Least Square Problems Gauss Newton Method



1 Least Square Problems

2 Gauss Newton Method

Least Square

Least-square Problem

$$egin{array}{rcl} oldsymbol{x}^* &=& rg\min_{oldsymbol{x}\in\mathbb{R}^n}f(oldsymbol{x})\ f(oldsymbol{x}) &=& rac{1}{2}\sum_{j=1}^m r_j(oldsymbol{x})^2 \end{array}$$

where $r_j(\boldsymbol{x}): \mathbb{R}^n \to \mathbb{R}, \ j=1,2,\cdots,m$ are smooth functions.

- $r_j(\boldsymbol{x})$, $j = 1, 2, \cdots, m$ are referred as *residuals*, ie $r_j(\boldsymbol{x}) = y_j \phi(\boldsymbol{x}; t_j)$; $\phi(\boldsymbol{x}; t_j)$ is a model
- It is assumed that $m \ge n$

Least Square Problems Gauss Newton Method

Least-square Problem: Example

Linear Least-square

- $r_j(\beta) = y_j \phi(\beta; t_j); \ j = 1, 2, \cdots, m$
- $\phi(\boldsymbol{\beta};t) = \beta_0 + \beta_1 t; \, \boldsymbol{\beta} = [\beta_0, \beta_1]^T.$
- $f(\boldsymbol{\beta}) = \frac{1}{2} \sum_{j=1}^{m} r_j(\boldsymbol{\beta})^2 = \frac{1}{2} \sum_{j=1}^{m} (y_j \beta_0 \beta_1 t_j)^2$

Least-square Problem: Example

Linear Least-square

•
$$\phi(\boldsymbol{\beta};t) = \beta_0 + \beta_1 t; \, \boldsymbol{\beta} = [\beta_0, \beta_1]^T.$$

•
$$y = 2 * t - 1 + \eta$$
; $\eta \sim \mathcal{N}(0, 2)$, with $t = 0, 1, \cdots, 20$



Non Least-square Problem: Example

Non Linear Least-square

•
$$\phi(\beta; t) = \beta_0 + \beta_1 t + \beta_2 t^2 + \beta_3 e^{-\beta_4 t}$$

•
$$r_j(\boldsymbol{\beta}) = t_j - \phi(\boldsymbol{\beta}; t_j); \ j = 1, 2, \cdots, m; \ \boldsymbol{\beta} = [\beta_0, \beta_1, \cdots, \beta_4]^T.$$



Least Square

Least-square Problem

Defining
$$\boldsymbol{r}(\boldsymbol{x}) = [r_1(\boldsymbol{x}), r_2(\boldsymbol{x}), \cdots, r_m(\boldsymbol{x})]^T$$

$$f(\boldsymbol{x}) = \frac{1}{2} \sum_{j=1}^{m} r_j(\boldsymbol{x})^2 = \frac{1}{2} \|\boldsymbol{r}(\boldsymbol{x})\|_2^2 = \frac{1}{2} \boldsymbol{r}(\boldsymbol{x})^T \boldsymbol{r}(\boldsymbol{x})$$

Least Square: Gradient

Then

$$Df(\boldsymbol{x}) = \frac{1}{2} \left(\boldsymbol{r}(\boldsymbol{x})^T D \boldsymbol{r}(\boldsymbol{x}) + \boldsymbol{r}(\boldsymbol{x})^T D \boldsymbol{r}(\boldsymbol{x}) \right) = \boldsymbol{r}(\boldsymbol{x})^T D \boldsymbol{r}(\boldsymbol{x})$$

$$\nabla f(\boldsymbol{x}) = D \boldsymbol{r}(\boldsymbol{x})^T \boldsymbol{r}(\boldsymbol{x}) = \mathbf{J}(\boldsymbol{x})^T \boldsymbol{r}(\boldsymbol{x})$$

where \mathbf{J} is the Jacobian of $\boldsymbol{r}:\mathbb{R}^n\to\mathbb{R}^m$ and

$$\begin{aligned} \mathbf{J}(\boldsymbol{x}) &= & [J_{ij}]_{\substack{i=1,\dots,m\\ j=1,\dots,n}} \\ &= & [\frac{\partial r_i(\boldsymbol{x})}{\partial x_j}]_{\substack{i=1,\dots,m\\ j=1,\dots,n}} \end{aligned}$$

