

Robust Parameter Estimation in Computer Vision*

Charles V. Stewart[†]

Abstract. Estimation techniques in computer vision applications must estimate accurate model parameters despite small-scale noise in the data, occasional large-scale measurement errors (outliers), and measurements from multiple populations in the same data set. Increasingly, robust estimation techniques, some borrowed from the statistics literature and others described in the computer vision literature, have been used in solving these parameter estimation problems. Ideally, these techniques should effectively ignore the outliers and measurements from other populations, treating them as outliers, when estimating the parameters of a single population. Two frequently used techniques are least-median of squares (LMS) [P. J. Rousseeuw, *J. Amer. Statist. Assoc.*, 79 (1984), pp. 871–880] and M-estimators [*Robust Statistics: The Approach Based on Influence Functions*, F. R. Hampel et al., John Wiley, 1986; *Robust Statistics*, P. J. Huber, John Wiley, 1981]. LMS handles large fractions of outliers, up to the theoretical limit of 50% for estimators invariant to affine changes to the data, but has low statistical efficiency. M-estimators have higher statistical efficiency but tolerate much lower percentages of outliers unless properly initialized.

While robust estimators have been used in a variety of computer vision applications, three are considered here. In analysis of range images—images containing depth or X , Y , Z measurements at each pixel instead of intensity measurements—robust estimators have been used successfully to estimate surface model parameters in small image regions. In stereo and motion analysis, they have been used to estimate parameters of what is called the “fundamental matrix,” which characterizes the relative imaging geometry of two cameras imaging the same scene. Recently, robust estimators have been applied to estimating a quadratic image-to-image transformation model necessary to create a composite, “mosaic image” from a series of images of the human retina. In each case, a straightforward application of standard robust estimators is insufficient, and carefully developed extensions are used to solve the problem.

Key words. computer vision, robust statistics, parameter estimation, range image, stereo, motion, fundamental matrix, mosaic construction, retinal imaging

AMS subject classification. 68T10

PII. S0036144598345802

1. Introduction. The goal of computer vision algorithms is to extract geometric, photometric, and semantic information from image data. This may include the position and identity of an object [2, 8, 21, 39, 50, 60], the motion of a camera attached to a car or an autonomous vehicle [1, 6, 19], the geometry of object surfaces [10, 15, 72], or the transformations necessary to build a large composite image (a mosaic) from a series of overlapping images of the same scene [37, 71]. Each of the processes used to extract this information requires some form of parameter estimation to describe

*Received by the editors October 8, 1998; accepted for publication February 10, 1999; published electronically July 27, 1999. This work was supported by National Science Foundation awards IRI-9217195 and IRI-9408700.

<http://www.siam.org/journals/sirev/41-3/34580.html>

[†]Department of Computer Science, Rensselaer Polytechnic Institute, Troy, NY 12180-3590 (stewart@cs.rpi.edu).

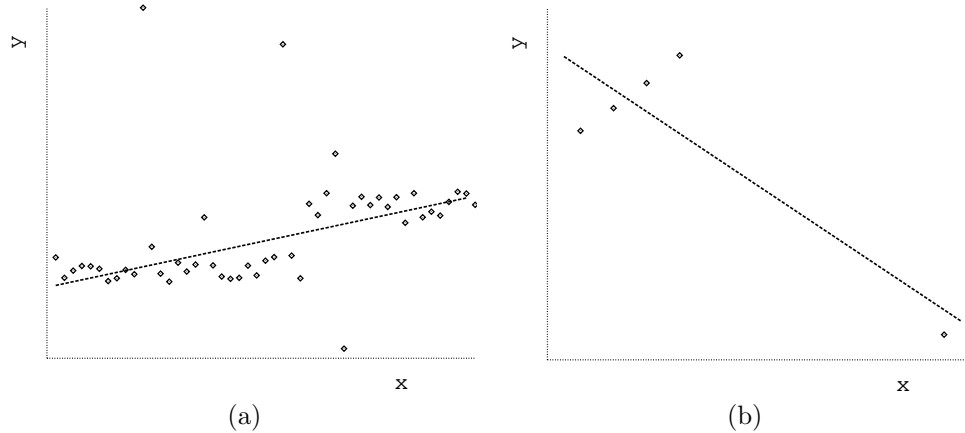


Fig. 1 Simple examples demonstrating the effects of (a) multiple image structures (data populations) plus outliers and (b) a single outlier (a “leverage point”) on linear least-squares fits. If the example in (a) is thought of as the plot of the cross section of a range image (section 3), then the x coordinate represents image position, the y coordinate represents the depth measurement, and the two populations correspond to measurements from two different surfaces, one closer to the sensor and perhaps occluding the other.

intensity edge curves; motion models; surface normals and curvatures; and Euclidean, affine, and projective transformation models [25].

An intensity image typically contains 250,000 pixels or more, with each pixel storing an 8-bit grey level vector or three 8-bit components of a color vector. Range images, where pixels record X, Y, Z scene coordinates as measured by special-purpose sensors, may contain as many measurement vectors. Each pixel measurement is subject to small-scale random variations (noise) caused by the processes of sensing and digitization. This huge volume of data implies that parameter estimation techniques in vision are heavily overconstrained, even for problems where low-level feature extraction, such as edge detection, is used as a preprocessing step. This in turn implies that parameter estimation problems in vision should be solved by least squares or, more generally, maximum likelihood estimation (MLE) techniques.

Unfortunately, computer vision data are rarely drawn from a single statistical population as required for effective use of MLE. Intensity and range images may contain light and depth measurements from multiple scene surfaces. A moving camera, which induces apparent image motion, may also image an independently moving object, inducing a second apparent motion. In both cases, multiple structures (populations) are represented in the image data. Additionally, some image data may be measurements that are difficult to assign to any population. These data are gross errors (“outliers”) which may be caused by specular highlights, saturation of sensors, or mistakes in feature extraction techniques such as edge and corner detectors. It is important to note that these gross errors may be arbitrarily large and therefore cannot be “averaged out,” as is typically done with small-scale noise. As illustrated in Figure 1, when image data are drawn from multiple populations or outliers, application of MLE can produce nonsense results. Much work in computer vision, therefore, has been directed at separating (segmenting) data into distinct populations prior to or during parameter estimation. Complete segmentation, however, is not possible without parameter estimation, because the process of assigning data points to populations depends, at least partially, on the parameters describing the structure of each population. A simple illustration of this is that it is impossible to know which points belong to a line until the line parameters are known.

This difficulty has sparked growing interest among the computer vision community in the use of robust estimation techniques, which have been developed over the last 25 years in both the computer vision and the statistics literatures [26, 28, 35, 36, 53, 63]. These techniques are attractive because they are specifically designed to accommodate data from multiple populations when estimating the parameters describing the dominant population. Ideally, the parameters estimated should not differ substantially from those estimated via MLE for the dominant population in isolation.

This paper provides a tutorial introduction to robust parameter estimation in computer vision. The paper starts with a summary of commonly used robust estimation techniques and then describes how they have been applied and extended. Three applications are considered in detail. The first is estimating the parameters describing low-order surface geometry from range data [7]. The second is fundamental matrix estimation—the problem of using corresponding points to establish the relationship between two different images of a scene taken with uncalibrated cameras [51, 81]. The third problem is building a “mosaic” image of a human retina by matching and combining overlapping retinal images into a single, larger image [5]. The last two problems are closely related.

In thinking about the techniques presented, it is important for the reader to note that robust estimators are not necessarily the only or even the best technique that can be used to solve the problems caused by outliers and multiple populations (structures) in all contexts. Specialized heuristics for handling occasional outliers appear throughout the computer vision literature, and in some cases, most notably when estimating surfaces from raw range or intensity images, multiple populations may be treated as an edge detection or segmentation problem (see textbook discussions in [29], for example). On the other hand, since robust estimation techniques have been designed to handle outliers and multiple populations and since these problems are pervasive, knowledge of robust estimation and its limitations is important in addressing parameter estimation problems in computer vision.

2. Robust Estimation. The first step in describing robust estimators is to state more clearly what is meant by robustness. Several measures of robustness are used in the literature. Most common is the *breakdown point* [63]—the minimum fraction of outlying data that can cause an estimate to diverge arbitrarily far from the true estimate. For example, the breakdown point of least squares is 0 because one bad point can be used to move the least squares fit arbitrarily far from the true fit. The theoretical maximum breakdown point is 0.5 because when more than half the data are outliers they can be arranged so that a fit through them will minimize the estimator objective function.

A second measure of robustness is the *influence function* [28, 35] which, intuitively, is the change in an estimate caused by insertion of outlying data as a function of the distance of the data from the (uncorrupted) estimate. For example, the influence function of the least squares estimator is simply proportional to the distance of the point from the estimate. To achieve robustness, the influence function should tend to 0 with increasing distance.

Finally, although not a measure of robustness, the efficiency of a robust estimator is also significant.¹ This is the ratio of the minimum possible variance in an estimate to the actual variance of a (robust) estimate [48], with the minimum possible variance

¹There is a potential confusion between the statistical notion of efficiency and the computational notion of efficiency, which is associated with the algorithm implementing a robust estimator. Where the meaning is not clear from the context, the phrases “statistical efficiency” and “computation efficiency” will be used.

being determined by a target distribution such as the normal (Gaussian) distribution. Efficiency clearly has an upper bound of 1.0. Asymptotic efficiency is the limit in efficiency as the number of data points tends to infinity. Robust estimators having a high breakdown point tend to have low efficiency, so that the estimates are highly variable and many data points are required to obtain precise estimates.

Robust estimators are usually defined and analyzed in terms of either linear regression or estimation of univariate or multivariate location and scatter [28, 35, 63]. (Most computer vision problems requiring robust estimation are similar to regression problems.) To set a general context, let $X = \{\mathbf{x}_i\}$ be a set of data points (vectors) and let \mathbf{a} be a k -dimensional parameter vector to be estimated. The objective functions used in robust estimation, like those used in MLE, are defined in terms of an error distance or residual function, denoted by $r_{i,\mathbf{a}} = r(\mathbf{x}_i; \mathbf{a})$. Ideally this should be a true geometric distance—e.g., the Euclidean distance between a point \mathbf{x}_i and the curve determined by \mathbf{a} —or better yet, a Mahalanobis distance if the covariance matrix of \mathbf{x}_i is known.

2.1. M-Estimators. While many variations on robust estimation have been proposed in the statistics literature, the two main techniques used in computer vision are M-estimators and least median of squares (LMS). M-estimators are generalizations of MLEs and least squares [28, 35]. In particular, the M-estimate of \mathbf{a} is

$$(1) \quad \hat{\mathbf{a}} = \operatorname{argmin}_{\mathbf{a}} \sum_{\mathbf{x}_i \in X} \rho(r_{i,\mathbf{a}}/\sigma_i),$$

where $\rho(u)$ is a robust loss function that grows subquadratically and is monotonically nondecreasing with increasing $|u|$. Also, σ_i^2 is the variance (scale) associated with the scalar value $r_{i,\mathbf{a}}$. Constraints on $\hat{\mathbf{a}}$ are easily incorporated.

