

UCLA

Department of Statistics Papers

Title

A Bias Bound for Applying Linear Regression to a General Linear Model

Permalink

<https://escholarship.org/uc/item/5pf1860c>

Authors

Li, Ker-Chau
Duan, Naihua

Publication Date

1991-07-16

Peer reviewed

A BIAS BOUND FOR LEAST SQUARES LINEAR REGRESSION

Naihua Duan and Ker-Chau Li

RAND Corporation and University of California

Abstract: Consider a general linear model $y=g(\alpha+\beta\mathbf{x})+\epsilon$, where the link function g is arbitrary and unknown. The maximal component of (α,β) that can be identified is the direction of β , which measures the substitutability of the components of \mathbf{x} . If $\zeta(\beta\mathbf{x})=E(\mathbf{x}|\beta\mathbf{x})$ is linear in $\beta\mathbf{x}$, the least squares linear regression of y on \mathbf{x} gives a consistent estimate for the direction of β , despite possible nonlinearity in the link function (Brillinger (1977, 1982)). If $\zeta(\beta\mathbf{x})$ is nonlinear, the linear regression might be inconsistent for the direction of β . We establish a bound for the asymptotic bias, which is determined from the nonlinearity in $\zeta(\beta\mathbf{x})$, and the multiple correlation coefficient R^2 for the least squares linear regression of y on \mathbf{x} . According to the bias bound, the linear regression is nearly consistent for the direction of β , despite possible nonlinearity in the link function, provided that the nonlinearity in $\zeta(\beta\mathbf{x})$ is small compared to R^2 . Our measure of nonlinearity in $\zeta(\beta\mathbf{x})$ is analogous to the maximal curvature studied by Cox and Small (1978). The bias bound is tight; we give the construction for the least favorable models which achieve the bias bound. The theory is applied to a special case for an illustration.

Key words and phrases: Lack of fit, link function, maximal curvature, nonlinearity, projection index, projection pursuit.

1. Main Result

Least squares linear regression is one of the most widely used statistical tools. It is based on the standard linear model:

$$y = \alpha + \beta\mathbf{x} + \epsilon, \quad \epsilon | \mathbf{x} \sim N(0, \sigma^2), \quad (1)$$

where y denotes a scalar outcome variable, and \mathbf{x} denotes a p -dimensional column vector of regressor variables. In empirical applications, it is unlikely for the standard linear model to hold exactly. Therefore we need to be concerned about possible violations of the model assumptions. For example, we might consider *distribution violation*: the error distribution might not be normal. There is a rich literature on robust methods for estimating the linear model in the presence of distribution violation; see, e.g., Huber (1981).

Another serious challenge to the standard linear model is the violation of the functional form. For example, the true model might be a power transformation model; the working model (1) might be based on a wrong transformation. More specifically, we assume the true model has the following form:

$$y = g(\alpha + \beta\mathbf{x}) + \epsilon, \quad E(\epsilon | \mathbf{x}) = 0, \quad \beta \neq \mathbf{0}, \quad (2)$$

where g is the *link function*, assumed to be arbitrary and unknown. Following Brillinger (1977, 1982), we call a model of form (2) a general linear model (GLM). To avoid trivialities, we assume in (2) that β is not null.

When the link function is arbitrary and unknown, any linear transformation of $\alpha + \beta\mathbf{x}$ can be absorbed into the link function. Therefore we cannot identify the intercept α , nor can we identify the length or the orientation of β . The most that we can identify is the direction of β , i.e., the collection of the ratios $\{\beta_j/\beta_k, j, k = 1, \dots, p\}$. The direction of β measures the substitutibility of the components of \mathbf{x} , and might be the main quantities of interest in many applications.

We will focus on estimating the direction of β , using the least squares linear regression of y on \mathbf{x} . We are concerned whether the linear regression still provides a valid estimate for the direction of β when the link function is nonlinear. We assume a mild condition on the GLM and the regressor \mathbf{x} :

(A.1) The regressor \mathbf{x} is sampled randomly from a probability distribution $Q(\mathbf{x})$; the following moments exist: $\mu = E(\mathbf{x})$, $\Sigma = \text{Cov}(\mathbf{x})$, Σ^{-1} , $E(y\mathbf{x}')$, $\sigma^2(\mathbf{x}) = \text{Var}(\epsilon | \mathbf{x})$, and $E[\sigma^2(\mathbf{x})\mathbf{x}\mathbf{x}']$.

