



# A fast procedure for calculating importance weights in bootstrap sampling

Hua Zhou<sup>a,\*</sup>, Kenneth Lange<sup>b</sup>

<sup>a</sup> Department of Human Genetics, University of California, Los Angeles, CA 90095-1766, United States

<sup>b</sup> Departments of Biomathematics, Human Genetics, and Statistics, University of California, Los Angeles, CA 90095-1766, United States

## ARTICLE INFO

### Article history:

Received 6 June 2009

Received in revised form 21 April 2010

Accepted 21 April 2010

Available online 6 May 2010

### Keywords:

Importance resampling

Bootstrap

Majorization

Quasi-Newton acceleration

## ABSTRACT

Importance sampling is an efficient strategy for reducing the variance of certain bootstrap estimates. It has found wide applications in bootstrap quantile estimation, proportional hazards regression, bootstrap confidence interval estimation, and other problems. Although estimation of the optimal sampling weights is a special case of convex programming, generic optimization methods are frustratingly slow on problems with large numbers of observations. For instance, interior point and adaptive barrier methods must cope with forming, storing, and inverting the Hessian of the objective function. In this paper, we present an efficient procedure for calculating the optimal importance weights and compare its performance to standard optimization methods on a representative data set. The procedure combines several potent ideas for large-scale optimization.

© 2010 Elsevier B.V. All rights reserved.

## 1. Introduction

Because it involves Monte Carlo estimation, the nonparametric bootstrap is an obvious candidate for importance sampling. To our knowledge, Johns (1988) and Davison (1988) first recognized the possibilities in the context of quantile estimation. The general idea is to sample cases with nonuniform weights. If the weights are carefully tuned to a given statistic, then importance sampling can dramatically reduce the variance of the bootstrap sample average estimating the mean of the statistic (Hinkley and Shi, 1989). Bootstrap importance sampling has expanded beyond quantile estimation to include proportional hazards regression, bootstrap confidence interval estimation, and many other applications (Do et al., 2001; Hall, 1992; Johns, 1988; Hu and Su, 2008).

Although estimation of the optimal sampling weights is a constrained optimization problem that yields to standard methods of convex programming, there is still room for improvement, particularly in problems with large numbers of observations. Interior point and adaptive barrier methods incur heavy costs in forming, storing, and inverting the Hessian of the objective function. In the current paper, we present an efficient procedure for calculating the optimal importance weights and compare its performance to standard optimization methods on a representative data set. The procedure combines several potent ideas for large-scale optimization. Briefly these include: (a) approximating the objective function by a quadratic, (b) majorizing the quadratic by a simple quadratic surrogate with parameters separated, (c) reparameterizing the surrogate so its minimization reduces to finding the closest point to a simplex, (d) mapping the simplex solution back to the original parameters, and (e) accelerating the entire scheme by a quasi-Newton improvement for finding the fixed point of a smooth algorithm map. The procedure sounds complicated, but each step is fast and straightforward to implement. On a test example with 1664 observations, our accelerated algorithm surpasses the performance of current standard methods for optimization.

\* Corresponding author.

E-mail addresses: [huazhou@ucla.edu](mailto:huazhou@ucla.edu) (H. Zhou), [klange@ucla.edu](mailto:klange@ucla.edu) (K. Lange).

Section 2 introduces the convex optimization problem defining the importance weights and derives our optimization procedure. Section 3 reviews and generalizes the clever simplex projection algorithm of Michelot. Section 4 summarizes our quasi-Newton acceleration; this scheme is specifically tailored to high-dimensional problems. Section 5 compares our new procedure, both unaccelerated and accelerated, to the standard methods of convex optimization on our sample problem. Finally, our discussion points readers to other applications of the design principles met here for high-dimensional optimization.

