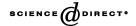


Available online at www.sciencedirect.com



Journal of COMPLEXITY

Journal of Complexity 20 (2004) 484-492

http://www.elsevier.com/locate/jco

On the complexity of curve fitting algorithms

N. Chernov,* C. Lesort, and N. Simányi

Department of Mathematics, University of Alabama at Birmingham, Birmingham, AL 35294, USA
Received 29 September 2003; accepted 15 January 2004

Abstract

We study a popular algorithm for fitting polynomial curves to scattered data based on the least squares with gradient weights. We show that sometimes this algorithm admits a substantial reduction of complexity, and, furthermore, find precise conditions under which this is possible. It turns out that this is, indeed, possible when one fits circles but not ellipses or hyperbolas.

© 2004 Elsevier Inc. All rights reserved.

Keywords: Least-squares fit; Curve fitting; Algebraic gradient weight fitting; Complexity

In many applications one needs to fit a curve described by a polynomial equation

$$P(x, y; \boldsymbol{\Theta}) = 0$$

(here Θ denotes the vector of unknown parameters) to experimental data (x_i, y_i) , i = 1, ..., n. In this equation P is a polynomial in x and y, and its coefficients are either unknown parameters or functions of unknown parameters. For example, a number of recent publications [5,6,9] are devoted to the problem of fitting quadrics $Ax^2 + Bxy + Cy^2 + Dx + Ey + F = 0$, in which case $\Theta = (A, B, C, D, E, F)$ is the parameter vector. The problem of fitting circles, given by equation $(x - a)^2 + (y - b)^2 - R^2 = 0$ with three parameters a, b, R, also arises in practice [4,8].

It is standard to assume that the data (x_i, y_i) are noisy measurements of some true (but unknown) points (\bar{x}_i, \bar{y}_i) on the curve, see [1,3,7,8] for details. The noise vectors $e_i = (x_i - \bar{x}_i, y_i - \bar{y}_i)$ are then assumed to be independent Gaussian vectors with zero mean and a scalar covariance matrix, $\sigma^2 I$. In this case the maximum likelihood estimate of Θ is given by the *orthogonal least-squares fit* (OLSF), which is based on

E-mail address: chernov@math.uab.edu (N. Chernov).

0885-064X/\$- see front matter © 2004 Elsevier Inc. All rights reserved. doi:10.1016/j.jco.2004.01.004

^{*}Corresponding author. Fax: 1-205-934-9025.

the minimization of the function

$$\mathscr{F}(\Theta) = \sum_{i=1}^{n} d_i^2, \tag{1}$$

where d_i denotes the distance from the point (x_i, y_i) to the curve $P(x, y; \Theta) = 0$.

Under these assumptions the OLSF is statistically optimal—it provides estimates of Θ whose covariance matrix attains its Rao-Cramer lower bound [3,7,8]. The OLSF is widely used in practice, especially when one fits simple curves such as lines or circles. However, for more general curves the OLSF becomes intractable, because the precise distance d_i is hard to compute. In those cases one resorts to various alternatives, and the most popular one is the *algebraic fit* (AF) based on the minimization of

$$\mathscr{F}_{a}(\Theta) = \sum_{i=1}^{n} w_{i} \left[P(x_{i}, y_{i}; \Theta) \right]^{2}, \tag{2}$$

where $w_i = w(x_i, y_i; \Theta)$ are suitably defined weights. The choice of the weight function $w(x, y; \Theta)$ is important. The AF is known [3] to provide a statistically optimal stimate of Θ (in the sense that the covariance matrix will attain its Rao–Cramer lower bound) if and only if the weight function satisfies

$$w(x, y; \Theta) = a(\Theta) / ||\nabla P(x, y; \Theta)||^2$$
(3)

for all points x, y on the curve, i.e. such that $P(x, y; \Theta) = 0$. Here $\nabla P = (\partial P/\partial x, \partial P/\partial y)$ is the gradient vector of the polynomial P, and $a(\Theta) > 0$ may be an arbitrary function of Θ (in practice, one simply sets $a(\Theta) = 1$). Any other choice of w will result in the loss of accuracy, see [3]. We call $w(x, y; \Theta)$ a gradient weight function if it satisfies (3) for all x, y on the curve $P(x, y; \Theta) = 0$. The AF (2) with a gradient weight function $w(x, y; \Theta)$ is commonly referred to as the gradient weighted algebraic fit (GRAF). It was introduced in the mid-1970s [14] and recently became standard for polynomial curve fitting, see, for example, [5,9,13].

