument can be applied to coefficients of those edges, so all coefficients remain positive throughout the implementation of the algorithm, and their indices are never removed from the sets J_k . From any J , the projection of y onto the face \mathcal{F}_J is easily found by averaging the y values between $j \in J$.

4.3. Shape-Restricted Regression Splines

The monotone regression estimator is a step function, and the convex regression estimator is piecewise linear—neither is satisfactory if f known to be smooth. For model (5), a

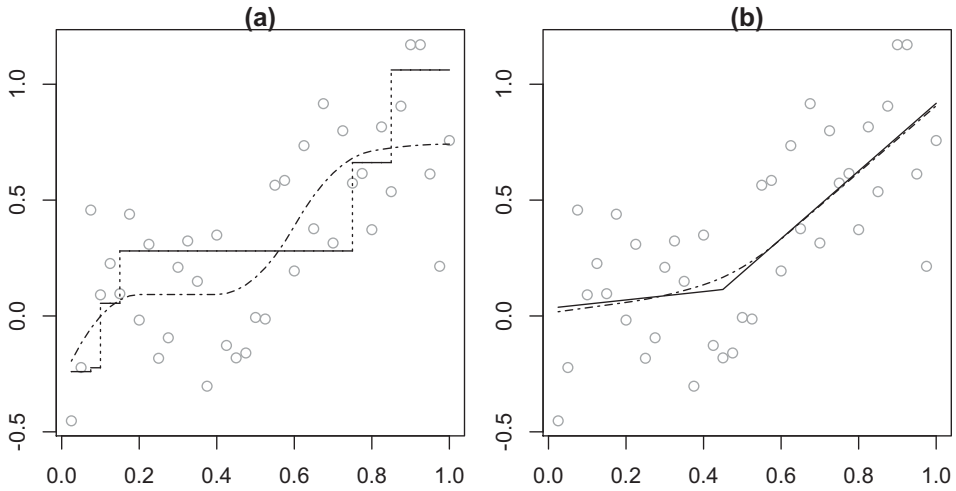


Figure 3. (a) Unsmoothed least-squares monotone fits a dataset (step function) and quadratic monotone regression spline fit (dot–dash curve). (b) Least-squares convex fit (solid) and cubic convex spline (dot–dash). Both spline fits use three interior knots.

method for estimating f with smoothed monotone regression using I -splines was given by Ramsay (1988) using regression splines, and Meyer (2008) extended the method to convex restrictions using C -splines. An alternative method for either is outlined here using the more familiar B -splines (see de Boor, 2001, for details). Given a set of k distinct knots over the range of the x values, a set of m piecewise degree- p spline basis functions are to be defined that span the space of such piecewise polynomials. The spline basis vectors contain the values of these functions at the observed x_i ; let the columns of the matrix \mathbf{X} contain the basis vectors. The unconstrained fit is obtained by minimizing $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$ over $\boldsymbol{\beta} \in \mathbb{R}^m$.

A quadratic ($p = 2$) spline function is increasing if and only if it is increasing at the knots; hence, if A_{ij} is the slope of the j th spline basis function at the i th knot, the monotone spline estimator minimizes the sum of squared residuals subject to $\mathbf{A}\boldsymbol{\beta} \geq \mathbf{0}$. Similarly, a cubic spline function is convex if and only if it is convex at the knots, and the cubic B -spline basis functions may be used with A_{ij} equal to the second derivative of the j th spline basis function at the i th knot. Examples of constrained fits to a scatterplot are shown in Fig. 3. In Fig. 3(a), the unsmoothed monotone regression is shown as the piecewise constant function, and the monotone quadratic spline with four interior knots is shown as the dot–dash curve. In Fig. 3(b), the unsmoothed convex regression function is shown as the solid piecewise linear function, and the convex cubic regression spline is the dot–dash curve. For either monotone or convex constraints, we can demonstrate that when the constraints hold, the constrained version of the estimator provides smaller squared error loss, by an argument similar to that for the parametric case (see Meyer, 2008).

