

Assessing the Accuracy of Outlier Tests in Nonlinear Regression

Myung Wook Kahng^{1,a}, Bu-Yong Kim^a

^aDept. of Statistics, Sookmyung Women's Univ.

Abstract

Given the specific mean shift outlier model, the standard approaches to obtaining test statistics for outliers are discussed. Accuracy of outlier tests is investigated using subset curvatures. These subset curvatures appear to be reliable indicators of the adequacy of the linearization based test. Also, we consider obtaining graphical summaries of uncertainty in estimating parameters through confidence curves. The results are applied to the problem of assessing the accuracy of outlier tests.

Keywords: Mean shift outlier model, outlier test, curvature measure, confidence curves.

1. Introduction

A standard paradigm for the analysis of nonlinear models is to assume that results for the linear model hold at least approximately for large enough sample sizes. Beginning with Beale (1960) and more importantly with the papers of Bates and Watts (1980, 1981), it has become increasingly clear that this approximation may not be at all adequate in some problems. Researchers have developed methods both to determine when linear approximations are inadequate and to give alternative procedures when the standard methods fail.

Section 2 considers the problem of testing outliers in nonlinear regression. Given the specific mean shift model, standard approaches to obtaining test statistics for outliers are discussed. The test based on linear approximation, namely the score test, is easy to calculate, but quite different from likelihood ratio test. Accuracy of the linear approximation is investigated using subset curvature measures in Section 3.

We consider obtaining graphical summaries of uncertainty in estimation in nonlinear regression models using confidence curves in Section 4. The results of confidence curves are applied not only to the problem of finding significance levels for outlier tests, but also to assessing the accuracy of the test based on linear approximation. We provide an example in Section 5.

2. Outlier Tests in Nonlinear Regression

The standard nonlinear regression model can be expressed as

$$y_i = f(x_i, \theta) + \epsilon_i, \quad i = 1, \dots, n$$

in which the i^{th} response y_i is related to the q -dimensional vector of known explanatory variables x_i through the known model function f , which depends on the p -dimensional unknown parameter $\theta \in \Theta$

This research was supported by the Sookmyung Women's University Research Grants 2007.

¹ Corresponding author: Professor, Department of Statistics, Sookmyung Women's University, Seoul 140-742, Korea.
E-mail: mwkahng@sm.ac.kr

and ϵ_i is error. We assume that f is twice continuously differentiable in θ and errors ϵ_i are *i.i.d.* normal random variables with mean 0 and variance σ^2 . In matrix notation we may write,

$$\mathbf{y} = \mathbf{f}(\mathbf{X}, \theta) + \epsilon, \quad (2.1)$$

where \mathbf{y} is an n -dimensional vector with elements y_1, \dots, y_n , \mathbf{X} is an $n \times q$ matrix with rows $\mathbf{x}_1^T, \dots, \mathbf{x}_n^T$, ϵ is an n -dimensional vector with elements $\epsilon_1, \dots, \epsilon_n$ and $\mathbf{f}(\mathbf{X}, \theta) = (f(x_1, \theta), \dots, f(x_n, \theta))^T = (f_1(\theta), \dots, f_n(\theta))^T = \mathbf{f}(\theta)$. Given the response vector \mathbf{y} , the least squares estimate of θ is denoted $\hat{\theta}$, the predicted response vector is $\hat{\mathbf{y}} = \mathbf{f}(\mathbf{X}, \hat{\theta}) = \mathbf{f}(\hat{\theta})$ and the residual vector is $\mathbf{e} = \mathbf{y} - \mathbf{f}(\mathbf{X}, \hat{\theta}) = \{e_i\}$. A tangent plane approximation to the expectation surface at $\hat{\theta}$ is used to make inferences about θ through the derived linear model $\mathbf{f}(\theta) = \mathbf{f}(\hat{\theta}) + \hat{\mathbf{V}}(\theta - \hat{\theta})$, where $\mathbf{V} = \mathbf{V}(\theta) = \partial \mathbf{f} / \partial \theta^T$ is the $n \times p$ matrix and $\hat{\mathbf{V}} = \mathbf{V}(\hat{\theta})$ is $\partial \mathbf{f} / \partial \theta^T$ evaluated at $\hat{\theta}$.

