

## Gradient Methods of Maximization

by Herman Chernoff  
Jean Bronfenbrenner

1. Introduction

In estimating a set of parameters by maximum likelihood methods, it is ordinarily necessary to find the point in the parameter space which maximizes a certain function (the likelihood function). Frequently the form of the function and the number of independent variables involved make it prohibitively difficult to determine the maximizing point directly, and methods of successive approximations are accordingly used.

Suppose we have an approximation  $x^{(0)} = (x_1^{(0)} \ x_2^{(0)} \ \dots \ x_n^{(0)})$  to the point  $c = (c_1 \ c_2 \ \dots \ c_n)$  which maximizes a function  $f(x)$ . Then if we expand  $f(x)$  in a Taylor series about  $x^{(0)}$  we have a means for obtaining an improved approximation to  $c$ . Presumably the more terms of the expansion we consider, the greater the speed with which we may expect to approach  $c$ . However, calculation of terms involving higher order derivatives increases considerably the computational cost of each iteration.

The methods of successive approximation to be discussed in this paper are gradient methods using the first order derivatives only and the Newton method which uses first and second order terms. In both cases it is possible to obtain from the successive approximations certain relevant information about terms of higher order than those actually computed, and to use this information to improve convergence.

2. Gradient Methods

Given an initial approximation  $x^{(0)} = (x_1^{(0)} \ x_2^{(0)} \ \dots \ x_n^{(0)})$  it is natural to choose the next approximation  $x^{(1)}$  in such a way that the step from  $x^{(0)}$  to

$x^{(1)}$  is in the direction of steepest ascent with respect to  $f(x)$ , i.e., in the direction of the gradient. The direction of steepest ascent depends, however, on the concept of distance used. In general there is no reason to assume that a unit of distance along the  $x_i$ -axis is in any sense equivalent to a unit of distance along the  $x_j$ -axis,  $i \neq j$ . The metric chosen implies a particular system for weighting these units.

Let  $\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)$ , and suppose that the distance from  $x^{(0)}$  to the point  $(x^{(0)} + \epsilon)$  be defined as

$$d = \left( \sum_{i,j} B_{ij} \epsilon_i \epsilon_j \right)^{\frac{1}{2}}$$

where  $B = \| B_{ij} \|$  is a positive definite symmetric matrix which may or may not depend on the point  $x^{(0)}$ . Then the locus of points  $(x^{(0)} + \epsilon)$  lying at a distance  $k$  from  $x^{(0)}$  is given by the ellipsoid

$$\sum B_{ij} \epsilon_i \epsilon_j = k^2$$

with center at  $x^{(0)}$ .

If we now let  $\epsilon_i = k \delta_i$  this ellipsoid becomes

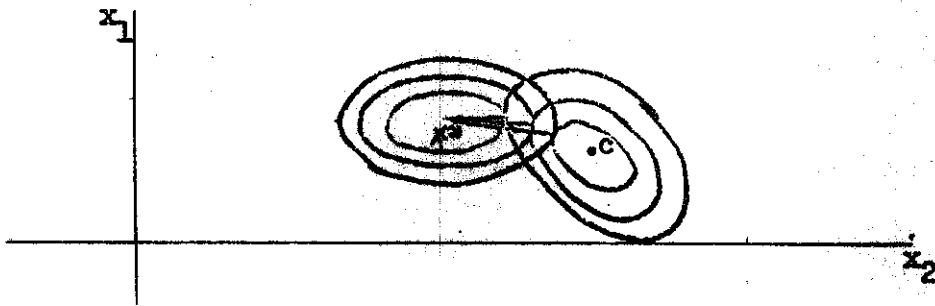
$$(1) \quad \sum B_{ij} \delta_i \delta_j = 1$$

Then the direction of steepest ascent in the  $k$ -neighborhood of  $x^{(0)}$  is the direction from  $x^{(0)}$  to the point of the ellipsoid for which the function is a maximum. In the limiting case as  $k$  approaches zero it is easy to approximate this direction by writing the Taylor expansion of  $f(x^{(0)} + k\delta)$ , maximizing subject to the restraint (1) and disregarding higher order terms in  $k$ .

For a particular value of  $k$ , the maximizing values of the  $\delta_i$  give the projection on the  $x_1 x_2 \dots x_n$  - space of the steepest secant that can be drawn through the  $(n+1)$  - dimensional point  $[x^{(0)}, f(x^{(0)})]$  and any point  $[x^{(0)} + k\delta, f(x^{(0)} + k\delta)]$  such that  $\delta$  satisfies (1). In the limit as  $k \rightarrow 0$ , the maximizing values of the  $\delta_i$  give the projection of the steepest

-3-

tangent line through  $[x^{(0)}, f(x^{(0)})]$ . The two dimensional case is illustrated in the graph below. The curves enclosing  $c$  are contour lines, each corresponding to  $f(x) = \text{constant}$ . The curves enclosing  $x^{(0)}$  are ellipses corresponding to different values of  $k$ . The arrows indicate the directions of the steepest secants for the respective values of  $k$ .