Under (A.1), the least squares estimate $\hat{\beta}_{LS}$ converges to

$$\beta_{LS} = \text{Cov}(g(\alpha + \beta\mathbf{x}), \mathbf{x}')\Sigma^{-1}. \quad (3)$$

When the link function g is nonlinear, β_{LS} might not have the same direction as β . Brillinger (1977, 1982) established that β_{LS} does have the same direction as β , despite possible nonlinearity in the link function, provided that \mathbf{x} is normally distributed. The result also holds under the weaker condition

(A.2) $\zeta(\beta\mathbf{x}) = E(\mathbf{x}|\beta\mathbf{x})$ is linear in $\beta\mathbf{x}$.

Theorem 1. *Assume that the random vector (y, \mathbf{x}') follows a GLM (2), and satisfies (A.1) and (A.2). We then have $\beta_{LS} \propto \beta$.*

For empirical applications, it is unlikely for condition (A.2) to hold exactly. When (A.2) fails, β_{LS} might not have the same direction as β . However, we

would expect the noncollinearity to be minor if $\zeta(\beta\mathbf{x})$ is nearly linear in $\beta\mathbf{x}$. We will quantify the magnitude of the noncollinearity between β_{LS} and β , using a measure of nonlinearity in $\zeta(\beta\mathbf{x})$ which is analogous to Cox and Small's (1978) maximal curvature.

We measure the noncollinearity between β and β_{LS} by the squared sine between the two vectors,

$$\sin^2(\beta, \beta_{LS}) = 1 - (\beta_{LS}\Sigma\beta')^2 / (\beta\Sigma\beta')(\beta_{LS}\Sigma\beta'_{LS}),$$

where the sine function is taken with respect to the inner product

$$(\mathbf{v}, \mathbf{w}) = \mathbf{v}\Sigma\mathbf{w}'. \tag{4}$$

Note that $1 - \sin^2(\beta, \beta_{LS})$ is the squared correlation between $\beta\mathbf{x}$ and $\beta_{LS}\mathbf{x}$. If β and β_{LS} are nearly collinear, the linear regression of y on \mathbf{x} does provide useful information on the approximate direction of β . If the noncollinearity is severe, the linear regression might be misleading and should be applied with caution.

Our main result is the following bias bound, which is proved in Section 2.

Theorem 2. *Assume the random vector (y, \mathbf{x}') follows a GLM (2), and satisfies (A.1). The noncollinearity between β and β_{LS} in (3) satisfies the following bias bound:*

$$\sin^2(\beta, \beta_{LS}) \leq \min \left\{ 1, \frac{\nu_2(\beta)/(1 - \nu_2(\beta))}{R^2/(1 - R^2)} \right\}, \tag{5}$$

where $R^2 = \text{Var}(\beta_{LS}\mathbf{x})/\text{Var}(y)$ is the usual R^2 for the least squares linear regression of y on \mathbf{x} , and $\nu_2(\beta)$ is given by

$$\nu_2(\beta) = \max_{\mathbf{b} \in R^p} \frac{\mathbf{b}T\mathbf{b}'}{\mathbf{b}\Sigma\mathbf{b}'} \mid \text{constraint } \mathbf{b}\Sigma\beta' = 0, \quad \text{where } T = \text{Cov}(\zeta(\beta\mathbf{x})).$$

The scalar $\nu_2(\beta)$ is the second eigenvalue for the spectral decomposition of T with respect to Σ . (The first eigenvector is β ; the first eigenvalue is one.) We now interpret $\nu_2(\beta)$ as a measure for the nonlinearity in $\zeta(\beta\mathbf{x})$. More specifically, $\nu_2(\beta)$ measures the deviation of $\zeta(\beta\mathbf{x})$ from the linear regression function

$$l(\beta\mathbf{x}) = \mu + \Sigma\beta'\beta(\mathbf{x} - \mu)/\beta\Sigma\beta',$$

where we take the linear regression of \mathbf{x} on $\beta\mathbf{x}$. Note that $\zeta(\beta\mathbf{x})$ is linear if and only if the regression of $\mathbf{b}\mathbf{x}$ on $\beta\mathbf{x}$ is linear for all $\mathbf{b} \in R^p$. For any $\mathbf{b} \in R^p$, consider the decomposition

