Statistical Science 2019, Vol. 34, No. 4, 545–565 https://doi.org/10.1214/18-STS694 @ Institute of Mathematical Statistics 2019

# Models as Approximations II: A Model-Free Theory of Parametric Regression<sup>1</sup>

Andreas Buja, Lawrence Brown, Arun Kumar Kuchibhotla, Richard Berk, Edward George and Linda Zhao

Abstract. We develop a model-free theory of general types of parametric regression for i.i.d. observations. The theory replaces the parameters of parametric models with statistical functionals, to be called "regression functionals," defined on large nonparametric classes of joint x-y distributions, without assuming a correct model. Parametric models are reduced to heuristics to suggest plausible objective functions. An example of a regression functional is the vector of slopes of linear equations fitted by OLS to largely arbitrary x-y distributions, without assuming a linear model (see Part I). More generally, regression functionals can be defined by minimizing objective functions, solving estimating equations, or with ad hoc constructions. In this framework, it is possible to achieve the following: (1) define a notion of "wellspecification" for regression functionals that replaces the notion of correct specification of models, (2) propose a well-specification diagnostic for regression functionals based on reweighting distributions and data, (3) decompose sampling variability of regression functionals into two sources, one due to the conditional response distribution and another due to the regressor distribution interacting with misspecification, both of order  $N^{-1/2}$ , (4) exhibit plug-in/sandwich estimators of standard error as limit cases of x-y bootstrap estimators, and (5) provide theoretical heuristics to indicate that x-y bootstrap standard errors may generally be preferred over sandwich estimators.

*Key words and phrases:* Ancillarity of regressors, misspecification, econometrics, sandwich estimator, bootstrap, bagging.

Andreas Buja is the Liem Sioe Liong/First Pacific Company Professor of Statistics, Lawrence Brown was the Miers Busch Professor of Statistics, Arun Kumar Kuchibhotla is Doctoral Student of Statistics, Richard Berk is Professor of Criminology and Statistics, Edward George is the Universal Furniture Professor of Statistics, Linda Zhao is Professor of Statistics, Statistics Department, The Wharton School, University of Pennsylvania, 400 Jon M. Huntsman Hall, 3730 Walnut Street, Philadelphia, Pennsylvania 19104-6340, USA (e-mail: lzhao@wharton.upenn.edu).

<sup>1</sup>Discussed in 10.1214/19-STS722, 10.1214/19-STS723, 10.1214/19-STS724, 10.1214/19-STS725, 10.1214/19-STS726, 10.1214/19-STS746, 10.1214/19-STS747, 10.1214/19-STS748, 10.1214/19-STS756; rejoinder at 10.1214/19-STS762.

"The hallmark of good science is that it uses models and 'theory' but never believes them." (J. W. Tukey, 1962, citing Martin Wilk)

#### 1. INTRODUCTION

We develop in this second article a model-free theory of parametric regression, assuming for simplicity i.i.d. x-y observations with quite arbitrary joint distributions. The starting point is the realization that regression models are approximations and should not be thought of as generative truths. A general recognition of this fact may be implied by the commonly used term "working model," but this vague term does not resolve substantive issues, created here by the fact that models are approximations and not truths. The primary issue

is that traditional model parameters define meaningful quantities only under conditions of model correctness. If the idea of models as approximations is taken seriously, one has to extend the notion of parameter from model distributions to basically arbitrary distributions. This is achieved by what is often called "projection onto the model," that is, finding for the actual data distribution the best approximating distribution within the model; one defines that distribution's parameter settings to be the target of estimation. Through such "projection" the parameters of a working model are extended to "statistical functionals," that is, mappings of largely arbitrary data distributions to numeric quantities. We have thus arrived at a functional point of view of regression, a view based on what we call regression functionals.

The move from traditional regression parameters in correctly specified models to regression functionals obtained from best approximations may raise fears of opening the gates to irresponsible data analysis where misspecification is of no concern. No such thing is intended here. Instead, we rethink the essence of regression and develop a new notion of well-specification of regression functionals, to replace the notion of correct specification of regression models. In the following bullets, we outline an argument in the form of simple postulates:

- The essence of regression is the asymmetric analysis of association: Variables with a joint distribution P are divided into response and regressors.
- Motivated by prediction and causation problems, interest focuses on properties of the conditional distribution of the response given the regressors.
- The goal or, rather, the hope is that the chosen quantities/functionals of interest are properties of the observed conditional response distribution, irrespective of the regressor distribution.
- Consequently, a regression functional will be called *well-specified* if it is a property of the observed conditional response distribution at hand, *irrespective of the regressor distribution*.

The first bullet is uncontroversial: asymmetric analysis is often natural, as in the contexts of prediction and causation. The second bullet remains at an intended level of vagueness as it explains the nature of the asymmetry, namely, the focus on the regressor-conditional response distribution. Intentionally there is no mention of regression models. The third bullet also steers clear of regression models by addressing instead quantities of interest, that is, regression functionals. In this and

the last bullet, the operational requirement is that the quantities of interest not depend on the regressor distribution. It is this constancy across regressor distributions that turns the quantities of interest into properties of the conditional response distribution alone.

All this can be made concrete with reference to the groundwork laid in Part I (Section 4) of this twopart series of articles. Consider the regression functional consisting of the coefficient vector obtained from OLS linear regression. It was shown in Part I that this vector does not depend on the regressor distribution (is "well-specified") if and only if the conditional response mean is a linear function of the regressors. This is the situation in which the coefficient vector fully describes the conditional mean function, but no other aspect of the conditional response distribution. Wellspecification of the OLS coefficient functional is therefore a weaker condition than correct specification of the linear model by setting aside homoskedasticity and Gaussianity which are linear model requirements not intimately tied to the slopes.

A desirable feature of the proposed definition of well-specification is that it generalizes to arbitrary types of parametric regression or, more precisely, to the statistical functionals derived from them. In particular, it applies to GLMs where the meaning of well-specified coefficients is again correct specification of the mean function but setting aside other model requirements. Well-specification further applies to regression functionals derived from optimizing general objective functions or solving estimating equations. Well-specification finally applies to any ad hoc quantities if they define regression functionals for joint *x-y* distributions.

The proposed notion of well-specification of regression functionals does not just define an ideal condition for populations but also lends itself to a tangible methodology for real data. A diagnostic for well-specification can be based on perturbation of the regressor distribution without affecting the conditional response distribution. Such perturbations can be constructed by reweighting the joint x-y distribution with weight functions that only depend on the regressors. If a regression functional is not constant under such reweighting, it is misspecified.

In practice, use of this diagnostic often works out as follows. Some form of misspecification will be detected for some of the quantities of interest, but the diagnostic will also aid in interpreting the specifics of the misspecification. The reason is that reweighting essentially localizes the regression functionals. For the

coefficients of OLS linear regression, for example, this means that reweighting reveals how the coefficients of the best fitting linear equation vary as the weight function moves across regressor space. Put this way, the diagnostic is related to nonparametric regression, but its advantage is that it focuses on the quantities of interest at all times, while switching from parametric to nonparametric regression requires a rethinking of the meaning of the original quantities in terms of the nonparametric fit. To guide users of the diagnostic to insightful choices of weight functions, we introduce a set of specific reweighting methodologies, complete with basic statistical inference.

Following these methodological proposals, we return to the inferential issues raised in Part I and treat them in generality for all types of well-behaved regression functionals. We show that sampling variation of regression functionals has two sources, one due to the conditional response distribution, the other due to the regressor distribution interacting with misspecification, where "misspecification" is meant in the sense of "violated well-specification" of the regression functional. A central limit theorem (CLT) shows that *both* sources, as a function of the sample size N, are of the usual order  $N^{-1/2}$ . Finally, it is shown that asymptotic plugin/sandwich estimators of standard error are limits of x-y bootstrap estimators, revealing the former to be an extreme case of the latter.

The present analysis becomes necessarily more opaque because algebra that worked out explicitly and lucidly for linear OLS in Part I is available in the general case only in the form of asymptotic approximation based on influence functions. Still, the analysis is now informed by the notion of well-specification of regression functionals, which gives the results a rather satisfactory form.

The article continues as follows. In Section 2, we discuss typical ways of defining regression functionals, including optimization of objective functions and estimating equations. In Section 3, we give the precise definition of well-specification and illustrate it with various examples. In Section 4, we introduce the reweighting diagnostic for well-specification, illustrated in Section 5 with specific reweighting methodologies applied to the LA homeless data (Part I). Section 6 shows for plug-in estimators of regression functionals how the sampling variability is canonically decomposed into contributions from the conditional response noise and from the randomness of the regressors. In Section 7, we state general CLTs analogous to the OLS versions

of Part I. In Section 8, we analyze model-free estimators of standard error derived from the M-of-N pairs bootstrap and asymptotic variance plug-in (often of the sandwich form). It holds in great generality that plug-in is the limiting case of bootstrap when  $M \to \infty$ . We also give some heuristics to suggest that bootstrap estimators might generally be preferred over plug-in/sandwich estimators. In Section 9, we summarize the path taken in these two articles.

REMARK. For notes on the history of model robustness, see Part I, Section 1. For the distinction between model robustness and outlier/heavy-tail robustness, see Part I, Section 13.

### 2. TARGETS OF ESTIMATION: REGRESSION FUNCTIONALS

This section describes some of the ways of constructing regression functionals, including those based on "working models" used as heuristics to suggest plausible objective functions. We use the following notation and assumptions throughout: At the population level, there are two random variables, the regressor  $\vec{X}$  with values in a measurable space  $\mathcal{X}$  and the response Y with values in a measurable space  $\mathcal{Y}$ , with a joint distribution  $P_{Y,\vec{X}}$ , a conditional response distribution  $P_{Y|\vec{X}}$  and a marginal regressor distribution  $P_{\vec{X}}$ . We express the connection between them using " $\otimes$ " notation:

$$(1) P_{Y,\vec{X}} = P_{Y|\vec{X}} \otimes P_{\vec{X}}.$$

Informally, this is expressed in terms of densities by  $p(y, \vec{x}) = p(y|\vec{x}) p(\vec{x})$ . In contrast to Part I, the regressor and response spaces  $\mathcal{X}$  and  $\mathcal{Y}$  are now entirely arbitrary. The typographic distinction between  $\vec{X}$  and Y is a hold-over from the OLS context of Part I. Both spaces,  $\mathcal{X}$  and  $\mathcal{Y}$ , can be of any measurement type, univariate or multivariate, or even spaces of signals or images.

Regression functionals need to be defined on universes of joint distributions that are sufficiently rich to grant the manipulations that follow, including the assumed existence of moments, influence functions and closedness for certain mixtures. The details are tedious, hence deferred to Appendix A (Buja et al., 2019) without claim to technical completeness. The treatment is largely informal so as not to get bogged down in distracting detail. Also, the asymptotics will be traditional in the sense that  $\mathcal X$  and  $\mathcal Y$  are fixed and  $N \to \infty$ . For more modern technical work on related matters; see Kuchibhotla et al. (2018). Readers comfortable with defining functionals by way of objective functions and estimating equations may want to continue with Section 2.3.