## 2. Optimization in importance resampling

In standard bootstrap resampling with  $n$  observations, each observation  $x_i$  is resampled uniformly with probability  $n^{-1}$ . As just argued, it is often helpful to implement importance sampling by assigning different resampling probabilities  $p_i$  to the different observations (Davison, 1988; Do and Hall, 1991; Hesterberg, 1996). For instance, with univariate observations  $(x_1, \dots, x_n)$ , we may want to emphasize one of the tails of the empirical distribution. If we elect to resample nonuniformly according to the multinomial distribution with proportions  $p = (p_1, \dots, p_n)^t$ , then the change of measure equality

$$\begin{aligned} \mathbf{E}[T(\mathbf{x}^*)] &= \mathbf{E}_p \left[ T(\mathbf{x}^*) \frac{\binom{n}{m_1^* \dots m_n^*} \left(\frac{1}{n}\right)^n}{\binom{n}{m_1^* \dots m_n^*} \prod_{i=1}^n p_i^{m_i^*}} \right] \\ &= \mathbf{E}_p \left[ T(\mathbf{x}^*) \prod_{i=1}^n (np_i)^{-m_i^*} \right] \end{aligned}$$

connects the uniform expectation and the importance expectation on the bootstrap resampling space. Here  $m_i^*$  represents the number of times sample point  $x_i$  appears in  $\mathbf{x}^*$ . Thus we can approximate the mean  $\mathbf{E}[T(\mathbf{x}^*)]$  by taking a bootstrap average

$$\frac{1}{B} \sum_{b=1}^B T(\mathbf{x}_b^*) \prod_{i=1}^n (np_i)^{-m_{b,i}^*}$$

with multinomial sampling relative to  $p$ . This Monte Carlo approximation has variance

$$\frac{1}{B} \left\{ \mathbf{E}_p \left[ T(\mathbf{x}^*)^2 \prod_{i=1}^n (np_i)^{-2m_{b,i}^*} \right] - \mathbf{E} [T(\mathbf{x}^*)]^2 \right\},$$

which achieves its minimum with respect to  $p$  when the theoretical second moment  $\mathbf{E}_p \left[ T(\mathbf{x}^*)^2 \prod_{i=1}^n (np_i)^{-2m_{b,i}^*} \right]$  is minimized.

Hall (1992) suggests approximately minimizing the second moment by taking a preliminary uniform bootstrap sample of size  $B_1$ . Based on the preliminary resample, we approximate  $\mathbf{E}_p \left[ T(\mathbf{x}^*)^2 \prod_{i=1}^n (np_i)^{-2m_{b,i}^*} \right]$  by the Monte Carlo average

$$\begin{aligned} s(p) &= \frac{1}{B_1} \sum_{b=1}^{B_1} T(\mathbf{x}_b^*)^2 \prod_{i=1}^n (np_i)^{-2m_{b,i}^*} \prod_{i=1}^n (np_i)^{m_{b,i}^*} \\ &= \frac{1}{B_1} \sum_{b=1}^{B_1} T(\mathbf{x}_b^*)^2 \prod_{i=1}^n (np_i)^{-m_{b,i}^*}. \end{aligned}$$

The function  $s(p)$  serves as a surrogate for  $\mathbf{E}_p \left[ T(\mathbf{x}^*)^2 \prod_{i=1}^n (np_i)^{-2m_{b,i}^*} \right]$ . It is possible to minimize  $s(p)$  on the open unit simplex by standard methods. Unfortunately, Newton’s method is hampered when the  $n$  is large by the necessity of evaluating, storing, and inverting the Hessian matrix at each iteration. This dilemma prompted our quest for a more efficient algorithm for minimizing  $s(p)$ .

Consider the optimization problem

$$\min_p s(p) \quad \text{subject to} \quad \sum_{i=1}^n p_i = 1, \quad p_i \geq \epsilon, \quad 1 \leq i \leq n.$$

Here the lower bound  $\epsilon > 0$  is imposed so that sampling does not entirely neglect some observations. In practice we take  $\epsilon = n^{-2}$  or  $n^{-3}$ . The gradient and second differential (Hessian) of  $s(p)$  are

$$\nabla s(p) = -\frac{1}{B_1} \sum_{b=1}^{B_1} T(\mathbf{x}_b^*)^2 \prod_{j=1}^n (np_j)^{-m_{b,j}^*} \begin{pmatrix} \frac{m_{b,1}^*}{p_1} \\ \vdots \\ \frac{m_{b,n}^*}{p_n} \end{pmatrix}$$

and

$$d^2s(p) = \frac{1}{B_1} \sum_{b=1}^{B_1} T(x_b^*)^2 \prod_{j=1}^n (np_j)^{-m_{bj}^*} \times \left[ \begin{array}{c} \left( \begin{array}{c} m_{b_1}^* \\ p_1 \\ \vdots \\ m_{b_n}^* \\ p_n \end{array} \right) \left( \begin{array}{cc} m_{b_1}^* & m_{b_n}^* \\ p_1 & p_n \end{array} \right) + \left( \begin{array}{c} m_{b_1}^* \\ p_1^2 \\ \vdots \\ m_{b_n}^* \\ p_n^2 \end{array} \right) \end{array} \right].$$

Because  $d^2s(p)$  is positive definite,  $s(p)$  is strictly convex. Evaluation of the gradient and Hessian requires  $O(nB_1)$  and  $O(n^2B_1)$  operations, respectively.

