

Function minimization

Volker Blobel – University of Hamburg
March 2005

1. Optimization
2. One-dimensional minimization
3. Search methods
4. Unconstrained minimization
5. Derivative calculation
6. Trust-region methods

Optimization in Physics data analysis

Objective function $F(\mathbf{x})$ is

- (negative log of a) Likelihood function in the Maximum Likelihood method, or
- sum of squares in a (nonlinear) Least Squares problem

Constraints $c_i(\mathbf{x})$ are

- equality constraints, expressing relations between parameters, and
- inequality constraints are limits of certain parameters, defining a restricted range of parameter values (e.g. $m_\nu^2 \geq 0$).

In mathematical terms:

$$\begin{array}{ll} \text{minimize} & F(\mathbf{x}) \\ \text{subject to} & c_i(\mathbf{x}) = 0, \quad i = 1, 2, \dots, m' \\ & c_i(\mathbf{x}) \geq 0, \quad i = m' + 1, \dots, m \end{array} \quad \mathbf{x} \in \mathcal{R}^n$$

Sometimes there are two objective functions: $F(\mathbf{x})$, requiring a good data description, and $G(\mathbf{x})$, requiring e.g. smoothness:

$$\text{minimize sum} \quad F(\mathbf{x}) + \tau \cdot G(\mathbf{x})$$

with *weight* parameter τ .

Aim of optimization

- Find the **global minimum** of the objective function $F(\mathbf{x})$,
- within the allowed range of parameter values,
- in a short time even for a large number of parameters and a complicated objective function,
- even if there are local minima.

Most methods will converge to the *next* minimum, which may be the global minimum or a local minimum (corresponding to *rapid* cooling), going immediately downhill as far as possible.

Search for **the global minimum** requires a special effort.

A special case: combinatorical minimization

The travelling salesman problem:

Find the shortest cyclical itenary for a travelling salesman who must visit each of N cities in turn.

Space over which the function is defined is a discrete, but very large, configuration space (e.g. set of possible orders of cities), with a number of elements in the configuration space which is factorial large (remember: $69! \approx 10^{100}$). The concept of a “direction” may have no meaning.

Objective function is e.g. the total path length.

Desired global extremum may be hidden among many, poorer, local extrema.

Nature is, for slowly cooled systems, able to find the minimum energy state of a crystal that is completely ordered over a large distance.

The method: [Simulated annealing](#), in analogy with the way metals cool and anneal, with **slow** cooling, using Boltzmanns probability distribution

$$\text{Prob}(E) \sim \exp(-E/k_B T)$$

for a system in equilibrium.

Even at low T there is a chance for a higher energy state, with a corresponding chance to get out of a local minimum.

The Metropolis algorithm: the elements

- A description of possible system configurations.
- A generator of random changes in the configuration (called “option”).
- An objective function E (analog of energy) whose minimization is the goal of the procedure.
- A control parameter T (analog of temperature) and an *annealing* schedule (for cooling down).

A random option will change the configuration from energy E_1 to energy E_2 . This option is **accepted** with probability

$$p = \exp \left[\frac{-(E_2 - E_1)}{k_b T} \right]$$

Note: if $E_2 < E_1$ the option is always accepted.

Example: the “effective” solution of the travelling salesman problem

One-dimensional minimization

Search for **minimum of function $f(x)$** of (scalar) argument x .

Important application in multidimensional minimization: robust minimization *along a line*

$$f(z) = F(\mathbf{x} + z \cdot \mathbf{\Delta x}) \quad \text{Line search}$$

Aim: robust, efficient and as fast as possible, because each function evaluation may require a large cpu time.

Standard method: **iterations**, starting from x_0 , with expression $\Phi(x)$:

$$x_{k+1} = \Phi(x_k) \quad \text{for } k = 0, 1, 2, \dots$$

with convergence to fixed point x^* with $x^* = \Phi(x^*)$

Search methods in n dimensions

Search methods in n dimensions **do not require any derivatives**, only **functions values**.

