

Chapter 11

Numerical methods

Most numerical gradient-based methods require an initial design \underline{x}_0 , a search direction \underline{d}_0 , and a step size α_0 . In this way, the improved design can be defined as $\underline{x}_1 = \underline{x}_0 + \alpha_0 \underline{d}_0$. Iteratively, these algorithms may be written as

$$\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k. \quad (11.1)$$

Gradient-based numerical methods are characterized by the different ways they determine the search direction vector \underline{d}_k . The step size α_k is determined by one-dimensional optimization methods, also referred to as *line search methods*.

11.1 Steepest descent method

11.1.1 Definition

This method was presented by Cauchy in 1847. In the original definition, a current design \underline{x}_k is driven to an improved design \underline{x}_{k+1} such that

$$f(\underline{x}_{k+1}) < f(\underline{x}_k); \quad (11.2)$$

this is,

$$f(\underline{x}_k + \alpha_k \underline{d}_k) < f(\underline{x}_k). \quad (11.3)$$

Using Taylor series expansion about \underline{x}_k , the objective function f can be written as

$$f(\underline{x}) = f(\underline{x}_k) + \nabla f(\underline{x}_k)^T (\underline{x} - \underline{x}_k) + R_1(\underline{x}). \quad (11.4)$$

Evaluating in $\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k$, one obtains that

$$f(\underline{x}_{k+1}) = f(\underline{x}_k) + \nabla f(\underline{x}_k)^T (\alpha_k \underline{d}_k) + R_1(\alpha_k^2). \quad (11.5)$$

For α_k sufficiently small, the residual term can be disregarded. Since $f(\underline{x}_{k+1}) - f(\underline{x}_k) < 0$, see (11.2), one

observes that

$$\alpha_k \nabla f(\underline{x}_k)^T \underline{d}_k < 0. \quad (11.6)$$

For any $\alpha_k > 0$,

$$\nabla f(\underline{x}_k)^T \underline{d}_k < 0, \quad (11.7)$$

where \underline{d}_k is a *descent direction* and (11.7) is referred to as *descent condition*. If

$$\underline{d}_k = -\nabla f(\underline{x}_k), \quad (11.8)$$

then

$$\nabla f(\underline{x}_k)^T \underline{d}_k = -\|\nabla f(\underline{x}_k)\|^2 < 0. \quad (11.9)$$

The search direction in (11.8) is referred to as *steepest descent direction*.

11.1.2 Significance

In \mathbb{R}^n the gradient $\nabla f(\underline{x}_k)$ is a normal vector to the tangent hyperplane of the function isovalue at \underline{x}_k . The gradient vector is aligned with the direction of maximum growth steepest ascent of the function f at \underline{x}_k .

Example. Consider the function

$$f(\underline{x}) = 2x_1^2 + x_2^2.$$

Determine if the search direction $\underline{d}_0 = (0.5, 0.8)^T$ is a descent direction in $\underline{x}_0 = (1, 2)^T$. Determine the steepest descent direction in this point.

The gradient of the function is

$$\nabla f(\underline{x}) = \begin{pmatrix} 4x_1 \\ 2x_2 \end{pmatrix}.$$

Evaluating at $\underline{x}_0 = (1, 2)^T$ yields

$$\nabla f(\underline{x}_0) = \begin{pmatrix} 4 \\ 4 \end{pmatrix};$$

therefore, $\underline{d}_s = -(4, 4)^T$ is the steepest descent direction of f at \underline{x}_0 . The descent condition for $\underline{d}_0 = (0.5, 0.8)^T$ can be verified with the sign of the directional derivative of f_0 along \underline{d}_0 . This is

$$\begin{pmatrix} 4 & 4 \end{pmatrix} \begin{pmatrix} 0.5 \\ 0.8 \end{pmatrix} = 5.2 \not< 0,$$

then, \underline{d}_0 is not a descent direction. ■

11.1.3 Line search

Once the search direction \underline{d}_k in the point \underline{x}_k is identified, the algorithm has to identify the step size α_k . This problem can be stated as,

$$\min_{\alpha} f(\alpha) = f(\underline{x}_k + \alpha \underline{d}_k). \quad (11.10)$$

The necessary and sufficient condition is obtained from

$$\frac{df(\alpha)}{d\alpha} = \frac{d(\underline{x}_k + \alpha \underline{d}_k)}{d\alpha} = 0. \quad (11.11)$$

This is

$$\frac{df(\underline{x}_{k+1})}{d\alpha} = \frac{\partial f(\underline{x}_{k+1})}{\partial \underline{x}} \frac{d\underline{x}_{k+1}}{d\alpha} = \nabla f(\underline{x}_{k+1})^T \underline{d}_k = 0. \quad (11.12)$$

In other words, the directional derivative of $f(\underline{x}_{k+1})$ along \underline{d}_k must be zero. Using (11.8), one can observe that

$$\underline{d}_{k+1} = -\nabla f(\underline{x}_{k+1}), \quad (11.13)$$

therefore,

$$\frac{df(\underline{x}_{k+1})}{d\alpha} = -\underline{d}_{k+1}^T \underline{d}_k = 0. \quad (11.14)$$

Interestingly, since two consecutive steepest descent directions are orthogonal, this method approaches the minimum in a “zig-zag” shape.

11.1.4 Convergence criteria

The necessary condition for optimality is satisfied in a particular point \underline{x}_k if

$$\|\nabla f(\underline{x}_k)\| \leq \varepsilon_G, \quad (11.15)$$

where ε_G is a tolerance on the gradient supplied by the user. Notice that the gradient can vanish at any stationary point (maximum, minimum, or saddle point). However, the chances to find a maximum or a saddle point with a steepest descent algorithm are remote.

The algorithm should also check successive reduction in the function value, for example,

$$|f(\underline{x}_{k+1}) - f(\underline{x}_k)| \leq \varepsilon_A + \varepsilon_R |f(\underline{x}_k)|, \quad (11.16)$$

where ε_A is the absolute tolerance on the change in the function value while ε_R is the relative tolerance. Suggested values might be $\varepsilon_G = 10^{-4}$, $\varepsilon_A = 10^{-6}$, and $\varepsilon_R = 10^{-2}$.

11.1.5 Algorithm

Step 1. Estimate a reasonable initial point \underline{x}_0 and convergence parameters ε_A , ε_G , and ε_R , where ε_A is the absolute tolerance, ε_G and ε_R tolerances for the gradient and the function.

Step 2. Determine $\nabla f(\underline{x}_k)$. Stop if

$$\|\nabla f(\underline{x}_k)\| \leq \varepsilon_G.$$

Otherwise, define the search direction $\underline{d}_k = -\nabla f(\underline{x}_k)$.

Step 3. Obtain α_k minimizing $f(\alpha) = f(\underline{x}_k + \alpha \underline{d}_k)$, for $\alpha > 0$. Update $\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k$.

Step 4. Evaluate $f(\underline{x}_{k+1})$. Stop if

$$|f(\underline{x}_{k+1}) - f(\underline{x}_k)| \leq \varepsilon_A + \varepsilon_R |f(\underline{x}_k)|$$

is satisfied in two consecutive iterations. Otherwise, update $k = k + 1$ and $\underline{x}_k = \underline{x}_{k+1}$ and go to Step 2.

