# Causal Inference with Corrupted Data: Measurement Error, Missing Values, Discretization, and Differential Privacy

Anish Agarwal
Columbia University

Rahul Singh*
Harvard University

Original draft: July 2021. This draft: February 2024.

## Abstract

The US Census Bureau will deliberately corrupt data sets derived from the 2020 US Census, enhancing the privacy of respondents while potentially reducing the precision of economic analysis. To investigate whether this trade-off is inevitable, we formulate a semiparametric model of causal inference with high dimensional corrupted data. We propose a procedure for data cleaning, estimation, and inference with data cleaning-adjusted confidence intervals. We prove consistency and Gaussian approximation by finite sample arguments, with a rate of $n^{-1/2}$ for semiparametric estimands that degrades gracefully for nonparametric estimands. Our key assumption is that the true covariates are approximately low rank, which we interpret as approximate repeated measurements and empirically validate. Our analysis provides nonasymptotic theoretical contributions to matrix completion, statistical learning, and semiparametric statistics. Calibrated simulations verify the coverage of our data cleaning-adjusted confidence intervals and demonstrate the relevance of our results for Census-derived data.

# 1 Introduction

The 2010 US Census inadvertently revealed too much information. In a simulated hack, researchers at the Census Bureau could re-identify between 52 and 179 million respondents from anonymous summary tables [Hawes, 2021]. To protect privacy, the Bureau will inject synthetic noise into summary tables of the 2020 Census and coarsen wage microdata in the Current Population Survey (CPS). Techniques like these, called privacy mechanisms in computer science, guarantee a property called *differential privacy* via deliberate data corruption [Dwork et al., 2006]. Differential privacy is widely implemented in the technology sector, e.g. Apple iOS and Google Chrome data. Due to its recent adoption in the government sector, several economists have warned of a looming trade-off: the privacy of respondents versus the precision of economic analysis [Abowd and Schmutte, 2019, Hotz et al., 2022].

We study differential privacy and discretization as modern challenges for causal inference. Economic data continue to suffer from classical data corruptions such as missing values and measurement error. Therefore, we analyze a class of data corruptions that encompasses both modern and classical issues *simultaneously*, while remaining agnostic about their relative magnitudes. Our research question is how (and when) it is possible to estimate typical causal parameters using high dimensional economic data that suffer from measurement error, missing values, discretization, and differential privacy mechanisms. An answer requires nonasymptotic analysis because differential privacy is defined as a finite sample property.

We study a broad class of causal parameters, including semiparametric scalars such as the average treatment effect, the local average treatment effect, and the average elasticity, as well as nonparametric functions such as heterogeneous treatment effects, in a nonlinear and high dimensional setting. Our main contribution is a procedure for automatic data cleaning, causal estimation, and inference with confidence intervals that account for the bias and variance consequences of data cleaning. The procedure is simple. It essentially combines principal component analysis, ordinary least squares, and sample splitting in new ways.

Our key assumption is that the true covariates are approximately low rank, which we validate for US Census-derived data and interpret from a causal perspective. We argue that covariates collected from the Census include approximate repeated measurements—e.g.

disability benefits, medical benefits, and unemployment benefits—which implies that they are approximately low rank. There are three key aspects of our contribution.

First, our simple procedure adapts to the *type* and *level* of data corruption. The same code works in a variety of settings, allowing for classical types such as measurement error and missing values as well as modern types such as discretization and differential privacy mechanisms, across variance levels. Crucially, the researcher does not need exact knowledge of the corruption distribution, e.g. its parametric form or covariance structure, and in this way we depart from the error-in-variable Lasso and Dantzig literatures; see Section 2. We depart from previous work on principal component regression by proposing new variants that involve "implicit" data cleaning—i.e. prediction on a test observation without cleaning it—and inference in nonlinear, heterogeneous causal models. We propose an error-in-variable balancing weight that adapts to the causal parameter of interest, which is a natural yet original solution for cross sectional data. Our theory of implicit data cleaning and our error-in-variable balancing weight appear to be new. The former is of independent interest.

Second, our theoretical analysis allows the rate of data cleaning to be slower than the rate of causal inference, so an analyst can use matrix completion for automatic data cleaning of covariates. We extend the classic semiparametric framework, where the goal is to obtain $n^{-1/2}$ convergence for the causal parameter despite a slower rate of convergence for a nonparametric "nuisance" regression. Our goal is to obtain $n^{-1/2}$ convergence for the causal parameter despite a slower rate of convergence for high dimensional data cleaning, which is a "nuisance" task. Since our data cleaning guarantees only hold on-average, we are unable to use previous semiparametric results; instead, we generalize semiparametric and nonparametric debiased machine learning theory to i.n.i.d. corrupted data, with new results on nominal and conservative variance estimation. Altogether, our framework translates slow, on-average data cleaning guarantees into fast causal estimation and inference guarantees.

Third, our empirical results suggest that there exist scenarios in which the trade-off between privacy and precision can be overcome, and others in which it cannot. We replicate and extend [Autor et al., 2013]'s seminal paper on the effect of import competition in US labor markets. To begin, we demonstrate the plausibility of our key assumption: Census data products contain many variables that are approximate repeated measurements. Next, we corrupt the data, injecting synthetic noise calibrated to the privacy level mandated for

the 2020 US Census. We implement differential privacy and discretization in a way that belongs to our class of data corruptions, which can therefore be cleaned and adjusted for in the confidence interval. We recover the main results of [Autor et al., 2013] without losing statistical precision. In this representative setting for economic research, it appears to be possible to achieve both privacy at the individual level and precision at the population level.

Section 2 situates our contributions within the context of related work. Section 3 formalizes our class of data corruptions and our key assumption. Section 4 proposes our procedure and demonstrates its performance in simulations. Section 5 theoretically justifies our procedure, and verifies the key assumption for nonlinear factor models. Section 6 presents the semi-synthetic exercise and discusses limitations. Section 7 concludes.

## 2   Related work

**Semiparametrics**. We use two insights from classic [Hasminskii and Ibragimov, 1979, Klaassen, 1987, Robinson, 1988, Bickel et al., 1993, Andrews, 1994, Newey, 1994, Robins and Rotnitzky, 1 Ai and Chen, 2003, Van der Laan and Rubin, 2006, Hahn and Ridder, 2013] and modern [Zheng and Van der Laan, 2011, Athey et al., 2018, Chernozhukov et al., 2018, Hirshberg and Wager, 202 Chernozhukov et al., 2022a, Chernozhukov et al., 2023] semiparametric theory. First, a causal parameter typically has regression and balancing weight representations, and both appear in the semiparametrically efficient asymptotic variance. We directly build on this insight: an error-in-variable regression and an error-in-variable balancing weight appear in our data cleaning-adjusted confidence intervals. Second, sample splitting eliminates restrictive conditions on the data generating process and estimation procedure. We combine these two classic ideas with implicit data cleaning, which appears to be a new idea.

**Error-in-variable regression**. We provide a framework to repurpose error-in-variable regression estimators for downstream causal inference. Error-in-variable regression has a vast literature spanning econometrics, statistics, and computer science studying the model

$$Y_i = \gamma_0(X_{i,\cdot}) + \varepsilon_i, \quad Z_{i,\cdot} = X_{i,\cdot} + H_{i,\cdot} \quad (X_{i,\cdot} \text{ is the } i\text{-th row of matrix } \boldsymbol{X} \text{ and so on}) \quad (1)$$

where $(X_{i,\cdot}, \varepsilon_i, H_{i,\cdot})$ are mutually independent and $(\varepsilon_i, H_{i,\cdot})$ are mean zero. We consider a generalization of this setting with missingness, and we define our causal parameter as a scalar

summary of nonlinear $\gamma_0$. Methods in econometrics typically assume auxiliary information for identification: repeated measurements [Hausman et al., 1991, Li and Vuong, 1998, Schennach, 2004], instrumental variables [Schennach, 2007, Hu and Schennach, 2008], and negative controls [Miao et al., 2018, Deaner, 2018]. Similar in spirit to repeated measurements, we assume $\boldsymbol{X}$ is approximately low rank. Methods in statistics extend the Lasso and Dantzig selector to high dimensional error-in-variable regression [Loh and Wainwright, 2012, Rosenbaum and Tsybakov, 2013, Datta and Zou, 2017]. However, these methods require linearity and exact sparsity of $\gamma_0$, as well as knowledge of the covariance of measurement error $H_{i,\cdot}$. By contrast, we assume $H_{i,\cdot}$ are subexponential; the analyst does not need to know the measurement error covariance, and therefore can be agnostic about the type and level of corruption. We propose new variants of principal component regression (PCR) for the error-in-variable regression and balancing weight. Previous work studies PCR for error-in-variable regression only, explicitly cleaning all observations [Stock and Watson, 2002a, Bai and Ng, 2006, Agarwal et al., 2020a]. We develop a technique of implicit data cleaning that avoids mixing together signal and noise across observations, which aids with downstream statistical inference of nonlinear models. Moreover, our error-in-variable balancing weight for cross sectional data appears to be new.

**PCA for large factor models**. The initial step of PCR is PCA. A vast literature studies the identification, estimation, and inference of latent factors $(\lambda_i, \mu_j)$ in models of the form

$$Z_{i,\cdot} = X_{i,\cdot} + H_{i,\cdot}, \quad X_{ij} = \lambda_i^T \mu_j \quad (X_{i,\cdot} \text{ is the } i\text{-th row of matrix } \boldsymbol{X} \text{ and so on}) \quad (2)$$

where $Z_{i,\cdot}$ is observed, the ambient dimension $dim(X_{i,\cdot})$ is high, and the latent dimension $dim(\lambda_i)$ is fixed [Bai, 2003, Bai and Ng, 2013]. Our interest in downstream causal inference allows us to bypass the issue of identifying latent factors, and to relax the linear factor model; instead, we require that the approximate rank of $\boldsymbol{X}$ diverges more slowly than $dim(X_{i,\cdot})$. The nonlinear factor model $X_{ij} = g(\lambda_i, \mu_j)$, where $dim(\lambda_j)$ may slowly diverge and $g$ is smooth, is *sufficient* but *unnecessary* for our analysis. Like the factor model literature, we allow weak correlation and heteroscedastity of measurement errors within units.

**Low rank causal models**. Whereas we study treatment effects, policy effects, and elasticities in cross sectional data, a rich literature studies treatment effects, in panel data,

via a low rank factor model for potential outcomes [Athey et al., 2021, Bai and Ng, 2019, Xiong and Pelger, 2023, Fernández-Val et al., 2021, Agarwal et al., 2020b, Feng, 2020]. By contrast, we study a more general class of causal parameters, in cross sectional data, when covariates are approximately low rank. The only previous work to consider both measurement error and missingness in cross sectional treatment effects appears to be [Kallus et al., 2018]. The authors study average treatment effect and prove consistency, without inference, for a parametric linear model where the true covariates are low dimensional Gaussians and the measurement error distribution is correctly specified. By contrast, we study a broad class of semiparametric and nonparametric causal parameters and provide inference, with data cleaning-adjusted confidence intervals. We do not require exact distributional knowledge of (high dimensional) true covariates or measurement error.

**Privacy in econometrics**. Our research question complements others in a recent literature on private econometrics. One strand considers how to disclose estimates obtained with private data access [Dwork and Lei, 2009, Smith, 2011, Komarova and Nekipelov, 2020]. Another considers how to conduct estimation after linking public and private records, where privacy considerations constrain linkage [Komarova et al., 2018]. We ask a complementary question motivated by empirical economic research using public, Census-derived data products: how to conduct causal inference in the presence of simultaneous data corruptions, including canonical privacy mechanisms applied to the data before estimation.

# 3   Model overview

**Causal parameter**. For readability, we focus on one causal parameter in the main text: the average treatment effect (ATE) with i.n.i.d. data $\theta_0 = \frac{1}{n}\sum_{i=1}^{n}\theta_i$, where $\theta_i = \mathbb{E}[Y_i^{(1)} - Y_i^{(0)}]$. Here, $Y_i^{(d)}$ is the potential outcome for unit $i$ under intervention $D = d$. $\theta_0$ is a sample average because different units may be drawn from different distributions—a challenging yet plausible scenario when data are corrupted. With i.i.d. data, $\theta_0$ simplifies to the familiar ATE. Appendix E considers a general class of semi- and nonparametric causal parameters e.g. the local average treatment effect, average elasticity, and heterogeneous treatment effects.

We denote the actual outcome by $Y_i \in \mathbb{R}$, the assigned treatment by $D_i \in \{0, 1\}$, and

the covariates that determine treatment assignment by $X_{i,\cdot} \in \mathbb{R}^p$. In order to express $\theta_0$ in terms of $(Y_i, D_i, X_{i,\cdot})$, we impose some additional structure on the problem. Generalizing a classic assumption in the literature on distribution shift, we assume that the conditional distributions $\mathbb{P}(Y_i | D_i, X_{i,\cdot})$ and $\mathbb{P}(D_i | X_{i,\cdot})$ are common across units; distribution shift is only in the marginal distributions of covariates $\mathbb{P}_i(X_{i,\cdot})$.

