

A Practical Approximation Algorithm for the LTS Estimator[☆]

David M. Mount^{a,1,*}, Nathan S. Netanyahu^{b,c}, Christine D. Piatko^d, Angela Y. Wu^e, Ruth Silverman^c

^aDepartment of Computer Science, University of Maryland, College Park, MD, USA

^bDepartment of Computer Science, Bar-Ilan University, Ramat-Gan 52900, Israel

^cCenter for Automation Research, University of Maryland, College Park, MD, USA

^dThe Johns Hopkins University Applied Physics Laboratory, Laurel, MD, USA

^eDepartment of Computer Science, American University, Washington, DC, USA

Abstract

The linear least trimmed squares (LTS) estimator is a statistical technique for fitting a linear model to a set of points. It was proposed by Rousseeuw as a robust alternative to the classical least squares estimator. Given a set of n points in \mathbb{R}^d , the objective is to minimize the sum of the smallest 50% squared residuals (or more generally any given fraction). There exist practical heuristics for computing the linear LTS estimator, but they provide no guarantees on the accuracy of the final result. Two results are presented. First, a measure of the numerical condition of a set of points is introduced. Based on this measure, a probabilistic analysis of the accuracy of the best LTS fit resulting from a set of random elemental fits is presented. This analysis shows that as the condition of the point set improves, the accuracy of the resulting fit also increases. Second, a new approximation algorithm for LTS, called Adaptive-LTS, is described. Given bounds on the minimum and maximum slope coefficients, this algorithm returns an approximation to the optimal LTS fit whose slope coefficients lie within the given bounds. Empirical evidence of this algorithm's efficiency and effectiveness is provided for a variety of data sets.

Keywords: robust estimation, linear estimation, least trimmed squares, approximation algorithms, computational geometry

1. Introduction

Consider an n -element point set $P = \{p_1, \dots, p_n\}$, where $p_i = (x_{i,1}, \dots, x_{i,d-1}, y_i) \in \mathbb{R}^d$. In standard linear regression with intercept, we assume the model

$$y_i = \sum_{j=1}^{d-1} \beta_j x_{i,j} + \beta_d + e_i, \quad \text{for } i = 1, \dots, n,$$

[☆]Dedicated to the memory of our dear friend and longtime colleague, Ruth Silverman

*Corresponding author

Email addresses: mount@cs.umd.edu (David M. Mount), nathan@{cs.biu.ac.il; cfar.umd.edu} (Nathan S. Netanyahu), christine.piatko@jhuapl.edu (Christine D. Piatko), awu@american.edu (Angela Y. Wu)

¹The work of this author has been partially supported by NSF grant CCF-1117259 and ONR grant N00014-08-1-1015.

Report Documentation Page

Form Approved
OMB No. 0704-0188

Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.

1. REPORT DATE 02 JUL 2015		2. REPORT TYPE		3. DATES COVERED 00-00-2015 to 00-00-2015	
4. TITLE AND SUBTITLE A Practical Approximation Algorithm for the LTS Estimator				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) University of Maryland, Department of Computer Science, College Park, MD, 20742				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution unlimited					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT The linear least trimmed squares (LTS) estimator is a statistical technique for fitting a linear model to a set of points. It was proposed by Rousseeuw as a robust alternative to the classical least squares estimator. Given a set of n points in Rd, the objective is to minimize the sum of the smallest 50% squared residuals (or more generally any given fraction). There exist practical heuristics for computing the linear LTS estimator, but they provide no guarantees on the accuracy of the final result. Two results are presented. First, a measure of the numerical condition of a set of points is introduced. Based on this measure, a probabilistic analysis of the accuracy of the best LTS fit resulting from a set of random elemental fits is presented. This analysis shows that as the condition of the point set improves, the accuracy of the resulting fit also increases. Second a new approximation algorithm for LTS, called Adaptive-LTS, is described. Given bounds on the minimum and maximum slope coefficients, this algorithm returns an approximation to the optimal LTS fit whose slope coefficients lie within the given bounds. Empirical evidence of this algorithm's efficiency and effectiveness is provided for a variety of data sets.					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT	18. NUMBER OF PAGES	19a. NAME OF RESPONSIBLE PERSON
a. REPORT unclassified	b. ABSTRACT unclassified	c. THIS PAGE unclassified			

where $(x_{i,1}, \dots, x_{i,d-1})$ are the given independent variables, the y_i 's are the given dependent variables, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_d)$ is the unknown coefficient vector, and the e_i 's are the errors. We refer to $(\beta_1, \dots, \beta_{d-1})$ as the *slope coefficients* and β_d as the *intercept*. Given an estimator $\hat{\boldsymbol{\beta}} \in \mathbb{R}^d$, define the *ith residual* to be $r_i(\hat{\boldsymbol{\beta}}, P) = y_i - (\sum_{j=1}^{d-1} \hat{\beta}_j x_{i,j} + \hat{\beta}_d)$. Let $r_{[i]}(\hat{\boldsymbol{\beta}}, P)$ denote the *ith smallest residual* in terms of absolute value.

Robust estimators (see, e.g., [1]) have been introduced in order to eliminate sensitivity to outliers, that is, points that fail to follow the linear pattern of the majority of the points. The basic measure of the robustness of an estimator is its *breakdown point*, that is, the fraction (up to 50%) of outlying data points that can corrupt the estimator arbitrarily. One of the most widely studied robust linear estimators is Rousseeuw's *least median of squares estimator* (LMS) [2], which is defined to be the estimator that minimizes the median squared residual. More generally, given an integer *coverage value* h , the objective is to find the hyperplane that minimizes the h th smallest squared residual. A number of papers, both practical and theoretical, have been devoted to computing this estimator in the plane and in higher dimensions (see, e.g., [3, 4, 5, 6, 7, 8, 9]).

It has been observed by Rousseeuw and Leroy [1] that LMS may not be the best estimator from the perspective of statistical properties. They argue in support of the *least trimmed squares* (or LTS) linear estimator [2]. Given an n -element point set P and a coverage $h \leq n$, it is defined to be the estimator that minimizes the *sum* (as opposed to the maximum) of the h smallest squared residuals. For the sake of preserving scale, we convert this into a quantity that more closely resembles a standard deviation. More formally, given a non-vertical hyperplane $\hat{\boldsymbol{\beta}}$ define its *LTS cost* with respect to P and h to be

$$\Delta_{\hat{\boldsymbol{\beta}}}(P, h) = \left(\frac{1}{h-1} \sum_{i=1}^h r_{[i]}^2(\hat{\boldsymbol{\beta}}, P) \right)^{1/2}.$$

Note that this measure is scale equivariant, and minimizing it is equivalent to minimizing the sum of squared residuals. The *LTS estimator* is the coefficient vector of minimum LTS cost, which we denote throughout by $\boldsymbol{\beta}^*(P, h)$. Let $\Delta^*(P, h)$ denote the associated LTS cost. When P and h are clear from context, we refer to these simply as $\boldsymbol{\beta}^*$ and Δ^* , respectively. The *LTS problem* is that of computing $\boldsymbol{\beta}^*$ from P and h . We refer to the points having the h smallest squared residuals as *inliers* and the remaining points as *outliers*. This generalizes the ordinary least squares estimator (when $h = n$). It is customary to set $h = \lfloor (n + d + 1)/2 \rfloor$ for outlier detection [1]. In practice, h may be set to some constant fraction of n based on the expected number of outliers. The statistical properties of LTS are analyzed in [1, 2].

The computational complexity of LTS is less well understood than that of LMS. Various exact algorithms have been presented, which are based on variants of branch-and-bound search and efficient incremental updates [10, 11, 12]. Unfortunately, these algorithms are practical only for fairly small point sets. Hössjer [13] presented an $O(n^2 \log n)$ algorithm for LTS in \mathbb{R}^2 based on plane sweep. In a companion paper [14] we presented an exact algorithm for LTS in \mathbb{R}^2 , which runs in $O(n^2)$ time. We also presented an algorithm for \mathbb{R}^d , for any $d \geq 3$, which runs in $O(n^{d+1})$ time. (Throughout, we assume that d is a fixed constant.) For large n these running times may be unacceptably high, even in spaces of moderate dimension.

Given these relatively high running times, it is natural to consider whether this problem can be solved approximately. There are a few possible ways to formulate LTS as an approximation problem, either by approximating the residual, by approximating the quantile, or both. The following formulations were introduced in [14]. The approximation parameters ε_r and ε_q denote the allowed residual and quantile errors, respectively.

Residual Approximation: The requirement of minimizing the sum of squared residuals is relaxed. Given $0 < \varepsilon_r$, an ε_r -residual approximation is any hyperplane β such that

$$\Delta_\beta(P, h) \leq (1 + \varepsilon_r) \Delta^*(P, h).$$

Quantile Approximation: Much of the complexity of LTS arises because of the requirement that *exactly* h points be covered. We can relax this requirement by introducing a parameter $0 < \varepsilon_q < h/n$ and requiring that the fraction of inliers used is smaller by ε_q . Let $h^- = h - \lfloor n\varepsilon_q \rfloor$. An ε_q -quantile approximation is any hyperplane β such that

$$\Delta_\beta(P, h^-) \leq \Delta^*(P, h).$$

Hybrid Approximation: The above approximations can be merged into a single approximation. Given ε_r and ε_q as in the previous two approximations, let h^- be as defined above. An $(\varepsilon_r, \varepsilon_q)$ -hybrid approximation is any hyperplane β such that

$$\Delta_\beta(P, h^-) \leq (1 + \varepsilon_r) \Delta^*(P, h).$$

Note that approximating the LTS cost does not imply that the optimum slope coefficients themselves are well approximated. Computing an approximation to the slope coefficients seems to be difficult. In particular, for some pathological point sets there may be many solutions that have nearly the same LTS costs but very different slope coefficients. Consider, for example, fitting a plane to a set of points uniformly distributed within a sphere. In an earlier paper [14], we presented an approximation algorithm for LTS whose execution time is roughly $O(n^d/h)$. In the same paper, we presented asymptotic lower bounds for computing the LTS and the related LTA (least trimmed absolute value) estimators. These results suggest that substantial improvements to these worst-case complexity bounds, either exact or approximate, are unlikely.

There are, however, simple and practical heuristic algorithms for LTS. One is the Fast-LTS heuristic of Rousseeuw and Van Driessen [15]. This algorithm is based on a combination of random sampling and local improvement. It involves repeatedly fitting a hyperplane to a small random sample of points, which is called an *elemental fit*, and then applying a small number of local improvements, called *concentration steps* or *C-steps*. (The algorithm will be presented in Section 2.) Each C-step transforms an existing fit into one of equal or lower LTS cost. In practice Fast-LTS performs quite well. However, it provides no assurance on the quality of the final fit. A more efficient implementation of this general approach was presented recently by Torti *et al.* [16]. Another example of a local search heuristic is Hawkins' feasible point algorithm [17]. It does not offer any formal performance guarantees either.

While Fast-LTS works well in practice, it does not represent a truly satisfactory algorithmic solution to the LTS problem. It implicitly assumes that the data are sufficiently well conditioned such that, with reasonably high probability, a random elemental fit is close enough to the optimal fit such that subsequent C-steps will converge to a fit that is very nearly optimal. If the data are poorly conditioned, however, Fast-LTS may require a very large number of random trials before it succeeds. This raises the question of what constitutes a well-conditioned point set and whether it is possible to derive an algorithm that provides results of guaranteed accuracy on all input sets.

We present two principal results in this paper. First, in Section 2 we introduce a measure of the numerical condition of a set of points and present a probabilistic analysis of the accuracy of the LTS fit resulting from a series of random elemental fits. Our results show that as the

condition of the point set improves, the probability of obtaining an accurate fit improves as well. Second, in Section 3 we present an algorithm that produces outputs of guaranteed approximate accuracy. Specifically, we present a hybrid approximation algorithm for LTS, called *Adaptive-LTS*. The input to the algorithm consists of P , h , bounds on the minimum and maximum slope coefficients, and approximation parameters ε_r and ε_q . It computes a hyperplane that is an $(\varepsilon_r, \varepsilon_q)$ -hybrid approximation to the optimal hyperplane whose slope coefficients lie within the specified bounds. The most complex aspect of the algorithm involves the solution of a variant of the 1-dimensional LTS problem for a set of input intervals, which we call the *interval LTS problem*. In Section 4 we present an $O(n \log n)$ time solution to this problem. We have implemented the Adaptive-LTS algorithm, and in Section 5, we present empirical studies contrasting these two algorithms.

2. Numerical Condition and Elemental Fits

In this section we consider the question of how the numerical conditioning of a point set influences the accuracy of an LTS fit based on random elemental fits. Consider a set of n points P in \mathbb{R}^d and a coverage value $h \leq n$. We assume that the points are in general position, so that any subset of d points defines a unique $(d - 1)$ -dimensional hyperplane. This can always be achieved through an infinitesimal perturbation. Any hyperplane generated in this manner is called an *elemental fit*. Elemental fits play an important role in Fast-LTS (described below). In this section we introduce a measure on the numerical condition of a point set and establish formal bounds on the accuracy of a random elemental fit in terms of this measure.

We begin with a brief presentation of Rousseeuw and van Driessen’s Fast-LTS algorithm. (Our description is slightly simplified, and we refer the reader to [15] for details.) The algorithm consists of a series of stages. Each stage starts by generating a random subset of P of size d and fitting a hyperplane β' to these points. It then generates a series of local improvements as follows. The h points of P that have the smallest squared residuals with respect to β' are determined. Let β'' be the least squares linear estimator of these points. Clearly, the LTS cost of β'' is not greater than that of β' . We set $\beta' \leftarrow \beta''$ and repeat the process. This is called a *C-step*. C-steps are repeated a small fixed number of times.

The overall algorithm operates by performing a large number of stages, where each stage starts with a new random elemental fit. After all stages have completed, a small number of results with the lowest LTS costs are saved. For each of these, C-steps are repeatedly applied until convergence. (It is easy to show that the result is at a local minimum of the LTS cost function.) The algorithm’s final output is the hyperplane with the lowest overall cost. Assuming constant dimension, each C-step can be performed in $O(n)$ time, and thus the overall running time is $O(kn)$, where k is the number of stages.

It is easy to see why Fast-LTS performs well under nominal circumstances. Suppose that we succeed in sampling d inliers, which, assuming at least h inliers, would be expected to occur with probability roughly $(h/n)^d$. Suppose as well that these points are well conditioned in the sense that they define a unique hyperplane even considering floating-point round-off errors. Then the resulting elemental fit should be sufficiently close to the optimum that the subsequent C-steps will converge to the true optimum. Thus, although there are no guarantees of good performance, if the data are well conditioned, then after a sufficiently large number of random starts, it is reasonable to believe that at least one of the resulting fits will lead to the true optimum.

The objective of this section is explore how we might formalize the notion of “well conditioned” data. We will ignore the effect of C-steps, and focus only on analyzing the performance

of repeated random elemental fits. (An analysis of the effect of C-steps would be a valuable addition to our results, but this seems to be a significantly harder problem.) Our aim is to provide sufficient conditions on the properties of point sets so that with constant probability, a random elemental fit will provide a reasonably good approximation to the optimum LTS hyperplane.

