

A Conspiracy of Random X and Model Violation against Classical Inference in Linear Regression

A. Buja, R. Berk, L. Brown, E. George, M. Traskin, K. Zhang, L. Zhao
The Wharton School, University of Pennsylvania

April 19, 2013

Abstract

Following the econometric literature on model misspecification, we examine statistical inference for linear regression coefficients β_j when the predictors are random and the linear model assumptions of first and/or second order are violated: $\mathbf{E}[Y|X_1, \dots, X_p]$ is not linear in the predictors and/or $\mathbf{V}[Y|X_1, \dots, X_p]$ is not constant. Such inference is meaningful if the linear model is seen as a useful approximation rather than part of a generative truth.

A difficulty well-known in econometrics is that the required standard errors under random predictors and model violations can be very different from the conventional standard errors that are valid when the linear model is correct. The difference stems from a synergistic effect between model violations and randomness of the predictors. We show that asymptotically the ratios between correct and conventional standard errors can range between infinity and zero. Also, these ratios vary between predictors within the same multiple regression.

This difficulty has consequences for statistics: It entails the breakdown of the classical ancillarity argument for predictors. If the assumptions of a generative regression model are violated, then the ancillarity argument for the predictors no longer holds, treating predictors as fixed is no longer valid, and standard inferences may lose their significance and confidence guarantees.

The standard econometric solution for consistent inference under misspecification and random predictors is based on the “sandwich estimator” of the covariance matrix of $\hat{\beta}$. A plausible alternative is the paired bootstrap which resamples predictors and response jointly. Discrepancies between conventional and bootstrap standard errors can be used as diagnostics for predictor-specific model violations, in analogy to econometric misspecification tests. The good news is that when model violations are sufficiently strong to invalidate conventional linear inference, their nature tends to be visible in graphical diagnostics.

Keywords: Ancillarity of predictors; First and second order incorrect models; Model misspecification

1 Introduction

Classical inference for linear regression is conditional on the observed predictors. This means that in a linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}_N, \sigma^2 \mathbf{I}_{N \times N}), \quad (1)$$

the standard errors $\text{SE}(\hat{\beta}_j)$ only account for the variability in the error vector $\boldsymbol{\epsilon}$ but not in the predictor data \mathbf{X} . In effect, the predictor data are treated as fixed known constants even when they are samples, as they are in observational data. The justification for conditioning on random \mathbf{X} is that the predictor distribution is ancillary

for the parameters β and σ of the model, hence conditioning on \mathbf{X} produces valid frequentist inference for these parameters (Cox and Hinkley 1974, Example 2.27).

A problem with the ancillarity argument is that it holds only when the linear model is correct. In practice, of course, whether a model is correct is never known, yet it is fitted just the same and also richly interpreted with statistical inference that assumes the truth of the model. An awareness of potential problems with such inference is occasionally handled with a reference to “robustness of validity”, meaning that conventional t - and F -tests are conservative (though not efficient) when the error distribution is heavier-tailed than normal. This argument, however, does not address model failure from non-linearity or heteroskedasticity. When first and/or second order assumptions are violated, there may arise pernicious problems for statistical inference because these violations can interact with the randomness of the predictors and invalidate conventional standard errors derived from \mathbf{X} -conditional linear models theory. As we will see, this interaction usually inflates and occasionally deflates the sampling variability of slope estimates relative to conventional \mathbf{X} -conditional standard errors.

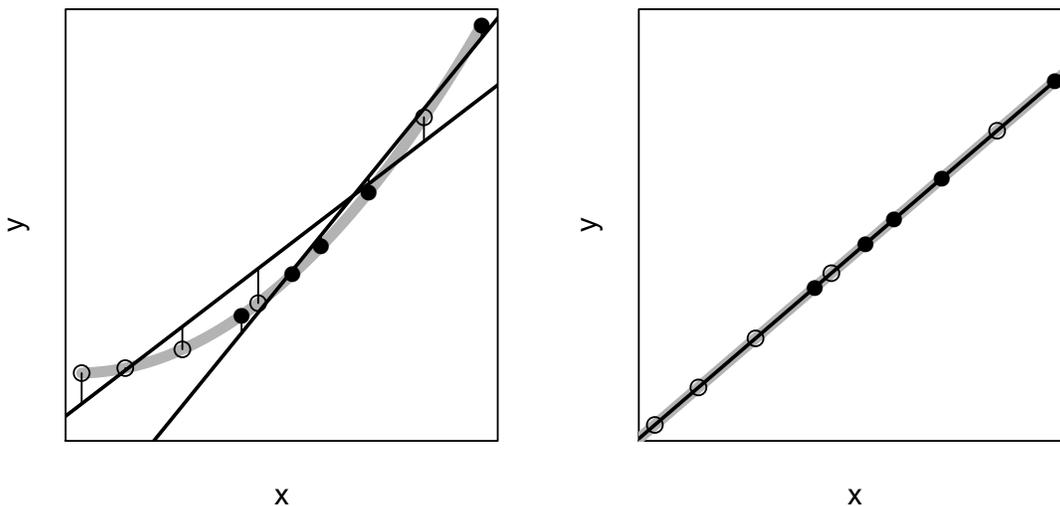


Figure 1: *Single Predictor, error-less Data: The filled and the open circles represent two “datasets” from the same population. The x -values are random; the y -values are a deterministic function of x : $y = \mu(x)$ (shown in gray). Left: The true response $\mu(x)$ (gray curve) is nonlinear; the open and the filled circles have different LS lines (black lines). Right: The true response $\mu(x)$ (gray line) is linear; the open and the filled circles have the same LS line (black overlaid on gray).*

The nature of the interaction between model violations and randomness of the predictors is most drastically illustrated with the artificial example of an error-free nonlinearity, that is, $y = \mu(x)$ is a deterministic but nonlinear function of a (single) random predictor x , as in the left hand plot of Figure 1. It shows two error-free (x, y) “datasets” of size $N = 5$ each, marked by open circles (“white data”) and filled circles (“black data”). The datasets are clearly shifted against each other, the white data being more to the left, the black data more to the right, even though both are i.i.d. random samples from the same predictor distribution. The relative shift implies that when a straight line is fitted to the “response” in both datasets the white dataset will “see” a smaller slope than the black dataset. This difference is solely due to the random differences between the two predictor samples. — This effect vanishes when the response function $\mu(x)$ is linear, as depicted in the right hand plot of Figure 1 for the same two sets of predictor values: obviously both sets “see” the same slope. Thus it is the combination

of randomness of the x -values and nonlinearity of the response function $\mu(x)$ that conspires to create sampling variability in LS estimates of slopes (and intercepts as well). This additional source of sampling variability is separate from the variability analyzed in \mathbf{X} -conditional linear models theory, which is error in the response. Furthermore, the additional variability is of the same order $1/\sqrt{N}$ as the conventional error-based variability.

In the econometric literature on misspecification it is known that standard errors can be invalidated not only by nonlinearity but by heteroskedasticity as well. Heteroscedasticity's effect on sampling variability is somewhat less transparent than nonlinearity's, yet in econometrics it is heteroscedasticity that appears to be more frequently debated. In the present context the relative importance of the two effects is not the issue; rather, it is the sheer fact that statistical inference is affected by an interaction of first and second order model violations with predictor randomness.

An immediate consequence for classical statistics is that for observational data \mathbf{X} -conditional inference must be considered as suspect, and so must the ancillarity argument that justifies it. The ancillarity argument for regression has been unquestioned by statisticians for many decades, and this acceptance may be partly supported by a fallacious intuition. The fallacy may run along the following lines: One starts from the premise that if the data do not satisfy the first and second order assumptions for the model at hand, they will do so very likely for some larger linear model. In that "correct" model the ancillarity argument applies, hence it applies in the smaller model. This very last step is of course the fallacy: ancillarity for the predictors can only be argued in the correct model, not its incorrect submodels.

A further consequence for classical statistics is that if \mathbf{X} -conditional inference is to be used, then checking first and second order correctness of the chosen model is essential. To address this need we will propose some simple diagnostics that provide indications of model failure in a predictor-specific manner. Such diagnostics have of course precursors in econometric misspecification tests. Our diagnostics will add a new level of interpretability to these established techniques [??? make more specific].

Standard errors that capture the variability generated by the interaction of random predictors and model violations have existed for quite some time: asymptotically correct standard errors can be obtained from so-called "sandwich estimators" of the covariance matrix of the estimates (White 1980a, 1980b, 1981, 1982). These have generated a considerable literature and have been applied to every conceivable type of parametric regression. They are better known for being "heteroskedasticity consistent" (White 1980b) than nonlinearity consistent (White 1980a).

Another and quite obvious approach to standard errors under random predictors and model violations is bootstrapping $(\vec{\mathbf{X}}_i, Y_i)$ pairs. It captures the joint, unconditional sampling variability of the response *and* the predictors without making model assumptions other than i.i.d. sampling. There exists a literature that shows this bootstrap to be asymptotically correct for linear regression under very general conditions, even if the linear model is incorrect; see for example Freedman's (1981) analysis of the "correlation model", and Mammen's (1993) analysis for increasing p . The paired bootstrap, however, has been a source of unease to some statisticians because it violates conditionality on the predictors. As we have seen above, such unease is ill-founded because conditionality relies on the ancillarity argument for predictors, which in turn assumes the correctness of the model. Such assumption creates exposure to potentially flawed statistical inference.

For linear models applied to observational data it should be recommended that all inference be accompanied by standard errors obtained from either sandwich estimators or a paired bootstrap simulation, if only as a diagnostic. If data analysts observe sizable discrepancies between conventional and bootstrap standard errors, it should alert them to potential problems on two levels: (1) the existence of first and/or second order model incorrectness and consequently a need to diagnose the nature of nonlinearities and/or heteroskedasticities, and (2) a need to use the

bootstrap standard errors for statistical inference if these are wider.

The reader may object at this point that something is amiss with this argument: If the functional form of the model is incorrect and if there is indeed a nonlinear association between predictors and response, what is the meaning of a slope and its standard error? Shouldn't the nonlinearities be taken care of first, for example, by fitting an additive model (Buja, Hastie, and Tibshirani 1989; Hastie and Tibshirani 1990) and/or transforming the response (Box and Cox 1964)? These arguments are of course valid: competent data analysts would discover critical nonlinearities with residual analysis, non-parametric fitting, nonlinear transformation of the response, and other methods of model building. Yet, the counterargument is that even when a linear model for observational data is first order incorrect, it can still be *theoretically meaningful* and *practically useful* (viz Box' (1979, second section heading) famous dictum). Both assertions can be made precise:

- Theoretically, a fitted linear model can always be interpreted as estimating the best linear approximation $\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$ to the conditional expectation $\mathbf{E}[Y|X_1, \dots, X_p]$ at the population. This view is in line with machine learning approaches which investigate best performance within a functional form without assuming it to contain the truth. The difference of our focus to that in machine learning is that we are not concerned with prediction but with valid inference for the parameters of the fitted equation.
- Practically, a fitted linear model, even if not strictly correct, can give information about the direction of the association — positive or negative — between Y and X_j in the presence of the other predictors. Also, it is possible to give meaning to LS slopes as weighted averages of “observed slopes”, and this carries over to populations without linear model assumptions. Because of its intuitive appeal we discuss this interpretation in Appendix A.

In conclusion, the problem of inference for “incorrect” linear models is not misguided; it can be both meaningful and useful.