Imposing these conditions as well as selection on $X_{i,\cdot}$, we recover two classic formulations of the treatment effect. The outcome formulation is in terms of the outcome mechanism $\gamma_0(D_i, X_{i,\cdot}) = \mathbb{E}[Y_i | D_i, X_{i,\cdot}]$, also called the regression, which is common across units: $\theta_i = \mathbb{E}[\gamma_0(1, X_{i,\cdot}) - \gamma_0(0, X_{i,\cdot})]$. The treatment formulation is in terms of the treatment mechanism $\mathbb{E}[D_i | X_{i,\cdot}]$, which is also common across units, and which appears in the denominator of the balancing weight $\alpha_0(D_i, X_{i,\cdot}) = \frac{D_i}{\mathbb{E}[D_i | X_{i,\cdot}]} - \frac{1-D_i}{1-\mathbb{E}[D_i | X_{i,\cdot}]}$: here, $\theta_i = \mathbb{E}[Y_i \cdot \alpha_0(D_i, X_{i,\cdot})]$. Our estimation and analysis combine both classic formulations.

**Data corruption**. The crux of our problem is that we observe $(Y_i, D_i, Z_{i,\cdot})$ instead:

$$Y_i = \gamma_0(D_i, X_{i,\cdot}) + \varepsilon_i, \quad Z_{i,\cdot} = (X_{i,\cdot} + H_{i,\cdot}) \odot \pi_{i,\cdot}. \tag{3}$$

Though the outcome $Y_i$ is generated from treatment $D_i$ and true covariates $X_{i,\cdot}$, we do not observe $X_{i,\cdot}$; instead, we observe the corrupted covariates $Z_{i,\cdot}$, which are the true covariates $X_{i,\cdot}$ plus conditionally mean zero corruption $H_{i,\cdot}$, multiplied elementwise by an independent masking vector $\pi_{i,\cdot} \in \{\texttt{NA}, 1\}^p$. Our concise model (3) generalizes the models (1) and (2), and it encompasses all four types of corruption. For example, to encode classical measurement error, let $Z_{i,\cdot}$ equal $X_{i,\cdot}$ plus a vector of Gaussian noise. To encode missing values, let $Z_{i,\cdot} = X_{i,\cdot} \odot \pi_{i,\cdot}$. In Appendix E, we accommodate corruption of the outcome $Y_i$ and treatment $D_i$, under restrictions.

Discretization is a process by which a continuous vector $X_{i,\cdot}$ maps to a discrete vector $Z_{i,\cdot}$, and our class encodes variants where $\mathbb{E}[Z_{i,\cdot} | X_{i,\cdot}] = X_{i,\cdot}$. For example, the covariate of interest may be a vector of probabilities $X_{i,\cdot}$, yet we observe actual occurrences $Z_{i,\cdot} \sim \text{Bernoulli}(X_{i,\cdot})$. Another example is randomized rounding, where continuous values are randomly rounded to nearby integers, e.g. $Z_{i,\cdot} = sign(X_{i,\cdot})\text{Poisson}(|X_{i,\cdot}|)$. Our class does not include deterministic rounding. Instead, it provides guidance on which types of rounding can be handled well in downstream causal inference. As such, it suggests alternative types of discretization for wage data in the CPS which are more favorable for economic research.

Differential privacy is a concept from computer science which means plausible deniability that any individual contributed their data to tabular summaries. The canonical mechanism that ensures differential privacy is to release $Z_{i,\cdot}$ equal to $X_{i,\cdot}$ plus a vector of Laplacian noise, calibrating the variance of the Laplacian to a priori bounds on the true values and other properties of the tabular summary statistics [Dwork et al., 2006]. Our framework allows for other canonical privacy mechanisms where $\mathbb{E}[Z_{i,\cdot}|X_{i,\cdot}] = X_{i,\cdot}$, e.g. discrete Gaussian, piece wise uniform, and bounded mechanisms. In the context of the Census, we consider adding Laplacian noise to data on aggregate units, which we formalize in Section 6. Injecting synthetic noise in this way helps to prevent the kind of attack simulated on the 2010 Census.

Across examples, $H_{i,\cdot}$ is subexponential, i.e. its tails are no worse than an exponential distribution's. So are compositions of various types of data corruption since the class of subexponential distributions is closed under addition. Therefore our class of data corruptions includes classical and modern issues *simultaneously*. It allows us to address the trade-off between privacy and precision in the context of heteroscedastic measurement error—a major aspect of the problem often overlooked [Chetty and Friedman, 2019, Steed et al., 2022].

What corruptions do we exclude? Our definition of $H_{i,\cdot}$ rules out nonseparable or endogeneous measurement error. Our definition of $\pi_{i,\cdot}$ rules out endogenous missingness, i.e. sample selection. We also rule out deterministic rounding, i.e. interval censoring. However, our class includes canonical privacy mechanisms as well as randomized rounding, itself a privacy mechanism. Therefore we study this class, which we formalize in Section 5. Importantly, we allow heteroscedastic corruptions that are dependent within a unit.

**Key assumption: Approximate repeated measurements**. Our key assumption is that the true covariates are approximately low rank: the rank of the matrix $\boldsymbol{X} \in \mathbb{R}^{n \times p}$ is approximately $r \ll (n, p)$. Among the $p$ covariates in the data set, there are approximately only $r$ latent types of covariates. For intuition, consider repeated measurements. In the classic repeated measurement model, we have two noisy measurements of one signal. In our model, we have $p$ noisy measurements $(Z_{i1}, ..., Z_{ip})$ that are approximately repeated measurements of only $r$ signals, where both $(r, p)$ grow with sample size $n$, yet $r \ll (n, p)$.

We place this assumption because it seems plausible in Census-derived data. Consider the commuting zone (CZ) level data set of [Autor et al., 2013]. Each CZ is a local economy with a vector of covariates $X_{i,\cdot} \in \mathbb{R}^{30}$ if we use variables from the authors' preferred specification

as well as additional variables from their appendix. The variables include average disability, unemployment, and medical benefits, which are not precisely repeated measurements but approximately so. We compute the singular value decomposition of $\boldsymbol{X}$ then visualize its singular values, also called its principal components, in Figure 1. We see that only about five principal components are significantly positive; $r = 5$ while $p = 30$.

Our key assumption admits a causal interpretation in the running example of ATE. Consider the special case where the true covariates are exactly low rank, i.e. $r = rank(\boldsymbol{X})$. The singular value decomposition is $\boldsymbol{X} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^T$ where $\boldsymbol{U} \in \mathbb{R}^{n\times r}$, $\boldsymbol{\Sigma} \in \mathbb{R}^{r\times r}$, and $\boldsymbol{V} \in \mathbb{R}^{p\times r}$. $\boldsymbol{V}$ consists of $r$ vectors in $\mathbb{R}^p$, called the right singular vec-
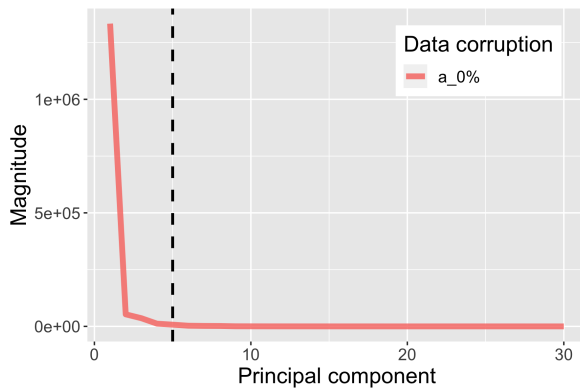


Figure 1: Key assumption in Census data

tors of $\boldsymbol{X}$, which are also the eigenvectors of the empirical covariance $n^{-1}\boldsymbol{X}^T\boldsymbol{X}$. The span of these vectors is an $r$ dimensional subspace of $\mathbb{R}^p$, i.e. a low dimensional subset of a high dimensional ambient space. In this scenario, we assume that treatment assignment is determined by the subspace. More generally, when covariates are approximately low rank, $\boldsymbol{X} = \boldsymbol{X}^{(\mathrm{LR})} + \boldsymbol{E}^{(\mathrm{LR})}$, where $\boldsymbol{X}^{(\mathrm{LR})} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^T$ is a rank $r$ approximation to $\boldsymbol{X}$, and $\boldsymbol{E}^{(\mathrm{LR})}$ is the approximation residual. We can either assume (i) selection is determined by $\boldsymbol{X}^{(\mathrm{LR})}$ only, i.e. the treatment assignment for unit $i$ depends on the *projection* of $X_{i,\cdot}$ onto the subspace spanned by $\boldsymbol{V}$; or (ii) selection is determined by both $\boldsymbol{X}^{(\mathrm{LR})}$ and $\boldsymbol{E}^{(\mathrm{LR})}$. To handle the latter, we keep track of $\Delta_E = \|\boldsymbol{E}^{(\mathrm{LR})}\|_{\max}$ in our theoretical analysis. Our analysis is robust to small violations of the exactly low rank assumption from statistical and causal perspectives.

# 4  Data cleaning-adjusted confidence interval

We would like a procedure that estimates parameters in nonlinear, heterogeneous causal models as if data were uncorrupted, yet adjusts for data cleaning in the confidence interval. Moreover, we would like a procedure that does not require knowledge of the corruption covariance structure in advance, departing from previous work. If such a procedure were to

exist, it would in some sense preempt the looming trade-off between privacy and precision.

**Why is inference hard**? We illustrate our procedure with an average treatment effect simulation. By construction, the treatment effect is $\theta_0 = 2.2$. We consider a data generating process (DGP) which satisfies our key assumption: one sample involves a matrix of covariates $\boldsymbol{X} \in \mathbb{R}^{100 \times 100}$ with rank $r = 5$. See Appendix L for details and for
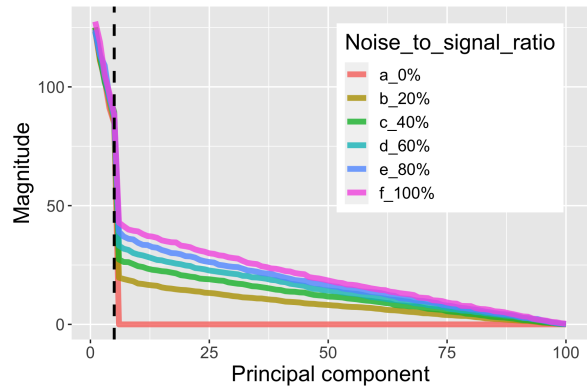


Figure 2: Key assumption in simulated data

similar results using alternative dimensions of $\boldsymbol{X}$. The DGP has nonlinear outcome and treatment mechanisms. Figure 2 plots the principal components of true covariates $\boldsymbol{X}$ in red. As expected, five principal components are nonzero and the rest are zero since $rank(\boldsymbol{X}) = 5$.

As a first pass, we implement ordinary least squares (OLS) of $Y_i$ on $(D_i, X_{i,\cdot})$. Running OLS on clean data 1000 times, the point estimates $\hat{\theta}$ (Figure 3a) center around the true value 2.2, and appear Gaussian. OLS works well in the absence of data corruption; there is nothing hidden in the DGP for clean data. We repeat this exercise introducing measurement error with variance that is 20% of the variance of the covariates. Inversion of the empirical covariance matrix $n^{-1}\boldsymbol{Z}^T\boldsymbol{Z}$ becomes numerically unstable, manifesting in point estimates that are erratic (Figure 3b) and standard errors that are typically `NA`'s. Notably, data corruption flips the sign about a quarter of the time, a phenomenon we verify for 2SLS and for settings closer to [Autor et al., 2013] in Appendix L. OLS is not well-suited to the combination of high dimensional covariates, (approximate) low rank, and measurement error. Indeed, any estimator that ignores covariate measurement error in a nonlinear, heterogeneous causal model suffers from bias of a complicated form [Battistin and Chesher, 2014].