Suppose we suspect in advance that m cases indexed by an m -dimensional vector $\mathbf{I} = (i_1, \dots, i_m)$ are outliers. A useful framework used to study outliers is the mean shift outlier model,

$$\mathbf{y} = \mathbf{f}(\mathbf{X}, \theta) + \mathbf{D}\delta + \epsilon, \quad (2.2)$$

where $\delta = (\delta_{i_1}, \dots, \delta_{i_m})^T$, $\mathbf{D} = (\mathbf{d}_1, \dots, \mathbf{d}_m)$ and \mathbf{d}_j is the i_j^{th} standard basis vector for \mathbf{R}^n . The testing of the hypothesis $\delta = \mathbf{0}$ is equivalent to testing whether cases in the set \mathbf{I} are outliers.

We denote the log-likelihood for model (2.2) by $L(\theta, \delta, \sigma^2)$ and obtain

$$\begin{aligned} L(\theta, \delta, \sigma^2) &= -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{f}(\mathbf{X}, \theta) - \mathbf{D}\delta)^T (\mathbf{y} - \mathbf{f}(\mathbf{X}, \theta) - \mathbf{D}\delta) \\ &= -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} S(\theta, \delta), \end{aligned} \quad (2.3)$$

where $S(\theta, \delta) = (\mathbf{y} - \mathbf{f}(\mathbf{X}, \theta) - \mathbf{D}\delta)^T (\mathbf{y} - \mathbf{f}(\mathbf{X}, \theta) - \mathbf{D}\delta)$. Given σ^2 , (2.3) is maximized with respect to $\phi = (\theta, \delta)$ when $S(\theta, \delta)$ is minimized at the least squares estimates $\hat{\phi} = (\hat{\theta}_{(I)}, \hat{\delta})$. Furthermore, $\partial L / \partial \sigma^2 = 0$ has solution $\sigma^2 = S(\theta, \delta) / n$, which gives a maximum for given ϕ as the second derivative is negative. This suggests that $\hat{\phi} = (\hat{\theta}_{(I)}, \hat{\delta})$ and $\hat{\sigma}_{(I)}^2 = S(\hat{\theta}_{(I)}, \hat{\delta}) / n$ are the maximum likelihood estimates. Under the null hypothesis $\delta = \mathbf{0}$, the maximum likelihood estimates are $\phi_0 = (\hat{\theta}, \mathbf{0})$ and $\hat{\sigma}^2 = S(\hat{\theta}, \mathbf{0}) / n$, which are the maximum likelihood estimates of model (2.1).

Now, we consider procedures for testing $H_0 : \delta = \mathbf{0}$ against $H_1 : \delta \neq \mathbf{0}$. Three asymptotically equivalent test methods are discussed by Kahng (1995), namely the likelihood ratio (LR) test introduced by Neyman and Pearson (1928), Wald (WD) test proposed by Wald (1943) and the score (S) test due originally to Rao (1947) and developed further by Silvey (1959). The three test statistics are defined as follows:

$$\begin{aligned} \text{LD} &= 2 \{L(\hat{\phi}) - L(\phi_0)\} \\ &= n \{ \log S(\hat{\theta}, \mathbf{0}) - \log S(\hat{\theta}_{(I)}, \hat{\delta}) \}, \\ \text{WD} &= \hat{\delta}^T \{ \text{Var}(\hat{\delta}) \}^{-1} \hat{\delta} \\ &= \frac{1}{\hat{\sigma}_{(I)}^2} \hat{\delta}^T (\mathbf{I} - \hat{\mathbf{H}}_I) \hat{\delta}, \\ S &= \mathbf{U}(\phi_0)^T \mathbf{I}(\phi_0)^{-1} \mathbf{U}(\phi_0) \\ &= \frac{1}{\hat{\sigma}^2} \mathbf{e}_I^T (\mathbf{I} - \hat{\mathbf{H}}_I)^{-1} \mathbf{e}_I, \end{aligned}$$

where \tilde{H}_I is $m \times m$ minor of $\tilde{H} = \tilde{V}(\tilde{V}^T \tilde{V})^{-1} \tilde{V}^T$ with rows and columns indexed by \mathbf{I} , $\tilde{V} = V(\hat{\theta}_{(I)})$ is $\partial f / \partial \theta^T$ evaluated at $\hat{\theta}_{(I)}$, $\mathbf{U}(\phi) = \partial L(\phi) / \partial \phi$, $\mathbf{I}(\phi_0) = \partial^2 L(\phi) / \partial \phi \partial \phi^T$, \hat{H}_I is the $m \times m$ minor of $\hat{H} = \hat{V}(\hat{V}^T \hat{V})^{-1} \hat{V}^T$ with rows and columns indexed by \mathbf{I} and \mathbf{e}_I is m -dimensional vector whose j^{th} element is e_{ij} .