$$\mathbf{b}\mathbf{x} = \mathbf{b}l + (\mathbf{b}\zeta - \mathbf{b}l) + (\mathbf{b}\mathbf{x} - \mathbf{b}\zeta),$$

where the first term is the linear regression of $\mathbf{b}\mathbf{x}$ on $\beta\mathbf{x}$, the second term is the lack of fit for this linear regression, and the third term is the pure error. We measure the nonlinearity in $\mathbf{b}\zeta(\beta\mathbf{x})$ by the proportion of the variance in $\mathbf{b}\mathbf{x}$ accounted for by the lack of fit term above:

$$LF(\beta, \mathbf{b}) = \text{Var}(\mathbf{b}\zeta - \mathbf{b}l) / \text{Var}(\mathbf{b}\mathbf{x}).$$

We then maximize $LF(\beta, \mathbf{b})$ over \mathbf{b} to measure the nonlinearity in $\zeta(\beta\mathbf{x})$. Since $\beta\mathbf{x} = \beta\zeta$, it is easy to verify that this nonlinearity measure coincides with $\nu_2(\beta)$:

$$\nu_2(\beta) = \max_{\mathbf{b} \in R^p} LF(\beta, \mathbf{b}).$$

The nonlinearity measure $\nu_2(\beta)$ is analogous to the maximum curvature studied in Cox and Small (1978). Cox and Small considered the quadratic regression of $\mathbf{b}\mathbf{x}$ on $\beta\mathbf{x}$, and measured the nonlinearity by the proportion of the variance in $\mathbf{b}\mathbf{x}$ accounted for by the quadratic term. They then maximized the nonlinearity over \mathbf{b} . This is analogous to our consideration of LF and its maximization over \mathbf{b} .

Theorem 1 follows immediately from Theorem 2: if $\zeta(\beta\mathbf{x})$ is linear in $\beta\mathbf{x}$, we have $\nu_2(\beta) = 0$, therefore the right hand side of (5) is zero, and β_{LS} is collinear with β . When $\zeta(\beta\mathbf{x})$ is nonlinear, Theorem 1 might not hold; however, the noncollinearity between β_{LS} and β is small if the nonlinearity in $\zeta(\beta\mathbf{x})$, $\nu_2(\beta)$, is small compared to R^2 .

The bias bound (5) is tight in the following sense: for a given β , a given $Q(\mathbf{x})$, and a given R^2 , we can find a GLM for which the noncollinearity $\sin^2(\beta, \beta_{LS})$ equals the right hand side of (5). The construction of such a GLM is given in Section 2.

For empirical applications, we need to estimate the right hand side of (5). If we have a good initial estimate $\hat{\beta}$ for the direction of β , we can estimate $\zeta(\beta\mathbf{x})$ by an appropriate nonparametric regression of \mathbf{x} on $\hat{\beta}\mathbf{x}$, then estimate T , $\nu_2(\beta)$, and the bias bound. If we don't have a good initial estimate for the direction of β , we can take a conservative approach and replace $\nu_2(\beta)$ in (5) by its maximum

$$\nu_2^{\text{sup}} = \max_{\beta \in R^p} \nu_2(\beta).$$

For each $\beta \in R^p$, we estimate $\zeta(\beta\mathbf{x})$, then estimate T and $\nu_2(\beta)$. We then search for the maximum of the estimated $\nu_2(\beta)$'s. This is a projection pursuit problem, with $\nu_2(\beta)$ as the projection index. Huber (1985 and discussions) gave a comprehensive review of the projection pursuit problem. Cox (1985) suggested using the maximum curvature in Cox and Small (1978) as the projection index. This is analogous to using $\nu_2(\beta)$ as the projection index.

We prove Theorem 2 in Section 2, then apply the theory to a special case in Section 3 for an illustration.

Remark 1. It is helpful to interpret the bias bound (5) in terms of the equivalent magnitude of estimation error. Under the standard linear model (1), the least squares estimate $\hat{\beta}_{LS}$ is unbiased for β . We measure the magnitude of its estimation error by the mean squared sine,

$$\text{MSS} = E[\sin^2(\beta, \hat{\beta}_{LS})],$$

where the sine function is taken with respect to the inner product (4). It is easy to verify that

$$\text{MSS} \cong n^{-1}(p-1)(1-R^2)/R^2,$$

where n is the sample size.

