

Nonlinear Least Squares

(Com S 477/577 Notes)

Yan-Bin Jia

Nov 27, 2018

1 Introduction

In least-squares fitting, we fit a model function $F(\mathbf{x}; \mathbf{c})$, controlled by the parameter vector $\mathbf{c} = (c_1, \dots, c_k)^T$, to n data points $(\mathbf{x}_1, f_1), \dots, (\mathbf{x}_n, f_n)$ of some unknown function f . This is achieved by minimizing the error function

$$\begin{aligned}\omega(\mathbf{c}) &= \sum_{i=1}^n (f_i - F(\mathbf{x}_i; \mathbf{c}))^2 \\ &= (\mathbf{f} - \boldsymbol{\alpha})^T (\mathbf{f} - \boldsymbol{\alpha}),\end{aligned}\tag{1}$$

where

$$\mathbf{f} = \begin{pmatrix} f_1 \\ \vdots \\ f_n \end{pmatrix} \quad \text{and} \quad \boldsymbol{\alpha}(\mathbf{c}) = \begin{pmatrix} F(\mathbf{x}_1; \mathbf{c}) \\ \vdots \\ F(\mathbf{x}_n; \mathbf{c}) \end{pmatrix}.$$

Linear least squares use a model function in the form of $F = c_1\phi_1(\mathbf{x}) + \dots + c_k\phi_k(\mathbf{x})$, where $\phi_1(\mathbf{x}), \dots, \phi_k(\mathbf{x})$ are basis functions.

It is more often that the model function assumes a nonlinear form. For example, to estimate the location of a target point \mathbf{x} , we take range measurements at n known locations \mathbf{p}_i , $1 \leq i \leq n$, to obtain its distances from these locations:

$$d_i = \|\mathbf{x} - \mathbf{p}_i\| + \nu_i,$$

where ν_i s are the measurement noises. Then, we use least squares to estimate \mathbf{x} as the location which minimizes the error

$$\sum_{i=1}^n (\|\mathbf{x} - \mathbf{p}_i\| - d_i)^2.$$

GPS works in the above manner.

When the model function is nonlinear in the parameters, we need to find their values numerically, starting with an initial guess $\mathbf{c}^{(0)}$. At each step, the parameter vector $\mathbf{c}^{(i)}$ is replaced by a new estimate $\mathbf{c}^{(i+1)} = \mathbf{c}^{(i)} + \boldsymbol{\delta}$, where $\boldsymbol{\delta} = (\delta_1, \dots, \delta_k)^T$. If we use the gradient descent method, then $\boldsymbol{\delta}$ should be along the negative gradient $\nabla\omega$, i.e., along

$$\boldsymbol{\delta}_g \equiv \frac{\partial \boldsymbol{\alpha}}{\partial \mathbf{c}} (\mathbf{f} - \boldsymbol{\alpha}).\tag{2}$$

The partial derivative $\partial\boldsymbol{\alpha}/\partial\mathbf{c}$ is the Jacobian J , which is an $n \times k$ matrix given below:

$$J = \begin{pmatrix} \frac{\partial F}{\partial c_1}(\mathbf{x}_1, \mathbf{c}) & \cdots & \frac{\partial F}{\partial c_k}(\mathbf{x}_1, \mathbf{c}) \\ \vdots & \ddots & \vdots \\ \frac{\partial F}{\partial c_1}(\mathbf{x}_n, \mathbf{c}) & \cdots & \frac{\partial F}{\partial c_k}(\mathbf{x}_n, \mathbf{c}) \end{pmatrix}.$$

Gradient descent has slow convergence since every two adjacent moving directions are perpendicular to each other, which results in the second movement undoing part of the progress made by the first movement. The step size also has to be carefully controlled, otherwise failure to converge may happen frequently.