The idea that models are approximations, not generative truths, is found in a literature too large to list. Aspects of approximation are explicit in work on robustness, on misspecified models, and on nonparametric fitting. Most famously the idea of models as approximations is implied by Box (1979), albeit in more catchy language. On this matter we follow Cox' (1995) opinion that “it does not seem helpful to say that all models are wrong. The very word model implies simplification and idealization,” which implies approximation. The theme runs through the discussion generated by Chatfield (1995), which is where we found Cox' comment. — Among the several meanings of “approximation” as they relate to statistical models, the only one we have in mind is that of approximating the conditional expectation $\mathbf{E}[Y|X_1, \dots, X_p]$ of a numerically interpreted response by a linear function $\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$ of the predictors. The possible causes of nonlinearity include incorrect functional form, omission of predictors, and measurement errors in predictors. Whatever the causes, known or unknown, we assume nevertheless that interest is in $\mathbf{E}[Y|X_1, \dots, X_p]$ through inference for the coefficients of the best linear approximant $\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$.

This article continues as follows: Section 2 shows a data example in which the conventional and bootstrap standard errors of some slopes are off by as much as a factor 2. Section 3 introduces notation for populations and samples, LS approximations, adjustment operations, nonlinearities, and decompositions of responses. Section 4 describes the sampling variation caused by random predictors and nonlinearities. Section 5 compares three types of standard errors and shows that both the (potentially flawed) conventional and the (correct) unconditional standard errors are inflated by nonlinearities, but Section 6 shows that asymptotically there is no limit to which inflation of the unconditional standard error can exceed inflation of the conventional standard error. Appendix A gives intuitive meaning to LS slopes in terms of weighted average slopes found in the data.

2 The Problem Illustrated with

We will use the well-known Boston Housing data for illustration, but we must begin with a caveat: These data (Harrison and Rubinfeld 1978) are an enumeration of the census tracts in the larger Boston area. They do **not** form a random sample; at a minimum, stochastic modeling would have to include consideration of spatial correlation. Yet, these data have been used to illustrate regression methods at least since the regression diagnostics book by Belsey, Kuh and Welsch (1980). In what follows we will treat the data as if they were a random sample to illustrate the issue of conditional versus unconditional statistical inference in linear regression. Because of peculiarities in the associations among the 14 variables in these data, they are uniquely suited to illustrate the phenomenon we describe here.

The following table shows in the first three numeric columns the results of a linear regression of MEDV (median value of owner occupied homes, standardized) on thirteen predictors (all standardized). As we can see, the predictors RM (average number of rooms per dwelling) and LSTAT (percent of lower status population) are the strongest in terms of t -ratios.

	Slope.est	SE.conv	t.val	SE.boot	SE.boot/SE.conv
CRIM	-0.099	0.031	-3.261	0.033	1.074
ZN	0.121	0.035	3.508	0.035	1.004
INDUS	0.017	0.046	0.382	0.038	0.843
CHAS	0.074	0.024	3.152	0.036	1.503 <?
NOX	-0.224	0.048	-4.687	0.048	1.003
RM	0.290	0.032	9.149 <<	0.065	2.049 <<
AGE	0.002	0.040	0.044	0.050	1.236
DIS	-0.344	0.045	-7.598 <	0.048	1.068
RAD	0.288	0.062	4.620	0.060	0.958
TAX	-0.233	0.068	-3.409	0.051	0.740
PTRATIO	-0.218	0.031	-7.126 <	0.026	0.865
B	0.092	0.026	3.467	0.027	1.036
LSTAT	-0.413	0.039	-10.558 <<	0.078	1.995 <<

Important for us is the comparison of the conventional standard errors (`SE.conv`, conditional on \mathbf{X}) and the paired bootstrap standard errors (`SE.boot`, sampling (x, y) pairs, not residuals) in the second and the fourth numeric column, respectively. The last column shows their ratio. We see that for the two strongest predictors, LSTAT and RM, the bootstrap standard errors are about twice the size of the conventional standard errors. As a consequence, these predictors seem to have their t -statistics of -10.558 and 9.149 inflated by about a factor of two. If so, they would have to cede their first and second ranks to the predictors DIS and PTRATIO whose bootstrap standard errors are in rough agreement with the traditional standard errors and who exhibit t -statistics of -7.598 and -7.126. So we must ask how a discrepancy can arise between bootstrap and traditional standard errors.

A pointer toward an answer is given in Figure 2: We are first shown residuals plotted against all predictors, augmented with smooths. The smooths indicate that the predictors LSTAT and RM are those that most strongly suggest nonlinear effects. Even more convincingly, a second series of plots shows the results of an additive model fit, $Y \sim \phi_1(X_1) + \dots + \phi_p(X_p)$, augmented with bootstrap bands. The transformations of the predictors LSTAT and RM suggest most strongly a nonlinear effect once again. Other predictors may also have nonlinear associations, but LSTAT and RM are the two most egregious cases.

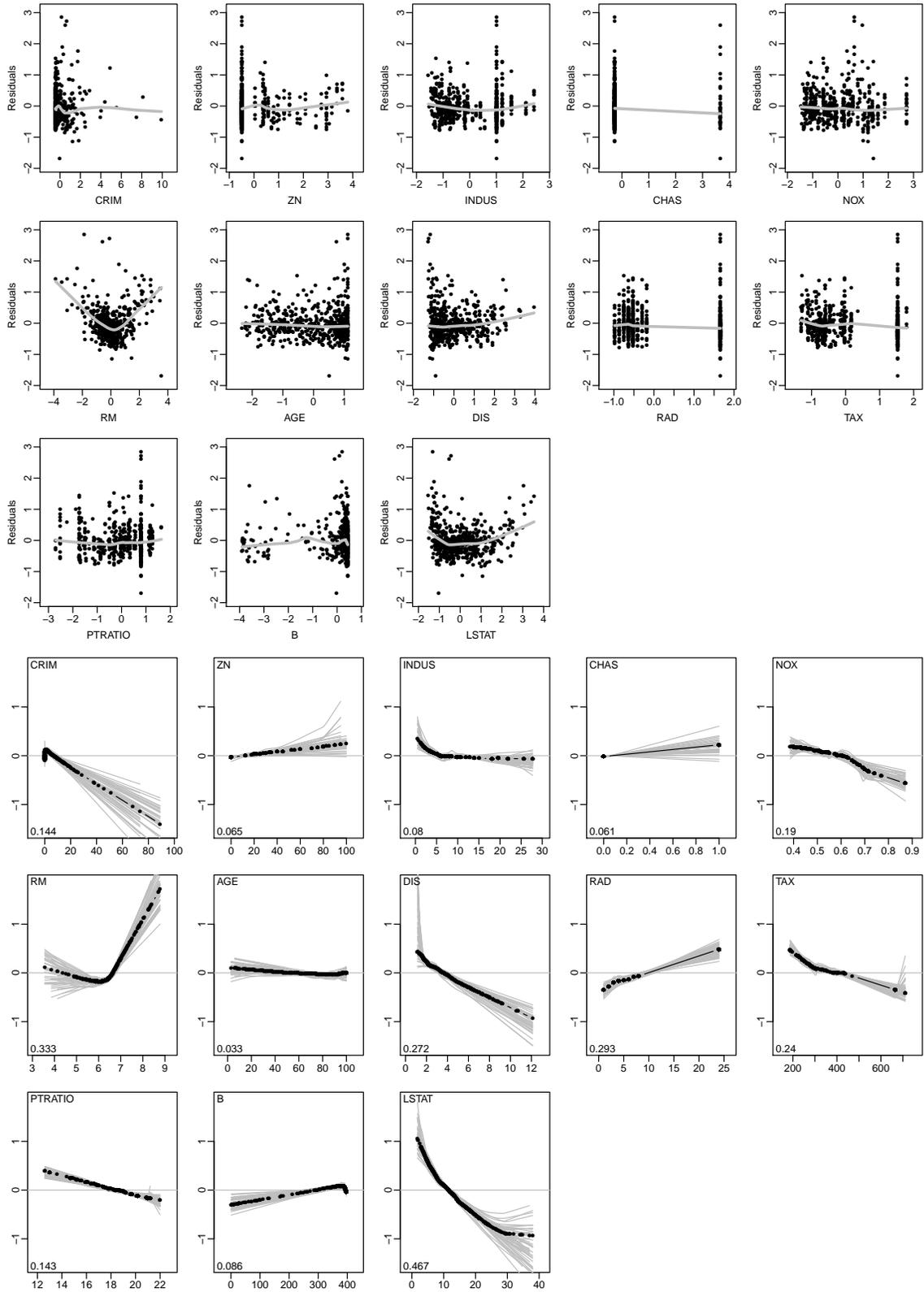


Figure 2: *Boston Housing data: The strongest two predictors, RM and LSTAT, show substantial nonlinearities. Top half: Residuals versus all predictors augmented with lowess smooths. Bottom half: ACE transformations for an additive model; the response is standardized but otherwise untransformed. The bottom left corners show the standard deviations of the transforms (= analogs of $|\hat{\beta}_j|$ when the $\text{sd}(\mathbf{x}_j) = 1$).*

3 Observational Regression Data: Estimation and its Target

Following Freedman (1981) and Mammen (1993), we introduce a framework in which linear LS is applied to data that are not described by a generative linear model but only by an i.i.d. sampling assumption. The framework therefore consists of (1) a largely arbitrary multivariate population in which one variable is singled out as the response, (2) the associated decomposition of the population response into a linear LS approximation, a nonlinearity, and error, (3) samples from the population as well as LS estimates obtained from the samples.

3.1 Quantities associated with the population

- The population will be described by random variables X_1, \dots, X_p and Y . At this point the only assumption is that these variables have a nontrivial joint distribution $\mathbf{P} = \mathbf{P}(dx_1, \dots, dx_p, dy)$ in the sense that second moments exist and the covariance matrix is of full rank.
- We write $\vec{\mathbf{X}} = (1, X_1, \dots, X_p)^T$ for the *column* random vector consisting of the predictor variables with a constant 1 prepended to accommodate an intercept term. We write any function $f(X_1, \dots, X_p)$ equivalently as $f(\vec{\mathbf{X}})$ as the prepended constant 1 is irrelevant.
- The “true response surface” $\mu(\vec{\mathbf{X}})$ is the conditional mean given the predictors $\vec{\mathbf{X}}$:

$$\mu(\vec{\mathbf{X}}) = \mathbf{E}[Y|\vec{\mathbf{X}}]. \quad (2)$$

No assumption of linearity as a function of $\vec{\mathbf{X}}$ is made.

- The **population linear LS approximation** to Y is the affine function defined by

$$\beta^T \vec{\mathbf{X}} = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p, \quad (3)$$

where the population intercept and slopes are collected in the $(p+1)$ -vector $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$ and β is defined by the following equivalent expressions:

$$\beta(\mathbf{P}) = \operatorname{argmin}_{\beta} \mathbf{E}[(Y - \beta^T \vec{\mathbf{X}})^2] = \operatorname{argmin}_{\beta} \left(-2\beta^T \mathbf{E}[Y \vec{\mathbf{X}}] + \beta^T \mathbf{E}[\vec{\mathbf{X}} \vec{\mathbf{X}}^T] \beta \right) \quad (4)$$

$$= \mathbf{E}[\vec{\mathbf{X}} \vec{\mathbf{X}}^T]^{-1} \mathbf{E}[Y \vec{\mathbf{X}}] \quad (5)$$

$$= \operatorname{argmin}_{\beta} \mathbf{E}[(\mu(\vec{\mathbf{X}}) - \beta^T \vec{\mathbf{X}})^2] \quad (6)$$

$$= \mathbf{E}[\vec{\mathbf{X}} \vec{\mathbf{X}}^T]^{-1} \mathbf{E}[\mu(\vec{\mathbf{X}}) \vec{\mathbf{X}}]. \quad (7)$$

Equalities (5)-(7) follow from the normal equations for the population linear LS problem in (4):

$$\mathbf{E}[\vec{\mathbf{X}} \vec{\mathbf{X}}^T] \beta = \mathbf{E}[Y \vec{\mathbf{X}}] = \mathbf{E}[\mu(\vec{\mathbf{X}}) \vec{\mathbf{X}}]. \quad (8)$$

The vector $\beta = \beta(\mathbf{P})$ is a population parameter and the target of estimation. Unlike in a linear model, it only describes the best linear approximation to $\mu(\vec{\mathbf{X}})$, not $\mu(\vec{\mathbf{X}})$ itself. There is no generative model assumed here, hence $\beta = \beta(\mathbf{P})$ is just a statistical functional defined on multivariate distributions $\mathbf{P} = \mathbf{P}(dx_1, \dots, dx_p, dy)$.

