Asymptotically Optimal Regression Prediction Intervals and Prediction Regions for Multivariate Data

David J. Olive¹

¹ Department of Mathematics, Southern Illinois University, Carbondale, IL, USA

Correspondence: David J. Olive, Department of Mathematics, Southern Illinois University, Carbondale, IL., USA. Tel: 1-618-453-6566. E-mail: dolive@siu.edu

Received: October 24, 2012 Accepted: January 13, 2013 Online Published: January 22, 2013 doi:10.5539/ijsp.v2n1p90 URL: http://dx.doi.org/10.5539/ijsp.v2n1p90

The research is financed by National Science Foundation grant DMS 0600933

Abstract

This paper presents asymptotically optimal prediction intervals and prediction regions. The prediction intervals are for a future response Y_f given a $p \times 1$ vector \mathbf{x}_f of predictors when the regression model has the form $Y_i = m(\mathbf{x}_i) + e_i$ where *m* is a function of \mathbf{x}_i and the errors e_i are iid from a continuous unimodal distribution. The prediction intervals have coverage near or higher than the nominal coverage for many techniques even for moderate sample size *n*, say n > 10(model degrees of freedom). The prediction regions are for a future vector of measurements \mathbf{x}_f from a multivariate distribution. The nonparametric prediction region developed in this paper has correct asymptotic coverage if the data $\mathbf{x}_1, ..., \mathbf{x}_n$ are iid from a distribution with a nonsingular covariance matrix. For many distributions, this prediction region appears to have good coverage for n > 20p, and this region is asymptotically optimal on a large class of elliptically contoured distributions. Hence the prediction intervals and regions perform well for moderate sample sizes as well as asymptotically.

Keywords: additive models, nonlinear regression, prediction intervals, prediction regions, regression

1. Introduction

This paper presents asymptotically optimal prediction intervals and prediction regions. The prediction regions are for a future vector of measurements x_f from a multivariate distribution, and are asymptotically optimal on a large class of elliptically contoured distributions. Regression is the study of the conditional distribution Y|x of the response Y given the $p \times 1$ vector of predictors x. The prediction intervals are for a future response Y_f given a vector x_f of predictors when the regression model has the form

$$Y_i = m(\boldsymbol{x}_i) + e_i \tag{1}$$

for i = 1, ..., n where *m* is a function of x_i and the errors e_i are iid from a continuous unimodal distribution. Many of the most important regression models have this form, including the multiple linear regression model and many time series, nonlinear, nonparametric and semiparametric models. If \hat{m} is an estimator of *m*, then the *i*th residual is $r_i = Y_i - \hat{m}(x_i) = Y_i - \hat{Y}_i$.

Olive (2007) showed how to form asymptotically optimal prediction intervals for model (1), but for many regression models and estimators, large n is needed for the intervals to perform well. Prediction intervals derived for multiple linear regression did perform well. This paper derives asymptotically optimal prediction intervals that perform well for many models for moderate n.

A large sample $100(1 - \delta)\%$ prediction interval (PI) has the form (\hat{L}_n, \hat{U}_n) where $P(\hat{L}_n < Y_f < \hat{U}_n) \xrightarrow{P} 1 - \delta$ as the sample size $n \to \infty$. Following Olive (2007), let ξ_{δ} be the δ percentile of the error e, i.e., $P(e \le \xi_{\delta}) = \delta$. Let $\hat{\xi}_{\delta}$ be the sample δ percentile of the residuals. Consider predicting a future observation Y_f given a vector of predictors \mathbf{x}_f where (Y_f, \mathbf{x}_f) comes from the same population as the past data (Y_i, \mathbf{x}_i) for i = 1, ..., n. Let $1 - \delta_2 - \delta_1 = 1 - \delta$ with $0 < \delta < 1$ and $\delta_1 < 1 - \delta_2$ where $0 < \delta_i < 1$. Then $P[Y_f \in (m(\mathbf{x}_f) + \xi_{\delta_1}, m(\mathbf{x}_f) + \xi_{1-\delta_2})] = 1 - \delta$.

Assume that \hat{m} is consistent: $\hat{m}(\mathbf{x}) \xrightarrow{P} m(\mathbf{x})$ as $n \to \infty$. Then $r_i = Y_i - \hat{m}(\mathbf{x}_i) \xrightarrow{P} Y_i - m(\mathbf{x}_i) = e_i$ and, under "mild"

regularity conditions, $\hat{\xi}_{\delta} \xrightarrow{P} \xi_{\delta}$. If $a_n \xrightarrow{P} 1$ and $b_n \xrightarrow{P} 1$, then

$$(\hat{L}_n, \hat{U}_n) = (\hat{m}(\mathbf{x}_f) + a_n \hat{\xi}_{\delta_1}, \hat{m}(\mathbf{x}_f) + b_n \hat{\xi}_{1-\delta_2})$$
(2)

is a large sample $100(1 - \delta)\%$ PI for Y_f .

According to regression folklore, the percentiles of the residuals are consistent estimators, $\hat{\xi}_{\delta} \xrightarrow{P} \xi_{\delta}$, under "mild" regularity conditions, and this consistency is the basis for using QQ plots. The folklore is true for linear models: sufficient conditions are $\hat{\beta} \xrightarrow{P} \beta$ and the x_i are bounded in probability. See Olive and Hawkins (2003), Welsh (1986) and Rousseeuw and Leroy (1987, p. 128).

