

# UC Berkeley

## UC Berkeley Previously Published Works

### Title

Bias Amplification and Bias Unmasking

### Permalink

<https://escholarship.org/uc/item/4bm6q6vf>

### Journal

Political Analysis, 24(3)

### ISSN

1047-1987

### Authors

Middleton, Joel A  
Scott, Marc A  
Diakow, Ronli  
et al.

### Publication Date

2016

### DOI

10.1093/pan/mpw015

Peer reviewed

# Bias Amplification and Bias Unmasking

Joel A. Middleton<sup>1</sup>, Marc A. Scott<sup>2</sup>, Ronli Diakow<sup>3</sup> and Jennifer L. Hill<sup>2</sup>

<sup>1</sup>Department of Political Science, University of California, Berkeley

<sup>2</sup>Humanities and Social Sciences in the Professions, New York University,  
Steinhardt

<sup>3</sup>New York City Department of Education

May, 2016

## ABSTRACT

In the analysis of causal effects in non-experimental studies, conditioning on observable covariates is one way to try to reduce unobserved confounder bias. However, a developing literature has shown that conditioning on certain covariates may increase bias, and the mechanisms underlying this phenomenon have not been fully explored. We add to the literature on bias-increasing covariates by first introducing a way to decompose omitted variable bias into three constituent parts: bias due to an unobserved confounder, bias due to *excluding* observed covariates, and bias due to amplification. This leads to two important findings. While instruments have been the primary focus of the bias amplification literature to date, we identify the fact that the popular approach of adding group fixed-effects can lead to bias amplification as well. This is an important finding because many practitioners think that fixed effects are convenient way to account for any and all group-level confounding and are at worst harmless. The second finding introduces the concept of bias *unmasking* and shows how it can be even more insidious than bias amplification in some cases. After introducing these new results analytically, we use constructed observational placebo studies to illustrate bias amplification and bias unmasking with real data. Finally, we propose a way to add bias decomposition information to graphical displays for sensitivity analysis to help practitioners think through the potential for bias amplification and bias unmasking in actual applications.

---

This research was partially supported by Institute of Education Sciences grants R305D110037 and R305B120017. For replication files see Middleton (2016).

## 1. INTRODUCTION

In the analysis of causal effects in non-experimental studies the key assumption necessary for unbiased estimation is that all confounders (pre-treatment variables that are part of the data generating mechanism for both treatment assignment and outcome) have been measured. In the social science literature this assumption is often referred to as “selection on observables” (Heckman and Robb 1985, 1986), “conditional independence” (Lechner 2001) or “ignorability” (Rubin 1978), and it is well-known that violation of this assumption leads to biased inference. However it is typically implausible to believe that we have measured all confounders, raising the question as to which of the available covariates should be adjusted for (conditioned upon) in practice.

Advice in the extant literature on which variables to condition on are contradictory. One recommendation has been that conditioning on more, rather than fewer, available (pretreatment) covariates is the best way to minimize bias associated with unobserved sources of heterogeneity (Rubin 2002; Rosenbaum 2002). Another recommendation says that those variables that are related to the treatment assignment mechanism should be included in the analysis (D’Agostino, Jr. 1998). Still other advice is to choose covariates based on their relationship to the outcome, rather than to the treatment (Brookhart et al. 2010; Austin et al. 2007; Hill 2007a). There is even some controversy within the controversy: in political science, Clarke (2005) contends that researchers are encouraged to include as many covariates as possible, without regard to the potential for increased bias; relatedly, in most intermediate econometrics texts, the choice between fixed and random effects is framed as a bias versus efficiency tradeoff, with the approach that includes more predictors, fixed effects, described as unbiased (under assumptions that we contend are overly optimistic; see (Greene 2000), for example).

There are, however, two notable classes of covariates that most agree should be *excluded* from the set of conditioning covariates. These are *bias inducers* and *bias amplifiers*. Bias inducers include posttreatment variables such as mediators and colliders (Cole et al. 2010; Pearl 2000;

Schisterman et al. 2009), and a particular group of pretreatment covariates (pretreatment colliders that lead to M-bias or butterfly-bias) (Ding and Miratrix 2014; Sjlander 2009; Pearl 2009). Such bias inducers may not be troublesome in practice, however, either because they can be identified for exclusion, as is sometimes the case for posttreatment variables, or because the bias they induce tends to be small (Ding and Miratrix 2014; Liu et al. 2012; Greenland 2002).

Bias amplifiers have received recent attention as variables that should be excluded from the conditioning set (Pearl 2010; Wooldridge 2009; Bhattacharya and Vogt 2007; Pearl 2011; Myers et al. 2011; Wyss et al. 2014). These covariates cannot induce bias where there is none, but they increase bias by modifying bias that is due to an unobserved confounder. Instruments, variables that only affect the outcome through their impact on the treatment, are the canonical example of a bias amplifying covariate. Conditioning on an instrument can hurt but can never help. On the one hand, this may seem like a trivial concern because it is unclear under what circumstances a researcher would be unaware that a variable was a true instrument for their treatment variable. However, even imperfect instruments can amplify bias (cf. Pearl 2010) and, as we will show below, even noninstruments can amplify bias.

The purpose of this paper is to clarify the relationship between predictor inclusion and bias through a novel decomposition that serves to unify separate strands in the literature, which purport to make recommendations for variable inclusion or exclusion. We then discuss specific conditions in which fixed effects for groups can be detrimental (bias increasing), even though they are largely regarded as addressing a form of omitted variable bias benignly. We show, using a placebo test in a constructed observational study, how this decomposition and an enhanced graphical display devised for sensitivity analysis can better inform variable selection in the context of bias from an unobserved confounder.

The outline of the paper is as follows. In Section 2 we derive our main results, showing how omitted variable bias may be decomposed into several constituent parts: bias attributable to the unobserved confounder, bias due to omitting observed covariates, and bias due to amplification.

These decompositions lead naturally to key insights regarding bias increasing covariates. First, we show that fixed effects for groups can act as pure bias amplifiers. This is novel because fixed effects are not instruments, which have been the focus of the bias amplification literature to date. Moreover, fixed effects are often thought to be a canonical technique for absorbing unmeasured group-level confounding, so demonstrating that they can increase bias in general runs counter to common understanding and practice. Second, we introduce the concept of bias “unmasking”, which explains why even variables that do not amplify bias per se may still lead to net increased bias. In Section 3, we examine two case studies in which the causal effect is known, and where confounding is likely to be present, to estimate and decompose biases into their constituent parts. In one case study amplification is a major contributor to net bias. In the other, the inclusion of covariates leads to a larger *net* bias due to “unmasking” of unobserved confounder bias. These examples reinforce prior advice to avoid inadvertently controlling for instruments when trying to infer causal effects from data where the causal variable was not randomized. They also illustrate why applied researchers might be concerned about arbitrary use of fixed effects for groups in non-experimental studies. In Section 4, we introduce an important addition to graphical displays for sensitivity analysis to help practitioners assess the potential for amplification and unmasking. We summarize our findings in Section 5.

## 2. BIAS DECOMPOSITION

In this section we establish the conditions under which a researcher would want to condition on a set of covariates,  $X$ , in estimating the effect of a treatment,  $Z$ , on an outcome,  $Y$ . To help make ideas concrete we will map this notation to the case study presented in Section 3. Therefore we conceptualize  $Z$  as a pre-recorded get-out-the-vote phone call and  $Y$  as voter turnout.

We start in Section 2.1 with a simple case where  $X$ , a matrix of one or more covariates, is being considered for inclusion or exclusion in its entirety. To fix ideas first consider  $X$  to include only a randomized assignment to receive the get-out-the-vote call, which acts as an instrument (for

more details see Section 3). Then in Section 2.2 we generalize those results to the case where we want to know whether to include *all* of the covariates in  $X$  in the conditioning set given that *some* of them will be included in the conditioning set. We start with the simpler case for clarity of exposition, and the generalization allows us to apply the bias decompositions to more realistic data applications. We then provide conditions specific to models that include fixed effects for groups to demonstrate how they may act as pure bias amplifiers.