- **Decompositions** of the response and response function:

Response:	Y		
Response function:	$\mu(\vec{X})$	$= \mathbf{E}[Y \vec{X}]$	
Population LS approximation:	$\beta^T \vec{X}$	(see (4)-(7))	
Nonlinearity, model bias:	$\eta(\vec{X})$	$= \mu(\vec{X}) - \beta^T \vec{X}$	(9)
Error, noise:	ϵ	$= Y - \mu(\vec{X})$	(10)
Population residual:	ξ	$= Y - \beta^T \vec{X} = \eta(\vec{X}) + \epsilon$	(11)
Error variance:	$\sigma^2(\vec{X})$	$= \mathbf{V}[Y \vec{X}] = \mathbf{E}[\epsilon^2 \vec{X}]$	(12)
Conditional MSE:	$\rho^2(\vec{X})$	$= \mathbf{E}[\xi^2 \vec{X}] = \eta^2(\vec{X}) + \sigma^2(\vec{X})$	(13)

The first three items are recapitulations, and the rest are definitions. Crucial are the *nonlinearity* or *model bias* function $\eta(\vec{X})$ to permit first order incorrectness, and the *error variance* function $\sigma^2(\vec{X})$ to permit second order model incorrectness. For lack of a better name we call ξ the *population residual*. The error ϵ satisfies $\mathbf{E}[\epsilon|\vec{X}] = 0$.

[Caution: Freedman (1981) and Mammen (1993) use “ ϵ ” to mean $\xi = \eta + \epsilon$ in our notation. They do not divide the population residual into error and nonlinearity.]

- **Orthogonalities/decorrelation** at the population: All of ξ , $\eta(\vec{X})$ and ϵ are orthogonal to the predictors:

$$\mathbf{E}[\xi \vec{X}] = \mathbf{E}[\eta(\vec{X}) \vec{X}] = \mathbf{E}[\epsilon \vec{X}] = \mathbf{0}. \quad (14)$$

The reason for ξ and $\eta(\vec{X})$ is that they are residuals in the population linear LS regressions of Y and $\mu(\vec{X})$, respectively, on \vec{X} . The error ϵ is orthogonal to \vec{X} because $\mathbf{E}[\epsilon|\vec{X}] = 0$. — As a consequence of the inclusion of an intercept entry “1” in \vec{X} , orthogonality of ξ and $\eta(\vec{X})$ to \vec{X} entails centering, whereas for ϵ it follows from $\mathbf{E}[\epsilon|\vec{X}] = 0$:

$$\mathbf{E}[\xi] = \mathbf{E}[\eta(\vec{X})] = \mathbf{E}[\epsilon] = 0. \quad (15)$$

The error ϵ satisfies much stronger orthogonalities: For any functions $f(\vec{X})$ and $g(\vec{X})$ (with suitable moments), $\epsilon f(\vec{X})$ is orthogonal to $g(\vec{X})$. This means $\mathbf{E}[(\epsilon f(\vec{X}))g(\vec{X})] = 0$, which is an immediate consequence of $\mathbf{E}[\epsilon|\vec{X}] = 0$. Yet ϵ is not generally independent of \vec{X} .

- **Adjustment formulas for the population:** We introduce the population version of “adjusted” predictors in order to express the multiple regression coefficient β_j as a simple regression coefficient. The adjusted predictor $X_{j\bullet}$ is defined as the residual of the regression of X_j (used as the response) on all other predictors. To this end we collect all other predictors by writing $\vec{X}_{\setminus j} = (1, X_1, \dots, X_{j-1}, X_{j+1}, \dots, X_p)^T$ and defining the p -vector $\beta_{j\bullet}$ by

$$\beta_{j\bullet} = \operatorname{argmin}_{\tilde{\beta}} \mathbf{E}[(X_j - \tilde{\beta}^T \vec{X}_{\setminus j})^2] = \mathbf{E}[\vec{X}_{\setminus j} \vec{X}_{\setminus j}^T]^{-1} \mathbf{E}[\vec{X}_{\setminus j} X_j].$$

The adjusted predictor is

$$X_{j\bullet} = X_j - \beta_{j\bullet}^T \vec{X}_{\setminus j}. \quad (16)$$

The representation of β_j as a simple regression coefficient is

$$\beta_j = \frac{\mathbf{E}[Y X_{j\bullet}]}{\mathbf{E}[X_{j\bullet}^2]} = \frac{\mathbf{E}[\mu(\vec{X}) X_{j\bullet}]}{\mathbf{E}[X_{j\bullet}^2]}. \quad (17)$$

This is a componentwise alternative to (5) above.

3.2 Notation for observational datasets

- When we use the term “observational data” we actually mean “cross-sectional data” consisting of i.i.d. cases $(X_{i,1}, \dots, X_{i,p}, Y_i)$ drawn from the multivariate distribution $\mathbf{P}(dx_1, \dots, dx_p, dy)$ ($i = 1, 2, \dots, N$). We collect the predictors of case i in a column $(p+1)$ -vector $\vec{\mathbf{X}}_i = (1, X_{i,1}, \dots, X_{i,p})^T$, prepended with 1 for an intercept.
- We stack the N samples to form random column N -vectors and a random predictor $N \times (p+1)$ -matrix:

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_N \end{bmatrix}, \quad \mathbf{X}_j = \begin{bmatrix} X_{1,j} \\ \vdots \\ X_{N,j} \end{bmatrix}, \quad \mathbf{X} = [\mathbf{1}, \mathbf{X}_1, \dots, \mathbf{X}_p] = \begin{bmatrix} \vec{\mathbf{X}}_1^T \\ \vdots \\ \vec{\mathbf{X}}_N^T \end{bmatrix}.$$

- Similarly we stack the values $\mu(\vec{\mathbf{X}}_i)$, $\eta(\vec{\mathbf{X}}_i)$, $\epsilon_i = Y_i - \mu(\vec{\mathbf{X}}_i)$, ξ_i , and $\sigma(\vec{\mathbf{X}}_i)$ to form random column N -vectors:

$$\boldsymbol{\mu} = \begin{bmatrix} \mu(\vec{\mathbf{X}}_1) \\ \vdots \\ \mu(\vec{\mathbf{X}}_N) \end{bmatrix}, \quad \boldsymbol{\eta} = \begin{bmatrix} \eta(\vec{\mathbf{X}}_1) \\ \vdots \\ \eta(\vec{\mathbf{X}}_N) \end{bmatrix}, \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_N \end{bmatrix}, \quad \boldsymbol{\xi} = \begin{bmatrix} \xi_1 \\ \vdots \\ \xi_N \end{bmatrix}, \quad \boldsymbol{\sigma} = \begin{bmatrix} \sigma(\vec{\mathbf{X}}_1) \\ \vdots \\ \sigma(\vec{\mathbf{X}}_N) \end{bmatrix}. \quad (18)$$

- The definitions (11) of ξ , (10) of ϵ , and (9) of $\eta(\vec{\mathbf{X}})$ translate to vectorized forms:

$$\boldsymbol{\xi} = \mathbf{Y} - \mathbf{X}\boldsymbol{\beta}, \quad \boldsymbol{\epsilon} = \mathbf{Y} - \boldsymbol{\mu}, \quad \boldsymbol{\eta} = \boldsymbol{\mu} - \mathbf{X}\boldsymbol{\beta}. \quad (19)$$

- It is important to keep in mind the distinction between population and sample properties, in particular:
 - The $(p+1)$ -vector $\boldsymbol{\beta} = \boldsymbol{\beta}(\mathbf{P})$ is *not* the LS estimate obtained from the sample; it is the population parameter and the target of estimation.
 - The N -vectors $\boldsymbol{\xi}$, $\boldsymbol{\epsilon}$ and $\boldsymbol{\eta}$ are *not* orthogonal to the predictor columns \mathbf{X}_j in the sample; in general $\langle \boldsymbol{\xi}, \mathbf{X}_j \rangle \neq 0$, $\langle \boldsymbol{\epsilon}, \mathbf{X}_j \rangle \neq 0$, $\langle \boldsymbol{\eta}, \mathbf{X}_j \rangle \neq 0$, even though the associated random variables are orthogonal to X_j in the population: $\mathbf{E}[\xi X_j] = \mathbf{E}[\epsilon X_j] = \mathbf{E}[\eta(\vec{\mathbf{X}})X_j] = 0$.

3.3 Quantities associated with LS estimation

- The **sample linear LS estimate** of $\boldsymbol{\beta}$ is the random column $(p+1)$ -vector

$$\hat{\boldsymbol{\beta}} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)^T = \operatorname{argmin}_{\boldsymbol{\beta}} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}. \quad (20)$$

Randomness stems from both the random response \mathbf{Y} and the random predictors in \mathbf{X} . Associated with $\hat{\boldsymbol{\beta}}$ are the following:

$$\text{Hat or projection matrix:} \quad \mathbf{H} = \mathbf{X}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \quad (21)$$

$$\text{LS fits:} \quad \hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{H}\mathbf{Y} \quad (22)$$

$$\text{Residuals:} \quad \mathbf{r} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} = (\mathbf{I} - \mathbf{H})\mathbf{Y} \quad (23)$$

[Note: Sample residual vector $\mathbf{r} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} \neq$ population residual vector $\boldsymbol{\xi} = \mathbf{Y} - \mathbf{X}\boldsymbol{\beta}$.]

- **Adjustment formulas for samples:** Again, we introduce adjusted variables to express a multiple regression as a collection of simple regressions. The sample version of adjustment is immediately useful in suggesting diagnostic plots, whereas the population version above served to provide targets of estimation. In a sample, we collect all predictor columns other than \mathbf{X}_j in a $N \times p$ matrix $\mathbf{X}_{\setminus j} = (\mathbf{1}, \dots, \mathbf{X}_{j-1}, \mathbf{X}_{j+1}, \dots)$, and we define

$$\hat{\beta}_{j\bullet} = \operatorname{argmin}_{\tilde{\beta}} \|\mathbf{X}_j - \mathbf{X}_{\setminus j}\tilde{\beta}\|^2 = (\mathbf{X}_{\setminus j}^T \mathbf{X}_{\setminus j})^{-1} \mathbf{X}_{\setminus j}^T \mathbf{X}_j.$$

The sample-adjusted predictor is

$$\mathbf{X}_{j\bullet} = \mathbf{X}_j - \hat{\beta}_{j\bullet}^T \mathbf{X}_{\setminus j} = (\mathbf{I} - \mathbf{H}_{\setminus j}) \mathbf{X}_j. \quad (24)$$

where $\mathbf{H}_{\setminus j} = \mathbf{X}_{\setminus j}(\mathbf{X}_{\setminus j}^T \mathbf{X}_{\setminus j})^{-1} \mathbf{X}_{\setminus j}^T$ is the associated projection or hat matrix. The j 'th slope estimate of the multiple linear regression of \mathbf{Y} on $\mathbf{X}_1, \dots, \mathbf{X}_p$ can then be expressed in the well-known manner as the slope estimate of the simple linear regression (without intercept) of \mathbf{Y} on $\mathbf{X}_{j\bullet}$:

$$\hat{\beta}_j = \frac{\langle \mathbf{Y}, \mathbf{X}_{j\bullet} \rangle}{\|\mathbf{X}_{j\bullet}\|^2}. \quad (25)$$

[An important distinction is that the vector $\mathbf{X}_{j\bullet}$ is sample-adjusted, whereas the random variable $X_{j\bullet}$ is population-adjusted. This makes the notation “ $X_{i,j\bullet}$ ” ambiguous: it could refer to the i 'th component of the sample-adjusted vector $\mathbf{X}_{j\bullet}$ or to the i 'th case drawn from the population-adjusted random variable $X_{j\bullet}$. These are not the same. We will therefore avoid using “ $X_{i,j\bullet}$ ”, which is possible without difficulty.]

