# 1 Nonlinear equations

## 1.1 Introduction

*Non-linear equations* (or *root-finding*) is a problem of finding a set of $n$ variables $\mathbf{x} = \{x_1, \ldots, x_n\}$ which satisfy a system of $n$ (generally) non-linear equations

$$f_i(x_1, ..., x_n) = 0 \Big|_{i=1,...,n} . \tag{1}$$

In matrix notation the system is written as

$$\mathbf{f}(\mathbf{x}) = 0 , \tag{2}$$

where $\mathbf{f}(\mathbf{x}) \doteq \{f_1(x_1, \ldots, x_n), \ldots, f_n(x_1, \ldots, x_n)\}$.

In the one-dimensional case it is generally possible to plot the function in the region of interest and see whether the graph crosses the $x$-axis. One can then be sure the root exists and even figure out its approximate position to start one's root-finding algorithm with. In multi-dimensions one generally does not know if the root exists at all, until it is found.

The root-finding algorithms generally proceed by iteration, starting from some approximate solution and making consecutive steps (hopefully) in the direction of the suspected root until some convergence criterion is satisfied. The procedure is generally not even guaranteed to converge unless starting from a point close enough to the sought root.

We shall only consider the multi-dimensional case here since i) the multi-dimensional root-finding is more difficult, and ii) the multi-dimensional routines can also be used in a one-dimensional case.

## 1.2 Newton's method

Newton's method (also reffered to as Newton-Raphson method, after Isaac Newton and Joseph Raphson) is a root-finding algorithm that uses the first term of the Taylor series of the functions $f_i$ to linearise the system (1) in the vicinity of a suspected root. It is one of the oldest and best known methods and is a basis of a number of more refined methods.

Suppose that the point $\mathbf{x} = \{x_1, \ldots, x_n\}$ is close to the root. The Newton's algorithm tries to find the step $\Delta\mathbf{x}$ which would move the point towards the root, such that

$$f_i(\mathbf{x} + \Delta\mathbf{x}) = 0 \Big|_{i=1,...,n} . \tag{3}$$

The first order Taylor expansion of (3) gives a system of linear equations,

$$f_i(\mathbf{x}) + \sum_{k=1}^{n} \frac{\partial f_i}{\partial x_k} \Delta x_k = 0 \Big|_{i=1,\ldots,n} , \tag{4}$$

or, in the matrix form,

$$\mathbf{J}\Delta\mathbf{x} = -\mathbf{f}(\mathbf{x}), \tag{5}$$

where $\mathbf{J}$ is the matrix of partial derivatives,

$$J_{ik} \doteq \frac{\partial f_i}{\partial x_k} , \tag{6}$$

called the *Jacobian matrix*. In practice, if derivatives are not available analytically, one uses finite differences,

$$\frac{\partial f_i}{\partial x_k} \approx \frac{f_i(x_1,\ldots,x_{k-1},x_k+\delta x, x_{k+1},\ldots,x_n) - f_i(x_1,\ldots,x_k,\ldots,x_n)}{\delta x} , \tag{7}$$

with $\delta x \ll s$ with $s$ being the typical scale of the problem at hand.

The solution $\Delta\mathbf{x}$ to the linear system (5) — called the Newton's step — gives the approximate direction and the step-size towards the solution.

The Newton's method converges quadratically if sufficiently close to the solution. Otherwise the full Newton's step $\Delta\mathbf{x}$ might actually diverge from the solution. Therefore in practice a more conservative step $\lambda\Delta\mathbf{x}$ with $\lambda < 1$ is usually taken. The strategy of finding the optimal $\lambda$ is referred to as *line search*.

It is typically not worth the effort to find $\lambda$ which minimizes $\|\mathbf{f}(\mathbf{x} + \lambda\Delta\mathbf{x})\|$ exactly, since $\Delta\mathbf{x}$ is only an approximate direction towards the root. Instead an inexact but quick minimization strategy is usually used, like the *backtracking line search* where one first attempts the full step, $\lambda = 1$, and then backtracks, $\lambda \leftarrow \lambda/2$, until the condition

$$\|\mathbf{f}(\mathbf{x} + \lambda\Delta\mathbf{x})\| < \left(1 - \frac{\lambda}{2}\right)\|\mathbf{f}(\mathbf{x})\| \tag{8}$$

is satisfied. If the condition is not satisfied for sufficiently small $\lambda_{\min}$ the step is taken with $\lambda_{\min}$ simply to step away from this diffictul place and try again.

Following is a typical algrorithm of the Newton's method with backtracking line search and condition (8).