Least Square: Jacobian

$$\mathbf{J}(\boldsymbol{x}) = [J_{ij}]_{\substack{i=1,\dots,m\\ j=1,\dots,n}}$$

$$= \left[\frac{\partial r_i(\boldsymbol{x})}{\partial x_j}\right]_{\substack{i=1,\dots,m\\ j=1,\dots,n}}$$

$$= \begin{bmatrix} \nabla r_1(\boldsymbol{x})^T\\ \nabla r_2(\boldsymbol{x})^T\\ \vdots\\ \nabla r_m(\boldsymbol{x})^T \end{bmatrix}$$

$$\mathbf{J}(\boldsymbol{x})^T = [\nabla r_1(\boldsymbol{x}), \nabla r_2(\boldsymbol{x}), \cdots, \nabla r_m(\boldsymbol{x})]$$

Least Square Problems Gauss Newton Method

Least Square: Gradient

 ∇

$$f(\boldsymbol{x}) = \mathbf{J}(\boldsymbol{x})^T \boldsymbol{r}(\boldsymbol{x})$$

= $[\nabla r_1(x), \nabla r_2(x), \cdots, \nabla r_m(x)] \begin{bmatrix} r_1(x) \\ r_2(x) \\ \vdots \\ r_m(x) \end{bmatrix}$
= $\sum_{i=1}^m r_i(x) \nabla r_i(x)$

Least Square: Hessian

Gradient

$$\nabla f(\boldsymbol{x}) = \sum_{i=1}^{m} r_i(x) \nabla r_i(x)$$

Hessian

$$\nabla^2 f(\boldsymbol{x}) = \sum_{i=1}^m \nabla r_i(x) \nabla r_i(x)^T + r_i(x) \nabla^2 r_i(x)$$
$$= \mathbf{J}(\boldsymbol{x})^T \mathbf{J}(\boldsymbol{x}) + \sum_{i=1}^m r_i(x) \nabla^2 r_i(x)$$
$$= \mathbf{J}(\boldsymbol{x})^T \mathbf{J}(\boldsymbol{x}) + S(\boldsymbol{x})$$

Linear Least Square

Least-square Problem

$$\begin{aligned} \mathbf{r}(\mathbf{x}) &= \mathbf{J}\mathbf{x} - \mathbf{b} \\ f(\mathbf{x}) &= \frac{1}{2} \|\mathbf{J}\mathbf{x} - \mathbf{b}\|_2^2 \\ \mathbf{J}(\mathbf{x}) &= D\mathbf{r}(\mathbf{x}) = \mathbf{J} \\ \nabla f(\mathbf{x}) &= \mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x}) = \mathbf{J}^T (\mathbf{J}\mathbf{x} - \mathbf{b}) = \mathbf{J}^T \mathbf{J}\mathbf{x} - \mathbf{J}^T \mathbf{b} \\ \nabla^2 f(\mathbf{x}) &= \mathbf{J}^T \mathbf{J} = \mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}) + \mathbf{0} \end{aligned}$$

Note:

- $\mathbf{J}(\boldsymbol{x}) = \mathbf{J}$ is a constant matrix
- $\sum_{k=1}^m r_k(x) \nabla^2 r_k(x) = 0$ due to $\nabla^2 r_k(x) = 0$, ie, $r_k(x)$ is affine.

Linear Least Square

Least-square Problem

As

$$\nabla f(\boldsymbol{x}) = \mathbf{J}^T \mathbf{J} \boldsymbol{x} - \mathbf{J}^T \boldsymbol{b}$$

the optimum x^{st} satisfies

$$\mathbf{J}^T \mathbf{J} \mathbf{x} = \mathbf{J}^T \mathbf{b}$$

known as normal equations.