Consider the multiple linear regression model $Y = X\beta + e$ where Y is an $n \times 1$ vector of dependent variables, X is an $n \times p$ matrix of predictors, β is a $p \times 1$ vector of unknown coefficients, and e is an $n \times 1$ vector of unknown iid zero mean errors e_i with variance σ^2 . Let the hat matrix $H = X(X^TX)^{-1}X^T$. Let $h_i = h_{ii}$ be the *i*th diagonal element of H for i = 1, ..., n. Then h_i is called the *i*th *leverage* and $h_i = x_i^T(X^TX)^{-1}x_i$. Suppose new data is to be collected with predictor vector x_f . Then the leverage of x_f is $h_f = x_f^T(X^TX)^{-1}x_f$.

For the multiple linear regression model, let $\hat{\xi}_{\delta}$ be the sample quantile of the residuals. Following Olive (2007), let

$$a_n = b_n = \left(1 + \frac{15}{n}\right) \sqrt{\frac{n}{n-p}} \sqrt{(1+h_f)}.$$
 (3)

Then a large sample semiparametric $100(1 - \delta)\%$ PI for Y_f is

$$(\hat{Y}_f + a_n \hat{\xi}_{\delta/2}, \hat{Y}_f + a_n \hat{\xi}_{1-\delta/2}).$$
 (4)

A PI is asymptotically optimal if it has the shortest asymptotic length that gives the desired asymptotic coverage. The PI (4) is asymptotically optimal on a large class of unimodal continuous symmetric error distributions. For more general distributions, an asymptotically optimal PI can be created by applying the shorth(*c*) estimator to the residuals where $c = \lceil n(1 - \delta) \rceil$ and $\lceil x \rceil$ is the smallest integer $\ge x$, e.g., $\lceil 7.7 \rceil = 8$. See Grübel (1988). That is, let $r_{(1)}, ..., r_{(n)}$ be the order statistics of the residuals. Compute $r_{(c)} - r_{(1)}, r_{(c+1)} - r_{(2)}, ..., r_{(n)} - r_{(n-c+1)}$. Let $(r_{(d)}, r_{(d+c-1)}) = (\tilde{\xi}_{\delta_1}, \tilde{\xi}_{1-\delta_2})$ correspond to the interval with the smallest length. Following Olive (2007), a 100 $(1 - \delta)\%$ PI for Y_f is

$$(\hat{Y}_f + a_n \tilde{\xi}_{\delta_1}, \hat{Y}_f + a_n \tilde{\xi}_{1-\delta_2}) \tag{5}$$

where a_n is given by (3). This prediction interval performs well for moderate n for multiple linear regression and several estimators, including least squares.

A problem with prediction intervals is choosing a_n and b_n so that the intervals have short length and coverage close to or higher than the nominal coverage for a wide variety of regression models when n is moderate. Section 2.1 shows how to modify (4) and (5) to achieve these goals while Section 2.2 covers prediction regions for a future vector of measurements x_f . Examples and simulations are in Section 3.

2. Method

The idea for finding the asymptotically optimal prediction intervals and regions is simple. Find the target population $100(1 - \delta)\%$ covering region. For small *n*, the coverage of the training data will be higher than that for the future case to be predicted. In simulations for a large group of models and distributions, the undercoverage could be as high as min(0.05, $\delta/2$). Let $q_n = \min(1 - \delta + 0.05, 1 - \delta + p/n)$ for $\delta > 0.1$ and

$$q_n = \min(1 - \delta/2, 1 - \delta + 10\delta p/n), \text{ otherwise.}$$
(6)

If $1 - \delta < 0.999$ and $q_n < 1 - \delta + 0.001$, set $q_n = 1 - \delta$. Then use the prediction interval or region that covers $100q_n\%$ of the training data. The coverage of the training data is $100q_n\%$ and converges to $100(1 - \delta)\%$ as $n \to \infty$, even if the model assumptions fail to hold.

2.1 Asymptotically Optimal Prediction Intervals

The technique used to produce asymptotically optimal PIs that perform well for moderate samples is simple. Find \hat{Y}_f and the residuals from the regression model. Since the leverage of x_i is closely related to the Mahalanobis distance of x_i from the sample mean \bar{x} of the *n* predictor vectors, leverage and extrapolation are useful for a wide

range of regression models. For a wide range of regression models, extrapolation occurs if $h_f > 2p/n$: if x_f is too far from the data $x_1, ..., x_n$, then the model may not hold and prediction can be arbitrarily bad. This result suggests replacing (3) by

$$a_n = b_n = \left(1 + \frac{15}{n}\right) \sqrt{\frac{n+2p}{n-p}}.$$
(7)

Let $\delta_n = 1 - q_n$ where q_n is given by (6). Then

$$(\hat{L}_n, \hat{U}_n) = (\hat{m}(\mathbf{x}_f) + b_n \hat{\xi}_{\delta_n/2}, \hat{m}(\mathbf{x}_f) + b_n \hat{\xi}_{1-\delta_n/2})$$
(8)

is a large sample $100(1 - \delta)$ % PI for Y_f that is similar to (2) and (4).

Let $c = \lceil nq_n \rceil$. Compute $r_{(c)} - r_{(1)}, r_{(c+1)} - r_{(2)}, ..., r_{(n)} - r_{(n-c+1)}$. Let $(r_{(d)}, r_{(d+c-1)}) = (\tilde{\xi}_{\delta_1}, \tilde{\xi}_{1-\delta_2})$ correspond to the interval with the smallest length. Then the asymptotically optimal 100 $(1 - \delta)$ % large sample PI for Y_f is

$$(\hat{m}(\boldsymbol{x}_f) + b_n \tilde{\xi}_{\delta_1}, \hat{m}(\boldsymbol{x}_f) + b_n \tilde{\xi}_{1-\delta_2}),$$
(9)

and is similar to (5).