Data corruption can derail causal inference, which motivates filling the `NA`'s, reigning in the extremes, and otherwise de-noising the values in $\boldsymbol{Z}$ in hopes of recovering $\boldsymbol{X}$. These are precisely the goals of matrix completion applied to the matrix $\boldsymbol{Z}$ [Candès and Recht, 2009, Candès and Tao, 2010, Keshavan et al., 2009, Hastie et al., 2015, Chatterjee, 2015] . Our goal is to automate data cleaning via matrix completion, then to adjust for data cleaning in the confidence interval. To select an appropriate matrix completion method, we return to Figure 2 to visualize the principal components of the corrupted covariates $\boldsymbol{Z} = \boldsymbol{X} + \boldsymbol{H}$

9

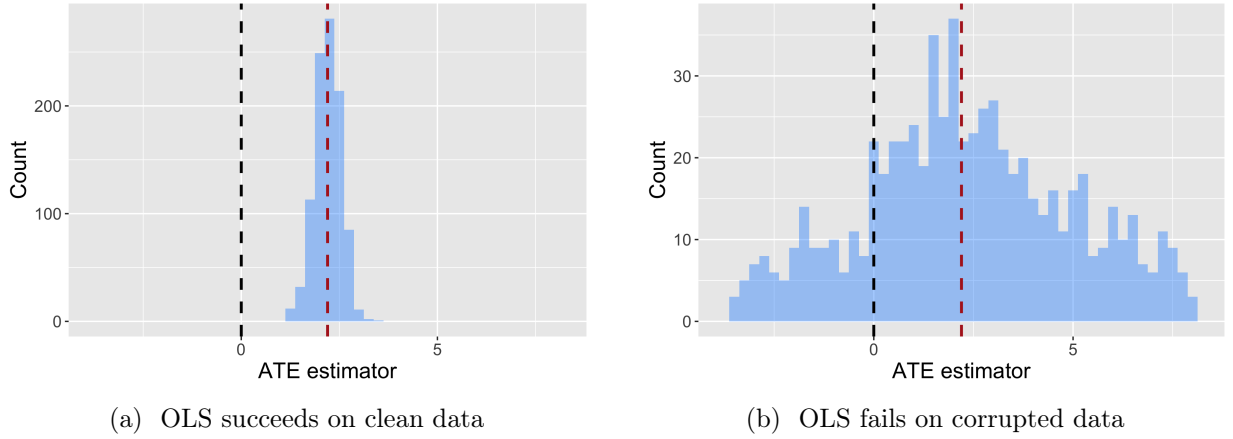(a) OLS succeeds on clean data   (b) OLS fails on corrupted data

Figure 3: First pass — OLS

for various noise-to-signal ratios (i.e. the noise variance divided by the signal variance). The initial five principal components remain virtually unchanged, while the lower principal components are amplified; signal remains spectrally concentrated while noise is spectrally diffuse. Therefore a natural way to clean the covariates is to discard the lower principal components—in essence, to perform principal component analysis (PCA).[1]

Why is inference hard after data cleaning? Several challenges arise. First, data cleaning may mix together signal and noise across observations; yet, we wish to prove Gaussian approximation via a central limit theorem. Our solution is to break dependence via both sample splitting, which is a classic idea [Klaassen, 1987], and implicit data cleaning, which is a new idea. Second, if we turn to automated data cleaning, the best rates of convergence to the true matrix $\boldsymbol{X}$ are slower than $n^{-1/2}$; yet, we wish to obtain a standard error of order $\hat{\sigma}n^{-1/2}$ for $\theta_0$. Our solution is to use a doubly robust estimating equation and to generalize double rate robustness [Chernozhukov et al., 2018, Van der Laan and Rose, 2018, Rotnitzky et al., 2021]. The third issue is a theoretical one to which we will return in Section 5: the best rates of matrix completion are not for recovering specific matrix entries but rather averages across matrix entries; yet, we wish to obtain downstream semiparametric inference. Our solution is to develop an algorithmic and analytic framework that forges a connection.

**Overview of the procedure**. Split the observations $(Y_i, D_i, Z_{i,\cdot})$ into equally sized TRAIN and TEST sets, each with $m = n/2$ observations. Our procedure consists of four steps, which we state at a high level before filling in the details: (i) data cleaning: $\hat{\boldsymbol{X}}$ using

---

[1]Alternative choices include canonical correlation analysis and partial least squares, which clean $\boldsymbol{Z}$ using $Y$. We leave these directions to future research.

TRAIN; (ii) error-in-variable regression: $\hat{\gamma}$ using TRAIN; (iii) error-in-variable balancing: $\hat{\alpha}$ using TRAIN; (iv) causal parameter: $\hat{\theta} \pm 1.96\hat{\sigma}n^{-1/2}$ using TEST. We opt for simplicity at each step, essentially combining PCA and OLS (albeit in new ways). We view these high level steps a template for more complex procedures in future work.

*Step 1: Data cleaning.* The automated data cleaning procedure is extremely simple: fill in missing values as zeros, scale appropriately, then perform PCA.

For any mathematical operation to be well defined, the NA's must be filled in somehow. To begin, we tally the likelihood of non-missingness for each covariate $j \in [p]$ in TRAIN: $\hat{\rho}_j = \max\{\frac{1}{m}\sum_{i \in \text{TRAIN}} \mathbb{1}(Z_{ij} \neq \text{NA}), \frac{1}{m}\}$, and $\hat{\boldsymbol{\rho}} = diag(\hat{\rho}_1, ..., \hat{\rho}_p) \in \mathbb{R}^{p \times p}$. Next, we fill in missing values with a FILL operator defined such that $\text{FILL}(Z_{ij}) = \frac{Z_{ij}}{\hat{\rho}_j}$ if $Z_{ij} \neq \text{NA}$ and $\text{FILL}(Z_{ij}) = 0$ if $Z_{ij} = \text{NA}$. Let $\boldsymbol{Z}^{\text{TRAIN}}$ be rows of $\boldsymbol{Z}$ where $i \in \text{TRAIN}$. The FILL operator may act on $\boldsymbol{Z}^{\text{TRAIN}}$ or $\boldsymbol{Z}^{\text{TEST}}$, but it always uses the likelihoods $\hat{\boldsymbol{\rho}}$ calculated from $\boldsymbol{Z}^{\text{TRAIN}}$.

**Proposition 4.1** (Filling with zeros is unbiased and simple). *For $i \in$ TEST,*

$$\mathbb{E}[\text{FILL}(Z_{ij})|X_{ij}, \text{TRAIN}] = X_{ij}\frac{\rho_j}{\hat{\rho}_j}.$$

*The alternative procedure of filling missing values with averages from TRAIN, denoted by $\bar{Z}_j^{\text{TRAIN}}$, gives*

$$\mathbb{E}[\text{FILL-AS-MEANS}(Z_{ij})|X_{ij}, \text{TRAIN}] = X_{ij}\rho_j + \bar{Z}_j^{\text{TRAIN}}(1 - \rho_j).$$

FILL-AS-MEANS gives a convex combination of the signal $X_{ij}$ and of the noisy average $\bar{Z}_j^{\text{TRAIN}}$. The noisy average introduces additional correlations that our procedure avoids.

After filling TRAIN, we project it onto its own principal subspace to calculate the cleaned training covariates $\hat{\boldsymbol{X}}$: $\text{FILL}(\boldsymbol{Z}^{\text{TRAIN}}) = \hat{\boldsymbol{U}}\hat{\boldsymbol{\Sigma}}\hat{\boldsymbol{V}}^T$, and $\hat{\boldsymbol{X}} = \hat{\boldsymbol{U}}_k\hat{\boldsymbol{\Sigma}}_k\hat{\boldsymbol{V}}_k^T$. We truncate the SVD of $\text{FILL}(\boldsymbol{Z}^{\text{TRAIN}})$ to include only the top $k$ principal components, where $k$ is a hyperparameter. Figure 2 suggests a choice of $k$. Below, we empirically verify that our results are robust to different choices of $k > r$. Future work may derive a data driven procedure $k = \hat{r}$ [Stock and Watson, 2002b, Bai and Ng, 2002, Onatski, 2009]. We preserve the ambient dimension $p$.

*Step 2: Error-in-variable regression.* Our error-in-variable regression is also simple: after cleaning TRAIN, perform ordinary least squares (OLS) on TRAIN, then use this OLS coefficient on the filled TEST for prediction. We only fill, and do not clean, the test set.

We introduce nonlinearity into the regression to allow treatment effect heterogeneity, which is crucial for causal inference. Appendix F characterizes what nonlinearity is allowed. Here, we focus on the interacted dictionary $b(D_i, \hat{X}_{i,\cdot}) = (D_i \hat{X}_{i,\cdot}, (1 - D_i) \hat{X}_{i,\cdot})$. Then the OLS coefficient $\hat{\beta} = \{(\hat{\beta}^{\text{TREAT}})^T, (\hat{\beta}^{\text{UNTREAT}})^T\}^T$ equals $[\{b(D^{\text{TRAIN}}, \hat{X})\}^T b(D^{\text{TRAIN}}, \hat{X})]^\dagger$ $[\{b(D^{\text{TRAIN}}, \hat{X})\}^T Y^{\text{TRAIN}}]$, where $\dagger$ means pseudoinverse.

The subtlety is in how predictions are constructed from $\hat{\beta}$. Out of sample prediction does *not* involve cleaning the test set: for $i \in \text{TEST}$, $\hat{\gamma}(D_i, Z_{i,\cdot}) = b\{D_i, \text{FILL}(Z_{i,\cdot})\}\hat{\beta}$.

**Proposition 4.2** (Implicit data cleaning preserves independence). *For $i \in$ TEST, $\hat{\gamma}(D_i, Z_{i,\cdot}) = b(D_i, Z_{i,\cdot})\tilde{\beta}$, where $\tilde{\beta} = \{(\hat{\boldsymbol{\rho}}^{-1} \hat{\beta}^{TREAT})^T, (\hat{\boldsymbol{\rho}}^{-1} \hat{\beta}^{UNTREAT})^T\}^T$ and we replace* `NA` *with 0 in $Z_{i,\cdot}$. Therefore for $(i,j) \in$ TEST, $\hat{\gamma}(D_i, Z_{i,\cdot}) \perp\!\!\!\perp \hat{\gamma}(D_j, Z_{j,\cdot}) | $ TRAIN.*

Remarkably, post-multiplying $b(D_i, Z_{i,\cdot})$ by $\tilde{\beta}$ handles the measurement error, missingness, discretization, and differential privacy of $Z_{i,\cdot}$ while also producing high quality nonlinear predictions of $Y_i$. We call this phenomenon "implicit" data cleaning. Moreover, since $\tilde{\beta}$ is learned exclusively from TRAIN, it is deterministic conditional on TRAIN, so predictions for observations $(i,j) \in$ TEST preserve their independence. This property of implicit data cleaning will be essential for our inferential theory.

Our new variant of PCR has broader use outside of causal inference. In online learning, a corrupted test observation $Z_{i,\cdot}$ does not need to be explicitly cleaned with respect to TEST or even TRAIN. Instead, it may be implicitly cleaned by post multiplying it with the coefficient $\tilde{\beta}$. For test observations, data cleaning and prediction can be combined into one step.

*Step 3: Error-in-variable balancing.* Our error-in-variable balancing weight generalizes our error-in-variable regression. It avoids the estimation and inversion of propensity scores, which may be numerically unstable in high dimensions. Pleasingly, it achieves exact balance for any finite sample size, in a sense that we formalize below. Moreover, it adapts to the causal parameter of interest, as we explain in Appendix I.

The only difference from the error-in-variable regression is that we replace the sufficient statistic $[\{b(D^{\text{TRAIN}}, \hat{X})\}^T Y^{\text{TRAIN}}] \in \mathbb{R}^{p'}$ with another sufficient statistic that we call the counterfactual moment $\hat{M} \in \mathbb{R}^{p'}$. The counterfactual moment resembles the expression $\theta_i = \mathbb{E}[\gamma_0(1, X_{i,\cdot}) - \gamma_0(0, X_{i,\cdot})]$, and it is the *only* aspect of the algorithm that changes

for different causal parameters. Formally, $\hat{\eta} = [\{b(D^{\text{TRAIN}}, \hat{\boldsymbol{X}})\}^T b(D^{\text{TRAIN}}, \hat{\boldsymbol{X}})]^{\dagger} \hat{M}$ and $\hat{M} = [\{b(1, \hat{\boldsymbol{X}})\}^T - \{b(0, \hat{\boldsymbol{X}})\}^T] \mathbb{1}_m$ where $\mathbb{1}_m \in \mathbb{R}^m$ is a vector of ones. As before, we do not clean the test set: for $i \in$ TEST, $\hat{\alpha}(D_i, Z_{i,\cdot}) = b\{D_i, \text{FILL}(Z_{i,\cdot})\}\hat{\eta}$.

**Proposition 4.3** (The balancing weight exactly balances covariates). *For any finite sample,*

$$\frac{1}{m} \sum_{i \in \text{TRAIN}} \hat{X}_{i,\cdot} = \frac{1}{m} \sum_{i \in \text{TRAIN}} D_i \hat{X}_{i,\cdot} \hat{\omega}_i^{\text{TREAT}} = \frac{1}{m} \sum_{i \in \text{TRAIN}} (1 - D_i) \hat{X}_{i,\cdot} \hat{\omega}_i^{\text{UNTREAT}},$$

*where $(\hat{\omega}_i^{\text{TREAT}}, \hat{\omega}_i^{\text{UNTREAT}}) \in \mathbb{R}$ are balancing weights computed from $\hat{\eta} = \{(\hat{\eta}^{\text{TREAT}})^T, (\hat{\eta}^{\text{UNTREAT}})^T\}^T$ as $\hat{\omega}_i^{\text{TREAT}} = \hat{X}_{i,\cdot} \hat{\eta}^{\text{TREAT}}$ and $\hat{\omega}_i^{\text{UNTREAT}} = -\hat{X}_{i,\cdot} \hat{\eta}^{\text{UNTREAT}}$.*

Deterministically, the error-in-variable balancing weight exactly balances the full population, the treated subpopulation, and the untreated subpopulation with respect to their cleaned covariates. It is precisely the reweighting that would ensure comparability of treated and untreated groups in a stratified sampling design. We articulate a more general balancing property for generic causal parameters in Appendix I. We also clarify the sense in which the error-in-variable regression and balancing weight coincide on TRAIN but not TEST.