Even though the GRAF is much cheaper than the OLSF, it is still a nonlinear problem requiring iterative methods. For example, in a popular *reweight procedure* [11,13] one uses the *k*th approximation $\Theta^{(k)}$ to compute the weights $w_i = w(x_i, y_i; \Theta^{(k)})$ and then finds $\Theta^{(k+1)}$ by minimizing (2) regarding the just computed w_i 's as constants. Note that if the parameters Θ are the coefficients of P, then (2), with fixed weights, becomes a quadratic function in Θ , and its minimum can be easily found. Another algorithm is based on solving the equation $\nabla_{\Theta} \mathcal{F}_a(\Theta) = 0$, i.e.

$$\sum P_i^2 \nabla_{\Theta} w_i + 2 \sum w_i P_i \nabla_{\Theta} P_i = 0 \tag{4}$$

for which various iterative schemes could be used. In the case of fitting quadrics, for example, the most advanced algorithms are the renormalization method [7], the heteroscedastic error-in-variables method [9] and the fundamental numerical scheme [5]. In all these algorithms, one needs to evaluate $\mathcal{O}(n)$ terms at each iteration. Therefore, the complexity of those algorithms is $\mathcal{O}(kn)$, where k is the number of iterations. Moreover, each algorithm requires access to individual coordinates x_i, y_i

of the data points at each iteration. These difficulties can be sometimes avoided in a remarkable way, as we show next.

Suppose, we need to fit circles given by equation

$$P(x, y) = (x - a)^{2} + (y - b)^{2} - R^{2} = 0.$$

Then, we have

$$||\nabla P(x, y; \Theta)||^2 = 4(x - a)^2 + 4(y - b)^2 = 4P(x, y) + 4R^2$$
(5)

hence $||\nabla P(x, y; \Theta)||^2 = 4R^2$ for all the points (x, y) lying on the circle P(x, y) = 0, and we can set $w(x, y; \Theta) = 1/R^2$. Therefore

$$\mathcal{F}_{a}(a,b,R) = \sum_{i=1}^{n} R^{-2} [x_{i}^{2} + y_{i}^{2} - 2ax_{i} - 2by_{i} + a^{2} + b^{2} - R^{2}]^{2}$$

$$= R^{-2} [z_{1} + az_{2} + bz_{3} + a^{2}z_{4} + b^{2}z_{5} + abz_{6} + cz_{7} + acz_{8}$$

$$+ bcz_{9} + c^{2}n], \tag{6}$$

where we denoted $c = a^2 + b^2 - R^2$ for brevity, and

$$z_1 = \sum (x_i^2 + y_i^2)^2, \quad z_2 = -4 \sum x_i (x_i^2 + y_i^2), \dots$$

are some expressions involving x_i and y_i only.

The minimization of (6) is still a nonlinear problem requiring iterative methods [4,2,10], but it has obvious advantages over the reweight procedure described above and other generic methods for solving Eq. (4). First of all, the values of $z_1, ..., z_9$ only need to be computed once, and then the cost of minimization of (6) will not depend on n anymore. Thus, the complexity of this algorithm is $\mathcal{O}(n) + \mathcal{O}(k)$, where $\mathcal{O}(n)$ is the cost of evaluation of $z_1, ..., z_9$ and $\mathcal{O}(k)$ is the cost of some k iterations spent on the subsequent minimization of $\mathcal{F}_a(a,b,R)$. Moreover, once the values of $z_1, ..., z_9$ are computed and stored, the coordinates x_i, y_i can be destroyed. Practically, $z_1, ..., z_9$ can be computed "on-line", when the data are collected. The minimization procedure per se can be implemented "off-line", without storage of (or access to) the data points. The quantities $z_1, ..., z_9$ here play the role of sufficient statistics.

