

# MULTIDIMENSIONAL SHARP QUADRATIC MAJORIZATION

JAN DE LEEUW

ABSTRACT. In a recent paper De Leeuw and Lange [2006] study sharp quadratic majorization for functions of a single variable. In this note we analyze a multivariate example and give some partial results.

## 1. INTRODUCTION

The problem we study in this paper is the minimization of a differentiable function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  using quadratic majorization. This means that in each iteration ( $k$ ) we find a positive definite matrix  $A(x^{(k)})$  such that

$$(1) \quad f(x) \leq g(x, x^{(k)}) = f(x^{(k)}) + (x - x^{(k)})' \mathcal{D}f(x^{(k)}) + \frac{1}{2}(x - x^{(k)})' A(x^{(k)})(x - x^{(k)})$$

for all  $x$ . We then define the algorithmic map by

$$(2) \quad x^{(k+1)} = \underset{x}{\operatorname{argmin}} g(x, x^{(k)}) = x^{(k)} - A(x^{(k)})^{-1} \mathcal{D}f(x^{(k)}).$$

Since

$$(3) \quad f(x^{(k+1)}) \leq g(x^{(k+1)}, x^{(k)}) = \min_x g(x, x^{(k)}) < g(x^{(k)}, x^{(k)}) = f(x^{(k)})$$

the algorithm generates a decreasing sequence of loss function values, and under the conditions given by Zangwill [1969] we have convergence to a local minimum at  $x_\infty$ . The linear convergence rate of the algorithm is the spectral radius  $\rho(I - A(x_\infty)^{-1} \mathcal{D}^2 f(x_\infty))$ , and consequently it is of interest to choose matrices  $A(x^{(k)})$  which are as large as possible.

This leads to the auxiliary problem of finding a “large”  $A$  satisfying the constraints

$$(4) \quad f(x) - f(x^{(k)}) - (x - x^{(k)})' \mathcal{D}f(x^{(k)}) - \frac{1}{2}(x - x^{(k)})' A(x - x^{(k)}) \leq 0$$

for all  $x$ . Equivalently (4) can be written as

$$(5) \quad \max_x f(x) - f(x^{(k)}) - (x - x^{(k)})' \mathcal{D}f(x^{(k)}) - \frac{1}{2}(x - x^{(k)})' A(x - x^{(k)}) = 0.$$

Note that (4) and (5) define a convex set of matrices  $A$ . Of course we need to have a specific definition of “large” to make this a well-defined optimization problem. One could use, for example, the trace, or the sum of squares, or the spectral radius.

## 2. MAJORIZATION OF MULTINOMIAL LOGISTIC

As an example consider the Multinomial Logistic with

$$\pi_j(x) = \frac{\exp(x_j)}{\sum_{\ell=1}^n \exp(x_\ell)}$$

and

$$f(x) = - \sum_{j=1}^m p_j \log \pi_j(x),$$

where  $p$  is a fixed probability vector.

Now

$$\mathcal{D}f(x) = \pi(x) - p,$$

and

$$\mathcal{D}^2 f(x) = \Pi(x) - \pi(x)\pi(x)'$$

Since  $\mathcal{D}^2 f(x) \lesssim \frac{1}{2}I$  uniform quadratic majorization is easy. The majorization function is just

$$g(x, y) = f(y) + (x - y)' \mathcal{D}f(y) + \frac{1}{4}(x - y)'(x - y),$$

leading to the algorithmic map

$$x^{(k+1)} = y - 2\mathcal{D}f(x^{(k)}),$$

which has convergence rate  $\rho(I - 2\mathcal{D}^2 f(\hat{x}))$

In this note we investigate if locally sharper quadratic majorizations are possible, i.e. we look for a  $\delta(y) \leq \frac{1}{2}$  such that

$$f(x) \leq f(y) + (x - y)' \mathcal{D}f(y) + \frac{1}{2}\delta(y)(x - y)'(x - y).$$

for all  $x$  and  $y$ . This will give a majorization algorithm with a faster convergence rate  $\rho(I - \delta(y)^{-1}\mathcal{D}^2 f(\hat{x}))$ .

