

1 Minimization

1.1 Introduction

Minimization (maximization) is the problem of finding the minimum (maximum) of a given—generally non-linear—real valued function $\phi(\mathbf{x})$ of an n -dimensional argument $\mathbf{x} \doteq \{x_1, \dots, x_n\}$. The function is often called the *objective function* or the *cost function*.

Minimization is a simpler case of a more general problem—*optimization*—which includes finding the best available values of the objective function within a given domain and/or subject to given constraints.

Minimization is not unrelated to root-finding: at the minimum all partial derivatives of the objective function vanish,

$$\frac{\partial \phi}{\partial x_i} = 0 \Big|_{i=1, \dots, n}, \quad (1)$$

and one can alternatively solve this system of (non-linear) equations.

1.2 Local minimization

1.2.1 Newton's methods

Newton's method is based on the quadratic approximation of the objective function $\phi(\mathbf{x})$ in the vicinity of the suspected minimum,

$$\phi(\mathbf{x} + \Delta \mathbf{x}) \approx \phi(\mathbf{x}) + \nabla \phi(\mathbf{x})^\top \Delta \mathbf{x} + \frac{1}{2} \Delta \mathbf{x}^\top \mathbf{H}(\mathbf{x}) \Delta \mathbf{x}, \quad (2)$$

where the vector $\nabla \phi(\mathbf{x})$ is the gradient of the objective function at the point \mathbf{x} ,

$$\nabla \phi(\mathbf{x}) \doteq \left\{ \frac{\partial \phi(\mathbf{x})}{\partial x_i} \right\}_{i=1, \dots, n}, \quad (3)$$

and $\mathbf{H}(\mathbf{x})$ is the *Hessian matrix*—a square matrix of second-order partial derivatives of the objective function at the point \mathbf{x} ,

$$\mathbf{H}(\mathbf{x}) \doteq \left\{ \frac{\partial^2 \phi(\mathbf{x})}{\partial x_i \partial x_j} \right\}_{i, j \in 1, \dots, n}. \quad (4)$$

The minimum of the quadratic form (2), as function of $\Delta \mathbf{x}$, is found at the point where its gradient with respect to $\Delta \mathbf{x}$ vanishes,

$$\nabla \phi(\mathbf{x}) + \mathbf{H}(\mathbf{x}) \Delta \mathbf{x} = 0. \quad (5)$$

This gives an approximate step towards the minimum, called the *Newton's step*,

$$\Delta \mathbf{x} = -\mathbf{H}(\mathbf{x})^{-1} \nabla \phi(\mathbf{x}). \quad (6)$$

The original Newton's method is simply the iteration,

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}(\mathbf{x}_k)^{-1} \nabla \phi(\mathbf{x}_k), \quad (7)$$

where at each iteration the full Newton's step is taken and the Hessian matrix is recalculated. In practice, instead of calculating \mathbf{H}^{-1} one rather solves the linear equation (5).

Usually the Newton's method is modified to take a smaller step \mathbf{s} ,

$$\mathbf{s} = \lambda \Delta \mathbf{x}, \quad (8)$$

with $0 < \lambda < 1$. The factor λ can be found by a backtracking algorithm similar to that in the Newton's method for root-finding. One starts with $\lambda = 1$ and then backtracks, $\lambda \leftarrow \lambda/2$, until the *Armijo condition*,

$$\phi(\mathbf{x} + \mathbf{s}) < \phi(\mathbf{x}) + \alpha \mathbf{s}^T \nabla \phi(\mathbf{x}), \quad (9)$$

is satisfied (or the minimal λ is reached, in which case the step is taken unconditionally). The parameter α can be chosen as small as 10^{-4} .

1.2.2 Quasi-Newton methods

Quasi-Newton methods are variations of the Newton's method which attempt to avoid recalculation of the Hessian matrix at each iteration, trying instead certain updates based on the analysis of the gradient vectors. The update $\delta \mathbf{H}$ is usually chosen to satisfy the condition

$$\nabla \phi(\mathbf{x} + \mathbf{s}) = \nabla \phi(\mathbf{x}) + (\mathbf{H} + \delta \mathbf{H}) \mathbf{s}, \quad (10)$$

called *secant equation*, which is the Taylor expansion of the gradient.