### 2.1 Regression Functionals from Optimization: ML and PS Functionals

In Part I, we described the interpretation of linear OLS coefficients as regression functionals. The expression "linear OLS" is used on purpose to avoid the expression "linear models" because no model is assumed. Fitting a linear equation using OLS is a procedure to achieve a best fit of an equation by the OLS criterion. This approach can be generalized to other objective functions  $\mathcal{L}(\theta; y, \vec{x})$ :

(2) 
$$\theta(P) = \underset{\theta \in \Theta}{\operatorname{argmin}} E_P [\mathcal{L}(\theta; Y, \vec{X})].$$

A common choice for  $\mathcal{L}(\theta; y, \vec{x})$  is the negative log-likelihood of a parametric regression model for  $Y|\vec{X}$ , defined by a parametrized family of conditional response distributions  $\{Q_{Y|\vec{X};\theta}:\theta\in\Theta\}$  with conditional densities  $\{q(y|\vec{x};\theta):\theta\in\Theta\}$ . The model is not assumed to be correctly specified, and its only purpose is to serve as a heuristic to suggest an objective function:

(3) 
$$\mathcal{L}(\boldsymbol{\theta}; y, \vec{\boldsymbol{x}}) = -\log q(y|\vec{\boldsymbol{x}}; \boldsymbol{\theta}).$$

In this case, the regression functional resulting from (2) will be called a maximum-likelihood functional or ML functional for short. It minimizes the Kullback-Leibler (K-L) divergence of  $P_{Y,\vec{X}} = P_{Y|\vec{X}} \otimes P_{\vec{X}}$  and  $Q_{Y|\vec{X}:\theta} \otimes P_{\vec{X}}$ , which is why one loosely interprets a ML functional as arising from a "projection of the actual data distribution onto the parametric model." ML functionals can be derived from major classes of regression models, including GLMs. Technically, they also comprise many M-estimators based on Huber  $\rho$ functions (Huber, 1964), including least absolute deviation (LAD,  $L_1$ ) as an objective function for conditional medians, and tilted  $L_1$  versions for arbitrary conditional quantiles, all of which can be interpreted as negative log-likelihoods of certain distributions, even if these may not usually be viable models for actual data. Not in the class of negative log-likelihoods are objective functions for M-estimators with redescending influence functions such as Tukey's biweight estimator (which also poses complications due to nonconvexity).

Natural extensions of ML functionals can be based on so-called "proper scoring rules" (Appendix B) which arise as cross-entropy terms of Bregman divergences. A special case is the expected negative log-likelihood arising as the cross-entropy term of K–L divergence. The optimization criterion is the proper scoring rule applied to the conditional response distribution  $P_{Y|\bar{X}}$  and model distributions  $Q_{Y|\bar{X};\theta}$ , averaged

over regressor space with  $P_{\vec{X}}$ . The resulting regression functionals may be called "proper scoring functionals" or simply *PS functionals*, a superset of ML functionals. All PS functionals, including ML functionals, have the important property of Fisher consistency: If the model is correctly specified, that is, if  $\exists \theta_0$  such that  $P_{Y|\vec{X}} = Q_{Y|\vec{X};\theta_0}$ , then the population minimizer is  $\theta_0$ :

(4) if 
$$P_{Y,\vec{X}} = Q_{Y|\vec{X}:\theta_0} \otimes P_{\vec{X}}$$
, then  $\theta(P) = \theta_0$ .

See Appendix B for background on proper scoring rules, Bregman divergences and some of their robustness properties to outliers and heavy tailed distributions.

Further objective functions are obtained by adding parameter penalties to existing objective functions:

(5) 
$$\tilde{\mathcal{L}}(\boldsymbol{\theta}; y, \vec{x}) = \mathcal{L}(\boldsymbol{\theta}; y, \vec{x}) + \lambda \mathcal{R}(\boldsymbol{\theta}).$$

Special cases are ridge and lasso penalties. Note that (5) results in one-parameter families of penalized functionals  $\theta_{\lambda}(P)$  defined for populations as well, whereas in practice  $\lambda = \lambda_N$  applies to finite N with  $\lambda_N \to 0$  as  $N \to \infty$ .

# 2.2 Regression Functionals from Estimating Equations: EE Functionals

Objective functions are often minimized by solving stationarity conditions that amount to estimating equations with the scores  $\psi(\theta; y, \vec{x}) = -\nabla_{\theta} \mathcal{L}(\theta; y, \vec{x})$ :

(6) 
$$E_{P}[\psi(\theta; Y, \vec{X})] = 0.$$

One may generalize and define regression functionals as solutions in cases where  $\psi(\theta; y, \vec{x})$  is not the gradient of an objective function; in particular it need not be the score function of a negative log-likelihood. Functionals in this class will be called *EE functionals*. For OLS, the estimating equations are the normal equations, as the score function for the slopes is

(7) 
$$\psi_{\text{OLS}}(\boldsymbol{\beta}; y, \vec{x}) = \vec{x}y - \vec{x}\vec{x}'\boldsymbol{\beta} = \vec{x}(y - \vec{x}'\boldsymbol{\beta}).$$

A seminal work that inaugurated asymptotic theory for general estimating equations is by Huber (1967). A more modern and rigorous treatment is in Rieder (1994).

An extension is the "Generalized Method of Moments" (GMM; Hansen, 1982). It applies when the number of moment conditions (the dimension of  $\psi$ ) is larger than the dimension of  $\theta$ . An important application is to causal inference based on numerous instrumental variables.

Another extension is based on "Generalized Estimating Equations" (GEE; Liang and Zeger, 1986). It applies to clustered data that have intracluster dependence, allowing misspecification of variances and intracluster dependences.

## 2.3 The Point of View of Regression Functionals and Its Implications

Theories of parametric models deal with the issue that a traditional model parameter has many possible estimators, as in the normal model  $\mathcal{N}(\mu, \sigma^2)$  where the sample mean is in various ways the optimal estimate of  $\mu$  whereas the median is a less efficient estimate of the same  $\mu$ . The comparison of estimates of the same traditional parameter has been proposed as a basis of misspecification tests (Hausman, 1978) and called "test for parameter estimator inconsistency" (White, 1982). In a framework based on regression functionals, the situation presents itself differently. Empirical means and medians, for example, are not estimators of the same parameter; instead, they represent different statistical functionals. Similarly, slopes obtained by linear OLS and linear LAD are different regression functionals. Comparing them by forming differences creates new regression functionals that may be useful as diagnostic quantities, but in a model-robust framework there is no concept of "parameter inconsistency" (White, 1982, p. 15), only a concept of differences between regression functionals.

A further point is that in a model-robust theory of observational (as opposed to causal) association, there is no concept of "omitted variables bias." There are only regressions with more or fewer regressor variables, none of which being "true" but some being more useful or insightful than others. Slopes in a larger regression are distinct from the slopes in a smaller regression. It is a source of conceptual confusion to write the slope of the jth regressor as  $\beta_i$ , irrespective of what the other regressors are. In more careful notation, one indexes slopes with the set of selected regressors M as well,  $\beta_{j\cdot M}$ , as is done of necessity in work on post-selection inference (e.g., Berk et al., 2013). Thus the linear slopes  $\beta_{j\cdot M}$  and  $\beta_{j\cdot M'}$  for the jth regressor, when it is contained in both of two regressor sets  $M \neq M'$ , should be considered as distinct regression functionals. The difference  $\beta_{j\cdot M'}-\beta_{j\cdot M}$  is not a bias but a difference between two regression functionals. If it is zero, it indicates that the difference in adjustment between M and M' is immaterial for the jth regressor. If  $\beta_{i.M'}$  and  $\beta_{i.M}$  are very different with opposite

signs, there exists a case of Simpson's paradox for this regressor.

Regression functionals generally depend on the full joint distribution  $P_{Y,\bar{X}}$  of the response *and* the regressors. Conventional regression parameters describe the conditional response distribution only under correct specification,  $P_{Y|\bar{X}} = Q_{Y|\bar{X};\theta}$ , while the regressor distribution  $P_{\bar{X}}$  is sidelined as ancillary. That the ancillarity argument for the regressors is not valid under misspecification was documented in Part I, Section 4. In the following sections, this fact will be the basis of the notion of well-specification of regression functionals.

Finally, we state the following to avoid misunderstandings: In the present work, the objective is not to recommend particular regression functionals, but to point out the freedoms we have in choosing them and the clarifications we need when using them.

### 3. MIS/WELL-SPECIFICATION OF REGRESSION FUNCTIONALS

Section 1 motivated a notion of well-specification for regression functionals, and this section provides the technical notations. The heuristic idea is that a regression functional is well-specified for a joint distribution of the regressors and the response if it does not depend on the marginal regressor distribution. In concrete terms, this means that the functional does not depend on where the regressors happen to fall. The functional is therefore a property of the conditional response distribution alone.

### 3.1 Definition of Well-Specification for Regression Functionals

Recall the notation introduced in (1):  $P_{Y,\vec{X}} = P_{Y|\vec{X}} \otimes P_{\vec{X}}$ . Here, a technical detail requires clarification: conditional distributions are defined only almost surely with regard to  $P_{\vec{X}}$ , but we will assume that  $\vec{x} \mapsto P_{Y|\vec{X}=\vec{x}}$  is a Markov kernel defined for all  $\vec{x} \in \mathcal{X}$ . With these conventions,  $P_{Y|\vec{X}}$  and  $P_{\vec{X}}$  uniquely determine  $P_{Y,\vec{X}} = P_{Y|\vec{X}} \otimes P_{\vec{X}}$  by (1), but not quite vice versa. Thus  $\theta(\cdot)$  can be written as

$$\theta(P) = \theta(P_{Y|\vec{X}} \otimes P_{\vec{X}}).$$

DEFINITION. The regression functional  $\boldsymbol{\theta}(\cdot)$  is well-specified for  $\boldsymbol{P}_{Y|\vec{X}}$  if

$$\theta(P_{Y|\vec{X}}\otimes P_{\vec{X}}) = \theta(P_{Y|\vec{X}}\otimes P_{\vec{X}}')$$

for all acceptable regressor distributions  $m{P}_{ec{X}}$  and  $m{P}_{ec{X}}'$ 

<sup>&</sup>lt;sup>1</sup>Thus we assume a "regular version" has been chosen, as is always possible on Polish spaces.

The term "acceptable" accounts for exclusions of regressor distributions such as those due to nonidentifiability when fitting equations, in particular, perfect collinearity when fitting linear equations (see Appendix A).

Importantly, the notion of well-specification is a *joint* property of a specific  $\theta(\cdot)$  and a specific  $P_{Y|\vec{X}}$ . A regression functional will be well-specified for some conditional response distributions but not for others.

The notion of well-specification represents an idealization, not a reality. Well-specification is never a fact, only degrees of misspecification are. Yet, idealizations are useful because they give precision and focus to an idea. Here, the idea is that a regression functional is intended to be a property of the conditional response distribution  $P_{Y|\bar{X}}$  alone, regardless of the regressor distribution  $P_{\bar{Y}}$ .

### 3.2 Well-Specification—Some Exercises and Special Cases

Before stating general propositions, here are some special cases to train intuitions.

1. The OLS slope functional can be written

$$\boldsymbol{\beta}(\boldsymbol{P}) = \boldsymbol{E}_{\boldsymbol{P}}[\vec{X}\vec{X}']^{-1}\boldsymbol{E}_{\boldsymbol{P}}[\vec{X}\mu(\vec{X})],$$

where  $\mu(\vec{x}) = E_P[Y|\vec{X} = \vec{x}]$  (Part I, Section 3.2). Thus  $\beta(P)$  depends on  $P_{Y|\vec{X}}$  only through the conditional mean function. The functional is well-specified if  $\mu(\vec{x}) = \beta_0'\vec{x}$  is linear, in which case  $\beta(P) = \beta_0$ . For the reverse, see Part I, Proposition 4.1.

2. A special case is regression through the origin, which we generalize slightly as follows. Let  $h(\vec{x})$  and g(y) be two nonvanishing real-valued square-integrable functions of the regressors and the response, respectively. Define

$$\theta_{h,g}(\mathbf{P}) = \frac{\mathbf{E}_{\mathbf{P}}[g(Y)h(\vec{X})]}{\mathbf{E}_{\mathbf{P}}[h(\vec{X})^2]}.$$

Then  $\theta_{h,g}(P)$  is well-defined for  $P_{Y|\vec{X}}$  if  $E_P[g(Y)|\vec{X}] = \beta \cdot h(\vec{X})$  for some  $\beta$ , hence  $\theta_{h,g}(P) = \beta$ .