2 The Gauss-Newton Method

The Gauss-Newton method [3] linearizes the model function $F(\mathbf{x}; \mathbf{c})$ using the Taylor series, for $i = 1, \dots, n$,

$$\begin{aligned} F(\mathbf{x}_i; \mathbf{c} + \boldsymbol{\delta}) &\approx F(\mathbf{x}_i; \mathbf{c}) + \sum_{j=1}^k \frac{\partial F}{\partial c_j}(\mathbf{x}_i; \mathbf{c}) \cdot \delta_j \\ &= F(\mathbf{x}_i; \mathbf{c}) + \left(\frac{\partial F}{\partial c_1}, \dots, \frac{\partial F}{\partial c_k} \right) \Big|_{(\mathbf{x}_i, \mathbf{c})} \boldsymbol{\delta}. \end{aligned}$$

Let us gather all k equations into one:

$$\boldsymbol{\alpha}(\mathbf{c} + \boldsymbol{\delta}) = \begin{pmatrix} F(\mathbf{x}_1; \mathbf{c}) \\ \vdots \\ F(\mathbf{x}_n; \mathbf{c}) \end{pmatrix} + J\boldsymbol{\delta}. \quad (3)$$

Now we approximate the least-squares error (1) as the parameter vector changes its value from \mathbf{c} to $\mathbf{c} + \boldsymbol{\delta}$:

$$\begin{aligned} \omega(\mathbf{c} + \boldsymbol{\delta}) &= (\mathbf{f} - \boldsymbol{\alpha}(\mathbf{c} + \boldsymbol{\delta}))^T (\mathbf{f} - \boldsymbol{\alpha}(\mathbf{c} + \boldsymbol{\delta})) \\ &\approx (\mathbf{f} - \boldsymbol{\alpha} - J\boldsymbol{\delta})^T (\mathbf{f} - \boldsymbol{\alpha} - J\boldsymbol{\delta}). \end{aligned}$$

The increment $\boldsymbol{\delta}$ is found by setting the partial derivative with respect to $\boldsymbol{\delta}$ to zero, namely,

$$-2(\mathbf{f} - \boldsymbol{\alpha} - J\boldsymbol{\delta})^T J = 0.$$

Thus, we obtain

$$\boldsymbol{\delta} = \boldsymbol{\delta}_t \equiv (J^T J)^{-1} J^T (\mathbf{f} - \boldsymbol{\alpha}), \quad (4)$$

where in the above the positive semidefinite $k \times k$ matrix $J^T J$ is invertible provided that all the columns of J are linearly independent. For later use we rewrite the above equation as

$$J^T J \boldsymbol{\delta}_t = J^T (\mathbf{f} - \boldsymbol{\alpha}). \quad (5)$$

In practice, \mathbf{c} is corrected by a fraction of $\boldsymbol{\delta}_t$ so as not to step out of the neighborhood within which the approximation (3) is good enough not to cause divergence. Failure to converge has been found to be not common with a step size $s\boldsymbol{\delta}_t$, $0 < s \leq 1$.

3 The Levenberg-Marquardt Method

Marquardt [6] observed that $\boldsymbol{\delta}_t$ is almost $\frac{\pi}{2}$ away from $\boldsymbol{\delta}_g$ in most problems because of the severe elongation of the surface $\omega(\mathbf{c})$. This could make the Gauss-Newton method still very slow. He realized that any improved method would in some sense interpolate between $\boldsymbol{\delta}_t$ and $\boldsymbol{\delta}_g$. At a suggestion from Levenberg to minimize ω locally, Marquardt [6] first modified the correction $\boldsymbol{\delta}_t$ given in (4) and use instead

$$\boldsymbol{\delta}_l = \left(J^T J + \lambda I_k \right)^{-1} J^T (\mathbf{f} - \boldsymbol{\alpha}), \quad (6)$$

where I_k is the $k \times k$ identity matrix. Note that the matrix $J^T J + \lambda I_k$ is positive definite for a positive λ , which is a damping factor. The iteration formula is thus

$$\mathbf{c}^{(i)} = \mathbf{c}^{(i-1)} + \left(J^T J + \lambda I_k \right)^{-1} J^T (\mathbf{f} - \boldsymbol{\alpha}), \quad (7)$$

where J and $\boldsymbol{\alpha}$ are evaluated at $\mathbf{c}^{(i-1)}$.