The minimization in (1) is solved by finding \mathbf{a} such that

$$\sum_{\mathbf{x}_i \in X} \psi(r_{i,\mathbf{a}}/\sigma_i) \frac{dr_{i,\mathbf{a}}}{d\mathbf{a}} \frac{1}{\sigma_i} = 0,$$

where $\psi(u) = \rho'(u)$. A common, although certainly not the only, next step [30], [33], [35, pp. 179–192] is to introduce a weight function w , where $w(u) \cdot u = \psi(u)$, and to solve

$$(2) \quad \sum_{\mathbf{x}_i \in X} w(r_{i,\mathbf{a}}/\sigma_i) \frac{1}{\sigma_i^2} \frac{dr_{i,\mathbf{a}}}{d\mathbf{a}} r_{i,\mathbf{a}} = 0.$$

This leads to a process known as “iteratively reweighted least squares” (IRLS), which alternates steps of calculating weights $w_i = w(r_{i,\mathbf{a}}/\sigma_i)$ using the current estimate of \mathbf{a} and solving (2) to estimate a new \mathbf{a} with the weights fixed. Initial estimates of \mathbf{a} may be obtained in a variety of manners, including nonrobust least squares or other robust estimators discussed below.

The many M-estimators that have been proposed differ in the shape of the functions $\rho(\cdot)$ and, as a result, $\psi(\cdot)$ and $w(\cdot)$. Three common functions are listed in Table 1, with weight functions plotted in Figure 2. The ψ functions, essentially, are proportional to the influence function. Hence, ψ functions tending to zero most quickly (“hard redescenders” in the terminology of [33]), such as the Beaton and Tukey function, allow the most aggressive rejection of outliers. This has been found to be important in computer vision problems [69] when outliers have small residual

Table 1 Three different robust loss functions, $\rho(u)$, and associated ψ functions. The “tuning parameters” a , b , and c are often tuned to obtain 95% efficiency [33].

Beaton and Tukey [4]	$\rho(u) = \begin{cases} \frac{a^2}{6} [1 - (1 - (\frac{u}{a})^2)^3], & u \leq a \\ \frac{a^2}{6}, & u > a \end{cases}$	$\psi(u) = \begin{cases} u [1 - (\frac{u}{a})^2]^2, & u \leq a \\ 0, & u > a \end{cases}$
Cauchy [33]	$\rho(u) = \frac{b^2}{2} \log[1 + (\frac{u}{b})^2]$	$\psi(u) = \frac{u}{1 + (u/b)^2}$
Huber [35, Chap. 7]	$\rho(u) = \begin{cases} \frac{1}{2}u^2, & u \leq c \\ \frac{1}{2}c(2 u - c), & c < u \end{cases}$	$\psi(u) = \begin{cases} u, & u \leq c \\ c \frac{u}{ u }, & u > c \end{cases}$

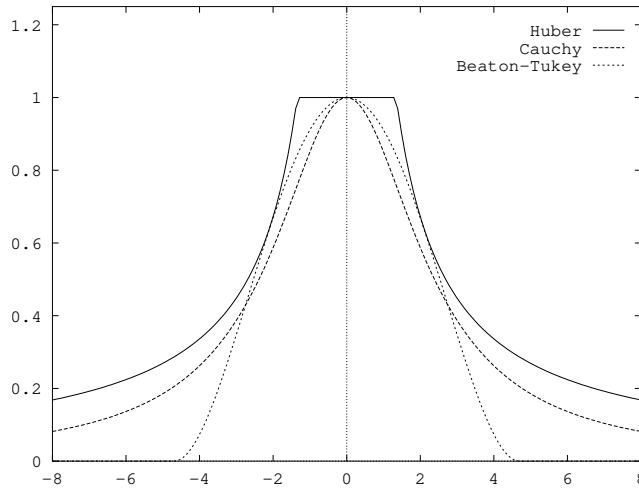


Fig. 2 Plots of the weight functions for the robust loss functions given in Table 1. The horizontal axis units are scale-normalized residuals $r_{i,\mathbf{a}}/\sigma_i$.

magnitudes such as in the range 4σ to 10σ . Unfortunately, when redescending ψ functions are used, the objective function $\sum_{\mathbf{x}_i \in X} \rho(r_{i,\mathbf{a}}/\sigma_i)$ is nonconvex, implying that IRLS will converge at local minima.

The scale values, σ_i , are used to normalize the error distances and to further weight the contribution of each point when scale varies with i . Scale values may be provided a priori from analysis of the sensor or process (such as edge detection) from which the data are obtained. In some cases, scale will not be known in advance and must be estimated from the data. In this case, σ_i is replaced by $\hat{\sigma}$ in the above equations. Equation (1) may be rewritten to jointly estimate \mathbf{a} and σ [35, Chap. 7]. Alternatively, σ may be estimated from an initial fit prior to IRLS and reestimated after each of the first few IRLS iterations. Robust scale estimation techniques such as the median absolute deviation [33] (4) are used. Scale must be fixed, however, before allowing IRLS to converge. The particularly difficult case of σ_i being both unknown and varying with i (heteroscedastic) will not be treated here.

An example of line parameter estimation using an IRLS implementation and the Beaton and Tukey [4] weight function is shown in Figure 3(a). Scale was reestimated after each of the first 3 iterations, and IRLS took 10 iterations to converge. Conver-

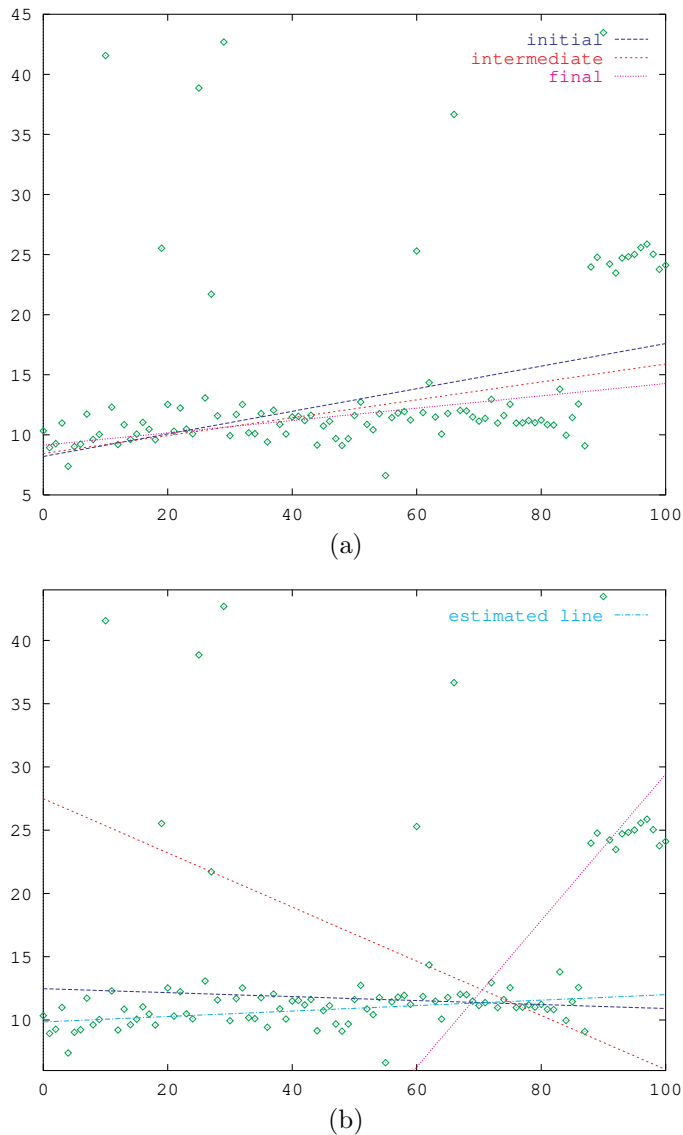


Fig. 3 Example robust line estimates: (a) shows initial, intermediate, and final estimated lines for an M -estimator using IRLS; (b) shows a number of lines hypothesized and tested during a random sampling implementation of LMS.

gence is typically faster with fewer iterations of scale estimation, but the result on this example is less accurate.

The breakdown point of standard M -estimators in regression is 0 because of the possibility of leverage points (Figure 1(b))—outliers positioned far from the remainder of the data in the independent variables as well as far from the fit to the uncorrupted data. A local minimum of (1) will generally occur for a fit rotated to pass through (or near) the leverage point(s), as in Figure 1(b). A generalization of M -estimates, GM-estimators (see discussion in [62]), which downgrade the influence of points based on both independent and dependent variables, can have breakdown points as high as $1/(k+1)$, where k is the length of \mathbf{a} .

2.2. Least-Median of Squares. Least-median of squares (LMS) is distinguished by having a breakdown point of 0.5, the highest value possible. In the notation here, the LMS estimate [61] is

$$(3) \quad \hat{\mathbf{a}} = \underset{\mathbf{a}}{\operatorname{argmin}} \operatorname{median}_{\mathbf{x}_i \in X} r_{i,\mathbf{a}}^2.$$

The intuition behind LMS is that up to half the data points can be arbitrarily far from the optimum estimate without changing the objective function value.

Since the median function is not differentiable, alternative search techniques are required. In the special case of fitting regression lines to points in \mathfrak{R}^2 , computationally efficient algorithms are known for solving (3) exactly [24, 67]. In more general settings, a random sampling technique, developed independently in the computer vision [26] and statistics [61] literatures, is required. The idea is to randomly select a number of k -point subsets of the N data points. A parameter vector, \mathbf{a}_s , is fit to the points in each subset, s [59]. Each \mathbf{a}_s is tested as a hypothesized fit by calculating the squared residual distance r_{i,\mathbf{a}_s}^2 of each of the $N - k$ points in $X - s$ and finding the median. The \mathbf{a}_s corresponding to the smallest median over S subsets is chosen as the estimate, $\hat{\mathbf{a}}$. Overall, this computation requires $O(SN)$ time using linear-time median finding techniques [22].

A crucial parameter in this algorithm is S , the number of subsets. S must be large enough to have a high probability of including at least one subset containing all “good” data points—points that are measurements from the structure of interest. If p is the minimum fraction of good points, then to a first approximation, the probability of a subset containing all good points is p^k . From this, it is easy to see that the probability that at least one of the S subsets contains all good points is

$$P_g = 1 - (1 - p^k)^S.$$

By choosing a desired value of P_g , the minimum value of S can be found. As an example with $P_g = 0.99$, if $k = 3$ and $p = 0.5$, then $S = 35$; if $k = 6$ and $p = 0.6$, then $S = 97$; and if $k = 6$ and $p = 0.5$, then $S = 293$. Clearly, the required number of subsets increases dramatically with increasing k and decreasing p . An example of line parameter estimation using a random sampling implementation of LMS is shown in Figure 3(b). Several hypothesized and tested lines are shown, including the final estimate.