3. An ad hoc estimate of a simple linear regression slope is

$$\boldsymbol{\theta}(\boldsymbol{P}) = \boldsymbol{E}_{\boldsymbol{P}}[(Y' - Y'')/(X' - X'') \mid |X' - X''| > \delta],$$

where (Y', X'),  $(Y'', X'') \sim P$  i.i.d. and  $\delta > 0$ . It is inspired by Part I, Section 10, and Gelman and Park (2009). It is well-specified if  $E_P[Y|X] = \beta_0 + \beta_1 X$ , in which case  $\theta(P) = \beta_1$ .

- 4. Ridge regression also defines a slope functional. Let  $\Omega$  be a symmetric nonnegative definite matrix and  $\boldsymbol{\beta}'\boldsymbol{\Omega}\boldsymbol{\beta}$  its quadratic penalty. Solving the penalized LS problem yields  $\boldsymbol{\beta}(\boldsymbol{P}) = (\boldsymbol{E}_{\boldsymbol{P}}[\vec{X}\vec{X}'] + \boldsymbol{\Omega})^{-1}\boldsymbol{E}_{\boldsymbol{P}}[\vec{X}\mu(\vec{X})]$ . This functional is well-specified if the conditional mean is linear,  $\mu(\vec{x}) = \boldsymbol{\beta}_0'\vec{x}$  for some  $\boldsymbol{\beta}_0$ , and  $\boldsymbol{\Omega} = c\boldsymbol{E}_{\boldsymbol{P}}[\vec{X}\vec{X}']$  for some  $c \geq 0$ , in which case  $\boldsymbol{\beta}(\boldsymbol{P}) = 1/(1+c)\boldsymbol{\beta}_0$ , causing uniform shrinkage across all regression coefficients.
- 5. Given a univariate response Y, what does it mean for the functional  $\theta(P) = E_P[Y]$  to be well-specified for  $P_{Y|\vec{X}}$ ? It looks as if it did not depend on the regressor distribution and is therefore always well-specified. This is a fallacy, of course. Because  $E_P[Y] = E_P[\mu(\vec{X})]$ , it follows that  $E_P[Y]$  is independent of  $P_{\vec{X}}$  iff the conditional response mean is constant:  $\mu(\vec{X}) = E_P[Y]$ .
- 6. Homoskedasticity: The average conditional variance functional  $\sigma^2(P) = E_P[V_P[Y|\vec{X}]]$  is well-specified iff  $V_P[Y|\vec{X} = \vec{x}] = \sigma_0^2$  is constant, in which case  $\sigma^2(P) = \sigma_0^2$ . A difficulty is that access to this functional assumes knowledge of a correctly specified mean function  $\mu(\vec{X}) = E_P[Y|\vec{X}]$ .
- 7. The average conditional MSE functional wrt linear OLS is  $E[(Y \beta(P)'\vec{X})^2] = E[m^2(\vec{X})]$  using the notation of Part I, end of Section 3.3. If it is well-specified, that is, if  $m^2(\vec{X}) = m_o^2$  is constant, then linear model-based inference is asymptotically justified (Part I, Lemma 11.4 (a)).
- 8. The correlation coefficient  $\rho(Y, X)$ , if interpreted as a regression functional in a regression of Y on X, is well-specified only in the trivial case when  $\mu(X)$  is constant and  $V_P[Y] > 0$ , hence  $\rho(Y, X) = 0$ .
- 9. Fitting a linear equation by minimizing least absolute deviations (LAD, the  $L_1$  objective function) defines a regression functional that is well-specified if there exists  $\beta_0$  such that median  $[P_{Y|\vec{X}}] = \beta_0'\vec{X}$ .
- 10. In a GLM regression with a univariate response and canonical link, the slope functional is given by

$$\beta(P) = \underset{\beta}{\operatorname{argmin}} E_{P} [b(\vec{X}'\beta) - Y\vec{X}'\beta],$$

where  $b(\theta)$  is a strictly convex function on the real line and  $\theta = \vec{x}' \beta$  is the "canonical parameter" modeled by a linear function of the regressors. The stationary/estimating equations are<sup>2</sup>

$$E_{P}[Y\vec{X}] = E_{P}[\partial b(\vec{X}'\beta)\vec{X}].$$

<sup>&</sup>lt;sup>2</sup>To avoid confusion with matrix transposition, we write  $\partial b$  instead of b' for derivatives.

This functional is well-specified iff  $E_P[Y|\vec{X}] = \partial b(\vec{X}'\beta)$  for  $\beta = \beta(P)$ . Well-specification of  $\beta(P)$  has generally no implication for  $V_P[Y|\vec{X}]$ , except in the next example.

- 11. Linear logistic regression functionals are a special case of GLM functionals where  $Y \in \{0, 1\}$  and  $b(\theta) = \log(1 + \exp(\theta))$ . Well-specification holds iff  $P[Y = 1|\vec{X}] = \varphi(\vec{X}'\beta)$  for  $\beta = \beta(P)$  and  $\varphi(t) = 1/(1+e^{-t})$ . Because the conditional response distribution is Bernoulli, the conditional mean of Y determines the conditional response distribution uniquely, hence well-specification of the regression functional  $\beta(P)$  is the same as correct specification of the logistic regression model.
- 12. If  $\theta(P)$  is well-specified for  $P_{Y|\vec{X}}$ , then so is the functional  $f(\theta(P))$  for any function  $f(\cdot)$ . An example in linear regression is the predicted value  $\beta(P)'\vec{x}$  at the regressor location  $\vec{x}$ . Other examples are contrasts such as  $\beta_1(P) \beta_2(P)$  where  $\beta_j(P)$  denotes the jth coordinate of  $\beta(P)$ .
- 13. A meaningless case of "misspecified functionals" arises when they do not depend on the conditional response distribution at all:  $\theta(P_{Y|\vec{X}} \otimes P_{\vec{X}}) = \theta(P_{\vec{X}})$ . Examples would be tabulations and summaries of individual regressor variables. They could not be well-specified for  $P_{Y|\vec{X}}$  unless they are constants.

### 3.3 Well-Specification of ML, PS and EE Functionals

The following proposition states that a regression functional defined by minimization of an objective function is well-specified if the value of the functional minimizes the objective function at all locations in regressor space.

PROPOSITION 3.1. If  $\theta_0$  minimizes

$$E_P[\mathcal{L}(Y|\vec{X};\theta)|\vec{X}=\vec{x}] \quad \forall \vec{x} \in \mathcal{X},$$

then the minimizer  $\theta(P)$  of

$$E_{P}[\mathcal{L}(Y|\vec{X};\theta)]$$

is well-specified for  $P_{Y|\vec{X}}$ , and  $\theta(P_{Y|\vec{X}} \otimes P_{\vec{X}}) = \theta_0$  for all acceptable regressor distributions  $P_{\vec{X}}$ .

The following is a corollary of Proposition 3.1 and follows from the fact that Bregman (and hence K-L) divergences are minimized when their two arguments are identical.

PROPOSITION 3.2. If  $\theta(\cdot)$  is a ML or PS functional for the working model  $\{Q_{Y|\vec{X};\theta}:\theta\in\Theta\}$ , it is well-specified for all model distributions  $P_{Y|\vec{X}}=Q_{Y|\vec{X};\theta}$ .

The next fact states that an EE functional is well-specified for a conditional response distribution if it satisfies the EE conditionally and globally across regressor space for one value  $\theta_0$ .

PROPOSITION 3.3. If  $\theta_0$  solves  $E_P[\psi(\theta_0; Y, \vec{X}) | \vec{X} = \vec{x}] = \mathbf{0}$  for all  $\vec{x} \in \mathcal{X}$ , then the EE functional defined by  $E_P[\psi(\theta; Y, \vec{X})] = \mathbf{0}$  is well-specified for  $P_{Y|\vec{X}}$ , and  $\theta(P_{Y|\vec{X}} \otimes P_{\vec{X}}) = \theta_0$  for all acceptable regressor distributions  $P_{\vec{X}}$ .

The proof is in Appendix D.

### 3.4 Well-Specification and Causality

The notion of well-specification for regression functionals relates to aspects of causal inference based on direct acyclic graphs (DAGs) and the Markovian structures they represent (e.g., Pearl, 2009). Given a DAG, the theory explains which choices of regressors X permit correct descriptions of causal effects for a given outcome variable Y. Focusing on one such choice of  $\hat{X}$ and Y, one is left with the task of describing interesting quantitative aspects of the conditional distribution  $P_{Y|\vec{X}}$ , which is thought to be unchanging under different manipulations and/or sampling schemes of the regressors X. Therefore, if a quantity of interest is to describe causal effects properly, it should do so irrespective of where the values of the causal variables  $\hat{X}$ have fallen. This is exactly the requirement of wellspecification for regression functionals. In summary, proper causal effects must arise as quantities of interest that are well-specified in the sense of Section 3.1.

Peters, Bühlmann and Meinshausen (2016, Section 1.1) discussed a related notion of "invariance" which can be interpreted as "invariance to regressor distributions." They propose this notion as a heuristic for causal discovery and inference based on multiple data sources with the same variables, one variable being singled out as the response Y. These multiple data sources are leveraged as follows: If for a subset of variables, X, the association  $X \to Y$  is causal, then the conditional distribution  $P_{Y|\vec{X}}$  will be the same across data sources. Subsets of causal variables  $\vec{X}$  with shared  $P_{Y|\vec{X}}$  across sources may therefore be discoverable if the sources differ in their regressor distributions and/or interventions on causal variables. For concreteness, the authors focus on a linear structural equation model (SEM), which allows us to reinterpret their proposals by abandoning the SEM assumption and to consider instead the regression functional consisting of the OLS regression coefficients resulting from the linear SEM. Thus the proposed method is at heart an approach to

detecting and inferring well-specified quantities, cast in a causal framework.

In the following section, we will introduce a diagnostic for well-specification that can be interpreted as emulating multiple data sources from a single data source. The proposal is to systematically reweight the data to synthetically create alternative datasets. Peters et al. (2016, Section 3.3) briefly mention the idea of conditioning as related to the idea of multiple data sources. Such conditioning is naturally achieved by locally reweighting the data, as will be shown next.

# 4. A REWEIGHTING DIAGNOSTIC FOR WELL-SPECIFICATION: TARGETS AND INFERENTIAL TOOLS

Well-specification of regression functionals connects naturally to reweighting, both of populations and of data. A concrete illustration of the basic idea can be given by again drawing on the example of linear OLS: The OLS slope functional is well-specified iff  $E_P[Y|\vec{X}] = \beta_0'\vec{X}$  for some  $\beta_0$ , in which case for any nonnegative weight function  $w(\vec{x})$  we have  $\beta_0 = \operatorname{argmin}_{\beta} E_P[w(\vec{X})(Y - \beta'\vec{X})^2]$ . Therefore, the reweighting of interest is with regard to weights that are functions of the regressors only. The general reason is that such weights affect the distribution of the regressors but not the conditional response distribution. Reweighting provides an intuitive basis for diagnosing well-specification of regression functionals. Because of the practical importance of the proposed reweighting diagnostic, we insert this material early, deferring estimation and inference to Section 6.

Reweighting has an extensive history in statistics, too rich to recount. The present purpose of reweighting is methodological: to diagnose the degree to which the null hypothesis of well-specification of a regression functional is violated. As an aid, we propose what we call a "tilt test." It provides evidence of whether a real-valued regression functional is likely to rise or fall (tilt up or down) from one extreme of reweighting to another. The conclusions from a rejection based on this test are simple and interpretable.

In practice, the majority of regression functionals of interest are regression slopes connected to specific regressors. A more interesting problem than detection of misspecification is another question: Does misspecification impinge on the statistical significance of a slope of interest? That is, would a slope have lost or gained statistical significance if the regressor distribution had been different? This is the primary question to be addressed by the reweighting diagnostic.