*Step 4: Causal estimation and inference.* The final step uses the error-in-variable regression $\hat{\gamma}$ and error-in-variable balancing weight $\hat{\alpha}$ learned from TRAIN, and evaluates them on TEST according to the doubly robust estimating equation: for $i \in$ TEST, $\hat{\psi}_i = \hat{\gamma}(1, Z_{i,\cdot}) - \hat{\gamma}(1, Z_{i,\cdot}) + \hat{\alpha}(D_i, Z_{i,\cdot})\{Y_i - \hat{\gamma}(D_i, Z_{i,\cdot})\}$ is the empirical influence of that observation. This process generates a vector $\hat{\psi} \in \mathbb{R}^m$. Reversing the roles of TRAIN and TEST, we generate another such vector. Slightly abusing notation, we concatenate the two to obtain a vector $\hat{\psi} \in \mathbb{R}^n$. We estimate the causal parameter as $\hat{\theta} = \text{MEAN}(\hat{\psi})$, its variance as $\hat{\sigma}^2 = \text{VAR}(\hat{\psi})$, and its data cleaning-adjusted confidence interval as $\text{CI} = \hat{\theta} \pm 1.96\hat{\sigma}n^{-1/2}$.

Our procedure deals with measurement error bias by cleaning the data. For the special case of ATE, the measurement error bias has a closed form solution in terms of the regression, propensity score, covariate density, and derivatives thereof [Battistin and Chesher, 2014]. We avoid estimation of the propensity scores, covariate density, and derivatives, which would be challenging in high dimensions. Instead, we simply combine PCA and OLS.

The way we impute missing values modifies multiple imputation [Rubin, 1976]. In multiple imputation, the analyst makes, say, two copies of the original data set, then imputes missing values (with some randomness so each imputation may be different).

Estimates and standard errors from each copy are then averaged. Our procedure splits the sample into two folds: TRAIN and TEST. We clean TRAIN and compute estimates and standard errors with TEST, then reverse the roles and take the average. We opt for sample splitting, rather than copying, and we additionally consider measurement error.

**Adapting to the type and level of corruption**. Next, we demonstrate that our four step procedure performs well in simulations with a broad variety of data corruptions. We run the same code in every setting; the procedure adapts to the *type* and *level* of data corruption, without prior knowledge of the corruption covariance structure.
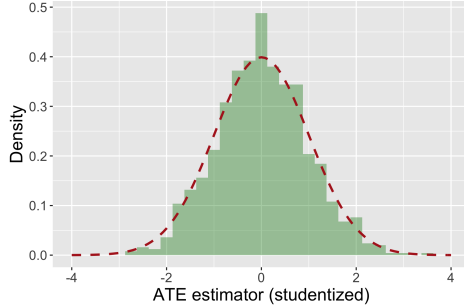
To begin, we consider measurement error $Z_{i,\cdot} = X_{i,\cdot} + H_{i,\cdot}$, where $H_{i,\cdot}$ is Gaussian noise, in the average treatment effect simulation described above. Recall that $\theta_0 = 2.2$, $\boldsymbol{X} \in \mathbb{R}^{100 \times 100}$, and $r = 5$. We implement our procedure on corrupted data 1000 times, collecting 1000 point estimates $\hat{\theta}$ and 1000 standard errors $\hat{\sigma}$. For a 20% noise-to-signal ratio, we visualize the studentized point estimates $(\hat{\theta} - \theta_0)/\hat{\sigma}$ in Figure 4a.i. Qualitatively, the histogram closely resembles the standard normal density.

We quantify performance in coverage tables. In Table 4a.ii, different rows correspond to different noise-to-signal ratios. Initially, we consider the oracle tuning of the PCA hyperparameter $k = r$. For each noise-to-signal ratio, we record the average point estimates, which are close to $\theta_0 = 2.2$. Next, we record the average standard errors, which adaptively increase in length to higher noise levels. Impressively, a 100% noise-to-signal ratio setting corresponds to a confidence interval that is only about 10% longer. These confidence intervals are the correct length, since about 950 of them include the true value $\theta_0 = 2.2$.

Table 4a.ii revisits the issue of tuning the hyperparameter $k$. This time, we fix the noise-to-signal ratio to 20%. Different rows correspond to different tunings: $k = r$, $k = r + 2$, and $k = r + 5$. Point estimates remain close to the true value $\theta_0 = 2.2$. The standard errors adaptively increase in length when $k$ deviates from $r$, though the length only increases about 10%. The confidence intervals are again the correct length, attaining nominal coverage.

We repeat this exercise with other types of data corruption: missing values (Figure 4b.i), discretization (Figure 4c.i), and differential privacy (Figure 4d.i). For missing values, $Z_{i,\cdot} = X_{i,\cdot} \odot \pi_{i,\cdot}$ and we consider non-response of 10%, 30%, and 50% of *all covariate entries*. In Census Bureau surveys, key variables such as income are missing 40% of the time. Fortunately, our procedure performs well even with this high level of missingness.

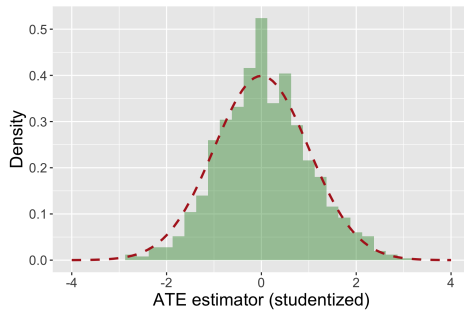For discretization, we consider randomized rounding $Z_{i,\cdot} = sign(X_{i,\cdot})\text{Poisson}(|X_{i,\cdot}|)$, which corresponds to a 33% noise-to-signal ratio. Finally, for differential privacy, $Z_{i,\cdot} = X_{i,\cdot} + H_{i,\cdot}$ where $H_{i,\cdot}$ is Laplacian noise, and we obtain results that are nearly identical to measurement error. Across settings, our results are robust to hyperparameter tuning.
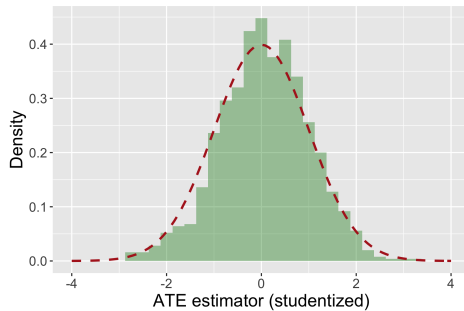


(a.i) Measurement error inference

| Meas. Err. | PCA | ATE | SE | 80% CI | 95% CI |
| --- | --- | --- | --- | --- | --- |
| 20% | k=5 | 2.22 | 0.35 | 0.81 | 0.96 |
| 60% | k=5 | 2.23 | 0.37 | 0.81 | 0.96 |
| 100% | k=5 | 2.28 | 0.39 | 0.82 | 0.95 |
| 20% | k=5 | 2.22 | 0.35 | 0.81 | 0.96 |
| 20% | k=7 | 2.21 | 0.36 | 0.84 | 0.96 |
| 20% | k=10 | 2.22 | 0.39 | 0.83 | 0.97 |

(a.ii) Measurement error coverage

(b.i) Missing values inference

| Miss. Val. | PCA | ATE | SE | 80% CI | 95% CI |
| --- | --- | --- | --- | --- | --- |
| 10% | k=5 | 2.20 | 0.35 | 0.81 | 0.96 |
| 30% | k=5 | 2.24 | 0.37 | 0.81 | 0.94 |
| 50% | k=5 | 2.35 | 0.41 | 0.79 | 0.94 |
| 10% | k=5 | 2.20 | 0.35 | 0.81 | 0.96 |
| 10% | k=7 | 2.19 | 0.37 | 0.81 | 0.95 |
| 10% | k=10 | 2.19 | 0.42 | 0.82 | 0.96 |

(b.ii) Missing values coverage

(c.i) Discretization inference

| Discret. | PCA | ATE | SE | 80% CI | 95% CI |
| --- | --- | --- | --- | --- | --- |
| 33% | k=5 | 2.23 | 0.36 | 0.81 | 0.96 |
| 33% | k=7 | 2.23 | 0.37 | 0.80 | 0.95 |
| 33% | k=10 | 2.23 | 0.41 | 0.81 | 0.95 |

(c.ii) Discretization coverage

(d.i) Differential privacy inference

| Diff. Priv. | PCA | ATE | SE | 80% CI | 95% CI |
| --- | --- | --- | --- | --- | --- |
| 20% | k=5 | 2.19 | 0.35 | 0.84 | 0.97 |
| 60% | k=5 | 2.23 | 0.37 | 0.81 | 0.96 |
| 100% | k=5 | 2.29 | 0.39 | 0.81 | 0.95 |
| 20% | k=5 | 2.19 | 0.35 | 0.84 | 0.97 |
| 20% | k=7 | 2.20 | 0.36 | 0.84 | 0.97 |
| 20% | k=10 | 2.19 | 0.39 | 0.86 | 0.97 |

(d.ii) Differential privacy coverage

Figure 4: Our approach adapts to the type and level of corruption.

# 5 Finite sample analysis

In the previous section, we articulate three reasons why inference after data cleaning is hard. First, data cleaning mixes signal and noise across observations. We introduce implicit data cleaning as an algorithmic solution, yet we still need to provide a theory of implicit data cleaning: why is it okay to not clean the test covariates? Second, the best rates of data cleaning are slower than $n^{-1/2}$. We incorporate the doubly robust estimating equation in the hope of achieving double rate robustness, yet we still need to prove that it works: how is causal inference still possible with standard errors of order $\hat{\sigma}n^{-1/2}$? Third, data cleaning recovers averages across matrix entries. How can we translate guarantees about recovering averages into guarantees about the coverage of data cleaning-adjusted confidence intervals? In this section, we answer these three theoretical questions with finite sample analysis.

We prove four theorems, each corresponding to a step in the procedure: (i) data cleaning: $\hat{X}$ converges to $X^{\text{TRAIN}}$; (ii) error-in-variable regression: $\hat{\gamma}$ converges to $\gamma_0$; (iii) error-in-variable balancing weight: $\hat{\alpha}$ converges to $\alpha_0$; (iv) causal parameter: $\mathbb{P}\{\theta_0 \in (\hat{\theta} \pm 1.96\hat{\sigma}n^{-1/2})\}$ converges to 0.95. We have already verified that our key assumption is reasonable in practice for US Census-derived data. In a corollary, we verify that it is reasonable in theory: it holds for a broad class of linear and nonlinear factor models.

**Step 1: Data cleaning**. For the data cleaning guarantee, we place four assumptions on the corrupted data. To lighten notation, we suppress indexing by TRAIN.

**Assumption 5.1** (Bounded signal). *There exists an absolute constant $\bar{A} < \infty$ such that for all $i \in [m]$ and $j \in [p]$, $|X_{ij}| \leq \bar{A}$.*

Bounded true values are standard in the matrix completion literature.

**Assumption 5.2** (Measurement error). *Each row of measurement error $H_{i,\cdot}$ is conditionally mean zero and subexponential, i.e. $\mathbb{E}[H_{i,\cdot}|X_{i,\cdot}] = 0$ and there exists $a \geq 1$ and $K_a < \infty$ such that $\|H_{i,\cdot}|X_{i,\cdot}\|_{\psi_a} \leq K_a$. Hence there exists $\kappa^2 > 0$ such that $\|\mathbb{E}[H_{i,\cdot}^T H_{i,\cdot}|X_{i,\cdot}]\|_{op} \leq \kappa^2$. We assume measurement error is independent across rows.*

Measurement error may be *dependent* within a given row. If each coordinate of $H_{i,\cdot} \in \mathbb{R}^p$ is independent, then $K_a$ and $\kappa^2$ are constants (i.e. they do not scale with $p$) [Vershynin, 2018,

Lemma 3.4.2]. More generally, $(K_a, \kappa)$ quantify the level of dependence among the entries of $H_{i,\cdot}$ within a row. Our model allows for a great deal of heteroscedasticity. In particular, the results to follow are conditional on $\boldsymbol{X}$, so the distribution of $H_{ij}$ may depend on $X_{ij}$ as long as it is conditionally mean zero and has tails no wider than those of an exponential distribution. Assumption 5.2 encompasses discretization and differential privacy.

**Assumption 5.3** (Missing values). *Each $\pi_{ij}$ is 1 with probability $\rho_j$ and NA otherwise. Identifying NA with 0, we assume there exists $\bar{K} < \infty$ such that $\|\pi_{i,\cdot} - (\rho_1, ..., \rho_p)|X_{i,\cdot}\|_{\psi_2} \leq \bar{K}$. Missingness $\pi_{i,\cdot}$ is independent of $H_{i,\cdot}$ given $X_{i,\cdot}$, and independent across rows.*

Our missingness model generalizes the standard missingness model in the PCR error-in-variable literature in two ways: (i) the missingness of one variable may depend on the missingness of another, and (ii) different variables may be missing with different probabilities. These additional degrees of flexibility are crucial for Census data, where non-responses for different variables are often correlated and where non-response rates of different variables can be vastly different. As with measurement error, missingness is independent across rows, but it may be *dependent* within a given row. If each coordinate of $\pi_{i,\cdot} \in \mathbb{R}^p$ is independent, then $\bar{K}$ is constant. More generally, $\bar{K}$ quantifies the level of dependence among the entries of $\pi_{i,\cdot}$ within a row. Our model allows for different probabilities of missingness for different variables, which may depend on the signal in a weak sense: the proof is conditional on $\boldsymbol{X}$, so the probability $\rho_j$ may depend on $X_{\cdot,j}$. We define the additional notation $\rho_{\min} := \min_{j \in [p]} \rho_j$ and $\boldsymbol{\rho} = diag(\rho_1, ..., \rho_p) \in \mathbb{R}^{p \times p}$.