The mean squared sine approximately equals the right hand side of the bias bound (5) if

$$n = \frac{p-1}{\nu_2(\beta)/(1-\nu_2(\beta))}. \quad (6)$$

If the sample size is much smaller than the right hand side of (6), bias is likely to be negligible compared with estimation error. If the sample size is much larger than the right hand side of (6), bias might dominate estimation error if the true link function is substantially nonlinear.

For the same nonlinearity measure $\nu_2(\beta)$, the right hand side of (6) is proportional to $p-1$. The asymptotic bias is less serious (compared to estimation error) for larger p 's.

Remark 2. The bias bound (5) and the discussion in Remark 1 deal with the *worst case* bias when the true link function is the least favorable. For a specific empirical study, the actual bias might be substantially smaller than the right hand side of (5) if the true link function is not the least favorable.

The results in this paper can be used as a screening device to diagnose whether we might have a serious bias *if* the true link function is substantially nonlinear. If the right hand side of (5) is big, the empirical scientist should be alerted to pay more attention to the goodness of the link function used in his working model. If the right hand side of (5) is small, say, compared to the MSS discussed in Remark 1, the goodness of the link function is less crucial.

2. Proof

The bias bound (5) is equivalent to

$$\frac{R^2 \sin^2(\beta, \beta_{LS})}{1 - R^2 \cos^2(\beta, \beta_{LS})} \leq \min(R^2, \nu_2(\beta)).$$

It is easy to see that the left hand side is bounded from above by R^2 . We now derive the other bound. In the derivation below, we consider several geometric measures such as norm and orthogonality; all of these are taken with respect to the inner product (4).

We decompose β_{LS} into

$$\beta_{LS} = c\beta + \beta^\perp, \quad \beta^\perp \perp \beta,$$

where c is a scalar. Since

$$\sin^2(\beta, \beta_{LS}) = \|\beta^\perp\|^2 / \|\beta_{LS}\|^2, \quad R^2 = \|\beta_{LS}\|^2 / \text{Var}(y),$$

we have

$$\frac{R^2 \sin^2(\beta, \beta_{LS})}{1 - R^2 \cos^2(\beta, \beta_{LS})} = \|\beta^\perp\|^2 / (\text{Var}(y) - c^2 \|\beta\|^2). \quad (7)$$

Let \mathbf{v} denote the slope vector for the least squares linear regression of $g(\alpha + \beta\mathbf{x})$ on ζ :

$$\mathbf{v}T = \text{Cov}(g(\alpha + \beta\mathbf{x}), \zeta(\beta\mathbf{x})').$$

If T does not have full rank, \mathbf{v} is not uniquely defined. We can take any version of \mathbf{v} . Consider the decomposition

$$\mathbf{v} = d\beta + \mathbf{w}, \quad \mathbf{w} \perp \beta, \quad (8)$$

where d is a scalar. We claim that

$$c = d, \quad \beta^\perp = \mathbf{w}T\Sigma^{-1}. \quad (9)$$

To see this, note that

$$\mathbf{v}T = \text{Cov}(g, \zeta') = \text{Cov}(g, \mathbf{x}') = \beta_{LS}\Sigma,$$

$$\beta T = \text{Cov}(\beta\zeta, \zeta') = \text{Cov}(\beta\mathbf{x}, \zeta') = \text{Cov}(\beta\mathbf{x}, \mathbf{x}') = \beta\Sigma.$$

It follows that

$$\beta_{LS}\Sigma = \mathbf{v}T = d\beta\Sigma + \mathbf{w}T,$$

which proves (9).