Example. Minimize the function

$$f(\underline{x}) = x_1^2 + x_2^2 - 2x_1x_2$$

using the steepest descent method from the point $\underline{x}_0 = (1, 0)^T$.

The search direction is $\underline{d}_0 = -(-2, 2)^T$. The step size minimizes $f(\alpha) = 18\alpha^2 - 8\alpha + 1$, this is, $\alpha_0 = 0.25$. The improved design is $\underline{x}_1 = (0.5, 0.5)^T$. In this point, the gradient of the function is $\underline{0}$ which satisfies the optimality condition and the convergence criterion. Therefore, this is the optimum point. ■

The steepest descent method is simple and robust to optimization. However, it has some drawbacks.

1. Even though this method is convergent, it is also slow even with positive definite quadratic forms.
2. The information from previous iterations is not used in the next steps.
3. In practice, the function is greatly reduced in the first iterations but it decreases mores slowly as the iterations continue.
4. The steepest descent direction makes sense from a local prospective (current point) but it can be improved in a more global sense.

11.1.6 Scaling

The order of convergence in this method is close to linear, even with some quadratic functions. That means

$$\lim_{k \rightarrow \infty} \frac{\|\underline{x}_{k+1} - \underline{x}^*\|}{\|\underline{x}_k - \underline{x}^*\|^p} < \infty, \quad (11.17)$$

where $p \rightarrow 1$. Using scaling techniques on the design variables it is possible to improve the order of convergence. A quadratic function can be scaled such that its Hessian is the identity matrix. In this case, the algorithm find the minimum in one iteration and the order of convergence is ∞ .

Example. Minimize the function

$$f(\underline{x}) = 25x_1^2 + x_2^2$$

using the steepest descent method form the point $\underline{x}_0 = (1, 1)^T$.

The gradient and the steepest descent direction is determined in every iteration. The line search is performed using a one-dimensional technique (e.g., golden section, polynomial approximations). Using the function `fminbnd` in MATLAB, this can be done as

```
function [xopt,fopt,aopt,ngradf]=afun(x0)
d0=-gfun('fun',x0);
[aopt,fopt] = fminbnd(@(a) fun(x0,d0,a),0,10);
xopt=x0+d0*aopt;
ngradf=norm(d0);

function f=fun(x0,d0,a)
if nargin==1
    d0=zeros(size(x0));
    a=0;
end
x=x0+a*d0;
f=25*x(1)^2+x(2)^2;
```

On the work space,

```
>> [x,f,a,n]=afun([1;1])
x =
    -0.0010
     0.9600
f =
     0.9215
a =
     0.0200
n =
    50.0400
```

Table 11.1 shows the iterative process.

Analytically, the Hessian is given by

$$\nabla^2 f(\underline{x}) = \begin{pmatrix} 50 & 0 \\ 0 & 2 \end{pmatrix}.$$

Table 11.1: Steepest descent, $f(x) = 25x_1^2 + x_2^2$

k	x_1	x_2	$f(\underline{x})$	α	$ \nabla f(\underline{x}) $
001	+1.0000E+00	+1.0000E+00	+2.6000E+01	+2.0033E-02	+5.0040E+01
002	-1.6444E-03	+9.5992E-01	+9.2154E-02	+4.7898E-02	+1.9216E+00
\vdots					
111	-2.3545E-06	+1.3752E-03	+1.8916E-02		+2.7531E-03

Introducing new variables $\hat{x} = (\hat{x}_1, \hat{x}_2)^T$ such that

$$\underline{x} = \underset{\sim}{D}\hat{x}, \quad (11.18)$$

where

$$\underset{\sim}{D} = \begin{pmatrix} \frac{1}{\sqrt{50}} & 0 \\ 0 & \frac{1}{\sqrt{2}} \end{pmatrix},$$

then $x_1 = \hat{x}_1/\sqrt{50}$, $x_2 = \hat{x}_2/\sqrt{2}$ and

$$f(\hat{x}) = \frac{1}{2}(\hat{x}_1^2 + \hat{x}_2^2).$$

The minimum of $f(\hat{x})$ can be located in one iteration using the steepest descent method. The point that minimizes f is $\hat{x}^* = (0, 0)^T$ in the scaled space. Then,

$$x_1^* = \hat{x}_1^*/\sqrt{50} = 0$$

and

$$x_2^* = \hat{x}_2^*/\sqrt{2} = 0.$$

■

Example. Minimize the function

$$f(\underline{x}) = 6x_1^2 - 6x_1x_2 + 2x_2^2 - 5x_1 + 4x_2 + 2$$

using the steepest descent method from the point $\underline{x}_0 = (-1, -2)^T$. Perform an appropriate scaling to improve the convergence rate of this method.

The Hessian of the function is given by

$$\nabla^2 f(\underline{x}) = \begin{pmatrix} 12 & -6 \\ -6 & 4 \end{pmatrix}.$$

The eigenvalues are $\lambda_1 = 0.7889$ and $\lambda_2 = 15.211$ with eigenvectors $\underline{e}_1 = (0.4718, 0.8817)^T$ and $\underline{e}_2 =$

$(-0.8817, 0.4718)^T$. The scaling is defined as

$$\underline{x} = \underline{Q} \underline{\hat{x}}, \quad (11.19)$$

where $\underline{Q} = (\underline{e}_1, \underline{e}_2)$, this is

$$\underline{Q} = \begin{pmatrix} 0.4718 & -0.8817 \\ 0.8817 & 0.4718 \end{pmatrix}.$$

Note that the Hessian of the scaled function is not the identity matrix yet. To accomplish this condition one needs a new change of variables

$$\underline{\hat{x}} = \underline{D} \underline{\hat{\hat{x}}}, \quad (11.20)$$

where

$$\underline{D} = \begin{pmatrix} \frac{1}{\sqrt{0.7889}} & 0 \\ 0 & \frac{1}{\sqrt{15.211}} \end{pmatrix}.$$

The Hessian of the new function $\underline{\hat{\hat{x}}}$ is, in fact, the identity matrix. Using the steepest descent method one reaches the optimum in one single iteration. The minimum is $\underline{\hat{\hat{x}}}^* = (-1.3158, -1.6142)^T$. Applying the linear transformations, the minimum in the original space is

$$\underline{x} = \underline{Q} \underline{D} \underline{\hat{\hat{x}}}, \quad (11.21)$$

and the solution is $\underline{x}^* = (-1/3, -2/3)^T$. ■

11.2 Conjugate gradient method

11.2.1 Definition

The conjugate gradient method was presented by Fletcher & Reeves (1964) as a quadratically convergent gradient method for locating an unconstrained local minimum of a function of several variables. With this method it is possible to locate the minimum of a quadratic function of n variables in n iterations. Consider the problem of minimizing a quadratic function

$$\min_{\underline{x}} \quad f(\underline{x}) = \frac{1}{2} \underline{x}^T \underline{A} \underline{x} + \underline{c}^T \underline{x} + b. \quad (11.22)$$