This improvement is essential for applications where one fits circular arc to data on a mass scale (often on-line) so that the speed of fitting algorithms is a factor. Such are, for example, modern experiments in high-energy physics (charged particles move along circular arcs in a homogeneous magnetic field, and computing the particle's energy requires estimation of the radius of the arc). Physicists process millions of images from nuclear accelerators on-line in a search of rare collisions (events), hence speedy fitting procedures are in high demand. For other experiments involving circle fitting see [2] and references therein.

Inspired by the above example, we might say that the problem of fitting a polynomial curve $P(x, y; \Theta) = 0$ admits a reduction of complexity if there are ℓ functions $z_j(x_1, y_1, ..., x_n, y_n)$, $1 \le j \le \ell$, with ℓ being independent of n and Θ , and a

gradient weight function $w(x, y; \Theta)$ such that

$$\mathscr{F}_{\mathbf{a}} = f(z_1, \dots, z_\ell; \Theta) \tag{7}$$

i.e. \mathcal{F}_a is a function of $z_1, ..., z_\ell$ and Θ only.

This definition does not suggest how to find the functions $z_1, ..., z_\ell$ in practical terms, though. Since \mathscr{F}_a is given by (2) with $P(x_i, y_i; \Theta)$ being a polynomial in x_i, y_i , then the most natural (if not the only) way to construct the functions $z_1, ..., z_\ell$ is to express the gradient weight function (3) in the form

$$w(x, y; \boldsymbol{\Theta}) = \sum_{k=1}^{K} C_k(\boldsymbol{\Theta}) D_k(x, y), \tag{8}$$

where C_k are functions of the parameter vector Θ alone, and D_k are functions of x and y only (here the number of terms, K, must be independent of Θ). Indeed, suppose that representation (8) is found. Since P^2 is a polynomial in x, y, we can expand it as

$$P^{2}(x,y) = \sum_{p,q} c_{p,q} x^{p} y^{q},$$

where $c_{p,q} = c_{p,q}(\Theta)$ denote its coefficients. Now the function \mathscr{F}_a can be evaluated as

$$\mathcal{F}_{a} = \sum_{k=1}^{K} \sum_{p,q} C_{k}(\Theta) c_{p,q}(\Theta) \sum_{i=1}^{n} x_{i}^{p} y_{i}^{q} D_{k}(x_{i}, y_{i})$$
$$= \sum_{k=1}^{K} \sum_{p,q} C_{k}(\Theta) c_{p,q}(\Theta) z_{k,p,q},$$

where

$$z_{k,p,q} = \sum_{i=1}^{n} x_i^p y_i^q D_k(x_i, y_i).$$

The values of $z_{k,p,q}$ depend on the data x_i, y_i only, hence we obtain the desired representation (7). Therefore, (8) implies (7). We believe that the converse is also true, i.e. the conditions (7) and (8) are actually equivalent, but we do not attempt to prove that.

Motivated by the above considerations, we adopt the following definition: the problem of fitting a polynomial curve $P(x, y; \Theta) = 0$ admits a reduction of complexity if the gradient weight function (3) can be expressed in form (8).

As we have seen, the problem of fitting circles admits a reduction of complexity (and so does the simpler problem of fitting lines). Now if the problem of fitting ellipses and/or hyperbolas admitted a reduction of complexity as defined above, we would be able to dramatically improve the known GRAF algorithms [5,7,9]. Unfortunately, this is impossible—there are deep mathematical reasons which prevent a reduction of complexity in the case of ellipses, hyperbolas, and parabolas.

In this paper, we find general conditions on the polynomial $P(x, y; \Theta)$ under which the problem of fitting the curve $P(x, y; \Theta) = 0$ allows a reduction of complexity. It turns out that lines and circles satisfy these conditions, but ellipses, hyperbolas, and

parabolas do not. Our results thus demonstrate (in a rigorous mathematical way) that fitting noncircular conics is an intrinsically more complicated problem than fitting circles or lines.

We will assume here that P is an irreducible polynomial, i.e. it cannot be expressed as a product $P = P_1(x, y; \Theta)P_2(x, y; \Theta)$ of polynomials of lower degrees (because otherwise fitting the curve P = 0 can be reduced to fitting simpler curves $P_1 = 0$ and $P_2 = 0$).