The “median absolute deviation” (MAD) scale estimate may be obtained from the estimate $\hat{\mathbf{a}}$ and used to gather data points to refine the fit. The scale estimate is

$$(4) \quad \hat{\sigma} = 1/\Phi^{-1}(0.75) \left(1 + \frac{5}{N - k} \right) \sqrt{\operatorname{median}_{\mathbf{x}_i \in X - s^*} r_{i,\hat{\mathbf{a}}}^2},$$

where s^* is the subset used to form $\hat{\mathbf{a}}$, Φ^{-1} is the inverse of the cumulative normal distribution, and $1 + 5/(N - k)$ is a finite sample correction factor [63]. The result is an unbiased estimate of $\hat{\sigma}$, which means that if all N points are sampled from a normal distribution with variance σ^2 , then $\hat{\sigma} \rightarrow \sigma$ as $N \rightarrow \infty$. The estimated scale may be used to refine $\hat{\mathbf{a}}$: the data points such that $r_{i,\hat{\mathbf{a}}}^2 < (\theta\hat{\sigma})^2$, for constant θ typically around 2.5, can be identified and used to calculate a least squares estimate of \mathbf{a} . This is important because LMS has low statistical efficiency.

LMS may be generalized in a number of ways based on the order statistics of the residual distances [62]. Some of these are discussed later.

2.3. Requirements for Robust Estimation. The nature of computer vision problems alters the performance requirements of robust estimators in a number of ways:

- The optimum breakdown point of 0.5 must be surpassed in some domains. While this is not possible in general, due to the definition of breakdown, it is possible in particular instances due to problem-specific considerations. For example, fewer than half the points are from any one line in Figure 1(a), but each line is clearly visible. A robust estimator should not fail to estimate a line that approximates one of these two.
- Having enough inliers to satisfy the breakdown limit will not guarantee satisfactory results since it only guarantees that the outliers will not cause the estimate to stray arbitrarily far. For example, the estimated line shown in Figure 1(a), which is approximately the LMS estimate when just more than 50% of the points are from the lower line [69], does not represent a breakdown of the estimator, but it is clearly wrong.
- The low breakdown point of M-estimators can be misleading. A local optimum fit passing near leverage points (Figure 1(b)) can be avoided by proper (robust) initialization of IRLS iterations. (This issue is explored in [47].) Use of an M-estimator having an influence function that tends to zero quickly, so that the weight of the leverage points is zero, will then cause outliers to be ignored.
- The emphasis is often, although not always, on tolerating large numbers of outliers rather than on statistical efficiency. One effect of this is on the choice of tuning parameters for M-estimators—they are often set much lower [9, 69] than suggested by efficiency experiments [33]. This tends to “narrow” the weight functions shown in Figure 2, reducing the influence of points on the tail of the inlier (noise) distribution (hence the loss of efficiency) but making the estimator less sensitive to outliers near the tail of the distribution.
- In some cases, computational efficiency can be of utmost concern. This could render impractical random sampling algorithms generating large numbers of samples.

2.4. Robust Estimation Techniques Developed in Computer Vision. Several robust estimation techniques have been developed in computer vision, either independent of the statistics literature or as an extension of techniques appearing there.² The two most important robust techniques developed independently in computer vision are Hough transforms [36, 43] and RANSAC (random sample consensus) [26, 59]. A Hough transform is a voting technique. The domain (the “parameter space”) of the parameter vector \mathbf{a} to be estimated is discretized, one voting bin is allocated per discrete parameter vector, and each data point \mathbf{x}_i “votes” for parameter vectors \mathbf{a} (i.e., the associated voting bin is incremented) for which the fit residual $r_{\mathbf{x}_i; \mathbf{a}}$ is small enough. The parameter space is searched after voting is complete to locate maxima or to locate and analyze clusters of large numbers of votes. In the former case, the objective function underlying the Hough transform is quite similar to that of an M-estimator and shares some of its limitations [69]. An advantage of Hough transforms in general is that a thorough sampling of the parameter space is obtained. A disadvantage is that the size of the voting space is exponential in the number of parameters, rendering it impractical for many applications.

²In addition, many published algorithms include heuristics for eliminating outliers; these are too numerous and too application-specific to be reviewed here.

RANSAC [26, 59] has similarities to both M-estimators and LMS. Like LMS, it is a minimal subset random sampling search technique, predating LMS by three years. The objective function to be maximized, however, is the number of data points (inliers) having absolute residuals smaller than a predefined value, θ . Equivalently, this may be viewed as minimizing the number of outliers, which may then be viewed as a binary robust loss function that is 0 for small (absolute) residuals, 1 for large absolute residuals, and has a discontinuous transition at θ . Interestingly, both RANSAC and Hough transforms, by virtue of the prespecified inlier band, can be used to find structures formed by substantially fewer than half the data. A cost of this is that small, random structures can also be found [69], implying that careful postprocessing analysis of structures estimated using RANSAC or Hough transforms is required.

Several extensions to LMS have been introduced in the computer vision literature [44, 54, 55]. The simplest is to alter the fraction of inliers used, moving up or down from the 0.5 of the median. In general, reducing the fraction will reduce the breakdown point, since a smaller fraction of outliers is required to corrupt the estimate, but this is advisable when there are multiple structures in the data, when there are relatively few true outliers, and when the minimum fraction of data corresponding to a single structure is known.

A more general approach, which does not rely on a fixed inlier fraction, is based on realizing that the scale estimate in (4) can be generalized to any inlier fraction p by replacing $\Phi^{-1}(0.75)$ with $\Phi^{-1}(0.5+p/2)$. This gives an unbiased scale estimate for the residuals to the correct fit when these residuals follow a normal distribution. Suppose, however, that the data consist of k points measured from a normal distribution and $N - k$ outliers. Here, at the correct fit, the scale estimate will be an overestimate for $p < k/N$, but it will be approximately constant. Significantly, the scale estimate will increase substantially for $p > k/N$. Also, the scale estimates will increase, especially for $p < k/N$, as the fit is moved away from the correct one. These observations have been built into two estimators, ALKS (adaptive least k th squares) [44] and MUSE (minimum unbiased scale estimate) [54, 55], which are based on finding the estimate and associated fraction p that minimize a measure tied to this scale estimate.

The final extension of LMS, an estimator referred to as MINPRAN (minimize probability of randomness), is based on an assumed outlier distribution, which may be known or hypothesized when properties of the sensing process are known [68]. The estimated parameter vector is the one whose chosen inlier set is *least likely* to be a set of outliers over all possible inlier sets. This inlier set is determined from an inlier bound, θ , as in RANSAC, but θ is varied over all possible values to minimize the probability of randomness. Somewhat surprisingly, the computational cost of this is only $O(SN \log N + N^2)$.

3. Surface Estimation from Range Data. The first application of robust estimation in computer vision to be considered here is the problem of estimating the parameters of surfaces from range data. This problem is closely related to regression since errors in range measurements tend to be concentrated along the optical axis of the cameras. This allows the depth measurement to be treated as the only independent variable.

3.1. Range Data and the Estimation Problem. Range data are sets of measurement vectors $(X, Y, Z)^T$ from surfaces in the field of view of a range sensor [10, 58]. Each data set may contain tens or hundreds of thousands of points. There are several different types of range sensor. Many use triangulation, either from two or more cameras at known positions, or by replacing one camera with a light source that shines

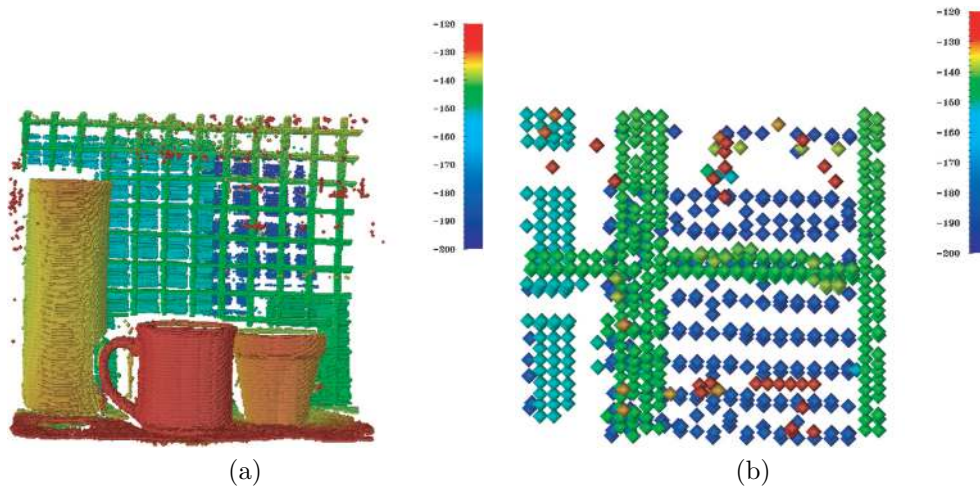


Fig. 4 A complicated range image data set. Image (a) shows the original range image using color to encode depth (red being closer to the viewer). The scene includes a tennis ball can, a coffee mug, and a planting pot in the foreground, a wire mesh table (on its side) behind these, and two different planar surfaces in the background. Image (b) shows a small window to which local surface estimation might be applied; points are from three different surfaces.

known planes of light into a scene, measuring $(X, Y, Z)^T$ from reflected positions in the remaining camera(s). Other techniques determine depth by measuring time of flight or phase shifts in radiation reflected from surfaces. In general, the $(X, Y, Z)^T$ measurements are usually recorded at the grid locations of an image, creating what is called a range image (Figure 4(a)). This range image may be either sparsely or densely populated with points. Accuracy and error of the measured points are around $1/250$ to $1/1000$ of the depth interval over which measurements are taken, which can vary from the order of centimeters to the order of meters. Some of the measurements may be gross errors (outliers) caused by specular highlights, by thresholding effects, or by interreflection of light from multiple surfaces. For more background on range sensors see [40, 58], while for background on the analysis and use of range data in computer vision see [2].

The first goal of range data analysis is to describe the surfaces from which the measurements were taken [15]. The descriptions may be polynomial models such as planar or quadratic models [10, 34, 45, 72], implicit functions such as quadrics and superquadrics [46], or even functions of unknown form, in which case the surfaces are ultimately defined discretely, at each grid location [14, 27], or using a finite-element model [73]. Two complications make the problem particularly difficult. First, each surface will be of limited spatial extent, implying that unknown surface boundaries may need to be located. Second, there will be multiple surfaces, some abutting or even overlapping in image locations and perhaps close in depth measurements. These difficulties suggest what might be termed a “local-to-global” approach, where low-order polynomials first are estimated in small and potentially overlapping image regions and then are gathered into more global, complete surface descriptions. Many range image analysis techniques take this approach [9, 16, 23, 34, 45, 66].