4 Decomposition of the LS Estimate according to its Sources of Variation

When the predictors are random and the linear model is treated as an approximation rather than a generative truth, the sampling variation of the LS estimate $\hat{\beta}$ can be additively decomposed into two components: one component due to error, and another component due to nonlinearity interacting with randomness of the predictors. This decomposition is a direct reflection of the decomposition of the population residual into error and nonlinearity, $\boldsymbol{\xi} = \boldsymbol{\epsilon} + \boldsymbol{\eta}$, according to (11). In this section we state asymptotic normality for each part of the decomposition. The relevance of the decomposition is that it explains what the paired bootstrap estimates, while the associated asymptotic normalities are necessary to justify the paired bootstrap.

In the conventional analysis of linear models, which is conditional on \mathbf{X} , the target of estimation is $\mathbf{E}[\hat{\beta}|\mathbf{X}]$. When \mathbf{X} is treated as random and nonlinearity is allowed, the target of estimation is the population LS solution β as defined in (4)-(7). In this case, $\mathbf{E}[\hat{\beta}|\mathbf{X}]$ is a random variable that sits between $\hat{\beta}$ and β :

$$\hat{\beta} - \beta = (\hat{\beta} - \mathbf{E}[\hat{\beta}|\mathbf{X}]) + (\mathbf{E}[\hat{\beta}|\mathbf{X}] - \beta) \quad (26)$$

This decomposition corresponds to the decomposition $\boldsymbol{\xi} = \boldsymbol{\epsilon} + \boldsymbol{\eta}$ as the following lemma shows. The lemma gives both vectorized and componentwise versions, as the former provide easy proofs, while the latter will result in simple diagnostics.

Lemma and Definition: *The following quantities will be called “Estimation Offsets” or “EO” for short, and*

they will be prefixed as follows:

$$\text{Total EO:} \quad \hat{\beta} - \beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\xi}, \quad \hat{\beta}_j - \beta_j = \frac{\langle \mathbf{X}_{j\bullet}, \boldsymbol{\xi} \rangle}{\|\mathbf{X}_{j\bullet}\|^2}. \quad (27)$$

$$\text{Error EO:} \quad \hat{\beta} - \mathbf{E}[\hat{\beta}|\mathbf{X}] = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}, \quad \hat{\beta}_j - \mathbf{E}[\hat{\beta}_j|\mathbf{X}] = \frac{\langle \mathbf{X}_{j\bullet}, \boldsymbol{\epsilon} \rangle}{\|\mathbf{X}_{j\bullet}\|^2}. \quad (28)$$

$$\text{Nonlinearity EO:} \quad \mathbf{E}[\hat{\beta}|\mathbf{X}] - \beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\eta}, \quad \mathbf{E}[\hat{\beta}_j|\mathbf{X}] - \beta_j = \frac{\langle \mathbf{X}_{j\bullet}, \boldsymbol{\eta} \rangle}{\|\mathbf{X}_{j\bullet}\|^2}. \quad (29)$$

Proof: It is sufficient to consider the vectorized versions. The lemma is a consequence of similar equations for \mathbf{Y} , $\boldsymbol{\mu}$ and $\mathbf{X}\boldsymbol{\beta}$:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad (30)$$

$$\mathbf{E}[\hat{\beta}|\mathbf{X}] = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\mu} \quad (31)$$

$$\beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\boldsymbol{\beta}) \quad (32)$$

Of these equations the first is the definition of $\hat{\beta}$, the middle follows from the first using $\mathbf{E}[Y_i|\vec{\mathbf{X}}_i] = \mu(\vec{\mathbf{X}}_i)$, and the last is a tautology. The equations of the lemma follow by forming the pairwise differences (30)-(32), (30)-(31), and (31)-(32), respectively, and recalling the vectorized definitions (19): $\boldsymbol{\xi} = \mathbf{Y} - \mathbf{X}\boldsymbol{\beta}$, $\boldsymbol{\epsilon} = \mathbf{Y} - \boldsymbol{\mu}$, $\boldsymbol{\eta} = \boldsymbol{\mu} - \mathbf{X}\boldsymbol{\beta}$. \square

The EOs are not actionable quantities in real data analysis, but they will be useful in comparing unconditional (bootstrap) and conditional (linear) inference.

Proposition: *The three EOs follow central limit theorems:*

$$N^{1/2} (\hat{\beta} - \beta) \xrightarrow{\mathcal{D}} \mathcal{N} \left(\mathbf{0}, \mathbf{E}[\vec{\mathbf{X}} \vec{\mathbf{X}}^T]^{-1} \mathbf{E}[\rho^2(\vec{\mathbf{X}}) \vec{\mathbf{X}} \vec{\mathbf{X}}^T] \mathbf{E}[\vec{\mathbf{X}} \vec{\mathbf{X}}^T]^{-1} \right) \quad (33)$$

$$N^{1/2} (\hat{\beta} - \mathbf{E}[\hat{\beta}|\mathbf{X}]) \xrightarrow{\mathcal{D}} \mathcal{N} \left(\mathbf{0}, \mathbf{E}[\vec{\mathbf{X}} \vec{\mathbf{X}}^T]^{-1} \mathbf{E}[\sigma^2(\vec{\mathbf{X}}) \vec{\mathbf{X}} \vec{\mathbf{X}}^T] \mathbf{E}[\vec{\mathbf{X}} \vec{\mathbf{X}}^T]^{-1} \right) \quad (34)$$

$$N^{1/2} (\mathbf{E}[\hat{\beta}|\mathbf{X}] - \beta) \xrightarrow{\mathcal{D}} \mathcal{N} \left(\mathbf{0}, \mathbf{E}[\vec{\mathbf{X}} \vec{\mathbf{X}}^T]^{-1} \mathbf{E}[\eta^2(\vec{\mathbf{X}}) \vec{\mathbf{X}} \vec{\mathbf{X}}^T] \mathbf{E}[\vec{\mathbf{X}} \vec{\mathbf{X}}^T]^{-1} \right) \quad (35)$$

Proof: It is sufficient to prove the first statement, say, as the others follow the same way.

$$\begin{aligned} N^{1/2} (\hat{\beta} - \beta) &= \left(\frac{1}{N} \mathbf{X}^T \mathbf{X} \right)^{-1} \left(\frac{1}{N^{1/2}} \mathbf{X}^T \boldsymbol{\xi} \right) \\ &\xrightarrow{\mathcal{D}} \mathbf{E}[\vec{\mathbf{X}} \vec{\mathbf{X}}^T]^{-1} \mathcal{N} \left(\mathbf{0}, \mathbf{E}[\rho^2(\vec{\mathbf{X}}) \vec{\mathbf{X}} \vec{\mathbf{X}}^T] \right) \\ &= \mathcal{N} \left(\mathbf{0}, \mathbf{E}[\vec{\mathbf{X}} \vec{\mathbf{X}}^T]^{-1} \mathbf{E}[\rho^2(\vec{\mathbf{X}}) \vec{\mathbf{X}} \vec{\mathbf{X}}^T] \mathbf{E}[\vec{\mathbf{X}} \vec{\mathbf{X}}^T]^{-1} \right), \end{aligned}$$

We used $\mathbf{E}[\boldsymbol{\xi}] = \mathbf{0}$ ($= \mathbf{E}[\boldsymbol{\epsilon}] = \mathbf{E}[\boldsymbol{\eta}]$) according to (15). It follows that

$$\mathbf{V}[\boldsymbol{\xi} \vec{\mathbf{X}}] = \mathbf{E}[\boldsymbol{\xi}^2 \vec{\mathbf{X}} \vec{\mathbf{X}}^T] = \mathbf{E}[\mathbf{E}[\boldsymbol{\xi}^2 | \vec{\mathbf{X}}] \vec{\mathbf{X}} \vec{\mathbf{X}}^T] = \mathbf{E}[\rho^2(\vec{\mathbf{X}}) \vec{\mathbf{X}} \vec{\mathbf{X}}^T] \quad \square$$

The asymptotic variance/covariance matrices in the proposition have the well-known sandwich or butterfly form. Asymptotic normality can also be expressed for each $\hat{\beta}_j$ separately using population adjustment:

Corollary: *Each component of an estimation offset follows a central limit theorem.*

$$N^{1/2}(\hat{\beta}_j - \beta_j) \xrightarrow{\mathcal{D}} \mathcal{N}\left(0, \frac{\mathbf{E}[\rho^2(\vec{\mathbf{X}})X_{j\bullet}^2]}{\mathbf{E}[X_{j\bullet}^2]^2}\right) \quad (36)$$

$$N^{1/2}(\hat{\beta}_j - \mathbf{E}[\hat{\beta}_j|\mathbf{X}]) \xrightarrow{\mathcal{D}} \mathcal{N}\left(0, \frac{\mathbf{E}[\sigma^2(\vec{\mathbf{X}})X_{j\bullet}^2]}{\mathbf{E}[X_{j\bullet}^2]^2}\right) \quad (37)$$

$$N^{1/2}(\mathbf{E}[\hat{\beta}_j|\mathbf{X}] - \beta_j) \xrightarrow{\mathcal{D}} \mathcal{N}\left(0, \frac{\mathbf{E}[\eta^2(\vec{\mathbf{X}})X_{j\bullet}^2]}{\mathbf{E}[X_{j\bullet}^2]^2}\right) \quad (38)$$

Proof: By the proposition above we know that each component is marginally asymptotically normal. Hence we only need to verify the form of the asymptotic variance. This in turn is achieved by specializing the asymptotic variance/covariance matrix of the proposition above to the single predictor case with $X_{j\bullet}$ as the random variable underlying the predictor. \square

Corollary: *The unconditional or marginal bias of $\hat{\beta}_j$ vanishes faster than $N^{-1/2}$:*

$$N^{1/2}(\mathbf{E}[\hat{\beta}_j] - \beta_j) \rightarrow 0 \quad (N \rightarrow \infty). \quad (39)$$

Proof: This follows from (36). \square

[xxxxxxxxx Preliminary note – might go elsewhere:

In econometrics it has long been exploited that (33) suggests a standard error estimate — the so-called “sandwich estimator” — that is asymptotically correct even when the linear model is incorrect. It is obtained by estimating $\rho^2(\vec{\mathbf{X}})$ with $(Y - \hat{\beta}^T \vec{\mathbf{X}})^2$, which is asymptotically unbiased: $\mathbf{E}[(Y - \hat{\beta}^T \vec{\mathbf{X}})^2 | \vec{\mathbf{X}}] \rightarrow \rho^2(\vec{\mathbf{X}})$. The sandwich estimator has come under criticism for being inefficient (Kauermann and Carroll 2001). This inefficiency problem might get corrected by drawing on the adjustment version (36) and noting that

$$\mathbf{E}[\rho^2(\vec{\mathbf{X}})X_{j\bullet}^2] = \mathbf{E}[\mathbf{E}[\rho^2(\vec{\mathbf{X}}) | X_{j\bullet}^2] X_{j\bullet}^2].$$

The expression $\mathbf{E}[\rho^2(\vec{\mathbf{X}}) | X_{j\bullet}^2]$ can be estimated nonparametrically in the sample by smoothing $\mathbf{r}^2 = (\mathbf{Y} - \mathbf{X}\hat{\beta})^2$ (= the vector of squared residuals) against $\mathbf{X}_{j\bullet}^2$ (= the vector of squared adjusted predictor values). Such smoothing should remove some of the noisiness and hence inefficiency of the sandwich estimator.]