Suppose that only the i^{th} case is suspected as being an outlier. In this single outlier case, the outlier model (2.2) becomes

$$\mathbf{y} = \mathbf{f}(\mathbf{X}, \boldsymbol{\theta}) + \delta \mathbf{d}_i + \boldsymbol{\epsilon}, \quad (2.4)$$

where \mathbf{d}_i is the i^{th} standard basis vector for \mathbf{R}^n . The three statistics take the form

$$\begin{aligned} \text{LD} &= n \left\{ \log S(\hat{\boldsymbol{\theta}}, 0) - \log S(\hat{\boldsymbol{\theta}}_{(i)}, \hat{\delta}) \right\}, \\ \text{WD} &= \frac{\hat{\delta}^2 (1 - \tilde{h}_{ii})}{\hat{\sigma}_{(i)}^2}, \\ S &= \frac{e_i^2}{\hat{\sigma}^2 (1 - \hat{h}_{ii})}, \end{aligned}$$

where \tilde{h}_{ii} and \hat{h}_{ii} are the i^{th} diagonal elements of \tilde{H} and \hat{H} , respectively. The three statistics differ in computational features. Unlike the other two tests, the score test requires only quantities calculated under the null hypothesis. Nevertheless, all three statistics are invariant under the reparametrization of $\boldsymbol{\theta}$ and have the same asymptotic distribution under the null hypothesis $H_0 : \delta = 0$.

3. Accuracy of Linear Approximation

Let's consider the mean shift outlier model (2.4). Since parameter δ is the parameter subset of this model, the subset curvatures for δ can be found by the methods described in Cook and Goldberg (1986).

Let $\mathbf{h}(\delta) = \mathbf{f}(\mathbf{X}, \tilde{\boldsymbol{\theta}}(\delta)) + \delta \mathbf{d}_i$, where $\tilde{\boldsymbol{\theta}}(\delta)$ denotes the p -dimensional vector-valued function that maximizes $L(\boldsymbol{\theta}, \delta, \sigma^2)$ over $\boldsymbol{\theta}$ given δ . To obtain precise expressions, we need the following definitions. Define $n \times p \times p$ array $\tilde{W} = W(\hat{\boldsymbol{\theta}}_{(i)})$ which is $\partial^2 f / \partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T$ evaluated at $\hat{\boldsymbol{\theta}}_{(i)}$. We consider the QR decomposition of $n \times (p+1)$ matrix V_ϕ , namely

$$V_\phi = \left. \frac{\partial \mathbf{f}}{\partial \boldsymbol{\phi}^T} \right|_{\phi=\hat{\phi}} = (\tilde{V}, \mathbf{d}_i) = \mathbf{Q}\mathbf{R},$$

where \mathbf{Q} is an $n \times (p+1)$ matrix with orthogonal columns and \mathbf{R} is a $(p+1) \times (p+1)$ upper triangular matrix. Now consider the transformation $\tilde{W}_\phi = \mathbf{R}^{-T} \tilde{W}_\phi \mathbf{R}^{-1}$, where the i^{th} face of \tilde{W}_ϕ is

$$(\tilde{W}_\phi)_i = \left. \frac{\partial^2 h_i}{\partial \boldsymbol{\phi} \partial \boldsymbol{\phi}^T} \right|_{\phi=\hat{\phi}} = \begin{bmatrix} \tilde{W}_i & \mathbf{0} \\ \mathbf{0}^T & 0 \end{bmatrix},$$

\tilde{W}_i is the i^{th} face of \tilde{W} and h_i is the i^{th} element of $\mathbf{h}(\delta)$.