3.2. Local Estimation of Surfaces. The goal of local estimation is to describe the surfaces represented in small image regions using planar or quadratic polynomials, identify and eliminate outliers, and perhaps generate an initial estimate of surface

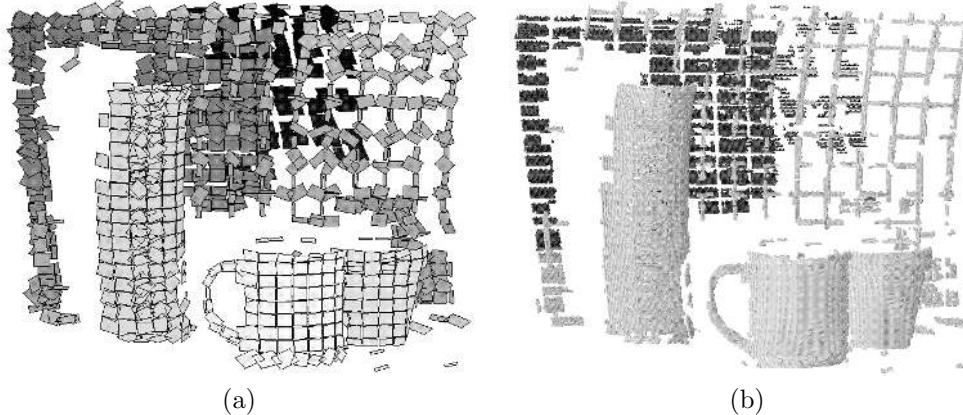


Fig. 5 Results (a slightly sideways view) of local surface estimation on the range image data shown in Figure 4. Planar surface patches with bounding rectangles on their inlier sets are shown in (a), and estimated inliers are shown in (b). Outliers have clearly been eliminated and most surfaces accurately estimated. (Use of the rectangular bounding boxes in (a) introduces some artifacts.) Some surfaces from the wire mesh, unfortunately, were missed because they did not occupy a large enough fraction of the inliers to form an acceptable MUSE fit.

boundaries. This is easy when only one surface is present but more difficult when there are two or more surfaces (Figure 4(b)). In the latter case, it is possible for fewer than 50% of the points to correspond to one surface. Ideally, all surfaces should be described.

The first work in using robust techniques for local estimation was done by Besl, Birch, and Watson [9] using M-estimators to robustly compute planar surface patches, and by Meer et al. [52, 53] using LMS in much the same way. Subsequently, Mirza and Boyer [56] explored a variety of M-estimators for local estimation, concluding that use of those with ψ (and weight) functions redescending to an asymptotic limit of zero at infinity yielded the best performance on data corrupted with random Gaussian outliers. Sinha and Schunck [66] used LMS locally to eliminate outliers and estimate local variations in surface orientation. Stewart [68] introduced the MINPRAN method discussed earlier to estimate multiple, potentially overlapping surface patches, not requiring any patch to contribute more than 50% of the points. Better performance at small magnitude discontinuities was obtained by Miller and Stewart [54, 55] using MUSE and by Lee, Meer, and Park [44] using ALKS (section 2.4). The former incorporates techniques from MINPRAN to determine when the MUSE estimate is distinctive enough to be considered nonrandom; this allows effective determination of the correct number of surface estimates to describe the data set taken from small image regions. An example result using MUSE [54, 55] is shown in Figure 5.

ALKS and MUSE produce accurate surface patches, including multiple surface patches, except when the depth change (discontinuity magnitude) between abutting surfaces is less than about 4.5σ [54], where σ is the standard deviation of the noise in the data. Most standard robust estimation techniques, including LMS and M-estimators, fail when the discontinuity magnitude is 7.5σ or higher [69]. (Figure 1(a) is generated from a model of a discontinuity where these estimators would fail.) The failure in all cases is an estimated surface patch that “bridges” the two real surfaces (see Figure 1(a)). This bridging surface is the global minimum of the objective function and is therefore not an artifact of the random sampling or IRLS search techniques

[69]. In summary, the best robust estimators yield accurate local fits except at extremely small magnitude depth discontinuities, which may be caused by small-scale surface structures.

3.3. Global Estimation of Surfaces. In global estimation of surfaces, which should yield complete descriptions of surfaces and their extents, robust estimation plays a reduced role. It must at least be used in combination with other techniques and is often not used at all in current systems. There are several illuminating reasons for this. First, the appropriate model for the surface is usually unknown, implying that models must be selected automatically [18, 75]. For example, the technique of Boyer, Mirza, and Ganguly [16] grows planar and quadratic surfaces from seed regions—regions of minimum variance in the data—estimated using an M-estimator. Growth is controlled using the same M-estimator, and the model is selected by a modified form of [17]. A related growth technique, which uses prediction intervals to detect finer magnitude discontinuities, is described in [54]. See [34] for a review and empirical comparison of other techniques that address this “segmentation” problem.

A second reason for the reduced role of robust estimation is that even if the surface model is known the extent of the surface must be determined. Without it, robust estimation alone has little chance of success, because a single surface will generally correspond to an extremely small percentage of the data in the entire range image and therefore the potential for “bridging fits,” as discussed above, is even more substantial than in local estimation. This is the reason for growing from seed regions, as in [10, 16, 72].

A third reason is that the surface may be too complicated to be described in a closed form and can only be described at discrete locations or using a spline or finite-element model [15]. Here, while robust estimation can help eliminate outliers [66], the more common use is to convert least squares error norms on the data and on the smoothness of the reconstructed surface into robust error norms based on M-estimator loss functions [13].

3.4. Final Comments on Surface Estimation. Two final comments about the nature of range data and surface estimation are important in order to raise issues that have been incompletely addressed in the literature but are important for practical applications. These affect the use of robust estimation and other techniques. First, most sensors produce quantized data, often with precisions as low as eight bits. At the extreme, any “noise” in the data is buried in the quantization. This differs significantly from the assumptions under which estimators, robust or not, are typically analyzed. Analytical results should therefore be used with caution when predicting estimator performance on this quantized range data. Second, the variance in (nonquantized) range data may vary spatially and with depth, even across a single surface [54, 55]. If this “heteroscedasticity” were known in advance, it could be built into the M-estimator objective function (as written in (1) above), but prior knowledge is difficult to obtain. The effect of such heteroscedasticity on high breakdown point estimators such as LMS has not been studied in computer vision.

4. Estimation of the Fundamental Matrix. Consider two different images of the same scene, taken from two different viewpoints. A great deal of information may be gleaned from these images about the relative positions of the camera(s) when the images were taken, about the structure of the scene, and, if the images were taken at different times, about changes in position of any scene objects. This information is encoded in differences between image positions of scene points. By extracting distinc-

tive features, such as intensity corners, in each image, and then matching them across the two images, a set of correspondences can be established. Each correspondence is a pair of image positions, one from each image, and the points in each pair are hypothesized to correspond to the same scene point. These correspondences are used to estimate image-to-image transformations, camera motion, and scene structure. What is actually estimated depends on what is known in advance about the camera, the motion, and the scene.

This section and the next consider two such problems of estimating image-to-image transformations based on image-to-image correspondences. In each case the correspondence set is generally contaminated by numerous outliers. Robust techniques are used to estimate the transformation parameters and identify the outlier (incorrect) correspondences. The estimated transformation may then be used to return to the initial feature set and identify new correspondences which are much more likely to be correct. This aspect of the problem differs from most applications of robust estimation where the input to the estimation process is fixed.

4.1. Introduction to the Fundamental Matrix. When two images are taken by a single uncalibrated camera or by two different, uncalibrated cameras, the relationship between the images can be captured in the “fundamental matrix” [51, 81]. This matrix imposes a linear constraint on the image positions of corresponding features. Once the fundamental matrix is known, the scene geometry can be reconstructed up to a projective transformation, with additional constraints (e.g., the corner of a room where three orthogonal faces meet) leading to affine or Euclidean reconstructions [25].

Let $\{\tilde{\mathbf{x}}_i\}$ be a set of feature point locations in image 1 and let $\{\tilde{\mathbf{x}}'_i\}$ be a set of corresponding feature locations in image 2, with each location vector written in homogeneous coordinates. The features are intensity corners [82] or other distinctive locations detected by an interest operator [64]. Correspondence for each $\tilde{\mathbf{x}}_i$ is found by searching an area of image 2 established by the range of scene depths and camera motions determined a priori. Within this area, which is sometimes as small as 30 pixels on a side [77] but could be much larger depending on what restrictions are placed on camera motion, $\tilde{\mathbf{x}}'_i$ is the location of the image feature most similar to that of $\tilde{\mathbf{x}}_i$, as decided by, for example, correlation of the surrounding image regions.

The fundamental matrix, \mathbf{F} , is a 3×3 , rank-2 matrix such that, for each i ,

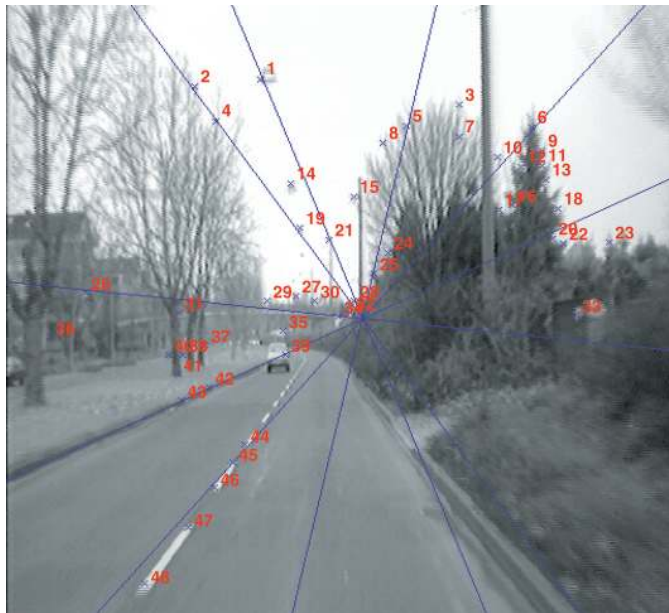
$$(5) \quad \tilde{\mathbf{x}}_i^T \mathbf{F} \tilde{\mathbf{x}}'_i = 0.$$

Because it is homogeneous and rank deficient, \mathbf{F} has seven degrees of freedom and therefore can be estimated from at least seven correspondence pairs. For any point, $\tilde{\mathbf{x}}'$, in image 2, $\mathbf{F}\tilde{\mathbf{x}}'$ defines a line in image 1 through which the point corresponding to $\tilde{\mathbf{x}}'$ must pass. Similarly, for any point, $\tilde{\mathbf{x}}$, in image 1, $\mathbf{F}^T\tilde{\mathbf{x}}$ defines a line in image 2 through which the point corresponding to $\tilde{\mathbf{x}}$ must pass. All such lines, known as “epipolar lines,” pass through the “epipoles,” which are the null spaces of \mathbf{F}^T and \mathbf{F} , respectively. Figure 6 shows a pair of images and epipolar lines. The epipole is the focus of expansion if the intercamera motion is purely translational motion.