Rousseeuw and Leroy provide an analysis of the number of elemental fits needed to obtain a good fit with a given probability (see Chapter 5 of [1]), but their analysis makes the tacit assumption that the inliers lie on the optimum LTS fit and thus determine a unique hyperplane. Such an argument ignores the potential effects of numerical instabilities. For example, if the inliers are clustered close to a low-dimensional flat,² then the resulting elemental fit will be numerically unstable.

Fix P and h throughout the remainder of this section, and let $\beta^* = \beta^*(P, h)$ denote the optimum LTS estimator. Define \mathbf{x}_i to be the d -element column vector $(x_{i,1}, \dots, x_{i,d-1}, 1)$ and β' to be the row vector $(\beta'_1, \dots, \beta'_d)$. The i th residual is just $y_i - \beta' \cdot \mathbf{x}_i$. Let $\mathcal{E}(P)$ denote the collection of all $\binom{n}{d}$ subsets of cardinality d of the index set $\{1, \dots, n\}$. Given $E = \{i_1, \dots, i_d\} \in \mathcal{E}$, let X_E denote the $d \times d$ matrix whose columns are the coordinate vectors of the corresponding independent variables:

$$X_E = \left[\mathbf{x}_{i_1} \mid \dots \mid \mathbf{x}_{i_d} \right].$$

Let \mathbf{y}_E be the d -element row vector of associated dependent variables $(y_{i_1}, \dots, y_{i_d})$. By our general position assumption, these vectors are linearly independent, and so there is a unique d -element row vector β' such that

$$\mathbf{y}_E = \beta' X_E, \quad \text{that is} \quad \beta' = \mathbf{y}_E X_E^{-1}$$

(see Fig. 1). For our subsequent analysis, we also define $\mathbf{y}_E^* = \beta^* X_E$ to be the vector of y -values resulting from evaluating the optimum LTS hyperplane at the points of the elemental set. Thus, we have $\beta^* = \mathbf{y}_E^* X_E^{-1}$.

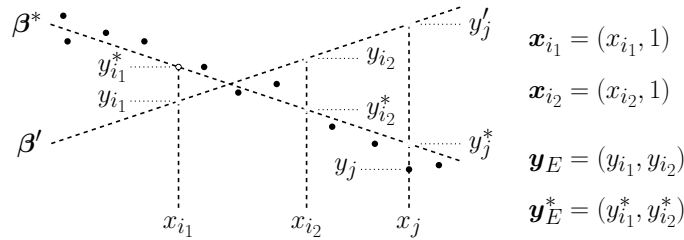


Figure 1: Estimated and optimal fits for $E = \{i_1, i_2\}$ and $\mathbf{y}_E = (y_{i_1}, y_{i_2})$. Solid points indicate points of P and hollow points are computed quantities.

We begin with a few standard definitions [19, 20]. Consider a d -vector \mathbf{x} and a $d \times d$ matrix X . Let $X_{i,j}$ denote the element in row i and column j of X , and let X_j denote the j th column

²This arises, for example, when fitting a plane to a set of points in d -dimensional space that are clustered very close to a k -dimensional flat, where $0 \leq k \leq d - 2$. One way to create such a set in dimension d is to select such an integer k , generate vectors whose first k coordinates are uniform over $[0, 1]$ and whose remaining $d - k$ coordinates are drawn from a Gaussian with a very small standard deviation, and then reorient the flat by applying a random rotation from the orthogonal group $SO(d)$ [18].

vector. The norms of \mathbf{x} and X (in the L_2 and Frobenius norms, respectively) are defined to be

$$\|\mathbf{x}\| = \left(\sum_i x_i^2 \right)^{1/2} \quad \text{and} \quad \|X\| = \left(\sum_j \sum_i X_{i,j}^2 \right)^{1/2} = \left(\sum_j \|X_j\|^2 \right)^{1/2}.$$

Let $\kappa(X)$ denote the *condition number* of X relative to this norm, that is, $\kappa(X) = \|X\| \cdot \|X^{-1}\|$. The concept of condition number is well known in matrix algebra as a means to characterize the stability of linear systems involving X [20]. It is a dimensionless quantity, where lower values mean greater stability.

Our condition measure is based on two properties of the independent variables of the inliers. Let $H \subseteq \{1, \dots, n\}$ denote the h -element index set of the inliers relative to P 's optimal LTS hyperplane. Intuitively, the first property states that a random elemental set of inliers provides a stable elemental fit. This is expressed by bounding the median condition number of X_E for $E \in \mathcal{E}(H)$. That is, we are taking the median over the collection of all d -element subsets of the h inliers. This property alone is not sufficient, since even a minuscule error in the placement of the hyperplane might be fatal if any inlier is located so far away that its associated residual dominates the overall cost. Such a point is called a *leverage point* [21]. Leverage points are problematic if they are few in number (and hence, have a low probability of being sampled) and they do not agree with the linear model determined by the remaining inliers (and hence, almost all elemental fits miss the leverage points). We will adopt a conservative approach of considering all leverage points to be harmful. The second property rules out leverage point inliers by bounding the ratio between the average squared norm of the points \mathbf{x}_i , for $i \in H$ (which is sensitive to leverage points) and the median squared norm of these same points (which is not sensitive). More formally, given two positive real parameters γ and λ , we say that P is (γ, λ) -well conditioned if

$$(i): \quad \text{med}_{E \in \mathcal{E}(H)} \kappa(X_E) \leq \gamma \quad \text{and} \quad (ii): \quad \frac{\sum_{j \in H} \|\mathbf{x}_j\|^2}{h \cdot \text{med}_{j \in H} \|\mathbf{x}_j\|^2} \leq \lambda^2. \quad (1)$$

Note that these assumptions constrain only the independent variables (not the dependent variables), and they apply only to the set of inliers of the optimal fit (and thus there are no assumptions imposed on the outliers). The use of the median in the two assumptions is not critical. At the expense of a constant factor in the analysis, we could have replaced the median with any constant quantile of h .

The main result of this section is given in the following theorem. It shows that depending on the condition of the point set, after a sufficiently large number of random elemental fits, we can achieve a constant-factor (residual) approximation for LTS with arbitrarily high confidence.

Theorem 1. *Let P be an n -element point set in \mathbb{R}^d that is (γ, λ) -well conditioned for some γ and λ . Let h be the coverage, and let $q = h/n$. Then for any $0 < \delta < 1$, with probability at least $1 - \delta$ at least one out of $4(2/q)^d \ln(1/\delta)$ random elemental fits, denoted β' , satisfies*

$$\Delta_{\beta'}(P, h) \leq (1 + \sqrt{2}\gamma\lambda)\Delta^*(P, h).$$

Assuming that the dimension d is a fixed constant, observe the number of random elemental fits grows on the order of $O(q^{-d} \ln(1/\delta))$, which is very sensitive to the dimension and relatively insensitive to the failure probability. The proof relies on the following lemma, which shows that with constant probability, a single elemental fit provides the desired approximation bound.

Lemma 1. *Given the preconditions of Theorem 1, let E be a random element of $\mathcal{E}(P)$, and let β' denote the associated random elemental fit. With probability at least $(q/2)^d/4$,*

$$\Delta_{\beta'}(P, h) \leq (1 + \sqrt{2}\gamma\lambda)\Delta^*(P, h).$$

Assuming for now that this lemma holds, we obtain Theorem 1 as follows. By repeating m independent random elemental fits, the probability of failing to achieve the stated approximation bound is at most $(1 - (q/2)^d/4)^m \leq \exp(-m(q/2)^d/4)$. By selecting $m \geq 4(2/q)^d \ln(1/\delta)$, the probability of failure is at most δ , from which Theorem 1 follows.

The remainder of the section is devoted to proving Lemma 1.

With probability at least q , a random index from $\{1, \dots, n\}$ is drawn from H , and so the probability that a random element of $\mathcal{E}(P)$ lies in $\mathcal{E}(H)$ approaches q^d for all sufficiently large n . (Note that we sample without replacement.) The remainder of the proof is conditioned on this assumption.

Recall that β^* is the optimal LTS hyperplane, and Δ^* is its cost. Let $\Delta' = \Delta_{\beta'}(P, h)$. Let H denote the set of indices of the h points of P whose residuals with respect to β^* are the smallest. We make use of the following technical result, which is proved in Appendix A.

Lemma 2. *Given Δ^* and Δ' as defined above,*

$$\frac{\Delta'}{\Delta^*} - 1 \leq \left(\frac{\|\mathbf{y}_E - \mathbf{y}_E^*\|^2}{\sum_{j \in H} (y_j - y_j^*)^2} \right)^{1/2} \cdot (\|X_E^{-1}\| \cdot \|X_E\|) \cdot \left(\frac{\sum_{j \in H} \|\mathbf{x}_j\|^2}{\|X_E\|^2} \right)^{1/2}.$$

Given this result, let us analyze each of these three factors separately. By definition, the middle factor is just $\kappa(X_E)$. To analyze the first factor, observe that $\|\mathbf{y}_E - \mathbf{y}_E^*\|^2 = \sum_{j \in E} (y_j - y_j^*)^2$. By our earlier assumption, E is a random subset of H , and so with probability approaching $1/2^d$ for large n , every element of $\{(y_j - y_j^*)^2 : j \in E\}$ is at most $\text{med}_{j \in H} (y_j - y_j^*)^2$. Since at least half of the elements of the multiset $\{(y_j - y_j^*)^2 : j \in H\}$ are at least as large as the median of the set, it follows that $\sum_{j \in H} (y_j - y_j^*)^2 \geq (h/2) \text{med}_{j \in H} (y_j - y_j^*)^2$. Thus, with probability approaching $1/2^d$ (and under the assumption that $E \subseteq H$) we have

$$\frac{\|\mathbf{y}_E - \mathbf{y}_E^*\|^2}{\sum_{j \in H} (y_j - y_j^*)^2} \leq \frac{\text{med}_{j \in H} (y_j - y_j^*)^2}{(h/2) \text{med}_{j \in H} (y_j - y_j^*)^2} \leq \frac{2}{h}.$$

To analyze the third factor, recall that from the definition of the Frobenius norm, $\|X_E\|^2 = \sum_{j \in E} \|\mathbf{x}_j\|^2$. By our earlier assumption, E is a random subset of H , and so with probability at least $1 - (1/2)^d \geq 1/2$, at least one element of $\{\|\mathbf{x}_j\|^2 : j \in E\}$ is at least as large as $\text{med}_{j \in H} \|\mathbf{x}_j\|^2$. If so, $\|X_E\|^2$ is at least as large as the median as well. Combining this with property (ii) of well-conditioned sets, with probability at least $1/2$ (and under the assumption that $E \subseteq H$) we have

$$\frac{\sum_{j \in H} \|\mathbf{x}_j\|^2}{\|X_E\|^2} \leq \frac{\lambda^2 h \cdot \text{med}_{j \in H} \|\mathbf{x}_j\|^2}{\text{med}_{j \in H} \|\mathbf{x}_j\|^2} \leq \lambda^2 h.$$

Combining our bounds on all three factors together with condition (i) in the definition of well-conditioned sets, we conclude that for large n , with probability approaching $q^d(1/2^d)(1/2) = (q/2)^d/2$, we have

$$\frac{\Delta'}{\Delta^*} - 1 \leq \left(\frac{2}{h} \right)^{1/2} \cdot \kappa(X_E) \cdot \left(\lambda \sqrt{h} \right) = \sqrt{2} \cdot \kappa(X_E) \cdot \lambda.$$

Since E is a random element of $\mathcal{E}(H)$, with probability at least $1/2$, $\kappa(X_E)$ is not greater than $\text{med}_{F \in \mathcal{E}(H)} \kappa(X_F)$, and so by condition (i) of well-conditioned sets we have $\kappa(X_E) \leq \gamma$. Therefore, with probability at least $(q/2)^d/4$,

$$\frac{\Delta'}{\Delta^*} \leq 1 + \sqrt{2} \cdot \kappa(X_E) \cdot \lambda \leq 1 + \sqrt{2}\gamma\lambda.$$

This completes the proof of Lemma 1 and establishes Theorem 1.

Note that the definition of well conditioning would appear to require the computation of medians over sets of size $O(h^d)$, which is quite large. It is clear from the proof, however, that replacing the medians with the r th quantile in the definition of well conditioning merely changes the factor $1/2^{d+2}$ of Lemma 1 to $1/r^{d+2}$. Thus, rather than computing the median exactly, it would suffice to employ random sampling to estimate its value, at the price of slightly increasing the number of elemental fits.

Of course, Theorem 1 does not completely resolve the question of how the numerical conditioning of the point set affects the convergence of Fast-LTS. First, it ignores the effect of C-steps. Second, since it applies only to the inliers, whose elements we do not normally know, it is not obvious how to turn this into an effective computational procedure. Nonetheless, if *a priori* knowledge about the inlier distribution is available, or if the outliers can be detected and removed *a posteriori*, this result can provide useful bounds on the expected accuracy of any approach to LTS based on repeated random elemental fits in terms of easily computable quantities. In Section 5.3 we will present some experimental analysis of these bounds.

3. An Adaptive Approximation Algorithm

In this section we present an approximation algorithm for LTS, which we call *Adaptive-LTS*. The input to the algorithm is an n -element point set P in \mathbb{R}^d , a coverage value $h \leq n$, a residual approximation parameter $\varepsilon_r \geq 0$, and a quantile approximation parameter ε_q , where $0 \leq \varepsilon_q < h/n$. The algorithm also accepts bounds on the minimum and maximum slope coefficients. If these are not given, the algorithm automatically estimates them from the point set. The algorithm returns a hyperplane that is an $(\varepsilon_r, \varepsilon_q)$ -hybrid approximation to the optimal hyperplane whose slope coefficients lie within the specified slope bounds. Let $h^- = h - \lfloor n\varepsilon_q \rfloor$ denote the *reduced coverage value*, as given in the quantile-approximation definition from Section 1. Recall that we assume that the dimension d is a fixed constant.

This algorithm is based on a general technique called *geometric branch-and-bound*, which involves a recursive subdivision of the solution space in search of the optimal solution. Geometric branch-and-bound has been used in other applications of geometric optimization, notably point pattern matching by Huttenlocher and Rucklidge [22, 23, 24] and Hagedoorn and Velkamp [25] and image registration by Mount *et al.* [26].

Before presenting the algorithm we provide a brief overview of the approach. Recall that the solution to the LTS problem can be represented as a d -dimensional vector $\boldsymbol{\beta} = (\beta_1, \dots, \beta_d)$, consisting of the coefficients of the estimated hyperplane. By fixing P and h , each point $\boldsymbol{\beta} \in \mathbb{R}^d$ is associated with its LTS cost, $\Delta_{\boldsymbol{\beta}} = \Delta_{\boldsymbol{\beta}}(P, h)$. This defines a scalar field over \mathbb{R}^d , called the *solution space*. The LTS problem reduces to computing the point $\boldsymbol{\beta} \in \mathbb{R}^d$ that minimizes this cost, that is, the lowest point in this field. This is the context in which geometric branch-and-bound is applied.

