

On perturbations of linear least squares problems

M.R.Osborne *

Abstract

The asymptotic behaviour of a class of least squares problems when subjected to structured perturbations is considered. It is permitted that the number of rows (observations) in the design matrix can be unbounded while the number of degrees of freedom (variables) is fixed. It is shown that for certain classes of random data the solution sensitivity depends asymptotically on the condition number of the design matrix rather than on its square which is the classical result for inconsistent systems. Extension of these results to the case where the perturbations are due to rounding errors is considered. This appears possible provided certain scaling problems can be resolved. These scales are not compatible with worst case perturbation theory.

1 Introduction

The linear least squares problem has the general form

$$\min_{\mathbf{x}} \mathbf{r}^T \mathbf{r}; \quad \mathbf{r} = A\mathbf{x} - \mathbf{b}. \quad (1)$$

where the design matrix $A : R^p \rightarrow R^n$, the residual and observation vectors $\mathbf{r}, \mathbf{b} \in R^n$, and the vector of model parameters $\mathbf{x} \in R^p$. Typically p will be fixed corresponding to a known model, while n will usually be assumed "large enough". Limiting processes will assume that p is fixed and $n \rightarrow \infty$.

*Mathematical Sciences Institute, Australian National University, ACT 0200, AUSTRALIA. <mailto:Mike.Osborne@anu.edu.au>

This paper is based on a presentation given to the meeting honouring the lives of Gene Golub and Ron Mitchell held at the Australian National University on February 28-29, 2008. It is dedicated to the memory of these two good friends.

This problem is a simple optimization problem subject to equality constraints. The necessary conditions for a minimum give

$$0 = \nabla_{\mathbf{x}} \mathbf{r}^T \mathbf{r} = \mathbf{2r}^T A. \quad (2)$$

Substituting for \mathbf{r} from (1) gives the *normal equations*

$$A^T A \mathbf{x} = A^T \mathbf{b}. \quad (3)$$

This system defines uniquely both the least squares estimator $\mathbf{x}^{(n)}$ and the corresponding residual vector $\mathbf{r}^{(n)}$ providing the design matrix A has full column rank p , and this condition is assumed.

An important modelling context which generates linear least squares problems is the following. Assume noisy observations are made on a system at a sequence of configurations labelled by a reference variable t which could be time. Let these be summarised by

$$b_i = y(t_i) + \varepsilon_i, \quad i = 1, 2, \dots, n, \quad (4)$$

where $y(t)$ is the error free signal (true model) which is assumed to be expressible in parametric form as

$$y(t) = \sum_{i=1}^p x_i^* \phi_i(t), \quad (5)$$

the x_i^* , $i = 1, 2, \dots, p$, are the (hypothesised) true parameter values, the $\phi_i(t)$, $i = 1, 2, \dots, p$ are basis functions specifying the model class, and the ε_i are random variables summarising the noise in the observations. A standard assumption would be that the ε_i are independent and normally distributed with mean 0, and standard deviation σ ($\varepsilon \sim N(0, \sigma^2 I)$). For this model

$$A_{ij} = \phi_j(t_i), \quad j = 1, 2, \dots, p, \quad i = 1, 2, \dots, n.$$

It is important to know how the estimate $\mathbf{x}^{(n)}$ of \mathbf{x}^* given by (3) behaves as the number of observations increases without limit because this permits statements to be made about rates of convergence. This is not just a theoretical point because it informs on how much data needs to be collected and so directly relates to the practicality of the measurement exercise. It is also important to know how the computational algorithm chosen to solve (1) will behave on large data sets. In this connection, the first point to make is that a *systematic* process capable of automation is required to generate the values of the reference variable t_i and record the observations b_i associated with large data sets if an asymptotic analysis is to be possible. The nature

of this recording process must depend on the nature of the system being observed. It is assumed that the system has the property that after a finite horizon, assumed to correspond to values $0 \leq t \leq 1$ for the labelling variable, no further information on model structure is available. One case corresponds to signals decaying to zero. However, the case of a finite observation window dictated by external factors is also included. Such systems may be required to be controlled, and may have quite complicated stability properties.