Let  $\bar{z}$  be a row vector whose  $i^{\text{th}}$  element is  $\frac{\partial f}{\partial x_i} [= f_{i1}]$  and let  $L$  be a matrix whose  $ij^{\text{th}}$  element is  $-\frac{\partial^2 f}{\partial x_i \partial x_j} [= -f_{ij}]$ . Then

$$f(x^{(0)} + k\delta) = f(x^{(0)}) + k \delta \bar{z}'(x^{(0)}) - \frac{k^2}{2} \delta L(x^{(0)}) \delta' + \dots \quad \underline{1/}$$

Maximizing subject to the restraint (1) we have

$$k \bar{z}'(x^{(0)}) - k^2 L(x^{(0)}) \delta' + \dots = 2 \lambda B \delta',$$

where  $\lambda$  is a Lagrange multiplier, and for very small  $k$

$$\delta' \approx B^{-1} \bar{z}'(x^{(0)})$$

except for a scale factor. This indicates in effect the proper direction for a step, in the limit as the size of the step approaches zero. The problem then arises as to how large a step may profitably be taken in the limiting direction indicated. This clearly depends on how fast the slope of  $f(x)$  changes as we depart from  $x^{(0)}$  in various directions and so cannot be determined without some consideration of higher order terms.

We observe that if

$$x^{(1)'} = x^{(0)'} + h\delta' = x^{(0)'} + h B^{-1} \bar{z}'(x^{(0)}),$$

---

1/ A vector  $x$  is considered to be a matrix with only one row and  $x'$ , the transpose of  $x$ , would therefore consist of one column. Thus  $\delta \bar{z}'$  is  $\sum \delta_i \bar{z}_i$  or the scalar product of  $\delta$  and  $\bar{z}$ .

then for small enough  $h > 0$  we necessarily have

$$f(x^{(1)}) > f(x^{(0)}).$$

For

$$\begin{aligned} f(x^{(1)}) &= f(x^{(0)} + h\delta) = f(x^{(0)}) + h\delta \dot{f}(x^{(0)}) + \dots \\ &= f(x^{(0)}) + h \dot{f}(x^{(0)}) B^{-1} \dot{f}(x^{(0)}) + \dots \end{aligned}$$

But

$$h \dot{f}(x^{(0)}) B^{-1} \dot{f}(x^{(0)}) > 0$$

since  $h > 0$  and  $B$  is positive definite.

It may further be shown that if a sequence of successive approximations  $x^{(0)}, x^{(1)}, \dots, x^{(m)}, \dots$  is formed, where

$$(2) \quad x^{(m+1)} = x^{(m)} + h B^{-1} \dot{f}(x^{(m)})$$

then if the initial approximation is in a sufficiently small neighborhood of  $c$  this sequence converges to  $c$  for sufficiently small values of  $h$ . To

establish this result we expand  $f(x)$  about  $c$ . Letting  $x - c = e$  we have

$$(3) \quad \dot{f}(x) = \dot{f}(c + e) = \dot{f}(c) - L(c) e + \dots = -L(c) e + \dots$$

since  $\dot{f}(c) = 0$ . If  $x^{(m)} - c = e^{(m)}$  we have from (2) and (3)

$$x^{(m+1)} - c = x^{(m)} - c - h B^{-1} L(c) e^{(m)} + \dots$$

$$(4) \quad e^{(m+1)} \approx [I - h B^{-1} L(c)] e^{(m)} \approx [I - h B^{-1} L(c)]^{m+1} e^{(0)}$$

Now if  $c$  is an isolated maximum,  $L(c)$  is a positive definite matrix and  $B^{-1}L(c)$  is therefore the product of two positive definite matrices and itself positive definite. Let  $\lambda_i$  and  $\mu_i$ ,  $i = 1, 2, \dots, n$  be respectively the characteristic values and vectors of  $B^{-1}L(c)$ . Then  $\lambda_i > 0$  for all  $i$ . Indeed we assume  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$ . We observe that the  $\mu_i$  are also the characteristic vectors of  $I - h B^{-1}L(c)$  while  $1 - h\lambda_i$ ,  $i = 1, 2, \dots, n$ , are the characteristic values. For by definition

$$\begin{aligned}
 [B^{-1} L(c) - \lambda_i I] \mu_i' &= 0 \\
 B^{-1} L(c) \mu_i' &= \lambda_i \mu_i' \\
 \mu_i' - h B^{-1} L(c) \mu_i' &= \mu_i' - h \lambda_i \mu_i' \\
 [I - h B^{-1} L(c) - (1 - h \lambda_i) I] \mu_i' &= 0
 \end{aligned}$$

For sufficiently small values of  $h$ , we have  $0 < 1 - h \lambda_i < 1^{1/2}$  since  $\lambda_i > 0$ .

We may decompose the vector  $(x^{(0)} - c) = e^{(0)}$  into components along the  $n$  orthogonal characteristic vectors  $\mu_i$ . Then (4) gives for  $m = 0$

$$\begin{aligned}
 e^{(1)'} &= [I - h B^{-1} L(c)] [k_1 \mu_1' + k_2 \mu_2' + \dots + k_n \mu_n'] + O(e^{(0)2}) \\
 &= (1 - h \lambda_1) k_1 \mu_1' + \dots + (1 - h \lambda_n) k_n \mu_n' + O(e^{(0)2})
 \end{aligned}$$