Let us begin with a generic overview of how geometric branch-and-bound works. It recursively subdivides the solution space into a collection of disjoint regions, called *cells*. The initial

cell covers any subset of the solution space that is known to contain the optimum solution. For a given cell C , let Δ_C denote the smallest value of Δ_β , for any $\beta \in C$. For the LTS problem, this is the minimum LTS cost under the constraint that the coefficient vector of the hyperplane is chosen from C . The algorithm then computes a lower bound and an upper bound on Δ_C . Letting Δ_C^- and Δ_C^+ denote these bounds, respectively, we have $\Delta_C^- \leq \Delta_C \leq \Delta_C^+$. (Note that Δ_C^+ is an upper bound on the *minimum* cost within the cell, not an upper bound on all the costs within the cell.) The algorithm also maintains the hyperplane β' associated with the smallest upper bound encountered in any cell that has been seen so far. If C 's lower bound exceeds the smallest upper bound seen so far, that is, if $\Delta_C^- > \Delta_{\beta'}$, we may safely eliminate C from further consideration, since any solution that it might provide cannot improve upon the current best solution. (The coefficient vector β' and $\Delta_{\beta'}$ will be defined below in the section on *Killing cells*.) If so, we say that C is *killed* when this happens. Otherwise, C is said to be *active*.

A geometric branch-and-bound algorithm runs in a sequence of *stages*. At the start of any stage, the algorithm maintains a collection of all the currently active cells. During each stage, one of the active cells is selected and is processed by subdividing it into two or more smaller cells, called its *children*. We compute the aforementioned upper and lower bounds for each child, and if possible, we kill them. Each surviving child cell is then added to the collection of active cells for future processing. The algorithm terminates when no active cells remain.

The precise implementation of any geometric branch-and-bound algorithm depends on the following elements:

- How are cells represented?
- What is the initial cell?
- How is a cell subdivided?
- How are the upper and lower bounds computed?
- Under what conditions is a cell killed?
- Among the active cells, which cell is selected next for processing?

Let us describe each of these elements in greater detail in the context of our approximation algorithm for LTS.

Cells and their representations. Although a solution to LTS is fully specified by the d -vector of coefficients, we will use a more concise representation. Recall that a hyperplane is represented by the equation $y = \sum_{j=1}^{d-1} \beta_j x_j + \beta_d$, where β_d is the intercept term. Rather than store all d -coefficients, we will store only the $d - 1$ slope coefficients $\beta = (\beta_1, \dots, \beta_{d-1}) \in \mathbb{R}^{d-1}$. We will show in Section 4 that given such a vector, it is possible to compute the optimum value of β_d in $O(n \log n)$ time by solving a 1-dimensional LTS problem. We shall show below that the total processing time for a cell is $O(n \log n)$, and so this does not affect the asymptotic time complexity. Because the number of cells that need to be processed tends to grow exponentially with dimension (see Section 5.2), this reduction in dimension can significantly reduce the algorithm's overall running time. Given $\beta \in \mathbb{R}^{d-1}$, let $\Delta_\beta(P, h)$ denote the minimum LTS cost achievable by extending β to a d -dimensional vector.

For our algorithm, a *cell* C is a closed, axis-parallel hyperrectangle in \mathbb{R}^{d-1} . It is represented by a pair of vectors $\beta^-, \beta^+ \in \mathbb{R}^{d-1}$, and consists of the Cartesian product of intervals $\prod_{i=1}^{d-1} [\beta_i^-, \beta_i^+]$. Let $\Delta_C(P, h) = \min_{\beta \in C} \Delta_\beta(P, h)$ denote the smallest LTS cost of any point of C .

Sampled elemental fits. We will employ random elemental fits to help guide the search algorithm. Before the algorithm begins, a user-specified number of randomly sampled elemental fits is generated. Each elemental fit is associated with a $(d - 1)$ -dimensional vector of slope coefficients. Let S_0 denote this set of vectors. For each active cell C , let $S(C)$ denote the elements of S_0 that lie within C . Each active cell C stores the elements of $S(C)$ in a list and the minimum axis-aligned hyperrectangle that contains $S(C)$.

Initial cell. The initial cell can be specified by the user or generated automatically by the algorithm from the sampled elemental fits. In the latter case, the algorithm generates the initial cell by computing an axis-aligned hyperrectangle in \mathbb{R}^{d-1} that encloses some or all of the elements of S_0 . Because any hyperrectangle enclosing S_0 will be heavily influenced by outliers, the algorithm provides an option to compute a subrectangle that contains a user-specified fraction of the points of S_0 .

This subrectangle is computed as follows. Let m denote the number of points of S_0 , and let $m' \leq m$ denote the user-specified number of points to lie within the subrectangle. Define $\rho = (m'/m)^{1/(d-1)}$. For $i = 1, \dots, d - 1$, project the points of S_{i-1} onto the i th coordinate axis and compute the smallest interval that contains a fraction of ρ of these projected points. (This can be done easily in $O(m \log m)$ time by sorting and sweeping through the projected points.) Remove all the points from S_{i-1} that lie outside of this interval and increment i . After doing this for all values of i , return the smallest axis-aligned rectangle that contains the final set S_{d-1} . Since each iteration keeps a fraction ρ the elements of S_{i-1} , it follows that $|S_i|$ is roughly $m\rho^i$. On completion, the number of points that remain is $m\rho^{d-1} = m'$, as desired. (Some additional care is needed when rounding reals to integers to obtain exactly m' points.)

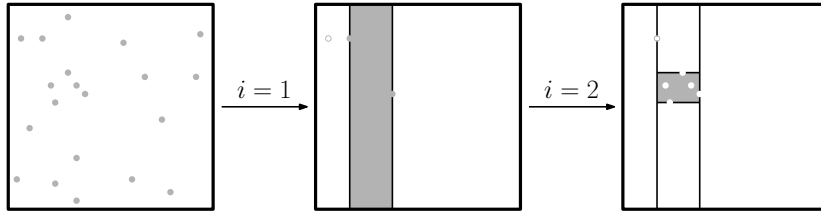


Figure 2: Computing the initial cell when $d = 3$, $m = 20$, and $m' = 5$, and $\rho = (5/20)^{1/2} = 1/2$. Each point represents a randomly sampled elemental fit, and each iteration computes the smallest interval containing a fraction ρ of the remaining points.

In our experiments (see Section 5) we found that the value $m' = cm(h/n)^d$ worked well, where $c \geq 1$ is a small constant. The rationale behind this choice is that, since each elemental fit is based on a sample of d points, and each point has a probability of roughly h/n of being an inlier, a fraction of roughly $(h/n)^d$ of the m elemental fits will be based entirely on inliers. Ideally, the slope coefficients of these fits will be well clustered, and hence for a suitable constant $c \geq 1$, the smallest hyperrectangle containing $cm(h/n)^d$ sampled slopes is a good candidate to include the LTS optimum.

Cell subdivision. We will discuss how active cells are selected for processing below. Once a cell C is selected, it is subdivided into two hyperrectangles by an axis parallel hyperplane. This hyperplane is chosen as follows. First, if $S(C)$ is not empty, then the algorithm determines the longest side length of the bounding rectangle containing $S(C)$ (with ties broken arbitrarily). It

then sorts the vectors of $S(C)$ along the axis that is parallel to this side, and splits the set by a hyperplane that is orthogonal to this side and passes through the median of the sorted sequence (see Fig. 3(a)). If $S(C)$ is empty, the algorithm bisects the cell through its midpoint by a hyperplane that is orthogonal to the cell's longest side length (see Fig. 3(b)). The resulting cells, denoted C_0 and C_1 in the figure, are called C 's children.

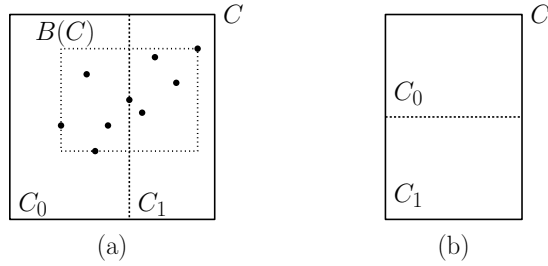


Figure 3: Subdividing a cell (a) when $S(C)$ is nonempty and (b) when $S(C)$ is empty.

An example of the subdivision process is shown below in Fig. 4 for a 2-dimensional solution space. In (a) we show a cell C and its samples $S(C)$. In (b) we show the subdivision that results by repeatedly splitting each cell through its median sample. In (c) we show the effect of two additional steps of subdivision by bisecting the cell's longest side. (This assumes that all the cells are selected for processing and none are killed.) Note that the subdivision has the desirable property that it tends to produce more cells in regions of space where the density of samples is higher.

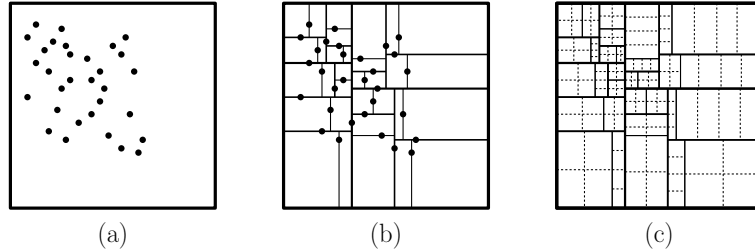


Figure 4: Example of the subdivision process: (a) initial cell and samples, (b) sample-based subdivision, (c) two more stages of midpoint subdivision.

Upper and lower bounds. Recall that the upper bound for a cell C , denoted $\Delta^+(C)$, is an upper bound on the smallest LTS cost among all hyperplanes whose slope coefficients lie within C . To compute such an upper bound, it suffices to sample any representative slope from within the cell and then compute the y -intercept of the hyperplane with this slope that minimizes the LTS cost. We shall show below that it suffices to use the reduced coverage value h^- defined earlier.

Our algorithm chooses the representative slope as follows. If $S(C)$ is nonempty, we take the median sample point that was chosen in the cell-subdivision process. Otherwise, we take the cell's midpoint. As mentioned above, this provides only the slope coefficients. In Section 4, we will show how to compute the optimum intercept value in $O(n \log n)$ time. To further improve

the upper bound, we then apply a user-specified number of C-steps. (We used two C-steps in our experiments.)

The computation of the lower bound for the cell C , denoted $\Delta^-(C)$, is more complex than the upper bound. Our approach is similar to the upper bound computation, but rather than computing the optimum intercept for a specific slope, we compute the optimum intercept under the assumption that the slope could be any point of C (using the full coverage value h , not h^-). In Section 4, we will show how to do this in $O(n \log n)$ time by solving a problem that we call the interval LTS problem.

Killing cells. Recall that our objective is to report a hyperplane β that satisfies

$$\Delta_{\beta}(P, h^-) \leq (1 + \varepsilon_r) \Delta^*(P, h).$$

Let β' denote the coefficient vector associated with the smallest value of $\Delta_{\beta'}(P, h^-)$ encountered so far in the search. (Note that this cost is computed using the reduced coverage value h^- .) We assert that a cell C can be killed if

$$\Delta^-(C) > \frac{\Delta_{\beta'}(P, h^-)}{1 + \varepsilon_r}. \quad (2)$$

To justify this assertion, observe that any coefficient vector β lying within C satisfies

$$(1 + \varepsilon_r) \Delta_{\beta}(P, h) \geq (1 + \varepsilon_r) \Delta^-(C) > \Delta_{\beta'}(P, h^-).$$

Thus, the LTS cost of β is sufficiently high that it cannot contradict the hypothesis that β' is a valid hybrid approximation. Therefore, it is safe to eliminate C from further consideration. (Note that future iterations might change β' , but only in a manner that would decrease $\Delta_{\beta'}(P, h^-)$. Thus, the decision to kill a cell never needs to be reconsidered.)

It is not hard to show that whenever the ratio between a cell's upper bound and lower bound becomes smaller than $(1 + \varepsilon_r)$, the cell will be killed. This is because β' was selected from the processed cell having the smallest value of Δ^+ , and so such a cell satisfies

$$\Delta^-(C) > \frac{\Delta^+(C)}{1 + \varepsilon_r} \geq \frac{\Delta_{\beta'}(P, h^-)}{1 + \varepsilon_r}.$$

This implies that C satisfies Eq. (2) and will be killed. It follows that the algorithm terminates once all active cells have been subdivided to such a small size that this condition holds.

Cell processing order. An important element in the design of an efficient geometric branch-and-bound algorithm is the order in which active cells are selected for processing. When a cell is selected for processing, there are two desirable outcomes we might hope for. First, the representative hyperplane from this cell will have a lower LTS cost than the best hyperplane β' seen so far. This will bring us closer to the optimum LTS cost, and it has the benefit that future cells are more likely to be killed by the condition of Eq. (2). Second, the lower bound of this cell will be so high that the cell will be killed, thus reducing the number of active cells that need to be considered. This raises the question of whether there are easily computable properties of the active cells that are highly correlated with either of these desired outcomes. We considered the following potential criteria for selecting the next cell:

- (1) *Max-samples:* Select the active cell that contains the largest number of randomly sampled elemental fits, that is, the cell C that maximizes $|S(C)|$.

- (2) *Min-lower-bound*: Select the active cell having the smallest lower bound, $\Delta^-(C)$.
- (3) *Min-upper-bound*: Select the active cell having the smallest upper bound, $\Delta^+(C)$.
- (4) *Oldest-cell*: Select the active cell that has been waiting the longest to be processed.

Based on our experience, none of these criteria alone is the best choice throughout the entire search. At the start of the algorithm, when many cells have a large number of samples, criterion (1) tends to work well because it favors cells that are representative of typical elemental fits. But as the algorithm proceeds, and cells are further subdivided, the number of samples per cell decreases. In such cases, this criteria cannot distinguish one cell from another. Criteria (2) and (3) are generally good in these middle phases of the algorithm. But when these two criteria are used exclusively, cells that are far from the LTS optimum are never processed. Adding criterion (4) helps remedy this unbalance.

Our approach to selecting the next active cell employs an adaptive strategy. The algorithm assigns a weight to each of the four criteria. Initially, all the criteria are given the same weight. At the start of each stage, the algorithm selects the next criteria randomly based on these weights. That is, if the current weights are denoted w_1, \dots, w_4 , criterion i is selected with probability $w_i/(w_1 + \dots + w_4)$. We select the next cell using this criterion. After processing this cell, we either reward or penalize this criterion depending on the outcome. First, we compute the ratio between the cell's new lower bound and the best LTS cost seen so far, that is, $\rho^- = \Delta^-(C)/\Delta_{\beta'}(P, h^-)$. Observe that ρ^- is nonnegative and increases as the lower bound approaches the current best LTS cost. With probability $\min(1, \rho^-)$ we reward this choice by increasing the current criterion's weight by 50%, and otherwise we penalize it by decreasing it by 10%. Second, we compute the inverse of the ratio between the cell's new upper bound and the best LTS cost seen so far, that is, $\rho^+ = \Delta_{\beta'}(P, h^-)/\Delta^+(C)$. Again, ρ^+ is nonnegative and increases as the upper bound approaches the best LTS cost. With probability $\min(1, \rho^+)$ we reward this choice by increasing the current criterion's weight by 50%, and otherwise we decrease it by 10%.

The complete algorithm. The algorithm begins by generating the initial sample S_0 of random elemental fits, and the initial cell C_0 is computed as described above and is labeled as active. We then repeat the following stages until no active cells remain. First, the above cell-selection strategy is applied to select and remove the next active cell C . By the subdivision process, we split this cell into two children cells C_0 and C_1 . The set of samples $S(C)$ is then partitioned among these two cells as $S(C_0)$ and $S(C_1)$, respectively. Next, for each of these children, the lower and upper bounds are computed as described earlier. If the upper bound cost is smaller than the current best LTS cost, then we update the current best LTS fit β' and its associated cost. For each child cell, we test whether it can be killed by applying the condition of Eq. (2). If the cell is not killed, it is added to the set of active cells. In either case, the selection criteria that was used to determine this cell is rewarded or penalized, as described above. The algorithm terminates when no active cells remain. On termination, we output the current best fit β' and its associated cost.