Under the above settings and by applying the result of Kahng (1995), we may define the maximum parameter-effects and intrinsic curvatures of \mathbf{h} at $\hat{\delta}$ as

$$\begin{aligned} \Gamma_s^r(\delta) &= \hat{\sigma}_{(i)} \max_{\|\mathbf{b}\|=1} \left\| \mathbf{b}^T \mathbf{A}_{22} \mathbf{b} \right\|, \\ \Gamma_s^l(\delta) &= 2\hat{\sigma}_{(i)} \max_{\|\mathbf{b}\|=1} \left\| [\mathbf{b}^T] [\mathbf{A}_{12}] \mathbf{b} \right\|, \end{aligned}$$

where A is the $(p + 1) \times (p + 1) \times (p + 1)$ parameter-effects curvature array $A = [Q^T][\ddot{W}_\phi]$ (Bates and Watts, 1981), A_{12} and A_{22} are the subarrays of A .

If $\Gamma_s^r(\delta)$ and $\Gamma_s^q(\delta)$ are sufficiently small, the likelihood and linear confidence regions for δ will be similar, otherwise we can expect these confidence regions to be dissimilar. Following Ratkowsky (1983, p.18) and Cook and Goldberg (1986), $1/[2\sqrt{F_\alpha(1, n - p - 1)}]$ may be used as a rough guide for judging the size of these curvatures. This method can be used to judge the adequacy of the test procedures which are based on the linear approximation. When $\Gamma_s^r(\delta)$ and $\Gamma_s^q(\delta)$ are greater than the guide, the linearization based test, namely the score test, is quite different from the likelihood ratio test.

4. Confidence Curves for Parameter δ

Cook and Weisberg (1990) give the graphical alternative to likelihood and Wald confidence intervals for a component of the parameter vector. The reason for using confidence curves is that likelihood intervals can have a different shape for each significance level $1 - \alpha$. Wald intervals are always symmetric, so if we have a 95% interval, we can always obtain the 90% and 99% intervals in a simple way. With the likelihood intervals, this is not so. Consequently, a graphical summary that does not depend on level is desirable.

The confidence curves include two curves - a likelihood confidence curve and Wald confidence curve. Applying the result of Kahng (2003), the likelihood confidence curves are the set of points defined by

$$\begin{cases} \sqrt{\{S(\tilde{\theta}(\delta), \delta) - S(\hat{\theta}_{(i)}, \hat{\delta})\}/s_{(i)}^2}, & \text{on the horizontal axis,} \\ \delta, & \text{on the vertical axis,} \end{cases} \quad (4.1)$$

where $s_{(i)}^2 = S(\hat{\theta}_{(i)}, \hat{\delta})/(n - p - 1)$. This is a modification of the graphical summary of the standard profile log-likelihood which has δ on the horizontal axis and $S(\tilde{\theta}(\delta), \delta)$ on the vertical axis. The plot (4.1) will be curves, with the amount of curvature giving information about the nonlinearity of the model. To this plot, two straight lines passing through $(0, \hat{\delta})$ with slope $\pm se(\hat{\delta})$ are added. These two straight lines represent the Wald interval. At any point on the horizontal axis of the confidence curve plot, the interval between the upper and lower curves provides a confidence interval for δ for some level of $1 - \alpha$. The confidence level can be determined from the calibrating distribution for either the Wald or likelihood procedure, which, in the scale of the plot, is $t(\nu)$, where ν is the corresponding degrees of freedom. If $t_\nu(u)^{-1}$ is the inverse of the t cumulative distribution function with ν degrees of freedom evaluated at u , then the confidence level at a point along the horizontal axis is $1 - 2t_\nu^{-1}\left[\sqrt{\{S(\tilde{\theta}(\delta), \delta) - S(\hat{\theta}_{(i)}, \hat{\delta})\}/s_{(i)}^2}\right]$.

The Wald and the likelihood regions are tangent at the maximum likelihood estimate $\hat{\phi} = (\hat{\theta}_{(i)}, \hat{\delta})$. If the likelihood is exactly quadratic, that is, if the log-likelihood is exactly normal, then the likelihood curves are the same as the Wald curves (two straight lines). Non-normality is reflected in the likelihood confidence curves failing to be straight.