To see that the PI (9) is asymptotically optimal, assume that the sample percentiles of the residuals converge to the population percentiles of the iid unimodal errors: $\hat{\xi}_{\delta} \xrightarrow{P} \xi_{\delta}$. Also assume that the population shorth $(\xi_{\delta_1}, \xi_{1-\delta_2})$ is unique and has length *L*. Since $b_n \rightarrow 1$, $\hat{m}(\mathbf{x}_f) \xrightarrow{P} m(\mathbf{x}_f)$, and $q_n = 1 - \delta$ for large enough *n*, it is enough to show that the shorth of the residuals converges to the population shorth of the $e_i: (\xi_{\delta_1}, \xi_{1-\delta_2}) \xrightarrow{P} (\xi_{\delta_1}, \xi_{1-\delta_2})$. Let L_n be the length of $(\xi_{\delta_1}, \xi_{1-\delta_2})$. Let $0 < \tau < 1$ and $0 < \epsilon < L$ be arbitrary. Assume *n* is large enough so that $q_n = 1 - \delta$. Then $P(L_n > L + \epsilon) \rightarrow 0$ since $(\hat{\xi}_{\delta_1}, \hat{\xi}_{1-\delta_2})$ covers 100 $(1 - \delta)$ % of the data and $L_n = \xi_{1-\delta_2} - \xi_{\delta_1} \leq \xi_{1-\delta_2} - \hat{\xi}_{\delta_1} \xrightarrow{P} L$ as $n \rightarrow \infty$ since the sample percentiles are consistent and the shorth is the smallest interval covering 100 $(1 - \delta)$ % of the data. If $P(L_n < L - \epsilon) > \tau$ eventually, then the shorth is an interval covering 100 $(1 - \delta)$ % of the cases that is shorter than the population shorth with positive probability τ . Hence at least one of $\hat{\xi}_{1-\delta_2}$ or $\hat{\xi}_{\delta_1}$ would not converge, a contradiction. Since ϵ and τ were arbitrary, $L_n \xrightarrow{P} L = \xi_{1-\delta_2} - \xi_{\delta_1} - \epsilon) > \tau$ eventually, then $P(\tilde{\xi}_{1-\delta_2} < \xi_{1-\delta_2} - \epsilon/2) > \tau$ eventually since $L_n = \tilde{\xi}_{1-\delta_2} - \tilde{\xi}_{\delta_1} \xrightarrow{P} L = \xi_{1-\delta_2} - \xi_{\delta_1}$. But such an interval (of length going to *L* in probability with left endpoint less than $\xi_{\delta_1} - \epsilon$ and right endpoint less than $\xi_{\delta_1} - \epsilon/2$ occurs interval covering $100(1 - \delta)$ % of the cases that is shorter than the poblability going to one since the population shorth is the unique shortest interval covering $100(1 - \delta)$ % of the cases that is shorter than the shorth is the population shorth is the unique shortest interval covering $100(1 - \delta)$ % of the cases with probability going to one since the population shorth is the unique shortest interval covering $100(1 - \delta)$ % of the cases tha

The above results show that PI (9) and the shorth of the residuals behave well when the sample percentiles are consistent. Even if these assumptions do not hold, the PI covers $100q_n\%$ of the training data, and often the coverage of the future case will be close to $100(1 - \delta)$ if the future case Y_f is similar to the training data.

For asymptotic optimality, can not have extrapolation. Also, even if the coverage converges to the nominal coverage, the length of the PI need not be asymptotically shortest unless the highest $1-\delta$ density region of the probability density function of the iid errors is an interval. The highest density region is an interval for unimodal distributions, but need not be an interval for multimodal distributions for all δ . Also see Cai, Tian, Solomon and Wei (2008).

Notice that the technique computes a PI for coverage $q_n \ge 1 - \delta$ which converges to the nominal coverage $1 - \delta$ as $n \to \infty$. Suppose $n \le 20p$. Then the nominal 95% PI uses $q_n = 0.975$ while the nominal 50% PI uses $q_n = 0.55$. Prediction distributions depend both on the error distribution and on the variability of the estimator \hat{m} . This variability is typically unknown but converges to 0 as $n \to \infty$. Also, residuals tend to underestimate the errors for small *n*. For small *n*, ignoring estimator variability and using $q_n = 1 - \delta$ resulted in undercoverage as high as min(0.05, $\delta/2$). Letting the "coverage" q_n decrease to the nominal coverage $1 - \delta$ inflates the length of the PI for small *n*, compensating for the unknown variability of \hat{m} .

The geometry of the "asymptotically optimal prediction region" is simple. The region is the area between two parallel lines with unit slope. Consider a plot of $m(x_i)$ versus Y_i on the vertical axis. The identity line with zero intercept and unit slope is $E(Y_i) = m(x_i)$. Let (L_i, U_i) be the asymptotically optimal population 95% prediction interval containing $m(x_i)$. For example, if the errors are iid $N(0, \sigma^2)$, then $Y_i|m(x_i) \sim N(m(x_i), \sigma^2)$, and $(L_i, U_i) =$

 $(m(\mathbf{x}_i) - 1.96\sigma, m(\mathbf{x}_i) + 1.96\sigma)$. Then the upper line has unit slope and passes through $(m(\mathbf{x}_i), U_i)$ while the lower line has unit slope and passes through $(m(\mathbf{x}_i), L_i)$.

The geometry of the "prediction region" for PI (9) is a natural sample analog of the population "asymptotically optimal prediction region". A response plot of $\hat{Y}_i = \hat{m}(x_i)$ versus Y_i has identity line $\hat{E}(Y_i) = \hat{m}(x_i)$. The region corresponding to pointwise prediction intervals is between two lines with unit slope passing through the points $(\hat{m}(x_i), \hat{U}_i)$ and $(\hat{m}(x_i), \hat{L}_i)$, respectively, where (\hat{L}_i, \hat{U}_i) is the asymptotically optimal prediction interval (9) for Y_f if $x_f = x_i$. For the multiple linear regression model, expect the points in the response plot to scatter in an evenly populated band for n > 5p. Other regression models, such as additive models, may need a much larger sample size n. See Section 3.1 for an example and simulations.