From this the proof of convergence is evident for an initial approximation close enough to  $c$ . We have furthermore

$$\begin{aligned}
 (5) \quad e^{(m+1)'} &= (1 - h \lambda_1)^{m+1} k_1 \mu_1' + \dots + (1 - h \lambda_n)^{m+1} k_n \mu_n' \\
 (6) \quad \delta^{(m)'} &= \frac{x^{(m+1)'} - x^{(m)'}}{h} \approx -B^{-1} L(c) e^{(m)'} \approx -B^{-1} L(c) [I - h B^{-1} L(c)]^m e^{(0)'} \\
 &= -\sum_i \lambda_i (1 - h \lambda_i)^m k_i \mu_i'
 \end{aligned}$$

### 3. Choice of $h$

We mentioned above that an intelligent choice of  $h$  depends on some information regarding higher order terms in the Taylor expansion of  $f(x)$ . It appears from (5) that the relevant information concerns the characteristic values of  $B^{-1} L(c)$ , for by letting  $h = \frac{1}{\lambda_i}$  we may eliminate the  $i^{\text{th}}$  terms on the right in (5). In particular it will be shown below how we may, without explicitly evaluating  $L$ , use successive iterations with any sufficiently small  $h$  to estimate  $\lambda_1$  and  $\lambda_n$ , the largest and the smallest of the characteristic values.

- 1/ If  $B$  is not strictly constant, we may assume that it approaches a certain value as  $x$  approaches  $c_1$ . Then  $\lambda_i$  behave in a similar fashion and the inequality above is valid for  $x$  sufficiently close to  $c$ .
- 2/  $y = O(x)$  indicates that  $y$  is at most of the order of magnitude of  $x$ . That is there exists a  $k$  so that  $|y| \leq k|x|$  as  $x \rightarrow 0$ . The extension of this notation to vectors is trivial.

Now if we let  $h = \frac{1}{\lambda_n}$  the last term in (5) vanishes but the coefficient of the first becomes  $(1 - \frac{\lambda_1}{\lambda_n})^{m+1}$  and this may be very large if  $\lambda_1$  is much greater than  $\lambda_n$ . On the other hand if  $h = \frac{1}{\lambda_1}$  the first term vanishes but the coefficient of the last becomes  $(1 - \frac{\lambda_n}{\lambda_1})^{m+1}$ , and this, while less than unity, may decrease very slowly as  $m$  increases. Therefore, if a constant  $h$  is to be used the value

$$h = \frac{2}{\lambda_1 + \lambda_n}$$

is recommended. This minimizes the maximum value of  $|1 - h \lambda_i|$  by making the largest positive and the largest negative value numerically equal. I.e., it makes

$$\begin{aligned} -(1 - h \lambda_1) &= 1 - h \lambda_n = 1 - \frac{2\lambda_n}{\lambda_1 + \lambda_n} \\ &= \frac{\lambda_1 - \lambda_n}{1 + \frac{\lambda_n}{\lambda_1}} \\ &= \frac{1 - \frac{\lambda_n}{\lambda_1}}{1 + \frac{\lambda_n}{\lambda_1}} \end{aligned}$$

However there is no reason why a constant value of  $h$  should be used in successive iterations. On the contrary it would seem desirable to vary  $h$  over the range between  $1/\lambda_1$  and  $1/\lambda_n$  in order to provide opportunities for coming close to the reciprocals of as many of the  $\lambda_i$  as possible. In this case we have instead of (5) and (6)

$$(7) \quad e^{(m+1)'} = \sum_i k_i \mu_i' \left[ \prod_{k=0}^m (1 - h_k \lambda_i) \right]$$

$$(8) \quad \delta^{(m)'} = \frac{e^{(m+1)'} - e^{(m)'}}{h_m} = -\sum_i \lambda_i \prod_{k=0}^{m-1} (1 - h_k \lambda_i) k_i \mu_i'$$

---

1/ Since this is less than unity, it follows from (5) that such an  $h$  is small enough to assure convergence if the original approximation is close enough to the maximum.

The procedure for estimating the smallest root  $\lambda_n$  from the successive approximations  $x^{(m)}$  is as follows, assuming  $\lambda_i > \lambda_n$  for all  $i \neq n$ . Let  $x_j^{(m)}$  be the  $j$ th element of  $x^{(m)}$ ,  $\delta_j^{(m)}$  the  $j$ th element of  $\delta^{(m)}$  and  $(\mu_i)_j$  the  $j$ th element of  $\mu_i$ . Then from (8)

$$\frac{\delta_j^{(m)}}{\delta_j^{(m-1)}} = \frac{-\sum_i k_i (\mu_i)_j \left[ \lambda_i \prod_{k=0}^{m-1} (1 - h_k \lambda_i) \right]}{-\sum_i k_i (\mu_i)_j \left[ \lambda_i \prod_{k=0}^{m-2} (1 - h_k \lambda_i) \right]}$$