5. Example

The data on the metabolism of tetracycline were presented in Bates and Watts (1988) and are reproduced in Table 1. In this experiment, a tetracycline compound was administered orally to a subject

Table 1: Tetracycline data

| i | x_i | y_i |
|-----|-------|-------|
| 1 | 1 | 0.7 |
| 2 | 2 | 1.2 |
| 3 | 3 | 1.4 |
| 4 | 4 | 1.4 |
| 5 | 6 | 1.1 |
| 6 | 8 | 0.8 |
| 7 | 10 | 0.6 |
| 8 | 12 | 0.5 |
| 9 | 16 | 0.3 |

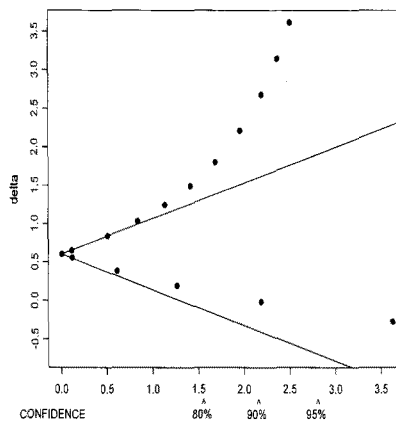


Figure 1: Confidence curves for δ

and the concentration(y_i) of tetracycline hydrochloride in the serum in micrograms per milliliter was measured over a period(x_i) of 16 hours.

We consider the outlier problem. Suppose that case 1 is suspected to be an outlier. The proposed model is the following:

$$f(x, \theta, \delta) = \theta_3[\exp\{-\theta_1(x - \theta_4)\} - \exp\{-\theta_2(x - \theta_4)\}].$$

The subset curvatures for δ are $\Gamma_s^r(\delta) = .5470$, $\Gamma_s^l(\delta) = .3653$. The corresponding guide is $1/[2\sqrt{F_{.05}(1, 4)}] = .180$, indicating inadequacy of the linear approximation.

Now we consider the confidence curves discussed in Section 4. The confidence curves for δ is given in Figure 1. In this figure, the likelihood confidence curves are severely curved and dissimilar to the Wald confidence curves. This indicates that estimations, inferences and diagnostics based on the linear approximation about δ may give misleading results.

References

Bates, D. M. and Watts, D. G. (1980). Relative curvature measures of nonlinearity (with discussion), *Journal of the Royal Statistical Society, Series B*, **42**, 1–25.
 Bates, D. M. and Watts, D. G. (1981). Parameter transformations for improved approximate confidence regions in nonlinear least squares, *The Annals of Statistics*, **9**, 1152–1167.

- Bates, D. M. and Watts, D. G. (1988). *Nonlinear Regression Analysis and Its Applications*, John Wiley & Sons, New York.
- Beale, E. M. L. (1960). Confidence regions in nonlinear estimation (with discussion), *Journal of the Royal Statistical Society, Series, B*, **22**, 41–76.
- Cook, R. D. and Goldberg, M. L. (1986). Curvatures for parameter subsets in nonlinear regression, *The Annals of Statistics*, **14**, 1399–1418.
- Cook, R. D. and Weisberg, S. (1990). Confidence curves in nonlinear regression, *Journal of the American Statistical Association*, **85**, 544–551.
- Kahng, M. W. (1995). Testing outliers in nonlinear regression, *Journal of the Korean Statistical Society*, **24**, 419–437.
- Kahng, M. W. (2003). Confidence curves for a function of parameters in nonlinear regression, *Journal of the Korean Statistical Society*, **32**, 1–10.
- Neyman, J. and Pearson, E. S. (1928). On the use and interpretation of certain test criteria for purposes of statistical inference, *Biometrika*, **20A**, 175–240 and 263–294.
- Rao, C. R. (1947). Large sample tests of statistical hypotheses concerning several parameters with applications to problems of estimation, *Proceedings of the Cambridge Philosophical Society*, **44**, 50–57.
- Ratkowsky, D. A. (1983). *Nonlinear Regression Modeling: A Unified Practical Approach*, Marcel Dekker, New York.
- Silvey, S. D. (1959). The Lagrangian multiplier test, *The Annals of Mathematical Statistics*, **30**, 389–407.
- Wald, A. (1943). Tests of statistical hypotheses concerning several parameters when the number of observations is large, *Transactions of the American Mathematical Society*, **54**, 426–482.

Received October 2008; Accepted November 2008