Define

$$\delta(x, y) = \frac{f(x) - f(y) - (x - y)' \mathcal{D}f(y)}{\frac{1}{2}(x - y)'(x - y)},$$

and

$$\hat{\delta}(y) = \max_x \delta(x, y).$$

We know that  $\hat{\delta}(y) \leq \frac{1}{2}$ , but computation suggests that  $\hat{\delta}(y)$  can actually be much smaller than  $\frac{1}{2}$ , especially if  $n$  is large and/or  $y$  has large variation.

*Example 1.* For  $n = 2$  and  $y_1 \neq y_2$  we have

$$\hat{\delta}(y) = \frac{\pi_2(y) - \pi_1(y)}{y_2 - y_1}.$$

attained for  $x_1 = y_2$  and  $x_2 = y_1$ . This is basically the same result as the sharp quadratic majorization of the logistic in De Leeuw and Lange [2006]. If  $y_2 \rightarrow y_1$  then  $\hat{\delta}(y) \rightarrow \frac{1}{2}$ .

*Example 2.* If all  $y$  are equal (without loss of generality we can take them all to be zero) then the optimal  $x$  has  $n - 1$  elements equal to  $-\frac{2}{n} \log(n - 1)$  and one element equal to  $2 \frac{n-1}{n} \log(n - 1)$ . Also

$$\hat{\delta}(y) = \frac{1}{2} \frac{n - 2}{(n - 1) \log(n - 1)}$$

By L'Hospital if  $n \rightarrow 2$  we have  $\hat{\delta}(y) \rightarrow \frac{1}{2}$ .

*Example 3.* We can use the fact that for all permutations  $P$  we have  $\pi(Px) = P(\pi(x))$ . Thus if  $\hat{x}$  is a permutation of the elements of  $y$  we must have  $(P - I)\pi(y) = \hat{\delta}(y)(P - I)y$ , i.e. there exist constants  $\alpha$  and  $\beta$  such that  $\pi_i(y) = \alpha + \beta y_i$ . Clearly we can solve for  $\alpha$  and  $\beta$  if  $y$  only takes two different values, generalizing Example 1. In this case  $\hat{\delta}(y)$  does not depend on the number of times the two values occur in  $y$ .

### 3. THE SHARP UPPER BOUND

For the computation  $\hat{\delta}(y)$  we could use  $\mathcal{D}_1 \delta(x, y) = 0$  if  $\pi(x) - \pi(y) = \delta(x, y)(x - y)$ , which suggests the iterative algorithm

$$x^{(k+1)} = y + \frac{1}{\delta(x^{(k)}, y)} (\pi(x^{(k)}) - \pi(y)).$$

This algorithm was discussed earlier in a unidimensional context by De Leeuw and Lange [2006]. Their proof that the algorithm is monotone and convergent for any convex function  $f$  with a bounded second derivative easily extends to the multivariate case. The R code for the function `deltaOpt()` is given in the Appendix.

At a stationary point  $\hat{x}$  we have

$$\hat{x} = y + \frac{1}{\delta(y)} (\pi(\hat{x}) - \pi(y)),$$

which implies  $\sum_{j=1}^m \hat{x}_j = \sum_{j=1}^m y_j$ .

Instead, we propose an alternative algorithm which generalizes more easily to situations in which  $\delta(y)$  is not necessarily a scalar. Let us start with some initial  $\delta^{(0)} \leq \hat{\delta}(y)$ . One easy initial estimate is

$$\delta^{(0)} = 2 \max_{i=1}^n f(y + e_i) - f(y) - \mathcal{D}_i f(y).$$

In iteration  $k$  we compute

$$\max_x \{f(x) - f(y) - (x - y)' \mathcal{D}f(y) - \frac{1}{2} \delta^{(k)} (x - y)' (x - y)\}.$$

If the maximum value is positive, attained say at  $x^{(k)}$ , then we set

$$\delta^{(k+1)} = \frac{f(x^{(k)}) - f(y) - (x^{(k)} - y)' \mathcal{D}f(y)}{\frac{1}{2} (x^{(k)} - y)' (x^{(k)} - y)},$$

and we conclude that  $\delta^{(k)} < \delta^{(k+1)} \leq \hat{\delta}(y)$ . If the maximum value is zero, we stop, because that implies  $\delta^{(k)} \geq \hat{\delta}(y)$ , and thus actually  $\delta^{(k)} = \hat{\delta}(y)$ .