Table 1: C-implementation of Newton's root-finding algorithm with back-tracking.

```
void newton_find_root
(void f(vector x,vector fx),vector x,float dx,float eps) {
 int n=x.size;
 matrix J=matrix_alloc(n,n);
 vector z=vector_alloc(n), Dx=vector_alloc(n);
 vector fx=vector_alloc(n), fz=vector_alloc(n);
 do{
  f(x,fx);
  for(int j=0;j<n;j++){                          /* Jacobian -> J     */
   vector_set(x,j,vector_get(x,j)+dx);
   f(x,fz); vector_add(fz,-1,fx);                /* fz-fx -> fz       */
   for(int i=0;i<n;i++) matrix_set(J,i,j,vector_get(fz,i)/dx);
   vector_set(x,j,vector_get(x,j)-dx);
  }
  givens_qr(J);
  vector_sum(0,fx,-1,fx,Dx);                     /* -f(x) -> Dx       */
  givens_qr_solve_inplace(J,Dx);                 /* -J^{-1} f(x) -> Dx */
  float lambda=2, normfx=vector_norm(fx);
  do{                                            /* back-tracking line-search: */
   lambda/=2; vector_sum(1,x,lambda,Dx,z); f(z,fz);
  }while(vector_norm(fz)>(1-lambda/2)*normfx && lambda>1/128.);
  vector_copy_left_to_right(z,x); vector_copy_left_to_right(fz,fx);
 }while(vector_norm(Dx)>dx && vector_norm(fx)>eps);
 matrix_free(J);
 vector_free(z); vector_free(Dx); vector_free(fx); vector_free(fz);
}
```

```
repeat
    solve JΔx = −f(x) for Δx
    λ ← 1
    while ( ‖f(x + λΔx)‖ > (1 − λ/2) ‖f(x)‖ and λ > 1/64 ) do λ ← λ/2
    x ← x + λΔx
until converged (e.g. ‖f(x)‖ < tolerance)
```

A somewhat more refined backtracking linesearch is based on an approximate minimization of the function

$$g(\lambda) \doteq \frac{1}{2}\|\mathbf{f}(\mathbf{x} + \lambda\Delta\mathbf{x})\|^2 \tag{9}$$

using interpolation. The values $g(0) = \frac{1}{2}\|\mathbf{f}(\mathbf{x})\|^2$ and $g'(0) = -\|\mathbf{f}(\mathbf{x})\|^2$ are already known (check this). If the previous step with certain $\lambda_{\text{trial}}$ was rejected,

we also have $g(\lambda_{\text{trial}})$. These three quantities allow to build a quadratic approximation,

$$g(\lambda) \approx g(0) + g'(0)\lambda + c\lambda^2 , \qquad (10)$$

where

$$c = \frac{g(\lambda_{\text{trial}}) - g(0) - g'(0)\lambda_{\text{trial}}}{\lambda_{\text{trial}}^2} . \qquad (11)$$

The minimum of this approximation (determined by the condition $g'(\lambda) = 0$),

$$\lambda_{\text{next}} = -\frac{g'(0)}{2c} , \qquad (12)$$

becomes the next trial step-size.

The procedure is repeated recursively until either condition (8) is satisfied or the step becomes too small (in which case it is taken unconditionally in order to simply get away from the difficult place).

## 1.3 Quasi-Newton methods

The Newton's method requires calculation of the Jacobian matrix at every iteration. This is generally an expensive operation. Quasi-Newton methods avoid calculation of the Jacobian matrix at the new point $\mathbf{x} + \Delta\mathbf{x}$, instead trying to use certain approximations, typically rank-1 updates.

### 1.3.1 Broyden's method

Broyden's algorithm estimates the Jacobian $\mathbf{J} + \Delta\mathbf{J}$ at the point $\mathbf{x} + \Delta\mathbf{x}$ using the finite-difference approximation,

$$(\mathbf{J} + \Delta\mathbf{J})\Delta\mathbf{x} = \Delta\mathbf{f} , \qquad (13)$$

where $\Delta\mathbf{f} \doteq \mathbf{f}(\mathbf{x} + \Delta\mathbf{x}) - \mathbf{f}(\mathbf{x})$ and $\mathbf{J}$ is the Jacobian at the point $\mathbf{x}$.

The matrix equation (13) is under-determined in more than one dimension as it contains only $n$ equations to determine $n^2$ matrix elements of $\Delta\mathbf{J}$. Broyden suggested to choose $\Delta\mathbf{J}$ as a rank-1 update, linear in $\Delta\mathbf{x}$,

$$\Delta\mathbf{J} = \mathbf{c}\,\Delta\mathbf{x}^T , \qquad (14)$$

where the unknown vector $\mathbf{c}$ can be found by substituting (14) into (13), which gives

$$\Delta\mathbf{J} = \frac{(\Delta\mathbf{f} - \mathbf{J}\Delta\mathbf{x})\Delta\mathbf{x}^T}{\Delta\mathbf{x}^T\Delta\mathbf{x}} . \qquad (15)$$

In practice if one wanders too far from the point where $\mathbf{J}$ was first calculated the accuracy of the updates may decrease significantly. In such case one might need to recalculate $\mathbf{J}$ anew. For example, two successive steps with $\lambda_{\min}$ might be interpreted as a sign of accuracy loss in $\mathbf{J}$ and subsequently trigger its recalculation.

### 1.3.2 Symmetric rank-1 update

The symmetric rank-1 update is chosen in the form

$$\Delta \mathbf{J} = \mathbf{u}\mathbf{u}^T , \tag{16}$$

where the vector $\mathbf{u}$ is found from the condition (13). The update is then given as

$$\Delta \mathbf{J} = \frac{(\Delta \mathbf{f} - \mathbf{J}\Delta \mathbf{x})(\Delta \mathbf{f} - \mathbf{J}\Delta \mathbf{x})^{\mathsf{T}}}{(\Delta \mathbf{f} - \mathbf{J}\Delta \mathbf{x})^{\mathsf{T}}\Delta \mathbf{x}} . \tag{17}$$