#### 4.1 Reweighting and Well-Specification

Consider reweighted versions of the joint distribution  $P = P_{Y,\vec{X}}$  with weight functions  $w(\vec{x})$  that depend only on the regressors, not the response, written as

$$\mathbf{P}_{Y,\vec{X}}^{w}(dy,d\vec{x}) = w(\vec{x})\mathbf{P}_{Y,\vec{X}}(dy,d\vec{x}) \quad \text{or}$$
$$p^{w}(y,\vec{x}) = w(\vec{x})p(y,\vec{x}),$$

where  $w(\vec{x}) > 0$  and  $E_P[w(\vec{X})] = 1$ , which turns  $P_{Y,\vec{X}}^w$  into a joint probability distribution for  $(Y,\vec{X})$  with the same support as  $P_{Y,\vec{X}}$ . At times, for specific weight functions, we will write  $w(\vec{X})P_{Y,\vec{X}}$  instead of  $P_{Y,\vec{X}}^w$ .

LEMMA 4.1. 
$$P_{Y|\vec{X}}^w = P_{Y|\vec{X}}$$
 and  $P_{\vec{X}}^w = w(\vec{X})P_{\vec{X}}$ .

The proof is elementary and simplest in terms of densities:

$$p^{w}(\vec{x}) = \int p^{w}(y, \vec{x}) \, dy$$

$$= \int w(\vec{x}) p(y, \vec{x}) \, dy = w(\vec{x}) \int p(y, \vec{x}) \, dy$$

$$= w(\vec{x}) p(\vec{x}),$$

$$p^{w}(y|\vec{x}) = p^{w}(y, \vec{x}) / p^{w}(\vec{x})$$

$$= (w(\vec{x}) p(y, \vec{x})) / (w(\vec{x}) p(\vec{x}))$$

$$= p(y, \vec{x}) / p(\vec{x}) = p(y|\vec{x}).$$

We obtain as an immediate consequence.

PROPOSITION 4.1. If the regression functional  $\theta(\cdot)$  is well-specified for  $P_{Y|\vec{X}}$ , it is unchanged under arbitrary  $\vec{X}$ -dependent reweighting:

$$\theta(P_{\vec{Y}\vec{X}}^w) = \theta(P_{\vec{Y}\vec{X}}).$$

REMARK. In fixed-X linear models theory, which assumes correct specification, it is known that reweighting the data with fixed weights grants unbiased estimation of coefficients. Translated to the current framework, this fact returns as a statement of invariance of well-specified functionals under  $\vec{X}$ -dependent reweighting.

Tests of misspecification based on reweighting were proposed by White (1980, Section 4) for linear OLS. The approach generalizes to arbitrary types of regression and regression functionals as follows: Given a weight function  $w(\vec{X})$  normalized for P, the null hypothesis is  $H_0: \theta(P^w) = \theta(P)$ . For the case that  $\theta(\cdot)$  is the vector of OLS linear regression coefficients, White

(1980, ibid., Theorem 4) proposes a test statistic based on plug-in estimates and shows its asymptotic null distribution to be  $\chi^2$ . The result is a Hausman test (1978) whereby (using model-oriented language) an efficient estimate under the model is compared to an inefficient but consistent estimate. Rejection indicates misspecification. We will not draw on White's results but instead use the x-y bootstrap as a basis of inference because (1) it directly applies to general types of regression under mild technical conditions, and (2) it lends itself to augmentation of visual displays that provide more informative diagnostics than vanilla tests. White (1980) did not develop a methodology for reweighting tests other than recommending experimentation with multiple weight functions. The present goal is to introduce highly interpretable one-parameter families of weight functions and to illustrate their practical use to gain insights into the nature of misspecifications.

### 4.2 The Well-Specification Diagnostic: Population Version

To obtain interpretable weight functions, we construct them as functions of a *univariate variable Z*. This variable will often be one of the real-valued regressors,  $Z = X_j$ . However, the variable Z may be any function of the regressors,  $Z = f(\vec{X})$ , as when  $Z = \beta' \vec{X}$  is the OLS fit of Y, or  $Z = X_{j \bullet}$  is  $X_j$  adjusted for all other regressors (Part I, Section 9).

Given a variable Z, consider for concreteness a univariate Gaussian weight function of Z, centered at  $\xi$  on the Z axis:

(8) 
$$w_{\xi}(z) = w_{\xi}^{*}(z) / E[w_{\xi}^{*}(Z)],$$

$$w_{\xi}^{*}(z) \propto \exp(-(z - \xi)^{2} / (2\gamma^{2})),$$

where  $\gamma$  is a user-specified bandwidth parameter (see Section 4.3 below).

Next, consider a one-dimensional regression functional  $\theta(P)$ , such as a linear regression slope. A graphical diagnostic is obtained by plotting  $\theta(\cdot)$  as a function of the reweighting centers  $\xi$ :

(9) 
$$\xi \mapsto \theta_{\xi}(\mathbf{P}) = \theta(w_{\xi}(Z)\mathbf{P}).$$

If the regression functional  $\theta(P)$  is well-specified for  $P_{Y|\vec{X}}$ , then  $\theta_{\xi}(P)$  is constant in  $\xi$  and equal to  $\theta(P)$ . Equivalently, if  $\theta_{\xi}(P)$  is not constant in  $\xi$ , then  $\theta(P)$  is misspecified. Thus nonconstancy is a sufficient criterion for misspecification. Insightful choices of traces of the form (9) will be proposed below.

#### 4.3 The Reweighting Diagnostic: Data Version

To make the diagnostic actionable on data, one obtains estimates

$$\hat{\theta}_{\xi} = \theta (\hat{w}_{\xi}(Z) \hat{P}),$$

where  $\hat{w}_{\xi}(x)$  is a weight function that is empirically normalized to unit mass,  $\hat{E}[\hat{w}_{\xi}(Z)] = 1$ , where  $\hat{E}[...]$  denotes the sample average. This means using weights for the observations of the form

$$w_i = \hat{w}_{\xi}(z_i) \propto \exp(-(z_i - \xi)^2 / (2\gamma^2)),$$
  
$$\frac{1}{N} \sum_i w_i = 1, \quad i = 1, \dots, N.$$

We parametrize the bandwidth  $\gamma = \alpha \hat{\sigma}(Z)$  in terms of the empirical standard deviation  $\hat{\sigma}(Z)$  of Z and a multiplier  $\alpha$ . In the examples, we use  $\alpha = 1$ .

In order to plot a discretized version of the trace  $\xi \mapsto \hat{\theta}_{\xi}$ , we obtain estimates  $\hat{\theta}_{\xi}$  for a grid of values  $\xi_{(1)} < \cdots < \xi_{(K)}$  on the Z axis, a simple choice being the interior deciles of the empirical Z distribution. Hence K=9, unless Z has numerous ties, causing some deciles to collapse. Finally, we plot  $\xi_{(k)} \mapsto \hat{\theta}_{\xi_{(k)}}$ . This is carried out in Figures 1–3 for the LA homeless data (see Section 5).

#### 4.4 Interpretations of the Reweighting Diagnostic

The reweighting diagnostic is likely to be accessible to practitioners of regression. One reason is that the restriction to weights as a function of a univariate variable Z permits a simple left-to-right comparison: Is  $\xi \mapsto \theta(w_{\xi}(Z)P)$  higher or lower on the right than on the left? In our experience, the dominant feature of such traces is indeed monotonicity. The intuitive appeal of reweighting is further helped by two mutually compatible interpretations:

- Data frequency: Reweighting mimics scenarios of datasets that contain more or fewer observations as a function of Z than the observed dataset. Thus it answers questions such as "what if there were more observations with low (or high) values of Z?" In this sense, reweighting mimics alternative data sources based on the data at hand.
- Conditioning: Reweighting can be seen as "soft conditioning on Z" in the sense that conditioning on "sharp inclusion" in an interval  $\xi c < Z < \xi + c$  is replaced by "soft inclusion" according to the weight function  $w_{\xi}(z)$ . In this sense reweighting localizes the regression functional. However, note that when  $Z = X_j$ , for example, the localization is of "codimension 1" in regressor space.

<sup>&</sup>lt;sup>3</sup>Mathematically, the restriction to weights as a function of univariate variables Z is no restriction at all because any  $w(\vec{x})$  can be trivially described as the identity function of  $Z = w(\cdot)$ .

In what follows, we use either of these interpretations depending on the context.

### 4.5 Inferential Features for Reweighting Diagnostics

Graphical diagnostics need inferential augmentation to answer questions of whether visually detected features are likely to be real. Presently the two main questions are:

- (1) Is the variation/nonconstancy in  $\xi_{(k)} \mapsto \hat{\theta}_{\xi_{(k)}}$  sufficiently strong to be statistically significant, and hence suggest misspecification of  $\theta(\cdot)$ ?
- (2) Where are the estimates  $\hat{\theta}_{\xi(k)}$  statistically significantly different from zero?

For regression slopes, question (2) may be more relevant than (1) because one usually cares about their statistical significance. Therefore, to answer question (2), we decorate the diagnostic plot with traces of bootstrapped estimates, as shown in the plots of Figures 1– 3. Bootstrap resampling is done from the actual, not the reweighted, data. The weight functions have the same centers  $\xi_{(k)}$ , but their bandwidth is based on bootstrapped standard deviations. In the figures, we show 199 bootstrap traces in gray color, amounting to a socalled "spaghetti plot." Along with the bootstrap replications we also show bootstrap error bars at the grid locations. Their widths are a lenient  $\pm 2$  bootstrap standard errors, not adjusted for multiplicity. (Such diagnostics are search expeditions and do not provide strict inference.)

As can be illustrated with Figures 1–3, statistical significance can feature a variety of patterns. Significance may exist...

- (2a) ... across the whole range of reweighting centers  $\xi_{(k)}$  and in the same direction, as in the top right plot of Figure 1;
- (2b) ... both on the left and the right but in opposite directions with a transition through insignificance in between, as is nearly the case in the center left plot of Figure 2;
- (2c) ... over part of the range, typically the left or the right side; such tendencies are seen in the two center plots of Figure 1;
- (2d) ... nowhere, as in the bottom right plot of Figure 2.

To answer question (1) regarding the presence of misspecification, we piggyback on the bootstrap exercise meant to answer question (2). Because most detections of misspecification arise from a monotone tilt in the trace  $\xi_{(k)} \mapsto \hat{\theta}_{\xi_{(k)}}$ , we construct a cheap test statistic by forming the difference between the two extreme points of the trace,  $\hat{\theta}_{\xi_{(K)}} - \hat{\theta}_{\xi_{(1)}}$ . We obtain its bootstrap distribution almost for free, hence we can perform a crude bootstrap test by placing the null value zero in the bootstrap distribution. The bootstrap p-value and the test statistic are shown near the top of each plot frame in Figures 1–3. For example, the top left frame of Figure 1 shows "Tilt: p=0.04 d=2.18," meaning that the difference of 2.18 is statistically significant with a (two-sided) p-value of 0.04.

Finally, we show on the left side of each frame a visual version of unweighted plain statistical significance of the quantity of interest in the form of a bootstrap confidence interval around the unweighted estimate  $\hat{\theta} \pm 2$  unweighted bootstrap standard errors. In addition, we show 199 bootstrap estimates (gray points horizontally jittered to reduce overplotting). The location on the horizontal axis has no meaning other than being 10% to the left of the range  $(\xi_{(1)}, \xi_{(K)})$  of the traces.