Examples:

Line search in one variable, sequentially in all dimensions

- this is usually rather inefficient

Simplex method by Nelder and Mead: simple, but making use of earlier function evaluations in an efficient way (“learning”)

Monte Carlo search: random search in n dimensions, using result as starting values for more efficient methods

- meaningful if several local minima may exist

In general search methods are acceptable initially (far from the optimum), but are inefficient and slow in the “end game”.

Simplex method in n dimensions

A **simplex** is formed by $n + 1$ points in n -dimensional space ($n = 2$ **triangle**)

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \mathbf{x}_{n+1}$$

sorted such that function values $F_j = F(\mathbf{x}_j)$ are in the order

$$F_1 \leq F_2 \leq \dots F_n \leq F_{n+1}$$

In addition: mean of best n points =

$$\text{center of gravity} \quad \mathbf{c} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$$

Method: sequence of cycles with new point in each cycle, replacing worst point \mathbf{x}_{n+1} , with new (updated) simplex in each cycle.

At the start of each cycle new test point \mathbf{x}_r , obtained by reflexion of worst point at the center of gravity:

$$\mathbf{x}_r = \mathbf{c} + \alpha \cdot (\mathbf{c} - \mathbf{x}_{n+1})$$

A cycle in the Nelder and Mead method

Depending on value $F(x_r)$:

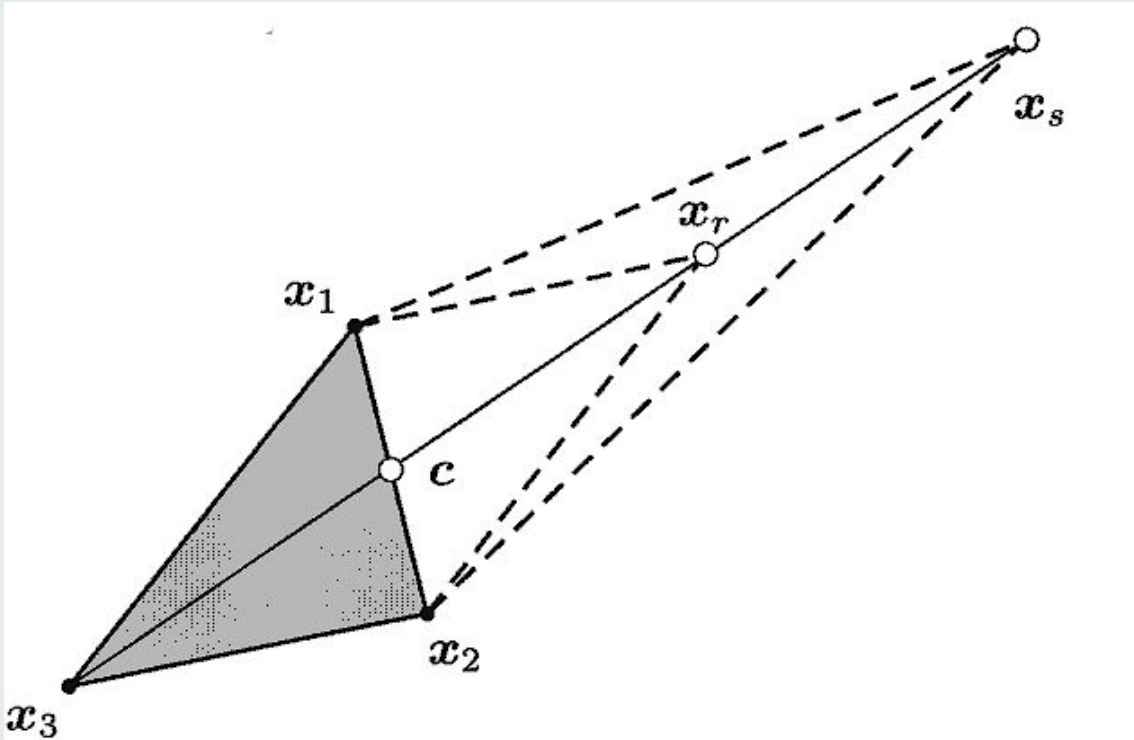
$[F_1 \leq F_r \leq F_n]$: Test point \mathbf{x}_r is middle point and is added, the previous worst point is removed.