#### 4. SHARP MAJORIZATION MATRIX

If we don't want to compute a sharp majorization constant, but a sharp majorization matrix, then that matrix  $\Delta$  should satisfy

$$(6) \quad h(x, y) = f(x) - f(y) - (x - y)' \mathcal{D}f(y) - \frac{1}{2} (x - y)' \Delta (x - y) \leq 0$$

for all  $x$ . Or, equivalently,  $\max_x h(x, y) = 0$ , with the maximum attained at  $x = y$ . We usually also require positive semi-definiteness, i.e.  $\Delta \succeq 0$ . This defines an infinite system of linear inequalities in  $\Delta$ , which means that for each  $y$  there is a convex set of matrices  $\Delta(y)$  satisfying these inequalities. In general this set could be empty, but we know that in our example every matrix  $\Delta \succeq \frac{1}{2} I$  will be in  $\Delta(y)$ .

If  $f$  is convex, then we can use majorization to check if  $\Delta \in \Delta(y)$ . The algorithm is

$$(7) \quad x^{(k+1)} = y + \Delta^{-1}(\mathcal{D}f(x^{(k)}) - \mathcal{D}f(y)) = y + \Delta^{-1}(\pi(x^{(k)}) - \pi(y)).$$

Thus will increase  $h(x, y)$  and as soon as  $h(x^{(k)}, y)$  becomes larger than zero, we know that  $\Delta \notin \Delta(y)$ . In fact, either  $h(x^{(k)}, y)$  converges to zero, and  $x^{(k)}$  converges to  $y$ , or  $h(x^{(k)}, y)$  becomes positive and then diverges. Note that  $p$  does not play a role in this algorithm. Again the R code to test a particular candidate matrix  $\Delta$  is in the Appendix.

Suppose  $\phi$  measures how large  $\Delta$  is. For instance,  $\phi$  could be the trace, or the sum-of-squares, or the sup-norm. Start with a  $\Delta$  which is not in  $\Delta(y)$ . Compute the  $x$  maximizing  $h(x,y)$ , say  $\hat{x}$ , and add the linear constraint  $h(\hat{x},y) \leq 0$  to the constraint set. Find  $\Delta$  minimizing  $\phi(\Delta)$  over the linear constraints accumulated so far. Find a new  $\hat{x}$ , and add the new constraint, and so on. This is a variation of the cutting plane method, and it will find the majorization matrix with the smallest value of  $\phi$ , in other words the majorization that is  $\phi$ -sharp.

To start we can use (6) to derive simple lower bounds for the diagonal elements of  $\Delta$ . We have

$$\delta_{ii} \geq 2\{f(y + e_i) - f(y) - \mathcal{D}_i f(y)\}$$

## APPENDIX A. CODE

```

1  delta<-function(x,y,p) {
2  pix<-exp(x)/sum(exp(x))
3  piy<-exp(y)/sum(exp(y))
4  fx<-sum(p*log(pix))
5  fy<-sum(p*log(piy))
6  gy<-piy-p
7  return((fx-fy-sum(gy*(x-y)))/(.5*sum((x-y)^2)))
8  }
9
10 deltaOpt<-function(y,p=rep(1/length(y),length(y)),itmax
    =100,eps=1e-10,verbose=TRUE) {
11     n<-length(y)
12     p<-runif(n)
13     p<-p/sum(p)
14     b<-bnds(y,p)
15     x<-y
16     x[which.max(b)]<-x[which.max(b)]+1
17     piy<-exp(y)/sum(exp(y))
18     fy<-sum(p*log(piy))
19     gy<-piy-p
20     itel<-1
21     repeat {
22     pix<-exp(x)/sum(exp(x))
23     fx<-sum(p*log(pix))
24     gx<-pix-p
25     lb<-sum(gy*(x-y))/(.5*sum((x-y)^2))
26     if (verbose) cat("D-Iteration: ",formatC(itel,
        digits=3,width=3),
27         "lb: ",formatC(lb,digits=6,width=10,
        format="f"),
28         "\n")
29     xnew<-y+(gx-gy)/lb
30     if ((max(abs(x-xnew)) < eps) | (itel == itmax))
        break()
31     x<-xnew