By the remarks made earlier, any cell that is killed cannot contradict the hypothesis that β' is an approximation to the optimal solution (subject to the assumption that the optimum lies within the initial slope bounds). The algorithm's correctness follows because on termination, when all the active cells have been killed, β' is a valid $(\varepsilon_r, \varepsilon_q)$ -hybrid approximation.

Unfortunately, it is not easy to provide a good asymptotic bound on the algorithm's overall running time. Unlike Fast-LTS, which runs for a fixed number of iterations irrespective of the structure of the data set, the branch-and-bound algorithm adapts its behavior in accordance with

the structural properties of the point set and the desired accuracy of the final result. In Section 5 we provide evidence that the algorithm terminates very rapidly for well-conditioned inputs, but it can take significantly longer when the data are poorly conditioned.

Even though we cannot bound the total number of stages, we can bound the running time of each stage. In the next section we will show that the upper and lower bounds, $\Delta^+(C)$ and $\Delta^-(C)$, can each be computed in $O(n \log n)$ time, where n is the number of points. Otherwise, the most time consuming part of each stage is the time needed to compute the next active cell. If m denotes the maximum number of active cells at any time, and we simply store them in a list, the running time of each stage is $O(m + n \log n)$. Through the use of more sophisticated data structures,³ it is possible to reduce the running time to $O(\log m + n \log n)$. Since we expect m to be smaller than n , we did not bother to implement this.

4. Computing Upper and Lower Bounds

In this section we describe how to compute the upper and lower bounds for a given cell C in the geometric branch-and-bound algorithm described in the previous section. Recall that P and h denote the point set and coverage, respectively, and that each hyperplane is represented by its slope coefficients $(\beta_1, \dots, \beta_{d-1})$. C is represented by a pair of vectors $\boldsymbol{\beta}^-, \boldsymbol{\beta}^+ \in \mathbb{R}^{d-1}$, and consists of the Cartesian product of intervals $\prod_{i=1}^{d-1} [\beta_i^-, \beta_i^+]$, which we denote by $[\boldsymbol{\beta}^-, \boldsymbol{\beta}^+]$. Let $\Delta_C(P, h) = \min_{\boldsymbol{\beta} \in C} \Delta_{\boldsymbol{\beta}}(P, h)$ denote the smallest LTS cost of any point of C . Our objective is to compute upper and lower bounds on this quantity.

In the previous section (see ‘‘Upper and lower bounds’’) we explained how the algorithm samples a representative vector of slope coefficients from C . To compute the cell’s upper bound, $\Delta^+(C)$, we compute the value of the y -intercept coefficient that, when combined with the representative vector of slope coefficients, produces the hyperplane that minimizes the LTS cost. It is straightforward to show that this can be reduced to solving a 1-dimensional LTS problem, which can then be solved in $O(n \log n)$ time [1]. We will present the complete algorithm, since it will be illustrative when we generalize this to the more complex problem of computing the cell’s lower bound.

Recall from Section 3 that, when computing $\Delta^+(C)$, we use the reduced coverage value h^- . To simplify notation, we will refer to this as h throughout this section. Let $p_i = (\mathbf{x}_i, y_i)$ denote the i th point of P , where $\mathbf{x}_i \in \mathbb{R}^{d-1}$ and $y_i \in \mathbb{R}$. The squared residual of p_i with respect to a hyperplane $\boldsymbol{\beta} = (\beta_1, \dots, \beta_d)$, which we denote by $r_i^2(\boldsymbol{\beta})$, is

$$r_i^2(\boldsymbol{\beta}) = \left(y_i - \left(\sum_{j=1}^{d-1} \beta_j x_{i,j} + \beta_d \right) \right)^2 = \left(\beta_d - \left(y_i - \sum_{j=1}^{d-1} \beta_j x_{i,j} \right) \right)^2. \quad (3)$$

Let $b_i = y_i - \sum_{j=1}^{d-1} \beta_j x_{i,j}$. We can visualize b_i as the intersection of the y -axis and the hyperplane passing through p_i whose slope coefficients match $\boldsymbol{\beta}$ (see Fig. 5(a)). Let $B = \{b_1, \dots, b_n\}$ denote this set of intercepts. For any hyperplane $\boldsymbol{\beta}$, the squared distance $(b_i - \beta_d)^2$ between intercepts is clearly equal to $r_i^2(\boldsymbol{\beta})$. Therefore, the desired intercept β_d of the LTS hyperplane whose slope coefficients match $\boldsymbol{\beta}$ ’s is the point on the y -axis that minimizes the sum of squared distances to its h closest points in B . That is, computing β_d reduces to solving the (1-dimensional) LTS problem on the n -element set B .

³Four priority queues are maintained, one for each of the cell-selection criteria. On creation, each cell is inserted into all four queues, and whenever a cell is selected for processing, it is removed from all four queues. By storing each priority queue as a balanced binary search tree [27] sorted according to the selection criterion together with cross-reference links, it is possible to perform each operation in $O(\log m)$ time.

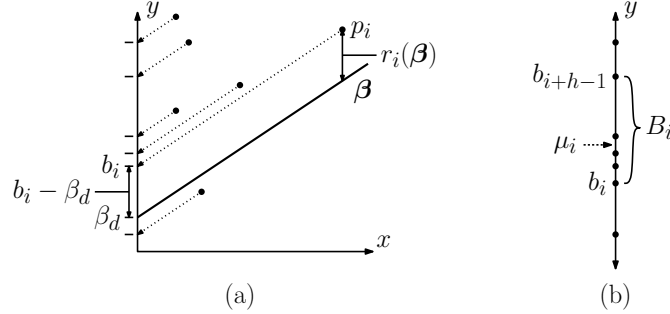


Figure 5: Computing the optimum y -intercept for a given slope β : (a) projecting the points and (b) the subset $\{b_i, \dots, b_{i+h-1}\}$.

Before presenting the algorithm for the 1-dimensional LTS problem, we begin with a few observations. First, let us assume that the points of B are sorted in nondecreasing order and have been renumbered so that $b_1 \leq \dots \leq b_n$. Clearly, the h -element subset B that minimizes the sum of squared distances to any point on the y -axis consists of h consecutive points of B . For $1 \leq i \leq n - h + 1$, define $B_i = \{b_i, \dots, b_{i+h-1}\}$ (see Fig. 5(b)). The desired upper bound is just the minimum among these $n - h + 1$ least-squares costs. This can be computed by brute force in $O(hn) = O(n^2)$ time, but we will present an $O(n)$ time incremental algorithm.

The point that minimizes the sum of squared deviates within each subset B_i is just the mean of the subset, which we denote by μ_i . The least-squares cost is the subset's standard deviation, which we denote by σ_i . Define $\text{SUM}(i) = \sum_{k=i}^{i+h-1} b_k$ and $\text{SSQ}(i) = \sum_{k=i}^{i+h-1} b_k^2$. Given these, it is well known that we can compute the mean and standard deviation as

$$\mu_i = \frac{\text{SUM}(i)}{h} \quad \text{and} \quad \sigma_i^2 = \left(\frac{\text{SSQ}(i) - (\text{SUM}(i)^2/h)}{h-1} \right)^{1/2}. \quad (4)$$

Thus, it suffices to show how to compute $\text{SUM}(i)$ and $\text{SSQ}(i)$, for $1 \leq i \leq n - h + 1$. We begin by computing the values of $\text{SUM}(1)$ and $\text{SSQ}(1)$ in $O(h)$ time by brute force. As we increment the value of i , we can update the quantities of interest in $O(1)$ time each. In particular,

$$\text{SUM}(i+1) \leftarrow \text{SUM}(i) + (b_{i+h} - b_i) \quad \text{and} \quad \text{SSQ}(i+1) \leftarrow \text{SSQ}(i) + (b_{i+h}^2 - b_i^2).$$

Thus, given the initial $O(n \log n)$ time to sort the points and the $O(h)$ time to compute $\text{SUM}(1)$ and $\text{SSQ}(1)$, the remainder of the computation runs in constant time for each i . Therefore, the overall running time is $O(n \log n + h + (n - h + 1)) = O(n \log n)$. The final LTS cost is just the minimum of the σ_i^2 's.

Next, we will show how to generalize this to the more complex task of computing the lower bound, $\Delta^-(C)$. Recall that the objective is to compute a lower bound on the LTS cost of *any* hyperplane β whose slope coefficients lie within the $(d-1)$ -dimensional hyperrectangle $[\beta^-, \beta^+]$. We will demonstrate that this can be reduced to the problem of solving a variant of the 1-dimensional LTS problem over a collection of n intervals, called the *interval LTS problem*, and we will present an $O(n \log n)$ time algorithm for this problem.

Let (x_i, y_i) denote the coordinates of $p_i \in P$, where $x_i \in \mathbb{R}^{d-1}$ and $y_i \in \mathbb{R}$. For any hyperplane β , recall the formula of Eq. (3) for the squared residual $r_i^2(\beta)$. We seek the value of β_d that

minimizes the sum of the h smallest squared residuals subject to the constraint that β 's slope coefficients lie within $[\beta^-, \beta^+]$. Although we do not know the exact values of these slope coefficients, we do have bounds on them. Using the fact that $\beta_j^- \leq \beta_j \leq \beta_j^+$, for $1 \leq j \leq d-1$, we can derive upper and lower bounds on the second term in the squared residual formula (Eq. (3)) by a straightforward application of interval arithmetic [28]. In particular, for $1 \leq i \leq n$,

$$b_i^+ = y_i - \sum_{j=1}^{d-1} \beta'_j x_{i,j}, \quad \text{where } \beta'_j = \begin{cases} \beta_j^- & \text{if } x_{i,j} \geq 0, \\ \beta_j^+ & \text{if } x_{i,j} < 0. \end{cases} \quad (5)$$

$$b_i^- = y_i - \sum_{j=1}^{d-1} \beta''_j x_{i,j}, \quad \text{where } \beta''_j = \begin{cases} \beta_j^+ & \text{if } x_{i,j} \geq 0, \\ \beta_j^- & \text{if } x_{i,j} < 0. \end{cases} \quad (6)$$

These quantities are analogous to the values b_i used in the upper bound, subject to the uncertainty in the slope coefficients. It is easy to verify that $b_i^- \leq b_i^+$, and for any β whose slope coefficients lie within $[\beta^-, \beta^+]$, the value of $r_i^2(\beta)$ is of the form $(\beta_d - b_i)^2$, for some $b_i \in [b_i^-, b_i^+]$. Analogous to the case of the upper bound computation, this is equivalent to saying that any hyperplane passing through p_i along the direction of β intersects the y -axis at a point in the interval $[b_i^-, b_i^+]$ (see Fig. 6). For the purposes of obtaining a lower bound on the LTS cost, we may take b_i to be the closest point in $[b_i^-, b_i^+]$ to the hypothesized intercept β_d . Thus, we have reduced our lower bound problem to computing the value of β_d that minimizes the sum of the smallest h squared distances to a collection of n intervals. This is a generalization of the 1-dimensional LTS problem, but where instead of points, we now have intervals.

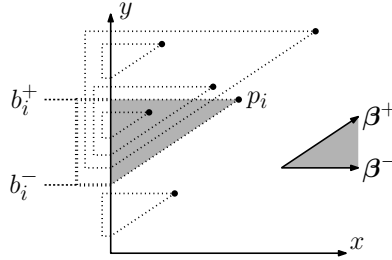


Figure 6: Reducing the lower bound to the interval LTS problem.

More formally, we define the *interval LTS problem* as follows. Given a point b and a closed interval, define the *distance* from b to this interval to be the minimum distance from b to any point of the interval. (The distance is zero if b is contained within the interval.) Given a set B of n nonempty, closed intervals on the real line and a coverage value h , for any real b , define $r_{[i]}(b, B)$ to be the distance from b to the i th closest interval of B . Also, define

$$\Delta_b(B, h) = \left(\frac{1}{h-1} \sum_{i=1}^h r_{[i]}^2(b, B) \right)^{1/2}.$$

The interval LTS problem is that of computing the minimum value of $\Delta_b(B, h)$, over all reals b , which we denote by $\Delta(B, h)$. In summary, we obtain the following reduction.

Lemma 3. Consider an n -element point set P in \mathbb{R}^d , a coverage h , and a cell C defined by the slope bounds $\beta^-, \beta^+ \in \mathbb{R}^{d-1}$. Let B denote the set of n intervals $[b_i^-, b_i^+]$ defined in Eqs. (5) and (6) above. Then for any $\beta \in C$, $\Delta(B, h) \leq \Delta_\beta(P, h)$.

It suffices to show how to solve the interval LTS problem. Our algorithm is a generalization of the algorithm for the upper-bound processing. First, we will identify a collection of $O(n)$ subsets of intervals of size h , such that the one with the lowest cost will be the solution to the interval LTS problem. Second, we will show how to compute the LTS costs of these subsets efficiently through an incremental approach. To avoid discussion of messy degenerate cases, we will make the general-position assumption that the interval endpoints of B are distinct. This can always be achieved through an infinitesimal perturbation of the endpoint coordinates.

Recall that in the upper-bound algorithm, we computed the LTS cost of each h consecutive points. In order to generalize this concept to the interval case, we first sort the left endpoints B in nondecreasing order, letting $b_{[j]}^-$ denote the j th smallest left endpoint. Similarly, we sort the right endpoints, letting $b_{[i]}^+$ denote the i th smallest right endpoint. Define I_i to be the (possibly empty) interval $[b_{[i]}^+, b_{[i+h-1]}^-]$, and define B_i to be the (possibly empty) subset of intervals of B whose right endpoints are at least $b_{[i]}^+$, and whose left endpoints are at most $b_{[i+h-1]}^-$. To avoid future ambiguity with the term “interval,” we will refer to any interval formed from the endpoints of B as a *span*, and each interval I_i is called an *h -span* (see Fig. 7). We next present a technical result about the sets B_i . The proof is given in Appendix A.

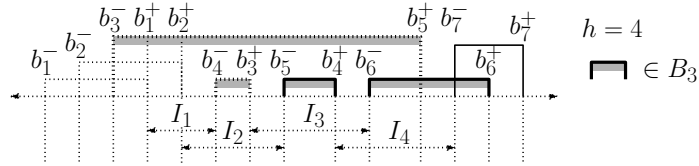


Figure 7: For $h = 4$, the h -spans I_i . The intervals of the subset B_3 corresponding to $I_3 = [b_{[3]}^+, b_{[6]}^-]$ are highlighted.

Lemma 4. *Let (B, h) be an instance of the interval LTS problem, and let B_i be the subsets of B defined above, for $1 \leq i \leq n - h + 1$. If for any i , $b_{[i]}^+ > b_{[i+h-1]}^-$, then $\Delta(B, h) = 0$. Otherwise, each set B_i consists of exactly h intervals of B .*

By this lemma we may assume henceforth that $b_{[i]}^+ > b_{[i+h-1]}^-$ for all i , since if we ever detect that this is not the case, we may immediately report that $\Delta(B, h)$ is zero and terminate the algorithm. Intuitively, the subsets B_i are chosen to be the “tightest” subsets of B having h elements. Our next lemma shows that for the purposes of solving the interval LTS problem, it suffices to consider just these subsets. For $1 \leq i \leq n - h + 1$, define $\Delta(B_i, h)$ to be the LTS cost of the interval LTS problem on only the h intervals of B_i .

