

Least-Squares Problems

Oscar Dalmau
dalmau@cimat.mx

Centro de Investigación en Matemáticas
CIMAT A.C. Mexico

March 2016

Outline

① Least Square Problems

② Gauss Newton Method

Least Square

Least-square Problem

$$\begin{aligned}\mathbf{x}^* &= \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \\ f(\mathbf{x}) &= \frac{1}{2} \sum_{j=1}^m r_j(\mathbf{x})^2\end{aligned}$$

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

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

Least-square Problem: Example

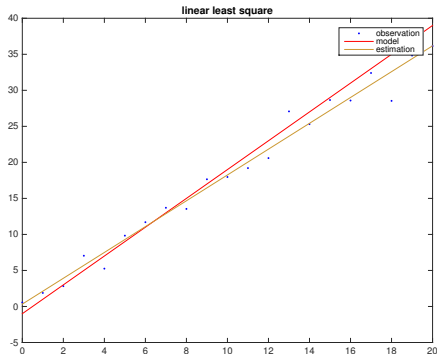
Linear Least-square

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

Least-square Problem: Example

Linear Least-square

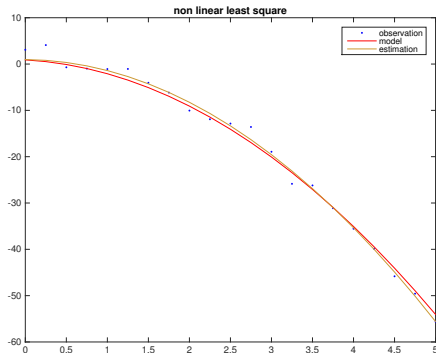
- $\phi(\boldsymbol{\beta}; t) = \beta_0 + \beta_1 t; \boldsymbol{\beta} = [\beta_0, \beta_1]^T$.
- $y = 2 * t - 1 + \eta; \eta \sim \mathcal{N}(0, 2)$, with $t = 0, 1, \dots, 20$



Non Least-square Problem: Example

Non Linear Least-square

- $\phi(\boldsymbol{\beta}; t) = \beta_0 + \beta_1 t + \beta_2 t^2 + \beta_3 e^{-\beta_4 t}$
- $r_j(\boldsymbol{\beta}) = t_j - \phi(\boldsymbol{\beta}; t_j); j = 1, 2, \dots, m; \boldsymbol{\beta} = [\beta_0, \beta_1, \dots, \beta_4]^T$.



Least Square

Least-square Problem

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

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

Least Square: Gradient

Then

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

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

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

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

Least Square: Jacobian

$$\begin{aligned}\mathbf{J}(\mathbf{x}) &= [J_{ij}]_{\substack{i=1,\dots,m \\ j=1,\dots,n}} \\ &= \left[\frac{\partial r_i(\mathbf{x})}{\partial x_j} \right]_{\substack{i=1,\dots,m \\ j=1,\dots,n}} \\ &= \begin{bmatrix} \nabla r_1(x)^T \\ \nabla r_2(x)^T \\ \vdots \\ \nabla r_m(x)^T \end{bmatrix} \\ \mathbf{J}(\mathbf{x})^T &= [\nabla r_1(x), \nabla r_2(x), \dots, \nabla r_m(x)]\end{aligned}$$

Least Square: Gradient

$$\begin{aligned}\nabla f(\mathbf{x}) &= \mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x}) \\ &= [\nabla r_1(x), \nabla r_2(x), \dots, \nabla r_m(x)] \begin{bmatrix} r_1(x) \\ r_2(x) \\ \vdots \\ r_m(x) \end{bmatrix} \\ &= \sum_{i=1}^m r_i(x) \nabla r_i(x)\end{aligned}$$

Least Square: Hessian

Gradient

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

Hessian

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

Linear Least Square

Least-square Problem

$$\mathbf{r}(\mathbf{x}) = \mathbf{J}\mathbf{x} - \mathbf{b}$$

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{J}\mathbf{x} - \mathbf{b}\|_2^2$$

$$\mathbf{J}(\mathbf{x}) = D\mathbf{r}(\mathbf{x}) = \mathbf{J}$$

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

$$\nabla^2 f(\mathbf{x}) = \mathbf{J}^T \mathbf{J} = \mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}) + \mathbf{0}$$

Note:

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

Linear Least Square

Least-square Problem

As

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

the optimum \mathbf{x}^* 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

$$\begin{aligned}\boldsymbol{x}^* &= \arg \min_{\boldsymbol{x} \in \mathbb{R}^n} f(\boldsymbol{x}) \\ f(\boldsymbol{x}) &= \frac{1}{2} \sum_{j=1}^m r_j(\boldsymbol{x})^2\end{aligned}$$

It exploits the structure of the Hessian $\nabla^2 f(\boldsymbol{x})$

Instead of the standard direction

$$\nabla^2 f(\mathbf{x}_k) \mathbf{d}_k^N = -\nabla f(\mathbf{x}_k)$$

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

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

Gauss Newton Method

- ① If $r_k(\mathbf{x}) \approx 0$ or $\nabla^2 r_k(\mathbf{x}) \approx 0, \forall k$ then

$$\nabla^2 f(\mathbf{x}) \approx \mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x})$$

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

- 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.