5 Standard Errors of Regression Coefficients

We compare conventional standard errors with the actual (marginal, unconditional) standard errors. The latter correctly include the additional variation due to the nonlinearity and randomness of \mathbf{X} . In what follows we distinguish between (squared) *standard error* $\mathbf{V}[\hat{\beta}_j]$, *conditional standard error* $\mathbf{V}[\hat{\beta}_j|\mathbf{X}]$, and *standard error estimate* $\hat{\mathbf{V}}[\hat{\beta}_j]$.

1. The (squared) *conventional standard error* of $\hat{\beta}_j$ is estimated by $\hat{\sigma}^2 / \|\mathbf{X}_{j\bullet}\|^2$ where one uses the usual estimate of σ^2 from the linear model as if the model were correct: $\hat{\sigma}^2 = \frac{1}{N-p-1} \|\mathbf{r}\|^2$. We calculate its conditional mean allowing for both nonlinearity and heteroskedasticity (see Appendix B.1):

$$\mathbf{E}[\hat{\sigma}^2 | \mathbf{X}] = \frac{1}{N-p-1} (\text{tr}((\mathbf{I} - \mathbf{H})\mathbf{D}_{\sigma^2}) + \|(\mathbf{I} - \mathbf{H})\boldsymbol{\eta}\|^2), \quad (40)$$

where \mathbf{D}_{σ^2} is the diagonal matrix with the values $\sigma^2(\vec{\mathbf{X}}_i)$ in its diagonal. Thus the conventional standard error estimate $\hat{\sigma}^2/\|\mathbf{X}_{j\bullet}\|^2$ is \mathbf{X} -conditionally an unbiased estimate of the following:

$$\text{SE}_{\text{conv}}^2[\hat{\beta}_j|\mathbf{X}] := \mathbf{E} \left[\frac{\hat{\sigma}^2}{\|\mathbf{X}_{j\bullet}\|^2} \middle| \mathbf{X} \right] = \frac{\frac{1}{N-p-1} \text{tr}((\mathbf{I} - \mathbf{H})\mathbf{D}_{\sigma^2})}{\|\mathbf{X}_{j\bullet}\|^2} + \frac{\frac{1}{N-p-1} \|(\mathbf{I} - \mathbf{H})\boldsymbol{\eta}\|^2}{\|\mathbf{X}_{j\bullet}\|^2} \quad (41)$$

On the right hand side, the first term accounts for error (possibly heteroscedastic), the second term for nonlinearity. Under homoscedasticity ($\mathbf{D}_{\sigma^2} = \sigma^2\mathbf{I}$) and linearity ($\boldsymbol{\eta} = \mathbf{0}$) it is true that $\text{SE}_{\text{conv}}^2[\hat{\beta}_j|\mathbf{X}] = \mathbf{V}[\hat{\beta}_j|\mathbf{X}]$, hence this is an appropriate conditional squared standard error for $\mathbf{E}[\hat{\beta}_j|\mathbf{X}]$, which under linearity equals β_j . In the presence of nonlinearity and/or heteroscedasticity, however, $\text{SE}_{\text{conv}}^2[\hat{\beta}_j|\mathbf{X}]$ is not a correct squared standard error for β_j , not even conditionally for $\mathbf{E}[\hat{\beta}_j|\mathbf{X}]$, as the next item will show.

2. The (squared) unconditional or *marginal standard error* is $\mathbf{V}[\hat{\beta}_j]$. Its theoretical relevance is that it is the true standard error that accounts for sampling variability due to error *and* due to randomness of the predictors. Its practical importance is that it can and should be estimated by sandwich estimators or the paired bootstrap. It can be decomposed as follows:

$$\text{SE}_{\text{marg}}^2[\hat{\beta}_j] := \mathbf{V}[\hat{\beta}_j] = \mathbf{E}[\mathbf{V}[\hat{\beta}_j|\mathbf{X}]] + \mathbf{V}[\mathbf{E}[\hat{\beta}_j|\mathbf{X}]] \quad (42)$$

$$= \mathbf{E} \left[\frac{\mathbf{X}_{j\bullet}^T \mathbf{D}_{\sigma^2} \mathbf{X}_{j\bullet}}{\|\mathbf{X}_{j\bullet}\|^4} \right] + \mathbf{V} \left[\frac{\langle \mathbf{X}_{j\bullet}, \boldsymbol{\eta} \rangle}{\|\mathbf{X}_{j\bullet}\|^2} \right]. \quad (43)$$

In the last equality we drew on (25).

Thus both the conventional and the marginal standard error can be additively decomposed into contributions due to $\sigma^2(\vec{\mathbf{X}})$ (error) and $\eta^2(\vec{\mathbf{X}})$ (nonlinearity). For a simplified asymptotic analysis, we form the scaled limits:

Proposition and Definition: *Defining $\text{SE}_{\text{conv}}^2[\hat{\beta}_j] := \mathbf{E}[\text{SE}_{\text{conv}}^2[\hat{\beta}_j|\mathbf{X}]]$ and $\text{SE}_{\text{marg}}^2[\hat{\beta}_j]$ as in (42), we have:*

$$N \text{SE}_{\text{conv}}^2(\hat{\beta}_j) \xrightarrow{N \rightarrow \infty} \underbrace{\frac{\mathbf{E}[\sigma^2(\vec{\mathbf{X}})]}{\mathbf{E}[X_{j\bullet}^2]}}_{\mathbf{AV}_{\text{conv}}^{(j)}(\sigma^2(\vec{\mathbf{X}}))} + \underbrace{\frac{\mathbf{E}[\eta^2(\vec{\mathbf{X}})]}{\mathbf{E}[X_{j\bullet}^2]}}_{\mathbf{AV}_{\text{conv}}^{(j)}(\eta^2(\vec{\mathbf{X}}))} = \underbrace{\frac{\mathbf{E}[\rho^2(\vec{\mathbf{X}})]}{\mathbf{E}[X_{j\bullet}^2]}}_{\mathbf{AV}_{\text{conv}}^{(j)}(\rho^2(\vec{\mathbf{X}}))} \quad (44)$$

$$N \text{SE}_{\text{marg}}^2(\hat{\beta}_j) \xrightarrow{N \rightarrow \infty} \underbrace{\frac{\mathbf{E}[\sigma^2(\vec{\mathbf{X}})X_{j\bullet}^2]}{\mathbf{E}[X_{j\bullet}^2]^2}}_{\mathbf{AV}_{\text{marg}}^{(j)}(\sigma^2(\vec{\mathbf{X}}))} + \underbrace{\frac{\mathbf{E}[\eta^2(\vec{\mathbf{X}})X_{j\bullet}^2]}{\mathbf{E}[X_{j\bullet}^2]^2}}_{\mathbf{AV}_{\text{marg}}^{(j)}(\eta^2(\vec{\mathbf{X}}))} = \underbrace{\frac{\mathbf{E}[\rho^2(\vec{\mathbf{X}})X_{j\bullet}^2]}{\mathbf{E}[X_{j\bullet}^2]^2}}_{\mathbf{AV}_{\text{marg}}^{(j)}(\rho^2(\vec{\mathbf{X}}))} \quad (45)$$

The proofs of (44) and (45) are obtained from the sum representations in (41) and (43). For details of (44) see Appendix B.3; essentially \mathbf{H} and $p+1$ in (41) become negligible for p fixed as $N \rightarrow \infty$. The limits in (45) are the asymptotic variances (36)-(38) of the preceding section.

Structurally all three terms $\mathbf{AV}_{\text{conv}}^{(j)}(\dots)$ have the same form, and so do the three terms $\mathbf{AV}_{\text{marg}}^{(j)}(\dots)$, justifying the notation. Only $\mathbf{AV}_{\text{marg}}^{(j)}(\dots)$ represent inferentially valid asymptotic variances, while $\mathbf{AV}_{\text{conv}}^{(j)}(\dots)$ are generally valid only when $\sigma^2(\vec{\mathbf{X}}) \equiv \sigma^2$ and $\eta \equiv 0$.

6 Asymptotic Comparison of Conventional and Marginal Standard Errors

We show that the conventional asymptotic variances can be too small or too large to unlimited degrees compared to the proper marginal asymptotic variances. A comparison of asymptotic variances can be done separately for $\sigma^2(\vec{X})$, $\eta^2(\vec{X})$ and $\rho^2(\vec{X})$. To this end we form the ratios $\mathbf{RAV}_j(\dots)$ as follows:

Definition and Lemma: *Ratios of marginal and conventional Asymptotic Variances*

$$\mathbf{RAV}_j(\sigma^2(\vec{X})) := \frac{\mathbf{AV}_{\text{marg}}^{(j)}(\sigma^2(\vec{X}))}{\mathbf{AV}_{\text{conv}}^{(j)}(\sigma^2(\vec{X}))} = \frac{\mathbf{E}[\sigma^2(\vec{X})X_{j\bullet}^2]}{\mathbf{E}[\sigma^2(\vec{X})]\mathbf{E}[X_{j\bullet}^2]} \quad (46)$$

$$\mathbf{RAV}_j(\eta^2(\vec{X})) := \frac{\mathbf{AV}_{\text{marg}}^{(j)}(\eta^2(\vec{X}))}{\mathbf{AV}_{\text{conv}}^{(j)}(\eta^2(\vec{X}))} = \frac{\mathbf{E}[\eta^2(\vec{X})X_{j\bullet}^2]}{\mathbf{E}[\eta^2(\vec{X})]\mathbf{E}[X_{j\bullet}^2]} \quad (47)$$

$$\mathbf{RAV}_j(\rho^2(\vec{X})) := \frac{\mathbf{AV}_{\text{marg}}^{(j)}(\rho^2(\vec{X}))}{\mathbf{AV}_{\text{conv}}^{(j)}(\rho^2(\vec{X}))} = \frac{\mathbf{E}[\rho^2(\vec{X})X_{j\bullet}^2]}{\mathbf{E}[\rho^2(\vec{X})]\mathbf{E}[X_{j\bullet}^2]} \quad (48)$$

The second equality on each line follows from (44) and (45). The ratios express by how much the conventional asymptotic variances need to be multiplied to match the proper marginal asymptotic variances. Among the three ratios the relevant one for the overall comparison of conventional and marginal inference is $\mathbf{RAV}_j(\rho^2(\vec{X}))$. For example, if $\mathbf{RAV}_j(\rho^2(\vec{X})) = 4$, say, then, for large sample sizes, the correct marginal standard error of $\hat{\beta}_j$ is about twice as large as the incorrect conventional standard error. In general \mathbf{RAV}_j expresses the following:

- if $\mathbf{RAV}_j(\rho^2(\vec{X})) = 1$, the conventional standard error for $\hat{\beta}_j$ is asymptotically correct;
- if $\mathbf{RAV}_j(\rho^2(\vec{X})) > 1$, the conventional standard error for β_j is asymptotically too small/optimistic;
- if $\mathbf{RAV}_j(\rho^2(\vec{X})) < 1$, the conventional standard error for β_j is asymptotically too large/pessimistic.