Our first step in minimizing the objective function  $s(p)$  is to approximate it by a quadratic around the current iterate  $p^k$ . According to Taylor's theorem, we have

$$\begin{aligned} s(p) &\approx s(p^k) + \nabla s(p^k)^t (p - p^k) + \frac{1}{2} (p - p^k)^t d^2s(p^k) (p - p^k) \\ &= s(p^k) - \frac{1}{B_1} \sum_b c_b v_b^t (p - p^k) + \frac{1}{2B_1} (p - p^k)^t \sum_b c_b (v_b v_b^t + D_b) (p - p^k) \\ &= r(p | p^k), \end{aligned}$$

where  $c_b = T(x_b^*)^2 \prod_{j=1}^n (np_j)^{-m_{bj}^*}$  and

$$v_b = \left( m_{b_1}^* p_1^{-1}, \dots, m_{b_n}^* p_n^{-1} \right)^t, \quad D_b = \text{diag} \left( m_{b_1}^* p_1^{-2}, \dots, m_{b_n}^* p_n^{-2} \right).$$

Our second step is to majorize the quadratic  $r(p | p^k)$  by a quadratic with parameters separated. If we set  $u = \sum_b c_b v_b$  and  $D = \sum_b c_b (\|v_b\|^2 I_n + D_b)$ , then application of the Cauchy-Schwarz inequality  $\|v_b\|^2 \|w\|^2 \geq (v_b^t w)^2$  yields the inequality

$$r(p | p^k) \leq s(p^k) - \frac{1}{B_1} [u^t + (p^k)^t D] p + \frac{1}{2B_1} p^t D p + c^k = q(p | p^k),$$

where  $c^k$  is an irrelevant constant that does not depend on  $p$ . Because equality holds in this last inequality whenever  $p = p^k$ , the function  $q(p | p^k)$  is said to majorize  $r(p | p^k)$ . The guiding principle of the MM algorithm (Hunter and Lange, 2004; Lange, 2004) is that minimizing  $q(p | p^k)$  drives  $r(p | p^k)$ , and presumably  $s(p)$ , downhill. Thus, we achieve a steady decrease in  $s(p)$ .

Our third step is to transform minimization of  $q(p | p^k)$  into a problem of finding the closest point to a truncated simplex. This step is effected by the reparameterization  $p^* = D^{1/2} p$ , where  $D^{1/2}$  is the matrix square root of  $D$ . Minimization of  $q(p | p^k)$  reduces to minimizing the squared distance

$$\frac{1}{2} \|p^* - (D^{-1/2} u + D^{1/2} p^k)\|^2$$

subject to the constraints  $\mathbf{1}^t D^{-1/2} p^* = 1$  and  $p^* \geq \epsilon D^{1/2} \mathbf{1}$ . Before we discuss how to project  $(D^{-1/2} u + D^{1/2} p^k)$  onto this truncated simplex, let us summarize our overall algorithm in pseudo-code.

---

#### Algorithm 1 Optimal Importance Weights

---

Initialize:  $p = (\frac{1}{n}, \dots, \frac{1}{n})$

**repeat**

$$c_b = T(x_b^*)^2 \prod_{j=1}^n (np_j)^{-m_{bj}^*}, \quad b = 1, \dots, B_1$$

$$v_b = \left( m_{b_1}^* p_1^{-1}, \dots, m_{b_n}^* p_n^{-1} \right)^t, \quad b = 1, \dots, B_1$$

$$D_b = \text{diag} \left( m_{b_1}^* p_1^{-2}, \dots, m_{b_n}^* p_n^{-2} \right), \quad b = 1, \dots, B_1$$

$$u = \sum_b c_b v_b$$

$$D = \sum_b c_b (D_b + \|v_b\|^2 I_n)$$

project  $x = D^{-1/2} u + D^{1/2} p$  onto

$$K = \{y \in \mathbb{R}^n : \sum_i d_i^{-1/2} y_i = 1, y_i \geq d_i^{1/2} \epsilon\}$$

$$p = D^{-1/2} P_K(x)$$

**until** convergence occurs

---

Several remarks are pertinent. (a) Projection onto the truncated simplex can be solved by a slight generalization of an efficient algorithm of Michelot (1986). The details spelled out in the next section show that projection requires at most