Since  $\lambda_n$  is the smallest of the  $\lambda_i$  it follows that in both numerator and denominator the last term in the summation becomes dominant after many iterations with small enough positive  $h$ . For if we divide both numerator and denominator by  $\prod_{k=0}^{m-2} (1 - h_k \lambda_n)$  then the  $i$ th term of each will contain as a factor

$$\frac{\prod_{k=0}^{m-2} (1 - h_k \lambda_i)}{\prod_{k=0}^{m-2} (1 - h_k \lambda_n)} = \frac{1 - h_k \lambda_i}{1 - h_k \lambda_n} \leq \left( \frac{1 - h \lambda_i}{1 - h \lambda_n} \right)^{m-1}$$

where  $h$  is the smallest of the  $h_k$ . For  $m$  sufficiently large this becomes arbitrarily small,  $i \neq n$ . Since the remaining factors are bounded, it follows that for large  $m$  we may neglect all terms except the last in both numerator and denominator. Thus we have

$$(9) \quad \frac{\delta_j^{(m)}}{\delta_j^{(m-1)}} \approx \frac{(k_n \mu_n)_j \prod_{k=0}^{m-1} (1 - h_k \lambda_n)}{(k_n \mu_n)_j \prod_{k=0}^{m-2} (1 - h_k \lambda_n)} = (1 - h_{m-1} \lambda_n)$$

Since only  $\lambda_n$  is unknown we may solve this equation for the desired estimate of  $\lambda_n$ .

The estimation of  $\lambda_1$  is more difficult and much less precise. Given an approximation to  $[B^{-1} L(c)]^m e^{(0)'}$ , we may make use of the fact that for large  $m$

$$\begin{aligned} [B^{-1} L(c)]^m e^{(0)' } &= [B^{-1} L(c)]^m (k_1 \mu_1' + k_2 \mu_2' + \dots + k_n \mu_n') \\ &= \lambda_1^m k_1 \mu_1' + \lambda_2^m k_2 \mu_2' + \dots + \lambda_n^m k_n \mu_n' \\ &\approx \lambda_1^m k_1 \mu_1' \end{aligned}$$

since  $\lambda_1$  is the largest characteristic value. Then

$$(10) \quad \frac{j^{\text{th}} \text{ element of } [B^{-1} L(c)]^m e^{(0)' }}{j^{\text{th}} \text{ element of } [B^{-1} L(c)]^{m-1} e^{(0)' }} \approx \frac{\lambda_1^m (\mu_1)_j}{\lambda_1^{m-1} (\mu_1)_j} = \lambda_1$$

The approximation to  $[B^{-1} L(c)]^m e^{(0)'}$  may be obtained by a recursive process. We have

$$(11a) \quad (x^{(1)} - x^{(0)})' = e^{(1)' } - e^{(0)' } \approx [I - h_0 B^{-1} L(c)] e^{(0)' } - e^{(0)' } = -h_0 B^{-1} L(c) e^{(0)' }$$

$$(11b) \quad \begin{aligned} (x^{(2)} - x^{(1)})' &\approx -h_1 B^{-1} L(c) e^{(1)' } = -h_1 B^{-1} L(c) [I - h_0 B^{-1} L(c)] e^{(0)' } \\ &= -h_1 B^{-1} L(c) e^{(0)' } + h_0 h_1 [B^{-1} L(c)]^2 e^{(0)' } \end{aligned}$$

$$(11c) \quad \begin{aligned} (x^{(3)} - x^{(2)})' &\approx -h_2 B^{-1} L(c) [I - h_1 B^{-1} L(c)] [I - h_0 B^{-1} L(c)] e^{(0)' } \\ &= -h_2 B^{-1} L(c) e^{(0)' } + h_2 (h_0 + h_1) [B^{-1} L(c)]^2 e^{(0)' } \\ &\quad - h_0 h_1 h_2 [B^{-1} L(c)]^3 e^{(0)' } \end{aligned}$$

Etc.

From these we obtain the following relations by repeated substitution.

$$(12a) \quad B^{-1} L(c) e^{(0)' } \approx -\frac{1}{h_0} (x^{(1)} - x^{(0)})'$$



$$(12b) \quad [B^{-1}L(c)]^2 e^{(0)'} \approx -\frac{1}{h_0^2 h_1} [h_1(x^{(1)} - x^{(0)}) - h_0(x^{(2)} - x^{(1)})]'$$

$$(12c) \quad [B^{-1}L(c)]^3 e^{(0)'} \approx \frac{1}{h_0^3 h_1^2 h_2} [-h_0^2 h_1(x^{(3)} - x^{(2)})' + h_0 h_2 (h_0 + h_1)(x^{(2)} - x^{(1)})' - h_1^2 h_2 (x^{(1)} - x^{(0)})']$$

Etc.

It should be noted that (12b) depends on both 11b and 12a; (12c) depends on (11c) (12a) and (12b); and so on. Since all these expressions neglect terms in higher powers of the elements of  $e^{(0)}$ , this provides opportunity for the cumulation of errors as  $m$  increases and limits the accuracy with which  $\lambda_1$  may be estimated by this process.

#### 4. Newton Method

We have not yet discussed the choice of the matrix  $B$ . There may be some a priori reason for weighting errors in one parameter more strongly than errors in another. Barring this situation it is desirable to choose  $B$  in such a way as to speed convergence as much as possible. This may involve varying  $B$  between iterations. If  $B^{-1}L(c)$  is close to the identity matrix then its characteristic values, the  $\lambda_i$ , will all be close to unity. Since from (5)

$$e^{(m+1)'} = \sum_{i=1}^n (1-h\lambda_i)^{m+1} k_i \mu_i' + O(e^{(m)2})$$

it is clear that a choice of  $h$  close to unity will then quickly reduce all the terms in the summation on the right, permitting  $x^{(m)}$  to converge rapidly to  $c$ .