**Assumption 5.4** (Concentrated signal). *Consider the approximation $\boldsymbol{X}^{(\text{LR})}$ to $\boldsymbol{X}$, with singular values $s_1, ..., s_r$. Assume that $s_1, ..., s_r \geq C\sqrt{\frac{mp}{r}}$, where $C$ is an absolute constant.*

Assumption 5.4 is analogous to incoherence-style conditions in econometrics and the notion of pervasiveness in matrix completion. Similar to a strong factor assumption, it ensures that the explanatory power of $\boldsymbol{X}^{(\text{LR})}$ dominates the explanatory power of various error terms. It requires signal to be spectrally concentrated. A natural setting in which Assumption 5.4 holds is when $X_{ij}^{(\text{LR})} = \Theta(1)$ and $s_1, ... s_r = \Theta(\tau)$. Then, for absolute constants $C, C', C'' > 0$, $C \cdot r \cdot \tau^2 = \sum_k s_k^2 = \|\boldsymbol{X}^{(\text{LR})}\|_{Fr}^2 = C' \cdot mp$ which implies $\tau = C''\sqrt{\frac{mp}{r}}$. Future work may extend our results to different spectral assumptions on $\boldsymbol{X}^{(\text{LR})}$.

**Remark 5.1.** *We parametrize our rates by the quality of low rank approximation.*

Without loss of generality, $\boldsymbol{X} = \boldsymbol{X}^{(\text{LR})} + \boldsymbol{E}^{(\text{LR})}$, where $\boldsymbol{X}^{(\text{LR})}$ is a low rank approximation to $\boldsymbol{X}$, and $\boldsymbol{E}^{(\text{LR})}$ is the approximation residual. The two key quantities are $r = rank\{\boldsymbol{X}^{(\text{LR})}\}$ and $\Delta_E = \|\boldsymbol{E}^{(\text{LR})}\|_{\max}$. It is *with* loss of generality that $r$ and $\Delta_E$ are simultaneously well behaved. Intuitively, as $r$ decreases, $\Delta_E$ increases (and vice-versa). Indeed, if $\boldsymbol{X}^{(\text{LR})} = \boldsymbol{X}$ then trivially $r \leq (m, p)$ and $\Delta_E = 0$; if $\boldsymbol{X}^{(\text{LR})} = 0$, then $r = 0$ but $\Delta_E = \bar{A}$. Our corollary shows that, under a nonlinear factor model, both $r$ and $\Delta_E$ behave well: $r \ll (m, p)$ and $\Delta_E \to 0$. Until that corollary, we parameterize rates by $(r, \Delta_E)$, which may be non-unique.[2]

**Theorem 5.1** (Finite sample data cleaning rate). *Suppose Assumptions 5.1, 5.2, 5.3, and 5.4 hold, $k = r$, and $\rho_{\min} > \frac{23 \log(mp)}{m}$. Then for an absolute constant $C > 0$,*

$$\frac{1}{m}\mathbb{E}\|\hat{\boldsymbol{X}} - \boldsymbol{X}\|_{2,\infty}^2 \leq C_1 \cdot \frac{r \ln^5(mp)}{\rho_{\min}^4}\left(\frac{1}{m} + \frac{1}{p} + \Delta_E^2\right),$$

*where $C_1 = C \cdot \bar{A}^4 (K_a + \bar{K})^2 (\kappa + K_a + \bar{K})^2$.*

The norm of convergence is called the $(2, \infty)$ norm: $\frac{1}{m}\|\hat{\boldsymbol{X}} - \boldsymbol{X}\|_{2,\infty}^2 = \max_{j \in [p]} \frac{1}{m}\|\hat{X}_{i,\cdot} - X_{\cdot,j}\|_2^2 = \max_{j \in [p]} \frac{1}{m}\sum_{i=1}^m (\hat{X}_{ij} - X_{ij})^2$ i.e. a maximum across columns and an average across rows. For any given variable $j \in [p]$, Theorem 5.1 guarantees that data cleaning performs well on average across observations $i \in [m]$. Our rate requires both $m$ and $p$ to increase: more repeated measurements improve the quality of data cleaning. For the bound to be meaningful, $(r, \Delta_E)$ must be simultaneously well behaved, which is our key assumption. Recall that $(K_a, \kappa, \bar{K})$ quantify the level of corruption dependence within a row. As long as the dependence is weak, e.g. $(K_a, \kappa, \bar{K})$ scale as some power of $\ln(mp)$, this dependence in negligible. Our downstream results for the error-in-variable regression and balancing weight build on this data cleaning guarantee. Signal is spectrally concentrated, while noise is spectrally diffuse, so we can concentrate out the noise.

**Step 2: Error-in-variable regression**. We place three additional assumptions.

**Assumption 5.5** (Response noise). *We have $\mathbb{E}[\varepsilon_i | X_{i,\cdot}] = 0$ and $\mathbb{V}[\varepsilon_i | X_{i,\cdot}] \leq \bar{\sigma}^2$. Response noise $\varepsilon_i$ is independent of $H_{i,\cdot}$ and $\pi_{i,\cdot}$ given $X_{i,\cdot}$, and independent across rows.*

---

[2]Since $r$ may be non-unique, there may be multiple valid choices of the hyperparameter $k$.

This condition permits measurement error and differential privacy of the outcome $Y_i$. Next we assume TRAIN and TEST each contains a sufficient variety of observations. For a matrix $\boldsymbol{M} \in \mathbb{R}^{m \times p}$, we define its row space as $\text{ROW}(\boldsymbol{M}) = span\{M_{i,\cdot}\}$.

**Assumption 5.6** (Row space inclusion). $\text{ROW}[b\{\boldsymbol{X}^{(LR),TRAIN}\}] = \text{ROW}[b\{\boldsymbol{X}^{(LR),TEST}\}]$.

This property permits $\boldsymbol{X}^{(\text{LR}),\text{TRAIN}} \neq \boldsymbol{X}^{(\text{LR}),\text{TEST}}$, and also permits the two matrices to have different SVDs. In Appendix H, we verify that Assumption 5.6 holds with high probability under weak auxiliary conditions. Finally, we place a weak technical condition.

**Assumption 5.7** (Well conditioned estimators). *Let $\hat{s}'_{k'}$ be the smallest non-zero singular value of $b(D^{TRAIN}, \hat{\boldsymbol{X}})$. Assume that $\hat{s}'_{k'} \gtrsim \frac{\bar{\varepsilon}}{\text{polynomial}(m,p)}$ where $\mathbb{E}[\varepsilon_i^8] \leq \bar{\varepsilon}^8$.*

For $(\hat{\beta}, \hat{\eta})$ to be well conditioned, the singular value $\hat{s}'_{k'}$ should not be too small. In particular, it must be bounded below by an arbitrary negative power of $m$ and $p$.

Before stating the result, we introduce a theoretical device $\beta^*$ as the coefficient of the best low rank nonlinear approximation to $\gamma_0$. In particular, we write $\gamma_0(D_i, X_{i,\cdot}) = b(D_i, X_{i,\cdot}^{(\text{LR})})\beta^* + \phi_i^{(\text{LR})}$ where $\phi_i^{(\text{LR})}$ is the approximation error. It will be convenient to keep track of this approximation error by defining $\phi_i := \gamma_0(D_i, X_{i,\cdot}) - b(D_i, X_{i,\cdot})\beta^*$. There will be a trade-off: a richer dictionary $b$ leads to a smaller approximation error in terms of $\|\phi\|_2^2$, but more compounding of measurement error and missingness. The following result helps to characterize how the compounded data corruption magnifies $(\rho_{\min}^{-1}, r, \Delta_E)$ but nothing else.

**Remark 5.2.** *Our results hold for a broad class of dictionaries, with the dictionary-specific constant $C'_b$ and the concise notation $(\rho'_{\min}, r', \Delta'_E)$ in Theorems 5.2 and 5.3. Appendix F proves that*

$$C'_b \leq 2^{d_{\max}} \cdot \bar{A}_{\max}^{2d_{\max}} \cdot \|\hat{\boldsymbol{X}}\|_{\max}^{2d_{\max}}, \quad \frac{1}{\rho'_{\min}} \leq \frac{d_{\max}\bar{A}^{d_{\max}}}{\rho_{\min}},$$

$$r' \leq r^{d_{\max}}, \quad and \quad \Delta'_E \leq C\bar{A}^{d_{\max}} \cdot d_{\max}\Delta_E,$$

*where $d_{\max}$ is the degree of the polynomial dictionary. Appendix F articulates restrictions on the class of dictionaries. For the interacted dictionary, $d_{\max} = 2$.*

**Remark 5.3.** *Under further incoherence-style assumptions, we bound $\|\hat{\boldsymbol{X}}\|_{\max} \leq C\sqrt{r}$ in Appendix G. Alternatively, one can bound*

$$\|\hat{\boldsymbol{X}}\|_{\max} \leq \|\hat{\boldsymbol{X}} - \boldsymbol{X}\|_{\max} + \|\boldsymbol{X}\|_{\max} \leq \|\hat{\boldsymbol{X}} - \boldsymbol{X}\|_{2,\infty}^2 + \bar{A}$$

*then appeal to Theorem 5.1 with high probability. Doing so for $C'_b$ does not affect the powers of $(m, p)$ but does increase the complexity of the pre-factors.*

**Theorem 5.2** (Finite sample error-in-variable regression rate)**.** *Suppose that the conditions of Theorem 5.1 hold, as well as Assumptions 5.5, 5.6, and 5.7. If we have that $\rho'_{\min} \gg \tilde{C} \sqrt{r'} \ln^{\frac{3}{2}}(mp) \left\{ \frac{1}{\sqrt{p}} \vee \frac{1}{\sqrt{m}} \vee \Delta_E \right\}$, where $\tilde{C} := C\bar{A}\left( \kappa + \bar{K} + K_a \right)$, then*

$$\mathcal{R}(\hat{\gamma}) \leq C'_b C_1 C_2 \cdot \bar{\sigma}^2 \cdot \frac{(r')^3 \ln^8(mp)}{(\rho'_{\min})^6} \|\beta^*\|_1^2 \left( \frac{1}{m} + \frac{p}{m^2} + \frac{1}{p} + \left(1 + \frac{p}{m}\right)(\Delta'_E)^2 + p(\Delta'_E)^4 \right)$$
$$+ C_2 \cdot \frac{(r')^2 \ln^3(mp)}{(\rho'_{\min})^2} \Delta_\phi \left(1 + (\Delta'_E)^2\right),$$

*where $\Delta_\phi = \frac{1}{m} \|\phi^{TRAIN}\|_2^2 \vee \frac{1}{m} \|\phi^{TEST}\|_2^2$,*

$$C_1 = C\bar{A}^4 (K_a + \bar{K})^2 (\kappa + \bar{K} + K_a)^2, \quad and \quad C_2 = C \cdot \bar{A}^4 (\kappa + \bar{K} + K_a)^2.$$

**Corollary 5.1** (Simplified regression rate)**.** *Suppose the conditions of Theorem 5.2 hold. Further suppose $\gamma_0$ is exactly linear in signal, which is exactly low rank. Then*

$$\mathcal{R}(\hat{\gamma}) \leq C_1 C_2 \cdot \bar{\sigma}^2 \cdot \frac{r^3 \ln^8(mp)}{\rho_{\min}^6} \|\beta^*\|_1^2 \left( \frac{1}{m} + \frac{p}{m^2} + \frac{1}{p} \right).$$

The norm of convergence is $\mathcal{R}(\hat{\gamma}) = \mathbb{E}\left[\frac{1}{m} \sum_{i \in \text{TEST}} \{\hat{\gamma}(D_i, Z_{i,\cdot}) - \gamma_0(D_i, X_{i,\cdot})\}^2\right]$, a relaxation of mean square error, where the expectation is over randomness in TRAIN and TEST. Two aspects of our problem necessitate this norm: (i) given the on-average data cleaning guarantee in Theorem 5.1, this is the best we can do; and (ii) for i.n.i.d. data, a population risk is otherwise not well defined.[3] Since the estimator $\hat{\gamma}$ does not involve cleaning TEST, Theorem 5.2 provides the theory of implicit data cleaning. The bound requires both $m$ and $p$ to increase, $p \ll m^2$, and $\rho_{\min} \gg p^{-1/2} \vee m^{-1/2} \vee \Delta_E$. For the bound to be meaningful, $(r, \Delta_E)$ must be simultaneously well behaved and the corruption dependence must be weak. Finally, the bound includes the nonlinear approximation error $\Delta_\phi$ and the size of the theoretical device $\|\beta^*\|_1$, which is well behaved if $\beta^*$ is approximately sparse. In summary, we keep track of the low rank approximation error $\Delta_E$ and the nonlinear sparse approximation error $\Delta_\phi$. To deal with $\Delta_E$, we demonstrate that nonlinear factor models admit low rank approximation below. Due to our doubly robust approach, estimation of the causal parameter $\theta_0$ is robust to non-vanishing $\Delta_\phi$—a discussion we revisit later.

---

[3]Interestingly, even with i.i.d. data, (i) necessitates this norm.