$O(n^2)$  operations and usually much fewer in practice. (b) Evaluation of  $u$  and  $D$  requires  $O(nB_1)$  operations. These represent potentially huge gains over Newton’s method if convergence occurs fast enough. Recall that Newton’s method needs  $O(nB_1)$  operations for evaluating the gradient,  $O(n^2B_1)$  operations for evaluating the Hessian matrix, and  $O(n^3)$  for inverting the Hessian matrix. (c) The boundary conditions and linear constraint are incorporated in the algorithm gracefully. (d) A side effect of majorization is the loss of the superlinear convergence enjoyed by Newton’s method. We therefore accelerate convergence by applying a general quasi-Newton scheme for fixed point problems. As discussed in Section 4, this scheme requires little extra computation per iteration and only  $O(n)$  storage. It is particularly attractive for high-dimensional problems. By contrast Newton’s method requires  $O(n^2)$  storage for manipulating the Hessian matrix.

### 3. Michelot algorithm

Michelot (1986) derived an efficient algorithm for projecting a point onto the unit simplex in  $\mathbb{R}^n$ . This algorithm converges in at most  $n$  iterations and often much sooner. We consider a trivial generalization that maps a point  $x \in \mathbb{R}^n$  to the closest point  $P_K(x)$  in the dilated and truncated simplex

$$K = \left\{ y \in \mathbb{R}^n : \sum_{i=1}^n \alpha_i y_i = c, y_i \geq \epsilon_i, 1 \leq i \leq n \right\}, \tag{1}$$

where the  $\alpha_i$  and  $\epsilon_i$  are strictly positive and together satisfy  $\sum_i \alpha_i \epsilon_i \leq c$ . The unit simplex is realized by taking  $c = 1$  and  $\alpha_i = 1$  and  $\epsilon_i = 0$  for all  $i$ . The revised algorithm cycles through the following steps.

---

**Algorithm 2** Michelot Algorithm: Project  $x \in \mathbb{R}^n$  onto the truncated simplex  $K = \{y \in \mathbb{R}^n : \sum_{i=1}^n \alpha_i y_i = c, y_i \geq \epsilon_i \text{ for all } i\}$

---

**repeat**

Project  $x$  onto the hyperplane  $H = \{x : \sum_i \alpha_i x_i = c\}$  via the map

$$P_H(x) = x - \frac{\alpha^T x - c}{\|\alpha\|^2} \alpha$$

For  $i = 1, \dots, n$ , if some  $x_i < \epsilon_i$ , then set  $x_i = \epsilon_i$  and eliminate  $x_i$  from further consideration

**until**  $x_i \geq \epsilon_i$  for all  $i$

---

The Michelot algorithm stops after a finite number of iterations because every iteration reduces the dimension  $n$  by at least 1. The first two steps of the algorithm are motivated by the following propositions, whose proofs are straightforward generalizations of those of Michelot (1986). Full validation of the revised algorithm follows from his further arguments.

**Proposition 3.1.** Suppose  $C$  is a closed convex set wholly contained within an affine subspace  $V$ . Then the projection  $P_C(x)$  onto  $C$  and the projection  $P_V(x)$  onto  $V$  satisfy  $P_C(x) = P_C \circ P_V(x)$ .

**Proof.** See Michelot’s paper (Michelot, 1986).  $\square$

**Proposition 3.2.** Suppose  $x \in \mathbb{R}^n$  satisfies  $\sum_i \alpha_i x_i = c$ , where  $\alpha_i > 0$  for all  $i$ . If  $x' \in \mathbb{R}^n$  has coordinates  $x'_i = \max\{x_i, \epsilon_i\}$ , then  $P_K(x) = P_K(x')$  for the truncated simplex (1).