The refinement process used to increase n is one in which independent trials are performed to obtain data for sequences of points $\{t_i^{(n)}, i = 1, 2, \dots, n\}$ for an increasing sequence of values of n . The use of the descriptor systematic is taken to mean that there is a limiting process such that

$$\frac{1}{n} \sum_{i=1}^n f(t_i^{(n)}) \rightarrow \int_0^1 f(t) dw(t), \quad n \rightarrow \infty, \quad (6)$$

holds for all sufficiently smooth $f(t)$ ($f(t) \in C[0, 1]$ for example) where $w(t)$ is a weight function characteristic of the sampling regime. The left hand side in (6) can be interpreted as a simple quadrature formula. For example, $w(t) = t$ in the two cases:

1. The t_i are equispaced. The corresponding quadrature error for smooth enough $f(t)$ is strictly $O(1/n)$.
2. The t_i are uniformly distributed in $[0, 1]$. The corresponding quadrature error is asymptotically normally distributed with variance $O(1/n)$.

Such samplings are called *regular* to stress that there is a sense in which the quadrature error is $o(1)$, $n \rightarrow \infty$. The associated convergence mode is denoted "r.e.". For example, $\xrightarrow[n \rightarrow \infty]{r.e.}$. The particular sense appropriate is implied.

Now assume that the design matrix in (1) is constructed for each n using a regular sampling procedure. Then the regularity condition gives

$$\begin{aligned} \frac{1}{n} A_{*i}^T A_{*j} &= \frac{1}{n} \sum_{k=1}^n \phi_i(t_k^{(n)}) \phi_j(t_k^{(n)}) \\ &\xrightarrow[n \rightarrow \infty]{r.e.} \int_0^1 \phi_i(t) \phi_j(t) dw(t) = G_{ij}. \end{aligned} \quad (7)$$

This states that the normal matrix in (3) scaled by $\frac{1}{n}$ approaches the Gram matrix $G : R^p \rightarrow R^p$ of the set $\Phi = \{\phi_j(t), j = 1, 2, \dots, p\}$ relative to the weight $w(t)$ with an error which is $o(1)$, $n \rightarrow \infty$. This rate is assumed fast

enough to ensure $\text{cond } A \rightarrow \text{cond } G$, $n \rightarrow \infty$. G is nonsingular and positive definite by assumption. It need not be well conditioned. For example, values of spectral condition numbers for Hilbert matrices corresponding to the case when elements of Φ are monomials and the t_i are equispaced are given in [4] for $2 \leq n \leq 16$. One consequence of the above discussion is that there is no real restriction in assuming that the subordinate matrix norm relative to the euclidean norm of the design matrix satisfies $\|A\| = \sqrt{n}$. This amounts to a rescaling of the design A by a quantity which is asymptotically constant.

The next section treats some consequences of perturbing the data of equation (1). The classic inequality of Golub and Wilkinson is derived and certain asymptotic properties for large n are explored. In particular, the influence of the stochastic components in the data vector \mathbf{b} is considered. This adds a somewhat different perspective to the usual worst case scenarios because here the law of large numbers [7] in the (extended) form

$$\frac{1}{n} \sum_{i=1}^n X_{ni} \varepsilon_{ni} \xrightarrow[n \rightarrow \infty]{a.s.} 0 \quad (8)$$

is available when the ε_{ni} are independent and of bounded variance for all n , and the constants X_{ni} are bounded. This result permits the influence of the “bad term” involving the square of the condition number of A to be ignored in certain circumstances. The resulting perturbation behaviour then becomes similar to that for consistent linear systems.

2 Perturbation of least squares problems

We consider the generic perturbed least squares problem (1) with data

$$\mathbf{r} = (A + \tau E) \mathbf{x} - (\mathbf{b} + \tau \mathbf{z})$$

where perturbations E , \mathbf{z} are fixed in the sense that they result from a well defined rule for each n . The perturbation E is assumed to be independent of any observational error. It is assumed that τ is a small parameter. The component-wise scale of the perturbations is fixed by requiring

$$\max_{i,j} |E_{ij}| = \eta \leq 1, \quad \|\mathbf{z}\|_{\infty} \leq 1. \quad (9)$$