# 5. THE REWEIGHTING DIAGNOSTIC FOR WELL-SPECIFICATION: METHODOLOGY AND EXAMPLES

The following subsections demonstrate three different purposes of the diagnostic. The quantities of interest are linear OLS slopes, though the approach generalizes to all types of regression that permit reweighting:

- Focal slope: Expose a slope  $\beta_k(P)$  of special interest to reweighting on each regressor in turn:  $Z = X_j$  for j = 1, ..., p (Section 5.1). This produces highly interpretable insights into interactions of regressor  $X_k$  with all other regressors  $X_j$ , without modeling these interactions directly.
- Nonlinearity detection: Expose each regression slope  $\beta_j(P)$  to reweighting on its own regressor,  $Z = X_j$  (Section 5.2). This produces insights into marginal nonlinear behaviors of response surfaces.
- Focal reweighting variable: Use a single reweighting variable of interest (here:  $Z = \beta' \vec{X}$ ) to diagnose well-specification of all components of a regression functional, here: slopes  $\beta_i(P)$  (Section 5.3).

<sup>&</sup>lt;sup>4</sup>This test statistic does not result in a Hausman (1978) test: both estimates are "inefficient under correct model specification." However, it quantifies an obvious visual feature of the traces.

<sup>&</sup>lt;sup>5</sup>For 199 bootstrap replicates, the lowest possible two-sided p-value is  $0.01 = 2 \cdot 1/(1 + 199)$ .

These diagnostics will be illustrated with the LA homeless data of Part I, Section 2 (Berk et al. 2008). The observations consist of a sample of 505 census tracts in the LA metropolitan area, and the variables are seven quantitative measures of the tracts with largely self-explanatory names: The response is the StreetTo-tal (count) of homeless people in a census tract, and the six regressors are: MedianIncome (of households, in \$1000s), PercMinority, and the prevalences of four types of lots: PercCommercial, PercVacant, PercResidential and PercIndustrial.

# 5.1 Diagnostics for a Focal Regression Coefficient of Interest (Figure 1)

One variable stands out as potentially accessible to intervention by public policies: PercVacant. Vacant lots could be turned into playgrounds, sports fields, parks, or offered as neighborhood gardens. It would therefore be of interest to check whether the regression coefficient of PercVacant possibly measures a causal effect, for which it is a necessary condition that it be well-specified (Section 3.4). To this end, Figure 1 shows diagnostics for the coefficient of PercVacant under reweighting on all six regressors.

As the plots show, statistical significance of the coefficient of PercVacant holds by and large under reweighting across the ranges of all six regressors. While this is comforting, there exists a weakening of significance in the extremes of the ranges of three regressors: high MedianIncome, low PercMinority and low PercResidential. With these qualitative observations, it is already indicated that wellspecification of the coefficient of PercVacant is doubtful, and indeed the tilt tests show statistical significance with 2-sided p-values of 0.01 and 0.02 for PercMinority and MedianIncome, respectively. The variable PercResidential also looks rather steep, but its tilt test has a weaker p-value around 0.1. Finally, a very weak indication is shown for larger effects at higher levels of PercVacant.

Does this indication of misspecification invalidate a causal effect of PercVacant? It does not. It only points to the likely possibility that the causal effect is not correctly described by a single linear regression coefficient; it is rather a more complex function of the regressors. Useful insight into the nature of the causal

effect (if this is what it is) can be gleaned from the diagnostic plots by using them to answer an obvious question: Where is the effect of PercVacant likely to be strong? An answer might indeed help in prioritizing interventions. Interpreting the plots of Figure 1 liberally, one could state that the effect of PercVacant looks strongest for census tracts with high PercMinority, followed by high PercResidential and low MedianIncome. These observations seem rather plausible and may indeed point to census tracts worth prioritizing for intervention with public policies.<sup>7</sup>

The insights gained so far point to the presence of interactions between PercVacant and other regressors because the slope of PercVacant varies at different levels of those other regressors. A natural next step would be more detailed modeling that includes interactions between PercVacant and the three interacting regressors, but the essential insights have already been gained.

# 5.2 Diagnostics for Slopes Reweighted by Their Own Regressors (Figure 2)

The top right plot in Figure 1 is a special case where the slope of interest is reweighted by its own regressor, PercVacant. It has a different interpretation, not related to interactions but to nonlinear effects. To get a better picture of the possibilities that can arise in real data, we show in Figure 2 the corresponding plots for all six regressors and their slopes.

Glancing at the six plots, we note some unpredictable effects of reweighting, both on the values and the estimation uncertainties of the slopes. We find examples of larger and smaller estimates as well as stronger and weaker statistical significances relative to their unweighted analogs:

- Bottom left plot for the regressor PercCommercial: The unweighted estimate of  $\beta_j(P)$  (on the left side of the plot) is weakly statistically significant (the lower end of the  $\pm 2$  standard error confidence interval touches zero). The reweighted estimates of  $\beta_j(w_\xi(X_j)P)$ , however, are closer to zero and nowhere statistically significant for any  $\xi$  in the range of PercCommercial.
- Top right plot for the regressor PercVacant: The unweighted estimate and the reweighted estimates are all statistically significant, but the reweighted ones are systematically larger and much more statistically significant.

<sup>&</sup>lt;sup>6</sup>Such programs have indeed been enacted in some cities. We abstain from commenting on the controversies surrounding such policies.

<sup>&</sup>lt;sup>7</sup>An application of this type of diagnostic to the Boston Housing data is in Appendix F.

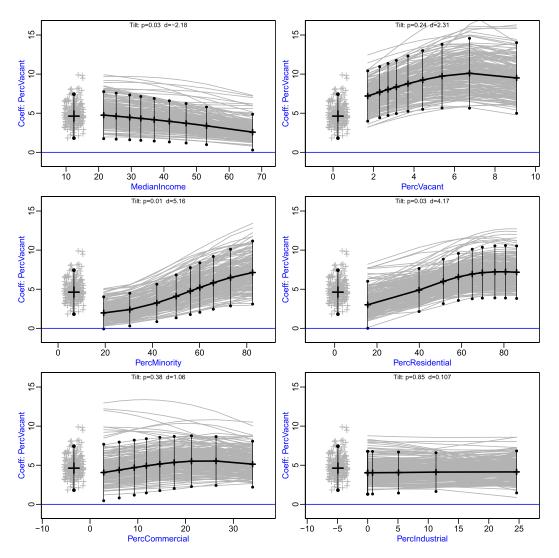


FIG. 1. Diagnostics for the slope of PercVacant; LA Homeless Data (see Section 2, Part I). Vertical axis = regression coefficient of PercVacant (in all frames); horizontal axes = regressors. If the vertical axis is interpreted causally as the effect size of PercVacant on the response StreetTotal (of homeless in a census tract), the following can be inferred: the effect size of PercVacant is greater for high values of PercMinority (center left frame) and low values of MedianIncome (top left frame), and possibly also for high values of PercResidential. Near the top margin of each frame are the p-values of a "Tilt" test for the difference between the right-most and left-most effect sizes.

Another noteworthy case of a different nature appears for the regressor PercMinority (Figure 2, center left plot). While the unweighted estimate is statistically insignificant, the locally reweighted estimates reveal a striking pattern:

- For low values of PercMinority  $\approx 20\%$ , the slope is negative and statistically significant: Incrementally, more minorities are associated with a lower StreetTotal of homeless.
- $\bullet$  For high values of PercMinority  $\approx 80\%$ , the slope is positive and (weakly) statistically signifi-

cant: Incrementally, more minorities are associated with a higher StreetTotal of homeless.

This finding (if real) represents a version of Simpson's paradox: In aggregate, there is no statistically significant association, but, conditional on low and high values of PercMinority, there is, and in opposite directions.

In Appendix E, we discuss some reasons for the unpredictable behaviors of slopes under reweighting wrt to their own regressors. We also mention a (weak) link to partial additive models (Hastie and Tibshirani, 1990) with one nonlinear term.

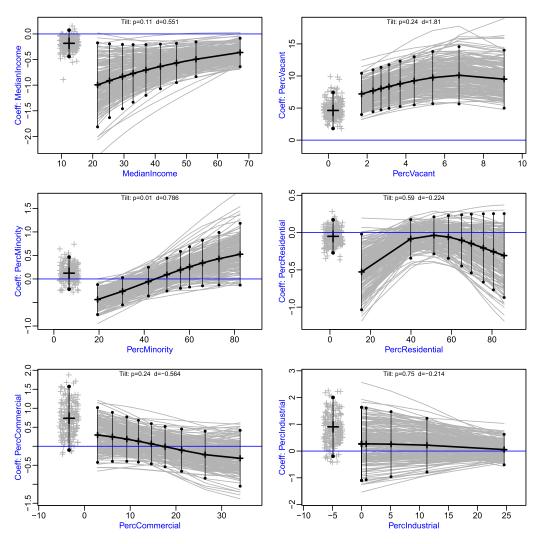


FIG. 2. Misspecification diagnostics: Slopes reweighted by their own regressors—indications of nonlinearity. The center left frame suggests that the slope of PercMinority reverses from a negative slope for low values of PercMinority to a positive slope for high values of PercMinority.

# 5.3 Diagnostics for a Focal Reweighting Variable of Interest (Figure 3)

Next, we illustrate a version of the diagnostics that subjects all slopes of a linear regression to a single reweighting variable of interest. The goal is to detect misspecification in any coefficient, and the hope is to do so by reweighting based on a variable Z that is both powerful and interpretable. Taking a cue from traditional residual diagnostics, we choose the OLS best approximation,  $Z = \beta' \vec{X}$ . The data version is based on reweighting as a function of the fitted values,  $z_i = \hat{y}_i = \hat{\beta}' \vec{x}_i$ . The question is whether any coefficient reveals

misspecification when comparing it on data with more low versus more high values of the linear approximation to the number of homeless. The expectation is that the gradient of the linear approximation should be a direction of high average response variation, and hence may have a higher promise of revealing misspecifications than other directions in regressor space.

Figure 3 shows this diagnostic applied to the LA homeless data, labeling the reweighting variable as Fitted. Some observations are as follows:

 The only slope with signs of misspecification is for MedianIncome (top left plot), whose tilt test has a p-value of 0.03. This slope achieves mild statistical significance for high values of Fitted, which would indicate that the "effect" (if any) of differ-

<sup>&</sup>lt;sup>8</sup>The estimated slope vector  $\hat{\beta}$  is frozen across bootstraps, ignoring a lower-order source of sampling variability.

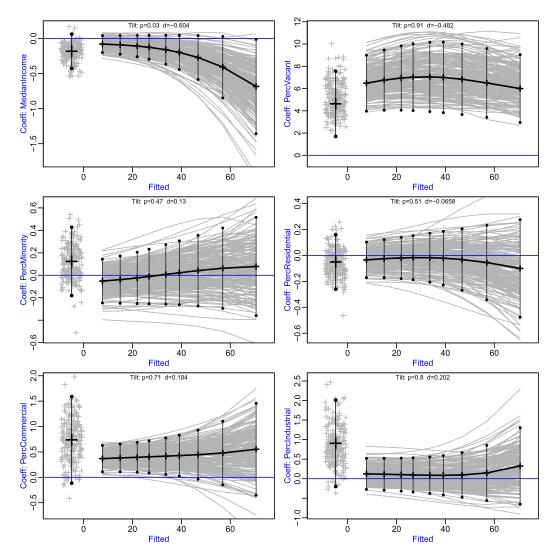


FIG. 3. Misspecification Diagnostics using one focal reweighting variable, the best linear approximation/prediction Fitted, for all slopes. The fan shapes left to right suggest that all slope estimates except the one for PercVacant have more sampling variability for higher fitted values.

ences in MedianIncome matter more for high values of Fitted.

- The slope of PercCommercial (bottom left plot) shows no signs of misspecification, but it is mildly statistically significant only for low values of Fitted due to the lower estimation uncertainty in that range.
- Five of the six plots feature a fan shape of the bootstrap spaghetti bands (exception: PercVacant).
   This indicates that these five slope estimates have greater estimation uncertainty for higher values of Fitted.

The last point illustrates that the diagnostic is not only informative about the average level of estimates but also about their estimation uncertainty.