With this in mind it would seem desirable to set  $B = L(x^{(m)})$  since in the neighborhood of  $c$  we may expect  $L^{-1}(x^{(m)})L(c)$  to be close to the

identity matrix. Then we have

$$x^{(m+1)'} = x^{(m)'} + L^{-1}(x^{(m)}) \pm'(x^{(m)})$$

or  
 (13) 
$$\delta^{(m)'} = x^{(m+1)'} - x^{(m)'} = L^{-1}(x^{(m)}) \pm'(x^{(m)})$$

But this is equivalent to the Newton method of successive approximations, in which we approximate  $f(x)$  to second order terms of the Taylor expansion (about  $x^{(m)}$ ) and then chose  $x^{(m+1)}$  so as to maximize the quadratic thus obtained. For if we write

$$f(x) = f(x^{(m)} + \delta^{(m)}) \approx f(x^{(m)}) + \delta^{(m)} \pm'(x^{(m)}) - \frac{1}{2} \delta^{(m)} L(x^{(m)}) \delta^{(m)'}$$

and maximize with respect to  $\delta^{(m)}$  we obtain

$$\pm'(x^{(m)}) - L(x^{(m)}) \delta^{(m)'} = 0$$

and (13) follows immediately.

It should be noted that speed of convergence increases as we move closer to the maximum. This is obvious since the characteristic values,  $\lambda_i$ , are pushed together as  $L^{-1}(x^{(m)}) L(c)$  approaches the identity matrix. In fact it may be shown that in the neighborhood of  $(c)$ , the elements of  $e^{(m+1)}$  are of the order of magnitude of the squares of the elements of  $e^{(m)}$ . This property makes the Newton method valuable when we are very close to the maximum.

However, the computational cost of finding the second order terms for each iteration may be considerable. To lessen this cost while still retaining some of the advantages of the Newton method, the following modification may be introduced. The metric  $B$  (i.e., the second order matrix  $L$ ) may be held constant for several iterations. That is, we may use  $L(x^{(r)})$  to compute  $\delta^{(r+1)}$ ,  $\delta^{(r+2)}$ , ...,  $\delta^{(r+k)}$ , thus avoiding calculation of  $L(x^{(r+1)})$ ,  $L(x^{(r+2)})$ , ...,  $L(x^{(r+k)})$ . However there will be no change in the  $\lambda_i$  and no increase in the speed of convergence so long as the metric is held constant.

### 5. The "Finagle Factor"

Much of our discussion of simple gradient methods centered around the size of the step to be taken in the chosen direction, i.e., the choice of  $h$ . The Newton method in its ordinary form implies a value of  $h$  equal to 1. So long as this is true it will be found that the iterations tend to undershoot or overshoot the mark in a systematic fashion depending on third order terms. By observing successive iterations we may correct for this systematic tendency without actually computing the third order terms. The following procedure is applicable to the modified Newton method. We deal first with the one-dimensional case.

Let  $x^{(r)}$  be the last point for which the second order term was evaluated. Then using the expansions of  $f$  and its derivatives about  $c$ , we have for  $\delta^{(m)}$

$$(14) \quad \delta^{(m)} = -\frac{f'(x^{(m)})}{f''(x^{(r)})} = -\frac{e^{(m)}f''(c) + \frac{1}{2}e^{(m)2}f'''(c) + O(e^{(m)3})}{f''(c) + e^{(r)}f'''(c) + O(e^{(r)2})}$$

If  $x^{(m+1)} = x^{(m)} + \delta^{(m)}$ , we have

$$(15) \quad e^{(m+1)} = e^{(m)} + \delta^{(m)} = \frac{\frac{1}{2}e^{(m)}(2e^{(r)} - e^{(m)})f'''(c) + O(e^{(r)2}e^{(m)})}{f''(c) + e^{(r)}f'''(c) + O(e^{(r)2})}$$

$$= e^{(m)}(2e^{(r)} - e^{(m)}) \frac{f'''(c)}{2f''(c)} [1 + O(e^{(r)})]$$

Since  $e^{(m)}$  is very close to  $-\delta^{(m)}$  it is clear that we could substantially reduce the right hand side of (15) by adding to it

$$\delta^{(m)}(2e^{(r)} - e^{(m)}) \frac{f'''(c)}{2f''(c)}$$

Thus convergence would be speeded if we formed  $x^{(m+1)}$  by adding to  $x^{(m)}$ , not  $\delta^{(m)}$ , but  $\delta^{(m)}$  multiplied by an appropriate "finagle factor":

$$\tilde{x}^{(m+1)} = x^{(m)} + (1 + \phi^{(m)}) \delta^{(m)}$$

where

$$(16) \quad \phi^{(m)} = (2e^{(r)} - e^{(m)}) \frac{f'''(c)}{2f''(c)}$$

Then

$$\tilde{e}^{(m+1)} = e^{(m)} + \delta^{(m)}(1 + \phi^{(m)}) = O(e^{(m)} e^{(r)2})$$

Thus far  $\phi^{(m)}$  is unknown. A relatively unsophisticated method of approximating this factor is the following. Compute  $\delta^{(r)}$ ,  $x^{(r+1)} = x^{(r)} + \delta^{(r)}$  and  $\delta^{(r+1)}$ .

$$(17) \quad -\delta^{(r+1)} \approx e^{(r+1)} \approx e^{(r)} \cdot e^{(r)} \frac{f'''(c)}{2f''(c)} \approx \delta^{(r)2} \frac{f'''(c)}{2f''(c)}$$

This can be used to estimate  $f'''(c)/2f''(c)$ .