$[F_r \leq F_1]$: Test point \mathbf{x}_r is best point, search direction seems to be effective. A new point $\mathbf{x}_s = \mathbf{c} + \beta \cdot (\mathbf{x}_r - \mathbf{c})$ (with $\beta > 1$) is determined and the function value is F_s evaluated. For $F_s < F_r$ the extra step is successful, \mathbf{x}_{n+1} is replaced by \mathbf{x}_s otherwise by \mathbf{x}_r .

$[F_r > F_n]$: the simplex is too big and has to be reduced in size. For $F_r < F_{n+1}$ the test point \mathbf{x}_r replaces the worst point \mathbf{x}_{n+1} . A new test point \mathbf{x}_s is defined by $\mathbf{x}_s = \mathbf{c} - \gamma(\mathbf{c} - \mathbf{x}_{n+1})$ (contraction) with $0 < \gamma < 1$. If this point \mathbf{x}_s with $F_s = F(\mathbf{x}_s) < F_{n+1}$ is an improvement, then \mathbf{x}_{n+1} is replaced by this point \mathbf{x}_s . Otherwise a new simplex is defined by replacing all points except \mathbf{x}_1 by $\mathbf{x}_j = \mathbf{x}_1 + \delta(\mathbf{x}_j - \mathbf{x}_1)$ for $j = 2, \dots, n+1$ with $0 < \delta < 1$, which requires n function evaluations.

Typical values are $\alpha = 1$, $\beta = 2$, $\gamma = 0.5$ und $\delta = 0.5$.

The Simplex method



A few steps of the Simplex method. Starting from the simplex (x_1, x_2, x_3) with center of gravity c . The points x_r and x_s are test points.

Unconstrained minimization

$$\boxed{\text{minimize } F(\mathbf{x}) \quad \mathbf{x} \in \mathcal{R}^n}$$

Taylor expansion:

$$\begin{array}{ll} \text{function} & F(\mathbf{x}_k + \Delta \mathbf{x}) \approx F_k + \mathbf{g}_k^T \Delta \mathbf{x} + \frac{1}{2} \Delta \mathbf{x}^T \mathbf{H}_k \Delta \mathbf{x} \\ \text{derivative} & \mathbf{g}(\mathbf{x}_k + \Delta \mathbf{x}) \approx \mathbf{g}_k + \mathbf{H}_k \Delta \mathbf{x} , \end{array}$$

Function value and derivatives are evaluated at $\mathbf{x} = \mathbf{x}_k$:

$$\begin{array}{ll} \text{function} & F_k = F(\mathbf{x}_k) \\ \text{gradient} & (g_k)_i = \left. \frac{\partial F}{\partial x_i} \right|_{\mathbf{x}=\mathbf{x}_k} \quad i = 1, 2, \dots, n \\ \text{Hesse matrix} & (H_k)_{ij} = \left. \frac{\partial^2 F}{\partial x_i \partial x_j} \right|_{\mathbf{x}=\mathbf{x}_k} \quad i, j = 1, 2, \dots, n \end{array}$$

Covariance matrix for statistical objective functions

Note: if the objective function is

- a sum of squares of deviations, defined according to the Method of Least Squares, or
- a negative log. Likelihood function, defined according to the Maximum Likelihood Method,

then the inverse Hessian \mathbf{H} at the minimum is a good estimate of the covariance matrix of the parameters \mathbf{x} :

$$\mathbf{V}_x \approx \mathbf{H}^{-1}$$

The Newton step

Step $\Delta \mathbf{x}_N$ determined from $\mathbf{g}_k + \mathbf{H}_k \Delta \mathbf{x} = 0$:

$$\Delta \mathbf{x}_N = -\mathbf{H}_k^{-1} \mathbf{g}_k .$$

For a quadratic function $F(\mathbf{x})$ the Newton step is, **in length and direction**, a step to the minimum of the function.