2.2 Prediction Regions

Asymptotically optimal prediction regions use ideas similar to those in the previous subsection. Some notation is needed. Let the *i*th case x_i be a $p \times 1$ random vector, and suppose the *n* cases are collected in an $n \times p$ matrix X with rows $x_1^T, ..., x_n^T$.

The classical estimator (\bar{x}, S) of multivariate location and dispersion is the sample mean and sample covariance matrix where

$$\overline{\boldsymbol{x}} = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{x}_i \text{ and } \boldsymbol{S} = \frac{1}{n-1} \sum_{i=1}^{n} (\boldsymbol{x}_i - \overline{\boldsymbol{x}}) (\boldsymbol{x}_i - \overline{\boldsymbol{x}})^{\mathrm{T}}.$$
 (10)

Some important joint distributions for x are completely specified by a $p \times 1$ population *location* vector μ and a $p \times p$ symmetric positive definite population *dispersion* matrix Σ . An important model is the elliptically contoured $EC_p(\mu, \Sigma, g)$ distribution with probability density function $f(z) = k_p |\Sigma|^{-1/2} g[(z - \mu)^T \Sigma^{-1} (z - \mu)]$ where $k_p > 0$ is some constant and g is some known function. The multivariate normal (MVN) $N_p(\mu, \Sigma)$ distribution is a special case.

Let the $p \times 1$ column vector T(X) be a multivariate location estimator, and let the $p \times p$ symmetric positive definite matrix C(X) be a dispersion estimator. Then the *i*th *squared sample Mahalanobis distance* is the scalar

$$D_i^2 = D_i^2(T(X), C(X)) = (x_i - T(X))^T C^{-1}(X)(x_i - T(X))$$
(11)

for each observation x_i . Notice that the Euclidean distance of x_i from the estimate of center T(X) is $D_i(T(X), I_p)$ where I_p is the $p \times p$ identity matrix. Often the data X will be suppressed. Then the classical Mahalanobis distance uses $(T, C) = (\bar{x}, S)$. Following Johnson (1987, pp. 107-108), the population squared Mahalanobis distance

$$U \equiv D^{2}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (\boldsymbol{x} - \boldsymbol{\mu})^{T} \boldsymbol{\Sigma}^{-1} (\boldsymbol{x} - \boldsymbol{\mu}),$$
(12)

and for elliptically contoured distributions, U has probability density function (pdf)

$$h(u) = \frac{\pi^{p/2}}{\Gamma(p/2)} k_p u^{p/2-1} g(u).$$
(13)

The volume of the hyperellipsoid

$$\{z: (z-\overline{x})^T S^{-1}(z-\overline{x}) \le h^2\} \text{ is equal to } \frac{2\pi^{p/2}}{p\Gamma(p/2)} h^p \sqrt{det(S)}, \tag{14}$$

see Johnson and Wichern (1988, pp. 103-104).

Note that if (T, C) is a \sqrt{n} consistent estimator of $(\mu, d \Sigma)$, then

$$D^{2}(T, C) = (x - T)^{T} C^{-1}(x - T) = (x - \mu + \mu - T)^{T} [C^{-1} - d^{-1} \Sigma^{-1} + d^{-1} \Sigma^{-1}](x - \mu + \mu - T)$$
$$= d^{-1} D^{2}(\mu, \Sigma) + O_{P}(n^{-1/2}).$$

Thus the sample percentiles of $D_i^2(T, C)$ are consistent estimators of the percentiles of $d^{-1}D^2(\mu, \Sigma)$. For multivariate normal data, $D^2(\mu, \Sigma) \sim \chi_p^2$.

Suppose $(T, C) = (\overline{x}_M, b S_M)$ is the sample mean and scaled sample covariance matrix applied to some subset of the data. For h > 0, the hyperellipsoid

$$\{z: (z-T)^T C^{-1}(z-T) \le h^2\} = \{z: D_z^2 \le h^2\} = \{z: D_z \le h\}$$
(15)

has volume equal to

$$\frac{2\pi^{p/2}}{p\Gamma(p/2)}h^p \sqrt{det(C)} = \frac{2\pi^{p/2}}{p\Gamma(p/2)}h^p b^{p/2} \sqrt{det(S_M)}$$
(16)

by (14). A future observation (random vector) \mathbf{x}_f is in region (15) if $D_{\mathbf{x}_f} \leq h$.

A large sample $(1 - \delta)100\%$ prediction region is a set \mathcal{A}_n such that $P(\mathbf{x}_f \in \mathcal{A}_n) \xrightarrow{P} 1 - \delta$. Let q_n be given by (6).

If (T, C) is a consistent estimator of $(\mu, d\Sigma)$, then (15) is a large sample $(1 - \delta)100\%$ prediction region if $h = D_{(up)}$ where $D_{(up)}$ is the q_n th sample quantile of the D_i . If $\mathbf{x}_1, ..., \mathbf{x}_n$ and \mathbf{x}_f are iid, then region (15) is asymptotically optimal on a large class of elliptically contoured distributions in that its volume converges in probability to the volume of the minimum volume covering region $\{\mathbf{z} : (\mathbf{z} - \mu)^T \Sigma^{-1} (\mathbf{z} - \mu) \le u_{1-\delta}\}$ where $P(U \le u_{1-\delta}) = 1 - \delta$ and U has pdf given by (13). The classical parametric multivariate normal large sample prediction region uses $D_{\mathbf{x}_f}(\overline{\mathbf{x}}, \mathbf{S}) \le \sqrt{\chi_{p,1-\delta}^2}$.