Let $\tau^2 = \text{Var}(y - \mathbf{v}\zeta)$. It follows from (8) and (9) that

$$\text{Var}(y) = \text{Var}(\mathbf{v}\zeta) + \tau^2 = c^2 \|\beta\|^2 + \mathbf{w}T\mathbf{w}' + \tau^2. \quad (10)$$

Combining (7), (9) and (10), we have

$$\frac{R^2 \sin^2(\beta, \beta_{LS})}{1 - R^2 \cos^2(\beta, \beta_{LS})} = \mathbf{w}T\Sigma^{-1}T\mathbf{w}' / (\mathbf{w}T\mathbf{w}' + \tau^2) \leq \mathbf{w}T\Sigma^{-1}T\mathbf{w}' / \mathbf{w}T\mathbf{w}'. \quad (11)$$

We want to maximize the right hand side of (11) under the constraint $\mathbf{w} \perp \beta$, thus we consider the spectral decomposition of $T\Sigma^{-1}T$ with respect to T . This is equivalent to the spectral decomposition of T with respect to Σ . If \mathbf{w} is an eigenvector with eigenvalue ν for the second spectral decomposition, i.e.,

$$\mathbf{w}T = \nu\mathbf{w}\Sigma, \quad (12)$$

then \mathbf{w} is also an eigenvector with the same eigenvalue for the first spectral decomposition. The first eigenvalue for the spectral decomposition (12) is one, with the corresponding eigenvector being $\mathbf{w} \propto \beta$. All other eigenvectors are orthogonal to β . It follows that the right hand side of (11) is maximized by taking \mathbf{w} to be the second eigenvector for the spectral decomposition (12). The maximum is $\nu_2(\beta)$, the second eigenvalue for (12). This completes the proof for Theorem 2.

We now establish that the bias bound (5) is tight. For a given β and a given $Q(\mathbf{x})$, the "least favorable" GLM's have no pure error, $\epsilon \equiv 0$, and have the following link function:

$$y = g(\alpha + \beta\mathbf{x}) = c\beta\mathbf{x} + \mathbf{w}^*\zeta(\beta\mathbf{x}), \quad (13)$$

where \mathbf{w}^* is the second eigenvector for (12), and c is a scalar to be determined from the given R^2 . For those GLM's, we have

$$R^2 = (c^2\|\beta\|^2 + \nu_2(\beta)^2\|\mathbf{w}^*\|^2) / (c^2\|\beta\|^2 + \nu_2(\beta)\|\mathbf{w}^*\|^2).$$

If $R^2 \geq \nu_2(\beta)$, we have

$$c^2 = (R^2 - \nu_2(\beta))\nu_2(\beta)\|\mathbf{w}^*\|^2 / (1 - R^2)\|\beta\|^2,$$

thus there exists a least favorable GLM of form (13) which has the given R^2 and achieves the right hand side of the bias bound (5).

If $R^2 < \nu_2(\beta)$, there does not exist a corresponding least favorable GLM of form (13). The right hand side of (5) can be achieved in this case with the GLM

$$y = \mathbf{w}^*\zeta(\beta\mathbf{x}) + \epsilon, \quad \epsilon | \mathbf{x} \sim N(0, \tau^2), \quad (13')$$

where

$$\tau^2 = (\nu_2(\beta) - R^2)\nu_2(\beta)\|\mathbf{w}^*\|^2/R^2.$$

The bias bound (5) is noninformative in this case: the noncollinearity between β and β_{LS} is unrestricted. For the GLM (13'), β_{LS} is orthogonal to β .

Remark 3. The bias bound allows the link function to be arbitrary. If we have prior information on the link function, it might be possible to sharpen the bound. For example, it might be known that the link function is monotonic. The link function for the least favorable GLM (13) might not be monotonic. Since

$$\text{Cov}(\mathbf{w}^*\zeta(\beta\mathbf{x}), \beta\mathbf{x}) = \mathbf{w}^*\Sigma\beta' = 0,$$

$\mathbf{w}^*\zeta(\beta\mathbf{x})$ is not monotonic in $\beta\mathbf{x}$. If c is small compared to $\|\mathbf{w}^*\|$, i.e., $R^2 - \nu_2(\beta)$ is small, the link function is not monotonic. If we restrict to monotonic link functions, it might be possible to improve upon the bias bound (5). An example is given in Section 3.

3. A Special Case

We illustrate the application of Theorem 2 with a special case which might be of interest in itself. We assume \mathbf{x} is uniformly distributed over the square $(-1 \leq x_1 \leq 1, -1 \leq x_2 \leq 1)$. Let $\beta = (1, t)$. Without loss of generality, we assume $0 \leq t \leq 1$.

