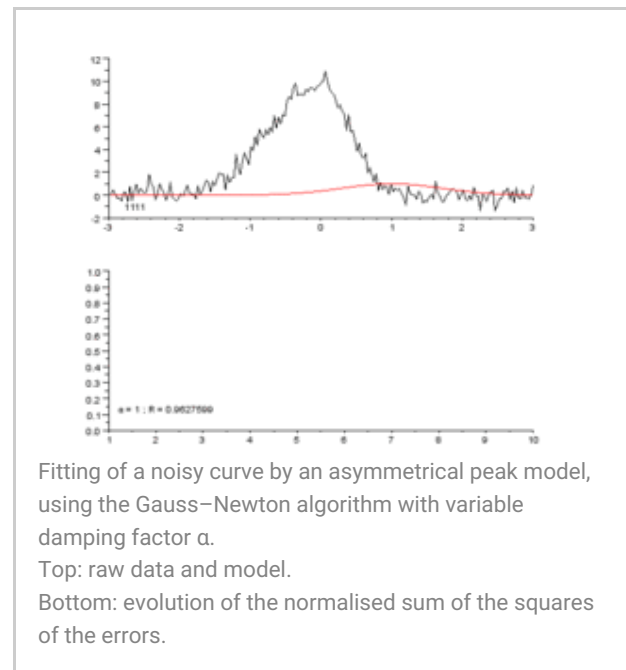


Gauss–Newton algorithm

en.wikipedia.org/wiki/Gauss–Newton_algorithm

The **Gauss–Newton algorithm** is used to solve non-linear least squares problems. It is a modification of Newton's method for finding a minimum of a function. Unlike Newton's method, the Gauss–Newton algorithm can only be used to minimize a sum of squared function values, but it has the advantage that second derivatives, which can be challenging to compute, are not required.

Non-linear least squares problems arise, for instance, in non-linear regression, where parameters in a model are sought such that the model is in good agreement with available observations.



The method is named after the mathematicians Carl Friedrich Gauss and Isaac Newton.

Description

Given m functions $\mathbf{r} = (r_1, \dots, r_m)$ (often called residuals) of n variables $\boldsymbol{\beta} = (\beta_1, \dots, \beta_n)$, with $m \geq n$, the Gauss–Newton algorithm iteratively finds the value of the variables that minimizes the sum of squares^[1]

$$S(\boldsymbol{\beta}) = \sum_{i=1}^m r_i^2(\boldsymbol{\beta}). \quad \{\displaystyle S(\boldsymbol{\beta}) = \sum_{i=1}^m r_i^2(\boldsymbol{\beta})\}$$

Starting with an initial guess $\boldsymbol{\beta}^{(0)}$ for the minimum, the method proceeds by the iterations

$$\boldsymbol{\beta}^{(s+1)} = \boldsymbol{\beta}^{(s)} - (\mathbf{J}_r^T \mathbf{J}_r)^{-1} \mathbf{J}_r^T \mathbf{r}(\boldsymbol{\beta}^{(s)}), \quad \{\displaystyle \boldsymbol{\beta}^{(s+1)} = \boldsymbol{\beta}^{(s)} - \left(\mathbf{J}_r^T \mathbf{J}_r \right)^{-1} \mathbf{J}_r^T \mathbf{r}(\boldsymbol{\beta}^{(s)})\}$$

where, if \mathbf{r} and $\boldsymbol{\beta}$ are column vectors, the entries of the Jacobian matrix are

$$(\mathbf{J}_r)_{ij} = \partial r_i(\boldsymbol{\beta}^{(s)}) / \partial \beta_j, \quad \{\displaystyle (\mathbf{J}_r)_{ij} = \frac{\partial r_i(\boldsymbol{\beta}^{(s)})}{\partial \beta_j}\}$$

and the symbol \mathbf{T} denotes the matrix transpose.

If $m = n$, the iteration simplifies to

$$\boldsymbol{\beta}^{(s+1)} = \boldsymbol{\beta}^{(s)} - (\mathbf{J}_r^T \mathbf{J}_r)^{-1} \mathbf{J}_r^T \mathbf{r}(\boldsymbol{\beta}^{(s)}), \quad \{\displaystyle \boldsymbol{\beta}^{(s+1)} = \boldsymbol{\beta}^{(s)} - \left(\mathbf{J}_r^T \mathbf{J}_r \right)^{-1} \mathbf{J}_r^T \mathbf{r}(\boldsymbol{\beta}^{(s)})\}$$

which is a direct generalization of Newton's method in one dimension.