Lemma 5. *Given an instance (B, h) of the interval LTS problem, $\Delta(B, h) = \Delta(B_i, h)$ for some i , where $1 \leq i \leq n - h + 1$.*

The proof is given in Appendix A. It follows that in order to solve the interval LTS problem, it suffices to compute the LTS cost of each of the subsets B_i , for $1 \leq i \leq n - h + 1$. The algorithm is a generalization of the upper bound algorithm, where intervals play the role of individual points. The details of the algorithm are relatively technical but straightforward. A complete presentation of the algorithm and its running time analysis are presented in Appendix B. The following theorem summarizes the main result of this section.

Theorem 2. *Given a collection B of n intervals on the real line and a coverage h , it is possible to solve the interval LTS problem on B in $O(n \log n)$ time.*

5. Empirical Analysis

In this section we present a number of experiments comparing the accuracy and running time of Adaptive-LTS and Fast-LTS on various input distributions. We consider input sets from various distributions, in various dimensions, and with various degrees of numerical conditioning.

Before presenting our results, let us remark on the structure that is common to all our experiments. We implemented both Rousseeuw and van Driessen’s Fast-LTS and our Adaptive-LTS in C++ (compiled by g++ 4.5.2 with optimization level “-O3” and running on a dual-core Intel Pentium-D processor under Ubuntu Linux 11.04). In all experiments we let the algorithm choose the initial cell bounds using the approach described under “Initial Cell” in Section 3.

Both Fast-LTS and Adaptive-LTS operate in a series of stages. Each stage of our implementation of Fast-LTS consists of a single random elemental fit followed by two C-steps. (This is the same number used by Rousseeuw and van Driessen in [15].) Although Adaptive-LTS has a termination condition, for the sake of comparison with Fast-LTS we disabled the termination condition in most of our experiments, and instead ran the two algorithms for the same fixed number of stages and compared their performance on a stage-by-stage basis. For the same reason, we omitted the final post-processing phase of Fast-LTS.

5.1. Convergence Rates for Synthetically Generated Data

For these experiments, we generated a data set consisting of 1000 points in \mathbb{R}^d , for $d \in \{2, 3\}$. In order to simulate point sets that contain outliers, each data set consisted of the union of points drawn from two distributions. A fixed fraction of points were sampled independently from an *inlier distribution*, which consists of points that lie close to a randomly generated hyperplane. The remaining points were sampled from some different distribution, called the *outlier distribution*. These distributions were varied to explore different numbers of dimensions and different degrees of numerical conditioning.

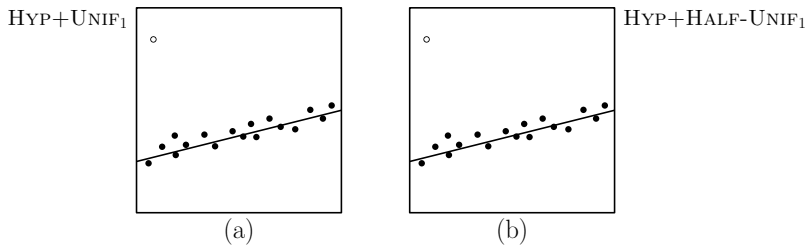


Figure 8: The distributions HYP+UNIFORM₁ and HYP+HALF-UNIF₁.

HYP+UNIFORM (2-dimensional): Our first experiment involved a 2-dimensional, well-conditioned distribution, consisting of 1000 points (one independent and one dependent variable). First, a line of the form $y = \beta_1 x + \beta_2$ was generated, where β_1 was sampled uniformly from $[-0.25, +0.25]$, and β_2 was sampled uniformly from $[-0.1, +0.1]$. Next, 550 points (55% inliers) were generated with the x -coordinates sampled uniformly from $[-1, +1]$, and each y -coordinate generated by computing the value of the line at these x -coordinates and adding a Gaussian deviate with mean 0 and standard deviation 0.01. The remaining 450 points (45% outliers) were sampled uniformly from the square $[-1, +1]^2$. Because the inliers were drawn from a hyperplane (actually a line in this case) and the outliers were drawn from a uniform distribution, we call this distribution

HYP+UNIFORM or generally HYP+UNIFORM $_k$ to denote the variant with k independent variables (see Fig. 8).

We ran both Fast-LTS and Adaptive-LTS on the resulting data set for 500 stages. Although in most of our experiments, we set the coverage value to 50%, to make this instance particularly easy to solve, we ran both algorithms with a relatively low coverage of $h = 0.1$ (that is, 10% inliers). In Fig. 9(a) and (b) we show the results for Fast-LTS and Adaptive-LTS, respectively, where the horizontal-axis indicates the stage number and the vertical-axis indicates the LTS cost. Note that the vertical-axis is on a logarithmic scale. Each (small square) point of each plot indicates the LTS cost of the sampled hyperplane. The upper (solid blue or “Best”) curve indicates the minimum LTS cost of any hyperplane seen so far. The lower (broken red or “Lower bound”) of Fig. 9(b) indicates the minimum lower bound of any active cell. The optimum LTS cost lies somewhere between these two curves, and therefore, the ratio between the current best and current lower bound limits the maximum relative error that would result by terminating the algorithm and outputting the current best fit.

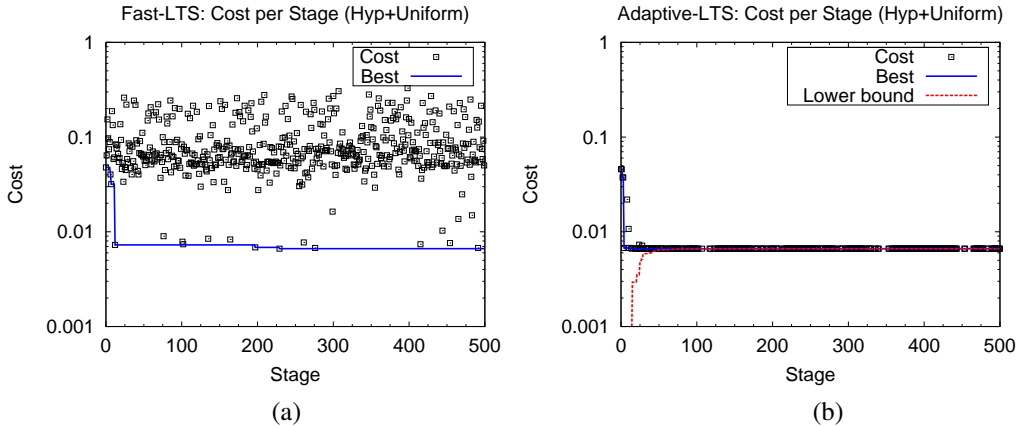


Figure 9: LTS cost versus stage number for the HYP+UNIFORM distribution ($d = 2$, $h = 0.1$) for (a): Fast-LTS and (b): Adaptive-LTS.

As mentioned earlier, this is a well-conditioned distribution, and therefore it not surprising that both algorithms converge very rapidly to a solution that is nearly optimal. Given only information on the best LTS cost, a user of Fast-LTS would not know when to terminate the search. In contrast, a user of Adaptive-LTS would know that after just 75 iterations, the relative difference between the best upper and lower bounds is smaller than 1%. Thus, terminating the algorithm at this point would imply an approximate error of at most 1%. It is also interesting to note that Fast-LTS generates elemental fits at random, and hence the pattern of plotted points remains consistently random throughout all 500 stages. In contrast, the sampled hyperplanes generated by Adaptive-LTS shows signs of converging to a subregion of the configuration space where the most promising solutions reside.

HYP+UNIFORM (3-dimensional): Our second experiment involved the same distribution, but in \mathbb{R}^3 (two independent and one dependent variable) and using a larger value of h . First, a plane of the form $y = \beta_1 x_1 + \beta_2 x_2 + \beta_3$ was generated, where β_1 and β_2 were sampled uniformly from $[-0.25, +0.25]$, and β_3 was sampled uniformly from $[-0.1, +0.1]$. As before, 550 points (55%

inliers) were generated with the x_1 - and x_2 -coordinates sampled uniformly from $[-1, +1]$, and the y -coordinates generated by computing the value of the plane at these x -coordinates and adding a 3-dimensional deviate whose coordinates are independent Gaussians of mean 0 and standard deviation 0.01. The remaining 450 points (45% outliers) were sampled uniformly from the cube $[-1, +1]^3$. We refer to this as `HYP+UNIFORM2`.

We ran both Fast-LTS and Adaptive-LTS on the data set for 500 stages with $h = 0.5$. In Fig. 10(a) and (b) we show the results for Fast-LTS and Adaptive-LTS, respectively. Again, both algorithms converge rapidly to a solution that is nearly optimal. However, as in the previous experiment, a user of Adaptive-LTS could use the ratio between the best upper and lower bounds to limit the maximum error in the LTS cost. For example, after roughly 250 stages the Adaptive-LTS relative error is within 20%, and after roughly 400 stages it is within 10%. Observe that the rate of convergence is slower than in the earlier 2-dimensional experiment. We conjecture that this is due to the increase in the dimension, which results in searching a larger solution space (see Section 5.2).

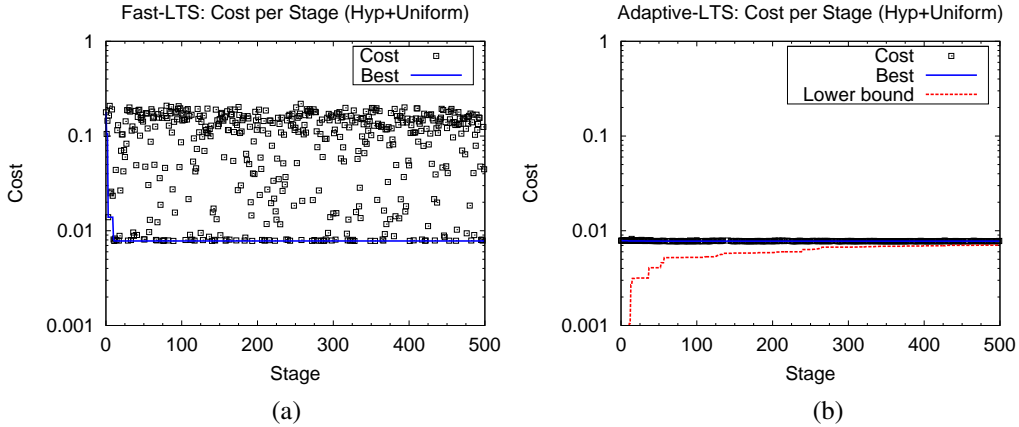


Figure 10: LTS cost versus stage number for the `HYP+UNIFORM` distribution ($d = 3$, $h = 0.5$) for (a): Fast-LTS and (b): Adaptive-LTS.

`HYP+HALF-UNIF`: This involved a data set consisting of 1000 points in \mathbb{R}^3 (two independent and one dependent variable), in which the outlier distribution was more skewed. The inliers were generated in the same manner as in the previous experiment. The outliers were sampled uniformly from the portion of the 3-dimensional hypercube $[-1, +1]^3$ that lies above the plane. (Fig. 8(b) illustrates this distribution in \mathbb{R}^2 .) As before, we ran both Fast-LTS and Adaptive-LTS on the data set for 500 stages with $h = 0.5$. In Fig. 11(a) and (b) we show the results for Fast-LTS and Adaptive-LTS, respectively. The convergence rates of the two algorithms are similar to the previous experiment. For example, after roughly 200 stages the relative error is within 20%, and after roughly 300 stages it is within 10%.

`FLAT+SPHERE`: For this experiment we designed a rather pathological distribution with the intention of challenging both algorithms. The inlier distribution was chosen so that almost all the inliers are degenerate (thus resulting in many unstable elemental fits) and the remaining few inliers are a huge distance away (thus penalizing any inaccurate elemental fits).

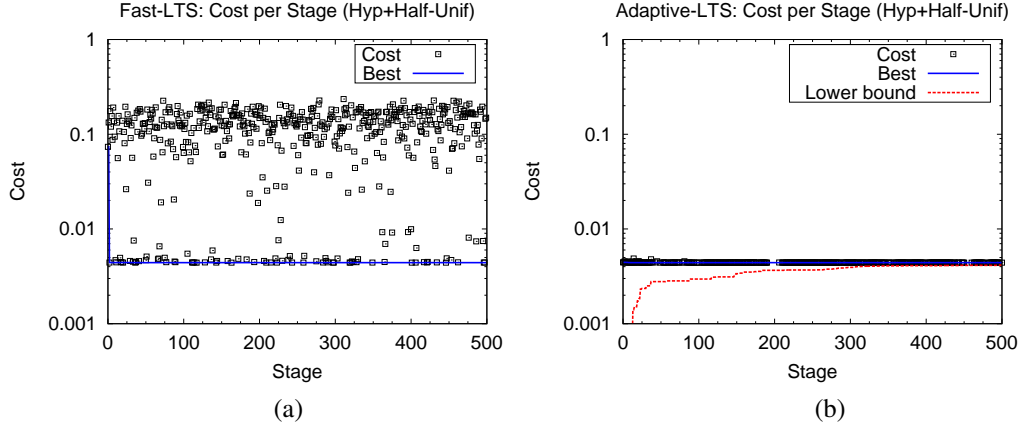


Figure 11: LTS cost versus stage number for the HYP+HALF-UNIF distribution ($d = 3, h = 0.5$) for (a): Fast-LTS and (b): Adaptive-LTS.

The point set consisted of 1000 points in \mathbb{R}^3 with 500 inliers and 500 outliers. Let S be a large sphere centered at the origin with radius 10,000, let β be a plane passing through the origin, and let ℓ be a line passing through the origin that lies on β (see Fig. 12). The inliers were drawn from two distributions. First, 490 *local inliers* were sampled uniformly along a line segment of length two centered at the origin and lying on ℓ , where each point was perturbed by a multivariate Gaussian with mean zero and standard deviation 0.10. The remaining 10 inliers, called *leverage inliers*, were sampled uniformly from the equatorial circle where β intersects S . We set the coverage to $h = 0.5$, implying that exactly 500 points are used in the computation of the LTS cost. Finally, the outliers consisted of 500 points that were sampled uniformly from S . We refer to this distribution as FLAT+SPHERE₂.

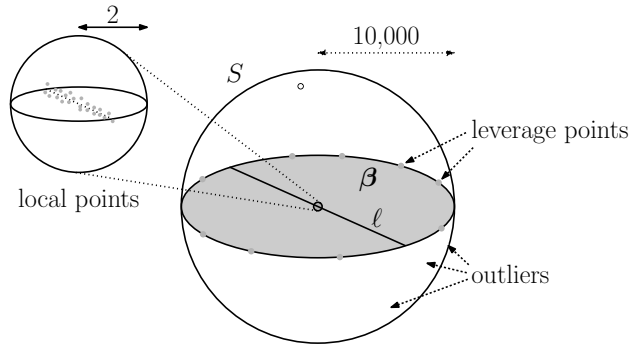


Figure 12: The FLAT+SPHERE distribution.