If $t = 0$ or $t = 1$, $\zeta(\beta\mathbf{x})$ is linear, thus β_{LS} is collinear with β . If $0 < t < 1$, then $\zeta(\beta\mathbf{x}) = (\zeta_1(\beta\mathbf{x}), \zeta_2(\beta\mathbf{x}))$ is nonlinear:

$$\zeta_1(\beta\mathbf{x}) = \begin{cases} (\beta\mathbf{x} - 1 + t)/2, & \text{if } \beta\mathbf{x} < -1 + t; \\ \beta\mathbf{x}, & \text{if } -1 + t \leq \beta\mathbf{x} \leq 1 - t; \\ (\beta\mathbf{x} + 1 - t)/2, & \text{if } \beta\mathbf{x} > 1 - t; \end{cases}$$

$$\zeta_2(\beta\mathbf{x}) = \begin{cases} (\beta\mathbf{x} + 1 - t)/2t, & \text{if } \beta\mathbf{x} < -1 + t; \\ 0, & \text{if } -1 + t \leq \beta\mathbf{x} \leq 1 - t; \\ (\beta\mathbf{x} - 1 + t)/2t, & \text{if } \beta\mathbf{x} > 1 - t; \end{cases}$$

thus β_{LS} and β might not be collinear. It is easy to verify that

$$\nu_2(\beta) = t(1 - t)^2/2 \leq 2/27,$$

where the maximum occurs at $t = 1/3$.

We have tabulated the bias bound (5) for $\beta = (1, 1/3)$. The second column in Table 1 gives the maximal angle between β and β_{LS} for R^2 's ranging from 10% to 90%. The asymptotic bias is substantial when R^2 is not close to one.

Table 1. Maximal angle between β_{LS} and $\beta = (1, 1/3)$

R^2	Maximal angle	
	Arbitrary g [1]	Monotonic g [2]
0.10	58.1°	12.5°
0.20	34.4°	12.5°
0.30	25.6°	12.5°
0.40	20.3°	12.5°
0.50	16.4°	12.5°
0.60	13.4°	12.5°
17/27	12.5°	12.5°
0.70	10.7°	10.7°
0.80	8.1°	8.1°
0.90	5.4°	5.4°

[1] No restrictions on the link function g .

[2] Link function g restricted to be monotonic.

For $p = 2$, $\nu_2(\beta) = 2/27$, the right hand side of (6) is 12.5, which is quite small. For most relevant sample sizes, bias might dominate estimation error if the true link function is substantially nonlinear.

It can be shown that the least favorable GLM (13) is monotonic in $\beta\mathbf{x}$ for $R^2 \geq 17/27$, and is not monotonic for $R^2 < 17/27$. If we restrict to monotonic link functions, we can sharpen the bias bound for $R^2 < 17/27$: the maximal angle is identical to that for $R^2 = 17/27$; see the third column in Table 1. For small R^2 's, the bias bound is sharpened substantially.

Acknowledgements

Naihua Duan's research is supported in part by a cooperative agreement between the RAND Corporation, SIMS, and U.S. Environmental Protection Agency, and in part by RAND corporate funds. Ker-Chau Li's research is supported by NSF Grant DMS86-02018. We appreciate the helpful comments from an anonymous referee.

References

- Brillinger, D. R. (1977). The identification of a particular nonlinear time series system. *Biometrika* **64**, 509-515.
- Brillinger, D. R. (1982). A general linear model with 'Gaussian' regressor variables. In *A*

Festschrift for Erich L. Lehmann (Edited by Bickel, P. J., Doksum, K. A. and Hodges, J. L.), 97-114. Wadsworth, Belmont, Calif., U.S.A.

Cox, D. R. (1985). Discussion of Huber, P. J.: Projection pursuit. *Ann. Statist.* **13**, 493-494.

Cox, D. R. and Small, N. J. H. (1978). Testing multivariate normality. *Biometrika* **65**, 263-272.

Huber, P. J. (1981). *Robust Statistics*. John Wiley, New York, U.S.A.

Huber, P. J. (1985). Projection pursuit (with discussion). *Ann. Statist.* **13**, 435-525.

RAND Corporation, Santa Monica, California 90406, U.S.A.

Department of Mathematics, University of California, Los Angeles, California 90024, U.S.A.

(Received August 1989; accepted June 1990)