```

```

32     itel<-itel+1
33     }
34     return(list(y=y,p=p,x=x,pix=pix,piy=piy,d=lb))
35 }
36
37 deltaMat<-function(y,x=sample(y,length(y)),d,itmax=1000,
    eps=1e-10,ops=Inf,verbose=TRUE) {
38     n<-length(y)
39     p<-runif(n)
40     p<-p/sum(p)
41     piy<-exp(y)/sum(exp(y))
42     fy<-sum(p*log(piy))
43     gy<-piy-p
44     itel<-1
45     repeat {
46         pix<-exp(x)/sum(exp(x))
47         fx<-sum(p*log(pix))
48         gx<-pix-p
49         hxy<-fx-(fy+sum((x-y)*gy)+sum((x-y)*colSums((x-y)
            *d))/2)
50         if (verbose) cat("H-Iteration: ",formatC(itel,
            digits=3,width=3),
51             "hxy: ",formatC(hxy,digits=6,width=10,
            format="f"),
52             "\n")
53         xnew<-y+solve(d,(gx-gy))
54         if ((max(abs(x-xnew)) < eps) | (itel == itmax))
            break()
55         if (hxy > ops) break()
56         x<-xnew
57         itel<-itel+1
58     }
59     return(list(h=hxy,x=x,cf=0.5*(x-y)^2,bf=fx-(fy+sum((x-y)
        *gy))))
60 }
61

```

```

62 bnds<-function(y,p) {
63   n<-length(y)
64   b<-rep(0,n)
65   piy<-exp(y)/sum(exp(y))
66   fy<-sum(p*log(piy))
67   gy<-piy-p
68   for (i in 1:n){
69     x<-y
70     x[i]<-x[i]+1
71     pix<-exp(x)/sum(exp(x))
72     fx<-sum(p*log(pix))
73     gx<-pix-p
74     b[i]<-2*(fx-(fy+gy[i]))
75   }
76   return(b)
77 }
78
79 altIter<-function(y,p=rep(1/length(y),length(y)),itmax
   =1000,eps=1e-10,verbose=FALSE) {
80   n<-length(y)
81   piy<-exp(y)/sum(exp(y))
82   fy<-sum(p*log(piy))
83   gy<-piy-p
84   b<-bnds(y,p)
85   iold<-which.max(b)
86   bold<-b[iold]
87   xold<-y
88   xold[iold]<-xold[iold]+1
89   itel<-1
90   repeat {
91     xupd<-deltaMat(y,x=xold,d=bold*diag(n),verbose=
       verbose)
92     xnew<-xupd$x
93     hval<-xupd$h
94     if (hval < eps) break()
95     ddst<-sum((xnew-y)^2)/2

```



```
96 pix<-exp(xnew)/sum(exp(xnew))
97 fx<-sum(p*log(pix))
98 bnew<-(fx-(fy+sum((xnew-y)*gy)))/ddst
99 cat("B-Iteration: ",formatC(itel,digits=3,width
    =3),
100     "bold: ",formatC(bold,digits=6,width=10,
    format="f"),
101     "bnew: ",formatC(bnew,digits=6,width=10,
    format="f"),
102     "hval: ",formatC(hval,digits=6,width=10,
    format="f"),
103     "ddst: ",formatC(ddst,digits=6,width=10,
    format="f"),
104     "\n")
105 if ((bnew-bold) < eps) | (itel == itmax)
    break()
106 bold<-bnew
107 itel<-itel+1
108 }
109 }
```

## REFERENCES

- J. De Leeuw and K. Lange. Sharp Quadratic Majorization in One Dimension. Preprint 503, UCLA Department of Statistics, 2006. URL [http://preprints.stat.ucla.edu/503/sharp\\_third.pdf](http://preprints.stat.ucla.edu/503/sharp_third.pdf).
- W. I. Zangwill. Convergence Conditions for Nonlinear Programming Algorithms. *Management Science*, 16:1–13, 1969.

DEPARTMENT OF STATISTICS, UNIVERSITY OF CALIFORNIA, LOS ANGELES, CA 90095-1554

*E-mail address*, Jan de Leeuw: [deleeuw@stat.ucla.edu](mailto:deleeuw@stat.ucla.edu)

*URL*, Jan de Leeuw: <http://gifi.stat.ucla.edu>