To provide theoretical basis, Marquardt established the following properties:

- i. $\omega(\mathbf{c} + \boldsymbol{\delta}_l) \leq \omega(\mathbf{c} + \boldsymbol{\delta})$ for all $\boldsymbol{\delta}$ satisfying $\|\boldsymbol{\delta}\| = \|\boldsymbol{\delta}_l\|$;
- ii. $\|\boldsymbol{\delta}_l(\lambda)\|$ decreases to zero monotonically as $\lambda \rightarrow \infty$;
- iii. $\boldsymbol{\delta}_l$ rotates from $\boldsymbol{\delta}_t$ to $\boldsymbol{\delta}_g$ monotonically as $\lambda \rightarrow \infty$.

The second insight Marquardt had was to replace the identity matrix I_k in (6) with a diagonal matrix $D = \text{diag}(d_1, \dots, d_k)$, where d_1, \dots, d_k are the diagonal elements of $J^T J$. This scaling is adapted to the curvature of the fitting error $\omega(\mathbf{c})$ so there is larger movement along the direction if the gradient is smaller, and it also makes λ dimensionless. More specifically, we denote $E = \text{diag}(\sqrt{d_1}, \dots, \sqrt{d_k})$ and introduce

$$\begin{aligned} K &= E^{-1} J^T J E^{-1}, \\ \mathbf{g} &= E^{-1} J^T (\mathbf{f} - \boldsymbol{\alpha}), \\ \boldsymbol{\sigma}_t &= E \boldsymbol{\delta}_t, \end{aligned}$$

and then transform (5) into

$$K \boldsymbol{\sigma}_t = \mathbf{g}.$$

In parallel to the change from (4) to (6), we modify the above equation to obtain a scaled correction:

$$\boldsymbol{\sigma}_m = (K + \lambda I_k)^{-1} \mathbf{g}.$$

Then, the increment to \mathbf{c} is

$$\begin{aligned} \boldsymbol{\delta}_m &= E^{-1} \boldsymbol{\sigma}_m \\ &= E^{-1} \left(E^{-1} J^T J E^{-1} + \lambda I_k \right)^{-1} E^{-1} J^T (\mathbf{f} - \boldsymbol{\alpha}) \\ &= \left(J^T J + \lambda D \right)^{-1} J^T (\mathbf{f} - \boldsymbol{\alpha}). \end{aligned} \quad (8)$$

This gives rise to the following iteration formula:

$$\mathbf{c}^{(i)} = \mathbf{c}^{(i-1)} + \left(J^T J + \lambda D \right)^{-1} J^T (\mathbf{f} - \boldsymbol{\alpha}). \quad (9)$$

A sufficiently large λ will rotate $\boldsymbol{\delta}_m$ towards the gradient $\boldsymbol{\delta}_g$ enough to decrease the error ω . However, a value of λ too large would make the method work like steepest descent, thereby inheriting its property of rapid initial progress followed by increasingly slower progress. Such a strategy would be a poor one globally. If λ is too big, $\mathbf{c}^{(i)}$ is too close to $\mathbf{c}^{(i-1)}$ to make fast progress. And if λ is too small, $\mathbf{c}^{(i)}$ is too far from $\mathbf{c}^{(i-1)}$, and linearization yields a poor approximation. The value of λ needs to be chosen small so minimization of ω is done in the maximum neighborhood over which the linearization (3) is adequate. Based on these considerations, Marquardt proposed an iterative algorithm below. Denote by $\lambda^{(i)}$, $\mathbf{c}^{(i)}$, and $\omega^{(i)}$ the values of λ , \mathbf{c} , and ω at the end of the i th iteration.