### 5.4 Summary Comments on Reweighting Diagnostics

The reweighting diagnostics proposed here are not meant to replace other types of diagnostics, typically based on residual analysis. They are, however, able to answer questions about quantities of interest and effects of regressors that residual analysis might not. They may also be able to provide insights into the nature of nonlinearities and interactions without explicitly modeling them. Furthermore, they are easily augmented with inferential features such as bootstrap spaghetti bands and (tentative) tests of misspecification with specific interpretations. Finally, they are able to localize regions in regressor space with high or low estimation uncertainty.

### 6. ESTIMATION OF REGRESSION FUNCTIONALS: CANONICAL DECOMPOSITION OF ESTIMATION OFFSETS

We return to the task of building a general framework of plug-in estimation of regression functionals based on i.i.d. data. We decompose sampling variability into its two sources, one due to the conditional response distribution, the other due to the randomness of the regressors interacting (conspiring) with misspecification. Along the way we find new characterizations of well-specification of regression functionals.

#### 6.1 Regression Data and Plug-in Estimation

We adopt some of the notations and assumptions from Part I, Section 5: Data consist of N i.i.d. draws  $(Y_i, \vec{X}_i) \sim P = P_{V\vec{X}}$ ; the responses  $Y_i$  are collected in a data structure  $Y = \{Y_i\}_i$ , and the regressors  $\vec{X}_i$  in another data structure  $X = \{\vec{X}_i\}_i$ , called "data frame" in programming languages such as R (2008). We avoid the terms "vector" and "matrix" because in a general theory of regression all variables—responses and regressors—can be of any type and of any dimension.<sup>9</sup> This is why not only X but Y is best thought of as a (random) "data frame." Regression of Y on X is any attempt at estimating aspects of the conditional distribution  $P_{V|\vec{X}}$ . We limit ourselves to regression functionals  $\boldsymbol{\theta}(\cdot)$  that allow plug-in estimation  $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}(\widehat{\boldsymbol{P}})$  where  $\widehat{\boldsymbol{P}} = \widehat{\boldsymbol{P}}_{Y,\vec{X}} = (1/N) \sum \delta_{(Y_i,\vec{X}_i)}$  is the joint empirical distribution. If necessary, we may write  $\hat{\boldsymbol{P}}_N$  for  $\hat{\boldsymbol{P}}$  and  $\hat{\theta}_N$  for  $\hat{\theta}$ . In addition, we will also need the empirical regressor distribution  $\hat{P}_{\vec{X}} = (1/N) \sum \delta_{\vec{X}_i}$ .

# 6.2 The Conditional Parameter of Model-Trusting Fixed-X Regression

We now define the important notion of a "conditional parameter" for arbitrary regression functionals, thereby providing the target of estimation for fixed-X theories. For OLS slopes, this target of estimation is  $\boldsymbol{\beta}(X) = \boldsymbol{E_P}[\hat{\boldsymbol{\beta}}|X]$  (Part I, Section 5). We use the idea that fixed-X theories condition on observed regressor observations  $\vec{X}_1, \ldots, \vec{X}_N$ , collected in the data frame X, and define a target of estimation by assuming that the population of Y-values at each  $\vec{X}_i$  is known:  $Y_i | \vec{X}_i \sim P_{Y|\vec{X}_i}$ . The joint distribution is then effectively  $P_{Y|\vec{X}} \otimes \hat{P}_{\vec{X}}$ , amounting to partial plug-in of  $\hat{P}_{\vec{X}}$ 

for  $P_{\vec{X}}$  in  $P_{Y,\vec{X}} = P_{Y|\vec{X}} \otimes P_{\vec{X}}$ . The conditional parameter for  $\theta(\cdot)$  is therefore defined as  $\theta(X) = \theta(P_{Y|\vec{X}} \otimes \widehat{P}_{\vec{X}})$ . We summarize notation, with emphasis on the second line:

$$\begin{split} \theta(\boldsymbol{P}) &= \theta(\boldsymbol{P}_{Y|\vec{X}} \otimes \boldsymbol{P}_{\vec{X}}), \\ \theta(\boldsymbol{X}) &= \theta(\boldsymbol{P}_{Y|\vec{X}} \otimes \boldsymbol{\hat{P}}_{\vec{X}}), \quad \boldsymbol{\hat{P}}_{\vec{X}} = (1/N) \sum \delta_{\vec{X}_i}, \\ \boldsymbol{\hat{\theta}} &= \theta(\boldsymbol{\hat{P}}). \end{split}$$

Note that X and  $\widehat{P}_{\vec{X}}$  contain the same information; the conditional response distribution  $P_{Y|\vec{X}}$  is implied and not shown in  $\theta(X)$ . The main points are:

- In model-trusting theories that condition on X, the target of estimation is  $\theta(X)$ . They assume  $\theta(X)$  is the same for all acceptable X.
- In model-robust theories that do not condition on X, the target of estimation is  $\theta(P)$ , whereas  $\theta(X)$  is a random quantity (Corollary 6.3 below).

The above definitions can be made more concrete by illustrating them with the specific ways of defining regression functionals of Section 2:

• Functionals defined through minimization of objective functions:

$$\begin{aligned} \boldsymbol{\theta}(\boldsymbol{P}) &= \operatorname{argmin}_{\boldsymbol{\theta}} \boldsymbol{E}_{\boldsymbol{P}} \big[ \mathcal{L}(\boldsymbol{\theta}; Y, \vec{\boldsymbol{X}}) \big], \\ \boldsymbol{\theta}(\boldsymbol{X}) &= \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i} \boldsymbol{E}_{\boldsymbol{P}} \big[ \mathcal{L}(\boldsymbol{\theta}; Y_{i}, \vec{\boldsymbol{X}}_{i}) \mid \vec{\boldsymbol{X}}_{i} \big], \\ \hat{\boldsymbol{\theta}} &= \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i} \mathcal{L}(\boldsymbol{\theta}; Y_{i}, \vec{\boldsymbol{X}}_{i}). \end{aligned}$$

• Functionals defined through estimating equations:

$$\theta(P): E_{P}[\psi(\theta Y, \vec{X})] = \mathbf{0},$$

$$\theta(X): \frac{1}{N} \sum_{i} E_{P}[\psi(\theta; Y_{i}, \vec{X}_{i}) | \vec{X}_{i}] = \mathbf{0},$$

$$\hat{\theta}: \frac{1}{N} \sum_{i} \psi(\theta; Y_{i}, \vec{X}_{i}) = \mathbf{0}.$$

Summary: Among the three cases in each bullet, the most impenetrable but also most critical case is the second one. It defines the "conditional parameter" through partial plug-in of the empirical regressor distribution. The conditional parameter is the target of fixed-X regression for arbitrary types of regression functionals.

#### 6.3 Estimation Offsets

The conditional parameter  $\theta(X)$  enables us to distinguish between two sources of estimation uncertainty: (1) the conditional response distribution and (2) the marginal regressor distribution. To this end, we defined

<sup>&</sup>lt;sup>9</sup>Recall that the typographic difference between Y and  $\vec{X}$  is a holdover from Part I, where the response and all regressors were assumed univariate quantitative.

in Part I for linear OLS what we call "estimation offsets." With the availability of  $\theta(X)$  for regression functionals, these can be defined in full generality:

$$Total EO = \hat{\theta} - \theta(P),$$

$$Noise EO = \hat{\theta} - \theta(X),$$

$$Approximation EO = \theta(X) - \theta(P).$$

The total EO is the offset of the plug-in estimate from its population target. The noise EO is the component of the total EO that is due to the conditional distribution  $Y|\vec{X}$ . The approximation EO is the part due to the randomness of  $\vec{X}$  under misspecification. These interpretations will be elaborated in what follows.

REMARK. We repeat an observation made in Part I, end of Section 5. The approximation EO  $\theta(X)$  –  $\theta(P)$  could be misinterpreted as a bias because it is the difference of two targets of estimation. This interpretation is *wrong*. In the presence of misspecification, the approximation EO is a nonvanishing random variable. It will be shown to contribute not a bias to  $\hat{\theta}$  but a  $N^{-1/2}$  term to the sampling variability of  $\hat{\theta}$ .

## 6.4 Well-Specification in Terms of Approximation EOs

The approximation EO lends itself for another characterization of well-specification.

PROPOSITION 6.1. Assume  $P_{\vec{X}} \mapsto \theta(P_{Y|\vec{X}} \otimes P_{\vec{X}})$  is continuous in the weak topology. Then  $\theta(\cdot)$  is well-specified for  $P_{Y|\vec{X}}$  iff  $\theta(X) - \theta(P) = 0$  for all acceptable X.

PROOF. If  $\theta(\cdot)$  is well-specified in the sense of Section 3, then

$$\theta(X) = \theta(P_{Y|\vec{X}} \otimes \widehat{P}_{\vec{X}}) = \theta(P_{Y|\vec{X}} \otimes P_{\vec{X}}) = \theta(P).$$

The converse follows because the empirical regressor distributions  $\hat{P}_{\vec{X}}$  (for  $N \to \infty$ ) form a weakly dense subset in the set of all regressor distributions, and the regression functional is assumed continuous in this argument.  $\square$ 

COROLLARY 6.1. Same assumptions as in Proposition 6.1.

- Fixed-X and random-X theories estimate the same target iff  $\theta(\cdot)$  is well-specified for  $P_{Y|\vec{X}}$ .
- $\theta(\cdot)$  is well-specified for  $P_{Y|\bar{X}}$  iff  $V_P[\bar{\theta}(X)] = 0$  for all acceptable  $P_{\bar{X}}$ .

The first bullet confirms that the notion of well-specification for regression functionals hits exactly the point of agreement between theories that condition on the regressors and those that treat them as random. The second bullet leads the way to the fact that a misspecified regression functional will incur sampling variability originating from the randomness of the regressors.

### 6.5 Deterministic Association Annihilates the Noise EO

While well-specification addresses a vanishing approximation EO, one can also consider the dual concept of a vanishing noise EO. Here is a sufficient condition under which the noise EO vanishes for all regression functionals.

PROPOSITION 6.2. If  $Y = f(\vec{X})$  is a deterministic function of  $\vec{X}$ , then  $\hat{\theta} - \theta(X) = 0$  for all regression functionals.

PROOF. The conditional response distribution is  $P_{Y|\vec{X}=\vec{x}} = \delta_{y=f(\vec{x})}$ , hence the joint distribution formed from  $P_{Y|\vec{X}=\vec{x}}$  and  $\hat{P}_{\vec{X}}$  is  $\hat{P}$ :  $P_{Y|\vec{X}} \otimes \hat{P}_{\vec{X}} = \hat{P}$ . It follows that  $\theta(X) = \theta(P_{Y|\vec{X}} \otimes \hat{P}_{\vec{X}}) = \theta(\hat{P}) = \hat{\theta}$ .

The proposition illustrates the fact that the noise EO is due to "noise," that is, variability of Y conditional on  $\vec{X}$ . Thus, although less transparent than in linear OLS, the conditional response distribution  $Y|\vec{X}$  is the driver of the noise EO.

#### 6.6 Well-Specification and Influence Functions

This section introduces influence functions for regression functionals which will prove useful for approximations in Section 6.7 and for asymptotic decompositions in Section 7. For background on influence functions see, for example, Hampel et al. (1986) and Rieder (1994).

The influence function is a form of derivative on the space of probability distributions, which makes it an intuitive tool to characterize well-specification of regression functionals: If  $\theta(P_{Y|\vec{X}}\otimes P_{\vec{X}})$  is constant in the argument  $P_{\vec{X}}$  at a fixed  $P_{Y|\vec{X}}$ , then this means intuitively that the "partial derivative" wrt  $P_{\vec{X}}$  vanishes.