Assume that \underline{A} is symmetric ($\underline{A}^T = \underline{A}$) and positive definite ($\underline{x}^T \underline{A} \underline{x} > 0, \forall \underline{x} \in \mathbb{R}^n, \underline{x} \neq \underline{0}$). Then, two search directions, \underline{d}_i and \underline{d}_j , are conjugate with respect to \underline{A} if

$$\underline{d}_i^T \underline{A} \underline{d}_j = 0, \quad i \neq j. \quad (11.23)$$

Consider an initial design \underline{x}_0 and a set of conjugate directions $\underline{d}_0, \underline{d}_1, \dots, \underline{d}_{n-1}$. Minimizing $f(\underline{x})$ along \underline{d}_0 one obtains \underline{x}_1 . From the point \underline{x}_1 , $f(\underline{x})$ is minimized along \underline{d}_1 to obtain \underline{x}_2 . The process continues until one reaches \underline{x}_n along \underline{d}_{n-1} . The point \underline{x}_n minimizes (11.22). Denoting

$$\nabla f_k = \nabla f(\underline{x}_k), \quad (11.24)$$

then

$$\nabla f_k = \underline{x}_k^T \underline{\underline{A}} + \underline{c}. \quad (11.25)$$

If the descent direction \underline{d}_k is known, then

$$\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k, \quad (11.26)$$

where α_k is the solution that minimizes $f(\alpha)$. This is

$$\min_{\alpha} f(\underline{x}_k + \alpha \underline{d}_k), \quad (11.27)$$

where

$$f(\underline{x}_k + \alpha \underline{d}_k) = \frac{1}{2}(\underline{x}_k + \alpha \underline{d}_k)^T \underline{\underline{A}}(\underline{x}_k + \alpha \underline{d}_k) + \underline{c}^T(\underline{x}_k + \alpha \underline{d}_k). \quad (11.28)$$

The derivative with respect to α can be expressed as

$$\frac{df(\underline{x}_k + \alpha \underline{d}_k)}{d\alpha} = \frac{1}{2}(\underline{x}_k + \alpha \underline{d}_k)^T \underline{\underline{A}} \underline{d}_k + \frac{1}{2} \left[\underline{d}_k^T \underline{\underline{A}}(\underline{x}_k + \alpha \underline{d}_k) \right]^T + \underline{c}^T \underline{d}_k \quad (11.29)$$

or

$$\frac{df(\underline{x}_k + \alpha \underline{d}_k)}{d\alpha} = (\underline{x}_k + \alpha \underline{d}_k)^T \underline{\underline{A}} \underline{d}_k + \underline{c}^T \underline{d}_k. \quad (11.30)$$

The necessary (and sufficient) condition for optimality is satisfied when (11.30) is equal to zero. This condition can be written as

$$\underline{x}_k^T \underline{\underline{A}} \underline{d}_k + \alpha \underline{d}_k^T \underline{\underline{A}} \underline{d}_k + \underline{c}^T \underline{d}_k = 0, \quad (11.31)$$

this is

$$\alpha \underline{d}_k^T \underline{\underline{A}} \underline{d}_k = -\underline{x}_k^T \underline{\underline{A}} \underline{d}_k - \underline{c}^T \underline{d}_k. \quad (11.32)$$

Comparing (11.32) with (11.25) and using the symmetry of $\underline{\underline{A}}$ yields

$$\alpha \underline{d}_k^T \underline{\underline{A}} \underline{d}_k = -\underline{d}_k^T \nabla f_k. \quad (11.33)$$

Finally, solving for α allows

$$\alpha_k = -\frac{\underline{d}_k^T \nabla f_k}{\underline{d}_k^T \underline{\underline{A}} \underline{d}_k}. \quad (11.34)$$

In this method the search direction \underline{d}_{k+1} has the form

$$\underline{d}_{k+1} = -\nabla f_{k+1} + \beta_k \underline{d}_k, \quad (11.35)$$

where $\beta_k \underline{d}_k$ represents the deflection of the search direction with respect to the steepest descent direction. If two consecutive search directions are conjugate with respect to \underline{A} then

$$\underline{d}_{k+1}^T \underline{A} \underline{d}_k = 0. \quad (11.36)$$

Using (11.34), the conjugate condition (11.36) can be written as

$$(-\nabla f_{k+1} + \beta_k \underline{d}_k)^T \underline{A} \underline{d}_k = 0, \quad (11.37)$$

this is

$$-\nabla f_{k+1}^T \underline{A} \underline{d}_k + \beta_k \underline{d}_k^T \underline{A} \underline{d}_k = 0. \quad (11.38)$$

Solving for β_k ,

$$\beta_k = \frac{\nabla f_{k+1}^T \underline{A} \underline{d}_k}{\underline{d}_k^T \underline{A} \underline{d}_k}. \quad (11.39)$$

From (11.26) one observes that

$$\underline{d}_k = \frac{\underline{x}_{k+1} - \underline{x}_k}{\alpha_k}.$$

Multiplying by \underline{A} one obtains

$$\underline{A} \underline{d}_k = \frac{\underline{A} \underline{x}_{k+1} - \underline{A} \underline{x}_k}{\alpha_k}. \quad (11.40)$$

Using (11.25),

$$\nabla f_k = \underline{A} \underline{x}_k + \underline{c}$$

and

$$\nabla f_{k+1} = \underline{A} \underline{x}_{k+1} + \underline{c},$$

then (11.40) can be written as

$$\underline{A} \underline{d}_k = \frac{\nabla f_{k+1} - \nabla f_k}{\alpha_k}. \quad (11.41)$$

Replacing (11.41) into (11.39) yields

$$\beta_k = \frac{\nabla f_{k+1}^T (\nabla f_{k+1} - \nabla f_k)}{\alpha_k \underline{d}_k^T \underline{A} \underline{d}_k}. \quad (11.42)$$

From (11.35) one observes that

$$\underline{d}_k = -\nabla f_k + \beta_{k-1} \underline{d}_{k-1}. \quad (11.43)$$

Multiplying (11.43) by ∇f_k ,

$$\underline{d}_k^T \nabla f_k = -\nabla f_k^T \nabla f_k + \beta_{k-1} \underline{d}_{k-1}^T \nabla f_k. \quad (11.44)$$

The condition $df/d\alpha = 0$ make \underline{d}_k and ∇f_{k+1} orthogonal, this is

$$\underline{d}_k^T \nabla f_{k+1} = 0, \quad (11.45)$$

or

$$\underline{d}_{k-1}^T \nabla f_k = 0.$$

In this way, (11.44) can be expressed as

$$\underline{d}_k^T \nabla f_k = -\nabla f_k^T \nabla f_k. \quad (11.46)$$

In consequence, (11.34) can be written as

$$\alpha_k = \frac{\nabla f_k^T \nabla f_k}{\underline{d}_k^T \underline{\underline{A}} \underline{d}_k}. \quad (11.47)$$

Substituting (11.47) into (11.39) yields

$$\beta_k = \frac{\nabla f_{k+1}^T (\nabla f_{k+1} - \nabla f_k)}{\nabla f_k^T \nabla f_k}. \quad (11.48)$$

This expression is used by the *Polak-Ribiere* algorithm. On the other hand, if one considers

$$\nabla f_{k+1}^T \nabla f_k = 0, \quad (11.49)$$

then

$$\beta_k = \frac{\nabla f_{k+1}^T \nabla f_{k+1}}{\nabla f_k^T \nabla f_k}. \quad (11.50)$$

This expression is used by the *Fletcher-Reeves* algorithm.