Notice that for the data $x_1, ..., x_n$, if C^{-1} exists, then $100q_n\%$ of the *n* cases are in the prediction region, and $q_n \rightarrow 1-\delta$ even if (T, C) is not a good estimator. Hence the coverage q_n of the data is robust to model assumptions. Of course the volume of the prediction region could be large if a poor estimator (T, C) is used or if the x_i do not come from an elliptically contoured distribution. Also notice that $q_n = 1 - \delta/2$ or $q_n = 1 - \delta + 0.05$ for $n \le 20p$ and $q_n \rightarrow 1 - \delta$ as $n \rightarrow \infty$. If $q_n \equiv 1 - \delta$, then (15) is a large sample prediction region, but taking q_n given by (6) improves the finite sample performance of the region. Taking $q_n \equiv 1 - \delta$ does not take into account variability of (T, C), and for small *n* the resulting prediction region tended to have undercoverage as high as min(0.05, $\alpha/2$). Using (6) helped reduce undercoverage for small *n* due to the unknown variability of (T, C).

The Olive and Hawkins (2010) RMVN estimator (T_{RMVN} , C_{RMVN}) is an easily computed \sqrt{n} consistent estimator of (μ , $c\Sigma$) under regularity conditions (E1) that include a large class of elliptically contoured distributions, and c = 1 for the $N_p(\mu, \Sigma)$ distribution. Also see Zhang, Olive and Ye (2012). The RMVN estimator also gives a useful estimate of (μ , Σ) for $N_p(\mu, \Sigma)$ data even when certain types of outliers are present.

Three new prediction regions will be considered. The nonparametric region uses the classical estimator $(T, C) = (\overline{x}, S)$ and $h = D_{(up)}$. The semiparametric region uses $(T, C) = (T_{RMVN}, C_{RMVN})$ and $h = D_{(up)}$. The parametric MVN region uses $(T, C) = (T_{RMVN}, C_{RMVN})$ and $h^2 = \chi^2_{p,q_n}$ where $P(W \le \chi^2_{p,q_n}) = q_n$ if $W \sim \chi^2_p$. All three regions are asymptotically optimal for $N_p(\mu, \Sigma)$ distributions with nonsingular Σ . The first two regions are asymptotically optimal for a large class of elliptically contoured distributions. For distributions with nonsingular covariance matrix $c_X \Sigma$, the nonparametric region is a large sample $(1 - \delta)100\%$ prediction region, but regions with smaller volume may exist. See Section 3.2 for examples and simulations.

3. Results

3.1 Regression



Figure 1. Pointwise prediction interval bands for Ozone data

Example 1 Chambers and Hastie (1993, pp. 251, 516) examine an environmental study that measured the four variables Y = ozone concentration, $x_1 = solar$ radiation, $x_2 = temperature$, and $x_3 = wind$ speed for n = 111

consecutive days. Figure 1 shows the response plot made in *Splus* with the pointwise large sample 95% PI bands for the additive model $Y = m(\mathbf{x}) + e$ where the additive predictor $m(\mathbf{x}) = \alpha + \sum_{j=1}^{3} S_j(x_j)$ for some functions S_j to be estimated. Here $\hat{m}(\mathbf{x}) =$ estimated additive predictor (EAP). Note that the plotted points scatter about the identity line in a roughly evenly populated band, and that 3 of the 111 PIs (9) corresponding to the observed data do not contain *Y*.

A small simulation study compares the PI lengths and coverages for sample sizes n = 50, 100 and 1000 for PIs (8) and (9). Values for PI (8) were denoted by scov and slen while values for PI (9) were denoted by ocov and olen. The five error distributions in the simulation were 1) N(0,1), 2) t_3 , 3) exponential(1) -1, 4) uniform(-1, 1) and 5) 0.9N(0, 1) + 0.1N(0, 100). The value $n = \infty$ gives the asymptotic coverages and lengths and does not depend on the model. So these values are same for multiple linear and nonlinear regression as well as additive models.

Software for the simulations is described in Section 4. The multiple linear regression model with $E(Y_i) = 1 + x_{i1} + \cdots + x_{i7}$ was used. The vectors $(x_1, ..., x_7)^T$ were iid $N_7(0, I_7)$ where I_p is the $p \times p$ identity matrix. Another regression model was $Y_i = m(\mathbf{x}_i) + e_i$, $E(Y_i) = m(\mathbf{x}_i) = \beta_1 x_{i1} + \beta_2 x_{i1}^2 + \beta_3 x_{i2} + \beta_4 x_{i2}^2 + \beta_5 x_{i3} + \beta_6 x_{i3}^2$. This model was fit as an additive model in x_1, x_2 , and x_3 . The model was also fit with nonlinear regression where the mean function is known up to the six parameters, although then the second order multiple linear regression model is appropriate. For the additive model, the additive predictor $m(\mathbf{x}_i) = \alpha + \sum_{j=1}^3 S_j(x_{ij})$. Both the nonlinear regression and additive model had the same mean function $m(\mathbf{x}_i) = x_{i1} + x_{i1}^2$. Thus $\boldsymbol{\beta} = (1, 1, 0, 0, 0, 0)^T$, $\alpha = 0$, $S_1(x_{i1}) = x_{i1} + x_{i1}^2$, $S_2(x_{i2}) = 0$ and $S_3(x_{i3}) = 0$. For these two models, the vectors $(x_1, x_2, x_3)^T$ were iid $N_3(0, I_3)$.