**Proof.** Consider minimizing the objective function  $y \mapsto \frac{1}{2} \|y - x\|^2$  subject to the linear constraint  $\sum_{i=1}^n \alpha_i y_i = c$  and boundary conditions  $y_i \geq \epsilon_i$  for every  $i$ . The Lagrangian function is

$$L(y, \lambda, \mu_i) = \frac{1}{2} \|y - x\|^2 + \lambda \left( \sum_i \alpha_i y_i - c \right) - \sum_i \mu_i (y_i - \epsilon_i).$$

Because this is a convex programming problem, the Karush–Kuhn–Tucker (KKT) optimality conditions are both necessary and sufficient. These conditions can be stated as

$$y_i - x_i + \lambda \alpha_i - \mu_i = 0 \tag{2}$$

$$(y_i - \epsilon_i) \mu_i = 0 \tag{3}$$

$$\mu_i \geq 0 \tag{4}$$

for multipliers  $\lambda$  and  $\mu_i$ . Multiplying both sides of equality (2) by  $\alpha_i$  and summing over  $i$  determines  $\lambda$  as the ratio

$$\lambda = \frac{\sum_i \mu_i \alpha_i}{\sum_i \alpha_i^2} \geq 0.$$

If  $x_i < \epsilon_i$ , then condition (2) implies  $\mu_i = y_i - x_i + \lambda \alpha_i > 0$ . But condition (3) now compels the equality  $y_i = \epsilon_i$ . Therefore we can replace  $x_i$  by  $\epsilon_i$  and  $\mu_i$  by  $\lambda \alpha_i$  and still maintain the KKT conditions at the point  $y$ . In other words,  $P_K(x) = P_K(x')$ .  $\square$

#### 4. A quasi-Newton acceleration scheme

In this section, we review a general quasi-Newton acceleration (Zhou et al., 2009) for fixed point problems. If  $F(x)$  is an algorithm map, then the idea of the scheme is to approximate Newton's method for finding a root of the equation  $\mathbf{0} = x - F(x)$ . Let  $G(x)$  now denote the difference  $G(x) = x - F(x)$ . Because  $G(x)$  has the differential  $dG(x) = I - dF(x)$ , Newton's method iterates according to

$$x^{k+1} = x^k - dG(x^k)^{-1}G(x^k) = x^k - [I - dF(x^k)]^{-1}G(x^k). \quad (5)$$

If we can approximate  $dF(x^k)$  by a low-rank matrix  $M$ , then we can replace  $I - dF(x^k)$  by  $I - M$  and explicitly form the inverse  $(I - M)^{-1}$ .

Quasi-Newton methods operate by secant approximations. We generate one of these by taking two iterates of the algorithm starting from the current point  $x^k$ . When we are close to the optimal point  $x^\infty$ , we have the linear approximation

$$F \circ F(x^k) - F(x^k) \approx M[F(x^k) - x^k],$$

where  $M = dF(x^\infty)$ . If  $v$  is the vector  $F \circ F(x^k) - F(x^k)$  and  $u$  is the vector  $F(x^k) - x^k$ , then the secant requirement is  $Mu = v$ . In fact, for the best results we require several secant approximations  $Mu_i = v_i$  for  $i = 1, \dots, q$ . These can be generated at the current iterate  $x^k$  and the previous  $q - 1$  iterates. The next proposition gives a sensible way of approximating  $M$ .

**Proposition 4.1.** *Let  $M = (m_{ij})$  be a  $n \times n$  matrix and  $\|M\|_F^2 = \sum_i \sum_j m_{ij}^2$  its squared Frobenius norm. Write the secant constraints  $Mu_i = v_i$  in the matrix form  $MU = V$  for  $U = (u_1, \dots, u_q)$  and  $V = (v_1, \dots, v_q)$ . Provided  $U$  has full column rank  $q$ , the minimum of the strictly convex function  $\|M\|_F^2$  subject to the constraints is attained by the choice  $M = V(U^t U)^{-1} U^t$ .*

**Proof.** See the Reference Zhou et al. (2009).  $\square$

To apply the proposition in our proposed quasi-Newton scheme, we must invert the matrix  $I - V(U^t U)^{-1} U^t$ . Fortunately, the explicit inverse

$$[I - V(U^t U)^{-1} U^t]^{-1} = I + V[U^t U - U^t V]^{-1} U^t$$

is a straightforward to check variant of the Sherman–Morrison formula. The  $q \times q$  matrix  $U^t U - U^t V$  is trivial to invert for  $q$  small even when  $n$  is large. This result suggest replacing the Newton update (5) by the quasi-Newton update

$$\begin{aligned} x^{k+1} &= x^k - [I - V(U^t U)^{-1} U^t]^{-1} [x^k - F(x^k)] \\ &= x^k - [I + V(U^t U - U^t V)^{-1} U^t] [x^k - F(x^k)] \\ &= F(x^k) - V(U^t U - U^t V)^{-1} U^t [x^k - F(x^k)]. \end{aligned}$$