### 2.1. *The Case of a Single Set of Conditioning Variables*

Mathematically, describing the magnitude of bias incurred by failure to satisfy the selection on observables assumption requires additional assumptions about the relationships between variables. We derive our results using the linear model as has been done in related work (Ding and Miratrix 2014; Pearl 2010; Clarke 2005, 2009) and which has the advantage of tying this work into more general results regarding omitted variable bias.<sup>1</sup>

To proceed, consider a linear model relating voter turnout,  $Y$  with the causal variable, contact via a get-out-the-vote phone call,  $Z$ ,

$$Y = Z\tau + X\beta^y + U\zeta^y + \epsilon^y. \tag{1}$$

In general  $X$  could represent a matrix of observed covariates; for our motivating example we will include only the randomized treatment assignment, which acts as an instrument.  $U$  is an unobserved confounder; in this setting  $U$  might represent degree of political engagement. As in Carnegie et al. (2014a) and Imbens (2003) we make the simplifying assumption that  $U \perp X$ . We can justify this assumption by conceptualizing  $U$  as the portion of the political engagement that is orthogonal to the randomized treatment assignment. For the sake of clarity, and without loss of

---

<sup>1</sup>We focus on the linear model; we expect that many of the results hold, broadly speaking, for GLM models. However, GLM models come with a host of their own problems with respect to bias. For example, coefficients from unadjusted and covariate adjusted logistic regression models are not comparable (VanderWeele and Arahc 2011; Freedman 2008; Breen et al. 2013), a problem sometimes referred to as “noncollapsibility” (cf. VanderWeele 2015). A discussion of the bias of adjusted and unadjusted GLM estimators is beyond the scope of this paper.

generality, we also assume that variables are mean centered.

To frame this model as a causal model we need to explicitly incorporate potential outcomes that formalize the counterfactual possibilities for the outcome under control and treatment conditions,  $Y(0)$  and  $Y(1)$ , respectively, (Rubin 1974). These represent the voting behavior (voted or did not) if not contacted and if contacted, respectively. Thus we write the causal version of the model as  $E[Y(Z)|Z, X, U] = Z\tau + X\beta^y + U\zeta^y$ . While we do not posit the functional form of the unconditional relationship between  $Z$  and  $U$ , we do not rule out dependence between these variables; in fact if they were independent then  $U$  would not be a confounder. For instance, in our example we expect that those who are contacted have a latent trait, a feature of their personality, that makes them more prone to answering the phone that is correlated with a willingness to vote.

In order to identify a causal effect with this model we need to assume ignorability holds conditional on everything in these models. Formally, this requires  $\{Y(z)\}_{z \in \mathcal{Z}} \perp Z | X, U$ , where  $\mathcal{Z}$  is the set of all possible levels of the treatment. In our context this means that among those who are in the treatment group if we subset on potential voters who are the same in terms of their willingness to vote, then whether or not contact was actually made with them is randomly assigned. For the remainder of the paper we will drop the potential outcome notation with the understanding that the potential outcomes are implicit in our modeling assumptions.

Next we utilize bias expressions to examine the cases where  $X$  is omitted from the conditioning set and when it is included in the conditioning set. We compare the two biases and show how to decompose these into constituent parts. Our derivation relies on a partition of the predictor set, in which letters  $S$  and  $O$ , stand for included (“seen”) and omitted sets, respectively. See Section A of the online supplementary materials for further details.

First suppose one calculates the unadjusted estimate of  $\tau$  in equation (1) by simply regressing  $Y$  on  $Z$  (ignoring known covariates  $X$ ). This translates to substituting  $S = [Z]$  and  $O = [X U]$  (A.2) in the online supplementary materials and yields an expression for the omitted variable bias of the crude estimator:

$$\begin{aligned}\text{Bias} [\widehat{\tau}_{Y|Z}] &= (Z'Z)^{-1} Z'X\beta^y + (Z'Z)^{-1} Z'U\zeta^y \\ &= \chi + v\end{aligned}\tag{2}$$

where  $\chi \equiv (Z'Z)^{-1} Z'X\beta^y$  is the bias due to omitting  $X$  and  $v \equiv (Z'Z)^{-1} Z'U\zeta^y$  is the bias due to omitting  $U$ .<sup>2</sup> These make sense because the bias due to omitting  $X$  is made large when  $X$  is a strong predictor of the outcome (as reflected in the magnitude of  $\beta^y$ ) conditional on the  $Z$ . The bias due to omitting  $U$  is made large when  $U$  is a strong predictor of the outcome conditional on  $Z$  (as reflected in  $\zeta^y$ ).<sup>3</sup> The symbols  $\chi$  and  $v$  are used as a shorthand to signify the constituent parts of the bias as these Greek letters most closely resemble their Latin alphabet counterparts.

To understand the impact of the randomized treatment assignment,  $X$ , on bias, we now calculate the bias when estimating  $\tau$  in a new model that includes  $X$  in the conditioning set,  $S$ . This translates to substituting  $S = [Z \ X]$  and  $O = [U]$  in (A.2) in the online supplementary materials. The bias can be written in partitioned notation as follows:

$$\text{Bias} \begin{bmatrix} \widehat{\tau}_{Y|ZX} \\ \widehat{\beta}^y \end{bmatrix} = \begin{bmatrix} Z'Z & Z'X \\ X'Z & X'X \end{bmatrix}^{-1} \begin{bmatrix} Z' \\ X' \end{bmatrix} U\zeta^y.\tag{3}$$

Using the inverse of the partition matrix and selecting off the element that corresponds to the coefficient on the causal variable  $Z$  (see Section A of the online supplementary materials), write

$$\begin{aligned}\text{Bias} [\widehat{\tau}_{Y|ZX}] &= \left( Z'Z - Z'X (X'X)^{-1} X'Z \right)^{-1} Z'U\zeta^y \\ &= \left( Z'Z - Z'X (X'X)^{-1} X'Z \right)^{-1} (Z'Z)(Z'Z)^{-1} Z'U\zeta^y \\ &= \left( Z'Z - Z'X (X'X)^{-1} X'Z \right)^{-1} (Z'Z) v \\ &= \frac{Z'Z}{(Z'Z - Z'H_X Z)} v = \frac{SST^*}{(SST^* - SSR^*)} v\end{aligned}$$

---

<sup>2</sup>Following (Greene 2000), these expressions assume that  $X$  and  $Z$  and  $U$  are nonstochastic. Generalizing to stochastic regressors involves taking the expectation of these expressions; these are omitted for the sake of clarity.

<sup>3</sup>In our example, this translates understanding whether the randomized treatment assignment is a stronger predictor of whether or not someone will vote than the unobserved willingness to vote characteristic.



$$\begin{aligned}
&= \left( \frac{1}{1 - r_{Z|X}^2} \right) v = v + \left( \frac{r_{Z|X}^2}{1 - r_{Z|X}^2} \right) v \\
&= v + \alpha
\end{aligned} \tag{4}$$

where  $H_X = X(X'X)^{-1}X'$  is the hat matrix associated with the regression; since we assume  $E(Z) = 0$ ,  $SST^* = SST + nE(Z)^2$  is the total sum of squares, and  $SSR^* = SSR + nE(Z)^2$  is the regression sum of squares, so that  $SSR/SST$  is  $r_{Z|X}^2$ , or the coefficient of determination, R-squared, in the regression of  $Z$  on  $X$ . Lastly, we make the identification  $\alpha \equiv \left( \frac{r_{Z|X}^2}{1 - r_{Z|X}^2} \right) v$ . The term  $\left( \frac{1}{1 - r_{Z|X}^2} \right)$  can be referred to as the amplification factor; importantly this term is identified. It reflects the extent to which the covariates, in this case the assignment to receive a get-out-the-vote phone call, predicts actually receiving the treatment, the phone call.

The amplification factor is particularly problematic if  $X$  accounts for a great deal of variation in  $Z$  as noted by Pearl (2010). So in our example, if most of the people assigned to receive a phone call actually received the phone call then this term would be large. The term  $\alpha \equiv \left( \frac{r_{Z|X}^2}{1 - r_{Z|X}^2} \right) v$  gives the change in bias attributable to amplification, call it the *net amplification bias*. See Section B in the online supplementary materials for intuition and context provided by another example.