The Gauss Newton Method is used to solve the problem

$$egin{array}{rcl} oldsymbol{x}^* &=& rg\min_{oldsymbol{x}\in\mathbb{R}^n}f(oldsymbol{x})\ f(oldsymbol{x}) &=& rac{1}{2}\sum_{j=1}^m r_j(oldsymbol{x})^2 \end{array}$$

It exploits the structure of the Hessian $abla^2 f({m x})$

Instead of the standard direction

$$abla^2 f(oldsymbol{x}_k) oldsymbol{d}_k^N = -
abla f(oldsymbol{x}_k)$$

one solves the following system of equation with respect to $oldsymbol{d}_k^{GN}$

$$\mathbf{J}(oldsymbol{x}_k)^T \mathbf{J}(oldsymbol{x}_k) oldsymbol{d}_k^{GN} = -\mathbf{J}(oldsymbol{x}_k)^T oldsymbol{r}(oldsymbol{x}_k)$$

(1) If $r_k({m x})pprox 0$ or $abla^2 r_k({m x})pprox 0$, orall k then

$$abla^2 f(\boldsymbol{x}) ~\approx~ \mathbf{J}(\boldsymbol{x})^T \mathbf{J}(\boldsymbol{x})$$

then, we do not require to compute the individual residual Hessians $\nabla^2 r_k(\boldsymbol{x}).$

2 There are many situation where $\mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x})$ dominates the second term. Therefore, $\mathbf{J}(\mathbf{x}_k)^T \mathbf{J}(\mathbf{x}_k)$ is a close approximation to $\nabla^2 f(\mathbf{x}_k)$ and the convergence rate of Gauss-Newton is similar to that of Newton's method.

• If \mathbf{J}_k has full rank and the gradient ∇f_k is nonzero, the direction d^{GN} is a descent direction, and therefore a suitable direction for a line search.

$$d^{GNT} \nabla f(\boldsymbol{x}) = d^{GNT} \mathbf{J}(\boldsymbol{x}_k)^T \boldsymbol{r}(\boldsymbol{x}_k)$$

= $-d^{GNT} \mathbf{J}(\boldsymbol{x}_k)^T \mathbf{J}(\boldsymbol{x}_k) d_k^{GN}$
= $-\|\mathbf{J}(\boldsymbol{x}_k) d_k^{GN}\|_2^2 \le 0$

What happens when $\mathbf{J}(\boldsymbol{x}_k)\boldsymbol{d}_k^{GN}=0?$

1 The final inequality is strict unless $\mathbf{J}(\boldsymbol{x}_k)\boldsymbol{d}_k^{GN} = \mathbf{0}$, in which case we have by the full rank of \mathbf{J}_k

$$\begin{aligned} \mathbf{J}(\boldsymbol{x}_k)^T \mathbf{J}(\boldsymbol{x}_k) \boldsymbol{d}_k^{GN} &= -\mathbf{J}(\boldsymbol{x}_k)^T \boldsymbol{r}(\boldsymbol{x}_k) \\ \mathbf{J}(\boldsymbol{x}_k)^T \mathbf{0} &= -\nabla f(\boldsymbol{x}_k) \\ \nabla f(\boldsymbol{x}_k) &= \mathbf{0} \end{aligned}$$

then x_k is a stationary point.

• The Gauss-Newton arises from the similarity between the equations

$$\mathbf{J}(oldsymbol{x}_k)^T \mathbf{J}(oldsymbol{x}_k) oldsymbol{d}_k^{GN} ~=~ - \mathbf{J}(oldsymbol{x}_k)^T oldsymbol{r}(oldsymbol{x}_k)$$

and the normal equations for the linear least-squares problem.
 The previous connection tells us that d_k^{GN} is in fact the solution of the linear least-squares problem

$$rgmin_{oldsymbol{d}} \|\mathbf{J}(oldsymbol{x}_k)oldsymbol{d} + oldsymbol{r}(oldsymbol{x}_k)\|^2$$

 If the QR (with column pivoting) or SVD-based algorithms are used to solve the corresponding linear system

$$\mathbf{J}(\boldsymbol{x}_k)^T \mathbf{J}(\boldsymbol{x}_k) \boldsymbol{d}_k^{GN} = -\mathbf{J}(\boldsymbol{x}_k)^T \boldsymbol{r}(\boldsymbol{x}_k)$$

there is no need to calculate the Hessian approximation $\mathbf{J}(\boldsymbol{x}_k)^T \mathbf{J}(\boldsymbol{x}_k)$ explicitly; we can work directly with the Jacobian $\mathbf{J}(\boldsymbol{x}_k)$.