The Olive (2007) PIs (4) and (5) are tailored for multiple linear regression but are liberal (too short) for moderate n for many other techniques. The new PIs (8) and (9) are meant to have coverage near or higher than the nominal coverage for moderate n and for a wide variety of techniques and are longer than PIs (4) and (5). For multiple linear regression, the new PIs (8) and (9) were conservative (too long with roughly 98% coverage for the 95% PI and 70% or 60% coverage for the 50% PI) for n = 50 and 100 compared to (4) and (5) for least squares, least absolute deviations L_1 and an *M*-estimator using the *Splus* functions llfit and rreg. See MathSoft (1999, pp. 293-295).

error		95%	PI	95%	PI	50%	PI	50%	PI
type	n	slen	olen	scov	ocov	slen	olen	scov	ocov
1	50	5.126	4.998	0.959	0.950	1.862	1.674	0.596	0.520
1	100	4.691	4.515	0.968	0.957	1.662	1.528	0.570	0.516
1	1000	3.994	3.944	0.954	0.949	1.379	1.351	0.514	0.505
1	∞	3.920	3.920	0.95	0.950	1.349	1.349	0.50	0.50
2	50	9.444	8.630	0.951	0.943	2.385	2.153	0.576	0.512
2	100	8.245	7.596	0.962	0.954	2.042	1.878	0.577	0.532
2	1000	6.523	6.388	0.950	0.946	1.584	1.553	0.499	0.489
2	∞	6.365	6.365	0.950	0.950	1.530	1.530	0.50	0.50
3	50	5.186	4.823	0.958	0.948	1.573	1.275	0.611	0.526
3	100	4.677	4.156	0.967	0.955	1.382	1.063	0.603	0.533
3	1000	3.771	3.227	0.954	0.952	1.112	0.774	0.509	0.512
3	∞	3.664	2.996	0.950	0.950	1.099	0.693	0.50	0.50
4	50	2.634	2.598	0.961	0.958	1.237	1.087	0.593	0.506
4	100	2.318	2.272	0.972	0.968	1.155	1.028	0.561	0.480
4	1000	1.936	1.926	0.959	0.954	1.014	0.969	0.499	0.486
4	∞	1.900	1.900	0.950	0.950	1.00	1.00	0.50	0.50
5	50	19.689	17.747	0.944	0.935	2.976	2.693	0.608	0.548
5	100	18.754	16.230	0.955	0.946	2.352	2.164	0.580	0.534
5	1000	13.855	12.930	0.946	0.943	1.602	1.569	0.510	0.504
5	∞	13.490	13.490	0.950	0.950	1.507	1.507	0.50	0.50

Table 1. PIs for additive models

The PIs (8) and (9) for nonlinear regression and additive models appear to have coverage near the nominal values in the simulations. For n = 50 and 100, the PIs for nonlinear regression were usually roughly 10% longer than those for additive models. The PIs for the additive model were computed using the *R* function gam. See Hastie

and Tibshirani (1990) and Wood (2006). The PI (8) is not asymptotically optimal with error type 3. It is not known whether \hat{m} is a consistent estimator of *m*, but the prediction intervals appear to have the correct asymptotic coverage and length. Some consistency results for the additive model and models of the form Y = m(x) + e where *m* is smooth are given in Müller, Schick and Wefelmeyer (2012) and Wang, Liu, Liang and Carroll (2011).

The simulation used 5000 runs and gave the proportion \hat{p} of runs where Y_f fell within the nominal $100(1 - \delta)\%$ PI. The count $m\hat{p}$ has a binomial $(m = 5000, p = 1 - \tau_n)$ distribution where $1 - \tau_n$ converges to the asymptotic coverage $(1 - \tau)$. The standard error for the proportion is $\sqrt{\hat{p}(1 - \hat{p})/5000} = 0.0031$ and 0.0071 for p = 0.05 and 0.5, respectively. Hence an observed coverage $\hat{p} \in (.941, .959)$ for 95% and $\hat{p} \in (.479, .521)$ for 50% PIs suggests that there is no reason to doubt that the PI has the nominal coverage.

Table 1 shows that for n = 1000, the coverages and lengths are near the asymptotic $n = \infty$ values. For the 95% PI (9), the coverages were in or near (.94, .96) while the 50% PI (9) was sometimes slightly conservative. The coverage for the 50% PI (8) was near 60% for n = 50. PI (9) is recommended since its asymptotic optimality does not depend on the symmetry of the error distribution.

3.2 Prediction Regions

Rousseeuw and Van Driessen (1999) introduce the DD plot of the classical Mahalanobis distances MD versus the robust distances RD. Olive (2002) shows that if consistent estimators are used and n is large, then the plotted points will follow the identity line with unit slope and zero intercept if the data distribution is multivariate normal, and the plotted points will follow some other line through the origin if the data distribution is from a large class of elliptically contoured distributions but not multivariate normal.

Example 2 Buxton (1920) gives five measurements on 87 men: *height, head length, nasal height, bigonal breadth* and *cephalic index*. The 5 outliers have *heights* that were recorded to be about 19mm and head lengths recorded as the heights. The DD plot of the classical Mahalanobis distances MD versus the RMVN distances RD can be used to visualize the prediction regions. Figure 2 shows the DD plot where points to the left of the vertical line are in the nonparametric large sample 90% prediction region. Points below the horizontal line are in the semiparametric region. The horizontal line at RD = 3.33 corresponding to the parametric MVN 90% region is obscured by the identity line. This region contains 78 of the cases. Since n = 87, the nonparametric and semiparametric regions used the 95th quantile. Since there were 5 outliers, this quantile was a linear combination of the largest clean distance and the smallest outlier distance. The semiparametric 90% region blows up unless the outlier proportion is small.

Figure 3 shows the DD plot and 3 prediction regions after the 5 outliers were removed. The classical and robust distances cluster about the identity line and the three regions are similar, with the parametric MVN region cutoff again at 3.33, slightly below the semiparametric region cutoff of 3.44.