In data fitting, where the goal is to find the parameters $\boldsymbol{\beta}$ such that a given model function $y = f(x, \boldsymbol{\beta})$ best fits some data points (x_i, y_i) , the functions r_i are the residuals:

$$r_i(\boldsymbol{\beta}) = y_i - f(x_i, \boldsymbol{\beta}). \quad \{\displaystyle r_i(\boldsymbol{\beta}) = y_i - f(x_i, \boldsymbol{\beta})\}$$

Then, the Gauss–Newton method can be expressed in terms of the Jacobian \mathbf{J}_f of the function f as

$$\boldsymbol{\beta}^{(s+1)} = \boldsymbol{\beta}^{(s)} + (\mathbf{J}_f^T \mathbf{J}_f)^{-1} \mathbf{J}_f^T \mathbf{r}(\boldsymbol{\beta}^{(s)}). \quad \{\displaystyle \boldsymbol{\beta}^{(s+1)} = \boldsymbol{\beta}^{(s)} + \left(\mathbf{J}_f^T \mathbf{J}_f \right)^{-1} \mathbf{J}_f^T \mathbf{r}(\boldsymbol{\beta}^{(s)})\}$$

Note that $(\mathbf{J}_f^T \mathbf{J}_f)^{-1} \mathbf{J}_f^T \quad \{\displaystyle \left(\mathbf{J}_f^T \mathbf{J}_f \right)^{-1} \mathbf{J}_f^T\}$

is the pseudoinverse of $\mathbf{J}_f \quad \{\displaystyle \mathbf{J}_f\}$.

Notes

The assumption $m \geq n$ in the algorithm statement is necessary, as otherwise the matrix $\mathbf{J}_r^T \mathbf{J}_r$ is not invertible and the normal equations cannot be solved (at least uniquely).

The Gauss–Newton algorithm can be derived by linearly approximating the vector of functions r_i . Using Taylor's theorem, we can write at every iteration:

$$\mathbf{r}(\boldsymbol{\beta}) \approx \mathbf{r}(\boldsymbol{\beta}^{(s)}) + \mathbf{J}_r(\boldsymbol{\beta}^{(s)}) \Delta \quad \{\displaystyle \mathbf{r}(\boldsymbol{\beta}) \approx \mathbf{r}(\boldsymbol{\beta}^{(s)}) + \mathbf{J}_r(\boldsymbol{\beta}^{(s)}) \Delta\}$$

with $\Delta = \boldsymbol{\beta} - \boldsymbol{\beta}^{(s)} \quad \{\displaystyle \Delta = \boldsymbol{\beta} - \boldsymbol{\beta}^{(s)}\}$

. The task of finding Δ minimizing the sum of squares of the right-hand side, i.e.,

$$\min \|\mathbf{r}(\boldsymbol{\beta}^{(s)}) + \mathbf{J}_r(\boldsymbol{\beta}^{(s)}) \Delta\|_2^2, \quad \{\displaystyle \min \|\mathbf{r}(\boldsymbol{\beta}^{(s)}) + \mathbf{J}_r(\boldsymbol{\beta}^{(s)}) \Delta\|_2^2\}$$

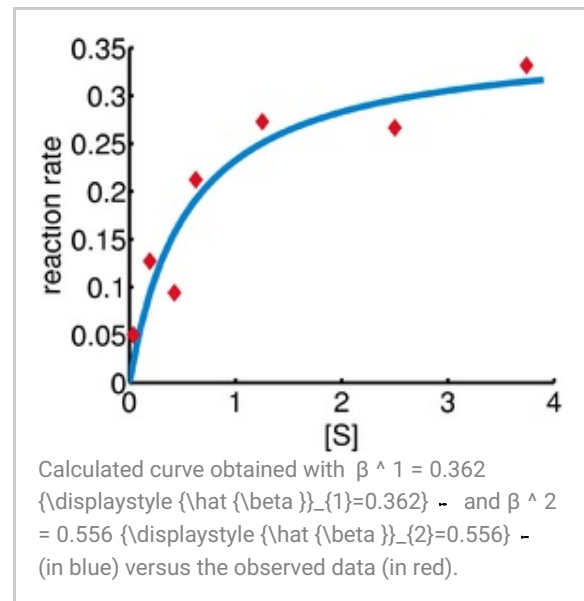
is a linear least-squares problem, which can be solved explicitly, yielding the normal equations in the algorithm.