1 The linear least-squares problem

$$rgmin_{oldsymbol{d}} \| \mathbf{J}(oldsymbol{x}_k) oldsymbol{d} + oldsymbol{r}(oldsymbol{x}_k) \|^2$$

can be viewed as the linear model for the the vector function $r(x_k + d) pprox r(x_k) + \mathbf{J}(x_k) d$ therefore

$$f(x_k + d) = \frac{1}{2} \|r(x_k + d)\|^2 \approx \frac{1}{2} \|\mathbf{J}(x_k)d + r(x_k)\|^2$$

2 Implementations of the Gauss-Newton method usually perform a line search in the direction d^{GN} .

Theorem 2.1

Suppose each residual function r_j is Lipschitz continuously differentiable in a neighborhood N of the bounded level set

 $\mathcal{L} = \{ \boldsymbol{x} | f(\boldsymbol{x}) \leq f(\boldsymbol{x}_0) \}$

where x_0 is the starting point for the algorithm, and that the Jacobians $\mathbf{J}(x)$ satisfy (the uniform full-rank condition) that there is a constant $\gamma > 0$ such that

 $\|\mathbf{J}(\boldsymbol{x})\boldsymbol{z}\| \geq \gamma \|\boldsymbol{z}\|$

for all x in a neighborhood \mathcal{N} of the level set \mathcal{L} . Then if the iterates x_k are generated by the Gauss-Newton method with step lengths α_k that satisfy the Wolfe conditions, we have

$$\lim_{k\to\infty}\mathbf{J}_k^T\boldsymbol{r}_k = 0.$$

Theorem 2.2

Let $\mathbf{r} : \mathbb{R}^n \to \mathbb{R}^m$, and let $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|^2$ be twice continuously differentiate in an open convex set Ω . Assume that $\mathbf{J}(\mathbf{x}) \in Lip_{\gamma}(\Omega)$ with $\|\mathbf{J}(\mathbf{x})\| \ge \alpha$ for all $\mathbf{x} \in \Omega$ and there exists $\mathbf{x}^* \in \Omega$ and $\lambda, \sigma \ge 0$ such that $\mathbf{J}(\mathbf{x}^*)^T \mathbf{r}(\mathbf{x}^*) = 0$, λ is the smallest eigenvalue of $\mathbf{J}(\mathbf{x}^*)^T \mathbf{J}(\mathbf{x}^*)$, and

$$\| (\mathbf{J}(x) - \mathbf{J}(x^*))^T r(x^*) \| \le \sigma \|x - x^*\|$$

for all $x \in \Omega$. If $\sigma < \lambda$ for any $c \in (1, \lambda/\sigma)$ there exists $\epsilon > 0$ such that for all $x_0 \in \mathcal{N}(x^*, \epsilon)$ the sequence generated by the Gauss-Newton method is well defined, converges to x^* , and obeys

$$\|oldsymbol{x}_{k+1}-oldsymbol{x}^*\|\leq rac{c\sigma}{\lambda}\|oldsymbol{x}_k-oldsymbol{x}^*\|+rac{clpha\gamma}{2\lambda}\|oldsymbol{x}_k-oldsymbol{x}^*\|^2$$

and $\|oldsymbol{x}_{k+1} - oldsymbol{x}^*\| \leq rac{c\sigma+\lambda}{2\lambda} \|oldsymbol{x}_k - oldsymbol{x}^*\| < \|oldsymbol{x}_k - oldsymbol{x}^*\|$

Corollary

Let the assumptions of the previous theorem be satisfied. If $r(x^*) = 0$, then there exists $\epsilon > 0$ such that for all $x_0 \in \mathcal{N}(x^*, \epsilon)$, the sequence $\{x_k\}$ generated by the Gauss-Newton method is well defined and converges quadratically to x^* .

Gauss Newton Method:Advantages

- 1 Locally quadratically convergent on zero-residual problems.
- Quickly locally q-linearly convergent on problems that aren't too nonlinear and have reasonably small residuals.
- **3** Solves linear least-squares problems in one iteration.

Gauss Newton Method: Disadvantages

- Slowly locally linearly convergent on problems that are sufficiently nonlinear or have reasonably large residuals.
- 2 Not locally convergent on problems that are very nonlinear or have very large residuals.
- **3** Not well defined if $\mathbf{J}(\boldsymbol{x}_k)$ doesn't have full column rank.
- 4 Not necessarily globally convergent.