The ratios $\mathbf{RAV}_j(\sigma^2(\vec{X}))$ and $\mathbf{RAV}_j(\eta^2(\vec{X}))$ express the degrees to which heteroscedasticity and/or nonlinearity contribute asymptotically to the defects of conventional standard errors. Structurally, the three ratios are inner products between the normalized squared quantities

$$\frac{\sigma^2(\vec{X})}{\mathbf{E}[\sigma^2(\vec{X})]}, \quad \frac{\eta^2(\vec{X})}{\mathbf{E}[\eta^2(\vec{X})]}, \quad \frac{\rho^2(\vec{X})}{\mathbf{E}[\rho^2(\vec{X})]}$$

on the one hand, and the normalized squared adjusted predictor

$$\frac{X_{j\bullet}^2}{\mathbf{E}[X_{j\bullet}^2]}$$

on the other hand. These inner products, however, are *not* correlations, and they are *not* bounded by +1. To analyze the range of \mathbf{RAV}_j and to take advantage of the identical form of \mathbf{RAV}_j in all three cases, we will write $f(\vec{X})$ for any of $\sigma^2(\vec{X})$, $\eta^2(\vec{X})$ and $\rho^2(\vec{X})$, and analyze the range of $\mathbf{RAV}_j(f^2(\vec{X}))$ over scenarios $f(\vec{X})$. It is easy to see that generally

$$\sup_f \mathbf{RAV}_j(f^2(\vec{X})) = \infty, \quad \inf_f \mathbf{RAV}_j(f^2(\vec{X})) = 0, \quad (49)$$

where $f(\vec{X})$ ranges over functions with suitable moments such that $\mathbf{RAV}_j(f^2(\vec{X}))$ exists. The limits (49) become fairly obvious from the following equivalent form of \mathbf{RAV}_j :

$$\mathbf{RAV}_j(f^2(\vec{X})) = \frac{\mathbf{E}[\mathbf{E}[f^2(\vec{X}) | X_{j\bullet}^2] X_{j\bullet}^2]}{\mathbf{E}[f^2(\vec{X})]\mathbf{E}[X_{j\bullet}^2]} \quad (50)$$

Thus $\mathbf{RAV}_j(f^2(\vec{X}))$ attains

- a large value if $\mathbf{E}[f^2(\vec{X}) | X_{j\bullet}^2]$ and $X_{j\bullet}^2$ are positively associated, which is the case for example if $\mathbf{E}[f^2(\vec{X}) | X_{j\bullet}^2]$ is a monotone increasing function of $X_{j\bullet}^2$;
- a value near zero if $\mathbf{E}[f^2(\vec{X}) | X_{j\bullet}^2]$ and $X_{j\bullet}^2$ are negatively associated, which is the case for example if $\mathbf{E}[f^2(\vec{X}) | X_{j\bullet}^2]$ is a monotone decreasing function of $X_{j\bullet}^2$.

We will illustrate the above with some simple scenario constructions.

Consider first nonlinearities $\eta(\vec{X})$: We construct a one-parameter family of nonlinearities $\eta_t(\vec{X})$ for which $\sup_t \mathbf{RAV}_j(\eta_t^2) = \infty$ and $\inf_t \mathbf{RAV}_j(\eta_t^2) = 0$. Generally in the construction of examples, it must be kept in mind that nonlinearities are orthogonal to (adjusted for) all other predictors: $\mathbf{E}[\eta(\vec{X})\vec{X}] = \mathbf{0}$. To avoid un insightful complications arising from adjustment due to complex dependencies among the predictors, we construct an example for simple linear regression with a single predictor $X_1 = X$ and an intercept $X_0 = 1$. W.l.o.g. we will further assume that X_1 is centered (population adjusted for X_0 , so that $X_{1\bullet} = X_1$) and standardized. In what follows we write X instead of X_1 , and the assumptions are $\mathbf{E}[X] = 0$ and $\mathbf{E}[X^2] = 1$. To make the example as simple as possible we adopt some additional assumptions on the distribution of X :

Proposition: *Define a one-parameter family of nonlinearities as follows:*

$$\eta_t(X) = \frac{1_{[|X|>t]} - p(t)}{\sqrt{p(t)(1-p(t))}}, \quad \text{where } p(t) = \mathbf{P}[|X| > t], \quad (51)$$

and we assume that $p(t) > 0$ and $1 - p(t) > 0 \quad \forall t > 0$. Assume further that the distribution of X is symmetric about 0, so that $\mathbf{E}[\eta_t(X)X] = 0$. Then we have:

$$\lim_{t \uparrow \infty} \mathbf{RAV}(\eta_t^2) = \infty;$$

$$\lim_{t \downarrow 0} \mathbf{RAV}(\eta_t^2) = 0 \text{ if the distribution of } X \text{ has no atom at the origin: } \mathbf{P}[X = 0] = 0.$$

(See Appendix B.5 for the simple proof). By construction these nonlinearities are centered and standardized, $\mathbf{E}[\eta_t(X)] = 0$ and $\mathbf{E}[\eta_t(X)^2] = 1$. They are also orthogonal to X , $\mathbf{E}[\eta_t(X)X] = 0$, due to the assumed symmetry of the distribution of X , $\mathbf{P}[X > t] = \mathbf{P}[X < -t]$, and the symmetry of the nonlinearities, $\eta_t(-X) = \eta_t(X)$.

Consider next heteroscedastic error variances $\sigma^2(\vec{X})$: The above construction for nonlinearities can be re-used. As with nonlinearities, for $\mathbf{RAV}(\sigma_t^2(X))$ to rise with no bound, the conditional error variance $\sigma_t^2(X)$ needs to place its large values in the unbounded tail of the distribution of X . For $\mathbf{RAV}(\sigma_t^2(X))$ to reach down to zero, $\sigma_t^2(X)$ needs to place its large values in the center of the distribution of X .

Proposition: *Define a one-parameter family of heteroscedastic error variances as follows:*

$$\sigma_t^2(X) = \frac{(1_{[|X|>t]} - p(t))^2}{p(t)(1-p(t))}, \quad \text{where } p(t) = \mathbf{P}[|X| > t], \quad (52)$$

and we assume that $p(t) > 0$ and $1 - p(t) > 0 \quad \forall t > 0$. Then we have:

$$\lim_{t \uparrow \infty} \mathbf{RAV}(\sigma_t^2) = \infty;$$

$\lim_{t \downarrow 0} \mathbf{RAV}(\sigma_t^2) = 0$ if the distribution of X has no atom at the origin: $\mathbf{P}[X = 0] = 0$.

An important difference between $\eta^2(\vec{\mathbf{X}})$ and $\sigma^2(\vec{\mathbf{X}})$ is that nonlinearities are constrained by orthogonalities to the predictors, whereas conditional error variances are not.

In the end, of course, it is $\rho^2(\vec{\mathbf{X}}) = \sigma^2(\vec{\mathbf{X}}) + \eta^2(\vec{\mathbf{X}})$ that determines the gap between conventional and marginal inferences. The following lemma can be used to examine when $\mathbf{RAV}_j(\rho^2(\vec{\mathbf{X}}))$ approaches ∞ or 0:

Lemma: *Define weights*

$$w_\sigma = \frac{\mathbf{E}[\sigma^2(\vec{\mathbf{X}})]}{\mathbf{E}[\rho^2(\vec{\mathbf{X}})]}, \quad w_\eta = \frac{\mathbf{E}[\eta^2(\vec{\mathbf{X}})]}{\mathbf{E}[\rho^2(\vec{\mathbf{X}})]}, \quad (53)$$

so that $w_\sigma + w_\eta = 1$. Then

$$\mathbf{RAV}_j(\rho^2(\vec{\mathbf{X}})) = w_\sigma \mathbf{RAV}_j(\sigma^2(\vec{\mathbf{X}})) + w_\eta \mathbf{RAV}_j(\eta^2(\vec{\mathbf{X}})). \quad (54)$$

Corollary: *If $\rho_t^2(\vec{\mathbf{X}}) = \sigma_t^2(\vec{\mathbf{X}}) + \eta_t^2(\vec{\mathbf{X}})$ is a one-parameter family of scenarios, and if $w_{\sigma_t} = w_\sigma > 0$ and $w_{\eta_t} = w_\eta > 0$ from (53) are fixed and strictly positive, then we have*

$$\lim_t \mathbf{RAV}_j(\rho_t^2(\vec{\mathbf{X}})) = \infty \text{ if } \lim_t \mathbf{RAV}_j(\sigma_t^2(\vec{\mathbf{X}})) = \infty \text{ or } \lim_t \mathbf{RAV}_j(\eta_t^2(\vec{\mathbf{X}})) = \infty;$$

$$\lim_t \mathbf{RAV}_j(\rho_t^2(\vec{\mathbf{X}})) = 0 \text{ if both } \lim_t \mathbf{RAV}_j(\sigma_t^2(\vec{\mathbf{X}})) = 0 \text{ and } \lim_t \mathbf{RAV}_j(\eta_t^2(\vec{\mathbf{X}})) = 0.$$

The corollary adds evidence that large values $\mathbf{RAV}(\rho^2(\vec{\mathbf{X}})) \gg 1$ seem more easily possible than small values $\mathbf{RAV}(\rho^2(\vec{\mathbf{X}})) \ll 1$, hence conventional inference seems more likely to be too optimistic/liberal than too pessimistic/conservative. The reason is that heteroscedasticity $\sigma^2(\vec{\mathbf{X}})$ or nonlinearity $\eta^2(\vec{\mathbf{X}})$ can each individually cause $\mathbf{RAV}(\rho^2(\vec{\mathbf{X}}))$ to explode, but each can provide a floor that prevents $\mathbf{RAV}(\rho^2(\vec{\mathbf{X}}))$ from dropping to 0.

7 Submodels

Following the notation in the PoSI article, we consider submodels given by indices $\mathbf{M} = \{j_1, \dots, j_m\}$ with predictor submatrices $\mathbf{X}_\mathbf{M} = (\mathbf{X}_{j_1}, \dots, \mathbf{X}_{j_m})$ and associated vectors of parameters and their LS estimates $\beta_\mathbf{M}$ and $\hat{\beta}_\mathbf{M}$, respectively. We know the asymptotic normality of $\hat{\beta}_\mathbf{M}$ for each \mathbf{M} , but we are now also interested in the joint asymptotic normality of $(\hat{\beta}_\mathbf{M}^T, \hat{\beta}_{\mathbf{M}'}^T)^T$ for two submodels \mathbf{M} and \mathbf{M}' . While the asymptotic variance for $\hat{\beta}_\mathbf{M}$ is

$$\mathbf{AV}(\hat{\beta}_\mathbf{M}) = \mathbf{E}[\vec{\mathbf{X}}_\mathbf{M} \vec{\mathbf{X}}_\mathbf{M}^T]^{-1} \mathbf{E}[\xi_\mathbf{M}^2 \vec{\mathbf{X}}_\mathbf{M} \vec{\mathbf{X}}_\mathbf{M}^T] \mathbf{E}[\vec{\mathbf{X}}_\mathbf{M} \vec{\mathbf{X}}_\mathbf{M}^T]^{-1},$$

using the submodel-specific ‘‘population residual’’ variable $\xi_\mathbf{M} = Y - \vec{\mathbf{X}}_\mathbf{M}^T \beta_\mathbf{M}$, the cross-asymptotic covariance between $\hat{\beta}_\mathbf{M}$ and $\hat{\beta}_{\mathbf{M}'}$ is

$$\mathbf{AV}(\hat{\beta}_\mathbf{M}, \hat{\beta}_{\mathbf{M}'}) = \mathbf{E}[\vec{\mathbf{X}}_\mathbf{M} \vec{\mathbf{X}}_\mathbf{M}^T]^{-1} \mathbf{E}[\xi_\mathbf{M} \xi_{\mathbf{M}'} \vec{\mathbf{X}}_\mathbf{M} \vec{\mathbf{X}}_{\mathbf{M}'}^T] \mathbf{E}[\vec{\mathbf{X}}_{\mathbf{M}'} \vec{\mathbf{X}}_{\mathbf{M}'}^T]^{-1}.$$

A The Meaning of Regression Slopes as Averages of Case-Based Slopes

In this section we give simple interpretations to slopes when the actual response function is not linear in the predictors. We first condition on the design and then consider the random design case. It will also be sufficient to consider simple linear regression only, keeping in mind that all formulas apply to multiple regression by replacing the single predictor with the j 'th predictor adjusted for all other predictors.