1. Let $\nu > 1$, and initialize $\lambda^{(0)} = 10^{-2}$.
2. Set λ to $\lambda^{(i-1)}/\nu$ and $\lambda^{(i-1)}$, respectively. Obtain two values of $\boldsymbol{\delta}_m$ according to (8), and compute $\omega(\mathbf{c}^{(i-1)} + \boldsymbol{\delta}_m)$.
3. Let ω_1 and ω_2 be the computed values of ω using $\lambda^{(i-1)}/\nu$ and $\lambda^{(i-1)}$, respectively.
 - (a) If $\omega_1 \leq \omega^{(i-1)}$, Set $\lambda^{(i)} \leftarrow \lambda^{(i-1)}/\nu$ and $\omega^{(i)} \leftarrow \omega_1$.
 - (b) If $\omega_1 > \omega^{(i-1)}$ but $\omega_2 \leq \omega^{(i-1)}$, set $\lambda^{(i)} \leftarrow \lambda^{(i-1)}$ and $\omega^{(i)} \leftarrow \omega_2$.
 - (c) If $\omega_1 > \omega^{(i-1)}$ and $\omega_2 > \omega^{(i-1)}$, start with $\lambda \leftarrow \lambda^{(i-1)}$ and repeatedly set $\lambda \leftarrow \lambda\nu$ until $\omega(\mathbf{c}^{(i-1)} + \boldsymbol{\delta}_m) \leq \omega^{(i-1)}$. Set $\lambda^{(i)}$ and $\omega^{(i)}$ to the ending values of λ and ω .
4. Set $\mathbf{c}^{(i)} \leftarrow \mathbf{c}^{(i-1)} + \boldsymbol{\delta}_m$, where $\boldsymbol{\delta}_m$ is evaluated from (8) with $\lambda = \lambda^{(i)}$.
5. Set $i \leftarrow i + 1$. Repeat steps 2 through 4.

The Levenberg-Marquardt method can be applied to other types of optimization and root finding problems. For example, suppose we want to solve a system of equations

$$f_i(\mathbf{x}) = 0, \quad i = 1, \dots, k.$$

We can turn it into a least-squares problem to minimize the error

$$\omega = \sum_{i=1}^k f_i^2(\mathbf{x}),$$

and apply the Levenberg-Marquardt method subsequently.¹

4 Curve and Surface Fitting

In computer vision, robotics, and geometric modeling, it is often important to derive some implicit representation of an object from raw shape data obtained using an imaging or touch sensing device. Such an implicit representation often takes the form of the zero set of a polynomial of even degree.²

¹Mathematica uses the Levenberg-Marquardt method in its command `FindMinimum` when given the `Method -> LevenbergMarquardt` option.

²It is known that the zero sets of polynomials of odd degrees always define unbounded curves and surfaces.

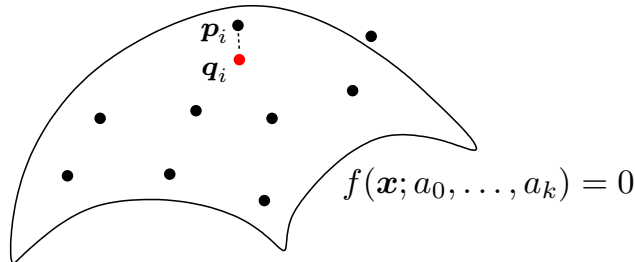


Figure 1: Fitting a surface over points.

In practice, quartics or superquartics [1] have often been the choices for describing surfaces of moderate geometric complexities.

An implicit representation in polynomial form has several advantages. First, it gives us a succinct shape description with few parameters so storage of a large amount of raw data is no longer necessary. Second, using the polynomial we can compute some algebraic invariants (which are expressions in terms of the polynomial's coefficients) to recognize the object from a set of known models. Third, we would be able to differentiate the shape to obtain local geometry indices such as curvature and torsion. Fourth, evaluating the polynomial at a point would easily tell us whether the point is inside, outside, or on the surface. Finally, a polynomial fit handles noisy data effectively in the sense that most shape irregularities due to noise are filtered out.