Let us use QR Factorization with Column Pivoting, ie

$$\mathbf{JP} = \mathbf{QR} = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{Q}_2 \end{bmatrix} \begin{bmatrix} \mathbf{R}_1 \\ 0 \end{bmatrix}$$

where $\mathbf{J} \in \mathbb{R}^{m \times n}$, $\mathbf{P} \in \mathbb{R}^{n \times n}$ is a permutation matrix, $\mathbf{Q} \in \mathbb{R}^{m \times m}$ and $\mathbf{R} \in \mathbb{R}^{m \times n}$, $\mathbf{Q}_1 \in \mathbb{R}^{m \times n}$ $\mathbf{Q}_2 \in \mathbb{R}^{m \times m-n}$ with

$$\mathbf{Q}\mathbf{Q}^T = \mathbf{Q}^T\mathbf{Q} = \mathbf{I}$$

and $\mathbf{R}_1 \in \mathbb{R}^{n imes n}$ is an upper triangular matrix with elements of the diagonal satisfying

$$|r_{11}| \ge |r_{22}| \ge \dots \ge |r_{nn}|$$

Considering $\|\mathbf{Q} \boldsymbol{x}\|_2 = \|\boldsymbol{x}\|_2$ if \mathbf{Q} is orthogonal, then

$$\begin{split} \|\mathbf{J}_k \boldsymbol{d}_k + \boldsymbol{r}_k\|^2 &= \|\mathbf{J}_k \mathbf{P} \mathbf{P}^T \boldsymbol{d}_k + \boldsymbol{r}_k\|^2 \\ &= \|\mathbf{Q} \mathbf{R} \mathbf{P}^T \boldsymbol{d}_k + \boldsymbol{r}_k\|^2 \\ &= \|\mathbf{R} \mathbf{P}^T \boldsymbol{d}_k + \mathbf{Q}^T \boldsymbol{r}_k\|^2 \\ &= \left\| \begin{bmatrix} \mathbf{R}_1 \\ 0 \end{bmatrix} \mathbf{P}^T \boldsymbol{d}_k + \begin{bmatrix} \mathbf{Q}_1^T \\ \mathbf{Q}_2^T \end{bmatrix} \boldsymbol{r}_k \right\|^2 \\ &= \|\mathbf{R}_1 \mathbf{P}^T \boldsymbol{d}_k + \mathbf{Q}_1^T \boldsymbol{r}_k\|^2 + \|\mathbf{Q}_2^T \boldsymbol{r}_k\|^2 \end{split}$$

From the last equation, ie,

$$\|\mathbf{J}_k \boldsymbol{d}_k + \boldsymbol{r}_k\|^2 = \|\mathbf{R}_1 \mathbf{P}^T \boldsymbol{d}_k + \mathbf{Q}_1^T \boldsymbol{r}_k\|^2 + \|\mathbf{Q}_2^T \boldsymbol{r}_k\|^2$$

Computing the gradient w.r.t d_k and due to \mathbf{P}, \mathbf{R}_1 have inverse,.. then

$$\begin{aligned} \mathbf{P} \mathbf{R}_1^T (\mathbf{R}_1 \mathbf{P}^T \boldsymbol{d}_k + \mathbf{Q}_1^T \boldsymbol{r}_k) &= 0 \\ \mathbf{R}_1 \mathbf{P}^T \boldsymbol{d}_k &= -\mathbf{Q}_1^T \boldsymbol{r}_k \end{aligned}$$

the previous system can be solved in two steps

Defining

$$egin{array}{rl} m{b} &=& - \mathbf{Q}_1^T m{r}_k \ m{z} &=& \mathbf{P}^T m{d}_k \end{array}$$

ie, $d_k = \mathbf{P} \boldsymbol{z}$. From

$$\mathbf{R}_1 \mathbf{P}^T \boldsymbol{d}_k = -\mathbf{Q}_1^T \boldsymbol{r}_k$$

we solve the following systems, first for z and then for d_k

$$egin{array}{rcl} {f R}_1 oldsymbol{z} &=& oldsymbol{b} \ {oldsymbol{d}}_k &=& {f P} oldsymbol{z} \end{array}$$