The definition of the full influence function of  $\theta(\cdot)$  is as follows:

(10) 
$$IF(y, \vec{x}) = \frac{d}{dt} \Big|_{t=0} \theta ((1-t)P + t\delta_{(y,\vec{x})}).$$

We omit  $\theta(\cdot)$  as well as  $P = P_{Y,\vec{X}}$  as arguments of  $IF(y,\vec{x})$  because both will be clear from the context, except for one occasion in Appendix C where we write

 $IF(y, \vec{x}; P)$ . More relevant is the following definition of the partial influence function of  $\theta(\cdot)$  with regard to the regressor distribution:

(11) 
$$IF(\vec{x}) = \frac{d}{dt}\Big|_{t=0} \theta(P_{Y|\vec{X}} \otimes ((1-t)P_{\vec{X}} + t\delta_{\vec{x}})).$$

For derivations of the following lemma and proposition, see Appendix C.

LEMMA 6.1. 
$$IF(\vec{x}) = E_P[IF(Y, \vec{X})|\vec{X} = \vec{x}].$$

PROPOSITION 6.3. A regression functional  $\theta(\cdot)$  with an influence function at  $P_{Y,\vec{X}}$  is well-specified for  $P_{Y|\vec{X}}$  iff  $IF(\vec{x}) = 0 \ \forall \vec{x}$ .

### 6.7 Approximating Estimation Offsets with Influence Functions

For linear OLS, the definition of EOS and the lemma in Section 5 of Part I exhibited an intuitive correspondence between the total, noise and approximation EO on the one hand and the population residual, the noise and the nonlinearity on the other hand. No such direct correspondence exists for general types of regression. The closest general statement about EOs is in terms of approximations based on influence functions. Assuming asymptotic linearity of  $\theta(\cdot)$ , the EOs have the following approximations to order  $o_P(N^{-1/2})$ , which will lead straight to the CLTs of the next section.

Total EO:

$$\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}(\boldsymbol{P}) \approx \frac{1}{N} \sum_{i} \boldsymbol{IF}(Y_i, \vec{\boldsymbol{X}}_i),$$

Noise EO:

$$\hat{\theta} - \theta(X) \approx \frac{1}{N} \sum_{i} (IF(Y_i, \vec{X}_i) - E_P[IF(Y, \vec{X}_i) | \vec{X}_i]),$$

Approximation EO:

$$\theta(X) - \theta(P) \approx \frac{1}{N} \sum_{i} E_{P}[IF(Y, \vec{X}_{i})|\vec{X}_{i}].$$

### 7. MODEL-ROBUST CENTRAL LIMIT THEOREMS DECOMPOSED

### 7.1 CLT Decompositions Based on Influence Functions

If the approximations of Section 6.7 hold, the EOs obey the following CLTs:

$$\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}(\boldsymbol{P})) \xrightarrow{\mathcal{D}} \\
\mathcal{N}(\boldsymbol{0}, V_{\boldsymbol{P}}[\boldsymbol{I}\boldsymbol{F}(\boldsymbol{Y}, \vec{\boldsymbol{X}})]), \\
\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}(\boldsymbol{X})) \xrightarrow{\mathcal{D}} \\
\mathcal{N}(\boldsymbol{0}, E_{\boldsymbol{P}}[V_{\boldsymbol{P}}[\boldsymbol{I}\boldsymbol{F}(\boldsymbol{Y}, \vec{\boldsymbol{X}})|\vec{\boldsymbol{X}}]]), \\
\sqrt{N}(\boldsymbol{\theta}(\boldsymbol{X}) - \boldsymbol{\theta}(\boldsymbol{P})) \xrightarrow{\mathcal{D}} \\
\mathcal{N}(\boldsymbol{0}, V_{\boldsymbol{P}}[E_{\boldsymbol{P}}[\boldsymbol{I}\boldsymbol{F}(\boldsymbol{Y}, \vec{\boldsymbol{X}})|\vec{\boldsymbol{X}}]]).$$

These are immediate consequences of the assumed asymptotic linearities. The asymptotic variances of the EOs follow the canonical decomposition,

$$V_{P}[IF(Y, \vec{X})] = E_{P}[V_{P}[IF(Y, \vec{X})|\vec{X}]] + V_{P}[E_{P}[IF(Y, \vec{X})|\vec{X}]],$$

the three terms being the asymptotic variance-covariance matrices of the total, the noise and the approximation EO, respectively. Implicit in this Pythagorean formula is that  $IF(Y, \vec{X}) - E_P[IF(Y, |\vec{X})]$  and  $E_P[IF(Y, |\vec{X})]$  are orthogonal to each other, which implies that the noise EO and the approximation EO are asymptotically orthogonal. Asymptotic orthogonalities based on conditioning are well-known in semi-parametric theory. For linear OLS, this orthogonality holds exactly for finite N due to  $\beta(X) = E[\hat{\beta}|X]$ :  $V_P[\hat{\beta} - \beta(X), \beta(X) - \beta(P)] = 0$ .

The following corollary is a restatement of Proposition 6.3, but enlightened by the fact that it relies on the asymptotic variance of the approximation EO.

COROLLARY 7.1. The regression functional  $\theta(\cdot)$  is well-specified for  $P_{Y|\bar{X}}$  iff the asymptotic variance of the approximation EO vanishes for all acceptable  $P_{\bar{X}}$ .

PROOF. Using careful notation, the condition says  $V_{P_{\vec{X}}}[E_{P_{Y|\vec{X}}}[IF(Y,\vec{X})|\vec{X}]] = \mathbf{0}$  for all acceptable  $P_{\vec{X}}$ . This in turn means  $E_{P_{Y|\vec{X}}}[IF(Y,\vec{X})|\vec{X}=\vec{x}] = \mathbf{0}$  for all  $\vec{x}$ , which is the condition of Proposition 6.3.  $\square$ 

#### 7.2 CLT Decompositions for EE Functionals

For EE functionals, the influence function is  $IF(y, \vec{x}) = \Lambda(\theta)^{-1} \psi(\theta; y, \vec{x})$  where  $\theta = \theta(P)$  and  $\Lambda(\theta) = \nabla_{\theta} E_{P}[\psi(\theta; Y, \vec{X})]$  is the Jacobian of size  $q \times q$ ,  $q = \dim(\psi) = \dim(\theta)$ . The CLTs specialize as follows:

$$\begin{split} & \sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \stackrel{\mathcal{D}}{\longrightarrow} \\ & \mathcal{N}\big(\boldsymbol{0}, \boldsymbol{\Lambda}(\boldsymbol{\theta})^{-1} \boldsymbol{V}_{P}\big[\boldsymbol{\psi}(\boldsymbol{\theta}; \boldsymbol{Y}, \vec{\boldsymbol{X}})\big] \boldsymbol{\Lambda}(\boldsymbol{\theta})^{\prime - 1}\big), \\ & \sqrt{N}\big(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}(\boldsymbol{X})\big) \stackrel{\mathcal{D}}{\longrightarrow} \\ & \mathcal{N}\big(\boldsymbol{0}, \boldsymbol{\Lambda}(\boldsymbol{\theta})^{-1} \boldsymbol{E}_{P}\big[\boldsymbol{V}_{P}\big[\boldsymbol{\psi}(\boldsymbol{\theta}; \boldsymbol{Y}, \vec{\boldsymbol{X}})|\vec{\boldsymbol{X}}\big]\big] \boldsymbol{\Lambda}(\boldsymbol{\theta})^{\prime - 1}\big), \\ & \sqrt{N}\big(\boldsymbol{\theta}(\boldsymbol{X}) - \boldsymbol{\theta}\big) \stackrel{\mathcal{D}}{\longrightarrow} \\ & \mathcal{N}\big(\boldsymbol{0}, \boldsymbol{\Lambda}(\boldsymbol{\theta})^{-1} \boldsymbol{V}_{P}\big[\boldsymbol{E}_{P}\big[\boldsymbol{\psi}(\boldsymbol{\theta}; \boldsymbol{Y}, \vec{\boldsymbol{X}})|\vec{\boldsymbol{X}}\big]\big] \boldsymbol{\Lambda}(\boldsymbol{\theta})^{\prime - 1}\big). \end{split}$$

The first line is Huber's (1967, Section 3) result. The three asymptotic variances have the sandwich form and are related according to the canonical decomposition,

$$V_{P}[\psi(\theta; Y, \vec{X})] = E_{P}[V_{P}[\psi(\theta; Y, \vec{X})|\vec{X}]] + V_{P}[E_{P}[\psi(\theta; Y, \vec{X})|\vec{X}]],$$

the terms again relating to the total EO, the noise EO and the approximation EO.

### 7.3 Implications of the CLT Decompositions

We address once again potential confusions relating to different notions of bias. Misspecification, in traditional parametric modeling, is sometimes called "model bias" which, due to unfortunate terminology, may suggest a connection to estimation bias,  $E_{P}[\hat{\theta}_{N}] - \theta(P)$ . Importantly, there is no connection between the two notions of bias. Estimation bias typically vanishes at a rate faster than  $N^{-1/2}$  and does not contribute to standard errors derived from asymptotic variances. Model bias, on the other hand, which is misspecification, generates in conjunction with the randomness of the regressors a contribution to the standard error, and this contribution is asymptotically of order  $N^{-1/2}$ , the same order as the better known contribution due to the conditional noise in the response. This is what the CLT decomposition shows. It also shows that the two sources of sampling variability are asymptotically orthogonal. In summary:

Model bias/misspecification does not create estimation bias; it creates sampling variability to the same order as the conditional noise in the response.

# 8. PLUG-IN/SANDWICH ESTIMATORS VERSUS M-OF-N BOOTSTRAP ESTIMATORS OF STANDARD ERROR

### 8.1 Plug-in Estimators Are Limits of M-of-N Bootstrap Estimators

In Part I, Section 8.2, it was indicated that for linear OLS there exists a connection between two ways of estimating asymptotic variance: the sandwich estimator for sample size N is the limit of the M-of-N bootstrap as  $M \to \infty$ , where bootstrap is the kind that resamples x-y cases rather than residuals. This connection holds at a general level: all plug-in estimators of standard error are limits of bootstrap in this sense.

The crucial observation of Part I goes through as follows: The M-of-N bootstrap is i.i.d. sampling of M observations from some distribution, hence there must hold a CLT as the resample size grows,  $M \to \infty$ . The distribution being (re)sampled is the empirical distribution  $\hat{P}_N = (1/N) \sum \delta_{(y_i, \vec{x}_i)}$ , where N is fixed but  $M \to \infty$  (causing ever more ties as M grows.) Therefore, the following holds for bootstrap resampling of any well-behaved statistical functional, be it in a regression context or not.

PROPOSITION 8.1. Assume the regression functional  $\theta(\cdot)$  is asymptotically normal for a sufficiently rich class of joint distributions  $P = P_{Y,\bar{X}}$  with acceptable regressor distributions  $P_{\bar{X}}$  as follows:

$$N^{1/2}(\hat{\boldsymbol{\theta}}_N - \boldsymbol{\theta}(\boldsymbol{P})) \xrightarrow{\mathcal{D}} \mathcal{N}(\boldsymbol{0}, \boldsymbol{AV}[\boldsymbol{P}; \boldsymbol{\theta}(\cdot)])$$
$$(N \to \infty).$$

Let  $\hat{\mathbf{P}}_N$  represent a fixed dataset of size N with acceptable regressors. Then a CLT holds for the M-of-N bootstrap resamples as  $M \to \infty$ , with an asymptotic variance obtained by plug-in. Letting  $\boldsymbol{\theta}_M^* = \boldsymbol{\theta}(\mathbf{P}_M^*)$  where  $\mathbf{P}_M^*$  is the empirical distribution of a resample of size M from  $\hat{\mathbf{P}}_N$ , we have

$$M^{1/2}(\boldsymbol{\theta}_{M}^{*} - \hat{\boldsymbol{\theta}}_{N}) \xrightarrow{\mathcal{D}} \mathcal{N}(\boldsymbol{0}, AV[\hat{\boldsymbol{P}}_{N}; \boldsymbol{\theta}(\cdot)])$$
  
 $(M \to \infty, \hat{\boldsymbol{P}}_{N} \text{ fixed}).$ 

The proposition contains its own proof. The following is the specialization to EE functionals where the asymptotic variance has the sandwich form.