Gauss Newton Method

- ① If \mathbf{J}_k has full rank and the gradient ∇f_k is nonzero, the direction \mathbf{d}^{GN} is a descent direction, and therefore a suitable direction for a line search.

$$\begin{aligned}\mathbf{d}^{GN T} \nabla f(\mathbf{x}) &= \mathbf{d}^{GN T} \mathbf{J}(\mathbf{x}_k)^T \mathbf{r}(\mathbf{x}_k) \\ &= -\mathbf{d}^{GN T} \mathbf{J}(\mathbf{x}_k)^T \mathbf{J}(\mathbf{x}_k) \mathbf{d}_k^{GN} \\ &= -\|\mathbf{J}(\mathbf{x}_k) \mathbf{d}_k^{GN}\|_2^2 \leq 0\end{aligned}$$

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

Gauss Newton Method

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

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

then \mathbf{x}_k is a stationary point.

Gauss Newton Method

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

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

and the *normal equations* for the linear least-squares problem.

- 2 The previous connection tells us that \mathbf{d}_k^{GN} is in fact the solution of the linear least-squares problem

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

Gauss Newton Method

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

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

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

Gauss Newton Method

- 1 The linear least-squares problem

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

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

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

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

Theorem 2.1

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

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

where \mathbf{x}_0 is the starting point for the algorithm, and that the Jacobians $\mathbf{J}(\mathbf{x})$ satisfy (the uniform full-rank condition) that there is a constant $\gamma > 0$ such that

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

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

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

Theorem 2.2

Let $\mathbf{r} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, 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 $\mathbf{J}(\mathbf{x}) \in \text{Lip}_\gamma(\Omega)$ with $\|\mathbf{J}(\mathbf{x})\| \geq \alpha$ for all $\mathbf{x} \in \Omega$ and there exists $\mathbf{x}^* \in \Omega$ and $\lambda, \sigma \geq 0$ 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

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

for all $\mathbf{x} \in \Omega$. If $\sigma < \lambda$ for any $c \in (1, \lambda/\sigma)$ there exists $\epsilon > 0$ such that for all $\mathbf{x}_0 \in \mathcal{N}(\mathbf{x}^*, \epsilon)$ the sequence generated by the Gauss-Newton method is well defined, converges to \mathbf{x}^* , and obeys

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \frac{c\sigma}{\lambda} \|\mathbf{x}_k - \mathbf{x}^*\| + \frac{c\alpha\gamma}{2\lambda} \|\mathbf{x}_k - \mathbf{x}^*\|^2$$

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

Corollary

Let the assumptions of the previous theorem be satisfied. If $\mathbf{r}(\mathbf{x}^*) = 0$, then there exists $\epsilon > 0$ such that for all $\mathbf{x}_0 \in \mathcal{N}(\mathbf{x}^*, \epsilon)$, the sequence $\{\mathbf{x}_k\}$ generated by the Gauss-Newton method is well defined and converges quadratically to \mathbf{x}^* .

Gauss Newton Method: Advantages

- ① Locally quadratically convergent on zero-residual problems.
- ② 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

- 1 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 $\mathbf{J}(\mathbf{x}_k)$ doesn't have full column rank.
- 4 Not necessarily globally convergent.

Implementation using QR

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

$$|r_{11}| \geq |r_{22}| \geq \cdots \geq |r_{nn}|$$

Implementation using QR

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

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

Implementation using QR

From the last equation, ie,

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

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

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

the previous system can be solved in two steps

Implementation using QR

Defining

$$\begin{aligned} \mathbf{b} &= -\mathbf{Q}_1^T \mathbf{r}_k \\ \mathbf{z} &= \mathbf{P}^T \mathbf{d}_k \end{aligned}$$

ie, $\mathbf{d}_k = \mathbf{P}\mathbf{z}$. From

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

we solve the following systems, first for \mathbf{z} and then for \mathbf{d}_k

$$\begin{aligned} \mathbf{R}_1 \mathbf{z} &= \mathbf{b} \\ \mathbf{d}_k &= \mathbf{P}\mathbf{z} \end{aligned}$$