The secant equation is under-determined in more than one dimension as it consists of only n equations for the n^2 unknown elements of the update $\delta \mathbf{H}$. Various quasi-Newton methods use different choices for the form of the solution of the secant equation.

In practice one typically uses the inverse Hessian matrix (often—but not always—denoted as \mathbf{B}) and applies the updates directly to the inverse matrix thus avoiding the need to solve the linear equation (5) at each iteration.

For the inverse Hessian matrix the secant equation (10) reads

$$(\mathbf{B} + \delta \mathbf{B}) \mathbf{y} = \mathbf{s}, \quad (11)$$

or, in short,

$$\delta \mathbf{B} \mathbf{y} = \mathbf{u}, \quad (12)$$

where $\mathbf{B} \doteq \mathbf{H}^{-1}$, $\mathbf{y} \doteq \nabla \phi(\mathbf{x} + \mathbf{s}) - \nabla \phi(\mathbf{x})$, and $\mathbf{u} \doteq \mathbf{s} - \mathbf{B} \mathbf{y}$.

One usually starts with the identity matrix as the zeroth approximation for the inverse Hessian matrix and then applies the updates.

Table 1: Quasi-newton minimisation algorithm with updates.

```

set the inverse Hessian matrix to unity, B = 1
repeat until converged (e.g.  $\|\nabla\phi\| < \text{tolerance}$ ) :
    calculate the Newton's step  $\Delta\mathbf{x} = -\mathbf{B}\nabla\phi$ 
    do linesearch starting with  $\lambda = 1$  :
        if  $\phi(\mathbf{x} + \lambda\Delta\mathbf{x}) < \phi(\mathbf{x})$  accept the step:
             $\mathbf{x} = \mathbf{x} + \lambda\Delta\mathbf{x}$ 
            update  $\mathbf{B} = \mathbf{B} + \delta\mathbf{B}$ 
            break linesearch
         $\lambda = \lambda/2$ 
    if  $\lambda$  is too small accept the step and reset B:
         $\mathbf{x} = \mathbf{x} + \lambda\Delta\mathbf{x}$ 
         $\mathbf{B} = 1$ 
        break linesearch
    continue linesearch

```

If the minimal λ is reached during the backtracking line-search—which might be a signal of lost precision in the approximate (inverse) Hessian matrix—it is advisable to reset the current inverse Hessian matrix to identity matrix.

Table 1.2.2 lists one possible algorithm of the quasi-newton method with updates.

Broyden's update The Broyden's update is chosen in the form

$$\delta\mathbf{B} = \mathbf{c}\mathbf{s}^T. \quad (13)$$

where the vector \mathbf{c} is found from the condition (12),

$$\mathbf{c} = \frac{\mathbf{u}}{\mathbf{s}^T\mathbf{y}}. \quad (14)$$

Sometimes the dot-product $\mathbf{s}^T\mathbf{y}$ becomes very small or even zero which results in serious numerical difficulties. One can avoid this by only performing update if the condition $|\mathbf{s}^T\mathbf{y}| > \epsilon$ is satisfied where ϵ is a small number, say 10^{-6} .

Symmetric Broyden's update The Broyden's update (13) is not symmetric (while the Hessian matrix should be) which is an obvious drawback. Therefore a better approximation might be the symmetric Broyden's update,

$$\delta\mathbf{B} = \mathbf{a}\mathbf{s}^T + \mathbf{s}\mathbf{a}^T. \quad (15)$$

The vector \mathbf{a} is again found from the condition (12),

$$\mathbf{a} = \frac{\mathbf{u} - \gamma\mathbf{s}}{\mathbf{s}^T\mathbf{y}}, \quad (16)$$

where $\gamma = (\mathbf{u}^T \mathbf{y}) / (2\mathbf{s}^T \mathbf{y})$.

Again one only performs the update if $|\mathbf{s}^T \mathbf{y}| > \epsilon$.