Sometimes large angle ($\approx 90^\circ$) between Newton direction $\Delta \mathbf{x}_N$ and \mathbf{g}_k (the direction of steepest descent).

Calculation of the distance ΔF to the minimum (called EDM in MINUIT):

$$d = -\mathbf{g}_k^T \cdot \Delta \mathbf{x}_N = \mathbf{g}_k^T \mathbf{H}_k^{-1} \mathbf{g}_k = \Delta \mathbf{x}_N^T \mathbf{H}_k \Delta \mathbf{x}_N > 0$$

if Hessian positive-definite. For a quadratic function the distance to the minimum is $d/2$.

General iteration scheme

Define initial value \mathbf{x}_0 , compute $F(\mathbf{x}_0)$ and set $k = 0$

1. **Test for convergence:** If the conditions for convergence are satisfied, the algorithm terminates with \mathbf{x}_k as the solution.
2. **Compute a search vector:** A vector $\Delta\mathbf{x}$ is computed as the new search vector. The Newton search vector is determined from $\mathbf{H}_k \Delta\mathbf{x}_N = -\mathbf{g}_k$.
3. **Line Search:** A one-dimensional minimization is done for the function $f(z) = F(\mathbf{x}_k + z \cdot \Delta\mathbf{x})$. and z_{\min} is determined.
4. **Update:** The point \mathbf{x}_{k+1} is defined by $\mathbf{x}_k + z_{\min} \cdot \Delta\mathbf{x}$, and k is increased by 1. Repeat steps.

Descent direction $\mathbf{g}_k^T \cdot \Delta\mathbf{x} < 0$ if Hessian positive-definite

Method of steepest descent

The search vector $\Delta \mathbf{x}$ is equal to the negative gradienten

$$\Delta \mathbf{x}_{\text{SD}} = -\mathbf{g}_k$$

- Step seems to be natural choice;
- only gradient required (no Hesse matrix) – good;
- no step size defined (in contrast to the Newton step) – bad;
- rate of convergence only linear:

$$\frac{F(\mathbf{x}_{k+1}) - F(\mathbf{x}^*)}{F(\mathbf{x}_k) - F(\mathbf{x}^*)} = c \approx \left(\frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}} \right)^2 = \left(\frac{\kappa - 1}{\kappa + 1} \right)^2 ,$$

(λ_{\max} and λ_{\min} are largest and smallest eigenvalue and κ condition number of Hesse matrix \mathbf{H} .
For a large value of κ c close to one and slow convergence – very bad.

Optimal step including size, if Hesse matrix known:

$$\Delta \mathbf{x}_{\text{SD}} = - \left(\frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{H}_k \mathbf{g}_k} \right) \mathbf{g}_k .$$

Problem: Hessian not positive definite

Diagonalization:

$$\mathbf{H} = \mathbf{U} \mathbf{D} \mathbf{U}^T = \sum_i \lambda_i \mathbf{u}_i \mathbf{u}_i^T,$$

vectors \mathbf{u}_i are the eigenvectors of \mathbf{H} , and \mathbf{D} is diagonal, with the eigenvalues λ_i of \mathbf{H} in diagonal. At least one of the eigenvalues λ_i is zero or less than zero.