11.2.2 Algorithm

Step 1. Estimate the initial design \underline{x}_0 . Select a convergence parameter ε . The first search direction is defined as the steepest descent direction,

$$\underline{d}_0 = -\nabla f_0.$$

If $\|\nabla f_0\| < \varepsilon$, then convergence is obtained and algorithm cannot improve the current design, otherwise go to Step 5.

Step 2. If $\|\nabla f_k\| < \varepsilon$ then convergence is obtained and the algorithm cannot improve the current design, otherwise go to the next Step.

Step 3. Determine the deflection factor using Polak-Ribiere equation (11.48) or Fletcher-Reeves equation (11.50). In the second case, the deflection factor can be expressed as

$$\beta_k = \frac{\nabla f_{k+1}^T \nabla f_{k+1}}{\nabla f_k^T \nabla f_k}$$

Step 4. Update the search direction

$$\underline{d}_{k+1} = -\nabla f_{k+1} + \beta_k \underline{d}_k.$$

Step 5. Determine the step size by performing a line search. In a quadratic function, the step size is given by

$$\alpha_k = \frac{\nabla f_k^T \nabla f_k}{\underline{d}_k^T \underline{A} \underline{d}_k}.$$

If the function is not quadratic, α is determined by minimizing $f(\alpha)$.

Step 6. Update the design

$$\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k$$

and go to Step 2.

If the minimum is not found after $n + 1$ iterations, it is recommended to restart the algorithm using the steepest descent direction.

Example. Consider the minimization of the function

$$f(\underline{x}) = x_1^2 + 2x_2^2 + 2x_3^2 + 2x_1x_2 + 2x_2x_3$$

from the initial design $\underline{x}_0 = (2, 4, 10)^T$. Using the conjugate gradient method, how many iterations would it take to find the minimum? Perform that many iterations and compare your numerical solution with the analytical solution.

Since this is a quadratic problem, we can use the Hessian to update the step size. The algorithm should converge in three iterations. The gradient of the function is given by

$$\nabla f(\underline{x}) = \begin{pmatrix} 2x_1 + 2x_2 \\ 2x_1 + 4x_2 + 2x_3 \\ 2x_2 + 4x_3 \end{pmatrix}$$

and the Hessian is

$$\underline{\underline{A}} = \nabla^2 f(\underline{x}) = \begin{pmatrix} 2 & 2 & 0 \\ 2 & 4 & 2 \\ 0 & 2 & 4 \end{pmatrix}$$

Iteration 1. The first improved point is $\underline{x}_1 = \underline{x}_0 + \alpha_0 \underline{d}_0$ where

$$\underline{d}_0 = -\nabla f_0 = \begin{pmatrix} -12 \\ -40 \\ -48 \end{pmatrix}$$

Since f is quadratic then the step size is

$$\alpha_0 = \frac{\nabla f_0^T \nabla f_0}{\underline{d}_0^T \underline{\underline{A}} \underline{d}_0} = \frac{253}{1594} = 0.1587.$$

Then,

$$\underline{x}_1 = \underline{x}_0 + \alpha_0 \underline{d}_0 = \begin{pmatrix} 0.0954 \\ -2.3488 \\ 2.3814 \end{pmatrix}$$

Iteration 2. The second improved point is $\underline{x}_2 = \underline{x}_1 + \alpha_1 \underline{d}_1$ where

$$\underline{d}_1 = -\nabla f_1 + \beta_0 \underline{d}_0$$

Using Fletcher-Reeves equation

$$\beta_0 = \frac{\nabla f_1^T \nabla f_1}{\nabla f_0^T \nabla f_0},$$

where

$$\nabla f_1 = \begin{pmatrix} -4.5069 \\ -4.4417 \\ 4.8281 \end{pmatrix};$$

then

$$\beta_0 = 0.01565$$

and

$$\underline{d}_1 = \begin{pmatrix} 4.3191 \\ 3.81566 \\ -5.5793 \end{pmatrix}.$$

Now, as before

$$\alpha_1 = \frac{\nabla f_1^T \nabla f_1}{\underline{d}_1^T \underline{\tilde{A}} \underline{d}_1} = 0.3155,$$

and

$$\underline{x}_2 = \underline{x}_1 + \alpha_1 \underline{d}_1 = \begin{pmatrix} 1.4578 \\ -1.1452 \\ 0.6214 \end{pmatrix}$$

Iteration 3. The third (hopefully last) design is determine in the same way:

$$\nabla f_2 = \begin{pmatrix} 0.62534 \\ -0.422104 \\ 0.195419 \end{pmatrix},$$

$$\beta_1 = \frac{\nabla f_2^T \nabla f_2}{\nabla f_1^T \nabla f_1} = 0.0096,$$

$$\underline{d}_2 = -\nabla f_2 + \beta_1 \underline{d}_1 = \begin{pmatrix} -0.5839 \\ 0.4587 \\ -0.2489 \end{pmatrix},$$

$$\alpha_2 = \frac{\nabla f_2^T \nabla f_2}{\underline{d}_2^T \underline{\tilde{A}} \underline{d}_2} = 2.4966,$$

and

$$\underline{x}^* = \underline{x}_3 = \underline{x}_1 + \alpha_1 \underline{d}_1 = \begin{pmatrix} 0.0000 \\ 0.0000 \\ 0.0000 \end{pmatrix}.$$

The minimum value of the function is $f(\underline{x}^*) = 0$.

Notice that the convergence criterion is satisfied, this is

$$\nabla f_3 = \begin{pmatrix} 0.0000 \\ 0.0000 \\ 0.0000 \end{pmatrix}.$$

The reader can easily prove that this corresponds to the analytical solution. ■

The conjugate gradient method is in essence the steepest descent with a deflected search direction. This simple modification does not require a considerably amount of computational work; however, it substantially improves the rate of convergence.

11.3 Newton's method

11.3.1 Definition

When second-order derivatives are available, the objective function can be better approximated. In this way, better search directions can be obtained and the convergence rate can be also improved. Newton's method makes use of the Hessian of the objective function and has quadratic convergence. When the method is applied to a quadratic function that is positive definite, it locates the minimum in one iteration.