To see why this distribution is particularly difficult, observe that barring lucky coincidences, the only way to obtain a low-cost LTS fit is to combine all the local inliers with all the leverage inliers to obtain a fit that lies very close to β . Thus, any elemental fit that involves even one outlier is unlikely to provide a good fit. Because of their degenerate arrangement along ℓ , any random sample of inliers that fails to include at least one leverage inlier is unlikely to lie close to

β . The probability of sampling three inliers such that at least one is a leverage inlier is roughly $(500/1000)^3 \cdot 10/500 = 0.025$. Note that by decreasing the fraction of leverage inliers relative to the total number of inliers, we can make this probability of a useful elemental fit arbitrarily small. This point set clearly fails both of the criteria given in Section 2 for well-conditioned point sets, since the vast majority of inliers are degenerately positioned, and there are a nontrivial number of leverage inliers.

Our experiments bear out the problems of this poor numerical conditioning (see Fig. 13). Both algorithms cast about rather aimlessly until hitting upon a good fit (around stage 230 for Fast-LTS and stage 50 for Adaptive-LTS). This experiment demonstrates the principal advantage of Adaptive-LTS over Fast-LTS. A user of Fast-LTS who saw only the results of the first 200 iterations might be tempted to think that the algorithm had converged to an optimal LTS cost of roughly 3.8, rather than the final cost of roughly 0.102. Terminating the algorithm prior to stage 200 would have resulted in a relative error of over 3000%. In contrast, a user of Adaptive-LTS would know by stage 200 that the relative error in the final LTS cost is less than 30% and after stage 300 it is close to 10%.

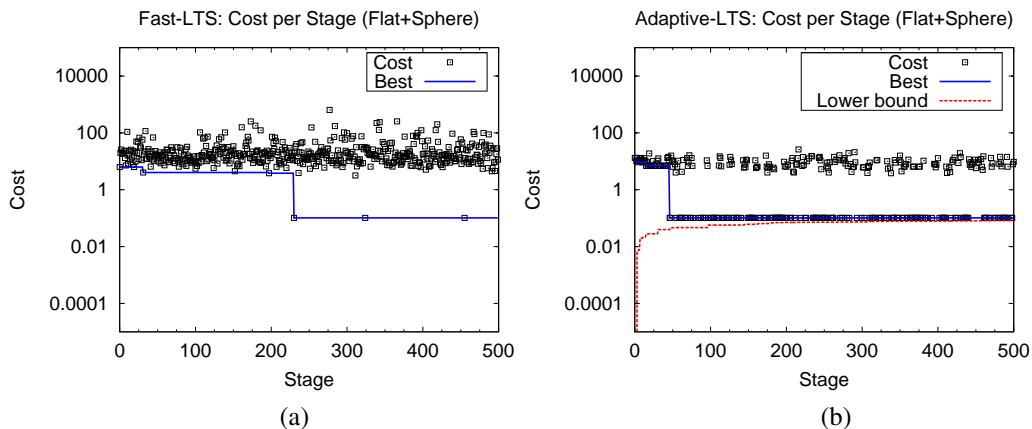


Figure 13: LTS cost versus stage number for the FLAT+SPHERE distribution ($d = 3$, $h = 0.5$) for (a): Fast-LTS and (b): Adaptive-LTS.

5.2. Execution Times and Dimension Dependence

So far we have analyzed only the LTS costs generated by the two algorithms. An important question is how the performance of Adaptive-LTS compares with Fast-LTS. Table 1 summarizes the execution times and LTS costs for both of these algorithms for 500 stages. In all cases, both algorithms generated essentially the same result and, hence, achieved the same LTS cost. While the execution times of Adaptive-LTS were consistently higher than those of Fast-LTS, it was only by a small constant factor.

Another issue is how the algorithms scale with dimension. Assuming that $n \gg d$, the running time of a fixed number k of iterations of Fast-LTS grows as $O(kdn)$. Of course, Theorem 1 predicts that the number of iterations needed to achieve an accurate fit is expected to grow exponentially with the dimension. In contrast, the running time of Adaptive-LTS is much harder to predict. This is because it depends on the effectiveness of the branch-and-bound pruning. By Lemma 2), the running time per cell grows as $O(n \log n)$, but it is reasonable to believe that the

Table 1: Summary of execution times and LTS costs for the synthetic experiments for 500 iterations of each algorithm. Running times are given in CPU seconds.

Distribution	Indep. Vars.	d	h	Fig.	Fast-LTS		Adaptive-LTS	
					Time	LTS Cost	Time	LTS Cost
HYP+UNIFORM ₁	1	2	0.1	9	0.14	0.00660	0.44	0.00660
HYP+UNIFORM ₂	2	3	0.5	10	0.47	0.00776	0.83	0.00776
HYP+HALF-UNIF ₂	2	3	0.5	11	0.48	0.00440	0.73	0.00440
FLAT+SPHERE ₂	2	3	0.5	13	0.81	0.10220	1.11	0.10220

number of cells will grow exponentially with the dimension of the solution space, which equals the number of independent variables. An intuitive explanation is that the accuracy of a cell's lower bound depends on its diameter, and the number of cells of a given diameter needed to cover a hypercube of a given side length (the initial cell) grows exponentially with the dimension of the hypercube.

We ran both Fast-LTS and Adaptive-LTS on the same 1000-point instances of HYP+UNIFORM, where the number of independent variables ranged from one to nine (that is, $2 \leq d \leq 10$). As in the earlier distributions, 55% of the points (inliers) had $d - 1$ independent variables sampled uniformly from $[-1, +1]^{d-1}$ and the remaining dependent variable was generated to lie on a random $(d - 1)$ -dimensional hyperplane plus a Gaussian error of mean 0 and standard deviation 0.01. The remaining 45% of points (outliers) were generated uniformly in a d -dimensional hypercube. Because it has $d - 1$ independent variables, we refer to this distribution as HYP+UNIFORM _{$d-1$} . In each case we requested a coverage of 50%.

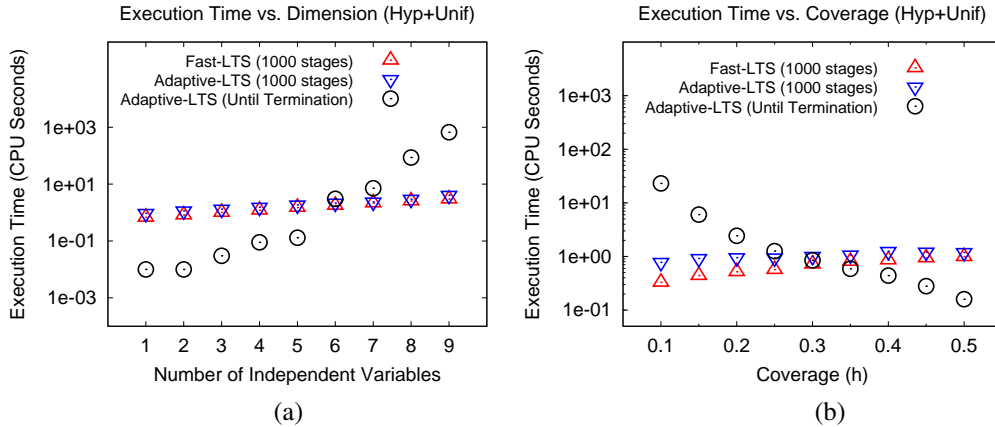


Figure 14: Execution times as a function of (a) the number of independent variables and (b) the coverage value.

We ran Fast-LTS and Adaptive-LTS each for 1000 iterations. We also ran Adaptive-LTS to its natural termination with $\varepsilon_r = 10$ and $\varepsilon_q = 0$. (That is, we allowed an error factor of 10 in the LTS cost, and allowed for no reduction in the coverage.) Fig. 14(a) shows the the execution times in CPU seconds (plotted on a log scale). All three converged to nearly the same final fit and had very similar LTS costs. As hypothesized, the execution times for the fixed-iteration versions are

relatively insensitive to variations in the dimension, but the execution times of full Adaptive-LTS show a strong dependence on dimension.

Next, we considered the impact of varying the coverage parameter h . We ran the same experiment as above but with the number of independent variables fixed at three ($d = 4$) and with h varying from 10% to 50%. Fig. 14(b) shows the the execution times in CPU seconds (plotted on a log scale). The execution times for the fixed-iteration versions are relatively insensitive to the coverage. The execution time of full Adaptive-LTS is strongly negatively correlated with h . This is rather surprising, since it would seem that decreasing the coverage should make the problem easier to solve. A detailed analysis reveals that the issue is the cell lower bounds used by the branch-and-bound algorithm. Smaller coverage values resulted in smaller lower bound values. The upper bounds, while smaller, were not as strongly affected. As a result, cells were not killed as efficiently, and this resulted in larger running times.

5.3. Empirical Analysis of Numerical Condition

The variation in convergence rates seen in Section 5.1 raises the question how these data sets compare to one another with respect to the numerical condition measures described in Section 2. To evaluate this, we ran a series of experiments involving both the low-dimensional data sets from Section 5.1 and two higher-dimensional data sets, `HYP+UNIFORM10` and `HYP+UNIFORM20`. These data sets and the associated coverage values used in the experiments are summarized in Table 2.

Table 2: Summary of data sets (inlier subsets)

Distribution	Indep. Vars.	Inliers	Residual Std. Dev.	q
<code>HYP+UNIFORM₁</code>	1	110	0.01	10%
<code>HYP+UNIFORM₂</code>	2	550	0.01	50%
<code>HYP+HALF-UNIF₂</code>	2	550	0.01	50%
<code>HYP+UNIFORM₁₀</code>	10	550	0.01	50%
<code>HYP+UNIFORM₂₀</code>	20	550	0.01	50%
<code>FLAT+SPHERE₂</code>	2	500	0.10	50%

For each data set, we computed the numerical condition measures γ and λ , as defined in Eq.(1) from Section 2. These are shown in Table 3. In particular, the matrix condition measure γ is the (estimated) median condition number of the set of elemental-fit matrices X_E . Its value was computed by randomly sampling 10,000 elemental subsets and computing the median (Frobenius) condition number of the resulting matrices. The leverage measure λ was computed (exactly) from Eq.(1) as the square root of the ratio of the average and the median of squared norms of the vectors of independent variables.

We also computed the number of elemental fits from Theorem 1 needed to guarantee (with probability at least $1 - \delta$) a fit that achieves a relative error of at most $\sqrt{2}\gamma\lambda$. This is shown in the fourth column of Table 3. Since Theorem 1 merely provides bounds, we ran an experiment to determine how well random elemental fits perform for these data sets. For each data set (consisting of 1000 points, including both inliers and outliers) we generated one million random elemental fits. In each case we computed the LTS cost of the resulting fit. In order to be faithful to Theorem 1 we did not apply any C-steps. We estimated the optimum LTS cost based on the noise standard deviation used when generating the inliers (shown in the fourth column of Table 2). Let

Table 3: Results of experiments on numerical condition for synthetic data sets.

Distribution	γ	λ	Elem-Fits	err \leq 1%	err \leq 10%	err \leq 100%
HYP+UNIFORM ₁	4.08	1.142	7369	0.78%	0.88%	1.75%
HYP+UNIFORM ₂	9.18	1.027	1179	3.80%	5.02%	11.20%
HYP+HALF-UNIF ₂	9.25	1.018	1179	9.18%	9.83%	13.39%
HYP+UNIFORM ₁₀	70.80	1.009	7.7×10^7	0.01%	0.03%	0.04%
HYP+UNIFORM ₂₀	195.08	1.006	8.1×10^{13}	0.00%	0.00%	0.00%
FLAT+SPHERE ₂	33.99	2606	1179	0.00%	0.01%	0.01%

Δ_i denote the LTS cost of the i th elemental fit, and let σ denote the this standard deviation. For $p \in \{0.01, 0.1, 1.0\}$, we declared that this elemental fit achieves a *relative error of at most $p\%$* if $\Delta_i \leq (1 + p/100)\sigma$. The fractions of fits satisfying these error bounds are shown in the last three columns of Table 3. Observe, for example, that for the data set HYP+UNIFORM₂, roughly 5% of the elemental fits achieved a relative error of at most 10%, suggesting that fits of this accuracy occur randomly at a rate of roughly one out of every twenty. These results qualitatively support Theorem 1, but the observed errors were quite a bit smaller than the bounds given by the theorem.

5.4. DPOSS Data Set

Our final experiment involved a point set from the Digitized Palomar Sky Survey (DPOSS), from the California Institute of Technology [29]. The data set consisted of 132,402 points in \mathbb{R}^6 , from which we sampled 10,000 at random for each of our experiments. Of the six coordinates, we considered the problem of estimating the first coordinate MAperF as a function of one or more of the coordinates csfF, csfJ, csfN. Fig. 15 shows a scatter plot of MAperF versus csfF, and also shows the LTS fit and an ordinary least square (OLS) fit. We used the same experimental setup as in the synthetic experiments, by running both Fast-LTS and Adaptive-LTS for 500 stages.

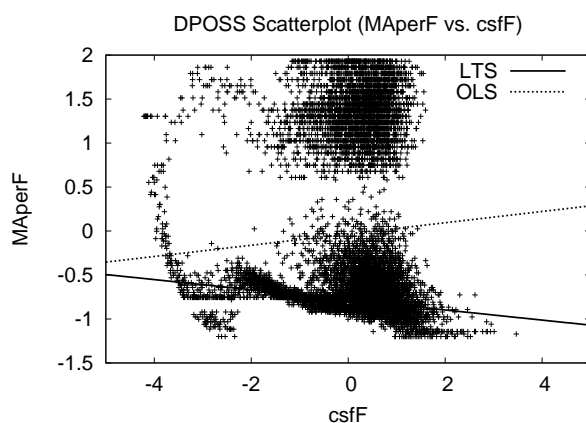


Figure 15: The DPOSS data set (MAperF vs. csfF).

In the first experiment, we considered a 2-dimensional instance of estimating MAperF as a function of csfF (see Fig. 16). Both algorithms rapidly converged to essentially the same solution

(shown in Fig. 15). Because it provides a lower bound, a user of Adaptive-LTS would know after just 30 stages that the relative error in the LTS cost of this solution is less than 5%.

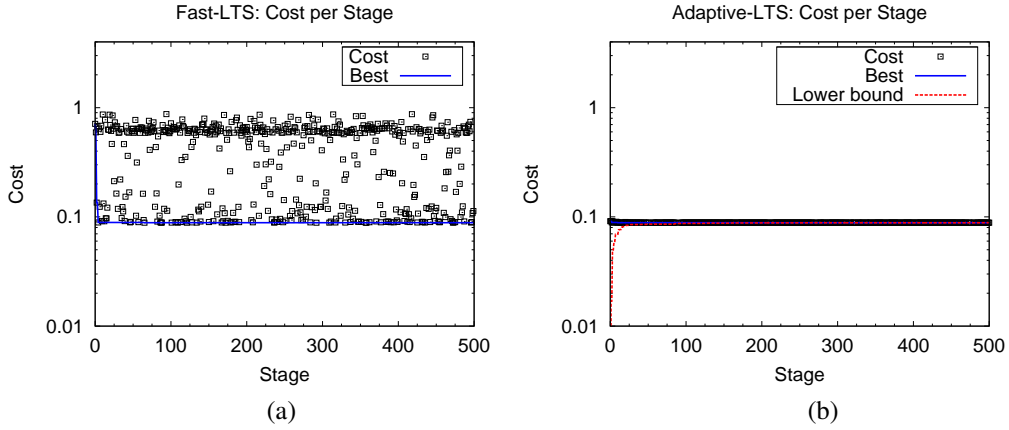


Figure 16: LTS cost versus stage number for estimating MAperF as function of csfF in the DPOSS data set ($d = 2$) for (a): Fast-LTS and (b): Adaptive-LTS.