4.4. Weighted and Iteratively Reweighted Least Squares

Constrained estimation for the regression models (4) or (5) may be accomplished when the errors have an arbitrary positive definite covariance matrix, through weighted regression. Specifically, for $\mathbf{y} = \boldsymbol{\theta} + \sigma\boldsymbol{\epsilon}$ with $\mathbf{A}\boldsymbol{\theta} \geq \mathbf{0}$, suppose $\text{cov}(\boldsymbol{\epsilon}) = \boldsymbol{\Sigma}$. The model equation is multiplied through by \mathbf{L}^{-1} , where $\mathbf{L}\mathbf{L}' = \boldsymbol{\Sigma}$, to get $\mathbf{z} = \boldsymbol{\phi} + \sigma\boldsymbol{\xi}$ where $\mathbf{z} = \mathbf{L}^{-1}\mathbf{y}$, $\boldsymbol{\phi} =$

$\mathbf{L}^{-1}\boldsymbol{\theta}$, and $\text{cov}(\boldsymbol{\xi})$ is the identity matrix. The new constraint matrix is $\mathbf{A}\mathbf{L}$. The projection is accomplished using the transformed model, and the reverse transformation provides the solution.

More general constrained maximum-likelihood estimation problems may be solved through iteratively reweighted cone projections. For example, we consider a constrained generalized regression model, where the response is a vector \mathbf{y} of independent observations from a distribution written in the form of an exponential family:

$$f(y_i) = \exp\{[y_i\theta_i - b(\theta_i)]/\tau^2 - c(y_i, \tau)\},$$

where the specifications of b and c determine the subfamily of models. See Silvapulle (1994) for examples including one-sided testing problems. Common examples are $b(\theta) = \log(1 + e^\theta)$ for the Bernoulli and $b(\theta) = \exp(\theta)$ for the Poisson model. The log-likelihood function

$$\ell(\boldsymbol{\theta}, \tau) = \sum_{i=1}^n \left[\frac{y_i\theta_i - b(\theta_i)}{\tau^2} \right]$$

is to be maximized over appropriate constraints on $\boldsymbol{\theta}$.

The vector $\boldsymbol{\mu} = E(\mathbf{y})$, can be observed to be $\mu_i = b'(\theta_i)$, and the variance of y_i is $b''(\theta_i)\tau$. The mean vector is related to the predictor variables through a link function $g(\mu_i) = \eta_i$. If $\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}$, constraints of the form $\mathbf{A}\boldsymbol{\beta} \geq \mathbf{d}$ may be considered. Further, if the link function is one-to-one, then monotonicity constraints imposed on the η values also constrain the mean function to be monotone. Constrained regression splines may be used to model $\boldsymbol{\eta}$, under assumptions of monotonicity and smoothness. In any of these applications, the log-likelihood function is maximized over $\boldsymbol{\eta} \in \mathcal{C}$ where \mathcal{C} is of the form (6). The algorithm involves iteratively reweighted cone projections and follows the same ideas for the generalized linear model as found in McCullagh and Nelder (1989). Starting with $\boldsymbol{\eta}^0 \in \mathcal{C}$, the estimate $\boldsymbol{\eta}^{k+1}$ is obtained from $\boldsymbol{\eta}^k$ by constructing \mathbf{z} :

$$z_i = \eta_i^k + (y_i - \mu_i^k) \left(\frac{d\eta}{d\mu} \right)_{ki},$$

where $\mu_i^k = g^{-1}(\eta_i^k)$ and the derivative of the link function is evaluated at μ_i^k . The weighted projection of \mathbf{z} onto \mathcal{C} is obtained with weight vector \mathbf{w} , where $1/w_i^k = (d\eta/d\mu)_k^2 V_k$, and V_k is the variance function evaluated at μ_i^k . This scheme can be shown to converge to the value of $\boldsymbol{\eta}$ that maximizes ℓ over $\boldsymbol{\eta} \in \mathcal{C}$. The proof is similar to that of theorem 1 in Meyer and Woodroffe (2004).

To illustrate, we use a data-set given by Ruppert et al. (2003), concerning incidence of bronchopulmonary dysplasia in $n = 223$ low birth weight infants. Suppose experts believe that the probability of the condition is a smooth, decreasing function of birth weight. The data are shown in Fig. 4 as tick marks at one for infants suffering from the condition and at zero otherwise. The ordinary logistic regression is shown as the dashed curve, but perhaps there is no reason to believe that the assumption of linear log odds holds. The unconstrained spline estimator with five equally spaced knots is shown as the dotted curve; this is unsatisfactory because it violates the assumption of decreasing probability. The constrained spline estimator using the same knots is shown as the solid curve; this shows a steeper descent than the fit assuming linear log odds, and a leveling off above zero. If the only valid assumptions are that the probability of the condition is a smooth, decreasing function of birth weight, the more flexible nonparametric fit may be preferred.