In surface fitting, a family of quartic polynomials with stably bounded zero sets is often considered. Choice of the family takes into account factors such as the effectiveness of representation and the computational efficiency of fitting. Polynomial coefficients are the parameters to be determined by, say, a least-squares method. For example, we may minimize the sum of squares of the distances from individual data points to the polynomial surface.

Suppose we need to fit a curve or surface over n data points $\mathbf{p}_1, \dots, \mathbf{p}_n$ in \mathbb{R}^2 or \mathbb{R}^3 . The family of polynomial curves or surfaces that we consider assumes the form

$$f(\mathbf{x}; a_0, \dots, a_k) = 0, \quad (10)$$

where $\mathbf{x} = (x, y)$ or (x, y, z) , and a_0, \dots, a_k are the coefficients. Denote by d_i the distance from \mathbf{p}_i to the surface (10), for $i = 1, \dots, n$. Then the fitting problem can be formulated in a least-squares fashion as

$$\min_{a_0, \dots, a_k} \sum_{i=1}^n d_i^2. \quad (11)$$

The issue, of course, is that we cannot determine d_i , $1 \leq i \leq n$, until a_0, \dots, a_k are known, which conversely depends on the knowledge of d_i according to (11). To get out of this chicken-and-egg situation, let us approximate d_i as follows. Denote by \mathbf{q}_i the closest point to \mathbf{p}_i on the surface to be determined. Obviously, $f(\mathbf{q}_i; a_0, \dots, a_k) = 0$. The best fitting f will place \mathbf{q}_i close enough to \mathbf{p}_i . So we can approximate the value of f at \mathbf{q}_i using Taylor's series at \mathbf{p}_i , discarding all terms of the second order and above:

$$f(\mathbf{q}_i; a_0, \dots, a_k) \approx f(\mathbf{p}_i; a_0, \dots, a_k) + \nabla f(\mathbf{p}_i; a_0, \dots, a_k) \cdot (\mathbf{q}_i - \mathbf{p}_i).$$

Since the left hand side of the above is zero, we have

$$\nabla f(\mathbf{p}_i; a_0, \dots, a_k) \cdot (\mathbf{q}_i - \mathbf{p}_i) \approx -f(\mathbf{p}_i; a_0, \dots, a_k).$$

Because \mathbf{q}_i and \mathbf{p}_i are very close to each other, the vector $\mathbf{q}_i - \mathbf{p}_i$ and the gradient $\nabla f(\mathbf{p}_i; a_0, \dots, a_k)$ are nearly parallel. The above equation yields an approximation:

$$d_i = \|\mathbf{q}_i - \mathbf{p}_i\| \approx \frac{|f(\mathbf{p}_i; a_0, \dots, a_k)|}{\|\nabla f(\mathbf{p}_i; a_0, \dots, a_k)\|}.$$

With a substitution of the above approximation into (11), we reformulate the fitting problem as

$$\min_{a_0, \dots, a_k} \frac{f^2(\mathbf{p}_i; a_0, \dots, a_k)}{\|\nabla f(\mathbf{p}_i; a_0, \dots, a_k)\|^2}.$$

The reader is referred to [2, 5, 7] for more on surface fitting techniques and experiments.

5 Surface Patch Reconstruction

A robotic hand can reconstruct an unknown surface patch by touch [4]. The idea is to track along three concurrent³ curves on the surface, obtaining tactile data points (x_k, y_k, z_k) , $1 \leq k \leq n$, in the meantime. Each curve is the intersection of the patch with a separate plane (called the sampling plane) within which the tracking motion is presently constrained. Denote by \mathbf{p} the intersection point of the three curves. Through fitting a parabola to the data points along each curve, its curvature at \mathbf{p} is estimated. The surface's normal curvature at \mathbf{p} in the tangent direction of this curve is the product of the (estimated) curvature with the cosine of the angle between the sampling plane and the normal plane containing the tangent direction.