For convenience, let us denote

$$Q(x, y; \Theta) := ||\nabla P(x, y; \Theta)||^2 = (\partial P/\partial x)^2 + (\partial P/\partial y)^2.$$

Clearly, $Q(x, y; \Theta)$ is itself a polynomial in x and y. Our subsequent arguments will involve some facts from complex analysis. We will treat x and y as complex, rather than real, variables.

Theorem. The problem of fitting curves $P(x, y; \Theta) = 0$ admits a reduction of complexity (as defined above) under the condition that the system of polynomial equations

$$P(x,y) = 0,$$

 $Q(x,y) = 0,$
(9)

has no solutions, real or complex, for any Θ .

Before we prove our theorem, we shall show how to use it. For the problem of fitting circles, we have already computed $Q = 4P + 4R^2$, see (5), hence system (9) has indeed no solutions for nondegenerate circles (for which $R \neq 0$).

When using the theorem, the following *invariance* property will be helpful. Let $(x,y)\mapsto (\tilde{x},\tilde{y})$ be a transformation of the xy plane that is a composition of translations, rotations, mirror reflections and similarities (the latter are defined by $(x,y)\mapsto (cx,cy)$ for some $c\neq 0$). Denote by $\tilde{P}(\tilde{x},\tilde{y})$ the polynomial P in the new coordinates \tilde{x},\tilde{y} . Then system (9) has a solution (real or complex) if and only if the corresponding system

$$\tilde{P}(\tilde{x}, \tilde{y}) = 0,$$

 $\tilde{Q}(\tilde{x}, \tilde{y}) = 0,$

has a solution, real or complex. Here $\tilde{Q} = ||\nabla \tilde{P}||^2$. This simple fact, which can be verified directly by the reader, allows us to simplify the polynomial P(x, y) before applying the theorem.

Consider the problem of fitting ellipses and hyperbolas. By using a translation and rotation of the xy plane we can always reduce the polynomial P to a canonical form $ax^2 + by^2 + c = 0$ (with $a \ne b$ and $abc \ne 0$). Then $Q = 4a^2x^2 + 4b^2y^2$ and we arrive at a system of equations

$$ax^2 + by^2 + c = 0,$$

 $a^2x^2 + b^2y^2 = 0.$

It is easy to see that it always has a solution

$$x = \pm \sqrt{\frac{bc}{a(a-b)}}, \quad y = \pm \sqrt{-\frac{ac}{b(a-b)}}$$

(note that x or y may be an imaginary number, which is allowed by our theorem). Therefore, the problem does not admit a reduction of complexity.

If our curve is a parabola, then we can use its canonical equation $y = cx^2$ for c > 0, hence $P = y - cx^2$ and $Q = 4c^2x^2 + 1$. Here again we have a common zero of P and Q at the point x = i/2c and y = -1/4c. Thus, no conic sections (except circles) satisfy the conditions of our theorem.

We now prove our theorem. Since $w(x, y; \Theta)$ must be a gradient weight function, the requirement (8) is equivalent to

$$\frac{1}{Q(x,y)} = \sum_{k=1}^{K} C_k(\Theta) D_k(x,y) \quad \text{whenever} \quad P(x,y) = 0$$
 (10)

(here we incorporated the factor $a(\Theta)$ into the coefficients $C_k(\Theta)$, for convenience). We emphasize that the left identity in (10) does not have to hold on the entire xy plane, it only has to hold on the curve P(x,y) = 0. If we denote that curve by \mathcal{L} , then (10) can be restated as

$$\frac{1}{Q(x,y)} = \sum_{k=1}^{K} C_k(\Theta) D_k(x,y) \quad \text{whenever } (x,y) \in \mathcal{L}.$$
(11)

The functions $D_k(x,y)$ in (10) cannot be arbitrary, they must be easily computable, i.e. available in the machine arithmetics. That is, they must be combinations of elementary functions—polynomials, exponentials, logarithms, trigonometric functions, etc. In that case $D_k(x,y)$ are analytic functions of x and y. Therefore, they have analytic extensions to the two-dimensional complex plane \mathbb{C}^2 . We note that they do not need be *entire functions*, i.e. analytic everywhere in \mathbb{C}^2 , they may have some singularities. For example, the function $(1+x^2+y^2)^{-1}$ is analytic in \mathbb{R}^2 but has singularities in \mathbb{C}^2 , e.g. the point x=i and y=0 is its singularity. Also, those extensions maybe multivalued functions (examples are $\ln x$ or \sqrt{x}).