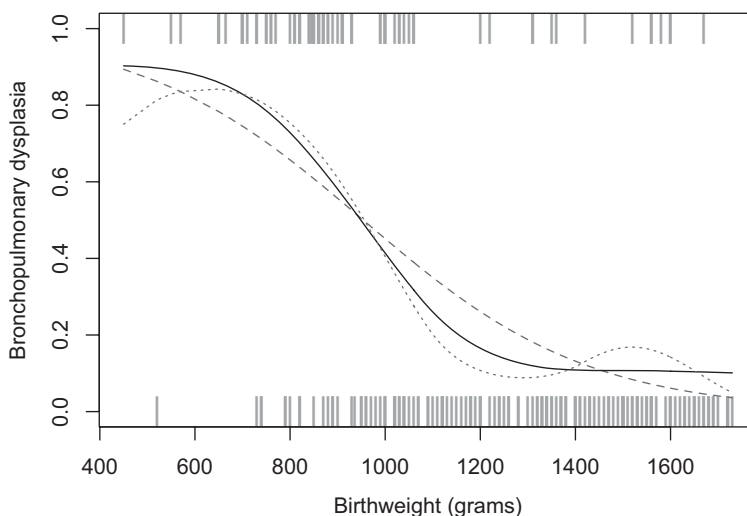


Figure 4. The estimated probability of bronchopulmonary dysplasia in low birth weight infants, as a function of birth weight. The dashed curve is the ordinary logistic regression, the dotted curve is the unconstrained regression spline estimate, and the solid curve is the monotone decreasing regression spline estimate.

5. Discussion

Constrained estimation is useful for both parametric and nonparametric function estimation. When constraints are valid, their use reduces squared error loss of the estimates and increases power of the tests for significance of the predictor. The hinge algorithm is intuitively simple and fast. It can be readily understood using ideas from linear models and does not require a background in optimization theory. The code provided in Appendix A is short and uses only basic matrix operations. More specific code for the other applications listed in this article, including the general quadratic programming problem, constrained parametric regression, constrained spline regression, and the generalized regression models, can be found at <http://www.stat.colostate.edu/~meyer/code.htm>.

Acknowledgment

This work was partially supported by NSF DMS 0905656.

Appendix A: Code

The hinge algorithm to solve the generic cone projection problem of minimizing $\|y - \theta\|^2$ subject to $A\theta \geq 0$ is provided here in R code, where amat is the irreducible constraint matrix A . Specifically, y is projected onto the polar cone to get $\hat{\rho}$, and Proposition 2.2 is used to get $\hat{\theta}$.

```
coneproj=function(y,amat){
  n=length(y);m=length(amat)/n
  sm=1e-8;h=1:m<0;obs=1:m;check=0
  delta=-amat;b2=delta%*%y
  if(max(b2)>sm){
```

```

    i=min(obs[b2==max(b2)])
    h[i]=TRUE
  }else{check=1;theta=1:n*0}
  while(check==0){
    xmat=matrix(delta[h,],ncol=n)
    a=solve(xmat*%*%t(xmat))*%*%xmat*%*%y
    if(min(a)<(-sm)){
      avec=1:m*0;avec[h]=a
      i=min(obs[avec==min(avec)])
      h[i]=FALSE;check=0
    }else{
      check=1
      theta=t(xmat)*%*%a
      b2=delta*%*(y-theta)
      if(max(b2)>sm){
        i=min(obs[b2==max(b2)])
        h[i]=TRUE;check=0}}
  }
  return(y-theta)}

```

Appendix B: Proofs

It is clear that if the hinge algorithm ends, it gives the correct solution. The algorithm ends because it does not choose the same set of edges twice. The sum of squared errors (SSE) $\|z - \hat{\phi}\|^2$ decreases for subsequent iterations with the same number of edges, so the algorithm cannot produce an infinite loop. First we show that the simplest type of loop does not occur.