Newton's method rely on a quadratic approximation of the function about a current design \underline{x}_k ,

$$f_Q(\underline{x}) = f(\underline{x}_k) + \nabla f(\underline{x}_k)^T(\underline{x} - \underline{x}_k) + \frac{1}{2}(\underline{x} - \underline{x}_k)^T \nabla^2 f(\underline{x}_k)(\underline{x} - \underline{x}_k) \quad (11.51)$$

and its gradient

$$\nabla f_Q(\underline{x}) = \nabla f(\underline{x}_k) + \nabla^2 f(\underline{x}_k)(\underline{x} - \underline{x}_k). \quad (11.52)$$

The minimum of f_Q satisfies that $\nabla f_Q(\underline{x}_Q^*) = \underline{0}$, this is

$$\nabla f_Q(\underline{x}_Q^*) = \nabla f(\underline{x}_k) + \nabla^2 f(\underline{x}_k) \underbrace{(\underline{x}_Q^* - \underline{x}_k)}_{\underline{d}_k} = \underline{0}, \quad (11.53)$$

then

$$\nabla^2 f(\underline{x}_k) \underline{d}_k = -\nabla f(\underline{x}_k) \quad (11.54)$$

Assuming that $\nabla^2 f(\underline{x}_k)$ is not singular, and therefore invertible, one observes that

$$\underline{d}_k = -[\nabla^2 f(\underline{x}_k)]^{-1} \nabla f(\underline{x}_k). \quad (11.55)$$

The descent condition (11.7) is satisfied if

$$-\nabla f_k^T [\nabla^2 f_k]^{-1} \nabla f_k < 0, \quad (11.56)$$

which is accomplished if and only if $\nabla^2 f(\underline{x}_k)$ is *positive definite*. Defining the search direction as

$$\underline{d}_k = \underline{x}_{k+1} - \underline{x}_k \quad (11.57)$$

the updating rule of the Newton's method can be expressed as

$$\underline{x}_{k+1} = \underline{x}_k - [\nabla^2 f(\underline{x}_k)]^{-1} \nabla f(\underline{x}_k). \quad (11.58)$$

Some drawbacks of the Newton's method include:

1. Each iteration in the Newton's method requires the Hessian of the objective function. In some func-

tions that might be impossible for certain points; however, when it is possible, the number of function calls increases substantially with respect to the previous first-order methods. For instance, if $f : \mathbb{R}^n \rightarrow \mathbb{R}$, then the algorithm requires $n(n+1)/2$ second-order derivatives.

2. The Hessian might be singular (or close to singular) in some points which makes impossible to determine a search direction.
3. The method does not use information collected in previous iterations.
4. The convergence of this method is not guaranteed.

11.3.2 Modified Newton's methods

When the function is highly nonlinear, then a quadratic approximation may not be good enough. In fact, it might happen that $f(\underline{x}_{k+1}) \not\leq f(\underline{x}_k)$ and convergence cannot be achieved. A modified Newton's method incorporates a line search to find a step size α_k over the search direction \underline{d}_k . This approach is referred to as *modified Newton's method*.

Marquardt (1964) suggested a modification on Newton's method by proposing a new way to find search directions. In this approach, Marquardt combines the advantages of the Newton's method with the steepest descent direction method. The proposed search direction is defined as

$$\underline{d}_k = -(\nabla^2 f(\underline{x}_k) + \gamma \underline{I})^{-1} \nabla f(\underline{x}_k), \quad (11.59)$$

where \underline{I} is the identity matrix and γ is a weighting value. When γ is big, then the direction corresponds to the steepest descent direction. When γ is small, then it corresponds to the Newton's direction. This approach is referred to as *Marquardt's modification*.

11.3.3 Algorithm

The algorithm for the modified Newton's method can be described as follows.

- Step 1. Estimate an initial design \underline{x}_0 and a convergence parameter ε .
- Step 2. Determine $\nabla f(\underline{x}_k)$. If $\|\nabla f(\underline{x}_k)\| < \varepsilon$, then convergence is obtained and algorithm cannot improve the current design, otherwise go to the next Step.
- Step 3. Determine the Hessian $\nabla^2 f(\underline{x}_k)$.
- Step 4. Determine the search direction

$$\underline{d}_k = -[\nabla^2 f(\underline{x}_k)]^{-1} \nabla f(\underline{x}_k).$$

Step 5. Update the design

$$\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k,$$

where α_k is the minimum of the function $f(\underline{x}_k + \alpha \underline{d}_k)$, then go to Step 2.

Example. Given the function

$$f(\underline{x}) = 10x_1^4 - 20x_1x_2 + 10x_2^2 + x_1^2 - 2x_1 + 5,$$

perform two iterations using the modified Newton's method from the point $\underline{x}_0 = (-1, 3)^T$.

This is not a quadratic function, therefore one cannot expect to find the minimum in one single iteration. The gradient of the function is given by

$$\nabla f(\underline{x}) = \begin{pmatrix} -2 + 2x_1 + 40x_1^3 - 20x_2 \\ -20x_1 + 20x_2 \end{pmatrix}$$

and the Hessian is given by

$$\nabla^2 f(\underline{x}) = \begin{pmatrix} 2 + 120x_1^2 & -20 \\ -20 & 20 \end{pmatrix}.$$

Iteration 1. The first improved point is $\underline{x}_1 = \underline{x}_0 + \alpha_0 \underline{d}_0$ where

$$\underline{d}_0 = -[\nabla^2 f(\underline{x}_0)]^{-1} \nabla f(\underline{x}_0) = \begin{pmatrix} 0.2353 \\ -3.7647 \end{pmatrix}$$

Minimizing $f(\alpha) = f(\underline{x}_0 + \alpha \underline{d}_0)$

$$\alpha_0 = 1.0045$$

then,

$$\underline{x}_1 = \underline{x}_0 + \alpha_0 \underline{d}_0 = \begin{pmatrix} -0.7637 \\ -0.7815 \end{pmatrix}$$

Iteration 2. Following the same procedure,

$$\underline{d}_1 = -[\nabla^2 f(\underline{x}_1)]^{-1} \nabla f(\underline{x}_1) = \begin{pmatrix} 0.1167 \\ 0.1346 \end{pmatrix}$$

$$\alpha_1 = 1.3424$$

$$\underline{x}_2 = \underline{x}_1 + \alpha_1 \underline{d}_1 = \begin{pmatrix} -0.6070 \\ -0.6008 \end{pmatrix}$$

The exact solution is given by $\underline{x}^* = (0.7207, 0.7207)^T$ and $f(\underline{x}^*) = 1.5818$.

11.4 Quasi-Newton methods

Quasi-Newton method's approximate the Hessian using first-order derivatives. For positive definite quadratic functions of n variables, these method achieve convergence in n iterations. If the minimum is not found after $n + 1$, then the algorithm restarts. In these methods the search direction is given by

$$\underline{d}_k = -\underline{\underline{H}}_k \nabla f(\underline{x}_k), \quad (11.60)$$

where $\underline{\underline{H}}_k$ is an approximation of $[\nabla^2 f(\underline{x}_k)]^{-1}$. Initially, this approximation is the identity matrix, $\underline{\underline{H}}_0 = \underline{\underline{I}}$.

11.4.1 Davidon-Fletcher-Powell (DFP) method

This method was proposed by Davidon (1959) and modified by Fletcher & Powell (1963). In the DFP method, the approximation of the inverse of the Hessian is given by

$$\underline{\underline{H}}_{k+1} = \underline{\underline{H}}_k - \frac{\underline{\underline{H}}_k \underline{y}_k \underline{y}_k^T \underline{\underline{H}}_k}{\underline{y}_k^T \underline{\underline{H}}_k \underline{y}_k} + \frac{\underline{s}_k \underline{s}_k^T}{\underline{s}_k^T \underline{y}_k}, \quad (11.61)$$

where $\underline{s}_k = \underline{x}_{k+1} - \underline{x}_k$ represents the change in the design and $\underline{y}_k = \nabla f(\underline{x}_{k+1}) - \nabla f(\underline{x}_k)$ represents the change in the gradient.