1. Strategy: define modified diagonal matrix $\overline{\mathbf{D}}$ with elements

$$\overline{\lambda}_i = \max(|\lambda_i|, \delta) \quad \delta > 0$$

(which requires diagonalization) and use modified **positive-definite** Hessian

$$\overline{\mathbf{H}} = \mathbf{U} \overline{\mathbf{D}} \mathbf{U}^T \quad \overline{\mathbf{H}}^{-1} = \mathbf{U} \overline{\mathbf{D}}^{-1} \mathbf{U}^T$$

2. Strategy: define modified Hessian by

$$\overline{\mathbf{H}} = \mathbf{H} + \alpha \mathbf{I}_n \quad \text{or} \quad \overline{\mathbf{H}}_{jj} = (1 + \lambda) \cdot \mathbf{H}_{jj} \quad (\text{Levenberg-Marquardt})$$

(simpler and faster compared to diagonalization).

- Large value of α (or λ) means search direction close to (negative) gradient (safe),
- small value of α (or λ) means search direction close to the Newton direction.

The factor α allows to interpolate between the methods of **steepest descent** and **Newton method**.

But what is large and what is small α : scale defined by

$$\alpha \approx \frac{|\mathbf{g}^T \mathbf{H} \mathbf{g}|}{\mathbf{g}^T \mathbf{g}}$$

suggested by step of steepest descent method.

Derivative calculation

The (optimal) Newton method requires

- first derivatives of $F(\mathbf{x})$: computation $\propto n$;
- second derivatives of $F(\mathbf{x})$: computation $\propto n^2$.
- Analytical derivatives may be impossible or difficult to obtain.
- Numerical derivatives require good step size δ for differential quotient

E.g. numerical derivative of $f(x)$ in one dimension:

$$\left. \frac{df}{dx} \right|_{x=x_0} \approx \frac{f(x_0 + \delta) - f(x_0 - \delta)}{2\delta}$$

Can the Newton (or quasi Newton) method be used without explicit calculation of the complete Hessian?

Variable metrik method

Calculation of Hessian (with $n(n+1)/2$ different elements) from sequence of first derivatives (gradients):
by update of estimate \mathbf{H}_k from change of gradient.

- Step $\Delta \mathbf{x}$ is calculated from $\mathbf{H}_k \Delta \mathbf{x} = -\mathbf{g}_k$
- After a **line search** with minimum at z_{\min} new value is $\mathbf{x}_{k+1} = \mathbf{x}_k + z_{\min} \Delta \mathbf{x}$,
with gradient \mathbf{g}_{k+1} .
- Update matrix: \mathbf{U}_k is added to \mathbf{H}_k to get new estimate of Hessian:

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \mathbf{U}_k$$

$$\mathbf{H}_{k+1} \boldsymbol{\delta} = \boldsymbol{\gamma} \quad \boldsymbol{\delta} = \mathbf{x}_{k+1} - \mathbf{x}_k \quad \boldsymbol{\gamma} = \mathbf{g}_{k+1} - \mathbf{g}_k$$

(with $\boldsymbol{\gamma}^T \boldsymbol{\delta} > 0$)

Note: an accurate **line search** is essential for the success.

Most effective update formula (Broydon/Fletcher/Goldfarb/Shanno (BFGS))

$$\mathbf{H}_{k+1} = \mathbf{H}_k - \frac{(\mathbf{H}_k \boldsymbol{\delta})(\mathbf{H}_k \boldsymbol{\delta})^T}{\boldsymbol{\delta}^T \mathbf{H}_k \boldsymbol{\delta}} + \frac{\boldsymbol{\gamma} \boldsymbol{\gamma}^T}{\boldsymbol{\gamma}^T \boldsymbol{\delta}} .$$

Initial matrix may be the unit matrix.

Also formulas for the update of the inverse \mathbf{H}_k^{-1} exist.

Properties: the method generates n linear independent search directions for a quadratic function and the estimated Hessian \mathbf{H}_{k+1} converges to the true Hessian \mathbf{H} .

Potential problems:

- no real convergence for "good" starting point;
- estimate "destroyed" for small, inaccurate steps (round-off errors).

Example: MINUIT from CERN

Several options can be selected:

Option MIGRAD: minimizes the objective function, calculates first derivatives numerically and uses the BFGS update formula for the Hessian – fast (if it works)

Option HESSE: calculates the Hesse matrix numerically – recommended after minimization

Option MINIMIZE: minimization by MIGRAD and HESSE calculation with checks.