We make several innovations relative to previous work on PCR. First, we propose an error-in-variable regression estimator that does not clean the test covariates, and we develop a new theory of implicit data cleaning. Second, we define a new norm of convergence which we subsequently use in causal inference. Appendix H compares our norm with those in previous work. Third, we allow for dependence of missingness across variables and for different probabilities of missingness across variables. This flexibility is realistic for Census data. Fourth, we consider a nonlinear regression function $\gamma_0$ that is approximated by a nonlinear dictionary of basis functions $b$. The dictionary of basis functions is important for causal inference because it allows for treatment effect heterogeneity, and it requires a novel characterization of which nonlinearities do not compound data corruption too much.

**Step 3: Error-in-variable balancing**. We place one additional assumption.

**Assumption 5.8** (Row space inclusion). $\hat{M} \in \text{ROW}\{b(D^{\text{TRAIN}}, \hat{\boldsymbol{X}})\}$.

Whereas Assumption 5.6 is about the low rank approximation of the signal across TRAIN and TEST, Assumption 5.8 is about the counterfactual moment in relation to TRAIN after cleaning. With $\hat{M} = [\{b(D^{\text{TRAIN}}, \hat{\boldsymbol{X}})\}^T Y^{\text{TRAIN}}]$, which reverts to error-in-variable regression, Assumption 5.8 immediately holds. In other cases, it limits the counterfactual queries that an analyst may ask. Because it concerns empirical quantities, it may be viewed as a diagnostic tool to determine whether the counterfactual can be extrapolated.

As before, we introduce a theoretical device $\eta^*$ as the coefficient of the best low rank nonlinear approximation to $\alpha_0$. In particular, we write $\alpha_0(D_i, Z_{i,\cdot}) = b(D_i, X_{i,\cdot}^{(\text{LR})})\eta^* + \zeta_i^{(\text{LR})}$ where $\zeta_i^{(\text{LR})}$ is the approximation error, and we study this approximation error by defining $\zeta_i := \alpha_0(D_i, Z_{i,\cdot}) - b(D_i, X_{i,\cdot})\eta^*$.[4] Again, there will be a trade-off: a richer dictionary $b$ leads to a smaller approximation error in terms of $\|\zeta\|_2^2$, but amplification of $(\rho_{\min}^{-1}, r, \Delta_E)$.

**Remark 5.4.** *Our results hold for a broad class of causal parameters, with parameter-specific constants $(C'_m, C''_m)$ in Theorem 5.3. Appendix I characterizes $(C'_m, C''_m)$ for several examples. For ATE with the interacted dictionary, $C'_m = 1$ and $C''_m = \bar{A}$.*

**Theorem 5.3** (Finite sample error-in-variable balancing weight rate). *Suppose the conditions of Theorem 5.1 hold, as well as Assumptions 5.6, 5.7, and 5.8. If $\rho'_{\min} \gg$*

---

[4]A further assumption that the treatment mechanism only depends on signal, i.e. $\mathbb{E}[D_i|X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}] = \mathbb{E}[D_i|X_{i,\cdot}]$, implies $\alpha_0(D_i, Z_{i,\cdot}) = \alpha_0(D_i, X_{i,\cdot}) = b(D_i, X_{i,\cdot}^{(\text{LR})})\eta^* + \zeta_i^{(\text{LR})}$.

$\tilde{C}\sqrt{r'}\ln^{\frac{3}{2}}(mp)\left\{\frac{1}{\sqrt{p}}\vee\frac{1}{\sqrt{m}}\vee\Delta_E\right\}$ *and* $\|\alpha_0\|_\infty\leq\bar{\alpha}$, *then*

$$\mathcal{R}(\hat{\alpha})\leq C_3\cdot\frac{(r')^5\ln^{13}(mp)}{(\rho'_{\min})^{10}}\|\eta^*\|_1^2\cdot\left\{\frac{1}{m}+\frac{1}{p}+\frac{p}{m^2}+\frac{m}{p^2}+\left(1+\frac{p}{m}+\frac{m}{p}\right)(\Delta'_E)^2\right.$$
$$\left.+(m+p)(\Delta'_E)^4+mp(\Delta'_E)^6\right\}+2\Delta_\zeta,$$

*where* $\Delta_\zeta=\frac{1}{m}\|\zeta^{\text{TRAIN}}\|_2^2\vee\frac{1}{m}\|\zeta^{\text{TEST}}\|_2^2$ *and*

$$C_3=C\bar{A}^{14}(C'_b+\sqrt{C'_m}+C''_m+\bar{\alpha}+\bar{A})^2(K_a+\bar{K})^4(\kappa+\bar{K}+K_a)^6.$$

**Corollary 5.2** (Simplified balancing weight rate). *Suppose the conditions of Theorem 5.3 hold. Further suppose $\alpha_0$ is exactly linear in signal, which is exactly low rank. Then*

$$\mathcal{R}(\hat{\alpha})\leq\tilde{C}_3\cdot\frac{r^5\ln^{13}(mp)}{\rho_{\min}^{10}}\|\eta^*\|_1^2\cdot\left\{\frac{1}{m}+\frac{1}{p}+\frac{p}{m^2}+\frac{m}{p^2}\right\},$$

*where* $\tilde{C}_3$ *is* $C_3$ *but with* $C'_b$ *replaced by* 1.

The norm of convergence $\mathcal{R}(\hat{\alpha})=\mathbb{E}\left[\frac{1}{m}\sum_{i\in\text{TEST}}\{\hat{\alpha}(D_i,Z_{i,\cdot})-\alpha_0(D_i,Z_{i,\cdot})\}^2\right]$ relaxes mean square error as before. Theorem 5.3 imposes a stronger condition on $p$ than Theorem 5.2: now, we need $m^{1/2}\ll p\ll m^2$. Once again, our bound keeps track of the low rank approximation error $\Delta_E$ and the nonlinear sparse approximation error $\Delta_\zeta$. Nonlinear factor models imply that the former vanishes, and our doubly robust approach allows the latter not to vanish, as we make precise below.

Theorem 5.3 innovates in several ways. Most importantly, it analyzes a new estimator for a new estimand: the error-in-variable balancing weight in cross sectional data. A rich literature proposes balancing weight estimators for causal inference with clean data, but to our knowledge, ours is the first error-in-variable balancing weight estimator for causal inference with corrupted cross sectional data. Appendix I shows that Theorem 5.3 holds for a broad class of counterfactual moments and therefore a broad class of causal parameters.

The counterfactual moment $\hat{M}=[\{b(D^{\text{TRAIN}},\hat{X})\}^TY^{\text{TRAIN}}]$ recovers error-in-variable regression. We choose not to simply subsume Theorem 5.2 by Theorem 5.3 for two reasons. First, doing so would require that $Y_i$ and $\varepsilon_i$ are bounded, which rules out differential privacy for the outcome. Second, Theorem 5.2 has lower powers of $(r,\rho_{\min}^{-1})$ and avoids the term $\frac{m}{p^2}$ so it is typically a tighter bound. If we are willing to accept these costs, then the application

of Theorem 5.3 to error-in-variable regression relaxes $\varepsilon_i \perp\!\!\!\perp H_{i,\cdot}, \pi_{i,\cdot}|X_{i,\cdot}$ in Assumption 5.5. Both approaches allow for heteroscedasticity of $\mathbb{V}[\varepsilon_i|X_{i,\cdot}]$ in the traditional sense.

**Step 4: Causal estimation and inference**. The corrupted data problem is an extended semiparametric problem. Let $W_{i,\cdot} = (D_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})$ concatenate the signal and the noise, so that $\mathbb{L}_2(\mathcal{W})$ consists of square integrable functions of the form $\gamma : (D_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) \to \mathbb{R}$. Both the true regression $\gamma_0(D_i, X_{i,\cdot})$ and our error-in-variable estimator $\hat{\gamma}(D_i, Z_{i,\cdot})$ belong to this space, which serves as our hypothesis space for semiparametric analysis.

**Assumption 5.9** (Distribution shift). *The extended outcome and treatment mechanisms, $\mathbb{E}[Y_i|D_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}]$ and $\mathbb{E}[D_i|X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}]$, do not vary across observations.*

Assumption 5.9 implies that $\gamma_0(W_{i,\cdot})$ and $\alpha_0(W_{i,\cdot})$ do not vary across observations, though the marginal distributions $\mathbb{P}_i(W_i)$ may vary. Our corruption model implies $\gamma_0(W_{i,\cdot}) = \gamma_0(D_i, X_{i,\cdot})$, and we are agnostic about whether $\alpha_0(W_{i,\cdot}) = \alpha_0(D_i, X_{i,\cdot})$ for the extended hypothesis space.[5] Our final assumption mildly strengthens common support.

**Assumption 5.10** (Bounded propensity). *The extended propensity score is bounded below and above, i.e. $1 - \bar{\phi} \leq \mathbb{E}[D_i|X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}] \leq \bar{\phi}$.[6]*

We introduce some additional notation to state the finite sample Gaussian approximation. Define the oracle influences $\psi_i = \psi(W_{i,\cdot}, \theta_i, \gamma_0, \alpha_0)$, where the influence function is

$$\psi(W_{i,\cdot}, \theta, \gamma, \alpha) = \gamma(1, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) - \gamma(0, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) + \alpha(W_{i,\cdot})\{Y_i - \gamma(W_{i,\cdot})\} - \theta.$$

$\mathbb{E}[\psi_i] = 0$ since $\mathbb{E}[\gamma_0(1, X_{i,\cdot}) - \gamma_0(0, X_{i,\cdot})] = \theta_i$ and $\mathbb{E}[\alpha_0(W_{i,\cdot})\{Y_i - \gamma_0(W_{i,\cdot})\}] = 0$ by law of iterated expectations. We define the higher moments and average higher moments by

$$\sigma_i^2 = \mathbb{E}[\psi_i^2], \qquad\qquad \xi_i^3 = \mathbb{E}[|\psi_i|^3], \qquad\qquad \chi_i^4 = \mathbb{E}[\psi_i^4];$$

$$\sigma^2 = \frac{1}{n}\sum_{i=1}^n \sigma_i^2, \qquad\qquad \xi^3 = \frac{1}{n}\sum_{i=1}^n \xi_i^3, \qquad\qquad \chi^4 = \frac{1}{n}\sum_{i=1}^n \chi_i^4.$$

---

[5] If $\mathbb{E}[D_i|X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}] = \mathbb{E}[D_i|X_{i,\cdot}]$, then $\alpha_0(W_{i,\cdot}) = \alpha_0(D_i, X_{i,\cdot})$.

[6] Our finite sample analysis allows $\bar{\phi} \uparrow 1$, and more generally $\bar{Q} \uparrow \infty$ in Remark 5.5, as the sample size $n \uparrow \infty$.

**Remark 5.5.** *Our results hold for a broad class of causal parameters, with parameter-specific constants $(\bar{Q}, \bar{q})$ in Theorems 5.4 and 5.5. For ATE, $\bar{Q} = 2\left(\frac{1}{\phi} + \frac{1}{1-\phi}\right)$ and $\bar{q} = 1$ under Assumptions 5.9 and 5.10. Appendix J characterizes $(\bar{Q}, \bar{q})$ for several other examples under generalizations of Assumptions 5.9 and 5.10. $\bar{Q}$ may be a diverging sequence.*

**Theorem 5.4** (Finite sample Gaussian approximation). *Suppose Assumptions 5.9 and 5.10 hold, $\mathbb{V}[\varepsilon_i \mid W_{i,\cdot}] \leq \bar{\sigma}^2$, $\|\alpha_0\|_\infty \leq \bar{\alpha}$, and for $(i,j) \in$ TEST, $\hat{\gamma}(W_{i,\cdot}) \perp\!\!\!\perp \hat{\gamma}(W_{j,\cdot}) \mid$ TRAIN and $\hat{\alpha}(W_{i,\cdot}) \perp\!\!\!\perp \hat{\alpha}(W_{j,\cdot}) \mid$ TRAIN. Then with probability $1 - \epsilon$,*

$$\sup_{z \in \mathbb{R}} \left| \mathbb{P}\left\{ \frac{n^{1/2}}{\sigma}(\hat{\theta} - \theta_0) \leq z \right\} - \Phi(z) \right| \leq 0.56 \left(\frac{\xi}{\sigma}\right)^3 n^{-\frac{1}{2}} + \frac{\Delta}{(2\pi)^{1/2}} + \epsilon,$$

*where $\Phi(z)$ is the standard Gaussian distribution function and*

$$\Delta = \frac{3L}{\epsilon\sigma}\left[(\bar{Q}^{1/2} + \bar{\alpha})\{\mathcal{R}(\hat{\gamma})\}^{\bar{q}/2} + \bar{\sigma}\{\mathcal{R}(\hat{\alpha})\}^{1/2} + \{n\mathcal{R}(\hat{\gamma})\mathcal{R}(\hat{\alpha})\}^{1/2}\right].$$

**Theorem 5.5** (Finite sample variance estimation). *Suppose Assumptions 5.9 and 5.10 hold, $\mathbb{V}[\varepsilon_i \mid W_{i,\cdot}] \leq \bar{\sigma}^2$, and $\|\hat{\alpha}\|_\infty \leq \bar{\alpha}'$. Then with probability $1 - \epsilon'$,*

$$\left|\hat{\sigma}^2 - (\sigma^2 + \text{BIAS})\right| \leq \Delta' + \Delta'' + 3\left[(\Delta')^{1/2}\{(\Delta'')^{1/2} + \sigma + \Delta_{OUT}^{1/2}\}\right.$$
$$\left. + (\Delta'')^{1/2}\{\Delta_{OUT}^{1/2} + (\Delta')^{1/4}\Delta_{OUT}^{1/4}\} + (\Delta')^{1/4}\Delta_{OUT}^{1/4}\sigma\right],$$