The normal equations are n simultaneous linear equations in the unknown increments Δ . They may be solved in one step, using Cholesky decomposition, or, better, the QR factorization of \mathbf{J}_r . For large systems, an iterative method, such as the conjugate gradient method, may be more efficient. If there is a linear dependence between columns of \mathbf{J}_r , the iterations will fail, as $\mathbf{J}_r^T \mathbf{J}_r$ becomes singular.

Example

In this example, the Gauss–Newton algorithm will be used to fit a model to some data by minimizing the sum of squares of errors between the data and model's predictions.

In a biology experiment studying the relation between substrate concentration $[S]$ and reaction rate in an enzyme-mediated reaction, the data in the following table were obtained.



i	1	2	3	4	5	6	7
[S]	0.038	0.194	0.425	0.626	1.253	2.500	3.740
rate	0.050	0.127	0.094	0.2122	0.2729	0.2665	0.3317

It is desired to find a curve (model function) of the form

$$\text{rate} = \frac{V_{\text{max}} [S]}{K_M + [S]}$$

that fits best the data in the least-squares sense, with the parameters V_{max} and K_M to be determined.

Denote by x_i and y_i the value of $[S]$ and the rate from the table, $i = 1, \dots, 7$. Let $\beta_1 = V_{\text{max}}$ and $\beta_2 = K_M$. We will find β_1 and β_2 such that the sum of squares of the residuals

$$r_i = y_i - \frac{\beta_1 x_i}{\beta_2 + x_i} \quad (i = 1, \dots, 7)$$

is minimized.

The Jacobian \mathbf{J}_r of the vector of residuals r_i with respect to the unknowns β_j is a 7×2

matrix with the i -th row having the entries

$$\frac{\partial r_i}{\partial \beta_1} = -x_i \beta_2 + x_i; \frac{\partial r_i}{\partial \beta_2} = \beta_1 x_i (\beta_2 + x_i)^2.$$

Starting with the initial estimates of $\beta_1 = 0.9$ and $\beta_2 = 0.2$, after five iterations of the Gauss–Newton algorithm the optimal values $\hat{\beta}_1 = 0.362$ and $\hat{\beta}_2 = 0.556$

are obtained. The sum of squares of residuals decreased from the initial value of 1.445 to 0.00784 after the fifth iteration. The plot in the figure on the right shows the curve determined by the model for the optimal parameters with the observed data.

Convergence properties

It can be shown^[2] that the increment Δ is a descent direction for S , and, if the algorithm converges, then the limit is a stationary point of S . However, convergence is not guaranteed, not even local convergence as in Newton's method, or convergence under the usual Wolfe conditions.^[3]

The rate of convergence of the Gauss–Newton algorithm can approach quadratic.^[4] The algorithm may converge slowly or not at all if the initial guess is far from the minimum or the matrix $\mathbf{J}_r^T \mathbf{J}_r$ is ill-conditioned. For example, consider the problem with $m = 2$ equations and $n = 1$ variable, given by

$$\begin{aligned} r_1(\beta) &= \beta + 1, \\ r_2(\beta) &= \lambda \beta^2 + \beta - 1. \end{aligned}$$

The optimum is at $\beta = 0$. (Actually the optimum is at $\beta = -1$ for $\lambda = 2$, because $S(0) = 1^2 + (-1)^2 = 2$, but $S(-1) = 0$.) If $\lambda = 0$, then the problem is in fact linear and the method finds the optimum in one iteration. If $|\lambda| < 1$, then the method converges linearly and the error decreases asymptotically with a factor $|\lambda|$ at every iteration. However, if $|\lambda| > 1$, then the method does not even converge locally.^[5]

Derivation from Newton's method

In what follows, the Gauss–Newton algorithm will be derived from Newton's method for function optimization via an approximation. As a consequence, the rate of convergence of the Gauss–Newton algorithm can be quadratic under certain regularity conditions. In general (under weaker conditions), the convergence rate is linear.^[6]

The recurrence relation for Newton's method for minimizing a function S of parameters β $\{\displaystyle {\boldsymbol {\beta }}\}$ is

$$\beta(s+1) = \beta(s) - H^{-1}g, \quad \{\displaystyle {\boldsymbol {\beta }}^{(s+1)} = {\boldsymbol {\beta }}^{(s)} - \mathbf{H}^{-1} \mathbf{g} \},$$

where \mathbf{g} denotes the gradient vector of S , and \mathbf{H} denotes the Hessian matrix of S . Since $S = \sum_{i=1}^m r_i^2$ $\{\displaystyle S = \sum_{i=1}^m r_i^2\}$