Restricted range of parameters possible: by transformation from a finite range (of external variables) to $-\infty \dots +\infty$ (of internal variables) using trigonometrical functions.

- allows to use minimization algorithm without range restriction;
- good (quadratic) function is transformed perhaps to a difficult function;
- range specification may be bad, e.g. $[0, 10^8]$ for optimum at 1.

Derivative calculation

Efficient minimization *must* take first and second order derivatives into account, which allows to apply Newton's principle.

Numerical methods for derivative calculation of $F(\mathbf{x})$ by finite differences

- require a good step length – not too small (roundoff errors) and not too large (influence of higher derivatives);
- require very many (time consuming) function evaluations.

Variable metric methods allow to avoid the calculation of second derivatives of $F(\mathbf{x})$, but still need the first derivative.

Are there other possibilities?

- the Gauss-Newton method, which reduces the derivative calculation; the Hessian can be calculated (estimated) without having to compute second derivatives;
- automatic derivative calculation.

Least squares curve fit

Assume a least squares fit of a curve $f(x_i, \mathbf{a})$ with parameters \mathbf{a} to Gaussian distributed data y_i with weight $w_i = 1/\sigma_i^2$:

$$y_i \cong f(x_i, \mathbf{a})$$

Negative Log of the likelihood function is equal to the sum of squares:

$$F(\mathbf{a}) = \sum_i w_i (f(x_i, \mathbf{a}) - y_i)^2$$

Derivatives of $F(\mathbf{a})$ are determined by derivatives of $f(x_i, \mathbf{a})$:

$$\begin{aligned}\frac{\partial F}{\partial a_j} &= \sum_i w_i \frac{\partial f}{\partial a_j} (f(x_i, \mathbf{a}) - y_i) \\ \frac{\partial^2 F}{\partial a_j \partial a_k} &= \sum_i w_i \frac{\partial f}{\partial a_j} \frac{\partial f}{\partial a_k} + \sum_i w_i \frac{\partial^2 f}{\partial a_j \partial a_k} (f(x_i, \mathbf{a}) - y_i)\end{aligned}$$

Usually the contributions from second derivatives of $f(x_i, \mathbf{a})$ are small and can be **neglected**; the approximate Hessian from first derivatives of $f(x_i, \mathbf{a})$ has positive diagonal elements and is expected to be positive-definite. The use of the approximate Hessian is called Gauss-Newton method.

Poisson contribution to objective function

Assume a maximum-likelihood fit of a function $f(x_i, \mathbf{a})$ with parameters \mathbf{a} to Poisson-distributed data y_i :

$$y_i \cong f(x_i, \mathbf{a})$$

Negative Log of the likelihood function:

$$F(\mathbf{a}) = \sum_i f(x_i, \mathbf{a}) - y_i \ln f(x_i, \mathbf{a})$$

or better
$$F(\mathbf{a}) = \sum_i (f(x_i, \mathbf{a}) - y_i) + y_i \ln \frac{y_i}{f(x_i, \mathbf{a})}$$

Derivatives of $F(\mathbf{a})$ are determined by derivatives of $f(x_i, \mathbf{a})$:

$$\frac{\partial F}{\partial a_j} = \sum_i y_i \frac{\frac{\partial f}{\partial a_j}}{f(x_i, \mathbf{a})} - \frac{\partial f}{\partial a_j}$$
$$\frac{\partial^2 F}{\partial a_j \partial a_k} = \sum_i y_i \frac{\frac{\partial f}{\partial a_j} \frac{\partial f}{\partial a_k} - \frac{\partial^2 f}{\partial a_j \partial a_k} f(x_i, \mathbf{a})}{f^2(x_i, \mathbf{a})} - \sum_i \frac{\partial^2 f}{\partial a_j \partial a_k}$$

Usually the contributions from second derivatives of $f(x_i, \mathbf{a})$ are small and can be **neglected**; the approximate Hessian from first derivatives of $f(x_i, \mathbf{a})$ has positive diagonal elements and is expected to be positive-definite. The use of the approximate Hessian is called Gauss-Newton method.