We note that if  $x^{(r)}$  is close to  $c$ ,  $e^{(r+1)}$ ,  $e^{(r+2)}$ ... will be small compared to  $e^{(r)}$  and while  $\phi^{(r)} = e^{(r)} \frac{f'''(c)}{2f''(c)}$ ,  $\phi^{(r+1)} \approx \phi^{(r+2)} \approx \dots \approx 2e^{(r)} \frac{f'''(c)}{2f''(c)}$ .

Approximating  $e^{(r)}$  by  $-\delta^{(r)}$  we have from 17

$$(18) \quad \phi^{(r+1)} \approx (-2\delta^{(r)}) - \frac{\delta^{(r+1)}}{\delta^{(r)2}} = \frac{2\delta^{(r+1)}}{\delta^{(r)}}$$

If  $x^{(r)}$  is not very close to  $c$ ,  $\phi^{(m)}$  will tend to change a little less rapidly from  $e^{(r)} \frac{f'''(c)}{2f''(c)}$  to  $2e^{(r)} \frac{f'''(c)}{f''(c)}$

For this reason one may find the following slightly more sophisticated approach useful especially in those cases where the work per iteration is considerable. This approach also makes it possible to estimate  $\frac{f'''(c)}{2f''(c)}$  without performing—as required by (17)—an iteration in which no "finable factor" is used.

Suppose that  $\hat{\phi}^{(m-1)}$  was used as an approximation to  $\phi^{(m-1)}$  in the  $m^{\text{th}}$  iteration. ( $\phi^{(m-1)}$  is never known exactly and may be approximated with varying degrees of accuracy.  $\hat{\phi}^{(m-1)}$  may even be zero if no previous knowledge concerning third order terms has been obtained). Thus  $x^{(m)} = x^{(m-1)} + \delta^{(m-1)}(1 + \hat{\phi}^{(m-1)})$

Then for  $m \geq r+1$ , we use (15) to obtain

$$(19) \quad e^{(m)} = e^{(m-1)} + \delta^{(m-1)} (1 + \hat{\phi}^{(m-1)}) e^{(m-1)} \left\{ [2e^{(r)} - e^{(m-1)}] \frac{f'''(c)}{2f''(c)} \right\} + \hat{\phi}^{(m-1)} \delta^{(m-1)} + O(e^{(r)2} e^{(m)})$$

Solving this to find an approximate value for  $\frac{f'''(c)}{2f''(c)}$  and substituting in

(16) we have

$$\phi^{(m)} \approx \frac{2e^{(r)} - e^{(m)}}{[2e^{(r)} - e^{(m-1)}]} \left[ \frac{e^{(m)} - \hat{\phi}^{(m-1)} \delta^{(m-1)}}{e^{(m-1)}} \right]$$

The error terms can be approximated by various methods of somewhat different accuracy. For example we may use (15) to note that

$$e^{(m)} = -\delta^{(m)} + O(e^{(m)} e^{(r)})$$

$$e^{(m-1)} = e^{(m-1)} - e^{(m)} + e^{(m)} = x^{(m-1)} - x^{(m)} - \delta^{(m)} + O(e^{(m)} e^{(r)})$$

$$2e^{(r)} - e^{(m)} = 2x^{(r)} - 2x^{(m)} - \delta^{(m)} + O(e^{(m)} e^{(r)})$$

$$2e^{(r)} - e^{(m-1)} = 2x^{(r)} - x^{(m)} - x^{(m-1)} - \delta^{(m)} + O(e^{(m)} e^{(r)})$$

Then

$$(20) \quad \phi^{(m)} \approx \hat{\phi}^{(m)} = \frac{2x^{(r)} - 2x^{(m)} - \delta^{(m)}}{2x^{(r)} - x^{(m)} - x^{(m-1)} - \delta^{(m)}} \cdot \frac{\delta^{(m)} + \hat{\phi}^{(m-1)} \delta^{(m-1)}}{x^{(m)} - x^{(m-1)} + \delta^{(m)}}$$

If after obtaining  $x^{(r_1)}$  ( $r_1 > r$ ), it is decided to recalculate second order terms we may proceed in the same manner to obtain

$$(20a) \quad \hat{\phi}^{(r_1)} = \frac{-\delta^{(r_1)}}{2x^{(r)} - x^{(r_1)} - x^{(r_1-1)} - \delta^{(r_1)}} \cdot \frac{\delta^{(r_1)} + \hat{\phi}^{(r_1-1)} \delta^{(r_1-1)}}{x^{(r_1)} - x^{(r_1-1)} + \delta^{(r_1)}}$$

In the  $n$ -dimensional case we are not in a position to determine the "finagle factor" so definitely or uniquely, since we have a vector rather than a scalar to eliminate. The best we can do is make use of an appropriate succession of  $h$ 's, which, as we might expect, are related to the characteristic roots of a certain matrix. Again let  $x^{(r)}$  be the last point for which the second order terms were evaluated.