A careful comparison of (4) and (2), reveals two key insights about adding  $X$  to the conditioning set of covariates. First, note the bias term in (2) associated with omitting  $X$ , namely,  $\chi \equiv (Z'Z)^{-1}Z'X\beta^y$ . Again this term is large when  $X$  is highly predictive of  $Y$ . This term is absent in (4) because  $X$  is adjusted for in this model. Second, the bias due to omitting  $U$  is modified from  $(Z'Z)^{-1}Z'U\zeta^y$  in (2) to become  $(Z'Z - Z'X(X'X)^{-1}X'Z)^{-1}Z'U\zeta^y$  in (4). The difference between these two terms results in the appearance of  $-Z'X(X'X)^{-1}X'Z$  in the denominator, a term which is necessarily less than or equal to zero because it is (-1 times) a quadratic form with positive definite matrix  $(X'X)^{-1}$  (cf. Greene 2000, Sections 2.8 and 2.8.1). This term is large if  $X$ , here the instrument, is highly predictive of the treatment receipt. Therefore, ex-

cept when the term  $-Z'X(X'X)^{-1}X'Z$  is zero (i.e.,  $X$  is not correlated with  $Z$ ), it shrinks the denominator in (4) relative to (2) resulting in *amplification* of the bias due to the unobserved  $U$ .

It is useful at this juncture to emphasize why instrumental variables (Angrist et al. 1996) have been a particular focus of attention in discussing bias amplifiers (Pearl 2010; Wooldridge 2009; Bhattacharya and Vogt 2007). When  $X$  is an instrument, as in our example,  $\beta^y = 0$  by definition. Therefore the bias due to omitting  $X$ ,  $\chi$ , in (2) equals 0 and there can be no benefit due to removing  $\chi$  bias when going from (2) to (4) and the only change in bias is an *increase* due to amplification,  $\alpha$ . In that sense, instruments can be referred to as *pure* amplifiers. When  $X$  is a pure amplifier it is necessarily true that  $v + \alpha$  is larger in magnitude than  $\chi + v$ .

Amplification is only part of the story concerning change in bias when going from an unadjusted to adjusted estimator, however. Whether conditioning on  $X$  increases or decreases the *net* bias depends on the magnitude of (2) relative to the magnitude of (4). Formally, a set of covariates,  $X$ , can be said to be net bias reducing only when

$$\left| (Z'Z)^{-1} Z'U\zeta^y + (Z'Z)^{-1} Z'X\beta^y \right| > \left| \left( Z'Z - Z'X(X'X)^{-1}X'Z \right)^{-1} Z'U\zeta^y \right|.$$

Or using our bias decomposition notation,  $|v + \chi| > |v + \alpha|$ . If  $v$  (the bias due to omitting  $U$ ) and  $\chi$  (the bias due to omitting  $X$ ) have the same sign then this implies that for  $X$  to be bias reducing  $|\chi| > |\alpha|$ . When  $v$  and  $\chi$  have opposite signs the requirement is  $|\chi| > |2v + \alpha|$ .

Clearly, conditioning on  $X$  can be net bias increasing in cases where the bias due to amplification,  $\alpha$ , is relatively large. This can happen, for instance, when the randomized treatment assignment is a strong predictor of actual contact via a get-out-the-vote message. However, conditioning on  $X$  can be net bias increasing even when  $r_{Z|X}^2 = 0$  (and, hence,  $\alpha=0$ ) if the bias due to omitting  $U$ ,  $v$ , and the bias due to omitting  $X$ ,  $\chi$ , have opposite signs and  $|\chi| < 2|v|$ . In that case, because  $\chi$  has an opposite sign to  $v$  but similar magnitude, it can be said to be *masking* (or canceling)  $v$  in (2). In that case the  $\chi$  is a “good” bias because it cancels with  $v$ , rendering the net

bias of the unadjusted estimator closer to zero than that of the adjusted estimator.

That bias due to omitting a known covariate  $X$  can be “good” bias (because its exclusion masks bias due to the unobserved confounder) is troubling because it implies that even when  $X$  is known to be predictive of  $Y$ , including it in the conditioning set of covariates may increase overall bias. To know whether removing  $\chi$  bias improves net bias or not, one must know something about  $v$  which is not identified. In light of this observation, it is clear that none of the existing recommendations for practice provide complete guidance on whether to condition on a covariate, or set of covariates.

## 2.2. The Case of Two Sets of Conditioning Variables

In this section we generalize the above results to the case in which we want to decide whether to include *all* of the covariates in  $X$  in the conditioning set given that *some* of them will be included in the conditioning set. Notationally, first partition the matrix of covariates such that  $X = [X_1 X_2]$ . Now assume that  $X_2$  will certainly be in the conditioning set and the question is whether to also include  $X_1$  in the conditioning set. This is a very common situation. The results for bias amplification are analogous to the prior case – we simply condition on  $X_2$  throughout the derivation. We omit unnecessary detail.

The model can now be written,

$$Y = Z\tau + X_1\beta_1^y + X_2\beta_2^y + U\zeta^y + \epsilon^y \quad (5)$$

where the  $U$  is independent of both  $X_1$  and  $X_2$ . Omitting  $X_1$  from the conditioning set leads to

$$\text{Bias} [\widehat{\tau}_{Y|ZX_2}] = \chi^* + v^* \quad (6)$$

where  $v^* \equiv \left( \frac{1}{1-r_{Z|X_2}^2} \right) v$ , likewise  $\chi^* \equiv \left( \frac{1}{1-r_{Z|X_2}^2} \right) \chi_1$  and  $r_{Z|X_2}^2$  is the R-squared in the regression of  $Z$  on  $X_2$ .

Including  $X_1$  in the conditioning set leads to

$$\text{Bias} [\widehat{\tau}_{Y|ZX_1X_2}] = v^* + \alpha^*. \quad (7)$$

Here the net amplification bias,  $\alpha^* \equiv \left( \frac{r_{Z|X_1X_2}^2 - r_{Z|X_2}^2}{1 - r_{Z|X_1X_2}^2} \right) v^*$ , is defined only slightly differently from  $\alpha$  above and the amplification factor can be written  $\left( \frac{1 - r_{Z|X_2}^2}{1 - r_{Z|X_1X_2}^2} \right)$ . Here  $r_{Z|X_1X_2}^2$  is the R-squared in the regression of  $Z$  on  $X_1$  and  $X_2$ .

So, conditioning on  $X$  is bias reducing when  $|v^* + \chi^*| > |v^* + \alpha^*|$ . When  $v^*$  and  $\chi^*$  have the same sign the requirement is that  $|\chi^*| > |\alpha^*|$ . When  $v^*$  and  $\chi^*$  have different signs the requirement is that  $|\chi^*| > |2v^* + \alpha^*|$ .

### 2.3. The Important Special Case of Fixed Effects

As mentioned above, pure bias amplifiers such as instruments can be particularly problematic because there cannot be any benefit to removing  $\chi$  from the bias equation since  $\chi = 0$ . In this section we derive the conditions in which fixed effects can be pure amplifiers – amplifying bias but providing no net improvement in bias due to removing  $\chi$  bias.

To consider fixed effects under the rubric presented above simply imagine  $X$  as a matrix of indicator variables representing groups. We ignore additional covariates in this development, but the results hold for that case as well. Starting from this point of view, the term  $\chi \equiv (Z'Z)^{-1} Z'X\beta^y$  in (2) can be written  $\chi \equiv \sum_{k=1}^K (Z'Z)^{-1} Z'X_k\beta^{yk}$  where  $X_k$  is the column vector from  $X$  associated with the  $k^{\text{th}}$  dummy variable, and  $\beta^{yk}$  is the corresponding coefficient for the  $k^{\text{th}}$  group.

Now consider the case where fixed effects are pure amplifiers – when the term  $\chi \equiv \sum_{k=1}^K (Z'Z)^{-1} Z'X_k\beta^{yk} = 0$ . Trivially, this term can be zero if the terms  $\beta^{yk}$  are all zero, i.e., if the fixed effects are instruments, but it can also be zero because the positive and negative terms sum to zero. When might those positive and negative terms net out to zero? To develop an intuition, consider a model for the treatment,  $Z$ ,

$$Z = \sum_{k=1}^K X_k\beta^{zk} + U\zeta^z + \epsilon^z \quad (8)$$

where  $U$  is defined as above. Now make the assumption, for the sake of simplifying the exposition, that  $Z$  has unit variance (in addition to having mean zero) and that the size of the  $K$  groups (associated with the fixed effects) are equal and thus  $E[X_k] = 1/K$ . Then we might write for the  $k^{\text{th}}$  dummy variable

$$\begin{aligned}
(Z'Z)^{-1} Z'X_k &= \frac{1}{n} Z'X_k \\
&= \text{Cov}(X_k, Z) = E[X_k Z] - E[X_k]E[Z] \\
&= E[(X_k)(\sum_j X_j \beta^{zj} + U\zeta^z + \epsilon^z)] - E[X_k]E[\sum_j X_j \beta^{zj} + U\zeta^z + \epsilon^z] \\
&= E[(X_k) \sum_j X_j \beta^{zj}] + E[(X_k)(U\zeta^z + \epsilon^z)] - E[X_k]E[\sum_j X_j \beta^{zj}] \\
&= E[(X_k)(X_k)\beta^{zk}] - \frac{1}{K} \sum_j E[X_j \beta^{zj}] \\
&= E[(X_k)]\beta^{zk} - \frac{1}{K} \sum_j \frac{1}{K} \beta^{zj} \\
&= \frac{1}{K} \beta^{zk} - \frac{1}{K^2} \sum_{j=1}^K \beta^{zj} \tag{9}
\end{aligned}$$

We rely on  $X_k X_k = X_k$  and  $X_k X_j = 0_n$ , for  $j \neq k$ , as well as  $E[X_k U] = E[X_k \epsilon^z] = 0$ , above.