For a sample (Y_i, X_i) ($i = 1, \dots, N$) of single-predictor/response data, the LS slope estimate can be represented as follows:

$$\hat{\beta} = \frac{\sum_i (Y_i - \bar{Y})(X_i - \bar{X})}{\sum_k (X_k - \bar{X})^2} = \frac{\sum_i \frac{Y_i - \bar{Y}}{X_i - \bar{X}} (X_i - \bar{X})^2}{\sum_k (X_k - \bar{X})^2} = \sum_i w_i b_i, \quad (55)$$

where

$$b_i = \frac{Y_i - \bar{Y}}{X_i - \bar{X}} \quad \text{and} \quad w_i = \frac{(X_i - \bar{X})^2}{\sum_k (X_k - \bar{X})^2}$$

are case-based slopes and weights. The slope for the i 'th case is based on the line segment that reaches from the center point (\bar{X}, \bar{Y}) to the data point (X_i, Y_i) . The weights grow with the square of the distance of X_i from \bar{X} . This is consistent with the variance of the case-based slopes b_i : $\mathbf{V}[b_i|\mathbf{X}] = \frac{n-1}{n} \sigma^2 / (X_i - \bar{X})^2$, assuming homoskedasticity only, $\mathbf{V}[Y_i|X_i] = \sigma^2$, but not linearity of $\mu(X_i) = \mathbf{E}[Y_i|X_i]$. Plausibly, $w_i \sim 1/\mathbf{V}[b_i] \sim (X_i - \bar{X})^2$, and $\hat{\beta}$ is therefore the LS solution of a weighted mean problem:

$$\hat{\beta} = \operatorname{argmin}_{\tilde{\beta}} \sum_i w_i (b_i - \tilde{\beta})^2,$$

which is the same as the unweighted LS problem

$$\hat{\beta} = \operatorname{argmin}_{\tilde{\beta}} \sum_i ((Y_i - \bar{Y}) - \tilde{\beta}(X_i - \bar{X}))^2.$$

[A well-known variant, used for example in Gelman and Park (2008), represents the slope estimate as follows:

$$\hat{\beta} = \sum_{i,i'} w_{i,i'} b_{i,i'}, \quad \text{where} \quad w_{i,i'} = \frac{(X_i - X_{i'})^2}{\sum_{k,k'} (X_k - X_{k'})^2} \quad \text{and} \quad b_{i,i'} = \frac{Y_i - Y_{i'}}{X_i - X_{i'}}$$

are pairwise weights and slopes.]

We can carry the preceding over to $\mathbf{E}[\hat{\beta}|\mathbf{X}]$, which is the target of LS estimation conditional on the design \mathbf{X} . We make no assumption of linearity about $\mu_i = \mathbf{E}[Y_i|\mathbf{X}]$, which can in fact be entirely arbitrary. Taking the conditional expectation in (55) given \mathbf{X} , we obtain

$$\mathbf{E}[\hat{\beta}|\mathbf{X}] = \frac{\sum_i (\mu_i - \bar{\mu}) \cdot (X_i - \bar{X})}{\sum_k (X_k - \bar{X})^2} = \sum_i w_i \beta_i, \quad (56)$$

where $\beta_i = (\mu_i - \bar{\mu}) / (X_i - \bar{X})$. Thus $\mathbf{E}[\hat{\beta}|\mathbf{X}]$, too, is an average of case-based slopes, for line segments from $(\bar{X}, \bar{\mu})$ to (X_i, μ_i) . [This is illustrated in Figure 3.]

Finally, we can extend the above to the population when the design \mathbf{X} is considered random, no conditioning is applied, and again no assumption of linearity about $\mu(X) = \mathbf{E}[Y|X]$ is made:

$$\begin{aligned} \beta &= \operatorname{argmin}_{\tilde{\beta}} \mathbf{E} \left[\left((Y - \mathbf{E}[Y]) - \tilde{\beta}(X - \mathbf{E}[X]) \right)^2 \right] \\ &= \frac{\mathbf{E} [(\mu(X) - \mathbf{E}[\mu(X)])(X - \mathbf{E}[X])]}{\mathbf{E}[(X - \mathbf{E}[X])^2]} \\ &= \mathbf{E}[w(X) \beta(X)], \end{aligned}$$

where

$$\beta(X) = \frac{\mu(X) - \mathbf{E}[\mu(X)]}{X - \mathbf{E}[X]} \quad \text{and} \quad w(X) = \frac{(X - \mathbf{E}[X])^2}{\mathbf{E}[(X - \mathbf{E}[X])^2]},$$

and, of course: $\mathbf{E}[\mu(X)] = \mathbf{E}[Y]$.

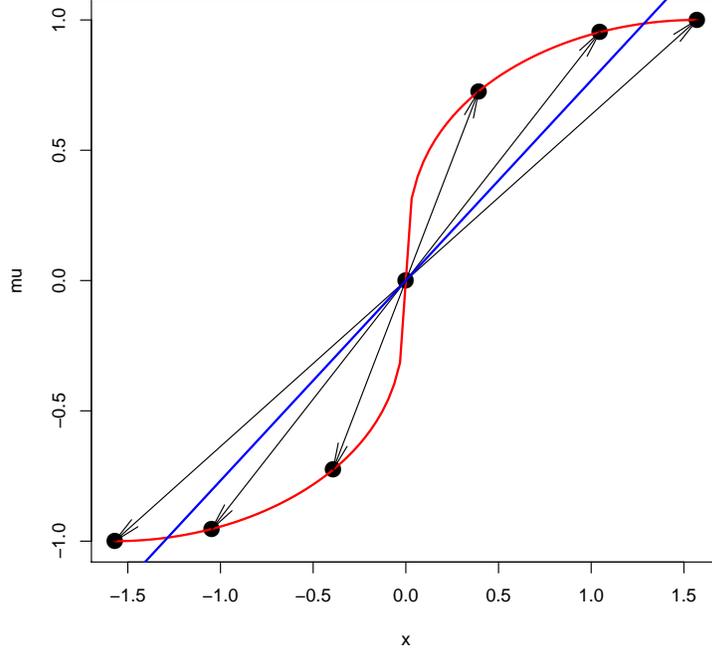


Figure 3: The dots show points (X_i, μ_i) that do not fall on a line; they fall on the red curve. The point $(0,0)$ is the center $(\bar{X}, \bar{\mu})$. The arrows show the segments from the center to the points (X_i, μ_i) ; their slopes are β_i . The average slope β is shown in blue. It is relatively flat because the weights are large for outlying values X_i .

B Proofs

B.1 Conditional Expectation of RSS

The conditional expectation of the RSS allowing for nonlinearity and heteroskedasticity:

$$\mathbf{E}[\|r\|^2 | \mathbf{X}] = \mathbf{E}[\mathbf{Y}^T (\mathbf{I} - \mathbf{H}) \mathbf{Y} | \mathbf{X}] \quad (57)$$

$$= \mathbf{E}[(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta} + \boldsymbol{\epsilon})' (\mathbf{I} - \mathbf{H}) (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta} + \boldsymbol{\epsilon}) | \mathbf{X}] \quad (58)$$

$$= \mathbf{E}[(\boldsymbol{\eta} + \boldsymbol{\epsilon})^T (\mathbf{I} - \mathbf{H}) (\boldsymbol{\eta} + \boldsymbol{\epsilon}) | \mathbf{X}] \quad (59)$$

$$= \text{tr}(\mathbf{E}[(\mathbf{I} - \mathbf{H}) (\boldsymbol{\eta} + \boldsymbol{\epsilon}) (\boldsymbol{\eta} + \boldsymbol{\epsilon})^T | \mathbf{X}]) \quad (60)$$

$$= \text{tr}((\mathbf{I} - \mathbf{H}) (\boldsymbol{\eta}\boldsymbol{\eta}^T + \mathbf{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T | \mathbf{X}])) \quad (61)$$

$$= \text{tr}((\mathbf{I} - \mathbf{H}) (\boldsymbol{\eta}\boldsymbol{\eta}^T + \mathbf{D}_{\sigma^2})) \quad (62)$$

$$= |(\mathbf{I} - \mathbf{H})\boldsymbol{\eta}|^2 + \text{tr}((\mathbf{I} - \mathbf{H})\mathbf{D}_{\sigma^2}) \quad (63)$$

B.2 Limit of Squared Adjusted Predictors

The asymptotic limit of $\|\mathbf{X}_{j\bullet}\|^2$:

$$\begin{aligned}
\frac{1}{N}\|\mathbf{X}_{j\bullet}\|^2 &= \frac{1}{N}\mathbf{X}_j^T(\mathbf{I} - \mathbf{H}_{\check{v}})\mathbf{X}_j \\
&= \frac{1}{N}(\mathbf{X}_j^T\mathbf{X}_j - \mathbf{X}_j^T\mathbf{H}_{\check{v}}\mathbf{X}_j) \\
&= \frac{1}{N}X_{i,j}^2 - \left(\frac{1}{N}\sum X_{i,j}\vec{\mathbf{X}}_{i,\check{v}}^T\right)\left(\sum_i\vec{\mathbf{X}}_{i,\check{v}}\vec{\mathbf{X}}_{i,\check{v}}^T\right)^{-1}\left(\sum_i\vec{\mathbf{X}}_{i,\check{v}}X_{i,j}\right) \\
&\xrightarrow{P} \mathbf{E}[X_j^2] - \mathbf{E}[X_j\vec{\mathbf{X}}_{\check{v}}]\mathbf{E}[\vec{\mathbf{X}}_{\check{v}}\vec{\mathbf{X}}_{\check{v}}^T]^{-1}\mathbf{E}[\vec{\mathbf{X}}_{\check{v}}X_j] \\
&= \mathbf{E}[X_{j\bullet}^2]
\end{aligned}$$

B.3 Conventional SE

We will use the notations $\Sigma_{\vec{\mathbf{X}}} = \mathbf{E}[\vec{\mathbf{X}}\vec{\mathbf{X}}^T]$ and $\Sigma_{\eta(\vec{\mathbf{X}}),\vec{\mathbf{X}}} = \mathbf{E}[\eta(\vec{\mathbf{X}})^2\vec{\mathbf{X}}\vec{\mathbf{X}}^T]$, and we will also need a p -dimensional normal random vector \mathbf{Z} for the following limit in distribution:

$$\frac{1}{N^{1/2}}\sum_{i=1}^N\eta(\vec{\mathbf{X}}_i)\vec{\mathbf{X}}_i^T \xrightarrow{\mathcal{D}} \mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \Sigma_{\eta(\vec{\mathbf{X}}),\vec{\mathbf{X}}}).$$