Example. Minimize the function

$$f(\underline{x}) = 5x_1^5 + 2x_1x_2 + x_2^2 + 7$$

from $\underline{x}_0 = (1, 2)^T$. Use $\varepsilon = 0.001$ as convergence parameter.

$$\underline{\underline{H}}_0 = \underline{\underline{I}} \text{ then}$$

$$\underline{d}_0 = -\nabla f_0 = \begin{pmatrix} 14 \\ 6 \end{pmatrix}.$$

The norm $\|\nabla f_0\| = 15.232 > \varepsilon$, then there is no convergence and the algorithm continues. The updated design is given by

$$\underline{x}_1 = \underline{x}_0 + \alpha_0 \underline{d}_0$$

where $\alpha_0 = 0.0988$. Por lo tanto

$$\underline{x}_1 = \begin{pmatrix} -1.046 \\ 2.042 \end{pmatrix}.$$

In the second iteration, $\|\nabla f_1\| = 2.29 > \varepsilon$, therefore the is no convergence. The algorithm determines

$\underline{s}_0 = \underline{x}_1 - \underline{x}_0 = (-1.386, -0.593)^T$ and $\underline{y}_0 = \nabla f_1 - \nabla f_0 = (-15.046, -3.958)^T$. Then,

$$\underline{\underline{H}}_1 = \begin{pmatrix} 0.148 & -0.211 \\ -0.211 & 0.950 \end{pmatrix}.$$

The new search direction is

$$\underline{d}_1 = \begin{pmatrix} 0.586 \\ -1.719 \end{pmatrix}.$$

The updated design is

$$\underline{x}_2 = \underline{x}_1 + \alpha_1 \underline{d}_1$$

where $\alpha_1 = 0.776$. Therefore

$$\underline{x}_2 = \begin{pmatrix} 0.069 \\ 0.073 \end{pmatrix}.$$

■

11.4.2 Broyden-Fletcher-Goldfarb-Shanno (BFGS) method

In the BFGS, the approximation of the inverse of the Hessian is given by

$$\underline{\underline{H}}_{k+1} = \underline{\underline{H}}_k - \frac{\underline{s}_k \underline{y}_k^T \underline{\underline{H}}_k + \underline{\underline{H}}_k \underline{y}_k \underline{s}_k^T}{\underline{s}_k^T \underline{y}_k} + \left(1 + \frac{\underline{y}_k^T \underline{\underline{H}}_k \underline{y}_k}{\underline{s}_k^T \underline{y}_k} \right) \frac{\underline{s}_k \underline{s}_k^T}{\underline{s}_k^T \underline{y}_k}, \quad (11.62)$$

where $\underline{s}_k = \underline{x}_{k+1} - \underline{x}_k$ represents the change in the design variables and $\underline{y}_k = \nabla f_{k+1} - \nabla f_k$ represents the change in the gradient.

Example. Minimize the function

$$f(\underline{x}) = 5x_1^5 + 2x_1x_2 + x_2^2 + 7$$

from $\underline{x}_0 = (1, 2)^T$. Use $\varepsilon = 0.001$ as convergence parameter.

The first iteration is just like in the previous example, this is $\underline{\underline{H}}_0 = \underline{\underline{I}}$ and

$$\underline{x}_1 = \begin{pmatrix} -1.046 \\ 2.042 \end{pmatrix}.$$

In the second iteration one obtains that

$$\underline{\underline{H}}_1 = \begin{pmatrix} 1.770 & -7.506 \\ -7.506 & 33.750 \end{pmatrix}$$

and the new search direction is

$$\underline{d}_1 = \begin{pmatrix} 17.20 \\ -76.77 \end{pmatrix}.$$

The updated design is given by

$$\underline{x}_2 = \underline{x}_1 + \alpha_1 \underline{d}_1$$

where $\alpha_1 = 0.0185$. Therefore

$$\underline{x}_2 = \begin{pmatrix} -0.0688 \\ -0.0098 \end{pmatrix}.$$

■

11.5 Trust regions methods

11.5.1 Definition

The basic idea under this methods is to approximate the objective function f with a simpler expression f_A that reasonably reflects its behavior in the neighborhood of a current design \underline{x}_k . This neighborhood is referred to as the trust region. If f is highly nonlinear the approximation f_A will be valid in a *trust region* Ω_k around \underline{x}_k . This region can be described as

$$\Omega_k = \{\underline{x} : \|\underline{x} - \underline{x}_k\|_\infty \leq h_k\}, \quad (11.63)$$

where h_k is the size of the trust region, which is dynamically adjusted. The sub-optimization problem can be stated as

$$\begin{aligned} \min_{\underline{d}} \quad & f_A(\underline{d}) \\ \text{s.a.} \quad & -h_k \leq d_i \leq h_k, \quad i = 1, \dots, n. \end{aligned} \quad (11.64)$$

The solution of this problem is \underline{d}_k . Notice that these methods do not require line search as the size of the trust region accounts for that.

If the approximation is linear, then the sub-optimization problem to locate \underline{d}_k within Ω_k can be described as

$$\begin{aligned} \min_{\underline{d}} \quad & f_L(\underline{d}) = f(\underline{x}_k) + \nabla f(\underline{x}_k)^T \underline{d} \\ \text{s.a.} \quad & -h_k \leq d_i \leq h_k, \quad i = 1, \dots, n. \end{aligned} \quad (11.65)$$

This type of problems with linear objective and linear constraints is known as *linear programming* (LP) problem.

If the approximation is quadratic, then the sub-optimization problem is

$$\begin{aligned} \min_{\underline{d}} \quad & f_Q(\underline{d}) = f(\underline{x}_k) + \nabla f(\underline{x}_k)^T \underline{d} + \frac{1}{2} \underline{d}^T \nabla^2 f(\underline{x}_k) \underline{d} \\ \text{s.a.} \quad & -h_k \leq d_i \leq h_k, \quad i = 1, \dots, n. \end{aligned} \quad (11.66)$$

This type of problems with a quadratic objective function and linear constraints is known as *quadratic programming* (QP) problem.

11.5.2 Reliability index

The reliability index is defined as

$$r_k = \frac{f(\underline{x}_k) - f(\underline{x}_k + \underline{d}_k)}{f(\underline{x}_k) - f_A(\underline{x}_k + \underline{d}_k)} \quad (11.67)$$

where f is the objective function to minimize and f_A is its approximation. If the numerator and the denominator are greater than zero and $r_k \approx 1$, then it is said that there is a good agreement between the function and its approximation. On the hand, if the numerator is negative, then r_k is also negative, which is not desired. In that case, the new point $\underline{x}_{k+1} = \underline{x}_k + \underline{d}_k$ does not decrease the function so the size of the trust region h_k has to be adjusted according to the value of r_k . A heuristic updating rule is as follows:

$$\begin{aligned} \text{If } r_k < 0.25 & \Rightarrow h_{k+1} = \frac{\|\underline{d}_k\|_\infty}{4} \\ \text{If } r_k > 0.75 \text{ y } h_k = \|\underline{d}_k\|_\infty & \Rightarrow h_{k+1} = 2h_k \\ \text{Otherwise} & \Rightarrow h_{k+1} = h_k \end{aligned} \quad (11.68)$$

11.5.3 Algorithm

Step 1. Estimate \underline{x}_0 , $h_0 = 1$ and a convergence parameter ε .