Now, the following function will also be analytic in \mathbb{C}^2 :

$$G(x, y) = 1 - Q(x, y) \sum_{k=1}^{K} C_k(\Theta) D_k(x, y)$$

since it is a combination of analytic functions. By (11), it vanishes on the curve \mathscr{L} in the real xy plane. Consider the subset $\mathscr{L}|\subset\mathbb{C}^2$ defined by the equation P(x,y)=0, where x and y are treated as complex variables. Note that \mathscr{L} is a curve on the two-dimensional manifold \mathscr{L} . We will prove that the function G(x,y) vanishes on the entire \mathscr{L} .

Since P(x, y) is an irreducible polynomial, \mathscr{Z} is a connected algebraic variety (a Riemann surface), hence it admits a complex parameterization (a complex coordinate, z), and the restriction of the function G onto \mathscr{Z} will be an analytic

function of z. It is known in complex analysis that if an analytic function $G(z), z \in \mathbb{C}$, vanishes on an infinite set having an accumulation point (in particular, on a one-dimensional curve in \mathbb{C}), then it is identically zero on \mathbb{C} , hence $G(z) \equiv 0$ for all $z \in \mathscr{Z}$. In our case, the curve on which G vanishes is \mathscr{L} (and we assume, of course, that it is a nondegenerate curve for all the relevant values of the parameter Θ). Hence, G vanishes on the entire \mathscr{L} , and therefore

$$G(x, y) = 0$$
 whenever $(x, y) \in \mathcal{Z}$. (12)

On the other hand, if the system of Eq. (9) has a complex solution (x, y), then (12) would be impossible, since any solution of (9) lies on the manifold \mathscr{L} (because P(x,y)=0), and at the same time Q(x,y)=0 implies G(x,y)=1. Therefore, if system (9) has a solution (real or complex), then representation (8) cannot possibly exist.

It remains to show that if system (9) has no solutions, then the representation (8) is possible, and hence our problem indeed admits a reduction of complexity. Assuming that (9) has no solutions, we will construct the representation (8) in the simplest, polynomial form

$$w(x, y; \Theta) = \sum_{p,q} w_{p,q}(\Theta) x^p y^q$$
(13)

the degree of this polynomial being independent of the parameter Θ . Consider a polynomial equation

$$P(x,y)U(x,y) + Q(x,y) W(x,y) = 1,$$
(14)

here U(x,y) and W(x,y) are unknown polynomials. A classical mathematical theorem, Hilbert's Nullstellensatz [15], says that Eq. (14) has polynomial solutions U(x,y) and W(x,y) if and only if P(x,y) and Q(x,y) have no common zeroes in \mathbb{C}^2 , i.e. whenever system (9) has no complex solutions, which is exactly what we have assumed. Note that since P and Q depend on Θ , then so do U and W, but we suppress this dependence in Eq. (14).

Now the polynomial W(x, y) solving (14) gives us the weight function $w(x, y; \Theta) = W(x, y)$, and it is easy to see that

$$W(x, y) = 1/Q(x, y)$$
 whenever $P(x, y) = 0$

Technically, the theorem is proved, but we make a further practical remark. Suppose we know that system (9) has no solutions, so that the problem admits a reduction of complexity. In this case, we need to find the polynomial W(x, y) solving (14) in an explicit form, in order to determine the weight function $w(x, y; \Theta)$. To this end, we describe a finite and relatively simple algorithm for computing the coefficients w_{pq} of the polynomial W. We substitute the expansions

$$W(x,y) = \sum_{p,q} w_{p,q} x_i^p y_i^q$$
 and $U(x,y) = \sum_{p,q} u_{p,q} x_i^p y_i^q$

into identity (14) and then equate the terms on the left-hand side and those on the right-hand side with the same degrees of the variables x, y. This gives a linear system of equations for the unknown coefficients w_{pq} and u_{pq} . This might be a large system

(its size depends on the degrees of U and W), but it is a linear system whose solution can be found by routine matrix methods. If the assumed degrees of U and V are high enough, then the above system is always solvable by the so called *effective Nullstellensatz*, see [12]. By solving that system we can obtain explicit formulas for the coefficients w_{pq} and u_{pq} . In fact, we only need the coefficients of W, not U, and those coefficients will be rational functions of the coefficients of the polynomial P(x,y), hence they will be easily computable. Lastly, we remark that solving this system, however complex, is *not* a part of the fitting algorithm—the system needs to be solved only once for every type of curves, and then its solution can be incorporated in the algorithm.