Figure 2. Prediction regions for Buxton data



Figure 3. Prediction regions for Buxton data without outliers

Example 3 Cook and Weisberg (1999, pp. 351, 433, 447) give a data set on 82 mussels sampled off the coast of New Zealand. The variables are $X_1 = \log(S)$, $X_2 = \log(M)$, $X_3 = L$, $X_4 = \log(W)$, and $X_5 = height$ where S is the *shell mass*, M is the *muscle mass* in grams, L is the *length* L, W is the *shell width* and H is the *height* of the shell in mm. Figure 4 shows a DD plot of the data with multivariate prediction regions added. This plot suggests that the data may come from an elliptically contoured distribution that is not multivariate normal. The semiparametric and nonparametric 90% prediction regions consist of the cases below the RD = 5.86 line and to the left of the MD = 4.41 line. These two lines intersect on a line through the origin that is followed by the plotted points. The parametric MVN prediction region is given by the points below the RD = 3.33 line and does not contain enough cases. Points to the left of a vertical line MD = 3.33 would give a modified classical MVN prediction region. Parametric prediction regions for multivariate normal data tend to have severe undercoverage if the data is not multivariate normal. This undercoverage problem becomes worse as p increases, since if the cutoff h is too small, then the volume of the prediction region depends on h^p by (14).



Figure 4. DD plot of the Mussels data

Simulations for the prediction regions used $\mathbf{x} = A\mathbf{w}$ where $\mathbf{A} = diag(\sqrt{1}, \sqrt{2}, ..., \sqrt{p}), \mathbf{w} \sim N_p(\mathbf{0}, \mathbf{I}_p), \mathbf{w} \sim LN(\mathbf{0}, \mathbf{I}_p)$ where the marginals are iid lognormal(0,1), or $\mathbf{w} \sim MVT_p(1)$, a multivariate t distribution with 1 degree of freedom so the marginals are iid Cauchy(0,1). All simulations used 5000 runs and $\delta = 0.1$.

Table 2. Coverages for 90% Prediction Regions

w dist	n	р	ncov	scov	mcov	voln	volm
MVN	600	30	0.906	0.919	0.902	0.503	0.512
MVN	1500	30	0.899	0.899	0.900	1.014	1.027
LN	1000	10	0.903	0.906	0.567	0.659	0 +
MVT(1)	1000	10	0.914	0.914	0.541	22634.3	0 +

For large *n*, the semiparametric and nonparametric regions are likely to have coverage near 0.90 because the coverage on the training sample is slightly larger than 0.9 and x_f comes from the same distribution as the x_i . For n = 10p and $2 \le p \le 40$, the semiparametric region had coverage near 0.9. The ratio of the volumes

$$\frac{h_i^p \sqrt{det(\boldsymbol{C}_i)}}{h_2^p \sqrt{det(\boldsymbol{C}_2)}}$$

was recorded where i = 1 was the nonparametric region, i = 2 was the semiparametric region, and i = 3 was the parametric MVN region. The volume ratio converges in probability to 1 for $N_p(\mu, \Sigma)$ data, and the ratio converges to 1 for i = 1 on a large class of elliptically contoured distributions. The parametric MVN region often had coverage much lower than 0.9 with a volume ratio near 0, recorded as 0+. The volume ratio tends to be tiny when the coverage is much less than the nominal value 0.9. For $10p \le n \le 20p$, the nonparametric region often had good coverage and volume ratio near 0.5.

Simulations and Table 2 suggest that for $N_p(\mu, \Sigma)$ data, the coverages (ncov, scov and mcov) for the 3 regions are near 90% for n = 20p and that the volume ratios voln and volm are near 1 for n = 50p. With fewer than 5000 runs, this result held for $2 \le p \le 80$. For the non-elliptically contoured LN data, the nonparametric region had voln well under 1, but the volume ratio blew up for $w \sim MVT_p(1)$.

4. Discussion

4.1 General Comments

There are not many practical competitors for the new prediction intervals and regions. Parametric prediction intervals and regions usually assume normality and tend to have severe undercoverage when the normality assumption does not hold. For confidence intervals and testing, misspecification of normality is sometimes not too important if the estimators are asymptotically normal, but for parametric prediction intervals and regions, correct specification of the parametric model is important. For example, do not use a parametric prediction region based on the multivariate normal distribution if the plotted points in the DD plot fail to cover the identity line.

Another competitor for regression is bootstrap prediction intervals. These PIs take hundreds of times longer to compute than PI (9), and convergence problems are greatly multiplied for models such as nonlinear regression models. Also bootstrap PIs may not be valid if a fixed number B of bootstrap samples are used. Di Bucchianico, Einmahl and Mushkudiani (2001) use the minimum volume ellipsoid (MVE) estimator to cover m out of n cases to produce MVE tolerance regions, but the technique can only be used on tiny data sets.

The location model is a special case of both the regression model (1) and of the multivariate location and dispersion

model. Let $a_n = \left(1 + \frac{15}{n}\right)\sqrt{\frac{n+1}{n-1}}$. Let $c = \lceil n(1-\delta) \rceil$. Let shorth $(c) = (Y_{(d)}, Y_{(d+c-1)})$. Let MED(n) be the sample median. If $Y_1, ..., Y_n$ are iid, then the recommended large sample $100(1-\delta)\%$ PI for Y_f is the closed interval $[L_n, U_n] = [(1-a_n)\text{MED}(n) + a_nY_{(d)}, (1-a_n)\text{MED}(n) + a_nY_{(d+c-1)}]$. This PI is (5) using the least absolute deviations estimator, but with a closed interval.