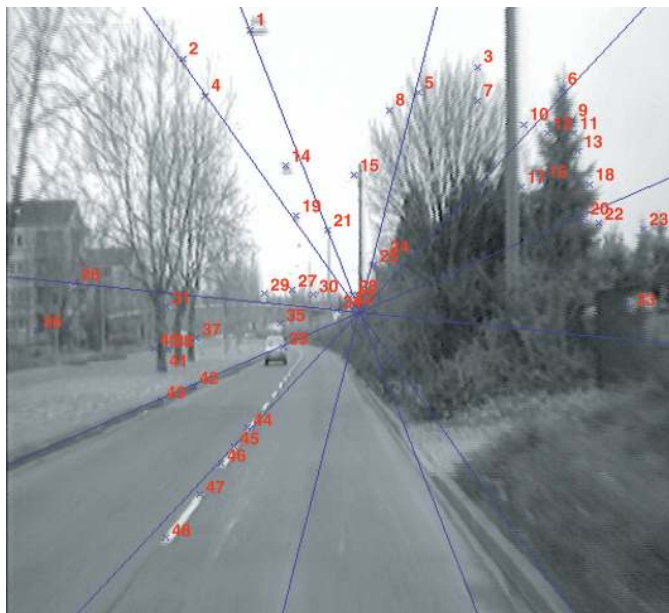
4.2. Estimating the Fundamental Matrix. Estimating the parameters of \mathbf{F} is first considered when there are no outliers and then again in the presence of outliers. The first step is to rewrite (5) $\tilde{\mathbf{x}}_i^T \mathbf{F} \tilde{\mathbf{x}}'_i$ as

$$\mathbf{z}_i^T \mathbf{f},$$

where if $\tilde{\mathbf{x}}_i = (x_i, y_i, 1)^T$ and $\tilde{\mathbf{x}}'_i = (x'_i, y'_i, 1)^T$, then $\mathbf{z}_i^T = (x_i x'_i, x_i y'_i, x_i, y_i x'_i, y_i y'_i, y_i, x'_i, y'_i, 1)$ and \mathbf{f} contains the appropriately ordered parameters of \mathbf{F} . This gives what is



(a)



(b)

Fig. 6 A pair of images and some of the epipolar lines resulting from the robust fundamental matrix technique of [82]. Points numbered are correspondences. Nonrobust estimation results in substantial skewing of these lines and the position of the epipole. Reprinted from *Artificial Intelligence* 78, Z. Zhang, R. Deriche, O. Faugeras, and Q. Luong, *A Robust Technique for Matching Two Uncalibrated Images through the Recovery of the Unknown Epipolar Geometry*, 1995, pp. 87–119, with permission from Elsevier Science.

typically called a fitting error [29, Chap. 11] or an algebraic distance [32]. Summing the squares of these algebraic distances over all i gives

$$\sum_i (\mathbf{z}_i^T \mathbf{f})^2 = \mathbf{f}^T \left(\sum_i \mathbf{z}_i \mathbf{z}_i^T \right) \mathbf{f} = \mathbf{f}^T \mathbf{Z} \mathbf{f}.$$

With appropriate centering and normalization of the data [31], \mathbf{f} , and therefore \mathbf{F} , may be estimated directly from this as the unit eigenvector corresponding to the minimum eigenvalue of \mathbf{Z} . Enforcing $\det(\hat{\mathbf{F}}) = 0$ is achieved by computing the singular value decomposition of $\hat{\mathbf{F}}$, setting the smallest singular value to zero, and recalculating $\hat{\mathbf{F}}$.

Some methods minimize an objective function closer to a geometric distance [77]. Letting $r_i = \mathbf{z}_i^T \mathbf{f}$ be the algebraic residual and $r_{i,x}$, $r_{i,y}$, $r_{i,x'}$, and $r_{i,y'}$ be its partial derivatives, then an approximation of the geometric distance of (x_i, y_i, x'_i, y'_i) to the manifold defined by \mathbf{f} is

$$\frac{r_i}{(r_{i,x}^2 + r_{i,y}^2 + r_{i,x'}^2 + r_{i,y'}^2)^{1/2}}.$$

This geometric distance value is incorporated into estimation of \mathbf{f} by scaling the algebraic residuals by the weights

$$w_i = \frac{1}{(r_{i,x}^2 + r_{i,y}^2 + r_{i,x'}^2 + r_{i,y'}^2)^{1/2}}.$$

This leads to the objective function

$$\sum_i (w_i \mathbf{z}_i^T \mathbf{f})^2,$$

which is minimized iteratively in exactly the same manner as IRLS (section 2.1).

4.3. Robust Estimation. Robust estimation of the fundamental matrix is important because the matching process is unreliable and the presence of an independently moving object in the field of view will induce a distinct fundamental matrix for the image projections of object points. Outlier percentages of 25% to 50% or more are not unrealistic. Application of robust estimation techniques appears straightforward, at first, but several issues arise. Some of these have been addressed in two prominent papers on the topic [77, 82]. These issues are discussed first for M-estimators and then for LMS.

M-estimators require a robust starting point and a weight function that tends to zero quickly. Even with an unrealistically small percentage of outliers (e.g., 10% to 15%), the initial estimate obtained from least squares does not allow M-estimation to converge to the correct estimate [77]. Because of the relatively small image search area for correspondences and because of the noise in estimated match positions, outlier correspondences tend to be close to correct epipolar lines. This implies that the domain of convergence for the correct estimate will be relatively small, in turn making it unlikely that the least squares fit used to initialize IRLS will be within this domain. The proximity of outlier correspondences to epipolar lines shows why the robust weight function must tend to zero quickly. These observations are all practical realizations of the theoretical properties of M-estimators: a zero breakdown point but also the ability to withstand outliers when robustly initialized and used with an appropriate weight function.

The difficulties associated with LMS are slightly different. The first problem is that the correct fundamental matrix is often nearly degenerate. (One example

of such a degeneracy is points from a planar surface.) This has two implications. First, this degeneracy should be avoided in the random samples themselves. Zhang et al. [82, 81] divided the image into nonoverlapping regions (buckets) and prevented a random sample from incorporating more than one match from any one bucket; spreading the random samples across the image reduces the likelihood of a degenerate sample. Torr and Murray [77] developed a rank test for degeneracy and considered degenerate samples no further. Second, LMS, with its low statistical efficiency, will often yield an unreliable estimate of F , even from a correct set of correspondences. Therefore, refinement of the LMS estimate is crucial. Zhang et al. [82] applied a least squares estimate to the matches falling within $2.5\hat{\sigma}$ of the epipolar lines in each image, $\hat{\sigma}$ being taken from the median of the summed squared epipolar distances. Torr and Murray [77] took the LMS estimate as the starting point for an M-estimator using a MAD scale estimate (4) refined using an expectation maximization (EM) algorithm.

The second problem in using LMS concerns the random sampling procedure itself. Zhang et al. [82] generated 8 correspondence samples from which F can be instantiated with the appropriate constraints. Torr and Murray [77] presented a method that instantiates F from 7 correspondence samples. Either way, the number of samples required is quite high, with 382 or 588 required for 7 point samples to have a 0.95 or 0.99 probability of obtaining a good sample and 766 or 1177 required for 8 point samples. This expense could be prohibitive for some applications.

The final issue in using LMS is that the number of outlier correspondences could in fact be greater than 50%. There are two interacting causes of the large number of outliers: (1) intensity corners are relatively weak features and are easily mismatched; (2) large interimage motions can lead to large matching search ranges and therefore more mismatches. (The results shown here are based on relatively small motions.) Torr and Murray [77] suggested the use of RANSAC, which has the concurrent limitation of requiring establishment of a prior inlier bound, when more than 50% of the matches may be outliers. The recent techniques [44, 54, 55] that surpass the 50% breakdown limit in a probabilistic sense without prior knowledge of an inlier bound have not yet been used in fundamental matrix estimation. More generally, the effects of large motions and large numbers of features require further exploration.

Example results of robust fundamental matrix estimation are shown in Figures 6 and 7. Figure 6, from [82], shows a pair of images and the robustly estimated epipolar lines. Figure 7, from [77], shows the motion vectors for matched corner features, and then separates them into those consistent (inliers) and inconsistent (outliers) with the fundamental matrix estimated from them. Notice that most matches from the independently moving person are eliminated as outliers and do not affect the fundamental matrix estimation. In general, however, unlike the local surface estimation problem, the conditions under which independent motions, producing distinct fundamental matrices, can be effectively handled are not yet fully understood [76, 78].

5. Construction of Retinal Image Mosaics. A new application of robust estimation is in the construction of image mosaics from overlapping images of the human retina. Starting from a series of images, the goal is to produce a composite that is much larger than any individual image and that shows the retina as a whole. Example images and a resulting mosaic are shown in Figure 8. This has numerous possible applications in ophthalmology.

The primary issue is calculating the transformation, T , mapping the coordinate system and therefore the pixels of a given image I_m to a reference image I_r which will form the “center” of the mosaic. Most mosaic construction techniques in com-



Fig. 7 *Fundamental matrix and motion estimation from an image sequence of a camera panning to follow a moving person: (a) shows one of these images, (b) shows the motion vectors obtained from correspondences superimposed on this image, (c) shows correspondences consistent with the estimated epipolar geometry (represented by the fundamental matrix), and (d) shows correspondences inconsistent with the estimated epipolar geometry. Reprinted from International Journal of Computer Vision 24, P. Torr and D. Murray, The Development and Comparison of Robust Methods for Estimating the Fundamental Matrix, 1997, pp. 271–300, with permission from Kluwer Academic Publishers.*

puter vision formulate T as an affine transformation in two dimensions [6] or as the model of the apparent (image) motion of a planar surface [1]. These yield 6-parameter and 8-parameter transformation models, respectively. Unfortunately, these models yield substantial mapping errors for the curved surface of the retina. As a result, a 12-parameter transformation model is required, which models the motion of a quadratic surface imaged using weak perspective (scaled orthographic) projection. The derivation of this model is similar to the derivation of the 8-parameter planar surface motion model [1] and is omitted here. To see how the transforma-