Utilizing the above, we have,

$$\begin{aligned}
\chi &\equiv \sum_{k=1}^K (Z'Z)^{-1} Z'X_k \beta^{yk} = \sum_{k=1}^K \text{Cov}(X_k, Z) \beta^{yk} \\
&= \sum_{k=1}^K \left( \frac{1}{K} \beta^{zk} - \frac{1}{K^2} \sum_{j=1}^K \beta^{zj} \right) \beta^{yk} \\
&= \frac{1}{K} \sum_{k=1}^K \beta^{zk} \beta^{yk} - \left( \frac{1}{K} \sum_{k=1}^K \beta^{yk} \right) \left( \frac{1}{K} \sum_{j=1}^K \beta^{zj} \right) \\
&= E_k [\beta^{zk} \beta^{yk}] - E_k [\beta^{zj}] E_k [\beta^{yk}] \\
&= \text{Cov}_k (\beta^{zk}, \beta^{yk}) \tag{10}
\end{aligned}$$

where we use the notation  $\text{Cov}_k$  to denote that covariance is to be taken *across* the  $K$  groups. Likewise,

$E_k$  is expectation across the  $K$  groups.

The derivation provides us with conditions in which fixed effects will be pure bias amplifiers. When  $Cov_k(\beta^{zk}, \beta^{yk}) = 0$ , pure amplification obtains, but clearly there should be concern when this condition is approximately met as well. One way to interpret this situation is that the group level structure in  $Y$  does not covary (or covaries extremely weakly) with the group level structure in  $Z$ . In other words, the average of the products of the group effects is zero.

If the expression in the last line of (10) can be estimated, then the fixed effects can be avoided when they are pure bias amplifiers (or close to it). Unfortunately, the term cannot be estimated unbiasedly, or even meaningfully bounded. To see this examine  $\text{Bias}[\widehat{\beta}^y]$  in (A.5) in the online supplementary materials. Instead, the usual regression estimator  $\widehat{\beta}^{yk}$  converges in probability to

$$\beta^{yk} - \beta^{zk} \left( Z'Z - Z'X [X'X]^{-1} X'Z \right)^{-1} Z'U\zeta^y. \quad (11)$$

Asymptotically then, a quantity that one might estimate is

$$\begin{aligned} \widehat{Cov}_k(\beta^{zk}, \beta^{yk}) &= Cov_k(\beta^{zk}, \widehat{\beta}^{yk}) \\ &= Cov_k\left(\beta^{zk}, \beta^{yk} - \beta^{zk} \left( Z'Z - Z'X [X'X]^{-1} X'Z \right)^{-1} Z'U\zeta^y\right) \\ &= Cov_k(\beta^{zk}, \beta^{yk}) - V_k(\beta^{zk}) \left( Z'Z - Z'X [X'X]^{-1} X'Z \right)^{-1} Z'U\zeta^y \\ &= Cov_k(\beta^{zk}, \beta^{yk}) - V_k(\beta^{zk}) (v + \alpha) \end{aligned} \quad (12)$$

where  $V_k(\beta^{zk})$  represents the variance of the values of  $\beta^{zk}$ . Because  $(v + \alpha)$  may take on a potentially wide range of values (12) is not a useful estimator of (10).<sup>4</sup>

It is certainly plausible that the relationship between group effects for the treatment and outcome are unrelated, net of controls. Since this cannot be determined empirically, a researcher must take seriously the potential for this to occur. Moreover, even if fixed effects are not pure bias am-

---

<sup>4</sup>That said, in a sensitivity analysis framework then, estimates for  $Cov_k(\beta^{zk}, \beta^{yk})$  might be computed for posited values of the bias term in (4).

plifiers, they may be bias increasing due to bias unmasking. Thus, it is quite possible that including fixed effects for groups will lead to increased absolute bias. The common rationale for including fixed effects is that they “cannot hurt; often help” is not supported by this decomposition of bias analysis.<sup>5</sup> Moreover, emphasizing the choice between fixed and random effects as a bias versus efficiency tradeoff subverts an important consideration, which is bias increasing under the inclusion of fixed effects in the presence of an unobserved confounder. While adjusting for group-level confounding, the fixed effects approach potentially introduces the two types of bias characterized in our decomposition.

### 3. CASE STUDIES

#### 3.1. *Case Study I: The effect of a Get-Out-the-Vote intervention*

In this subsection we repurpose the data from a study of the effect of prerecorded get-out-the-vote phone calls on voter turnout (Shaw et al. 2012) to illustrate the phenomenon of bias amplification. While the original study was a randomized experiment, we use the data to create a constructed observational study.

In the original experiment units were assigned to a condition that received a prerecorded telephone message encouraging them to vote or to a “no message” condition. In 1,597 precincts randomization was at the precinct level. In another 5,838 precincts, households were randomly assigned to treatment or control within precinct. We use the combined data file of 463,489 subjects.

Of interest in this study was the effect of contact,  $Z$ , on voter turnout,  $Y$ . However, individuals who were actually contacted may be different who were not contacted in ways that make them more likely to vote – for instance they were less likely to have died or moved. Therefore, naively regressing turnout on contact is likely to violate the selection on observables assumption and thus

---

<sup>5</sup>In work that supports this contention, Clarke (2009) concludes that despite some awareness of the potential for bias, the common practice in Political Science is to include as many predictors as possible. The author does not specifically name fixed effects in the admonitions and instead uses simulation studies to characterize absolute bias differences under inclusion and exclusion of single predictors.

yield a biased estimate of the effect of contact. Instrumental variables regression, using treatment assignment as the instrument, is the typical remedy in a situation like this with randomized assignment to treatment and strong evidence that the only pathway between the randomized assignment and the outcome is through treatment receipt. However, we are interested in illustrating bias so we do not use instrumental variables. and instead we make comparisons between those contacted and those not contacted. Moreover, we construct a placebo test, using turnout in *prior* elections as outcome measures. Since we know that contact in 2006 cannot affect the turnout in prior elections, the true treatment effect must be zero. Estimates that deviate from zero thus reveal the bias inherent in the estimator.<sup>6</sup>

Within the context of our constructed observational placebo study we can test whether two types of variables act as bias amplifiers when included as covariates in the specified model. In Section 3.2 we consider treatment assignment (an instrument for contact) as a bias amplifier. In Section 3.3 we consider fixed effects for precinct as bias amplifiers. In both sections, estimates from models that include the potential bias amplifier are compared to a simple regression of turnout on contact to see which yields an estimate closer to the true parameter value of zero. If a model with the potential bias amplifier yields an estimate that is further from zero, then this is evidence that the potential bias amplifier caused a net increase in bias. Furthermore, because the causal parameter is known to be zero, the constituent components of bias –  $\chi$ ,  $v$  and  $\alpha$  – are also identified. So for each outcome we can see if bias amplification is the cause of the net increase in bias.

### 3.2. *Instrument as Bias Amplifier: Analysis and Results*

Table 1, panel A shows the results for the analysis of the effect of conditioning on an instrument, randomized treatment assignment, on bias.<sup>7</sup>

---

<sup>6</sup>Another option would have been to use 2006 election turnout as the outcome and compare our observational estimates to the experimental benchmark created by the instrumental variables estimate. The downside of this approach is that this benchmark is itself noisy making it more difficult to precisely partition the bias. We prefer using the sharp 0 of our placebo tests as a comparison.

<sup>7</sup>For the data and replication files for all tables and figures herein, see Middleton (2016).