Our next experiment considered a 3-dimensional instance of estimating MAperF as a function of both csfF and csfJ (see Fig. 17). Again, both algorithms rapidly converged to essentially the same solution. A user of Adaptive-LTS would know that after roughly 250 stages the Adaptive-LTS relative error is within 20%, and after roughly 500 stages it is within 10%.

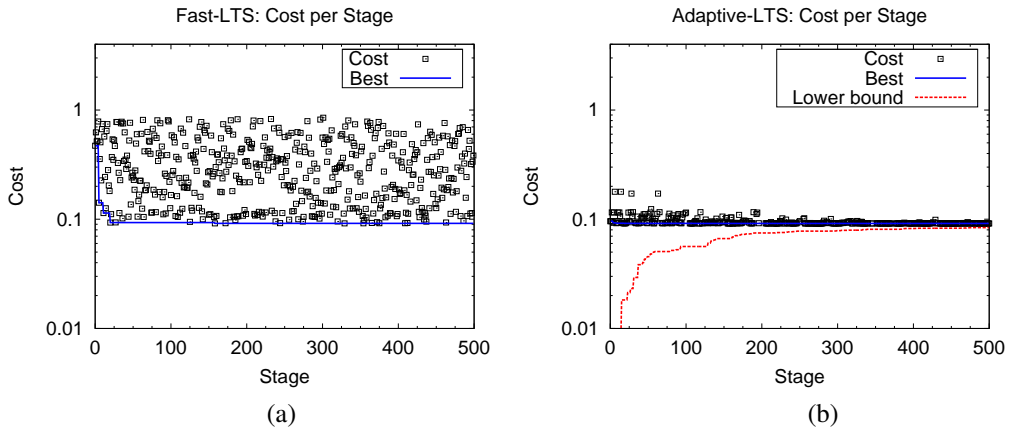


Figure 17: LTS cost versus stage number for the DPOSS data set ($d = 3$) for (a): Fast-LTS and (b): Adaptive-LTS.

Our final experiment demonstrates one of the practical limitations of Adaptive-LTS. We considered a 4-dimensional instance of estimating MAperF as a function of csfF, csfJ, and csfN (see Fig. 18).

Again, both algorithms rapidly converged to essentially the same solution. However, due to the increase in dimension, the convergence in the error bound is much slower. After 500 stages, the relative error bound for Adaptive-LTS is over 100%. Although it is not evident from the plot,

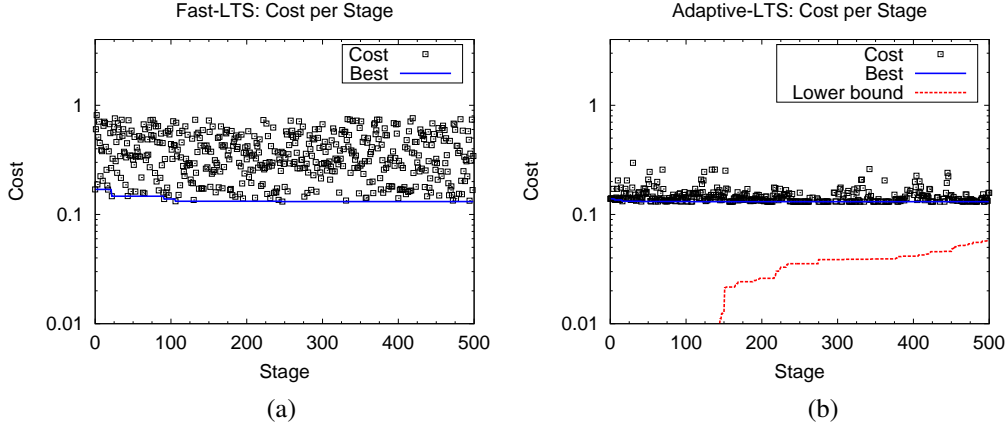


Figure 18: LTS cost versus stage number for the DPOSS data set ($d = 4$) for (a): Fast-LTS and (b): Adaptive-LTS.

it is only after roughly 5,000 stages that the error bound falls below 20%, and after 10,000 stages it falls below 10%.

As in our synthetic experiments, we compared the relative performance of Adaptive-LTS and Fast-LTS. In Table 4 we summarize the execution times and LTS costs for both of these algorithms for 500 stages. Both algorithms generated essentially the same fitting hyperplane and, hence achieved the same LTS cost. Because of its higher overhead, Adaptive-LTS did require more execution time for the same number of stages, but in all instances the running time was less than 50% greater than Fast-LTS.

Table 4: Summary of execution times and LTS costs for the DPOSS experiments for 500 iterations of each algorithm. Running times are given in CPU seconds.

d	h	Fig.	Fast-LTS		Adaptive-LTS	
			Time	LTS Cost	Time	LTS Cost
2	0.5	16	5.02	0.0884	7.31	0.0884
3	0.5	17	5.77	0.0915	7.40	0.0915
4	0.5	18	7.72	0.1308	8.54	0.1311

6. Concluding Remarks

In this paper we have presented two results related to the computation of the linear least trimmed squares estimator. First, we introduced a measure of the numerical condition of a point set. We showed that if a point set is well conditioned, then it is possible to bound the accuracy (as a function of the conditioning parameters) of the LTS fit resulting from a sufficiently large number of random elemental fits. Second, we presented an approximation algorithm for LTS, called Adaptive-LTS. Given a point set, a coverage value, and bounds on the maximum and minimum slope coefficients, this algorithm computes a hybrid approximation to the optimum LTS fit assuming the slope coefficients are drawn from these bounds. We implemented this

algorithm and presented empirical evidence on a variety of data sets of this algorithm's efficiency and effectiveness. In contrast to existing practical approaches, this algorithm computes a lower bound on the optimum LTS cost. Therefore, when the algorithm terminates it provides an upper bound on the approximation error. The software is freely available from the first author's web page, <http://www.cs.umd.edu/~mount>.

There are a number of possible directions for future research. First, our analysis of well-conditioned point sets considered only the results of random elemental fits, but it ignored the effect of possible C-steps. It would be interesting to consider the effect of the numeric condition on the convergence rate of the C-steps. Second, an obvious weakness of our Adaptive-LTS algorithm is its reliance on the requirement that the initial bounds on the slope coefficients be given. Although the program can estimate these bounds from random elemental fits, we cannot guarantee that the optimum fit lies within these bounds. It would be nice to have the algorithm compute the initial bounds in such a manner that the optimum LTS cost is guaranteed to lie within them. Finally, although we showed that each stage of the algorithm runs in $O(m + n \log n)$ (where n is the number of points and m is the number of active cells), we do not have good bounds on the number of stages until termination. Deriving asymptotic time bounds on the overall running time of the algorithm would be useful.

7. Acknowledgments

We would like to thank Peter Rousseeuw for providing us with the DPOSS data set. Also, we are very grateful to the anonymous referees for their informative and helpful comments.

References

- [1] P. J. Rousseeuw, A. M. Leroy, *Robust Regression and Outlier Detection*, Wiley, New York, 1987.
- [2] P. J. Rousseeuw, Least median-of-squares regression, *Journal of the American Statistical Association* 79 (1984) 871–880.
- [3] D. L. Souvaine, J. M. Steele, Time- and space- efficient algorithms for least median of squares regression, *Journal of the American Statistical Association* 82 (1987) 794–801.
- [4] H. Edelsbrunner, D. L. Souvaine, Computing median-of-squares regression lines and guided topological sweep, *Journal of the American Statistical Association* 85 (1990) 115–119.
- [5] D. M. Mount, N. S. Netanyahu, C. Piatko, R. Silverman, A. Y. Wu, Quantile approximation for robust statistical estimation and k -enclosing problems, *International Journal of Computational Geometry & Applications* 10 (2000) 593–608.
- [6] D. M. Mount, N. S. Netanyahu, E. Zuck, Analyzing the number of samples required for an approximate Monte-Carlo LMS line estimator, in: M. Hubert, G. Pison, A. Struyf, S. V. Aelst (Eds.), *Theory and Applications of Recent Robust Methods*, Statistics for Industry and Technology, Birkhäuser, Basel, 2004, pp. 207–219.
- [7] D. M. Mount, N. S. Netanyahu, K. R. Romanik, R. Silverman, A. Y. Wu, A practical approximation algorithm for the LMS line estimator, *Computational Statistics & Data Analysis* 51 (2007) 2461–2486.
- [8] T. Bernholt, Computing the least median of squares estimator in time $O(n^d)$, in: *Proc. Intl. Conf. on Computational Science and Its Applications*, Vol. 3480 of Springer LNCS, Springer-Verlag, Berlin, 2005, pp. 697–706.
- [9] J. Erickson, S. Har-Peled, D. M. Mount, On the least median square problem, *Discrete & Computational Geometry* 36 (2006) 593–607.
- [10] J. Agulló, New algorithms for computing the least trimmed squares regression estimator, *Computational Statistics & Data Analysis* 36 (2001) 425–439.
- [11] M. Hofmann, E. J. Kontoghiorghes, Matrix strategies for computing the least trimmed squares estimation of the general linear and sur models, *Computational Statistics & Data Analysis* 54 (2010) 3392–3403.
- [12] M. Hofmann, C. Gatu, E. J. Kontoghiorghes, An exact least trimmed squares algorithm for a range of coverage values, *J. Comput. and Graph. Stat.* 19 (2010) 191–204.
- [13] O. Hössjer, Exact computation of the least trimmed squares estimate in simple linear regression, *Computational Statistics & Data Analysis* 19 (1995) 265–282.

- [14] D. M. Mount, N. S. Netanyahu, C. Piatko, R. Silverman, A. Y. Wu, On the least trimmed squares estimator, *Algorithmica* 69 (2014) 148–183.
- [15] P. J. Rousseeuw, K. van Driessen, Computing LTS regression for large data sets, *Data Mining and Knowledge Discovery* 12 (2006) 29–45.
- [16] F. Torti, D. Perrotta, A. C. Atkinson, M. Riani, Benchmark testing of algorithms for very robust regression: FS, LMS and LTS, *Computational Statistics & Data Analysis* 56 (2012) 2501–2512.
- [17] D. M. Hawkins, The feasible solution algorithm for least trimmed squares regression, *Computational Statistics & Data Analysis* 17 (1994) 185–196.
- [18] P. Diaconis, M. Shahshahani, The subgroup algorithm for generating uniform random variables, *Prob. in the Engr. and Inf. Sci.* 1 (1987) 15–32.
- [19] S. H. Friedberg, A. J. Insel, L. E. Spence, *Linear Algebra*, 4th Edition, Prentice Hall, 2002.
- [20] G. H. Golub, C. F. V. Loan, *Matrix Computations*, The Johns Hopkins University Press, Baltimore, Maryland, 1996.
- [21] P. J. Rousseeuw, A diagnostic plot for regression outliers and leverage points, *Computational Statistics & Data Analysis* 11 (1991) 127–129.
- [22] D. P. Huttenlocher, W. J. Rucklidge, A multi-resolution technique for comparing images using the Hausdorff distance, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, New York, 1993, pp. 705–706.
- [23] W. J. Rucklidge, Efficient visual recognition using the Hausdorff distance, no. 1173 in *Lecture Notes in Computer Science*, Springer-Verlag, Berlin, 1996.
- [24] W. J. Rucklidge, Efficiently locating objects using the Hausdorff distance, *International Journal of Computer Vision* 24 (1997) 251–270.
- [25] M. Hagedoorn, R. C. Veltkamp, Reliable and efficient pattern matching using an affine invariant metric, *International Journal of Computer Vision* 31 (1999) 103–115.
- [26] D. M. Mount, N. S. Netanyahu, J. Le Moigne, Efficient algorithms for robust point pattern matching, *Pattern Recognition* 32 (1999) 17–38.
- [27] T. H. Cormen, C. E. Leiserson, R. L. Rivest, C. Stein, *Introduction to Algorithms*, 3rd Edition, MIT Press, Cambridge, MA, 2009.
- [28] L. Jaulin, M. Kieffer, O. Didrit, E. Walter, *Applied Interval Analysis*, Springer-Verlag, 2001.
- [29] S. G. Djorgovski, R. R. Gal, S. C. Odewahn, de R. R. Carvalho, R. Brunner, G. Longo, R. Scaramella, The Palomar digital sky survey (DPOSS), in: *Wide Field Surveys in Cosmology (14th IAP meeting)*, Editions Frontieres, Paris, 1998, p. 89.

Appendix A. Proofs of Technical Lemmas

PROOF OF LEMMA 2. By definition of the LTS cost we have

$$(\Delta^*)^2(h-1) = \sum_{j \in H} r_j^2(\beta^*, P) = \sum_{j \in H} (y_j - \beta^* x_j)^2.$$

Let H' denote the indices of the h points of P whose squared residuals with respect to β' are the smallest. Therefore,

$$(\Delta')^2(h-1) = \sum_{j \in H'} r_j^2(\beta', P) = \sum_{j \in H'} (y_j - \beta' x_j)^2 \leq \sum_{j \in H} (y_j - \beta' x_j)^2.$$

Thus, we have

$$(\Delta' - \Delta^*) \sqrt{h-1} \leq \left(\sum_{j \in H} (y_j - \beta' x_j)^2 \right)^{1/2} - \left(\sum_{j \in H} (y_j - \beta^* x_j)^2 \right)^{1/2}.$$

We may interpret $\Delta' \sqrt{h-1}$ as the Euclidean distance between two points in \mathbb{R}^h , particularly (y_1, \dots, y_h) and $(\beta' x_1, \dots, \beta' x_h)$. Similarly, we may interpret $\Delta^* \sqrt{h-1}$ as the Euclidean distance

between (y_1, \dots, y_h) and $(\beta^* \mathbf{x}_1, \dots, \beta^* \mathbf{x}_h)$. Let $\|\mathbf{u} - \mathbf{v}\|$ denote the Euclidean distance between \mathbf{u} and \mathbf{v} . By the triangle inequality and the above inequalities regarding norms, we have

$$\begin{aligned} (\Delta' - \Delta^*) \sqrt{h-1} &\leq \sum_{j \in H} (\|y_j - \beta' \mathbf{x}_j\| - \|y_j - \beta^* \mathbf{x}_j\|) \leq \sum_{j \in H} \|\beta' \mathbf{x}_j - \beta^* \mathbf{x}_j\| \\ &= \left(\sum_{j \in H} (\beta' \mathbf{x}_j - \beta^* \mathbf{x}_j)^2 \right)^{1/2} = \left(\sum_{j \in H} ((\beta' - \beta^*) \mathbf{x}_j)^2 \right)^{1/2} \\ &\leq \left(\sum_{j \in H} \|\beta' - \beta^*\|^2 \cdot \|\mathbf{x}_j\|^2 \right)^{1/2} \\ &\leq \|\beta' - \beta^*\| \left(\sum_{j \in H} \|\mathbf{x}_j\|^2 \right)^{1/2}, \end{aligned}$$

where the second to last step follows from the Cauchy-Schwarz inequality.