Expanding  $f_i(x)$  and  $f_{ij}(x)$  in Taylor series about  $c$  we have

$$f_i(x) = \sum_{\mathbf{l}} e_{\mathbf{l}} f_{i\mathbf{l}}(c) + \frac{1}{2} \sum_{\mathbf{n}, \mathbf{p}} e_{\mathbf{n}} e_{\mathbf{p}} f_{i\mathbf{n}\mathbf{p}}(c) + \dots$$

$$f_{ij}(x) = f_{ij}(c) + \sum_{\mathbf{n}} e_{\mathbf{n}} f_{ij\mathbf{n}}(c) + \dots$$

We have furthermore <sup>1/</sup>

$$f_i^{\mathbf{r}}(x) = f_i^{\mathbf{r}}(c) - \sum_{\mathbf{k}} f_i^{\mathbf{r}\mathbf{k}}(c) \sum_{\mathbf{q}} e_{\mathbf{q}} \sum_{\mathbf{s}} f_{\mathbf{k}\mathbf{q}\mathbf{s}}(c) f_i^{\mathbf{s}\mathbf{j}}(c) + \dots$$

$$\delta_j^{(m)} = - \sum_{\mathbf{l}} f_i(x^{(m)}) f_{ij}(x^{(r)})$$

$$= - \sum_{\mathbf{l}, \mathbf{l}'} e_{\mathbf{l}}^{(m)} f_{i\mathbf{l}}(c) f_{ij}(c) - \frac{1}{2} \sum_{\mathbf{i}, \mathbf{n}, \mathbf{p}} e_{\mathbf{n}}^{(m)} e_{\mathbf{p}}^{(m)} f_{i\mathbf{n}\mathbf{p}}(c) f_{ij}(c)$$

$$+ \sum_{\mathbf{l}, \mathbf{l}'} e_{\mathbf{l}}^{(m)} f_{i\mathbf{l}}(c) \sum_{\mathbf{k}} f_i^{\mathbf{r}\mathbf{k}}(c) \sum_{\mathbf{q}} e_{\mathbf{q}}^{(r)} \sum_{\mathbf{s}} f_{\mathbf{k}\mathbf{q}\mathbf{s}}(c) f_i^{\mathbf{s}\mathbf{j}}(c) + \dots$$

$$= - e_j^{(m)} - \frac{1}{2} \sum_{\mathbf{i}, \mathbf{n}, \mathbf{p}} e_{\mathbf{n}}^{(m)} e_{\mathbf{p}}^{(m)} f_{i\mathbf{n}\mathbf{p}}(c) f_{ij}(c) + \sum_{\mathbf{k}, \mathbf{q}, \mathbf{s}} e_{\mathbf{k}}^{(m)} e_{\mathbf{q}}^{(r)} f_{\mathbf{k}\mathbf{q}\mathbf{s}}(c) f_i^{\mathbf{s}\mathbf{j}}(c) + \dots$$

$$= - e_j^{(m)} + \frac{1}{2} \sum_{\mathbf{i}, \mathbf{n}, \mathbf{p}} e_{\mathbf{n}}^{(m)} (2e_{\mathbf{p}}^{(r)} - e_{\mathbf{p}}^{(m)}) f_{i\mathbf{n}\mathbf{p}}(c) f_{ij}(c) + \dots$$

$$e_j^{(m+1)} = e_j^{(m)} + h_m \delta_j^{(m)}$$

$$= (1-h_m) e_j^{(m)} + \frac{h_m}{2} \sum_{\mathbf{i}, \mathbf{n}, \mathbf{p}} e_{\mathbf{n}}^{(m)} (2e_{\mathbf{p}}^{(r)} - e_{\mathbf{p}}^{(m)}) f_{i\mathbf{n}\mathbf{p}}(c) f_{ij}(c) + \dots$$

Suppose  $m = r$ . Then

$$\delta_j^{(r)} = -e_j^{(r)} + \frac{1}{2} \sum_{\mathbf{i}, \mathbf{n}, \mathbf{p}} e_{\mathbf{n}}^{(r)} e_{\mathbf{p}}^{(r)} f_{i\mathbf{n}\mathbf{p}}(c) f_{ij}(c) + \dots$$

$$e_j^{(r+1)} = (1-h_r) e_j^{(r)} + \frac{h_r}{2} \sum_{\mathbf{i}, \mathbf{n}, \mathbf{p}} e_{\mathbf{n}}^{(r)} e_{\mathbf{p}}^{(r)} f_{i\mathbf{n}\mathbf{p}}(c) f_{ij}(c) + \dots$$

If we let  $U = \|U_{jn}\|$  where

$$U_{jp} = \sum_{\mathbf{i}, \mathbf{n}} e_{\mathbf{n}}^{(r)} f_{i\mathbf{n}\mathbf{p}}(c) f_{ij}(c)$$

then we may write in vector notation

<sup>1/</sup> Here we are making use of the following facts. If  $L(x) = \|f_{ij}(x)\|$  then  $L(c+e) = L(c) + D(e)$  where  $D(e) = \| \sum_{\mathbf{n}} e_{\mathbf{n}} f_{ij\mathbf{n}}(c) + O(e_{\mathbf{p}} e_{\mathbf{q}}) + \dots \|$

For  $(e)$  sufficiently small

$$L^{-1}(x) = L^{-1}(c) - L^{-1}(c) D(e) L^{-1}(c) + O(D^2)$$

then we may write in vector notation

$$\begin{aligned} e^{(r+1)'} &= (I - h_r I) e^{(r)'} + \frac{h_r}{2} U e^{(r)'} + \dots \\ &= \left[ I - h_r \left( I - \frac{U}{2} \right) \right] e^{(r)'} + \dots \end{aligned}$$