Proposition B.1. *The algorithm does not remove the edge that it just added.*

Proof. Suppose at the start of some iteration of the algorithm, the set of edges is J_k , and at the end it is $J_{k+1} = J_k \cup \{l\}$, so δ^l is the most recently added edge. The coefficient of δ^l is

$$b_l = \frac{\langle z, \tilde{\delta}^l \rangle}{\|\tilde{\delta}^l\|^2},$$

where $\tilde{\delta}^l$ is the residual from the regression of δ^l on the other regressors $\{\delta^j, j \in J_k\}$. The numerator of the right-hand side is equivalent to $\langle z - \hat{\phi}^k, \delta^l \rangle$, because of orthogonality:

$$\langle z, \tilde{\delta}^l \rangle = \langle z - \hat{\phi}^k, \tilde{\delta}^l \rangle = \langle z - \hat{\phi}^k, \delta^l \rangle > 0.$$

The first equality is because $\hat{\phi}^k \perp \tilde{\delta}^l$ and the second because $\delta^l - \tilde{\delta}^l \perp z - \hat{\phi}^k$.

The idea for proving that the algorithm stops is to show that the SSE at a given iteration with n_h edges is less than that of the last iteration with n_h edges. Suppose that at the beginning of some iteration of the algorithm, we have the solution:

$$\hat{\phi}^B = \sum_{j \in J_B} b_j^B \delta^j,$$

and that this solution satisfies the constraints. Suppose that it is not optimal and the algorithm adds the vector δ^l to the set of regressors. The least-squares fit produced is then

$$\hat{\phi}^M = \sum_{j \in \mathbf{J}_B} b_j^M \delta^j + b_l^M \delta^l.$$

Further, suppose that this $\hat{\phi}^M$ does not satisfy the constraints, so b_l^M , say, is negative. The algorithm will then remove δ^l from the set of regressors in step 2 and go to step 1 to refit the data. The next proposition shows that the new solution

$$\hat{\phi}^N = \sum_{j \in \mathbf{J}_N} b_j^N \delta^j,$$

where $\mathbf{J}_N = \mathbf{J}_B \cup \{l\} - \{i\}$, has $\text{SSE}(\hat{\phi}^N) < \text{SSE}(\hat{\phi}^B)$.

Proposition B.2. *If the algorithm replaces a edge, then the SSE after is less than the SSE before.*

Proof. Let

$$\tilde{\delta}^l = \delta^l - \sum_{j \in \mathbf{J}_B} \alpha_j^l \delta^j,$$

where the second term is the projection of δ^l onto the space spanned by $\{\delta^j, j \in \mathbf{J}_B\}$. Then, $\tilde{\delta}^l \perp \hat{\phi}^B$, so we can write

$$\begin{aligned} \hat{\phi}^M &= \hat{\phi}^B + b_l^M \tilde{\delta}^l \\ &= \sum_{j \in \mathbf{J}_B} (b_j^B - \alpha_j^l b_l^M) \delta^j + b_l^M \delta^l. \end{aligned}$$

We know that $b_i^B > 0$, since $\hat{\phi}^B$ satisfies the constraints, and $b_l^M > 0$, by Proposition B.1. Further,

$$b_i^M = b_i^B - \alpha_i^l b_l^M < 0,$$

so $\alpha_i^l > 0$. Let

$$\phi(x) = \sum_{j \in \mathbf{J}_B} (b_j^B - \alpha_j^l x) \delta^j + x \delta^l.$$

Note that $\phi(0) = \hat{\phi}^B$ and $\phi(b_l^M) = \hat{\phi}^M$. When $x = b_l^M$, the coefficient of δ^l in $\phi(x)$ disappears. Further, $0 < b_i^B / \alpha_i^l < b_l^M$, since $b_i^B - \alpha_i^l 0 > 0$ and $b_i^B - \alpha_i^l b_l^M < 0$. Since b_l^M minimizes $\|z - \phi(x)\|^2$, we have

$$\|z - \phi(0)\|^2 > \|z - \phi\left(\frac{b_i^B}{\alpha_i^l}\right)\|^2 > \|z - \phi(b_l^M)\|^2.$$

Further,

$$\|z - \hat{\phi}^N\|^2 < \|z - \phi\left(\frac{b_i^B}{\alpha_i^l}\right)\|^2$$

since $\hat{\phi}^N$ is the least-squares fit with the same regressors. So, finally,

$$\|z - \hat{\phi}^N\|^2 < \|z - \hat{\phi}^B\|^2.$$

The proof that the sum of squares decreases if *two* edges are replaced is similar in idea but more tedious, and can be found at <http://www.stat.colostate.edu/~meyer/quadprog.htm>.