To describe the model specification we refer back to model (1).  $Y$  is an  $(n \times 1)$  vector of voter turnout indicators,  $Z$  is an  $(n \times 1)$  vector of indicators for contact,  $X$  is an  $(n \times 1)$  vector of indicators of treatment assignment.  $U$  is the omitted confounder, assumed to have unit variance. Each row of the table conducts the analysis for a different election. The column labeled ‘OLS’ presents the estimated coefficient on  $Z$  when regressing  $Y$  on  $Z$  only. The column labeled ‘Inst.’ presents the estimated coefficient on  $Z$  when regressing  $Y$  on  $X$  and  $Z$ . The column ‘Diff’ presents the difference between the two estimates along with a bootstrapped standard error. In the columns labeled  $v$ ,  $\alpha$ , and  $\chi$  the observed bias is decomposed into constituent parts.

For the general election 2004, the OLS estimate exhibits a bias of 0.138 while the model controlling for treatment is much more biased at 0.478.<sup>8,9</sup> The bias increase of 0.339 (an increase of 244%) is entirely due to bias amplification,  $\alpha$ . That there is essentially no contribution to the bias through  $\chi$  is expected given that instruments are known to be pure amplifiers. Not surprisingly then, the unadjusted estimator is better than one that adjusts for an instrument.

Results in Table 1, Panel A, from other election years show substantively similar results. As in the case of the General Election 2004, adding the treatment indicator to the conditioning set of covariates leads to increased bias. The increase in bias is attributable to bias amplification.

Results in Table 1, Panel B, repeat this analysis for models that include additional covariates (turnout in prior elections) in the conditioning set. To describe this model specification we refer back to model (5).  $Y$  and  $Z$  are defined as in Panel A.  $X_1$  is the instrument while  $X_2$  is an  $(n \times k)$  matrix of indicators for turnout in  $k$  prior elections. For the General 2004 election outcome,  $X_2$  included General 2002 turnout, General 2000 turnout, Primary 2004 turnout, Primary 2002 turnout and Primary 2000 turnout. For the Primary 2004 election outcome,  $X_2$  included General 2002 turnout, General 2000 turnout, Primary 2002 turnout and Primary 2000 turnout. The column labeled ‘OLS’ presents the estimated coefficient on  $Z$  when regressing  $Y$  on  $Z$  and  $X_2$ . In the next

---

<sup>8</sup>The standard errors are so small as to suggest that the bias is measured with great precision.

<sup>9</sup>This is a tremendous amount of bias when one considers that the outcome is a binary, 0-1, outcome.

column, labeled ‘Inst.’, is the estimated coefficient on  $Z$  when regressing  $Y$  on  $Z$ ,  $X_2$  and also  $X_1$ . The remaining columns give the difference between the two estimates and the bias decomposition.

Overall the biases are smaller in Panel B. For example, in Panel A, for General 2004 the bias due to omitting  $U$ ,  $v$ , is estimated to be 0.140, or 14 percentage points. In Panel B, in contrast, the bias due to omitting  $U$ ,  $v^*$  is estimated to be 0.021 or two percentage points. This seems to suggest that the prior vote history variables are themselves covariates that reduce bias in this example. However, we note that 2.1 percentage points is still a substantively large bias. Moreover, the ratio  $\frac{\alpha^*}{v^*}$  is similar to  $\frac{\alpha}{v}$  above; bias due to amplification,  $\alpha^*$ , is over 250% the size of the bias due to omitting  $U$ ,  $v^*$ .

A. No covariates	OLS (SE)	Inst. (SE)	Diff (SE)	$v$	$\alpha$	$\chi$
General 2004	0.138 (0.004)	0.478 (0.004)	0.339 (0.003)	0.140	0.337	-0.002
General 2002	0.135 (0.005)	0.451 (0.004)	0.315 (0.004)	0.132	0.318	0.003
General 2000	0.131 (0.005)	0.451 (0.004)	0.32 (0.004)	0.132	0.318	-0.001
Primary 2004	0.094 (0.005)	0.285 (0.004)	0.192 (0.004)	0.084	0.202	0.01
Primary 2002	0.093 (0.004)	0.288 (0.003)	0.195 (0.004)	0.085	0.204	0.009
Primary 2000	0.113 (0.004)	0.37 (0.004)	0.257 (0.003)	0.109	0.261	0.004
B. With covariates	OLS (SE)	Inst. (SE)	Diff (SE)	$v^*$	$\alpha^*$	$\chi^*$
General 2004	0.017 (0.002)	0.077 (0.002)	0.06 (0.002)	0.021	0.056	-0.005
Primary 2004	0.025 (0.005)	0.062 (0.003)	0.037 (0.004)	0.017	0.045	0.008

Table 1: GOTV Example with Instrument as Potential Bias Amplifier. Results are displayed for estimates of the effect of the get-out-the-vote intervention on a number of pre-treatment outcomes thus creating placebo tests. Column 1 reveals that linear regression results suffer from bias due to selection on unobservables. Column 2 displays results from an extension of this analysis that could exacerbate the selection bias by including the indicator for the initial randomization, which in this case acts as an instrument. The third column presents the raw difference between columns 1 and 2. The final three columns decompose the bias into the constituent parts (see Sections 2 and 2.2).

### 3.3. *Fixed Effects as Bias Amplifiers: Analysis and Results*

Next, consider the implications for bias when adding fixed effects for precinct to the model specification. Table 2 presents these results. In Panel A, referring back to model (1),  $Y$  and  $Z$  are defined as the turnout indicators and contact indicators, as above, while  $X$  is now an  $(n \times K)$  matrix of dummy variable indicators for the  $K$  precincts.

Examining the results for 2004 election turnout, the fixed effects model is much more biased than the model without fixed effects; when regressing  $Y$  on  $Z$  only the estimate is 0.138 compared to 0.272 when regressing  $Y$  on  $Z$  and  $X$ . The net increase in bias is 97%. Here again, the major factor in the bias difference is bias amplification,  $\alpha$ . Fixed effects are essentially pure bias amplifiers as evidenced by the fact there is a virtually no bias associated with omitting them ( $\chi = 0.001$ ). We reiterate: group fixed effects have the potential to increase absolute bias by way of pure bias amplification.

Results in Table 2, Panel A from other election years show substantively similar results for adding fixed effects to the model specification. Adding the fixed effects to the conditioning set of covariates leads to increased bias due to bias amplification.

Panel B of Table 2 presents the analysis where additional covariates are included in the specification. Again, refer back to model (5) to see the model specification.  $Y$  and  $Z$  are specified as in Panel A. Here  $X_1$  is a matrix of dummy variables for precinct and  $X_2$  includes the prior election turnout indicators as in the Table 1, Panel B.

Results again show that fixed effects have amplified bias. While the amount of bias starts off lower for these models, the amplification factor is about the same, roughly doubling the bias of the estimate.

### 3.4. *Case Study II: The Effect of Selecting a Disadvantaged Village Council President*

In the previous case study we demonstrated a situation where bias amplification resulted from either adding an instrument or adding fixed effects to the conditioning set of covariates. This

A. No covariates	OLS (SE)	FE (SE)	Diff (SE)	$v$	$\alpha$	$\chi$
General 2004	0.138 (0.004)	0.272 (0.004)	0.134 (0.004)	0.137	0.135	0.001
General 2002	0.135 (0.005)	0.257 (0.004)	0.123 (0.005)	0.129	0.128	0.006
General 2000	0.131 (0.005)	0.257 (0.004)	0.126 (0.004)	0.129	0.128	0.002
Primary 2004	0.094 (0.005)	0.172 (0.003)	0.077 (0.005)	0.087	0.086	0.007
Primary 2002	0.093 (0.004)	0.174 (0.003)	0.081 (0.004)	0.087	0.086	0.006
Primary 2000	0.113 (0.004)	0.217 (0.004)	0.104 (0.003)	0.109	0.108	0.004
B. With covariates	OLS (SE)	FE (SE)	Diff (SE)	$v^*$	$\alpha^*$	$\chi^*$
General 2004	0.017 (0.002)	0.037 (0.001)	0.02 (0.002)	0.018	0.019	-0.001
Primary 2004	0.025 (0.005)	0.037 (0.002)	0.011 (0.005)	0.018	0.019	0.007

Table 2: GOTV Example with Set of Fixed Effects as Potential Bias Amplifier. The columns are otherwise similar to those in Table 1.

amplification occurred whether or not there were additional conditioning covariates specified in the model. In this section we consider a case study that repurposes data from another study (Dunning et al. 2013) to provide an example where bias amplification *per se* is not a major concern but where fixed effects nonetheless lead to a large increase in net bias as the bias due to  $v$  is “unmasked.”