SR1 update The symmetric-rank-1 update (SR1) is chosen in the form

$$\delta \mathbf{B} = \mathbf{v} \mathbf{v}^T, \quad (17)$$

where the vector \mathbf{v} is again found from the condition (10), which gives

$$\delta \mathbf{B} = \frac{\mathbf{u} \mathbf{u}^T}{\mathbf{u}^T \mathbf{y}}. \quad (18)$$

Again, one only performs the update if denominator is not too small, that is, $|\mathbf{u}^T \mathbf{y}| > \epsilon$.

Other popular updates The wikipedia article “Quasi-Newton method” lists several other popular updates.

1.2.3 Downhill simplex method

The *downhill simplex method* [1] (also called “Nelder-Mead” or “amoeba”) is a commonly used minimization algorithm where the minimum of a function in an n -dimensional space is found by transforming a simplex—a polytope with $n+1$ vertices—according to the function values at the vertices, moving it downhill until it converges towards the minimum.

The advantages of the downhill simplex method is its stability and the lack of use of derivatives. However, the convergence is relatively slow as compared to Newton’s methods.

In order to introduce the algorithm we need the following definitions:

- **Simplex:** a figure (polytope) represented by $n+1$ points, called vertices, $\{\mathbf{p}_1, \dots, \mathbf{p}_{n+1}\}$ (where each point \mathbf{p}_k is an n -dimensional vector).
- **Highest point:** the vertex, \mathbf{p}_{hi} , with the highest value of the function: $\phi(\mathbf{p}_{\text{hi}}) = \max_k \phi(\mathbf{p}_k)$.
- **Lowest point:** the vertex, \mathbf{p}_{lo} , with the lowest value of the function: $\phi(\mathbf{p}_{\text{lo}}) = \min_k \phi(\mathbf{p}_k)$.
- **Centroid:** the center of gravity of all points, except for the highest: $\mathbf{p}_{\text{ce}} = \frac{1}{n} \sum_{(k \neq \text{hi})} \mathbf{p}_k$

The simplex is moved downhill by a combination of the following elementary operations:

1. **Reflection:** the highest point is reflected against the centroid, $\mathbf{p}_{\text{hi}} \rightarrow \mathbf{p}_{\text{re}} = \mathbf{p}_{\text{ce}} + (\mathbf{p}_{\text{ce}} - \mathbf{p}_{\text{hi}})$.

Table 2: Downhill simplex (Nelder-Mead) algorithm

```

REPEAT :
  find highest, lowest, and centroid points of the simplex
  try reflection
  IF  $\phi(\text{reflected}) < \phi(\text{lowest})$  :
    try expansion
    IF  $\phi(\text{expanded}) < \phi(\text{reflected})$  :
      accept expansion
    ELSE :
      accept reflection
  ELSE :
    IF  $\phi(\text{reflected}) < \phi(\text{highest})$  :
      accept reflection
    ELSE :
      try contraction
      IF  $\phi(\text{contracted}) < \phi(\text{highest})$  :
        accept contraction
      ELSE :
        do reduction
UNTIL converged (e.g. size(simplex) < tolerance)

```

2. Expansion: the highest point reflects and then doubles its distance from the centroid, $\mathbf{p}_{\text{hi}} \rightarrow \mathbf{p}_{\text{ex}} = \mathbf{p}_{\text{ce}} + 2(\mathbf{p}_{\text{ce}} - \mathbf{p}_{\text{hi}})$.
3. Contraction: the highest point halves its distance from the centroid, $\mathbf{p}_{\text{hi}} \rightarrow \mathbf{p}_{\text{co}} = \mathbf{p}_{\text{ce}} + \frac{1}{2}(\mathbf{p}_{\text{hi}} - \mathbf{p}_{\text{ce}})$.
4. Reduction: all points, except for the lowest, move towards the lowest points halving the distance. $\mathbf{p}_{k \neq \text{lo}} \rightarrow \frac{1}{2}(\mathbf{p}_k + \mathbf{p}_{\text{lo}})$.

Table 2 shows one possible algorithm for the downhill simplex algorithm.

References

- [1] J.A.Nelder and R.Mead. A simplex method for function minimization. *Computer Journal*, 7:308-313, 1965.