It is assumed that τ is small enough for both A and $A + \tau E$ to have their full rank p . The necessary conditions for the perturbed and unperturbed least squares problems are

$$(A + \tau E)^T \hat{\mathbf{r}} = 0, \quad A^T \mathbf{r}^{(n)} = 0$$

where the $\hat{\cdot}$ indicates the solution of the perturbed problem. Subtracting gives

$$(A + \tau E)^T (\hat{\mathbf{r}} - \mathbf{r}^{(n)}) + \tau E^T \mathbf{r}^{(n)} = 0,$$

and substituting for the residual vectors gives the basic relation

$$(A + \tau E)^T (A + \tau E) (\hat{\mathbf{x}} - \mathbf{x}^{(n)}) = \tau \left\{ (A + \tau E)^T (\mathbf{z} - E\mathbf{x}^{(n)}) - E^T \mathbf{r}^{(n)} \right\}. \quad (10)$$

For small enough τ and each fixed n this gives

$$\begin{aligned} \hat{\mathbf{x}} - \mathbf{x}^{(n)} &= \tau \left\{ (A^T A)^{-1} (A^T (\mathbf{z} - E\mathbf{x}^{(n)}) - E^T \mathbf{r}^{(n)}) \right\} + O(\tau^2), \\ &= \tau \left\{ \begin{array}{l} \left(\frac{1}{\sqrt{n}} U \right)^{-1} \frac{1}{\sqrt{n}} Q_1^T (\mathbf{z} - E\mathbf{x}^{(n)}) \\ - \left(\frac{1}{n} A^T A \right)^{-1} \frac{1}{n} E^T \mathbf{r}^{(n)} \end{array} \right\} + O(\tau^2), \end{aligned} \quad (11)$$

where A possesses the orthogonal Q times upper triangular U factorization

$$A = Q \begin{bmatrix} U \\ 0 \end{bmatrix} = [Q_1 \quad Q_2] \begin{bmatrix} U \\ 0 \end{bmatrix} = Q_1 U,$$

and Q_1 corresponds to the first p columns of Q . There are two ways of looking at this relation. The first considers n fixed and worries about the size of $\text{cond}(A) = \frac{\sigma_1}{\sigma_p}$, the ratio of the largest to smallest singular values of A . This leads to the basic inequality

$$\|\hat{\mathbf{x}} - \mathbf{x}^{(n)}\| \leq \tau \left\{ \frac{\text{cond}(A)}{\sqrt{n}} \|\mathbf{z} - E\mathbf{x}^{(n)}\| + \frac{\text{cond}(A)^2}{n} \|E^T \mathbf{r}^{(n)}\| \right\} + O(\tau^2), \quad (12)$$

where the assumption that $\|A\| = \sqrt{n} = \sigma_1$ has been used. The original form of this result is due to [3]. Equation (12) reveals the possible dominance of the term $\text{cond}(A)^2$. This is likely if $\frac{1}{n} \|E^T \mathbf{r}^{(n)}\|$ is not small. The importance of the inequality (12) is that it is a generic result highlighting what is best possible. For this reason computational algorithms in which the error takes this form are said to have *optimal error structure*. Development of such optimal algorithms based on the use of orthogonal transformations goes back to [5], [2], and [1]. It follows from (11) that

$$\begin{aligned} \hat{\mathbf{r}} - \mathbf{r}^{(n)} &= -\tau \left\{ (I - P) \mathbf{z} + PE\mathbf{x}^{(n)} + A (A^T A)^{-1} E^T \mathbf{r}^{(n)} \right\} + O(\tau^2), \\ &= -\tau \left\{ (I - P) \mathbf{z} + PE\mathbf{x}^{(n)} + Q_1 U^{-1} E^T \mathbf{r}^{(n)} \right\} + O(\tau^2) \end{aligned} \quad (13)$$

where P is the orthogonal projection $A (A^T A)^{-1} A^T$ onto the range of A . Thus the result of the perturbation is a change of magnitude $O(\text{cond}(A))$

in the residual showing that a more satisfactory result is possible if the computed residual is the required quantity.