The original paper examined the effect of having a village council president from a disadvantaged group (scheduled cast or scheduled tribe) on programmatic spending in India. In certain locations in India, council seats are reserved for disadvantaged groups on a rotating basis. Villages were assigned to have a reserved seat by first creating a list of councils within each district sorted by size of the population of the target disadvantaged group. Then, councils above a certain cutoff had their presidencies reserved for a disadvantaged group. In subsequent elections, the list was rotated so that a different set of villages had reserved seats. The original study capitalized on the list rotation scheme to conduct a quasi-experimental study that compared cities just above the cutoff to cities just below.

### 3.5. Analysis and Results

We reuse these data in a way not intended by the original study in order to induce confounding and study the resulting bias. We induce confounding by using the entire data set, not just the quasi-experimental pairs. Including data from all villages introduces confounding because villages higher on the list are not valid counterfactual cases for those further down given that they were sorted by the population of the disadvantaged groups.

Next, because outcome data exist for a time period *before* the assignment of the treatment, we were able once again to conduct a placebo test, whereby the effect of the treatment on the outcomes in a prior time period could be analyzed.<sup>10</sup> As above, since the treatment cannot affect outcomes in the past, the true value of the parameter is known to be zero. Estimates from the data can be compared to the true benchmark of zero and deviations from zero can be considered evidence of bias.

Our analysis compared models with and without fixed effects for district for each of a number of outcome measures. The estimates from the two models can be compared to see which is closer to the true parameter value of zero. If the fixed effects model is further from zero, then this is evidence that fixed effects cause a net increase in bias.

The outcome measures reflect seven government programs. Table 3 provides the names of the programs. Outcomes are measured in thousands of rupies for the first five outcomes and in number of latrines for the last two.

Table 3 presents our results. Again referring back to the model in (1),  $Y$  is an  $(n \times 1)$  vector of expenditures (or number of latrines for the last two outcomes),  $Z$  is an  $(n \times 1)$  vector of indicators for treatment assignment (reserved council presidency) and  $X$  is an  $(n \times k)$  matrix of indicators of district (taluk). In Panel A, Table 3 the column labeled ‘OLS’ gives the estimated coefficient on  $Z$  when regressing  $Y$  on  $Z$  only. The column labeled ‘FE’ gives the estimated coefficient on  $Z$  when

---

<sup>10</sup>We examined the “effect” of seats reserved in the 2007 election on outcomes from 2006. We also limited the data set to those villages that did not have a reserved presidency in 2005-2006 election years.

regressing  $Y$  on  $Z$  and  $X$ .

Results in Panel A of Table 3 show that in three of seven cases (Ashraya, Latrines and Community Latrines), fixed effects appear to be moving estimates in the direction of zero. In two other cases (IAY Scheme and Ambedkar), the estimates are moving away from zero but only slightly so. In the case of MGNREGA, the result is a tossup with bias moving from 0.5 to -0.5 with the addition of the fixed effects.

In the case of the Water Infrastructure, however, the fixed effects estimate is much further from zero compared to the unadjusted estimate (0.3 compared to -10.2). The bootstrapped standard error suggests that this is a statistically significant difference.

In the last three columns of the table the observed biases are decomposed into constituent parts: the bias due to the unobserved confounder ( $v$ ), the bias due to amplification when controlling for  $X$  ( $\alpha$ ) and the bias due to omitting  $X$  from the conditioning set ( $\chi$ ). In decomposing this bias in the case of Water Infrastructure we can examine what is happening; the bias due to omitting  $X$  is roughly of the same magnitude as the bias due to omitting  $U$  but they have opposite signs. When neither  $X$  nor  $U$  are controlled for in the model, the two biases cancel. In this sense,  $\chi$  can be said to be “good” bias which is masking  $v$  in the unadjusted model. Bias amplification ( $\alpha$ ) plays a role, albeit a smaller one, accounting for about 9% (-0.6/-6.7) of the move away from zero when going from the unadjusted to the adjusted estimator.

Panel B of Table 3 shows the results for the model that includes several additional covariates in the conditioning set. Referring to model (5),  $Y$  and  $Z$  are defined as in Panel A.  $X_1$  is a matrix of dummy variables for district.  $X_2$  is an  $(n \times 7)$  matrix of covariates including village expenditures for the year, village population, population of scheduled caste members, population of scheduled tribe members, size of the literate population, and size of the working population.

The values in the OLS column of Panel B, Table 3 are the estimated coefficient on  $Z$  when regressing  $Y$  on  $Z$  and  $X_2$ . The values in the FE column are the estimated coefficient on  $Z$  when regressing  $Y$  on  $Z$ ,  $X_1$  and  $X_2$ .

Results confirm the main finding for Water Infrastructure. Including fixed effects in the model unmask the bias due to  $U, v^*$ , making the total bias worse than when fixed effects are not included in the model.

Interestingly, results in Panel B also show that when controlling for these other covariates,  $X_2$ , fixed effects actually increase bias for 6 out of 7 outcomes compared to 2 out of 7 in Panel A. That additional covariates,  $X_2$ , can alter whether fixed effects help or hurt greatly complicates the question of what to include in the conditioning set for practitioners.

A. No Covariates	OLS (SE)	FE (SE)	Diff (SE)	$v$	$\alpha$	$\chi$
Ashraya	3.4 (2.6)	0.4 (1.9)	-3.0 (2.8)	0.4	0.1	3.0
IAY	33.4 (12.4)	-34.0 (8.9)	-67.5 (9.5)	-30.0	-4.1	63.4
Ambedkar	-0.6 (0.6)	-0.8 (0.8)	-0.2 (0.2)	-0.7	-0.1	0.1
MGNREGA	0.5 (3.5)	-0.5 (2.7)	-1.0 (1.0)	-0.5	0.0	1.0
Water Infrastructure	0.3 (3.9)	-10.2 (4.0)	-10.5 (2.6)	-9.0	-1.2	9.3
Latrines	-12.3 (5.8)	-5.7 (5.3)	6.7 (4.0)	-5.0	-0.6	-7.3
Community Latrines	0.2 (0.2)	0.1 (0.1)	-0.2 (0.1)	0.0	0.0	0.2
B. With Covariates	OLS (SE)	FE (SE)	Diff (SE)	$v^*$	$\alpha^*$	$\chi^*$
Ashraya	1.1 (1.7)	1.5 (2.2)	0.4 (1.8)	1.3	0.2	-0.3
IAY	-5.9 (7.5)	-10.0 (6.8)	-4.0 (5.4)	-8.9	-1.1	2.9
Ambedkar	-0.8 (0.8)	-1.6 (1.5)	-0.8 (0.8)	-1.5	-0.1	0.6
MGNREGA	-0.3 (3.6)	-0.6 (3.2)	-0.3 (0.7)	-0.6	0.0	0.3
Water Infrastructure	1.8 (3.5)	-4.1 (2.8)	-5.9 (1.8)	-3.6	-0.4	5.4
Latrines	-5.1 (5.0)	-1.8 (5.5)	3.4 (3.3)	-1.6	-0.2	-3.6
Community Latrines	0.1 (0.2)	0.2 (0.2)	0.1 (0.1)	0.2	0.0	-0.1

Table 3: Village Council Presidency Example with Fixed Effects

#### 4. A SENSITIVITY ANALYSIS FRAMEWORK

For the case studies that we have examined, we have identified situations where adding an instrument or fixed effects to a set of conditioning variables increases bias. We can see this increase in bias because we have constructed these studies as placebo tests, whereby the true parameter value is known to be zero because the outcomes occurred before the treatment. However, our case studies provide little consolation to practitioners who do not know the true value of the parameter. The question we consider in this section is whether sensitivity analysis could be used to alert a practitioner to the potential for increases in bias that we have demonstrated.

Sensitivity analysis has been proposed as a way to visualize the potential for an unobserved confounder to bias results of an analysis (c.f. Imbens 2003; Rosenbaum and Rubin 1983; Clarke 2005, 2009). These approaches posit the attributes of an unobserved confounder,  $U$ , (usually its association with treatment  $Y$  and outcome  $Z$ ) that would be sufficient (in addition to observed confounders) to satisfy the selection on observable assumption. Then they calculate the amount of bias induced by failing to include  $U$  in the conditional set. Typically a full sensitivity analysis repeats this exercise across a range of possible attributes for  $U$  and the results can be visually displayed. If the estimated outcome changes very little except in the face of very extreme confounding by  $U$ , the results are said to be insensitive to omitted confounder bias. Similar attributes of observed covariates (for example their associations with treatment and outcome) can be used as benchmarks to help understand the range of plausible attribute values for “typical” covariates in that setting.