*where*

$$\text{BIAS} = \Delta_{OUT} + 2\Delta_{OUT}^{1/2}\sigma, \quad \Delta_{OUT} = \frac{1}{n}\sum_{i=1}^{n}[(\theta_i - \theta_0)^2],$$

$$\Delta' = 4(\hat{\theta} - \theta_0)^2 + \frac{24L}{\epsilon'}\left[\{\bar{Q} + (\bar{\alpha}')^2\}\mathcal{R}(\hat{\gamma})^{\bar{q}} + \bar{\sigma}^2\mathcal{R}(\hat{\alpha})\right], \quad \text{and} \quad \Delta'' = \left(\frac{2}{\epsilon'}\right)^{1/2}\chi^2 n^{-\frac{1}{2}}.$$

**Corollary 5.3** (Confidence interval coverage). *Suppose the conditions of Theorems 5.4 and 5.5 hold. Further assume (i) moment regularity: $\{(\xi/\sigma)^3 + \chi^2\}n^{-\frac{1}{2}} \to 0$; (ii) error-in-variable regression rate: $(\bar{Q}^{1/2} + \bar{\alpha}/\sigma + \bar{\alpha}')\{\mathcal{R}(\hat{\gamma})\}^{\bar{q}/2} \to 0$; (iii) error-in-variable balancing weight rate: $\bar{\sigma}\{\mathcal{R}(\hat{\alpha})\}^{1/2} \to 0$; (iv) product of rates is fast: $\{n\mathcal{R}(\hat{\gamma})\mathcal{R}(\hat{\alpha})\}^{1/2}/\sigma \to 0$. Then $\hat{\theta} \xrightarrow{p} \theta_0$, $\hat{\sigma}^2 \xrightarrow{p} \sigma^2 + \text{BIAS}$, and $\mathbb{P}\{\theta_0 \in (\hat{\theta} \pm 1.96\hat{\sigma}n^{-1/2})\} \to 0.95 + c$ where $\text{BIAS}, c \geq 0$. If in addition $\Delta_{OUT} \to 0$, i.e. there are not too many outliers, then $\hat{\theta} \xrightarrow{p} \theta_0$, $\hat{\sigma}^2 \xrightarrow{p} \sigma^2$, and $\mathbb{P}\{\theta_0 \in (\hat{\theta} \pm 1.96\hat{\sigma}n^{-1/2})\} \to 0.95$.*

**Remark 5.6.** *Corollary 5.3 holds for a broad class of semiparametric estimands such as the average elasticity and nonparametric estimands such as heterogeneous treatment effects.*

*Moreover, it holds for not only the data cleaning and estimation procedure that we propose, but for any data cleaning and estimation procedure satisfying its weak conditions.*

The rate conditions $\mathcal{R}(\hat{\gamma}) \to 0$, $\mathcal{R}(\hat{\alpha}) \to 0$ , and $\{n\mathcal{R}(\hat{\gamma})\mathcal{R}(\hat{\alpha})\}^{1/2} \to 0$ suffice for Gaussian approximation with standard deviation $\sigma n^{-1/2}$, generalizing the main result in [Chernozhukov et al., 2023] to the harder setting with corrupted and i.n.i.d. data. These rate conditions are in terms of a more general norm than previous work because of matrix completion in the data cleaning step. Nonetheless, we recover a familiar product rate condition from semiparametric theory. The conditions solve the two remaining theoretical challenges. First, they provide a framework to translate an on-average data cleaning guarantee into a data cleaning-adjusted confidence interval for the causal parameter, by using generalized norms. Second, they ensure that the standard deviation is $\sigma n^{-1/2}$ as long as the *product* of error-in-variable rates (and hence the product of data cleaning rates) is of order $n^{-1/2}$. In summary, they allow for causal inference at rates faster than matrix completion, which is essential to achieving precision for the population while maintaining privacy for individuals.

A technical innovation is semiparametric variance estimation in the i.n.i.d. setting, which is essential to the validity of confidence intervals. We define $\Delta_{\text{OUT}}$ to quantify the frequency of outliers. Since $\theta_i = \mathbb{E}[\gamma_0(1, X_{i,\cdot}) - \gamma_0(0, X_{i,\cdot})]$, $\Delta_{\text{OUT}}$ quantifies the shift in the marginal distributions of true covariates $\mathbb{P}_i(X_{i,\cdot})$. At best, $\Delta_{\text{OUT}} = 0$ in the i.i.d. case. At worst, $\Delta_{\text{OUT}}$ is a constant (when individual treatment effects are bounded). The condition $\Delta_{\text{OUT}} \to 0$, i.e. relatively few outliers, suffices for consistent variance estimation and nominal confidence intervals. When $\Delta_{\text{OUT}} \not\to 0$, our variance estimator is asymptotically biased upwards by $\text{BIAS} = \Delta_{\text{OUT}} + 2\Delta_{\text{OUT}}^{1/2}\sigma$, implying conservative confidence intervals. At worst, our confidence intervals are valid but conservative by a theoretically quantifiable amount.

Our exact characterization of BIAS may have broader consequences for design-based inference. Future work may study properties of our procedure in randomized experiments.

Data corruption only appears in the asymptotic variance $\sigma^2$ via the error-in-variable balancing weight $\alpha_0$. In the ATE example, noise appears in the asymptotic variance when the treatment mechanism depends on both signal and noise. If the treatment mechanism depends on signal alone, then our causal estimator implemented on corrupted data is asymptotically as efficient as our causal estimator implemented on clean data.

**Key assumption holds for nonlinear factor models**. Finally, we tie together our various results and revisit our key assumption that covariates are approximately low rank. We show that nonlinear factor models (i) encode the intuition of approximate repeated measurements; (ii) imply that covariates are approximately low rank; and (iii) satisfy the rate conditions for causal inference. In a nonlinear factor model, $X_{ij} = g(\lambda_i, \mu_j)$ where $(\lambda_i, \mu_j)$ are latent factors corresponding to units and covariates, respectively. We assume that the function $g$ is smooth in its second argument, formalizing the repeated measurement intuition.

**Assumption 5.11** (Generalized factor model). *Assume $\boldsymbol{X}$ is generated as $X_{ij} = g(\lambda_i, \mu_j)$, where $\lambda_i, \mu_j \in [0,1)^q$ and $g(\lambda_i, \cdot) \in \mathcal{H}(q, S, C_H)$. Here, $\mathcal{H}(q, S, C_H)$ is the Hölder class of functions $g : [0,1)^q \to \mathbb{R}$ whose partial derivatives satisfy*

$$\sum_{s : |s| = \lfloor S \rfloor} \frac{1}{s!} |\nabla_s g(\mu) - \nabla_s g(\mu')| \leq C_H \|\mu - \mu'\|_{\max}^{S - \lfloor S \rfloor}, \quad \forall \mu, \mu' \in [0,1)^q,$$

*where $\lfloor S \rfloor$ is the largest integer below $S$.*

A linear factor model is a special case where $g(\lambda_i, \mu_j) = \lambda_i^T \mu_j$, satisfying Assumption 5.11 for all $S \in \mathbb{N}$ and some $C_H = C < \infty$. Assumption 5.11 also allows for smooth nonlinear factor models, and it implies joint control over $(r, \Delta_E)$ as desired. Intuitively, as latent dimension $q$ increases, the rank $r$ increases. As smoothness $S$ increases, the approximation error $\Delta_E$ decreases. Our final result demonstrates that, as long as the ratio $q/S$ is small enough, the data cleaning adjusted confidence intervals are valid.

**Remark 5.7.** *Our results hold for a broad class of dictionaries, with the concise notation $q'$ in Corollary 5.4. Appendix K proves that $q' \leq d_{\max} q$, where $d_{\max}$ is the degree of the polynomial dictionary. For the interacted dictionary, $d_{\max} = 2$.*

**Corollary 5.4.** *Suppose the conditions of Theorems 5.2, 5.3, 5.4 and 5.5 hold, as well as Assumption 5.11. For simplicity, consider the semiparametric case where $\sigma, \bar{\sigma}, \bar{\alpha}, \bar{\alpha}', \bar{Q}$ are bounded above and $\bar{q} = 1$. Suppose in addition (i) moment regularity: $\{(\xi/\sigma)^3 + \chi^2\}n^{-\frac{1}{2}} \to 0$; (ii) weak dependence: $(K_a, \kappa, \bar{K}, \rho_{\min}^{-1})$ scale polynomially in $\ln(np)$; (iii) nonlinear sparse approximation: $m\Delta_\phi \leq \|\beta^*\|_1^2 < \infty$ and $m\Delta_\zeta \leq \|\eta^*\|_1^2 < \infty$; (iv) enough repeated measurements: $n^{\frac{2}{3}} \lesssim p \lesssim n^{\frac{3}{2}}$, i.e. $n = p^\upsilon$ or $p = n^\upsilon$ for $\upsilon \in [1, \frac{3}{2}]$; (v) small latent dimension to smoothness ratio: $\frac{q'}{S} < \frac{3}{4} - \frac{\upsilon}{2}$. Then the conclusions of Corollary 5.3 hold.*

26

In summary, we allow either $n > p$ or $p > n$ as long as $(n, p)$ increase at similar rates. Given $(n, p)$, the ratio of the latent dimension $q$ over smoothness $S$ in the generalized factor model must be sufficiently low. For example, if $n = p$ and $q' = q$ then we require $q < \frac{S}{4}$: the latent dimension must be less than a quarter of the smoothness. A sufficiently low $\frac{q}{S}$ ratio ensures sufficiently fast learning rates $\mathcal{R}(\hat{\gamma})$ and $\mathcal{R}(\hat{\alpha})$ for causal inference with standard error $\hat{\sigma}n^{-1/2}$. For the special case of a linear factor model, the $\frac{q}{S}$ ratio constraint becomes vacuous, and there is no restriction on the latent dimension $q$. The same is true for a polynomial factor model where $g(\lambda_i, \mu_j) = \text{polynomial}(\lambda_i, \mu_j)$. The doubly robust framework allows us to slightly relax the conditions stated above and still obtain consistent estimation for $\theta_0$: either $\Delta_\phi \nrightarrow 0$ or $\Delta_\zeta \nrightarrow 0$, i.e. $\gamma_0$ or $\alpha_0$ may be incorrectly specified.

# 6 Case study: Effect of import competition using Census data

**Can we recover the same effects with data corruption**? Equipped with theoretical guarantees, we return to the motivating real world issue: measurement error, missing values, discretization, and differential privacy in US Census data. We replicate a seminal paper in labor economics [Autor et al., 2013] that uses Census-derived data to ask: what is the effect of import competition on local labor markets in the US? We ask an additional question: can we recover the same effects after introducing various types and levels of synthetic corruption? In particular, we implement differential privacy at a level calibrated to the 2020 Census. Our empirical results represent a realistic use case of Census-derived data, yet an idealized data setting where the corruptions belong to our class. The causal parameter is the partially linear instrumental variable regression parameter described in Appendix E.

[Autor et al., 2013] use Census data at the commuting zone (CZ) level. A CZ is an aggregate unit interpretable as a local economy, and 722 CZs make up the mainland US. CZ data are constructed from individual microdata published by the US government. The outcome $Y_i$ is percent change in US manufacturing employment; the treatment $D_i$ is percent change in imports from China; the instrument $U_i$ is percent change in imports from China to other countries; and the covariates $X_{i,\cdot}$ are CZ characteristics. In the augmented

specification, the covariates $X_{i,\cdot} \in \mathbb{R}^{30}$ include approximate repeated measurements such as average disability, medical, and unemployment benefits, and appear to be approximately low rank in Figure 1.

Figures 5a, 5b, and 5c present our initial semi-synthetic exercises. For reference, we visualize in red the 2SLS point estimate and confidence interval of [Autor et al., 2013], using clean data. Immediately next to [Autor et al., 2013]'s results, we visualize our own point estimate and confidence interval with clean data. We recover essentially the same point estimate and a somewhat smaller confidence interval. The true covariates are approximately low rank, our procedure exploits this fact, and therefore it has an advantage. Subsequently, we implement our procedure with increasing levels of measurement error: 20%, 40%, 60%, 80%, and 100% noise-to-signal ratio. Our point estimates remain stable, and the confidence intervals subtly increase in length. We obtain similar results with missing values and discretization: point estimates remain stable and the confidence intervals adaptively increase in length for higher noise-to-signal ratios, similar to Figure 4. Appendix L confirms similar results when standardizing the true covariates before the semi-synthetic exercises.