However, there is an alternative way of considering this result which is important when n is large and $\boldsymbol{\varepsilon}$ is a random vector. The Gram matrix G (7) is used to write a limiting form of (11) as $n \xrightarrow{r.e.} \infty$. Contributions from quadrature error terms in this approximation have been ignored (Lemma 2 shows they contribute at most $\tau(o(1))$ for large n given regular sampling), and $G^{1/2}$ is written for the large n approximation to $\frac{1}{\sqrt{n}}U$ in (11).

$$\widehat{\mathbf{x}} - \mathbf{x}^{(n)} = \tau \left\{ G^{-1/2} \frac{1}{\sqrt{n}} Q_1^T (\mathbf{z} - E\mathbf{x}^{(n)}) - G^{-1} \frac{1}{n} E^T \mathbf{r}^{(n)} \right\} + O(\tau(o(1)), \tau^2). \quad (14)$$

We have the following bounds for the interesting terms in this equation.

Lemma 1

$$\begin{aligned} \frac{1}{\sqrt{n}} \|Q_1^T (\mathbf{z} - E\mathbf{x}^{(n)})\| &\leq \|\mathbf{z} - E\mathbf{x}^{(n)}\|_\infty, \\ \frac{1}{n} \|E^T \mathbf{r}^{(n)}\| &\leq \sqrt{\frac{p}{n}} \eta \|\mathbf{r}^{(n)}\|_\infty. \end{aligned}$$

Proof. The first part applies the standard inequality relating the 2 and ∞ norms. The second part follows in similar fashion from the inequality

$$\|E\| \leq \sqrt{np}\eta, \quad (15)$$

where the right hand side is a simple bound for the Frobenius norm of A . ■

Also it is important when fixing the order of dependence on τ that the n dependence of the $O(\tau)$ terms is appropriately bounded as $n \xrightarrow{r.e.} \infty$. The key is the following result.

Lemma 2

$$\frac{1}{n} (A + \tau E)^T (A + \tau E) \xrightarrow[n \rightarrow \infty]{r.e.} G + O(\tau).$$

It follows that the normal matrix associated with the perturbed least squares problem has a suitably bounded inverse under regular sampling.

Proof.

$$\begin{aligned} \frac{1}{n} (A + \tau E)^T (A + \tau E) &= \frac{1}{n} U^T \{ I + \tau \{ Q_1^T E U^{-1} + U^{-T} E^T Q_1 \} \\ &\quad + \tau^2 U^{-T} E^T E U^{-1} \} U. \end{aligned} \quad (16)$$

To show that the terms multiplying both τ and τ^2 in this expression are $O(1)$, $n \xrightarrow{r.e.} \infty$, requires a bound for $\|EU^{-1}\|$ valid for large n . Note $\|EU^{-1}\| \geq \|Q_1^T EU^{-1}\|$. The required bound can be constructed as follows:

$$\begin{aligned} \|U^{-T} E^T E U^{-1}\| &= \sup_{\mathbf{v}} \frac{\mathbf{v}^T U^{-T} E^T E U^{-1} \mathbf{v}}{\mathbf{v}^T \mathbf{v}}, \\ &= \sup_{\mathbf{w}} \frac{\mathbf{w}^T E^T E \mathbf{w}}{\mathbf{w}^T U^T U \mathbf{w}}, \\ &\leq \frac{\|E\|^2}{n \sigma_{\min} \left\{ \frac{1}{n} A^T A \right\}}, \\ &\leq \frac{p\eta^2}{\sigma_{\min} \{G\}} + o(1), \quad n \xrightarrow{r.e.} \infty, \end{aligned}$$

where the estimate of $\|E\|$ given in the previous Lemma has been used. Thus

$$\left\| \frac{1}{n} (A + \tau E)^T (A + \tau E) - G \right\| \leq \tau \left\{ 3 \|G\|^{1/2} \sqrt{p \text{cond}(G)} + o(1) \right\}, \quad n \xrightarrow{r.e.} \infty.$$

The last step uses $\tau^2 \|EU^{-1}\|^2 \leq \tau \|EU^{-1}\|$ when $\tau \|EU^{-1}\| \leq 1$. ■

Remark 3 *This result has the consequence that all terms in the basic relation (14) have the orders claimed as $n \xrightarrow{r.e.} \infty$. More can be said if the law of large numbers (8) can be applied to estimate $E^T \mathbf{r}^{(n)}$. It follows from the necessary conditions that*

$$A^T \mathbf{r}^{(n)} = 0 \Rightarrow Q_2 Q_2^T \mathbf{r}^{(n)} = \mathbf{r}^{(n)}.$$