, the gradient is given by

$$g_j = 2 \sum_{i=1}^m r_i \frac{\partial r_i}{\partial \beta_j}. \quad \{\displaystyle g_j = 2 \sum_{i=1}^m r_i \frac{\partial r_i}{\partial \beta_j} \}.$$

Elements of the Hessian are calculated by differentiating the gradient elements, g_j $\{\displaystyle g_j\}$, with respect to β_k $\{\displaystyle \beta_k\}$:

$$H_{jk} = 2 \sum_{i=1}^m \left(\frac{\partial r_i}{\partial \beta_j} \frac{\partial r_i}{\partial \beta_k} + r_i \frac{\partial^2 r_i}{\partial \beta_j \partial \beta_k} \right). \quad \{\displaystyle H_{jk} = 2 \sum_{i=1}^m \left(\frac{\partial r_i}{\partial \beta_j} \frac{\partial r_i}{\partial \beta_k} + r_i \frac{\partial^2 r_i}{\partial \beta_j \partial \beta_k} \right) \}$$

The Gauss–Newton method is obtained by ignoring the second-order derivative terms (the second term in this expression). That is, the Hessian is approximated by

$$H_{jk} \approx 2 \sum_{i=1}^m J_{ij} J_{ik}, \quad \{\displaystyle H_{jk} \approx 2 \sum_{i=1}^m J_{ij} J_{ik} \}$$

where $J_{ij} = \frac{\partial r_i}{\partial \beta_j}$ $\{\displaystyle J_{ij} = \frac{\partial r_i}{\partial \beta_j} \}$ are entries of the Jacobian \mathbf{J} . The gradient and the approximate Hessian can be written in matrix notation as

$$\mathbf{g} = 2 \mathbf{J}^T \mathbf{r}, \quad \mathbf{H} \approx 2 \mathbf{J}^T \mathbf{J}. \quad \{\displaystyle \mathbf{g} = 2 \mathbf{J}^T \mathbf{r}, \quad \mathbf{H} \approx 2 \mathbf{J}^T \mathbf{J} \}$$

These expressions are substituted into the recurrence relation above to obtain the operational equations

$$\beta(s+1) = \beta(s) + \Delta; \quad \Delta = -(\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \mathbf{r}. \quad \{\displaystyle {\boldsymbol {\beta }}^{(s+1)} = {\boldsymbol {\beta }}^{(s)} + \Delta; \quad \Delta = -(\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \mathbf{r} \}$$

Convergence of the Gauss–Newton method is not guaranteed in all instances. The approximation

$$\left| r_i \frac{\partial^2 r_i}{\partial \beta_j \partial \beta_k} \right| \ll \left| \frac{\partial r_i}{\partial \beta_j} \frac{\partial r_i}{\partial \beta_k} \right| \quad \text{or} \quad \left| \frac{\partial^2 r_i}{\partial \beta_j \partial \beta_k} \right| \ll \left| \frac{\partial r_i}{\partial \beta_j} \frac{\partial r_i}{\partial \beta_k} \right|$$

that needs to hold to be able to ignore the second-order derivative terms may be valid in two cases, for which convergence is to be expected:^[7]

1. The function values r_i are small in magnitude, at least around the minimum.
2. The functions are only "mildly" nonlinear, so that $\frac{\partial^2 r_i}{\partial \beta_j \partial \beta_k}$ is relatively small in magnitude.

Improved versions

With the Gauss–Newton method the sum of squares of the residuals S may not decrease at every iteration. However, since Δ is a descent direction, unless $S(\beta^s)$ is a stationary point, it holds that $S(\beta^s + \alpha \Delta) < S(\beta^s)$ for all sufficiently small $\alpha > 0$. Thus, if divergence occurs, one solution is to employ a fraction α of the increment vector Δ in the updating formula:

$$\beta^{s+1} = \beta^s + \alpha \Delta.$$

In other words, the increment vector is too long, but it points in "downhill", so going just a part of the way will decrease the objective function S . An optimal value for α can be found by using a line search algorithm, that is, the magnitude of α is determined by finding the value that minimizes S , usually using a direct search method in the interval $0 < \alpha < 1$.