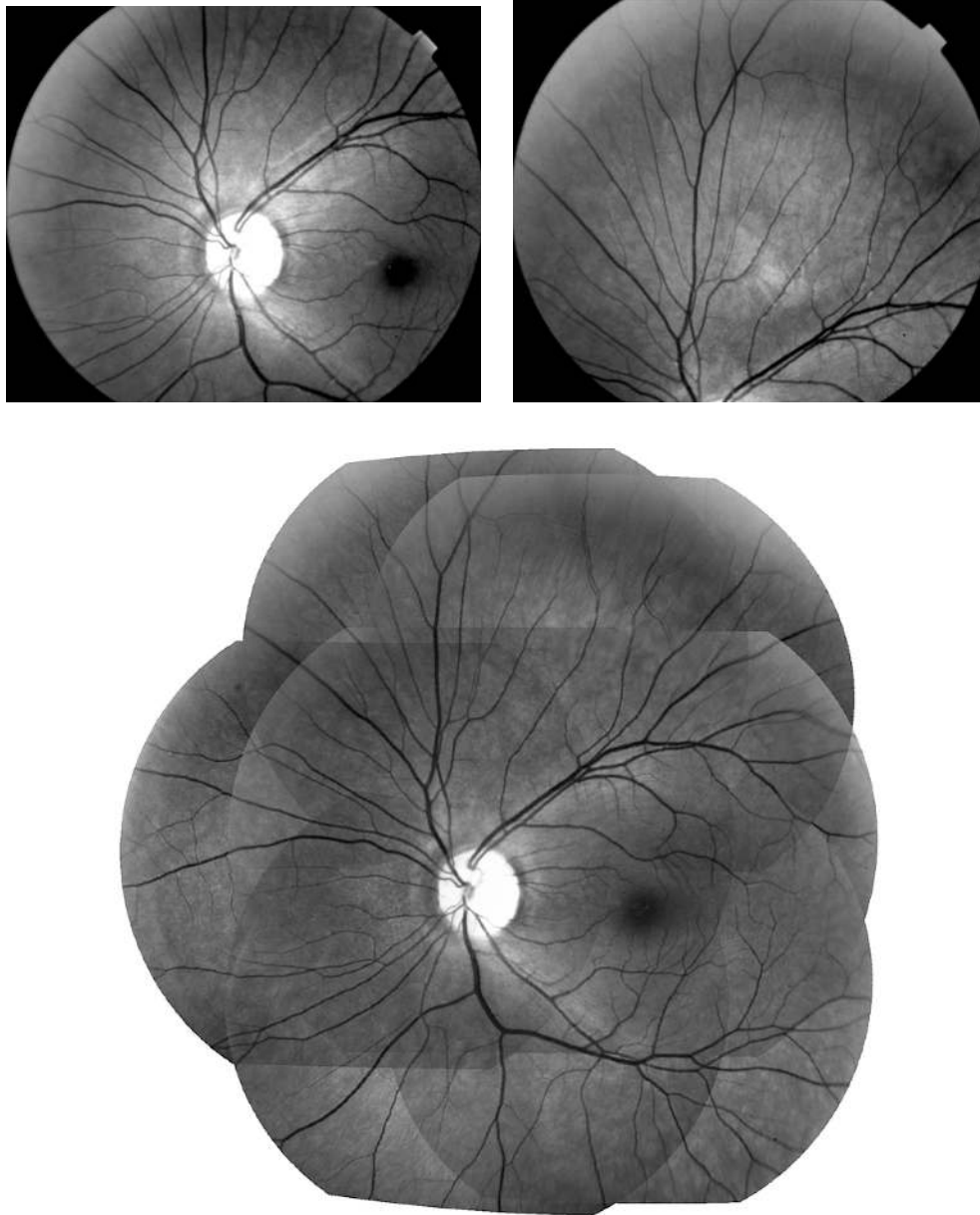


Fig. 8 *Two retinal images and a mosaic of many images estimated using the robust mapping technique described here.*

tion works, let $\mathbf{x}_m = (x_m, y_m)^T$ be the coordinates of a point in I_m , and define $\mathbf{X}(\mathbf{x}_m) = (x_m^2, y_m^2, x_m y_m, x_m, y_m, 1)$. Then the transformed point location in I_r is

$$\mathbf{x}_r = \mathbf{T}\mathbf{X}(\mathbf{x}_m),$$

where \mathbf{T} is a 2×6 matrix. For optimal estimates of \mathbf{T} , this yields transformation errors averaging about one pixel.

Similar to fundamental matrix estimation, \mathbf{T} is estimated by establishing correspondence between image features located in I_m and I_r . Features are bifurcations in

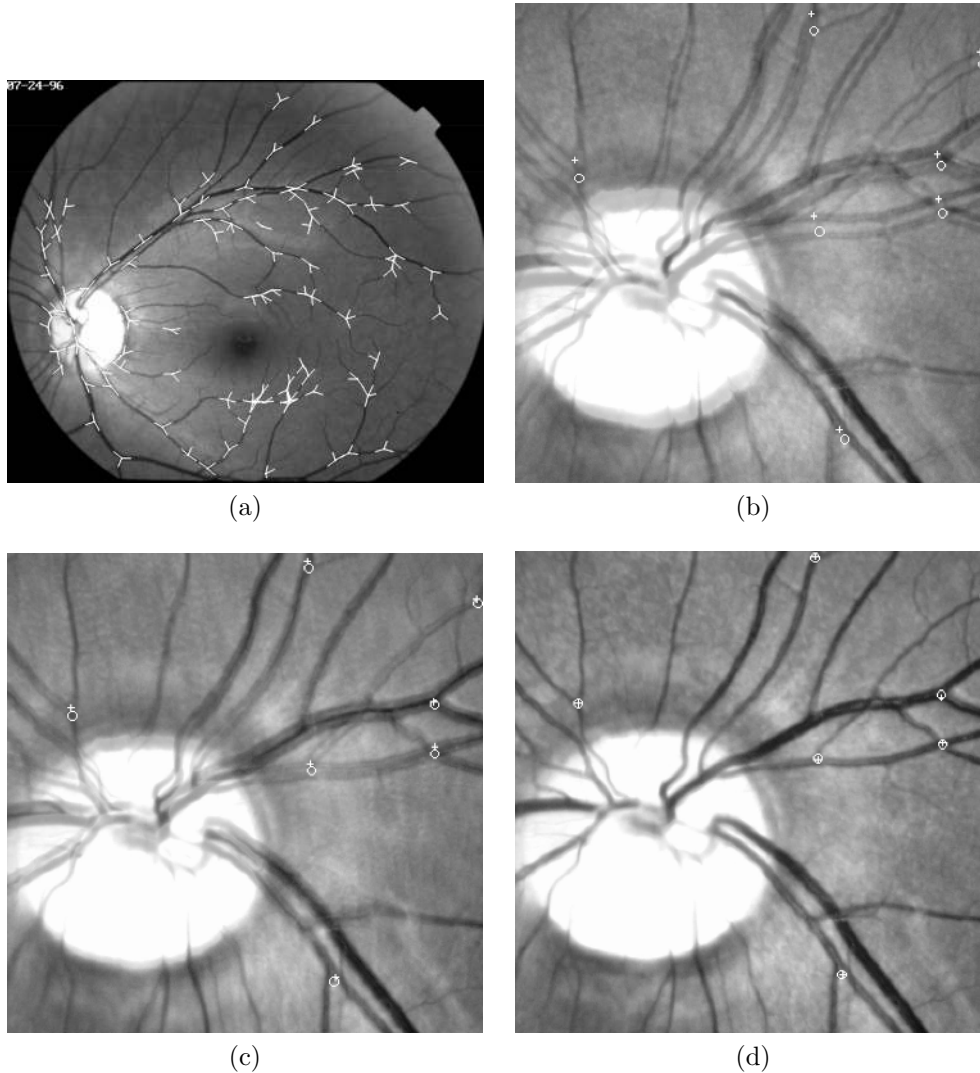


Fig. 9 Hierarchical retinal image transformation estimation results. Image (a) shows features extracted from one image, (b) shows the results of applying the estimated zeroth-order, translation-only model, (c) shows the results of applying the estimated first-order, 6-parameter, affine model, and (d) shows the results of applying the final, estimated second-order, 12-parameter quadratic model. In each of (b), (c), and (d), a small region around the optic disk is shown to emphasize modeling errors.

blood vessels, found through a tracing procedure [20]. Example images are shown in Figure 8 and extracted features are shown in Figure 9(a). Establishing correspondence between features, unfortunately, is substantially more difficult than in fundamental matrix estimation, because the search range is not known in advance and the degree of overlap between I_m and I_r may vary substantially. The result is that initially nearly every possible feature pair must be considered as a correspondence, and this correspondence set must be culled during estimation of T .

5.1. Hierarchical Transform Estimation. The combination of the 12-parameter transformation model and the vast majority of matches being incorrect is too difficult

to handle directly. Instead, low-order approximations to T are estimated to reduce the correspondence set and focus the final robust estimation of T . This hierarchy of transformation models differs from standard hierarchical (multiresolution) techniques in motion estimation, which start from low-resolution versions of images and gradually increase this resolution [6].

The first approximation to T is a simple translation:

$$\mathbf{x}_m = \mathbf{x}_r + \mathbf{t}.$$

Each possible match produces a hypothesized image translation vector, $\mathbf{t} = \mathbf{x}_m - \mathbf{x}_r$, and these vectors are recorded in a two-dimensional histogram with overlapping bins. The radius of each bin equals the experimentally determined maximum position error that could occur using the translation model. For each hypothesized \mathbf{t} , all bins overlapping \mathbf{t} are found. A weight determined by an image similarity measure (e.g., cross correlation) between matching features is added to each overlapping bin. At the end, the translation vector corresponding to the peak bin is taken as the translation estimate and all matches falling into this bin are retained for estimating the next finer approximation to the transformation. For some features, several matches could be retained, while for others, particularly those outside the overlap between images, there may be no matches. Figure 9(b) shows two images aligned using the translation model.

T is next approximated using a six-parameter affine model [6, 38]:

$$\mathbf{x}_m = \mathbf{A}\mathbf{x}_r + \mathbf{t},$$

where \mathbf{A} is a 2×2 matrix with no imposed constraints. Here, LMS is used to robustly estimate the parameters of \mathbf{A} and \mathbf{t} . Features outside the region of image overlap determined by the translation model (plus modeling error) are not considered during LMS. In selecting the matches to form a random sample, image features are first selected randomly and then a correspondence for each selected feature is chosen at random from its correspondences that were retained following translation estimation. Other than this, application of LMS is straightforward. Figure 9(c) shows two images aligned using an estimated affine model.

The final stage is estimation of the full 12-parameter model using an IRLS implementation of an M-estimator. The scale and residuals used in the first iteration of IRLS are taken from the optimal LMS affine transformation. In effect, the affine model estimated by LMS is used to initialize the M-estimate of the full model. This represents a novel twist on robust initialization of an M-estimator, since it is done with the robust estimate of a lower order model. The scale value taken from the affine model is reestimated following the first iteration of IRLS and then fixed for the remaining iterations.

