1 Non-linear Optimization: Newton — Gauss Method

The Newton–Gauss method is a second-order optimization technique for quadratic functions, utilizing a linear approximation of the optimized function at each step. It is applied to solve nonlinear least squares problems, effectively reducing them to a sequence of linear least squares problems.

Gradient and Hessian of the Loss Function:

Given the quadratic loss function

$$Q(\mathbf{x}) = \sum_{\mathbf{x} \in X^{\ell}} (a(\mathbf{x}, \boldsymbol{\theta}) - y(\mathbf{x}))^2$$
(1)

we can express the gradient and Hessian of the function in terms of the model's parameters:

1. The gradient components are

$$egin{aligned} Q'_j &= rac{\partial Q}{\partial heta_j} \ &= 2\sum_{\mathbf{x} \in X^\ell} (a(\mathbf{x}, oldsymbol{ heta}) - y(\mathbf{x})) \cdot rac{\partial a(\mathbf{x}, oldsymbol{ heta})}{\partial heta_j} \end{aligned}$$

2. The Hessian components are

$$\begin{split} Q_{i,j}'' &= \frac{\partial^2 Q}{\partial \theta_i \partial \theta_j} \\ &= 2 \sum_{\mathbf{x} \in X^{\ell}} \frac{\partial a(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_i} \frac{\partial a(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_j} - 2 \sum_{\mathbf{x} \in X^{\ell}} (a(\mathbf{x}, \boldsymbol{\theta}) - y(\mathbf{x})) \cdot \frac{\partial^2 a(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j}. \end{split}$$

Linear Approximation of the Algorithm:

Apply a Taylor series expansion of the algorithm up to the linear term near the current approximation of the parameter vector $\hat{\theta}$:

$$a(\mathbf{x},\boldsymbol{\theta}) = \underbrace{a(\mathbf{x},\hat{\boldsymbol{\theta}})}_{\text{const}} + \sum_{j} \underbrace{\frac{\partial a(\mathbf{x},\hat{\boldsymbol{\theta}})}{\partial \theta_{j}}}_{\text{const}_{j}} \underbrace{(\theta_{j} - \hat{\theta}_{j})}_{\delta \theta_{j}} + O\Big(\left\|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\right\|^{2}\Big), \tag{4}$$

 $a(\mathbf{x}, \hat{\boldsymbol{\theta}})$ is constant, and the linear term is the sum of the partial derivatives of $a(\mathbf{x}, \hat{\boldsymbol{\theta}})$ with respect to the parameters θ_j . The higher-order terms are negligible and will be omitted below.

Differentiate the linear approximation of the algorithm:

$$\frac{\partial}{\partial \theta_j} a(\mathbf{x}, \theta) \approx 0 + \underbrace{\frac{\partial a(\mathbf{x}, \hat{\theta})}{\partial \theta_j}}_{\text{const}_k} \cdot 1 + \underbrace{O(\|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\|^2)}_{\text{const}_j}$$
(5)
= const_j

The components of the sum depending on $\theta_{j\neq k}$ was zeroed out in the differentiation over θ_k . Substitute the obtained derivative into the expression for the Hessian:

$$Q_{i,j}'' \approx 2 \sum_{\mathbf{x} \in X^{\ell}} \underbrace{\frac{\partial a(\mathbf{x}, \hat{\boldsymbol{\theta}})}{\partial \theta_{i}}}_{\text{const}_{i}} \underbrace{\frac{\partial a(\mathbf{x}, \hat{\boldsymbol{\theta}})}{\partial \theta_{j}}}_{\text{const}_{j}} - 2 \sum_{\mathbf{x} \in X^{\ell}} \underbrace{(a(\mathbf{x}, \boldsymbol{\theta}) - \boldsymbol{y}(\mathbf{x})) \cdot \mathbf{0}}_{\boldsymbol{y}(\mathbf{x})}$$
(6)

The linear term will be zeroed out in the second differentiation and will not enter the Hessian.

Matrix Formulation of the Optimization Step:

Introduce the matrix of first partial derivatives and the algorithm's response vector at the current approximation of the parameters $\hat{\theta}$:

gradient is the column vector: $\nabla f(\boldsymbol{x}) \coloneqq \begin{pmatrix} \frac{\partial f(\boldsymbol{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\boldsymbol{x})}{\partial x_1} \end{pmatrix},$

(2)

and f'_j denotes *j*th component of the column

$$D \coloneqq \left\{ \frac{\partial a(\mathbf{x}_i, \hat{\boldsymbol{\theta}})}{\partial \theta_j} \right\}_{i,j}, \quad \boldsymbol{a} \coloneqq \begin{pmatrix} a(\mathbf{x}_1, \hat{\boldsymbol{\theta}}) \\ \vdots \\ a(\mathbf{x}_\ell, \hat{\boldsymbol{\theta}}) \end{pmatrix}$$

matrix D and vector a depend on the point of expansion $\hat{\theta}$ and are recalculated at each optimization step.

The gradient and Hessian (at each step) are calculated using the matrix D:

$$Q' = D^{\mathsf{T}} (\boldsymbol{a} - \boldsymbol{y}), \quad Q'' = D^{\mathsf{T}} D(\boldsymbol{a} - \boldsymbol{y})$$
(8)

The optimization step of the Newton — Rafson method is also expressed in terms of the matrix $D\!\!:$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma \cdot \underbrace{\left(\boldsymbol{D}^{\mathsf{T}} \; \boldsymbol{D} \right)^{-1} \boldsymbol{D}^{\mathsf{T}} \; (\boldsymbol{a} - \boldsymbol{y})}_{\boldsymbol{D}^{+} \qquad \boldsymbol{\varepsilon}} \tag{9}$$

The optimization step vector at each iteration can be determined from the linear system in any of these formulations:

$$\underbrace{\varepsilon}_{y} = D \cdot \underbrace{\delta \theta}_{\beta} \quad \Leftrightarrow \quad \delta \theta = D^{+} \varepsilon \quad \Leftrightarrow \quad \|D \cdot \delta \theta - \varepsilon\|^{2} \to \min_{\beta} \tag{10}$$

Newton — Rafson method is a second-order optimization technique that provides fast convergence. Newton–Gauss method is an approximate second-order method that uses a linear approximation of the optimized function at each step.

(7)

The nonlinear optimization problem is reduced to a sequence of linear problems: at each iteration, a linear expansion of the function is made, matrices are calculated, and a (new) system of linear equations is solved.

The method is a second-order approximation method, providing fast convergence and slightly inferior accuracy compared to the Newton–Raphson method.