We modify for our purposes a new sensitivity analysis package available in R: `treatSens` (Carnegie et al. 2014b,a). The tool takes a dual-parameter approach similar to that of Imbens (2003).

For a given combination of values of the sensitivity parameters (the coefficients on  $U$  in the  $Y$  and  $Z$  models:  $\zeta^z$  and  $\zeta^y$ , respectively) an estimate of the treatment effect,  $\tau$ , can be generated by first drawing candidate values of  $U$ , denoted  $\hat{U}$ , from the distribution implied by the sensitivity



parameters and then estimating the parameters of the model regressing  $Y$  on  $Z$ ,  $X_1$ ,  $X_2$  and  $\dot{U}$  using OLS. Call the estimate of the coefficient on  $Z$ ,  $\hat{\tau}(\zeta^z, \zeta^y)$ . An average of this parameter estimate is taken across 20 draws of  $\dot{U}$  to reduce the uncertainty associated with the random draws from the distribution of  $U$ .

The algorithm proceeds by considering a range of possible values of  $\zeta^z$  and  $\zeta^y$  in a grid. Values of  $\hat{\tau}(\zeta^z, \zeta^y)$  can be computed for each cell in the grid. The values in the grid can be represented on a plot with axes  $\zeta^z$  and  $\zeta^y$  and contours drawn representing constant values of  $\hat{\tau}(\zeta^z, \zeta^y)$ .

Note that for a given contour all of the bias terms in our equations are identified. Therefore, we modify the sensitivity analysis currently available in the `treatSens` package to label each contour in the grid with the values of  $v^*$ ,  $\alpha^*$ , and  $\chi^*$ , which are themselves implied by the values of the sensitivity parameters,  $\zeta^z$  and  $\zeta^y$ , in addition to  $\hat{\tau}(\zeta^z, \zeta^y)$ . Additionally, we place a contour demarcating the area in which given fixed group effects would increase bias, rather than decreasing it. This modification allows the user to identify whether the areas of the parameter space where bias increases due to this addition represent manifestations of the unobserved variable  $U$  that are plausible.

To calibrate the strength of the sensitivity parameters, we follow Imbens (2003) in plotting the coefficient estimates on the (standardized) observed covariates,  $X_2$  in the framework of Section 2.2, in the data.

In the next section, we present a sensitivity analysis plot for the voter turnout experiment. Section C in the online supplementary materials presents a figure for village expenditures in India.

#### 4.1. *Sensitivity Analysis of GOTV Outcomes*

The sensitivity analysis for the GOTV outcomes, in Figure 1, examines the potential for fixed effects to bias the estimates for the effect of contact on General 2004 turnout, presented in Panel B of Table 2.

In interpreting Figure 1, consider the point (0.1, 0.05). It falls approximately on the line labeled

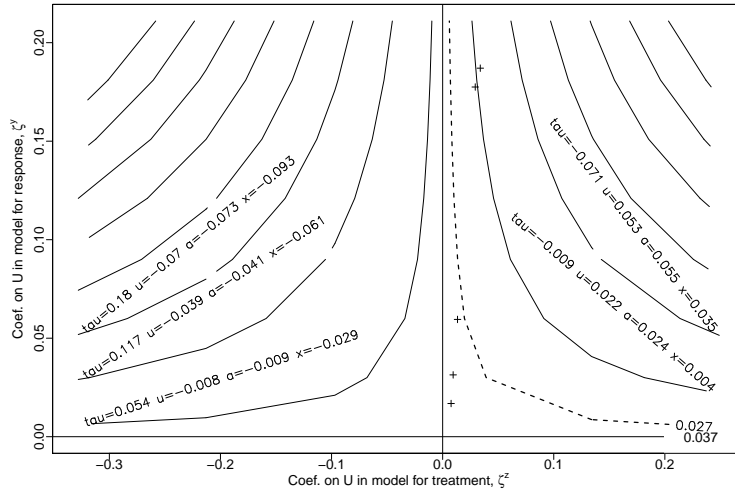


Figure 1: Sensitivity Plot For Shaw et al. (2012) Data

$\tau = -0.007$ . The figure implies that if  $\zeta^z = 0.1$  and  $\zeta^y = 0.05$ , then the true effect would be about  $-0.007$ . The line also provides the decomposed bias  $v = 0.021$ ,  $\alpha = 0.023$  and  $\chi = 0.003$ . From this we can conclude that, if  $\zeta^z = 0.1$  and  $\zeta^y = 0.05$ , then the bias of the estimator without fixed effects,  $v + \chi = 0.024$ , is smaller in magnitude than the bias when adjusting for fixed effects,  $v + \alpha = 0.044$ . The figure also alerts us that the net amplification bias,  $\alpha$ , is relatively large in this case being roughly 100% of the value of the omitted confounder bias,  $v$ , throughout the figure. As a helpful summary, the dashed line represents the threshold separating the region where fixed effects are bias increasing from the region in which the fixed effects are bias reducing. Above and to the right of the line, fixed effects are bias increasing; for all other values of  $\zeta^z$  and  $\zeta^y$  the fixed effects are bias reducing.

The plus signs in the figure represent estimated coefficients on the (standardized) covariates,  $X_2$ , from the outcome and treatment models. As in Carnegie et al. (2014a) and Imbens (2003) we interpret the plus signs as providing benchmarks that help the researcher assess the plausibility of an omitted confounder with similar properties. For example the mark furthest from the origin, at about  $(0.03, 0.19)$ , is plotting the coefficients on the (standardized) indicator of turnout in the 2000

general election. One interpretation is that it is likely that the sensitivity parameters corresponding to the omitted confounder could have properties similar to that of the indicator for turnout in the 2000 general election. Reassuringly, Figure 1 indeed would have provided a researcher with a warning to be wary of fixed effects for this dataset.

## 5. DISCUSSION

We have discussed the ways in which additional control covariates can increase bias, including bias amplification and bias unmasking. This decomposition made apparent a special case of pure bias amplification, in particular, when fixed effects amplify bias. The canonical example of a (pure) bias amplifying covariate in the literature to date has been instruments. However, we have shown that fixed effects can be pure bias amplifiers even though they do not act as instruments and even though they absorb heterogeneity in (and are causally related to) the outcome. Since fixed effects are often characterized as simply inefficient, rather than biased, in the econometric framework, the potential for amplification has been obscured, to date. We then presented a method of visualizing the conditions under which fixed effects are bias increasing (either via unmasking or amplification) to guide the researcher via a calibrated sensitivity analysis.

Our bias decomposition delineates two scenarios in which bias increases upon inclusion of particular covariates. In the case of amplification, the bias formulas provided in this paper help us to better understand the circumstances under which covariates may act as bias amplifiers or bias unmaskers. Examining  $\alpha \equiv \left( \frac{r_{Z|X}^2}{1-r_{Z|X}^2} \right) v$  provides some reassurance that amplification may not be a major concern in practice. It is only greater than the bias due to omitting  $U, v$ , when  $r_{Z|X}^2 > \frac{1}{2}$ . In words,  $X$  would have to account for more than half of the variability in the assignment mechanism for amplification to have the bias to be as large as the bias due to the unobserved confounder.<sup>11</sup> Fortunately  $r_{Z|X}^2$  is identified, a fact that should give us some idea of whether bias amplification should be a particular concern.

---

<sup>11</sup>In the case where  $X$  is a matrix of dummy variables for group, this condition is equivalent to saying that the intraclass correlation (ICC) is greater than 0.5.

However, concern over the phenomenon of bias unmasking should rival concern over bias amplification. In the second case study, for example, the Water Infrastructure outcome reveals that the bias due to omitting fixed effects,  $\chi$ , can be large, but of opposite sign and similar magnitude compared to the bias due to the unobserved confounder,  $v$ . Omitting both the fixed effects and the unobserved confounder was preferable to adjusting for the fixed effects precisely because the two biases counterbalanced one another in the unadjusted estimate. In practice, a researcher is unlikely to know whether adjusting for covariates will unmask unobserved confounder bias. Similar observations have led to somewhat pessimistic assessment of observational analysis, for example, in Clarke (2005) and Frisell et al. (2012) (but see also Clarke 2009).