Automatic differentiation

Can automatic differentiation be used?

GNU libmatheval is a library comprising several procedures that makes it possible to create in-memory tree representation of *mathematical functions* over single or multiple variables and later use this representation to *evaluate function for specified variable values*, to create corresponding tree for *function derivative* over specified variable or to get back textual representation of in-memory tree.

Calculated derivative is in mathematical sense correct no matters of fact that derivation variable appears or not in function represented by evaluator.

Supported elementary functions are: exponential (exp), logarithmic (log), square root (sqrt), sine (sin), cosine (cos), tangent (tan), cotangent (cot), secant (sec), cosecant (csc), inverse sine (asin), inverse cosine (acos), inverse tangent (atan), inverse cotangent (acot), inverse secant (asec), inverse cosecant (acsc), hyperbolic sine (sinh), cosine (cosh), hyperbolic tangent (tanh), hyperbolic cotangent (coth), hyperbolic secant (sech), hyperbolic cosecant (csch), hyperbolic inverse sine (asinh), hyperbolic inverse cosine (acosh), hyperbolic inverse tangent (atanh), hyperbolic inverse cotangent (acoth), hyperbolic inverse secant (asech), hyperbolic inverse cosecant (acsch), absolute value (abs), Heaviside step function (step) with value 1 defined for $x = 0$ and Dirac delta function with infinity (delta) and not-a-number (nandelta) values defined for $x = 0$.

Trust-region methods

If the derivatives of $F(\mathbf{x})$ are available the quadratic model

$$\tilde{F}(\mathbf{x}_k + \mathbf{d}) = F_k + \mathbf{g}_k^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \mathbf{H}_k \mathbf{d}$$

can be written at the point \mathbf{x}_k with the gradient $\mathbf{g}_k \equiv \nabla F$ and the Hessian $\mathbf{H}_k \equiv \nabla^2 F$ (matrix of second derivatives). From the approximation the step \mathbf{d} is calculated from

$$\mathbf{g}_k + \mathbf{H}_k \mathbf{d} = 0$$

which is

$$\mathbf{d} = -\mathbf{H}_k^{-1} \mathbf{g}_k .$$

In the **trust-region method** the model is **restricted** depending on a parameter Δ , which forces the next iterate into a sphere centered at \mathbf{x}_k and of radius Δ , which defines the *trust region*, where the model is assumed to be *acceptable*. The vector \mathbf{d} is determined by solving

$$\text{minimize } \tilde{F}(\mathbf{x}_k + \mathbf{d}), \quad |\mathbf{d}| \leq \Delta ,$$

a constrained quadratic minimization problem. The equation can be solved efficiently even if the Hessian is *not positive-definite*.

If λ_1 is the first eigenvalue of a not-positive definite Hessian a parameter λ with $\lambda \geq \max(0, -\lambda_1)$ is selected such that $(\mathbf{H}_k + \lambda \mathbf{I})$ is positive-definite. The minimum point of a Lagrangian

$$L(\mathbf{d}, \lambda) = F(\mathbf{x}_k + \mathbf{d}) + \frac{\lambda}{2} (|\mathbf{d}|^2 - \Delta^2)$$

for fixed λ is given by

$$\mathbf{d}_\Delta = -(\mathbf{H}_k + \lambda \mathbf{I})^{-1} \mathbf{g}_k$$

Now the merit function is $F(\mathbf{x}_k + \mathbf{d}_\Delta)$.

The trajectory $\{\mathbf{x}_k + \mathbf{d}_\Delta\}$ becomes a curve in the space R^n , which has a finite length if \tilde{F} has a minimum point at finite distance.

The actual method to determine \mathbf{d} is discussed later.