Step 2. Determine $\nabla f(\underline{x}_k)$ and $\nabla^2 f(\underline{x}_k)$. If $\|\nabla f(\underline{x}_k)\| \leq \varepsilon$ then convergence has been achieved and the minimum has been found. Otherwise, continue to the next Step.

Step 3. Obtain \underline{d}_k solving the sub-optimization problem (11.64).

Step 4. Evaluar $f(\underline{x}_{k+1})$ and obtain r_k according to (11.67).

Step 5. Update the size of the trust region h_{k+1} according to (11.68).

Step 6. Determine the new point \underline{x}_{k+1} according to the following rule:

$$\begin{aligned} \text{If } r_k \leq 0 & \Rightarrow \underline{x}_{k+1} = \underline{x}_k \text{ and go to Step 3} \\ \text{Otherwise} & \Rightarrow \underline{x}_{k+1} = \underline{x}_k + \underline{d}_k \text{ and go to Step 2} \end{aligned}$$

11.6 Least square error

Consider the system of linear equations

$$\underline{\underline{A}}\underline{\underline{x}} = \underline{\underline{b}}. \quad (11.69)$$

where the matrix $\underline{\underline{A}} \in \mathbb{R}^{m \times n}$ applies the vector $\underline{\underline{x}} \in \mathbb{R}^n$ into $\underline{\underline{b}} \in \mathbb{R}^m$. When the system is inconsistent, this is $\text{rank}(\underline{\underline{A}}) < \text{rank}(\underline{\underline{A}}, \underline{\underline{b}})$, it does not have a solution. However, one can attempt to provide an approximate point that minimizes $\|\underline{\underline{A}}\underline{\underline{x}} - \underline{\underline{b}}\|$. This is,

$$\min_{\underline{\underline{x}}} f(\underline{\underline{x}}) = (\underline{\underline{A}}\underline{\underline{x}} - \underline{\underline{b}})^T(\underline{\underline{A}}\underline{\underline{x}} - \underline{\underline{b}}) \quad (11.70)$$

The necessary condition for optimality can be expressed as $\nabla f(\underline{\underline{x}}^*) = \underline{\underline{0}}$, where

$$\begin{aligned} \nabla f(\underline{\underline{x}}) &= \underline{\underline{A}}^T \underline{\underline{A}}\underline{\underline{x}} - \underline{\underline{A}}^T \underline{\underline{b}} + \underline{\underline{A}}^T \underline{\underline{A}}\underline{\underline{x}} - \underline{\underline{A}}^T \underline{\underline{b}} \\ &= 2\underline{\underline{A}}^T \underline{\underline{A}}\underline{\underline{x}} - 2\underline{\underline{A}}^T \underline{\underline{b}}. \end{aligned}$$

In the optimal point,

$$\underline{\underline{A}}^T \underline{\underline{A}}\underline{\underline{x}}^* = \underline{\underline{A}}^T \underline{\underline{b}}. \quad (11.71)$$

Solving for $\underline{\underline{x}}^*$ yields

$$\underline{\underline{x}}^* = (\underline{\underline{A}}^T \underline{\underline{A}})^{-1} \underline{\underline{A}}^T \underline{\underline{b}}. \quad (11.72)$$

This solution is referred to as *least square error*. Notice that this is a convex problem, therefore $\underline{\underline{x}}^*$ is a global minimum.

Example. Solve for $\underline{\underline{x}}$ in the following system of linear equations,

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}.$$

Since $\text{rank}(\underline{\underline{A}}) = 2$ and $\text{rank}(\underline{\underline{A}}, \underline{\underline{b}}) = 3$ the system is inconsistent and does not have a solution. However, the solution that minimizes the quadratic error is obtained using (11.72), this is

$$\begin{pmatrix} x_1^* \\ x_2^* \end{pmatrix} = \left[\begin{pmatrix} 1 & 3 & -1 \\ 2 & 4 & 1 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 3 & 4 \\ -1 & 1 \end{pmatrix} \right]^{-1} \begin{pmatrix} 1 & 3 & -1 \\ 2 & 4 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \quad (11.73)$$

$$= \begin{pmatrix} -0.4516 \\ 0.6129 \end{pmatrix}. \quad (11.74)$$

■

One of the most important applications of this method is the identification of the coefficients of a function that better matches a cloud of points.

Example. An experiment yields the data in the following table.

t	x
0.00	1.001
0.10	1.089
0.23	1.240
0.70	1.604
0.90	1.738
1.50	2.020
2.65	1.412
3.00	1.241

Find the least square best fit coefficients a , b , and c if the assumed functional form is

(a) $x = a + bt + ct^2$,

(b) $x = a + b \sin(t) + c \sin(2t)$.

Which best fit estimate has the smallest least square error?

(a) The system of linear equations can be expressed as

$$\begin{pmatrix} 1.000 & 0.000 & 0.000 \\ 1.000 & 0.100 & 0.010 \\ 1.000 & 0.230 & 0.053 \\ 1.000 & 0.700 & 0.490 \\ 1.000 & 0.900 & 0.810 \\ 1.000 & 1.500 & 2.250 \\ 1.000 & 2.650 & 7.023 \\ 1.000 & 3.000 & 9.000 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 1.001 \\ 1.089 \\ 1.240 \\ 1.604 \\ 1.738 \\ 2.020 \\ 1.412 \\ 1.241 \end{pmatrix}.$$

Using (11.72), the solution can be expressed as

$$\begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 0.982 \\ 1.205 \\ -0.379 \end{pmatrix}.$$

The function is $x(t) = 0.982 + 1.205t - 0.379t^2$. The quadratic error is 0.0224.

(b) In this case, the system can be expressed as

$$\begin{pmatrix} 1.000 & 0.000 & 0.000 \\ 1.000 & 0.100 & 0.199 \\ 1.000 & 0.228 & 0.444 \\ 1.000 & 0.644 & 0.985 \\ 1.000 & 0.783 & 0.974 \\ 1.000 & 0.997 & 0.141 \\ 1.000 & 0.472 & -0.832 \\ 1.000 & 0.141 & -0.279 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 1.001 \\ 1.089 \\ 1.240 \\ 1.604 \\ 1.738 \\ 2.020 \\ 1.412 \\ 1.241 \end{pmatrix}.$$

Using (11.72), the solution can be expressed as

$$\begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 1.017 \\ 0.960 \\ -0.011 \end{pmatrix}.$$

The function is $x(t) = 1.017 + 0.960 \sin(t) - 0.011 \sin(2t)$. The quadratic error is 0.0157, therefore, it is a better approximation. ■

In some occasions, the data are weighted according to their reliability. In this case, one makes use of a weighting matrix $\underline{\underline{W}}$, such that

$$\underline{\underline{W}}\underline{\underline{A}}\underline{\underline{x}} = \underline{\underline{b}}. \quad (11.75)$$

In this case, the optimum solution is given by

$$(\underline{\underline{W}}\underline{\underline{A}})^T \underline{\underline{W}}\underline{\underline{A}}\underline{\underline{x}}^* = (\underline{\underline{W}}\underline{\underline{A}})^T \underline{\underline{b}} \quad (11.76)$$

$$\underline{\underline{x}}^* = [(\underline{\underline{W}}\underline{\underline{A}})^T \underline{\underline{W}}\underline{\underline{A}}]^{-1} (\underline{\underline{W}}\underline{\underline{A}})^T \underline{\underline{b}}. \quad (11.77)$$

This method is referred to as *weighted least square error*.