COROLLARY 8.1. The plug-in sandwich estimator for an EE functional is the asymptotic variance estimated by the M-of-N bootstrap in the limit  $M \to \infty$  for a fixed sample of size N.

# 8.2 Arguments in Favor of M-of-N Bootstrap over Plug-in Estimators

A natural next question is whether the plug-in/sandwich estimator is to be preferred over M-of-N bootstrap estimators, or whether there is a reason to prefer some form of M-of-N bootstrap. In the latter case, the follow-up question would be how to choose the resample size M. While we do not have any recommendations for choosing a specific M, there exist various arguments in favor of some M-of-N bootstrap over plug-in/sandwich estimation of standard error.

A first argument is that bootstrap is more flexible in that it lends itself to various forms of confidence interval construction that grant higher order accuracy of coverage. See, for example, Efron and Tibshirani (1993) and Hall (1992).

A second argument is related to the first but in a different direction: Bootstrap can be used to diagnose whether the sampling distribution of a particular functional  $\theta(\cdot)$  is anywhere near asymptotic normality for a given sample size N. This can be done by applying normality tests to simulated bootstrap values  $\theta_b^*$  (b = 1, ..., B), or by displaying these values in a normal quantile plot.

A third argument is that there exists theory that shows bootstrap to work for very small M compared to N in some situations where even conventional N-of-N bootstrap does not work. (See Bickel, Götze and van Zwet, 1997, following Politis and Romano, 1994, on subsampling.) It seems therefore unlikely that the limit  $M \to \infty$  for fixed N will yield any form of superiority to bootstrap with finite M.

A fourth argument derives from a result by Buja and Stuetzle (2001, 2016), which states that so-called "M-bagged functionals" have low complexity in a certain sense, the lower the smaller the resample size M is. The limit  $M \to \infty$  is therefore the most complex choice. The connection to the issue of "bootstrap versus plug-in/sandwich estimators" is that M-of-N bootstrap standard errors are simple functions of M-bagged functionals, hence the complexity comparison carries over to standard errors.

It appears that multiple arguments converge on the conclusion that the M-of-N bootstrap is to be preferred over plug-in/sandwich standard errors. (Also recall that both are to be preferred over the residual bootstrap.)

#### 9. SUMMARY AND CONCLUSION

This article completes important aspects of the program set out in Part I. It pursues the idea of model robustness to its conclusion for arbitrary types of regression based on i.i.d. observations. The notion of model robustness coalesces into a model-free theory where all quantities of interest are statistical functionals, called "regression functionals," and models take on the role of heuristics to suggest objective functions whose minima define regression functionals defined on largely arbitrary joint  $(Y, \vec{X})$  distributions. In this final section, we recount the path that makes the definition of well-specification for regression functionals compelling.

To start, an important task of the present article has been to extend the two main findings of Part I from linear OLS to arbitrary types of regression. The findings are that nonlinearity and randomness of the regressors interact ("conspire")

- (1) to cause the target of estimation to depend on the regressor distribution;
- (2) to cause  $N^{-1/2}$  sampling variability to arise that is wholly different from the sampling variability caused by the conditional noise in the response.

It was intuitively clear that these effects would somehow carry over from linear OLS to all types of regression, but it was not clear what would take the place of "nonlinearity," a notion of first-order misspecification peculiar to fitting linear equations and estimating linear slopes. In attempting to generalize Part I, a vexing issue is that one is looking for a framework free of specifics of fitted equations and additive stochastic components of the response. Attempts at directly generalizing the notions of "nonlinearity" and "noise" of Part I lead to dead ends of unsatisfactory extensions that are barely more general than linear OLS. This raises the question to a level of generality in which there is very little air to breathe: the objects that remain are a regression functional  $\theta(\cdot)$  and a joint distribution  $P_{Y,\vec{X}}$ . Given these two objects, what do mis- and well-specification mean? An answer, maybe the answer, is arrived at by casting regression in the most fundamental way possible: Regression is the attempt to describe the conditional response distribution  $P_{Y|\vec{X}}$ . This interpretation sweeps away idiosyncratic structure of special cases. It also suggests taking the joint distribution  $P_{V \vec{X}}$  apart and analyzing the issue of mis- and well-specification in terms of  $P_{Y|\vec{X}}$  and  $P_{\vec{X}}$ , as well as  $\theta(\cdot)$ , the quantities of interest. The solution, finally, to

- establishing a compelling notion of mis- and wellspecification at this level of generality, and
- extending (1) and (2) above to arbitrary types of regression,

is to look no further and use the "conspiracy effect" (1) as the definition: Misspecification means dependence of the regression functional on the regressor distribution. Conversely, well-specification means the regression functional does not depend on the regressor distribution; it is a property of the conditional response distribution alone.

The "conspiracy effect" (2) above is now a corollary of the definition: If the functional is not constant across regressor distributions, it will incur random variability on empirical regressor distributions, and this at the familiar rate  $N^{-1/2}$ .

The link between the proposed definition and conventional ideas of misspecification is as follows: Because most regressions consist of fitting some functional form of the regressors to the response, misspecification of the functional form is equivalent to misspecification of its parameters viewed as regression functionals: depending on where the regressors fall, the misspecified functional form needs adjustment of its

<sup>&</sup>lt;sup>10</sup>The term "bagging" was coined by Breiman (1996) for bootstrap averaging.

parameters to achieve the best approximation over the distribution of the regressors.

Well-specification of regression functionals being an ideal, in reality we always face degrees of misspecification. Acknowledging the universality of misspecification, however, does *not* justify carelessness in practice. It is mandatory to perform diagnostics and, in fact, we proposed a type of diagnostic in Sections 4 and 5 tailored to the present notion of mis/well-specification. The diagnostic consists of checking the dependence of regression functionals on the regressor distribution by systematically perturbing the latter, not by shifting or otherwise moving it, but by reweighting it. Reweighting has the considerable advantage over other forms of perturbation that it applies to all variable types, not just quantitative ones.

While the reality of misspecification imposes a duty to perform diagnostics, there is also an argument to be made to feel less guilty about choosing simpler models over more complex ones. One reason is that the reweighting diagnostic permits localization of models and thereby enables a systematic exploration of local best approximations, always in terms of model parameters interpreted as regression functionals. As shown in Sections 5.1–5.3, this possibility vastly extends the power of models beyond that of a single model fit.

Finally, there is an argument to be made in favor of using statistical inference that is model-robust, and to this end one can use x-y bootstrap estimators or plug-in/sandwich estimators of standard errors. Between the two, one can give arguments in favor of bootstrap over plug-in/sandwich estimators. Importantly, both approaches to inference are in accord with the insight that misspecification forces us to treat regressors as random.

#### **ACKNOWLEDGMENTS**

We are grateful to Gemma Moran and Bikram Karmakar for their help in the generalizations of Section 6.2, and to Hannes Leeb for pointing out the source of the Tukey quote shown before the Introduction.

L. Brown, L. Zhao., and A. Buja were supported in part by NSF Grants DMS-10-07657 and DMS-1310795. E. George was supported in part by NSF Grant DMS-14-06563.

### SUPPLEMENTARY MATERIAL

Supplement to "Models as Approximations II: A Model-Free Theory of Parametric Regression" (DOI: 10.1214/18-STS694SUPP; .pdf). This supplement contains Appendices A-G.

#### **REFERENCES**

- BERK, R., KRIEGLER, B. and YLVISAKER, D. (2008). Counting the homeless in Los Angeles county. In *Probability and Statistics: Essays in Honor of David A. Freedman* (D. Nolan and S. Speed, eds.). *Inst. Math. Stat.* (*IMS*) *Collect.* **2** 127–141. IMS, Beachwood, OH. MR2459951
- BERK, R., BROWN, L., BUJA, A., ZHANG, K. and ZHAO, L. (2013). Valid post-selection inference. *Ann. Statist.* **41** 802–837. MR3099122
- BICKEL, P. J., GÖTZE, F. and VAN ZWET, W. R. (1997). Resampling fewer than *n* observations: Gains, losses, and remedies for losses. *Statist. Sinica* **7** 1–31. MR1441142
- BREIMAN, L. (1996). Bagging predictors. *Mach. Learn.* **24** 123–140.
- BUJA, A. and STUETZLE, W. (2001,2016). Smoothing effects of bagging: Von Mises expansions of bagged statistical functionals. Available at arXiv:1612.02528.
- BUJA, A., BROWN, L., KUCHIBHOTLA, A. K., BERK, R., GEORGE, E. and ZHAO, L. (2019). Supplement to "Models as Approximations II: A Model-Free Theory of Parametric Regression." 10.1214/18-STS694SUPP.
- EFRON, B. and TIBSHIRANI, R. J. (1993). An Introduction to the Bootstrap. Monographs on Statistics and Applied Probability 57. CRC Press, New York. MR1270903
- GELMAN, A. and PARK, D. K. (2009). Splitting a predictor at the upper quarter or third and the lower quarter or third. *Amer. Statist.* **63** 1–8. MR2655696
- HALL, P. (1992). The Bootstrap and Edgeworth Expansion. Springer Series in Statistics. Springer, New York. MR1145237
- HAMPEL, F. R., RONCHETTI, E. M., ROUSSEEUW, P. J. and STA-HEL, W. A. (1986). Robust Statistics: The Approach Based on Influence Functions. Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics. Wiley, New York. MR0829458
- HANSEN, L. P. (1982). Large sample properties of generalized method of moments estimators. *Econometrica* 50 1029–1054. MR0666123
- HASTIE, T. J. and TIBSHIRANI, R. J. (1990). Generalized Additive Models. Monographs on Statistics and Applied Probability 43. CRC Press, London. MR1082147
- HAUSMAN, J. A. (1978). Specification tests in econometrics. Econometrica 46 1251–1271. MR0513692
- HUBER, P. J. (1964). Robust estimation of a location parameter. Ann. Math. Stat. 35 73–101. MR0161415
- HUBER, P. J. (1967). The behavior of maximum likelihood estimates under nonstandard conditions. In *Proc. Fifth Berkeley Sympos. Math. Statist. and Probability (Berkeley, Calif.*, 1965/66), *Vol. I: Statistics* 221–233. Univ. California Press, Berkeley, CA. MR0216620
- KUCHIBHOTLA, A. K., BROWN, L. D. and BUJA, A. (2018). Model-free study of ordinary least squares linear regression. Available at arXiv:1809.10538.
- LIANG, K. Y. and ZEGER, S. L. (1986). Longitudinal data analysis using generalized linear models. *Biometrika* 73 13–22. MR0836430
- PEARL, J. (2009). Causality: Models, Reasoning, and Inference, 2nd ed. Cambridge Univ. Press, Cambridge. MR2548166
- PETERS, J., BÜHLMANN, P. and MEINSHAUSEN, N. (2016). Causal inference by using invariant prediction: Identification

- and confidence intervals. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **78** 947–1012. MR3557186
- POLITIS, D. N. and ROMANO, J. P. (1994). Large sample confidence regions based on subsamples under minimal assumptions. *Ann. Statist.* **22** 2031–2050. MR1329181
- R DEVELOPMENT CORE TEAM (2008). R: A Language and Environment for Statistical Computing. *R Foundation for Statistical Computing*, Vienna, Austria. Available at http://www.R-project.org.
- RIEDER, H. (1994). Robust Asymptotic Statistics. Springer Series in Statistics. Springer, New York. MR1284041
- TUKEY, J. W. (1962). The future of data analysis. *Ann. Math. Stat.* **33** 1–67. MR0133937
- WHITE, H. (1980). Using least squares to approximate unknown regression functions. *Internat. Econom. Rev.* 21 149– 170. MR0572464
- WHITE, H. (1982). Maximum likelihood estimation of misspecified models. *Econometrica* **50** 1–25. MR0640163