By the earlier observations relating β' and β^* to \mathbf{y}_E and \mathbf{y}_E^* , and another application of the Cauchy-Schwarz inequality we obtain

$$\begin{aligned} (\Delta' - \Delta^*) \sqrt{h-1} &\leq \|\mathbf{y}_E X_E^{-1} - \mathbf{y}_E^* X_E^{-1}\| \left(\sum_{j \in H} \|\mathbf{x}_j\|^2 \right)^{1/2} \\ &\leq \|(\mathbf{y}_E - \mathbf{y}_E^*)\| \cdot \|X_E^{-1}\| \left(\sum_{j \in H} \|\mathbf{x}_j\|^2 \right)^{1/2}. \end{aligned}$$

For each $\mathbf{x}_j \in P$, let $y_j^* = \beta^* \mathbf{x}_j$. We can express Δ^* as $\sum_{j \in H} (y_j - y_j^*)^2$, from which it follows that

$$\frac{(\Delta' - \Delta^*)}{\Delta^*} \leq \|\mathbf{y}_E - \mathbf{y}_E^*\| \cdot \|X_E^{-1}\| \left(\sum_{j \in H} \|\mathbf{x}_j\|^2 \right)^{1/2} / \left(\sum_{j \in H} (y_j - y_j^*)^2 \right)^{1/2}.$$

By multiplying and dividing by $\|X_E\|$, we obtain

$$\frac{\Delta'}{\Delta^*} - 1 \leq \left(\frac{\|\mathbf{y}_E - \mathbf{y}_E^*\|^2}{\sum_{j \in H} (y_j - y_j^*)^2} \right)^{1/2} \cdot (\|X_E^{-1}\| \cdot \|X_E\|) \cdot \left(\frac{\sum_{j \in H} \|\mathbf{x}_j\|^2}{\|X_E\|^2} \right)^{1/2},$$

which completes the proof.

PROOF OF LEMMA 4. Observe that (by our general-position assumption) $i-1$ intervals of B lie strictly to the left of $b_{[i]}^+$, and $n - (i+h-1)$ intervals lie strictly to the right of $b_{[i+h-1]}^-$. If $b_{[i]}^+ > b_{[i+h-1]}^-$, let k denote the number of intervals of B that contain the span $[b_{[i+h-1]}^-, b_{[i]}^+]$. Every interval of B is either among the $i-1$ intervals to the left of $b_{[i]}^+$, and/or the $n - (i+h-1)$ intervals that lie to the right of $b_{[i+h-1]}^-$, or the k that contain $[b_{[i+h-1]}^-, b_{[i]}^+]$. Therefore,

$$n \leq (i-1) + (n - (i+h-1)) + k \leq n - h + k,$$

from which we conclude that $k \geq h$. Thus, every point of the span $[b_{[i+h-1]}^-, b_{[i]}^+]$ is contained within h intervals of B . This implies that the distance from any such point to all of its h closest intervals is zero. Therefore, $\Delta(B_i, h) = 0$, and so $\Delta(B, h) = 0$.

On the other hand, if $b_{[i]}^+ \leq b_{[i+h-1]}^-$, then the $i-1$ intervals of B to the left of $b_{[i]}^+$ are disjoint from the $n - (i+h-1)$ intervals to the right of $b_{[i+h-1]}^-$. Since the remaining n intervals of B are in B_i , it follows that $|B_i| = n - (i-1) - (n - (i+h-1)) = h$, as desired.

PROOF OF LEMMA 5. Let B' denote the subset of B of size h that achieves the minimum LTS cost. Suppose towards a contradiction that $\Delta(B', h) < \min_i \Delta(B_i, h)$. Let $b_{[i]}^+$ and $b_{[j]}^-$ denote the leftmost right endpoint and rightmost left endpoint in B' , respectively. Since the span from $b_{[i]}^+$ to $b_{[j]}^-$ must overlap at least h intervals of B , by the same reasoning as in Lemma 4, we have $j \geq i + h - 1$. Further, we may assume that $j > i + h - 1$, for otherwise B' would be equal to B_i . Thus, B' spans at least $h + 1$ intervals, which implies that there exists at least one “unused” interval of B that is in B_i but not in B' . Let I denote any such interval (see Fig. A.19(a)).

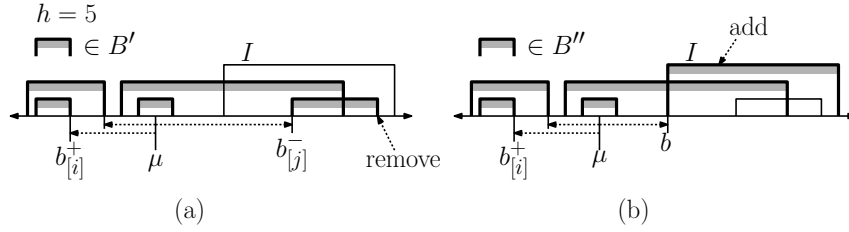


Figure A.19: Proof of Lemma 5. By replacing the interval whose left endpoint is anchored to $b_{[j]}^-$ with the unused interval I anchored to b , we decrease the LTS cost.

First, we assert that $b_{[i]}^+ < b_{[j]}^-$. To see why, observe that otherwise, $b_{[i]}^+$ is contained within $h + 1$ intervals (by the same reasoning as in Lemma 4), implying that $\Delta(B_i, h) = 0$, which would violate our hypothesis that none of the B_i 's is optimal. Given this assertion, let μ denote the point that minimizes the sum of squared distances to the intervals of B' (see Fig. A.19(a)). Clearly, $b_{[i]}^+ \leq \mu \leq b_{[j]}^-$. By left-right symmetry, we may assume without loss of generality that μ is closer to $b_{[i]}^+$ than it is to $b_{[j]}^-$. (More specifically, if we negate all the interval endpoints, we may apply the proof with the roles of $b_{[i]}^+$ and $b_{[j]}^-$ reversed.) Because $b_{[j]}^-$ is the rightmost left endpoint in B' , the distance from μ to $b_{[j]}^-$ is the largest among all the intervals of B whose closest endpoint to μ lies in the span $[b_{[i]}^+, b_{[j]}^-]$, and in particular this is true for the interval I . Consider the set B'' formed by taking B' and replacing the interval anchored at $b_{[j]}^-$ with I (see Fig. A.19(b)). Since I is closer to μ than $b_{[j]}^-$ is, the LTS cost of B'' is smaller than the LTS cost of B' . This contradicts our hypothesis that B' achieves the minimum LTS cost.

Appendix B. Complete Algorithm for the Interval LTS Problem

In this appendix we present the complete algorithm for the interval LTS problem. The input consists of a set B of n nonempty, closed intervals on the real line and a coverage value h . The objective is to determine the point $\mu \in \mathbb{R}$ that minimizes the average squared distances to the closest h intervals and the cost function associated with this point.

Recall from Section 4 the definition of the cost function $\Delta(B, h)$ for a given set B of intervals. Also recall the interval sets B_i and the h -spans I_i , for $1 \leq i \leq n - h + 1$, defined just prior to Lemma 4. Let μ_i denote the point that minimizes the sum of squared distances from itself to all of the intervals of B_i . Our approach will be to compute the initial values, μ_1 and $\Delta(B_1, h)$, explicitly by brute force in $O(h)$ time. Then, for $i = 2, \dots, n - h + 1$, we will show how to update the subsequent values μ_i and $\Delta(B_i, h)$, each in constant time. By Lemma 5, the μ_i yielding the smallest such cost is the desired result.

In the upper bound algorithm of Section 4, we observed that the value that minimizes the sum of squared distances to a set of numbers is just the mean of the set. However, generalizing

this simple observations to a set of intervals is less obvious. Given an arbitrary interval $[b^-, b^+]$ of B , the function that maps any point x on the real line to its squared distance to this interval is clearly convex. (It is not strictly convex because the function is zero for all points that lie within the interval.) By the remarks made after Lemma 4, we may assume that no h intervals of B contain a common point (for otherwise the cost is trivially zero). Since $\Delta(B_i, h)$ is the minimum of the sum of h such functions, and by our assumption that no h intervals of B contain in a single point, it follows that this function is strictly convex. Therefore, μ_i is this function's unique local minimum. Given an initial estimate on the location of μ_i , we can exploit convexity to determine the direction in which to move. As the algorithm transitions from B_{i-1} to B_i , we remove one interval from the left (the one anchored at $b_{[i-1]}^+$) and add one interval to the right (the one anchored at $b_{[i+h-1]}^-$). Since we remove an interval from the left end and add one on the right end, it follows that $\mu_i \geq \mu_{i-1}$. Thus, the search for μ_i proceeds monotonically from left to right.

To make this more formal, let $\{b_1, \dots, b_{2n}\}$ denote the sorted sequence of the interval endpoints of B in nondecreasing order. This implicitly defines $2n - 1$ intervals of the form $[b_{k-1}, b_k]$, for $1 < k \leq 2n$ where μ_i might lie. We refer to these intervals as *segments* (see Fig. B.20(a)). All the points within the interior of any given segment have the property that they lie to the left/right/contained within the same subset of intervals of B . Our algorithm operates by maintaining three indices, i , j , and k . We will set $j = i + h - 1$ so that $I_i = [b_{[i]}^+, b_{[j]}^-]$ is the current h -span of interest. The index k specifies the *current segment*, $[b_k, b_{k+1}]$, which represents the hypothesized segment that contains μ_i . The algorithm repeatedly increases k as long as μ_i is determined to lie to the right of b_{k+1} .

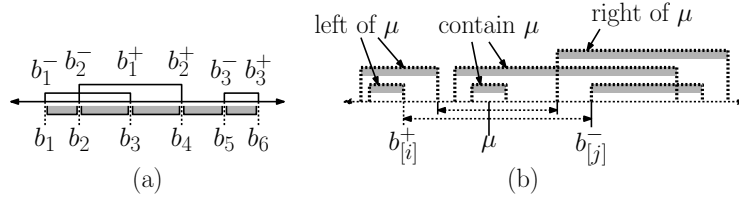


Figure B.20: Segments and contributions: (a) The sorted endpoints $\langle b_1, \dots, b_{2n} \rangle$ and the associated segments (shaded), (b) the mean μ_i and the distance contributions from the intervals of B_i .

How do we determine whether the current segment contains μ_i ? First, we classify the every interval of B_i as lying to the left, right, or containing the segment's interior. The intervals of B_i that contain the segment's interior contribute zero to the squared distance, and so they may be ignored. For those intervals that lie entirely to the segment's left (resp., right), the interval's right (resp., left) endpoint determines the squared distance (see Fig. B.20(b)). Let us call this endpoint the *relevant endpoint* of the associated interval. As in the computation of the upper bound, we will maintain two quantities *SUM* and *SSQ*, which will be updated throughout the algorithm. The first stores the sum of the relevant endpoints and the second stores the sum of squares of the relevant endpoints. As in Eq. (4), the value μ_i and the LTS cost $\Delta(B_i, h) = \sigma_i$ will be derived from these quantities.

Each time we advance the current segment by incrementing k , we are now either entering an interval of B (if b_k is a left endpoint) or leaving an interval (if b_k is a right endpoint). In the first case the interval we are entering was contributing its left endpoint to (*SUM* and *SSQ*) and is now making no contribution (see Fig. B.21(a)). In the second case, the interval we are leaving was

making no contribution and is now contributing its right endpoint (see Fig. B.21(b)).

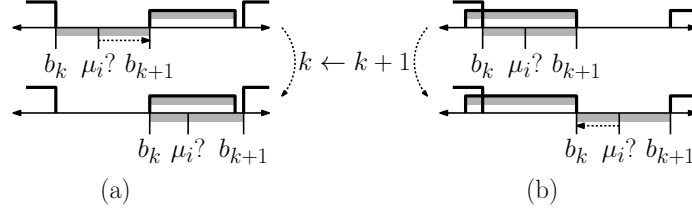


Figure B.21: Contribution updates when advancing the current segment.

We can now present the complete algorithm. Throughout, the algorithm maintains Δ_{best} , which stores the minimum LTS cost of any B_i encountered thus far.

- (1) [Initializations:]
 - (a) [Sort:] Sort the left endpoints of B , denoted $\langle b_{[1]}^-, \dots, b_{[n]}^- \rangle$, the right endpoints, denoted $\langle b_{[1]}^+, \dots, b_{[n]}^+ \rangle$, and all the endpoints, denoted $\langle b_1, \dots, b_{2n} \rangle$.
 - (b) [Initialize costs:] Set $\text{SUM} \leftarrow \sum_{1 \leq j \leq h} b_{[j]}^-$, $\text{SSQ} \leftarrow \sum_{1 \leq j \leq h} (b_{[j]}^-)^2$, and $\Delta_{\text{best}} \leftarrow +\infty$.
 - (c) [Initialize indices:] Let $i \leftarrow k \leftarrow 1$ and $j \leftarrow h$.
- (2) [Main loop:] While $i \leq n - h - 1$ do:
 - (a) [Check for h -fold interval overlap:] If $b^+[i] \geq b^-[j]$, set $\Delta_{\text{best}} \leftarrow 0$ and return (by the remarks made after Lemma 4).
 - (b) [Update μ :] Recalling Eq. (4), set $\mu \leftarrow \text{SUM}/h$.
 - (c) [Find the segment that contains μ :] While $b_{k+1} < \mu$ do:
 - (i) [Advance to the next segment:] Set $k \leftarrow k + 1$.
 - (ii) [Entering an interval?] If b_k is the left endpoint of some interval of B , set $\text{SUM} \leftarrow \text{SUM} - b_k$ and $\text{SSQ} \leftarrow \text{SSQ} - (b_k)^2$ (see Fig. B.21(a)).
 - (iii) [Exiting an interval?] Otherwise, b_k is the right endpoint of some interval of B . Set $\text{SUM} \leftarrow \text{SUM} + b_k$ and $\text{SSQ} \leftarrow \text{SSQ} + (b_k)^2$ (see Fig. B.21(b)).
 - (iv) [Update μ :] Set $\mu \leftarrow \text{SUM}/h$.
 - (d) [Update LTS cost:] Recalling Eq. (4), set $\sigma \leftarrow \left(\frac{\text{SSQ} - (\text{SUM}^2/h)}{h-1} \right)^{1/2}$. If $\sigma < \Delta_{\text{best}}$ then set $\Delta_{\text{best}} \leftarrow \sigma$.
 - (e) [Remove contribution of leftmost interval $b_{[i]}^+$:] $\text{SUM} \leftarrow \text{SUM} - b_{[i]}^+$ and $\text{SSQ} \leftarrow \text{SSQ} - (b_{[i]}^+)^2$.
 - (f) [Advance to next subset:] Set $i \leftarrow i + 1$ and $j \leftarrow j + 1$.
 - (g) [Add contribution of rightmost interval $b_{[j]}^-$:] $\text{SUM} \leftarrow \text{SUM} + b_{[j]}^-$ and $\text{SSQ} \leftarrow \text{SSQ} + (b_{[j]}^-)^2$.
- (3) Return Δ_{best} as the final LTS cost

The algorithm's correctness has already been justified. To establish the algorithm's running time, observe that Step (1a) takes $O(n \log n)$ time, and Step (1b) takes $O(h)$ time. All the other steps take only constant time. Each time the loop of Step (2c) is executed, the value of k is increased. Since k cannot go beyond the last segment (since clearly $\mu \leq b_{2n}$), the total number of iterations of this loop throughout the entire algorithm is at most $2n$. Since the loop of Step (2) is repeated at most $n - h + 1$ times, the overall running time is $O(n \log n + h + 2n + (n - h - 1)) = O(n \log n)$, as desired.