Sensitivity analysis when considering unobserved confounders has been previously considered elsewhere (e.g. Clarke 2009; Carnegie et al. 2014a; Imbens 2003). We proposed an important modification to sensitivity plots aimed at increasing the information available about the potential for bias amplification. Plotting the bias decompositions ( $\alpha$ ,  $v$ , and  $\chi$ ) on each of the contours should help practitioners to consider bias amplification and bias unmasking.

While better study designs are always the best way to address concerns over the dangers caused by failing to control for all potential confounders in an observational studies, the reality is that many questions of interest are difficult or impossible to study using randomized experiments. In the absence of controlled or natural experiments we need more tools to help applied researchers make the best choices regarding how to perform their analyses. Thoughtful consideration about the potential for bias amplification and unmasking should be part of this process. We hope that the methodology presented in this paper can assist the researcher and makes these ideas more concrete and fully contextualized.

## REFERENCES

Angrist, J.D., Imbens, G. and Rubin, D. (1996). Identification of Causal Effects Using Instrumental Variables. *Journal of the American Statistical Association* 91(434): 444-455.

- Angrist, J.D. and Pischke, J. (2009). *Mostly Harmless Econometrics*. Princeton University Press, Princeton.
- Austin, P. Grootendorst, P. and Anderson, G. (2007) A Comparison of the Ability of Different Propensity Score Models to Balance Measured Variables Between Treated and Untreated Subjects: A Monte Carlo Study. *Statistics in Medicine* 26: 734-753.
- Breen, R. K. Karlson, and A. Holm (2013). Total, direct, and indirect effects in logit and probit models. *Sociological Methods and Research*, 42(2), 164191.
- Brookhart, M., Sturmer, T., Glynn, R., Rassen, J. and Schneeweiss, S. (2010). Confounding Control in Healthcare Database Research. *Medical Care* 48: S114-S120.
- Bhattacharya, J. and Vogt, W. (2007). Do Instrumental Variables Belong in Propensity Scores? *NBER Working Paper* 343, National Bureau of Economic Research, MA.
- Carnegie, NB, J Hill, and M Harada. (2014) Assessing sensitivity to unmeasured confounding using simulated potential confounders. (unpublished)
- Carnegie, NB, J Hill, and M Harada. (2014) Package: TreatSens, <http://www.R-project.org>.
- Clarke, K. A., (2005) The Phantom Menace. *Conflict Management and Peace Science* 22:341352.
- Clarke, K. A., (2009) Return of the Phantom Menace. *Conflict Management and Peace Science* 26: 46-66.
- Cole, SR, R W Platt, E. F. Schisterman, H. Chu, D. Westreich, D. Richardson, and C. Poole (2010). Illustrating bias due to conditioning on a collider. *Int J Epidemiol* 39(2): 417420.
- D'Agostino, Jr. R. (1998). Propensity Score Methods for Bias Reduction in the Comparison of Treatment to Non-Randomized Control Group. *Statistics in Medicine* 17 314-316.

- Ding, Peng and Miratrix, Luke. (2014). To Adjust or Not to Adjust? Sensitivity Analysis of M-bias and Butterfly-Bias. *Journal of Causal Inference* 2: 2-17.
- Dunning, T. and Nilekani, J. (2013). Ethnic Quotas and Political Mobilization: Caste, Parties, and Distribution in Indian Village Councils. *American Political Science Review* 107: 35-56.
- Freedman, D. A. (2008). Randomization Does Not Justify Logistic Regression. *Statistical Science* 23(2): 237-249.
- Frisell, T., Oberg, S., Kuja-Halkola, R., and Sjolander, A. (2012). Sibling Comparison Designs: Bias From Non-Shared Confounders and Measurement Error *Epidemiology* 23(5): 713-720.
- Greene, WH. (2000). *Econometric Analysis* Prentice Hall (4th Edition).
- Greenland, S. (2002). Quantifying biases in causal models: classical confounding vs. collider-stratification bias. *Epidemiology* 14: 300-306.
- Hill, J. (2007a). Discussion of research using propensityscore matching: Comments on A critical appraisal of propensityscore matching in the medical literature between 1996 and 2003 by Peter Austin, *Statistics in Medicine* 27(12): 2055-2061
- Hill, J. (2007b). "Discussion of 'Can Nonrandomized Experiments Yield Accurate Answers? A Randomized Experiment Comparing Random and Nonrandom Assignments'" *Journal of the American Statistical Association* 103(484): 1346-1350.
- Heckman, J. and R. Robb. (1985). "Alternative Methods for Estimating the Impact of Interventions" in J.J. Heckman and B. Singer (Eds.) *Longitudinal Analysis of Labor Market Data*: Cambridge University Press.
- Heckman, J. and R. Robb. (1986). "Alternative Methods for Solving the Problem of Selection bias in Evaluating the Impact of Treatments on Outcomes" in H. Wainer (Ed.) *Drawing Inferences from Self-selected Samples*: New Jersey: Lawrence Erlbaum Associates.

- Imbens, Guido, W. (2003). Sensitivity to Exogeneity Assumptions in Program Evaluation. *Recent Advances in Econometric Methodology* 93(2): 126-132.
- Lechner, M., (2001), Identification and Estimation of Causal Effects of Multiple Treatments under the Conditional Independence Assumption, in Lechner and Pfeiffer (Eds.), *Econometric Evaluations of Active Labor Market Policies in Europe*, Heidelberg, Physica.
- Liu, W., M. A. Brookhart, S. Schneeweiss, X. Mi, and S. Setoguchi (2012). Implications of m-bias in epidemiologic studies: a simulation study. *American Journal of Epidemiology* 176: 938-948.
- Middleton, J.A., (2016), "Replication Data for: Bias Amplification and Bias Unmasking", <http://dx.doi.org/10.791/DVN/UO5WQ4>, Harvard Dataverse
- Myers J. A., J.A. Rassen, J.J. Gagne, K.F. Huybrechts, S. Schneeweiss, K.J. Rothman, M.M. Joffe, R.J. Glynn. (2011). Effects of adjusting for instrumental variables on bias and precision of effect estimates. *American Journal of Epidemiology* 174(11):1213-22.
- Pearl, Judea. (2000) *Causality* Cambridge.
- Pearl J. (2009). Myth, confusion, and science in causal analysis. technical report.
- Pearl, Judea. (2010). On a Class of Bias-Amplifying Variables that Endanger Effect Estimates. *Proceedings of UAI*, 417-424.
- Pearl, J. (2011) Invited Commentary: Understanding Bias Amplification. *American Journal of Epidemiology* 174(11): 1223-1227.
- Rosenbaum, P.R. (2002). *Observational Studies*, Springer.
- Rosenbaum, P. R. and Rubin, D. B. (1983), Assessing sensitivity to an unobserved binary covariate in an observational study with binary outcome. *Journal of the Royal Statistical Society. Series B (Methodological)* 45: 212-218

- Rubin, D.B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66: 688.
- Rubin, D. B. (1978). Bayesian Inference for Causal Effects: The Role of Randomization *The Annals of Statistics* 6(1) 34-58.
- Rubin, D. B. (2002). Using Propensity Scores to Help Deign Observational Studies: Application to the Tobacco Litigation *Health Services and Outcomes Research Methodology* 2 169-188.
- Shaw, D. R., D. P. Green, J. G. Gimpel, and A. S. Gerber. (2012) Do Robotic Calls from Credible Sources Influence Voter Turnout or Vote Choice? Evidence from a Randomized Field Experiment. *Journal of Political Marketing* 11 (4): 231-245.
- Schisterman, E.F., Cole, S. R., and R. W. Platt. (2009). Overadjustment Bias and Unnecessary Adjustment in Epidemiologic Studies. *Epidemiology* 20:488-495.
- Sjlinder, A. (2009). Propensity scores and M-structures. *Stat Med* 28:141620.
- Sobel, M.E., (2006). What Do Randomized Studies of Housing Mobility Demonstrate?: Causal Inference in the Face of Interference. *Journal of the American Statistical Association* 101: 1398-1407.
- Wooldridge, J. (2009). Should instrumental variables be used as matching variables? *Unpublished*.
- Wyss, R., M. Lunt, M. A. Brookhart, R. J. Glynn, T. Strumer. (2014) *Journal of Causal Inference* 2(2): 131-146.
- VanderWeele, T.J. (2015). *Explanation in Causal Inference: Methods for Mediation and Interaction* Oxford University Press.
- VanderWeele, T.J., and O.A. Arahc. (2011) Unmeasured Confounding for General Outcomes, Treatments, and Confounders: Bias Formulas for Sensitivity Analysis. *Epidemiology*. 22(1): 4252.