This means that the necessary conditions give

$$\mathbf{r}^{(n)} = Q_2 Q_2^T (A (\mathbf{x}^{(n)} - \mathbf{x}^*) - \boldsymbol{\varepsilon}) = -Q_2 Q_2^T \boldsymbol{\varepsilon},$$

so that

$$\begin{aligned} \frac{1}{n} E^T \mathbf{r}^{(n)} &= -\frac{1}{n} E^T Q_2 Q_2^T \boldsymbol{\varepsilon} \xrightarrow[n \rightarrow \infty]{a.s.} 0, \\ &\Rightarrow \frac{1}{n} G^{-1} E^T \mathbf{r}^{(n)} \xrightarrow[n \rightarrow \infty]{a.s.} 0. \end{aligned} \tag{17}$$

by the law of large numbers. This provides a sense in which the term in $G^{-1/2}$ dominates in (14) for large n .

3 What about rounding error?

There are several important distinctions between the development in the previous section and the asymptotic properties of the important class of problems in which E , \mathbf{z} , τ are determined by computational procedure and the characteristics of floating point arithmetic.

1. The scale τ is determined by the requirement that the component-wise scaling conditions (9) are satisfied. If these are set by reference to worst case error analysis (for example, [4], Theorem 19.3), then this gives $\tau(n) = \gamma_n u$ where u is unit roundoff and $\gamma_n = O(n)$. This is not compatible with the previous asymptotic results.
2. The application of the law of large numbers requires that ε be independent of E , \mathbf{z} . This cannot be strictly true here as right hand side values must influence rounding behaviour.
3. The values of E , \mathbf{z} depend on the detail of the particular algorithm implemented.

Putting aside the setting of τ for the moment, some progress can be made on the other matters. Consider the Golub orthogonal factorization algorithm based on Householder transformations [2]. A suitable form of error analysis is given in [6]. This shows that the potential $\text{cond}(A)^2$ comes from a term

$$\Delta = U^{-1} \delta Q \mathbf{r}^{(n)} = \left(\frac{1}{\sqrt{n}} U \right)^{-1} \delta Q \frac{1}{\sqrt{n}} \mathbf{r}^{(n)},$$

where δQ is the error in the computed orthogonal transformation. The computation of the factorization matrix Q does not involve the problem right hand side so the rounding error/stochastic error interactions can only contribute to potential $\text{cond}(A)$ terms. Now Δ can be estimated using the law of large numbers provided the individual elements of δQ have an $O\left(\frac{1}{\sqrt{n}}\right)$ estimate, a magnitude typical of the elements of an $n \times n$ orthogonal matrix.

This suggests that the analysis of the preceding section can be applied here provided τ is small. This requires systematic cancellation not allowed for in setting $\tau = \gamma_n u$ based on worst case analysis. There is more hope from informal observations which would seem to suggest that the cumulative effects of rounding errors prove relatively small in large scale computations. This could indicate something like a weak-mixing form of a law of large numbers (weak mixing because there is certainly some local rounding error interaction). Such a law need not depend on the precise statistics of individual rounding errors, and could be compatible with worst case error analysis

in the sense of allowing certain exceptional cases by analogy with almost sure convergence.

4 acknowledgement

The author is appreciative of the interest shown by the referees and by Nick Higham. Their comments have led to important improvements.

References

- [1] Å. BJÖRCK, *Solving linear least squares problems by Gram-Schmidt orthogonalization*, BIT, 7 (1967), pp. 1–22.
- [2] G. H. GOLUB, *Numerical methods for solving least squares problems*, Num. Math., 7 (1965), pp. 206–216.
- [3] G. H. GOLUB AND J. H. WILKINSON, *Iterative refinement of least squares solutions*, Num. Math., 9 (1966), pp. 189–198.
- [4] N. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, SIAM, 2002. Second Edition.
- [5] A. HOUSEHOLDER, *Unitary triangularization of a nonsymmetric matrix*, J. ACM, 6 (1958), pp. 339–342.
- [6] L. S. JENNINGS AND M. R. OSBORNE, *A direct error analysis for least squares*, Num. Math., 22 (1974), pp. 325–332.
- [7] K. SEN AND J. SINGER, *Large Sample Methods in Statistics*, Chapman and Hall, 1993.