Simulations were done in *Splus* and *R*. See R Development Core Team (2008). The Buxton data and programs in the collection of functions *rpack.txt* are available at (www.math.siu.edu/olive/ol-bookp.htm). For multiple linear regression, the function pisim simulates PIs (4) and (5) while the *Splus* function pisim4 simulates PIs (8) and (9) using OLS, L_1 and *M*-estimators. The function pisim3 was used to create Table 1 while pisim5 uses nls to simulate PIs for nonlinear regression. Care is needed when using pisim5 since for some versions of *R/Splus*, the nls function will fail to converge for some runs. Using nruns = 500 is less likely to cause an error than nruns=5000. The function *predsim* was used for Table 2. The function *ddplot4* was used to produce Figures 2, 3 and 4. The function *lpisim* simulates the PI for the location model while *covrmvn* computes the RMVN estimator.

4.2 Conclusions

Parametric prediction intervals and regions are notorious for severe undercoverage. The new techniques are designed to have good coverage at the training data, even if the model assumptions fail to hold. The Olive (2007) PIs (4) and (5) are tailored for multiple linear regression but are too short for many other techniques for moderate n. PIs (8) and (9) are generally longer than PIs (4) and (5) and have coverage near or higher than the nominal value for many techniques even for moderate n, say n > 10 (model degrees of freedom). PIs (8) and (9) are quite conservative for multiple linear regression for moderate n. These PIs are useful since the error distribution does not need to be known.

The new nonparametric and semiparametric prediction regions appear to have good coverage for n > 20p and may be the first easily computed prediction regions that are effective when the underlying multivariate distribution is unknown.

For the prediction regions, use the DD plot to check the multivariate normality assumption and to check for the presence of outliers. If n > 20p and the plotted points cluster tightly about a line through the origin, then the nonparametric and semiparametric prediction regions may have good coverage. For regression with additive errors, if *n* is large and the plotted points cluster about the identity line in the response plot, then the new prediction intervals may have good coverage.

Acknowledgements

The author thanks the editor and two referees for their comments that improved this article.

References

- Buxton, L. H. D. (1920). The anthropology of Cyprus. *The Journal of the Royal Anthropological Institute of Great Britain and Ireland*, 50, 183-235. http://dx.doi.org/10.2307/2843379
- Cai, T., Tian, L., Solomon, S. D., & Wei, L. J. (2008). Predicting future responses based on possibly misspecified working models. *Biometrika*, 95, 75-92. http://dx.doi.org/10.1093/biomet/asm078
- Chambers, J. M., & Hastie, T. J. (Eds.) (1993). Statistical models in S. New York, NY: Chapman & Hall.
- Cook, R. D., & Weisberg, S. (1999). Applied regression including computing and graphics. New York, NY: John Wiley & Sons.
- Di Bucchianico, A., Einmahl, J. H. J., & Mushkudiani, N. A. (2001). Smallest nonparametric tolerance regions. *The Annals of Statistics*, 29, 1320-1343. http://dx.doi.org/10.1214/aos/1013203456
- Grübel, R. (1988). The length of the shorth. *The Annals of Statistics*, 16, 619-628. http://dx.doi.org/10.1214/aos/1176350823
- Hastie, T. J., & Tibshirani, R. J. (1990). Generalized additive models. London, UK: Chapman & Hall.
- Johnson, M. E. (1987). Multivariate statistical simulation. New York, NY: John Wiley & Sons.
- Johnson, R. A., & Wichern, D. W. (1988). *Applied multivariate statistical analysis* (2nd ed.). Englewood Cliffs, NJ: Prentice Hall.
- MathSoft. (1999). S-plus 2000 guide to statistics, Vol. 1. Seattle, WA: MathSoft.
- Müller, U. U., Schick, A., & and Wefelmeyer, W. (2012). Estimating the error distribution function in semiparametric additive regression models. *Journal of Statistical Planning and Inference*, 142, 552-566. http://dx.doi.org/10.1016/j.jspi.2011.08.013
- Olive, D. J. (2002). Applications of robust distances for regression. *Technometrics*, 44, 64-71. http://dx.doi.org/10.1198/004017002753398335
- Olive, D. J. (2007). Prediction intervals for regression models. *Computational Statistics and Data Analysis*, 51, 3115-3122. http://dx.doi.org/10.1016/j.csda.2006.02.006
- Olive, D. J., & Hawkins, D. M. (2003). Robust regression with high coverage. *Statistics and Probability Letters*, 63, 259-266. http://dx.doi.org/10.1016/S0167-7152(03)00090-7
- Olive, D. J., & Hawkins, D. M. (2010). Robust multivariate location and dispersion. Retrieved from http://www.math.siu.edu/olive/preprints.htm

- R Development Core Team. (2008). R: a language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. Retrieved from http://www.R-project.org
- Rousseeuw, P. J., & Leroy, A. M. (1987). *Robust regression and outlier detection*. New York, NY: John Wiley & Sons.
- Rousseeuw, P. J., & Van Driessen, K. (1999). A fast algorithm for the minimum covariance determinant estimator. *Technometrics*, *41*, 212-223. http://dx.doi.org/10.1080/00401706.1999.10485670
- Wang, L., Liu, X., Liang, H., & Carroll, R. J. (2011). Estimation and variable selection for generalized additive partial linear models. *The Annals of Statistics*, *39*, 1827-1851. http://dx.doi.org/10.1214/11-AOS885SUPP
- Welsh, A. H. (1986). Bahadur representation for robust scale estimators based on regression residuals. *The Annals of Statistics*, 14, 1246-1251. http://dx.doi.org/10.1214/aos/1176350064
- Wood, S. N. (2006). Generalized additive models: an introduction with R. Boca Rotan, FL: Chapman & Hall/CRC.
- Zhang, J., Olive, D. J., & Ye, P. (2012). Robust covariance matrix estimation with canonical correlation analysis. *International Journal of Statistics and Probability*, 1, 119-136. http://dx.doi.org/10.5539/ijsp.v1n2p119