The quasi-Newton method is clearly feasible for high-dimensional problems. It takes two ordinary iterates to generate a secant condition and a corresponding quasi-Newton update. If a quasi-Newton update fails to send the objective function in the right direction, then one can always revert to the second iterate  $F \circ F(x^k)$ . For a given  $q$ , the obvious way to proceed is to do  $q$  initial ordinary updates and form  $q - 1$  secant pairs. At that point quasi-Newton updating can commence. After each accelerated update, one should replace the earliest retained secant pair by the new secant pair. The whole scheme is summarized in Algorithm 3. Note that the effort per iteration is relatively light: two ordinary iterates and some matrix times vector multiplications. Most of the entries of  $U^t U$  and  $U^t V$  can be computed once and used over multiple iterations. The scheme is also consistent with linear constraints. Thus, if the parameter space satisfies a linear constraint  $w^t x = a$  for all feasible  $x$ , then the quasi-Newton iterates also satisfy  $w^t x^k = a$  for all  $k$ . This claim follows from the equalities  $w^t F(x) = a$  and  $w^t V = \mathbf{0}$  in the above notation.

Earlier quasi-Newton accelerations (Jamshidian and Jennrich, 1997; Lange, 1995) focus on approximating the Hessian of the objective function rather than the differential of the algorithm map. In our recent paper (Zhou et al., 2009), we demonstrate that the current quasi-Newton acceleration significantly boosts the convergence rate of a variety of optimization algorithms. We apply it in the next section to importance sampling.

#### 5. Example

Our numerical example, borrowed from chapter 14 of the book (Moore and McCabe, 2005), contains a random sample of  $n = 1664$  repair times for Verizon's telephone customers. It is evident from the histogram displayed in Fig. 1 that the distribution of repair times has a long right tail and is far from normal. The median is 3.59 h, but the mean is 8.41 h, and the maximum is 191.6 h. For purposes of illustration, we focus on the probability that the repair time of a Verizon customer exceeds 100 h. The statistic of interest  $T(\mathbf{x}) = \frac{1}{n} \sum_i 1_{\{x_i > 100\}}$  is strongly influenced by extreme repair times. To estimate optimal importance weights, we took a preliminary bootstrap sample of size  $B_1 = 1000$  and executed our estimation procedure in MATLAB. We also performed three forms of interior point optimization in MATLAB's Optimization Toolbox as part of the `fmincon` function. The three standard methods use the exact Hessian of the objective function, a BFGS quasi-Newton approximation to it, and a limited-memory version (LBFGS) of the BFGS approximation. The LBFGS algorithm depends as

**Algorithm 3** Quasi-Newton Acceleration of an Algorithm Map  $F$  for Minimizing the Objective Function  $O$ 


---

```

Initialize:  $x^0, q$ 
for  $i=1$  to  $q + 1$  do
   $x^i = F(x^{i-1})$ 
end for
 $u_i = x^i - x^{i-1}, v_i = x^{i+1} - x^i, i = 1, \dots, q$ 
 $U = [u_1 \dots u_q], V = [v_1 \dots v_q]$ 
 $x^n = x^{q+1}$ 
repeat
   $x_1 = F(x^n)$ 
  if  $[O(x_1) - O(x^n)] \leq \epsilon[|O(x^n)| + 1]$  then
    break
  end if
   $x_2 = F(x_1)$ 
  update the oldest  $u$  in  $U$  by  $u = x_1 - x^n$ 
  update the oldest  $v$  in  $V$  by  $v = x_2 - x_1$ 
   $x_{qn} = x_1 + V(U^t U - U^t V)^{-1} U^t u$ 
  if  $x_{qn}$  falls outside the feasible region then
    project  $x_{qn}$  onto feasible region
  end if
  if the objective function satisfies  $O(x_{qn}) < O(x_2)$  then
     $x^{n+1} = x_{qn}$ 
  else
     $x^{n+1} = x_2$ 
  end if
until convergence occurs

```

---

**Table 1**

Comparison of algorithms for calculating importance weights for the Verizon repair time data set. There are  $n = 1664$  observations and  $B_1 = 1000$  bootstrap replicates.