The success of a step \mathbf{d}_k is measured by the ratio

$$\rho_k = \frac{F_k - F(\mathbf{x}_k + \mathbf{d}_k)}{F_k - \widetilde{F}(\mathbf{x}_k + \mathbf{d}_k)}$$

of *actual to predicted* decrease of the objective function. Ideally ρ_k is close to 1.

- The trust-region step \mathbf{d}_k is accepted when ρ_k is not much smaller than 1.
- If ρ_k is however much smaller than 1, then the trust-region radius Δ is too large.
- In addition a prescription to increase Δ in case of successful steps is needed for future steps.

A proposed strategy:

if $\rho_k > 0.9$	[very successful]	$x_{k+1} := x_k + d_k$	$\Delta_{k+1} := 2 \times \Delta_k$
otherwise if $\rho_k > 0.1$	[successful]	$x_{k+1} := x_k + d_k$	$\Delta_{k+1} := \Delta_k$
otherwise if $\rho_k \leq 0.1$	[unsuccessful]	$x_{k+1} := x_k$	$\Delta_{k+1} := 1/2 \times \Delta_k$

First case: if \mathbf{H} is positive definite and the solution \mathbf{d} of

$$\mathbf{H}\mathbf{d} = -\mathbf{g}$$

satisfies $||\mathbf{d}_\Delta(\lambda)|| = \Delta$, then the problem is solved.

Second case (otherwise): make the spectral decomposition (diagonalization)

$$\mathbf{H} = \mathbf{U}^T \mathbf{D} \mathbf{U}$$

with \mathbf{U} = matrix of orthogonal eigenvectors, \mathbf{D} = diagonal matrix made up of eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$. One has to calculate a value of λ with $\lambda \geq -\lambda_1$.

$$\mathbf{s}(\lambda) = -(\mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{g}$$

$$||\mathbf{s}(\lambda)||_2^2 = ||\mathbf{U}^T (\mathbf{D} + \lambda \mathbf{I})^{-1} \mathbf{U} \mathbf{g}||_2^2 = \sum_{i=1}^n \left(\frac{\gamma_i}{\lambda_i + \lambda} \right)^2 = \Delta^2$$

where $\gamma_i = \mathbf{e}_i^T \mathbf{U} \mathbf{g}$. This equation has to be solved for λ , which is extremely difficult.

It is far better to solve the equivalent equation

$$\frac{1}{||\mathbf{s}(\lambda)||} - \frac{1}{\Delta} = 0$$

which is an analytical equation, which has no poles and can be solved by Newton's method.

Properties of trust-region methods

Opinions about trust-region methods:

“The trust-region method is actually extremely important and might supersede line-searches, sooner or later.”

“Trust-region methods are not superior to methods with step-length strategy; at least the number of iterations is higher. ... they are not more robust.”

Parameters can be of very different order of magnitude and precision. Preconditioning with appropriate *scaling* of the parameters may be essential for the trust-region method.

Linesearch-methods are naturally *optimistic*, while trust-region methods appear to be naturally *pessimistic*. Trust-region methods are expected to be more robust, with convergence (perhaps slower) even in difficult cases.

Function minimization

Optimization in Physics data analysis	2
Aim of optimization	3
A special case: combinatorical minimization	4
The Metropolis algorithm: the elements	6
One-dimensional minimization	7
Search methods in n dimensions	8
Simplex method in n dimensions	9
A cycle in the Nelder and Mead method	10
The Simplex method	11
Unconstrained minimization	12
Covariance matrix for statistical objective functions	13
The Newton step	14
General iteration scheme	15
Method of steepest descent	16
Problem: Hessian not positive definite	17
Derivative calculation	19
Variable metrik method	20
Example: MINUIT from CERN	22
Derivative calculation	23
Least squares curve fit	24
Poisson contribution to objective function	25
Automatic differentiation	26
Trust-region methods	27
Trust region models	28
Solution of the trust-region problem	30
Properties of trust-region methods	31