References

- Brunk, H. D. (1955). Maximum likelihood estimates of monotone parameters. *Annals of Mathematical Statistics* 26(4):607–616.
- de Boor, C. (2001). *A Practical Guide to Splines*. Rev. ed. New York: Springer.
- Dykstra, R. J. (1983). An algorithm for restricted least squares regression. *Journal of the American Statistical Association* 78:837–842.
- Fang, S. C., Puthenpura, S. (1993). *Linear Optimization and Extensions. Theory and Algorithms*. Englewood Cliffs NJ: Prentice Hall.
- Fraser, D. A. S., Massam, H. (1989). A mixed primal-dual bases algorithm for regression under inequality constraints. Application to convex regression. *Scandinavian Journal of Statistics* 16:65–74.
- Goldfarb, D., Idnani, A. (1983). A numerically stable dual method for solving strictly convex quadratic programs. *Mathematical Programming* 27:1–33.
- Gouriéroux, C., Holly, A., Monfort, A. (1982). Likelihood ratio test, Wald test, and Kuhn-Tucker test in linear models with inequality constraints on the regression parameters. *Econometrica* 50(1):63–80.
- Gourieroux, C., Monfort, A. (1995). *Statistics and Econometric Models*. Cambridge: Cambridge University Press.
- Grömping, U. (2010). Inference with linear equality and inequality constraints using R: The package ic.infer. *Journal of Statistical Software* 33(10):1–31.
- Hawkins, D. M. (1994). Fitting monotonic polynomials to data. *Computational Statistics* 9:233–247.
- Judge, G. G., Takayama, T. (1966). Inequality restrictions in regression analysis. *Journal of the American Statistical Association* 61:166–181.
- Karmarkar, N. (1984). A new polynomial time algorithm for linear programming. *Combinatorica* 4:373–395.
- Lawson, C. L., Hanson, R. J. (1974). *Solving Least Squares Problems*. Englewood Cliffs NJ: Prentice Hall.
- Liew, C. K. (1976). Inequality constrained least-squares estimation. *Journal of the American Statistical Association* 71:746–751.
- McCullagh, P., Nelder, J. A. (1989). *Generalized Linear Models*. 2nd ed. New York: Chapman & Hall.
- Meyer, M. C. (1999). An extension of the mixed primal-dual bases algorithm to the case of more constraints than dimensions. *Journal of Statistical Planning and Inference* 81:13–31.
- Meyer, M. C. (2008). Inference using shape-restricted regression splines. *Annals of Applied Statistics* 2(3):1013–1033.
- Meyer, M. C., Woodroffe, M. (2004). Estimation of a unimodal density using shape restrictions. *Canadian Journal of Statistics* 32(1):85–100.
- Ramsay, J. O. (1988). Monotone regression splines in action. *Statistical Science* 3(4):425–461.
- Raubertas, R. F., Lee, C. I. C., Nordheim, E. V. (1986). Hypothesis tests for normal means constrained by linear inequalities. *Communications in Statistics—Theory and Methods* 15(9):2809–2833.

- Robertson, T., Wright, F. T., Dykstra, R. L. (1988). *Order Restricted Statistical Inference*. New York: Wiley.
- Ruppert, D., Wand, M. P., Carroll, R. J. (2003). *Semiparametric Regression (Cambridge Series in Statistical and Probabilistic Mathematics)*. Cambridge: Cambridge University Press.
- Silvapulle, M. J. (1994). On tests against one-sided hypotheses in some generalized linear models. *Biometrics* 50(3):853–858.
- Silvapulle, M. J., Sen, P. K. (2005). *Constrained Statistical Inference*. New York: Wiley.
- van Eeden, C. (2006). *Restricted Parameter Space Estimation Problems*. New York: Springer.
- Wilhelmsen, D. R. (1976). A nearest point algorithm for convex polyhedral cones and applications to positive linear approximation. *Mathematics of Computation* 30(133):48–57.