Algorithm	Iters	$s(p^*) (\times 10^{-6})$	Time
Naive	11586	4.666920	2023.1764
$q = 1$	24	4.665277	11.7226
$q = 2$	19	4.665277	8.9670
$q = 3$	16	4.665277	7.1562
$q = 4$	16	4.665277	6.9169
$q = 5$	17	4.665276	7.0566
$q = 6$	17	4.665277	6.7434
$q = 7$	18	4.665277	6.8271
$q = 8$	19	4.665277	7.0829
$q = 9$	20	4.665277	7.2667
$q = 10$	21	4.665277	7.4299
Int-Pt (Hessian)	18	4.665276	1247.8343
Int-Pt (BFGS)	69	4.665276	149.7225
Int-Pt (LBFGS)		Refer to Table 2	

well on the number  $q$  of secant conditions selected. The active set and trust region methods also implemented in MATLAB are ignored here. The first is noticeably slower than the interior point methods, and the second cannot handle equality constraints and boundary conditions.

The exact Hessian method takes only 18 iterations but 1248 s to converge. In practice, 99% of the execution time is spent on evaluating the Hessian matrix, which requires  $O(n^2 B_1)$  operations per iteration, and 1% of the time is spent on factoring the Hessian matrix, which requires  $O(n^3)$  operations per iteration. This example illustrates the extreme speed of MATLAB's matrix operations. The interior point method with BFGS updates takes many more iterations but much less time per iteration because it dispenses with evaluating and factoring the Hessian matrix. Part of the slow convergence of the BFGS method may be attributed to the boundary conditions. In contrast, our algorithm takes  $O(n B_1)$  operations per iteration and converges quickly under acceleration. As shown in Table 1, the accelerated algorithm with  $q \geq 1$  secant conditions is a clear winner, giving massive improvements in execution time over the naive MM algorithm and the Hessian and BFGS variants of Newton's method. Although our accelerated algorithm also beats the LBFGS algorithm for all choices of  $q$ , Table 2 shows that the later algorithm is highly competitive on large-scale problems. All comparisons listed in the two tables involve stringent stopping criteria, which are adjusted to give the same number of significant digits for the converged values of the objective function. Running times are recorded in seconds.

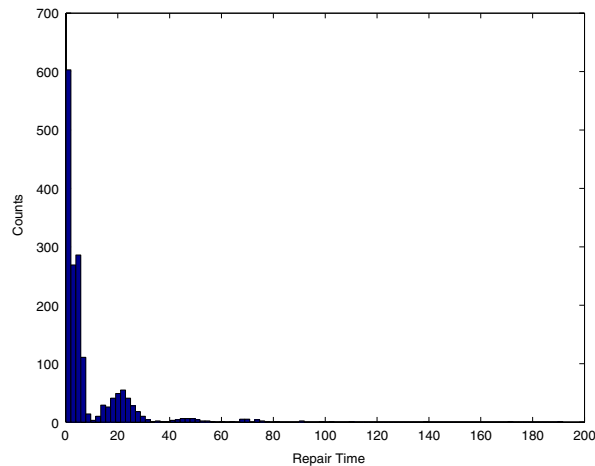


Fig. 1. Histogram of 1664 repair times for Verizon's customers.

Table 2

Comparison of our quasi-Newton acceleration and the LBFGS methods on the Verizon data set.

$q$	Quasi-Newton			Int-Pt (LBFGS)		
	Iters	$s(p^*)(\times 10^{-6})$	Time	Iters	$s(p^*)(\times 10^{-6})$	Time
1	21	4.66528164	11.3508	149	4.66527741	18.6839
2	17	4.66527843	8.5021	90	4.66527742	12.1505
3	15	4.66527728	7.0560	75	4.66527741	9.9297
4	15	4.66527699	6.6824	89	4.66527742	12.2690
5	16	4.66527679	7.1644	78	4.66527741	11.3548
6	16	4.66527712	7.3846	83	4.66527742	12.1189
7	17	4.66527710	6.3872	76	4.66527742	12.1928
8	18	4.66527704	6.9555	74	4.66527742	11.7458
9	19	4.66527699	6.8392	67	4.66527741	10.6362
10	20	4.66527703	7.0096	67	4.66527742	10.2277

## 6. Discussion

In summary, our procedure uses: (a) quadratic approximation of the objective function, (b) majorization by a second quadratic with parameters separated, (c) the Michelot algorithm to project points onto a truncated simplex, and (d) acceleration by quasi-Newton approximation of the algorithm map. Each of these ideas has other applications.