Three normal curvatures are thus obtained. From them and the relative orientations of the corresponding tangent vectors at \mathbf{p} , we solve for the principal curvatures κ_1 and κ_2 and the Darboux frame at the point \mathbf{p} . Under this frame, the surface patch locally takes the form

$$z(x, y) = \frac{1}{2}(\kappa_1 x^2 + \kappa_2 y^2) + \sum_{3 \leq i+j \leq d} a_{ij} x^i y^j, \quad (12)$$

where the terms of degrees above two are added to describe a larger area of the surface. The coefficients of the polynomial, gathered into a vector \mathbf{a} , are determined in a least-squares sense

$$\min_{\mathbf{a}} f(\mathbf{a}) \quad (13)$$

where

$$f(\mathbf{a}) = \frac{1}{n} \sum_{k=1}^n (z(x_k, y_k) - z_k)^2. \quad (14)$$

However, the resulting surface described by $z(x, y)$ may “fold” many times unrealistically, as indicated by a high total absolute Gaussian curvature:

$$g(\mathbf{a}) = \iint_D |K(x, y)| \cdot \sqrt{1 + z_x^2 + z_y^2} dx dy. \quad (15)$$

To get rid of unnecessary “folds”, we minimize the total absolute Gaussian curvature subject to the constraint that the surface fit should “pass through” the sampled data points:

$$\min_{\mathbf{a}} g(\mathbf{a}) \quad \text{subject to } f(\mathbf{a}) = \mathbf{0}.$$

³intersecting at one point

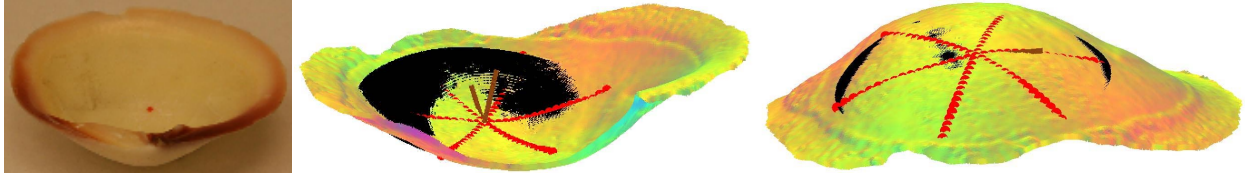


Figure 2: Reconstruction of the surface of a shell from three lines of data points

The technique of Lagrange multipliers is then applied.

Figure 2 displays a shell and two views (inside-out and outside-in) of the reconstructed patch⁴ overlaid onto its mesh model generated by NextEngine’s desktop 3-D scanner (accuracy 0.127mm). The data points from tracking are shown in red while the estimated Darboux frame in brown. The average distance of the mesh vertices in the overlay area to the patch is 0.2283mm, indicating a close match.

References

- [1] A. H. Barr. Superquadrics and angle-preserving transformations. *IEEE Computer Graphics and Applications*, 1(1):11–23, 1981.
- [2] M. M. Blane, Z. Lei, H. Civi, D. B. Cooper. The 3L algorithm for fitting implicit polynomial curves and surfaces to data. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(3):298–313, 2000.
- [3] H. O. Hartley. The modified G-N method for the fitting of nonlinear regression functions by least squares. *Technometrics*, 3(2):269–280, 1961.
- [4] Y.-B. Jia and J. Tian. Surface patch reconstruction from “one-dimensional” tactile data. *IEEE Transactions on Automation Science and Engineering*, 7(2):400–407, 2010.
- [5] D. Keren and D. Cooper. Describing complicated objects by implicit polynomials. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(1):38–53, 1994.
- [6] D. W. Marquardt. An algorithm for least-squares estimation of nonlinear parameters. *Journal of the Society for Industrial and Applied Mathematics*, 11(2):431–441, 1963.
- [7] G. Taubin, F. Cukierman, S. Sullivan, J. Ponce, and D. J. Kriegman. Parametrized families of polynomials for bounded algebraic curve and surface fitting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(3):287–303, 1994.

⁴In the reconstruction, d was set to be 4.