The most interesting aspect of the computation in the final stage is in the weighting. Let $(\mathbf{x}_{m,i}, \mathbf{x}_{r,j})$ be a match, let $r_{i,j} = \|\mathbf{x}_{r,j} - T\mathbf{X}(\mathbf{x}_{m,i})\|$ be the fit residual (the Euclidean distance between the transformed $\mathbf{x}_{m,i}$ and $\mathbf{x}_{r,j}$), let $w_{i,j}$ be the robust weight, and let $s_{i,j}$ be the match similarity measure. The weight calculation requires two steps. First, for each match the robust weights and similarity measures are multiplied to obtain

$$w_{i,j}^* = w_{i,j}s_{i,j}.$$

Second, these weights are scaled based on how they stack up against competing matches. Competing matches are other matches (in I_r) for feature point i (which

is from I_m) or other matches (in I_m) for feature point j (which is from I_r). The scaling factor is

$$\eta_{i,j} = \frac{w_{i,j}^*}{w_{i,j}^* + \sum_{k \neq i} w_{k,j}^* + \sum_{m \neq j} w_{i,m}^*}.$$

This scaling factor will be at or near 1.0 if no other matches for i or j have large weights and will be much less than 1.0 if they do. The point is to downgrade the influence of ambiguous matches. The final weight used in IRLS will be $\eta_{i,j} w_{i,j}^*$. One benefit of this is that it allows matches that were rejected early to reenter the computation. If they produce extreme residuals, their weights will be zero and the computation will proceed as though they didn't exist.

Figure 9(d) shows the final alignment between two images based on the robust estimate of the 12-parameter model. Figure 8 shows a retinal mosaic of many different images. The effects of the nonlinear mapping are most easily seen on the boundaries of the images.

6. Discussion and Conclusions. This paper has summarized several important robust parameter estimation techniques and outlined their use in three applications in computer vision. Many other applications have been considered in the computer vision literature as well, including optical flow and motion estimation [11, 23, 65], edge and feature detection [49, 57, 59], pose estimation [30, 42], and tracking [12]. Most applications of robust estimators, like the three emphasized here, are based on LMS, M-estimators, and their variations.

The observations about robust parameter estimation in the foregoing discussion can be summarized in three important points.

1. The theoretical limit of the 50% breakdown point can and must be surpassed in certain applications. This can be done through the use of RANSAC or Hough transforms if a prior error bound is available, through adaptive techniques based on scale estimates such as ALKS [44] and MUSE [54, 55], or by special-purpose techniques such as the hierarchy of models used in retinal mosaic construction. Care must be taken to ensure that the results are meaningful, especially when the estimated structure includes only a small percentage of the data as inliers. For example, MUSE [54, 55] incorporates a randomness test based on MINPRAN [68] to ensure that the estimated structures are significant.
2. Without a robust initial estimate, as provided by a high breakdown estimator such as LMS, M-estimation is likely to yield poor results. This reflects the low breakdown point of M-estimators. Hence, robust initialization is especially important when outliers are numerous. These outliers cause a nonrobust least squares estimate to be far from correct. Their presence requires that a weight function such as that of Beaton and Tukey [4], which tends to zero quickly, be used. When using such weight functions, lower tuning parameters than recommended in the statistics literature should be used, emphasizing outlier resistance over statistical efficiency. In a practical although not a theoretical sense, robust initialization is less important when few outliers are expected and, especially, when leverage points are nonexistent.
3. The nature of the outliers that might arise in the data should be considered carefully. Most important, of course, is the fact that outliers to one correctly estimated structure (population) are often inliers to (one or more) other structures. In this case, successful use of robust estimation depends

on the distribution of points between structures, the physical proximity of the point sets, and the similarity between the structures themselves. For example, small magnitude depth and orientation discontinuities in range data can lead to skewed surface estimates—“bridging fits”—when using current robust estimators. Similar difficulties are likely in fundamental matrix estimation when images include proximate objects undergoing slightly different motions; the extent of the problem has not yet been fully explored.

The last of these problems is the most difficult. The current best approach to addressing it is to use mixture model formulations [74] in which multiple structures are simultaneously and robustly estimated, and data are dynamically assigned to different structures. This has been studied most heavily in motion and fundamental matrix estimation [3, 41, 78, 79, 80] but should be used increasingly elsewhere. Aside from the added complexity, an important limitation of mixture models is that they are most effective when all structures in the data can be appropriately modeled. This is sometimes difficult or impossible and perhaps shouldn't be necessary when only a single estimate is required; e.g., when estimating a planar surface model for a roof, models for the leaves and branches of trees near the roof should not be required. Unmodeled structures must be treated as outliers to all models, which reduces the mixture model estimation problem to the original robust parameter estimation problem. Further work is clearly needed.

A new approach currently under investigation combines robust parameter estimation with boundary estimation in a single objective function [70]. This makes explicit the fact that structures in computer vision, e.g., object surfaces, have limited spatial extent. The objective function effectively treats points far outside boundaries as outliers, but allows structures to “grow” toward regions of data that are close to the extrapolated estimate. All of this results from the objective function minimization process. While still in the early stages of development, this technique has already shown promise of eliminating some of the problems in local and global surface estimation discussed earlier.

In summary, robust parameter estimation is an important though incompletely solved problem in computer vision. Computer vision data are numerous, corrupted by noise and outliers, and usually drawn from multiple statistical populations (structures). Because robust estimation techniques are designed to handle corrupted and incompletely modeled data, they provide an attractive set of tools for computer vision estimation problems. These tools cannot be applied to vision problems successfully, however, without careful thought about the populations (structures) to be estimated and the frequency and nature of outliers. Successful application can lead to computer vision techniques that accommodate substantial variations in the data. This is an important development in extending the overall reliability of computer vision technology.

Acknowledgments. The author thanks Dr. Kishore Bubna for feedback on an earlier draft of this paper and thanks the anonymous reviewers for their thoughtful comments.

REFERENCES

- [1] G. ADIV, *Determining 3-d motion and structure from optical flow generated by several moving objects*, IEEE Trans. Pattern Anal. Machine Intelligence, 7 (1985), pp. 384–401.
- [2] F. ARMAN AND J. AGGARWAL, *Model-based object recognition in dense range images*, Comput. Surveys, 25 (1993), pp. 5–43.

- [3] S. AYER AND H. SAWHNEY, *Layered representation of motion video using robust maximum likelihood estimation of mixture models and MDL encoding*, in Proceedings, IEEE International Conference on Computer Vision, 1995, pp. 777–784.
- [4] A. E. BEATON AND J. W. TUKEY, *The fitting of power series, meaning polynomials, illustrated on band-spectroscopic data*, *Technometrics*, 16 (1974), pp. 147–185.
- [5] D. E. BECKER, A. CAN, J. N. TURNER, H. L. TANENBAUM, AND B. ROYSAM, *Image processing algorithms for retinal montage synthesis, mapping and real-time location determination*, *IEEE Trans. Biomed. Engrg.*, 45 (1998), pp. 105–118.
- [6] J. BERGEN, P. ANANDAN, K. HANNA, AND R. HINGORANI, *Hierarchical model-based motion estimation*, in Proceedings, Second European Conference on Computer Vision, 1992, Lecture Notes in Comput. Sci. 588, Springer, New York, pp. 237–252.
- [7] P. BESL, *Surfaces in Range Image Understanding*, Springer-Verlag, New York, 1988.
- [8] P. BESL AND R. JAIN, *Three-dimensional object recognition*, *Comput. Surveys*, 17 (1985), pp. 75–145.
- [9] P. J. BESL, J. B. BIRCH, AND L. T. WATSON, *Robust window operators*, in Proceedings, IEEE International Conference on Computer Vision, 1988, pp. 591–600.
- [10] P. J. BESL AND R. C. JAIN, *Segmentation through variable-order surface fitting*, *IEEE Trans. Pattern Anal. Machine Intelligence*, 10 (1988), pp. 167–192.
- [11] M. J. BLACK AND P. ANANDAN, *A framework for the robust estimation of optical flow*, in Proceedings, IEEE International Conference on Computer Vision, 1993, pp. 231–236.
- [12] M. J. BLACK AND A. D. JEPSON, *EigenTracking: Robust matching and tracking of articulated objects using a view-based representation*, *Internat. J. Comput. Vision*, 26 (1998), pp. 63–84.
- [13] M. J. BLACK AND A. RANGARAJAN, *On the unification of line processes, outlier rejection, and robust statistics with applications in early vision*, *Internat. J. Comput. Vision*, 19 (1996), pp. 57–91.
- [14] A. BLAKE AND A. ZISSERMAN, *Visual Reconstruction*, MIT Press, Cambridge, MA, 1987.
- [15] R. BOLLE AND B. VEMURI, *On three-dimensional surface reconstruction methods*, *IEEE Trans. Pattern Anal. Machine Intelligence*, 13 (1991), pp. 1–13.
- [16] K. L. BOYER, M. J. MIRZA, AND G. GANGULY, *The robust sequential estimator: A general approach and its application to surface organization in range data*, *IEEE Trans. Pattern Anal. Machine Intelligence*, 16 (1994), pp. 987–1001.
- [17] H. BOZDOGAN, *Model selection and Akaike's information criterion (AIC): The general theory and its analytic extension*, *Psychometrika*, 52 (1987), pp. 345–370.
- [18] K. BUBNA AND C. V. STEWART, *Model selection and surface merging in reconstruction algorithms*, in Proceedings, IEEE International Conference on Computer Vision, 1998, pp. 895–902.
- [19] P. J. BURT, J. R. BERGEN, R. HINGORANI, R. KOLCZYNSKI, W. A. LEE, A. LEUNG, J. LUBIN, AND H. SHVAYTSER, *Object tracking with a moving camera: An application of dynamic motion analysis*, in Proceedings, IEEE Workshop on Visual Motion, 1989, pp. 2–12.
- [20] A. CAN, H. SHEN, J. N. TURNER, H. L. TANENBAUM, AND B. ROYSAM, *Rapid automated tracing and feature extraction from retinal fundus images using direct exploratory algorithms*, *IEEE Trans. Technology Biomed.*, 6 (1999).
- [21] R. CHIN AND C. DYER, *Model-based recognition in robot vision*, *Comput. Surveys*, 18 (1986), pp. 67–108.
- [22] T. H. CORMEN, C. E. LEISERSON, AND R. L. RIVEST, *Introduction to Algorithms*, McGraw-Hill, New York, 1990.
- [23] T. DARRELL AND A. PENTLAND, *Cooperative robust estimation using layers of support*, *IEEE Trans. Pattern Anal. Machine Intelligence*, 17 (1995), pp. 474–487.
- [24] H. EDELSBRUNNER AND D. L. SOUVAINÉ, *Computing least median of squares regression lines and guided topological sweep*, *J. Amer. Statist. Assoc.*, 85 (1990), pp. 115–119.
- [25] O. FAUGERAS, *Stratification of 3-dimensional vision: Projective, affine, and metric representations*, *J. Opt. Soc. Amer. A*, 12 (1995), pp. 465–484.
- [26] M. A. FISCHLER AND R. C. BOLLES, *Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography*, *Comm. ACM*, 24 (1981), pp. 381–395.
- [27] W. GRIMSON, *From Images to Surfaces: A Computational Study of the Human Early Visual System*, MIT Press, Cambridge, MA, 1981.
- [28] F. R. HAMPEL, P. J. ROUSSEEUW, E. RONCHETTI, AND W. A. STAHEL, *Robust Statistics: The Approach Based on Influence Functions*, John Wiley, New York, 1986.
- [29] R. HARALICK AND L. SHAPIRO, *Computer and Robot Vision*, Vol. 1, Addison-Wesley, Reading, MA, 1992.