Fitting in 3D. Even though we only discussed here the problem of fitting 2D curves to planar data, our analysis easily generalizes to fitting algebraic surfaces $P(x, y, z; \Theta) = 0$, where P is a polynomial in x, y, z, to three-dimensional data points (x_i, y_i, z_i) . The conclusion is similar: the reduction of complexity is possible if and only if the system of equations

$$P(x, y, z) = 0,$$

 $Q(x, y, z) = 0,$
(15)

with

$$Q(x, y, z) := ||\nabla P(x, y, z)||^2 = (\partial P/\partial x)^2 + (\partial P/\partial y)^2 + (\partial P/\partial z)^2$$

has no solutions, real or complex, for any Θ . For example, if we fit quadratic surfaces, then in a properly chosen coordinate frame, $P(x,y,z) = ax^2 + by^2 + cy^2 + d$ (we require that $d \neq 0$ and at least two of the coefficients a,b,c do not vanish). Then (15) takes form

$$ax^{2} + by^{2} + cy^{2} + d = 0,$$

$$a^{2}x^{2} + b^{2}y^{2} + c^{2}y^{2} = 0.$$

This system has no solutions (thus allowing a reduction of complexity) in only two cases: a = b = c (a sphere) and (up to a permutation of variables) a = b, c = 0 (a cylinder).

N. Chernov is partially supported by NSF Grant DMS-0098788 and N. Simányi is partially supported by NSF Grant DMS-0098773.

References

- [1] M. Berman, D. Culpin, The statistical behaviour of some least squares estimators of the centre and radius of a circle, J.R. Statist. Soc. B 48 (1986) 183–196.
- [2] N. Chernov, C. Lesort, Fitting circles and lines by least squares: theory and experiment, preprint, available at http://www.math.uab.edu/cl/cl1.
- [3] N. Chernov, C. Lesort, Statistical efficiency of curve fitting algorithms, Comput. Statist. Data Anal. to appear.
- [4] N.I. Chernov, G.A. Ososkov, Effective algorithms for circle fitting, Comput. Phys. Comm. 33 (1984) 329–333.

- [5] W. Chojnacki, M.J. Brooks, A. van den Hengel, Rationalising the renormalisation method of Kanatani, J. Math. Imaging Vision 14 (2001) 21–38.
- [6] W. Gander, G.H. Golub, R. Strebel, Least squares fitting of circles and ellipses, BIT 34 (1994) 558–578.
- [7] K. Kanatani, Statistical Optimization for Geometric Computation: Theory and Practice, Elsevier Science, 1996.
- [8] K. Kanatani, Cramer-Rao lower bounds for curve fitting, Graph. Models Image Proc. 60 (1998) 93-99.
- [9] Y. Leedan, P. Meer, Heteroscedastic regression in computer vision: problems with bilinear constraint, Internat. J. Comput. Vision 37 (2000) 127–150.
- [10] V. Pratt, Direct least-squares fitting of algebraic surfaces, Comput. Graphics 21 (1987) 145–152.
- [11] P.D. Sampson, Fitting conic sections to very scattered data: an iterative refinement of the Bookstein algorithm, Comput. Graphics Image Proc. 18 (1982) 97–108.
- [12] J.R. Shoenfield, Mathematical Logic, Addison-Wesley, Reading, MA, 1967, p. 100, Example 18 (e).
- [13] G. Taubin, Estimation of planar curves, surfaces and nonplanar space curves defined by implicit equations, with applications to edge and range image segmentation, IEEE Trans. Pattern Anal. Machine Intell. 13 (1991) 1115–1138.
- [14] K. Turner, Computer perception of curved objects using a television camera, Ph.D. Thesis, Department of Machine Intelligence, University of Edinburgh, 1974.
- [15] O. Zariski, P. Samuel, Commutative Algebra, Vol. 2, Van Nostrand, Princeton, NJ, 1958–1960, p. 164.