The following are the ingredients for the limiting behaviors of $\|\boldsymbol{\eta}\|^2$ and $\|\mathbf{H}\boldsymbol{\eta}\|^2$:

$$\begin{aligned}
N(\mathbf{X}^T\mathbf{X})^{-1} &= \left(\frac{1}{N}\sum_{i=1}^N\vec{\mathbf{X}}_i\vec{\mathbf{X}}_i^T\right)^{-1} \\
&\xrightarrow{P} \Sigma_{\vec{\mathbf{X}}}^{-1} \\
\frac{1}{N}\|\boldsymbol{\eta}\|^2 &= \frac{1}{N}\sum_{i=1}^N\eta(\vec{\mathbf{X}}_i)^2 \\
&\xrightarrow{P} \mathbf{E}[\eta(\vec{\mathbf{X}})^2] = \mathbf{V}[\eta(\vec{\mathbf{X}})] \\
\|\mathbf{H}\boldsymbol{\eta}\|^2 &= \boldsymbol{\eta}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\boldsymbol{\eta} \\
&= \left(\frac{1}{N^{1/2}}\sum_{i=1}^N\eta(\vec{\mathbf{X}}_i)\vec{\mathbf{X}}_i^T\right)\left(\frac{1}{N}\sum_{i=1}^N\vec{\mathbf{X}}_i\vec{\mathbf{X}}_i^T\right)^{-1}\left(\frac{1}{N^{1/2}}\sum_{i=1}^N\eta(\vec{\mathbf{X}}_i)\vec{\mathbf{X}}_i\right) \\
&\xrightarrow{\mathcal{D}} \mathbf{Z}^T\Sigma_{\vec{\mathbf{X}}}^{-1}\mathbf{Z} \\
\frac{1}{N}\|\mathbf{H}\boldsymbol{\eta}\|^2 &\xrightarrow{P} 0
\end{aligned}$$

For large N and fixed p , as $N/(N-p-1) \rightarrow 1$ we have

$$N\frac{\frac{1}{N-p-1}\|(\mathbf{I} - \mathbf{H})\boldsymbol{\eta}\|^2}{\|\mathbf{X}_{j\bullet}\|^2} \xrightarrow{P} \frac{\mathbf{E}[\eta(\vec{\mathbf{X}})^2]}{\mathbf{E}[X_{j\bullet}^2]}. \quad (64)$$

B.4 Asymptotic Normality in Terms of Adjustment

We gave the asymptotic limit of the conditional bias in vectorized form after (33)-(35). Here we derive the equivalent element-wise limit using adjustment to show (36)-(38). The variance of the conditional bias is the marginal inflator of SE.

$$N^{1/2}(\mathbf{E}[\hat{\beta}_j|\mathbf{X}] - \beta_j) = N^{1/2} \frac{\langle \mathbf{X}_{j\bullet}, \boldsymbol{\eta} \rangle}{\|\mathbf{X}_{j\bullet}\|^2} = \frac{\frac{1}{N^{1/2}} \mathbf{X}_j^T \boldsymbol{\eta} - \frac{1}{N^{1/2}} \mathbf{X}_j^T \mathbf{H}_{\setminus j} \boldsymbol{\eta}}{\frac{1}{N} \|\mathbf{X}_{j\bullet}\|^2}$$

$$\begin{aligned} \frac{1}{N^{1/2}} \mathbf{X}_j^T \mathbf{H}_{\setminus j} \boldsymbol{\eta} &= \frac{1}{N^{1/2}} \mathbf{X}_j^T \mathbf{X}_{\setminus j} (\mathbf{X}_{\setminus j}^T \mathbf{X}_{\setminus j})^{-1} \mathbf{X}_{\setminus j}^T \boldsymbol{\eta} \\ &= \left(\frac{1}{N} \sum_i X_{i,j} \vec{\mathbf{X}}_{i,\setminus j}^T \right) \left(\frac{1}{N} \sum_i \vec{\mathbf{X}}_{i,\setminus j} \vec{\mathbf{X}}_{i,\setminus j}^T \right)^{-1} \left(\frac{1}{N^{1/2}} \sum_i \vec{\mathbf{X}}_{i,\setminus j} \eta(\vec{\mathbf{X}}_i) \right) \\ &\stackrel{\mathcal{D}}{\approx} \mathbf{E}[X_j \vec{\mathbf{X}}_{\setminus j}] \mathbf{E}[\vec{\mathbf{X}}_{\setminus j} \vec{\mathbf{X}}_{\setminus j}^T] \left(\frac{1}{N^{1/2}} \sum_i \vec{\mathbf{X}}_{i,\setminus j} \eta(\vec{\mathbf{X}}_i) \right) \\ &= \boldsymbol{\beta}_{j\cdot}^T \left(\frac{1}{N^{1/2}} \sum_i \vec{\mathbf{X}}_{i,\setminus j} \eta(\vec{\mathbf{X}}_i) \right) \\ &= \frac{1}{N^{1/2}} \sum_i (\boldsymbol{\beta}_{j\cdot}^T \vec{\mathbf{X}}_{i,\setminus j}) \eta(\vec{\mathbf{X}}_i) \end{aligned}$$

$$\begin{aligned} \frac{1}{N^{1/2}} (\mathbf{X}_j^T \boldsymbol{\eta} - \mathbf{X}_j^T \mathbf{H}_{\setminus j} \boldsymbol{\eta}) &\stackrel{\mathcal{D}}{\approx} \frac{1}{N^{1/2}} \sum_i (X_{i,j} - \boldsymbol{\beta}_{j\cdot}^T \vec{\mathbf{X}}_{i,\setminus j}) \eta(\vec{\mathbf{X}}_i) \\ &\stackrel{\mathcal{D}}{\rightarrow} \mathcal{N} \left(0, \mathbf{V}[(X_j - \boldsymbol{\beta}_{j\cdot}^T \vec{\mathbf{X}}_{\setminus j}) \eta(\vec{\mathbf{X}})] \right) \\ &= \mathcal{N} \left(0, \mathbf{V}[X_{j\bullet} \eta(\vec{\mathbf{X}})] \right) \end{aligned}$$

$$N^{1/2}(\mathbf{E}[\hat{\beta}_j|\mathbf{X}] - \beta_j) \stackrel{\mathcal{D}}{\rightarrow} \mathcal{N} \left(0, \frac{\mathbf{V}[X_{j\bullet} \eta(\vec{\mathbf{X}})]}{\mathbf{E}[X_{j\bullet}^2]^2} \right)$$

B.5 A Family of Nonlinearities with Extreme *RAV*

We refer to the example constructed in Section 6 and recall that both $\mathbf{E}[X^2] = 1$ and $\mathbf{E}[\eta_t(X)^2] = 1$. We abbreviate $\bar{p}(t) = 1 - p(t)$.

$$\begin{aligned}
\mathbf{RAV}(\eta_t) &= \mathbf{E} [\eta_t(X)^2 X^2] \\
&= \frac{1}{p(t)\bar{p}(t)} \mathbf{E} \left[(1_{\{|X|>t\}} - p(t))^2 X^2 \right] \\
&= \frac{1}{p(t)\bar{p}(t)} \mathbf{E} \left[(1_{\{|X|>t\}} - 2 \cdot 1_{\{|X|>t\}} p(t) + p(t)^2) X^2 \right] \\
&= \frac{1}{p(t)\bar{p}(t)} \mathbf{E} \left[(1_{\{|X|>t\}}(1 - 2p(t)) + p(t)^2) X^2 \right] \\
&= \frac{1}{p(t)\bar{p}(t)} (\mathbf{E} [1_{\{|X|>t\}} X^2]) (1 - 2p(t) + p(t)^2) \\
&\geq \frac{1}{p(t)\bar{p}(t)} (p(t) t^2 (1 - 2p(t)) + p(t)^2) \quad \text{for } p(t) \leq \frac{1}{2} \\
&= \frac{1}{\bar{p}(t)} (t^2 (1 - 2p(t)) + p(t)) \\
&\geq t^2 (1 - 2p(t)) + p(t) \\
&\sim t^2 \quad \text{as } t \uparrow \infty.
\end{aligned}$$

For the following we note $1_{\{|X|>t\}} - p(t) = -1_{\{|X|\leq t\}} + \bar{p}(t)$:

$$\begin{aligned}
\mathbf{RAV}(\eta_t) &= \mathbf{E} [\eta_t(X)^2 X^2] \\
&= \frac{1}{p(t)\bar{p}(t)} \mathbf{E} \left[(1_{\{|X|\leq t\}} - \bar{p}(t))^2 X^2 \right] \\
&= \frac{1}{p(t)\bar{p}(t)} \mathbf{E} \left[(1_{\{|X|\leq t\}} - 2 \cdot 1_{\{|X|\leq t\}} \bar{p}(t) + \bar{p}(t)^2) X^2 \right] \\
&= \frac{1}{p(t)\bar{p}(t)} \mathbf{E} \left[(1_{\{|X|\leq t\}}(1 - 2\bar{p}(t)) + \bar{p}(t)^2) X^2 \right] \\
&= \frac{1}{p(t)\bar{p}(t)} (\mathbf{E} [1_{\{|X|\leq t\}} X^2] (1 - 2\bar{p}(t)) + \bar{p}(t)^2) \\
&\leq \frac{1}{p(t)\bar{p}(t)} (\bar{p}(t) t^2 (1 - 2\bar{p}(t)) + \bar{p}(t)^2) \quad \text{for } \bar{p}(t) \leq \frac{1}{2} \\
&= \frac{1}{p(t)} (t^2 (1 - 2\bar{p}(t)) + \bar{p}(t)) \\
&\sim t^2 + \bar{p}(t) \quad \text{as } t \downarrow 0,
\end{aligned}$$

assuming $\bar{p}(0) = P[X = 0] = 0$.

References

- [1] ALLISON, P. D. (1995). The Impact of Random Predictors on Comparisons of Coefficients Between Models: Comment on Clogg, Petkova, and Haritou. *American Journal of Sociology* **100** (5), 1294–1305.
- [2] BELSELEY, D. A., KUH, E. and WELSCH, R. E. (1980). *Regression diagnostics: identifying influential data and sources of collinearity*. Wiley series in probability and mathematical statistics. Hoboken, NJ: John Wiley & Sons, Inc.

- [3] BOX, G. E. P. (1979). Robustness in the Strategy of Scientific Model Building. in *Robustness in Statistics: Proceedings of a Workshop* (Launer, R. L., and Wilkinson, G. N., eds.) xxx–xxx, Place: Publisher.
- [4] BOX, G. E. P. and COX, D. R. (1964). An analysis of transformations. *Journal of the Royal Statistical Society, Series B* **26** (2), 211-252.
- [5] BREIMAN, L. and SPECTOR, P. (1992). Submodel Selection and Evaluation in Regression. The X-Random Case. *International Statistical Review* **50** (3), 291–319.
- [6] BUJA, A., HASTIE, T. and TIBSHIRANI, R. (1989). Linear Smoothers and Additive Models (with discussions and rejoinder). *The Annals of Statistics*, **17** (2), 453–555.
- [7] CHATFIELD, C. (1995). Model Uncertainty, Data Mining and Statistical Inference (with discussion). *Journal of the Royal Statistical Society, Series A* **158** (3), 419-466.
- [8] CLOGG, C. C., PETKOVA, E. and HARITOU, A. (1995). Statistical Methods for Comparing Regression Coefficients Between Models. *American Journal of Sociology* **100** (5), 1261–1293.
- [9] COX, D. R. and HINKLEY, D. V. (1974). *Theoretical Statistics*, London: Chapman & Hall.
- [10] COX, D.R. (1995). Discussion of Chatfield (1995). *Journal of the Royal Statistical Society, Series A* **158** (3), 455-456.
- [11] FREEDMAN, D. A. (1981). Bootstrapping Regression Models. *The Annals of Statistics* **9** (6), 1218–1228.
- [12] HARRISON, X. and RUBINFELD, X. (1978) Hedonic prices and the demand for clean air. *Journal of Environmental Economics and Management* **5**, 81–102.
- [13] HASTIE, T. J. and TIBSHIRANI, R. J. (1990). *Generalized Additive Models*, London: Chapman & Hall/CRC Monographs on Statistics & Applied Probability.
- [14] KAUERMANN, G. and CARROLL, R. J. (2001). *A note on the efficiency of sandwich covariance matrix estimation*, *Journal of the American Statistical Association* **96**(456), 1387-1396.
- [15] MAMMEN, E. (1993). Bootstrap and Wild Bootstrap for High Dimensional Linear Models. *The Annals of Statistics* **21** (1), 255–285.