- [30] R. M. HARALICK, H. JOO, C.-N. LEE, X. ZHUANG, V. G. BAIDYA, AND M. B. KIM, *Pose estimation from corresponding data*, IEEE Trans. Systems Man Cybernetics, 19 (1989), pp. 1426–1446.
- [31] R. HARTLEY, *In defence of the 8-point algorithm*, in Proceedings, IEEE International Conference on Computer Vision, 1995, pp. 1064–1070.
- [32] R. HARTLEY, *Minimizing algebraic error in geometric estimation problems*, in Proceedings, IEEE International Conference on Computer Vision, 1998, pp. 469–476.
- [33] P. W. HOLLAND AND R. E. WELSCH, *Robust regression using iteratively reweighted least-squares*, Comm. Statist. Theory Methods, A6 (1977), pp. 813–827.
- [34] A. HOOVER, G. JEAN-BAPTISTE, X. JIANG, P. FLYNN, H. BUNKE, D. GOLDFOG, K. BOWYER, D. EGGERT, A. FITZGIBBON, AND R. FISHER, *An experimental comparison of range image segmentation algorithms*, IEEE Trans. Pattern Anal. Machine Intelligence, 18 (1996), pp. 673–689.
- [35] P. J. HUBER, *Robust Statistics*, John Wiley, New York, 1981.
- [36] J. ILLINGWORTH AND J. KITTLER, *A survey of the Hough transform*, Computer Vision, Graphics, and Image Processing (CVGIP), 44 (1988), pp. 87–116.
- [37] M. IRANI, P. ANANDAN, AND S. HSU, *Mosaic based representations of video sequences and their applications*, in Proceedings, IEEE International Conference on Computer Vision, 1995, pp. 605–611.
- [38] M. IRANI, B. ROUSSO, AND S. PELEG, *Computing occluding and transparent motions*, Internat. J. Comput. Vision, 12 (1994), pp. 5–16.
- [39] D. JACOBS, *Matching 3-d models to 2-d images*, Internat. J. Comput. Vision, 21 (1997), pp. 123–153.
- [40] R. JARVIS, *A perspective on range finding techniques for computer vision*, IEEE Trans. Pattern Anal. Machine Intelligence, 5 (1983), pp. 122–139.
- [41] S. JU, M. BLACK, AND A. JEPSON, *Skin and bones: Multi-layer, locally affine, optical flow and regularization with transparency*, in Proceedings, IEEE Conference on Computer Vision and Pattern Recognition, 1996, pp. 307–314.
- [42] R. KUMAR AND A. R. HANSON, *Robust methods for estimating pose and a sensitivity analysis*, CVGIP: Image Understanding, 60 (1994), pp. 313–342.
- [43] V. LEAVERS, *Survey: Which Hough transform?*, CVGIP, 58 (1993), pp. 250–264.
- [44] K.-M. LEE, P. MEER, AND R.-H. PARK, *Robust adaptive segmentation of range images*, IEEE Trans. Pattern Anal. Machine Intelligence, 20 (1998), pp. 200–205.
- [45] A. LEONARDIS, A. GUPTA, AND R. BAJCSY, *Segmentation of range images as the search for geometric parametric models*, Internat. J. Comput. Vision, 14 (1995), pp. 253–277.
- [46] A. LEONARDIS, F. SOLINA, AND A. MACERL, *A direct recovery of superquadric models in range images using recover-and-select paradigm*, in Proceedings, Third European Conference on Computer Vision A, 1994, Springer-Verlag, pp. 309–318.
- [47] S. LI, *Robustizing robust M-estimation using deterministic annealing*, Pattern Recognition, 29 (1996), pp. 159–166.
- [48] B. W. LINDGREN, *Statistical Theory*, Chapman and Hall, London, 1993.
- [49] L. LIU, B. G. SCHUNCK, AND C. C. MEYER, *On robust edge detection*, in Proceedings of the International Workshop on Robust Computer Vision, 1990, pp. 261–286.
- [50] D. LOWE, *Three-dimensional object recognition from single two-dimensional images*, Artificial Intelligence, 31 (1987), pp. 355–395.
- [51] Q. LUONG AND O. FAUGERAS, *The fundamental matrix: Theory, algorithms, and stability analysis*, Internat. J. Comput. Vision, 17 (1996), pp. 43–75.
- [52] P. MEER, D. MINTZ, AND A. ROSENFELD, *Least median of squares based robust analysis of image structure*, in Proceedings of the DARPA Image Understanding Workshop, 1990, pp. 231–254.
- [53] P. MEER, D. MINTZ, A. ROSENFELD, AND D. Y. KIM, *Robust regression methods for computer vision: A review*, Internat. J. Comput. Vision, 6 (1991), pp. 59–70.
- [54] J. V. MILLER, *Regression-Base Surface Reconstruction: Coping with Noise, Outliers, and Discontinuities*, Ph.D. thesis, Rensselaer Polytechnic Institute, Troy, NY, 1997.
- [55] J. V. MILLER AND C. V. STEWART, *MUSE: Robust surface fitting using unbiased scale estimates*, in Proceedings, IEEE Conference on Computer Vision and Pattern Recognition, 1996, pp. 300–306.
- [56] M. J. MIRZA AND K. L. BOYER, *Performance evaluation of a class of M-estimators for surface parameter estimation in noisy range data*, IEEE Trans. Robotics Automat., 9 (1993), pp. 75–85.
- [57] N. NETANYAHU, V. PHILOMIN, A. ROSENFELD, AND A. STROMBERG, *Robust detection of straight and circular road segments in noisy aerial images*, Pattern Recognition, 30 (1997), pp. 1673–1686.

- [58] D. NITZAN, *Three-dimensional vision structure for robot applications*, IEEE Trans. Pattern Anal. Machine Intelligence, 10 (1988), pp. 291–309.
- [59] G. ROTH AND M. D. LEVINE, *Extracting geometric primitives*, CVGIP: Image Understanding, 58 (1993), pp. 1–22.
- [60] C. A. ROTHWELL, *Object Recognition through Invariant Indexing*, Oxford Science Publications, Oxford, UK, 1995.
- [61] P. J. ROUSSEEUW, *Least median of squares regression*, J. Amer. Statist. Assoc., 79 (1984), pp. 871–880.
- [62] P. J. ROUSSEEUW, *Introduction to positive-breakdown methods*, in G. Maddala and C. Rao, eds., Handbook of Statistics, Vol. 15: Robust Inference, Elsevier–North Holland, Amsterdam, 1997, pp. 101–121.
- [63] P. J. ROUSSEEUW AND A. M. LEROY, *Robust Regression and Outlier Detection*, John Wiley, New York, 1987.
- [64] C. SCHMID, R. MOHR, AND C. BAUCKHAGE, *Comparing and evaluating interest points*, in Proceedings, IEEE International Conference on Computer Vision, 1998, pp. 230–235.
- [65] B. G. SCHUNCK, *Image flow: Fundamentals and algorithms*, in W. N. Martin and J. K. Aggarwal, eds., Motion Understanding: Robot and Human Vision, Kluwer Academic Publishers, Norwell, MA, 1988, pp. 23–80.
- [66] S. S. SINHA AND B. G. SCHUNCK, *A two-stage algorithm for discontinuity-preserving surface reconstruction*, IEEE Trans. Pattern Anal. Machine Intelligence, 14 (1992), pp. 36–55.
- [67] D. L. SOUVAINE AND J. M. STEELE, *Time- and space-efficient algorithms for least median of squares regression*, J. Amer. Statist. Assoc., 82 (1987), pp. 794–801.
- [68] C. V. STEWART, *MINPRAN: A new robust estimator for computer vision*, IEEE Trans. Pattern Anal. Machine Intelligence, 17 (1995), pp. 925–938.
- [69] C. V. STEWART, *Bias in robust estimation caused by discontinuities and multiple structures*, IEEE Trans. Pattern Anal. Machine Intelligence, 19 (1997), pp. 818–833.
- [70] C. V. STEWART, K. BUBNA, AND A. PERERA, *Estimating Model Parameters and Boundaries by Minimizing a Joint, Robust Objective Function*, Proceedings, IEEE Conference on Computer Vision and Pattern Recognition, Fort Collins, CO, 1999, to appear.
- [71] R. SZELISKI, *Video mosaics for virtual environments*, IEEE Computer Graphics Appl., 16 (1996), pp. 22–30.
- [72] 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 Intelligence, 13 (1991), pp. 1115–1138.
- [73] D. TERZOPOULOS, *The computation of visible-surface representations*, IEEE Trans. Pattern Anal. Machine Intelligence, 10 (1988), pp. 417–438.
- [74] D. M. TITTERINGTON, A. F. M. SMITH, AND U. E. MAKOV, *Statistical Analysis of Finite Mixture Distributions*, John Wiley, New York, 1985.
- [75] P. TORR, *An assessment of information criteria for motion model selection*, in Proceedings, IEEE Conference on Computer Vision and Pattern Recognition, 1997, pp. 47–52.
- [76] P. TORR AND D. MURRAY, *Statistical detection of independent movement from a moving camera*, Image Vision Comput., 11 (1993), pp. 180–187.
- [77] P. TORR AND D. MURRAY, *The development and comparison of robust methods for estimating the fundamental matrix*, Internat. J. Comput. Vision, 24 (1997), pp. 271–300.
- [78] P. TORR AND A. ZISSERMAN, *Concerning Bayesian motion segmentation, model averaging, matching and the trifocal tensor*, in Proceedings, Fifth European Conference on Computer Vision, 1998, Springer-Verlag, pp. 511–527.
- [79] J. Y. A. WANG AND E. H. ADELSON, *Layered representation for motion analysis*, in Proceedings, IEEE Conference on Computer Vision and Pattern Recognition, 1993, pp. 361–366.
- [80] Y. WEISS, *Smoothness in layers: Motion segmentation using nonparametric mixture estimation*, in Proceedings, IEEE Conference on Computer Vision and Pattern Recognition, 1997, pp. 520–526.
- [81] Z. ZHANG, *Determining the epipolar geometry and its uncertainty: A review*, Internat. J. Comput. Vision, 27 (1998), pp. 161–195.
- [82] Z. ZHANG, R. DERICHE, O. FAUGERAS, AND Q. LUONG, *A robust technique for matching two uncalibrated images through the recovery of the unknown epipolar geometry*, Artificial Intelligence, 78 (1995), pp. 87–119.