In cases where the direction of the shift vector is such that the optimal fraction α is close to zero, an alternative method for handling divergence is the use of the Levenberg–Marquardt algorithm, also known as the "trust region" method.^[1] The normal equations are modified in such a way that the increment vector is rotated towards the direction of steepest descent,

$$(J^T J + \lambda D) \Delta = -J^T r,$$

where D is a positive diagonal matrix. Note that when D is the identity matrix I and $\lambda \rightarrow +\infty$, then $\lambda \Delta = \lambda (J^T J + \lambda I)^{-1} (-J^T r) = (I - J^T J / (\lambda + \dots)) (-J^T r) \rightarrow -J^T r$

$\mathbf{J}^T \mathbf{r}$ approaches the direction of the negative gradient $-\mathbf{J}^T \mathbf{r}$. Therefore the direction of Δ approaches the direction of the negative gradient $-\mathbf{J}^T \mathbf{r}$.

The so-called Marquardt parameter λ may also be optimized by a line search, but this is inefficient, as the shift vector must be recalculated every time λ is changed. A more efficient strategy is this: When divergence occurs, increase the Marquardt parameter until there is a decrease in S . Then retain the value from one iteration to the next, but decrease it if possible until a cut-off value is reached, when the Marquardt parameter can be set to zero; the minimization of S then becomes a standard Gauss–Newton minimization.

Large-scale optimization

For large-scale optimization, the Gauss–Newton method is of special interest because it is often (though certainly not always) true that the matrix $\mathbf{J}^T \mathbf{J}$ is more sparse than the approximate Hessian $\mathbf{J}^T \mathbf{J} \mathbf{J}^T$. In such cases, the step calculation itself will typically need to be done with an approximate iterative method appropriate for large and sparse problems, such as the conjugate gradient method.

In order to make this kind of approach work, one needs at least an efficient method for computing the product

$$\mathbf{J}^T \mathbf{J} \mathbf{p}$$

for some vector \mathbf{p} . With sparse matrix storage, it is in general practical to store the rows of \mathbf{J} in a compressed form (e.g., without zero entries), making a direct computation of the above product tricky due to the transposition. However, if one defines \mathbf{c}_i as row i of the matrix \mathbf{J} , the following simple relation holds:

$$\mathbf{J}^T \mathbf{J} \mathbf{p} = \sum_i \mathbf{c}_i (\mathbf{c}_i \cdot \mathbf{p}),$$

so that every row contributes additively and independently to the product. In addition to respecting a practical sparse storage structure, this expression is well suited for parallel computations. Note that every row \mathbf{c}_i is the gradient of the corresponding residual r_i ; with this in mind, the formula above emphasizes the fact that residuals contribute to the problem independently of each other.

Related algorithms

In a quasi-Newton method, such as that due to Davidon, Fletcher and Powell or Broyden–Fletcher–Goldfarb–Shanno (BFGS method) an estimate of the full Hessian $\frac{\partial^2 S}{\partial \beta_j \partial \beta_k}$ is built up numerically using first derivatives $\frac{\partial r_i}{\partial \beta_j}$ only so that after n refinement cycles the method closely approximates to Newton's

method in performance. Note that quasi-Newton methods can minimize general real-valued functions, whereas Gauss–Newton, Levenberg–Marquardt, etc. fits only to nonlinear least-squares problems.

Another method for solving minimization problems using only first derivatives is gradient descent. However, this method does not take into account the second derivatives even approximately. Consequently, it is highly inefficient for many functions, especially if the parameters have strong interactions.

Notes

1. ^a ^b Björck (1996)
2. ^a Björck (1996), p. 260.
3. ^a Mascarenhas (2013), "The divergence of the BFGS and Gauss Newton Methods", *Mathematical Programming*, **147** (1): 253–276, [arXiv:1309.7922](https://arxiv.org/abs/1309.7922), [doi:10.1007/s10107-013-0720-6](https://doi.org/10.1007/s10107-013-0720-6)
4. ^a Björck (1996), p. 341, 342.
5. ^a Fletcher (1987), p. 113.
6. ^a <http://www.henley.ac.uk/web/FILES/maths/09-04.pdf>
7. ^a Nocedal (1999), p. 259.

References

- Björck, A. (1996). *Numerical methods for least squares problems*. SIAM, Philadelphia. ISBN 0-89871-360-9.
- Fletcher, Roger (1987). *Practical methods of optimization* (2nd ed.). New York: John Wiley & Sons. ISBN 978-0-471-91547-8.
- Nocedal, Jorge; Wright, Stephen (1999). *Numerical optimization*. New York: Springer. ISBN 0-387-98793-2.

External links

Implementations

Artelys Knitro is a non-linear solver with an implementation of the Gauss–Newton method. It is written in C and has interfaces to C++/C#/Java/Python/MATLAB/R.