Quadratic approximation lies at the heart of Newton's method and its many spinoffs. Our outer-product majorization balances separation of parameters against poor approximation of the Hessian. Separation of parameters is often the key to solving high-dimensional problems. The loss of quadratic convergence is largely remedied by acceleration. This combination of tactics also has potential in fitting generalized linear models (GLM). In this setting, Fisher's scoring method uses the expected information matrix

$$J(\beta) = \sum_{i=1}^n \frac{1}{\sigma_i^2(\beta)} q'(x_i^t \beta)^2 x_i x_i^t$$

rather than the observed information matrix. Here  $\beta$  is the parameter vector,  $q(\cdot)$  is the inverse link function,  $x_i$  is the predictor vector for case  $i$ ,  $y_i$  the response for case  $i$ , and  $\sigma_i^2(\beta) = \mathbf{Var}(y_i)$ . Note that  $J(\beta)$  is again a sum of outer products. This fact suggests a combination of majorization and acceleration on high-dimensional GLM problems. In many such problems, it is prudent to also impose a ridge or lasso penalty. The ridge penalty preserves separation of parameters by a quadratic surrogate. The lasso penalty also preserves separation of parameters, but not by a quadratic surrogate. Lasso penalized maximum likelihood estimation is amenable to cyclic coordinate ascent because the lasso is linear on either side of 0. Our recent work on penalized ordinary and logistic regression (Wu and Lange, 2008; Wu et al., 2009) illustrates some of the possibilities. Finally, our paper (Zhou et al., 2009) amply illustrates the virtues of acceleration by quasi-Newton approximation of an algorithm map.

## References

- Davison, A.C., 1988. Discussion of paper by D.V. Hinkley. *J. Roy. Statist. Soc. Ser. B* 50, 356–357.  
 Do, K.A., Wang, X., Broom, B.M., 2001. Importance bootstrap resampling for proportional hazards regression. *Comm. Statist. Theory Methods* 30 (10), 2173–2188.  
 Do, K.A., Hall, P., 1991. On importance resampling for the bootstrap. *Biometrika* 78 (1), 161–167.

- Hall, P., 1992. The Bootstrap and Edgeworth Expansion. Springer-Verlag, New York.
- Hesterberg, T., 1996. Control variates and importance sampling for efficient bootstrap simulations. *Statist. Comput.* 6, 147–157.
- Hinkley, D.V., Shi, S., 1989. Importance sampling and the nested bootstrap. *Biometrika* 76 (3), 435–446.
- Hu, J., Su, Z., 2008. Bootstrap quantile estimation via importance resampling. *Comput. Statist. Data Anal.* 52 (12), 5136–5142.
- Hunter, D.R., Lange, K.L., 2004. A tutorial on MM algorithms. *Amer. Statist.* 58, 30–37.
- Jamshidian, M., Jennrich, R.I., 1997. Acceleration of the EM algorithm by using quasi-Newton methods. *J. Roy. Statist. Soc. Ser. B* 59 (3), 569–587.
- Johns, M.V., 1988. Importance sampling for bootstrap confidence intervals. *J. Amer. Statist. Assoc.* 83 (403), 709–714.
- Lange, K.L., 1995. A quasi-Newton acceleration of the EM algorithm. *Statist. Sinica* 5 (1), 1–18.
- Lange, K.L., 2004. Optimization. Springer-Verlag, New York.
- Michelot, C., 1986. A finite algorithm for finding the projection of a point onto the canonical simplex of  $\mathbf{R}^n$ . *J. Optim. Theory Appl.* 50 (1), 195–200.
- Moore, D.S., McCabe, G.P., 2005. Introduction to the Practice of Statistics. W.H. Freeman.
- Wu, T.T., Lange, K.L., 2008. Coordinate descent algorithms for lasso penalized regression. *Ann. Appl. Stat.* 2, 224–244.
- Wu, T.T., Chen, Y.F., Hastie, T., Sobel, E.M., Lange, K.L., 2009. Genomewide association analysis by lasso penalized logistic regression. *Bioinformatics* 25, 714–721.
- Zhou, H., Alexander, D., Lange, K.L., 2009. A quasi-Newton acceleration for high-dimensional optimization algorithms. *Stat. Comput.* doi:10.1007/s11222-009-9166-3.