The application of least square error is not limited to linear systems. In fact, one can find minimize the sum of squares deviations between a set of given values and predicted values. In any case, the least square error problem is given by

$$\min_{\underline{\underline{x}}} f(\underline{\underline{x}}) = \sum_i (g_i - p_i(\underline{\underline{x}}))^2 \quad (11.78)$$

where g_i represent the given data and p_i are the predicted data. To illustrate this concept, let us consider the following example.

Example. An experiment yields the data in the following table.

Find the least square best fit coefficients a , b , and c if the assumed functional form is

(a) $f(\underline{\underline{x}}) = ax_1 + bx_2 + c,$

f	x_1	x_2
2.2	5.0	10.0
9.5	3.0	1.0
23.6	0.6	0.6
74.3	0.1	2.0
6.3	3.0	1.8

(b) $f(\underline{x}) = ax_1^b x_2^c$.

(a) In this case we have a system of linear equations given by

$$\begin{pmatrix} 5.0 & 10.0 & 1.0 \\ 3.0 & 1.0 & 1.0 \\ 0.6 & 0.6 & 1.0 \\ 0.1 & 2.0 & 1.0 \\ 3.0 & 1.8 & 1.0 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 2.2 \\ 9.5 \\ 23.6 \\ 74.3 \\ 6.3 \end{pmatrix}$$

Using (11.72), the solution can be expressed as

$$\begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} -18.17 \\ 4.34 \\ 52.35 \end{pmatrix}.$$

The least square error is 707.2.

(b) In this case, the function is nonlinear and (11.72) cannot be used. In this case, the general approach in (11.78) is utilized. The optimization problem is given by

$$\min_{a,b,c} (2.2 - a5.0^b 10.0^c)^2 + \cdots + (6.3 - a3.0^b 1.8^c)^2$$

The solution of this nonlinear, unconstrained problem is given by

$$\begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 15.72 \\ -0.72 \\ -0.14 \end{pmatrix}.$$

The least square error is 8.08 which indicates that it is a better fit than the previous functional. ■

Chapter 12

Problems

12.1 Exercises

1. Consider the following functions:

- $f(\underline{x}) = (x_1 - 1)^2 + (x_2^2 - 1)^2$
- $f(\underline{x}) = x_1 + \frac{1}{x_1 x_1} + x_2$
- $f(\underline{x}) = \frac{x_1 + x_2}{1 + x_1^2 + x_2^2 + x_1 x_2}$
- $f(\underline{x}) = (x_1 - 1)^2 + x_1 x_3 + x_2 x_3$

Determine (a) all stationary points and (b) check if they are strict local minima using the sufficient conditions.

2. Determine if the given direction at the point is that of descent for the following functions (show all the calculations). In each case, determine the steepest descent direction.

- $f(\underline{x}) = 10(x_2 - x_1^2)^2 + (1 - x_1)^2$; $\underline{d} = (162, -40)^T$ at $\underline{x} = (2, 2)^T$.
- $f(\underline{x}) = (x_1 - 1)^2 + (x_2 - 2)^2 + (x_3 - 3)^2 + (x_4 - 4)^2$; $\underline{d} = (2, -2, 2, -2)^T$ at $\underline{x} = (2, 1, 4, 3)^T$.
- $f(\underline{x}) = \sum_{i=1}^{i=4} [(1 - x_i)^2 + 100(x_{i+1} - x_i^2)^2]$; $\underline{d} = (-1, 1, 2, 1, -1)^T$ at $\underline{x} = (-1, 1, 2, 1, -1)^T$.
- $f(\underline{x}) = x_1 \sin(x_2 x_3^2)$; $\underline{d} = (1, -1, -1)^T$ at $\underline{x} = (1, 2, 1)^T$.
- $f(\underline{x}) = \log_{x_1}(x_2)$; $\underline{d} = (-1, 2)^T$ at $\underline{x} = (10, 10)^T$.

3. Consider the following functions, search directions, and current points:

- $f(\underline{x}) = 4x_1^2 + x_2^2 - 4x_1 x_2$; $\underline{d} = (-1, 2)^T$ at $\underline{x} = (-1, 0)^T$
- $f(\underline{x}) = x_1^2 + x_2^2 + x_3^2$; $\underline{d} = (-2, -4, 2)^T$ at $\underline{x} = (1, 2, -1)^T$
- $f(\underline{x}) = (x_1 - 1)^2 + (x_2 - 2)^2 + (x_3 - 3)^2 + (x_4 - 4)^2$; $\underline{d} = (-2, 2, -2, 2)^T$ at $\underline{x} = (2, 1, 4, 3)$.

In each case, (a) derive the function of one variable (line search function) and (b) determine optimum step size (show all calculations).

4. Consider the function $f(\underline{x}) = x_1 x_2^2$ at the point $\underline{x}_0 = (1, 2)^T$. (a) Obtain the expression quadratic approximation f_Q and make a plot of f and f_Q . (b) Analytically, obtain the minimum of f_Q . This is equivalent to one iterations of Newton's method. (c) Improve this minimum by performing a line search. This is equivalent to one iteration of a modified Newton's method.
5. For the function $f(\underline{x}) = x_1 x_2$, determine the expression for $f(\alpha)$ along the line $x_1 = x_2$ and along the line joining $(0, 1)^T$ to $(1, 0)^T$.
6. Perform two iterations (by hand) using steepest descent, Fletcher-Reeves, Newtons's, and BFGS for the following functions:
 - $f(\underline{x}) = (x_1 + x_2 - 6)^2 + (2x_1 + x_2 - 5)^2$ at $\underline{x}_0 = (0, 0)^T$.
 - $f(\underline{x}) = (x_1 - 1)^2 + 2x_2^2 + 2x_3^2 + 2x_1 x_2 + 2x_2 x_3$ at $\underline{x}_0 = (0, 0, 0)^T$.
7. A cardboard box is to be designed to have a volume of 1 m^3 . Determine the optimal values of length, width, and height to minimize the amount of cardboard material. Hint: The problem can be formulated as unconstrained in terms of two design variables.

12.2 Mini-projects

1. Develop a MATLAB function to solve unconstrained optimization problems using:
 - Steepest descent method.
 - Fletcher-Reeves.
 - Polak-Ribiere.
 - Modified Newton's method (line search).
 - BFGS
 - DFP
 - Trust region method with linear approximation.

Makes sure that your algorithm works for two and three dimensional problems.