Since the elements of  $e^{(r+1)'}$  are small compared with those of  $e^{(r)'}$  when we are close to the maximum we have for the next iteration

$$\begin{aligned} e_j^{(r+2)'} &\approx (1 - h_{r+1}) e_j^{(r+1)'} + \frac{h_{r+1}}{2} \sum_{i,n,p} e_n^{(r+1)'} \frac{\partial^2 e^{(r)'}}{\partial x_i \partial x_p} f_{ij}^{(c)} + \dots \\ e^{(r+2)'} &\approx \left[ I - h_{r+1} (I - U) \right] e^{(r+1)'} + \dots \end{aligned}$$

Similarly, for  $m > r+1$ , so long as  $L(x^{(r)})$  is used

$$e^{(m+1)'} \approx \left[ I - h_m (I - U) \right] e^{(m)'} \approx \prod_{p=r+1}^m \left[ I - h_p (I - U) \right] e^{(r+1)'}$$

In this case we may knock out the major part of  $e^{(m+1)'}$ , thus speeding convergence, if we choose  $h$ 's close to the reciprocals of the characteristic values of  $I-U$ , where  $U=O(e^{(r)'})$  is a matrix involving third order terms of the expansion of  $f$ . The procedure for determining appropriate values of  $h$  from successive iterations is very similar to that discussed earlier for other gradient methods. Whenever second order terms are recomputed  $U$  will, of course, change, and the relevant characteristic values will change accordingly. This was to be expected since, in the forms used in our discussion of gradient methods, a recomputation of second order terms means a change in the metric  $B$ . We have the additional information that the subsequent iterations without recalculation of second order terms have the effect of spreading the relevant characteristic values in a systematic fashion. For if  $\nu_i, i = 1, 2, \dots, n$  are the characteristic values of  $U$ , then  $1 - \frac{\nu_i}{2}$  are the characteristic values of  $I - \frac{U}{2}$  as compared with  $1 - \nu_i$  for  $I-U$ .

6. Concluding Remarks

The method proposed to find the smallest  $\lambda$  can be modified to eliminate dominating characteristic components with small  $\lambda$ . If several comparatively small h's are used the relative importance of characteristic vectors with large  $\lambda$  in the error terms is diminished. If initially the component with the smallest  $\lambda$  is very small it may be that after a few such iterations the main components will be those with  $\lambda$  close to a certain value, say  $\lambda_0$ . This will be apparent by noting that  $\delta^{(m)}$  is very close to a multiple of  $\delta^{(m-1)}$ .

$$\delta^{(m)} \approx (1 - h_{m-1} \lambda_0) \delta^{(m-1)}$$

Then it would be wise to use  $h_m = \frac{1}{\lambda_0}$ . This would tend to eliminate the largest components of the error. If this h is quite large, it may tend to revive components with large  $\lambda$ 's for  $1 - h_m \lambda$  would be highly negative.

It is therefore wise to follow a large h with several small ones. A

heuristic method of estimating  $\lambda_0$  is illustrated by the following. If

$\delta^{(m)} = .6 \delta^{(m-1)}$ ,  $h_{m-1}$  was only .4 as large as it should have been, therefore,  $h_m = \frac{10}{4} h_{m-1}$ . This method is equivalent to the above for if

$1 - h_{m-1} \lambda_0 = .6$ ,  $\lambda_0 = .4 / h_{m-1}$ ,  $h_m = 1 / \lambda_0 = \frac{10}{4} h_{m-1}$ . If this method is

applied to the Modified Newton Method, one should allow for the spread of the  $\lambda$  about 1 immediately after the second order terms are calculated.

For example, if  $h = 1$  is used when the second order terms are calculated

and one obtains  $\delta^{(r+1)} = .1 \delta^{(r)}$ ;  $h_r = 1$  is only .9 as large as it should

be; thus, we may prefer to use  $h_{r+1} = \frac{10}{9}$ . However, to allow for the spread  $h_{r+1}$  close to  $1 + \frac{2}{9}$  would be preferable.

The method of finding the largest characteristic root of a matrix on which is based the method of finding the smallest  $\lambda$  is a well known method. The Newton method and the gradient method with Euclidean Metric



$B = I$  have been known and used frequently in the past.<sup>1/</sup> Indeed Curry has noted that the iterations with different scales along the axes would give different results. Generalizations of these methods have been applied to functions defined on spaces which are not Euclidean.<sup>2/</sup>

The Cowles Commission has been using these methods for several years in Statistical problems of Econometrics.<sup>3/</sup> Indeed in these problems one can make use of the fact that for reasonably large samples, certain matrices which are relatively easy to compute are practically equivalent to the matrix of second order terms.

---

<sup>1/</sup> H. B. Curry, "The Method of Steepest Descent for Non-Linear Maximization Problems," Quarterly of Applied Mathematics, Oct. 1944. In this paper Curry also refers to papers of Cauchy, Courant, and Hadamard in which these methods have been used.

<sup>2/</sup> Wolfowitz has applied the generalization of the Newton method to a problem in Calculus of Variations in an unpublished paper.

<sup>3/</sup> Cowles Commission Monograph 10, Section 4 (to be published).