(a) Measurement error

(b) Missing values

(c) Discretization

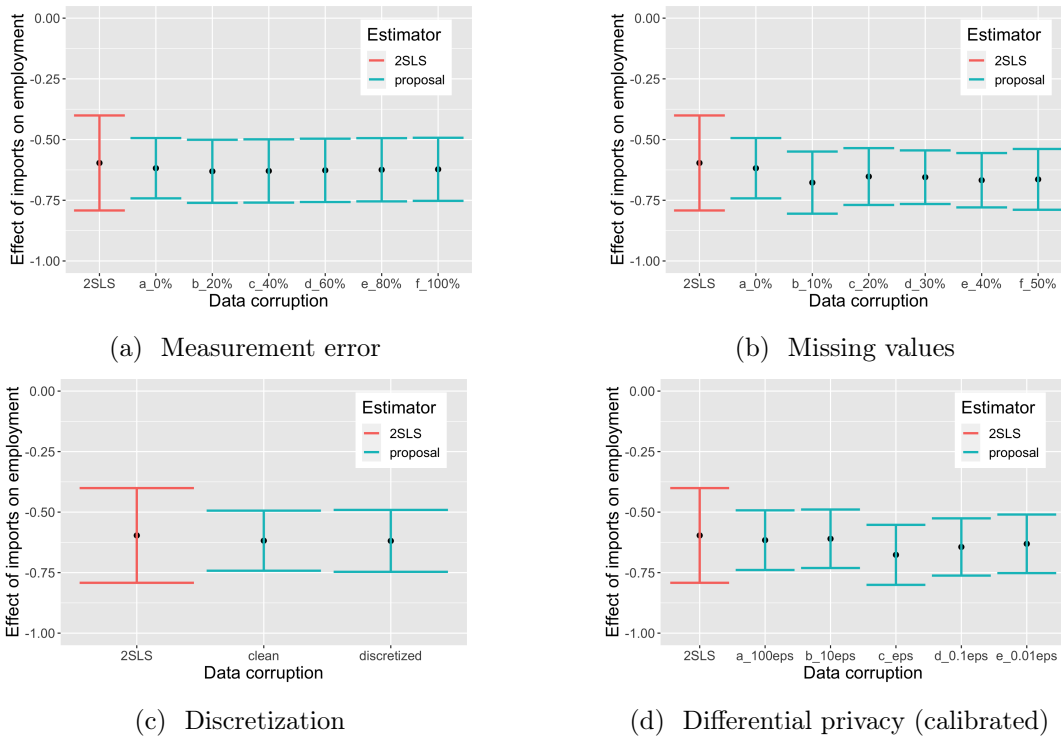(d) Differential privacy (calibrated)

Figure 5: Synthetic corruption

**Formalizing privacy**. Next, we calibrate our semi-synthetic exercise to privacy levels mandated by the US Census Bureau. To do so, we clarify how our model of causal inference

28

with corrupted data accommodates differential privacy mechanisms. With these formal results, we calibrate the variance of the Laplacian noise appropriately. In what follows, we focus on a one-off data release (formally called the non-interactive setting).

We maintain the following thought experiment: we are the Census Bureau, and our goal is to publish [Autor et al., 2013]'s CZ-level aggregated data set while protecting the privacy of individuals within CZs. In particular, we have access to the individual-level microdata, which we will *not* publicly share; we will only publish the CZ-level summaries for aggregate units. Consider a particular commuting zone $i \in [n]$ with $L_i$ individuals, and denote its individual-level microdata by $\boldsymbol{M}^{(i)} \in \mathbb{R}^{L_i \times p}$. We wish to publish $p$ summary statistics $X_{i,\cdot}$ for this CZ, where $X_{ij} = \frac{1}{L_i} \sum_{\ell=1}^{L_i} M_{\ell j}^{(i)}$, however we wish to maintain plausible deniability that each individual $\ell \in [L_i]$ contributed their data. The simulated attack on the 2010 Census found that Census block summary tables did not maintain this plausible deniability.

**Definition 6.1** (Differential privacy for summary tables)**.** *A randomized mechanism $\mathcal{M}$ confers differential privacy with privacy loss $\epsilon$ if and only if for any two possible individual-level data sets $\boldsymbol{M}$ and $\boldsymbol{M}'$ differing in a single row, and for all events $E$ in the range of $\mathcal{M}$, $\frac{\mathbb{P}(\mathcal{M}(\boldsymbol{M}) \in E)}{\mathbb{P}(\mathcal{M}(\boldsymbol{M}') \in E)} \leq e^\epsilon$ where the randomness is with respect to $\mathcal{M}$.*

The canonical mechanism that achieves differential privacy is to publish $\mathcal{M}(\boldsymbol{M}^{(i)}) = X_{i,\cdot} + H_{i,\cdot}$ instead of $X_{i,\cdot}$, where $H_{i,\cdot}$ is Laplacian noise with an appropriately calibrated variance.[7] In addition to the Laplace mechanism, the discrete Gaussian, piece wise uniform, and bounded mechanisms induce measurement error that is subexponential and conditionally mean zero, which fits within our framework. For simplicity, we focus on the Laplace mechanism when relating privacy to our theoretical results.[8]

**Proposition 6.1** (Strong protections for aggregate data)**.** *Suppose (i) each entry of microdata is bounded, i.e. $|M_{\ell j}^{(i)}| \leq \bar{A}_i$; (ii) no individual appears in two commuting zones. Then the mechanism $Z_{ij} = X_{ij} + H_{ij}$ where $H_{ij} \overset{i.i.d.}{\sim} Laplace\left(\frac{2\bar{A}_i}{\epsilon} \frac{p}{L_i}\right)$ confers $\epsilon$ differential privacy and the measurement error parameters satisfy $K_a, \kappa \leq \max_{i \in [n]} \frac{2^{3/2} \bar{A}_i}{\epsilon} \frac{p}{L_i}$. This privacy guarantee is immune to data cleaning.*

---

[7]More precisely, $\mathcal{M}_i : \mathbb{R}^{L_i \times p} \to \mathbb{R}^p$.

[8]Bureau policies mandate "zero concentrated" differential privacy, which is a closely related privacy concept for which our implementation suffices. See Appendix M for details.

**Corollary 6.1** (Safety in numbers)**.** *Suppose the conditions of Proposition 6.1 hold and* $\max_{i \in [n]} \frac{p}{L_i} \lesssim \ln(np)$. *Then the measurement error parameters satisfy* $K_a, \kappa \lesssim \ln(np)$, *and therefore our rates of data cleaning and error-in-variable estimation translate into data cleaning adjusted confidence intervals.*

In summary, the calibrated Laplacian variance depends on the privacy loss $\epsilon$, the most extreme true value $\bar{A}_i$, the number of covariates $p$, and number of individuals $L_i$ per aggregate unit. The auxiliary condition $\frac{p}{L_i} \lesssim \ln(np)$ is a practical diagnostic: roughly speaking, the number of published covariates should not greatly exceed the number of individuals per aggregate unit. It sheds light on limitations because it is plausible for CZs, but implausible for Census blocks. Much empirical economic research studies CZs, which we study in our semi-synthetic exercise. Future research may empirically investigate, through simulated attacks, how vulnerable various data releases may be for different $\frac{p}{L_i}$ regimes.

Figure 5d implements differential privacy for [Autor et al., 2013]'s CZ-level aggregated data set while protecting the privacy of individuals within CZs. We calibrate the Laplacian variance according to Proposition 6.1, where $\epsilon$ is based on Bureau memos, $p = 30$ in the augmented specification, and $(\bar{A}_i, L_i)$ are calculated from the microdata for each CZ; see Appendix M for details. To study the robustness of our results to the privacy loss parameter, we consider $(100\epsilon, 10\epsilon, \epsilon, 0.1\epsilon, 0.01\epsilon)$, which corresponds to privacy below and above the mandated level. Across levels, our point estimates and confidence intervals remain stable.

# 7 Conclusion

Recent developments in how the US Census Bureau publishes economic data motivate us to study a class of corruptions that is rich enough to encompass classical types of corruption, such as measurement error and missingness, as well as modern types, such as discretization and differential privacy mechanisms. Abstractly, our goal is to learn parameters in nonlinear, heterogeneous causal models from corrupted data; concretely, our goal is to characterize some scenarios in which it is possible to achieve both privacy and precision. To do so, we propose new data cleaning-adjusted confidence intervals that are computationally simple, statistically rigorous, and empirically robust in settings calibrated to empirical economic research. We build a framework to use matrix completion as data cleaning for downstream

causal inference, bridging two rich literatures. Future work may extend our results to confounded noise and sample selection bias.

# References

[Abowd et al., 2022] Abowd, J. M., Ashmead, R., Cumings-Menon, R., Garfinkel, S., Heineck, M., Heiss, C., Johns, R., Kifer, D., Leclerc, P., Machanavajjhala, A., Moran, B., Secton, W., Spence, M., and Zhuravlev, P. (2022). The 2020 census disclosure avoidance system TopDown algorithm. *Harvard Data Science Review*, (Special Issue 2).

[Abowd and Schmutte, 2019] Abowd, J. M. and Schmutte, I. M. (2019). An economic analysis of privacy protection and statistical accuracy as social choices. *American Economic Review*, 109(1):171–202.

[Agarwal et al., 2020a] Agarwal, A., Shah, D., and Shen, D. (2020a). On model identification and out-of-sample prediction of principal component regression: Applications to synthetic controls. *arXiv:2010.14449*.

[Agarwal et al., 2020b] Agarwal, A., Shah, D., and Shen, D. (2020b). Synthetic interventions. *arXiv:2006.07691*.

[Agarwal et al., 2021] Agarwal, A., Shah, D., Shen, D., and Song, D. (2021). On robustness of principal component regression. *Journal of the American Statistical Association*, 116(536):1731–1745.

[Ai and Chen, 2003] Ai, C. and Chen, X. (2003). Efficient estimation of models with conditional moment restrictions containing unknown functions. *Econometrica*, 71(6):1795–1843.

[Andrews, 1994] Andrews, D. W. K. (1994). Asymptotics for semiparametric econometric models via stochastic equicontinuity. *Econometrica*, 62(1):43–72.

[Athey et al., 2021] Athey, S., Bayati, M., Doudchenko, N., Imbens, G., and Khosravi, K. (2021). Matrix completion methods for causal panel data models. *Journal of the American Statistical Association*, 116(536):1716–1730.

[Athey et al., 2018] Athey, S., Imbens, G. W., and Wager, S. (2018). Approximate residual balancing: Debiased inference of average treatment effects in high dimensions. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 80(4):597–623.

[Autor et al., 2013] Autor, D. H., Dorn, D., and Hanson, G. H. (2013). The China syndrome: Local labor market effects of import competition in the United States. *American Economic Review*, 103(6):2121–2168.

[Bai, 2003] Bai, J. (2003). Inferential theory for factor models of large dimensions. *Econometrica*, 71(1):135–171.

[Bai and Ng, 2002] Bai, J. and Ng, S. (2002). Determining the number of factors in approximate factor models. *Econometrica*, 70(1):191–221.

[Bai and Ng, 2006] Bai, J. and Ng, S. (2006). Confidence intervals for diffusion index forecasts and inference for factor-augmented regressions. *Econometrica*, 74(4):1133–1150.

[Bai and Ng, 2013] Bai, J. and Ng, S. (2013). Principal components estimation and identification of static factors. *Journal of Econometrics*, 176(1):18–29.

[Bai and Ng, 2019] Bai, J. and Ng, S. (2019). Matrix completion, counterfactuals, and factor analysis of missing data. *arXiv:1910.06677*.

[Battistin and Chesher, 2014] Battistin, E. and Chesher, A. (2014). Treatment effect estimation with covariate measurement error. *Journal of Econometrics*, 178(2):707–715.

[Ben-Michael et al., 2021] Ben-Michael, E., Feller, A., Hirshberg, D. A., and Zubizarreta, J. R. (2021). The balancing act in causal inference. *arXiv:2110.14831*.

[Bickel et al., 1993] Bickel, P. J., Klaassen, C. A. J., Ritov, Y., and Wellner, J. A. (1993). *Efficient and Adaptive Estimation for Semiparametric Models*. Johns Hopkins University Press.

[Bun and Steinke, 2016] Bun, M. and Steinke, T. (2016). Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Theory of Cryptography Conference*, pages 635–658. Springer.

[Candès and Recht, 2009] Candès, E. J. and Recht, B. (2009). Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772.

[Candès and Tao, 2010] Candès, E. J. and Tao, T. (2010). The power of convex relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory*, 56(5):2053–2080.

[Chatterjee, 2015] Chatterjee, S. (2015). Matrix estimation by universal singular value thresholding. *The Annals of Statistics*, 43(1):177–214.

[Chernozhukov et al., 2018] Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W. K., and Robins, J. (2018). Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68.

[Chernozhukov et al., 2022a] Chernozhukov, V., Newey, W. K., and Singh, R. (2022a). Automatic debiased machine learning of causal and structural effects. *Econometrica*, 90(3):967–1027.

[Chernozhukov et al., 2022b] Chernozhukov, V., Newey, W. K., and Singh, R. (2022b). Debiased machine learning of global and local parameters using regularized Riesz representers. *The Econometrics Journal*, 25(3):576–601.

[Chernozhukov et al., 2023] Chernozhukov, V., Newey, W. K., and Singh, R. (2023). A simple and general debiased machine learning theorem with finite-sample guarantees. *Biometrika*, 110(1):257–264.

[Chetty and Friedman, 2019] Chetty, R. and Friedman, J. N. (2019). A practical method to reduce privacy loss when disclosing statistics based on small samples. In *AEA Papers and Proceedings*, volume 109, pages 414–20.

[Datta and Zou, 2017] Datta, A. and Zou, H. (2017). CoCoLasso for high-dimensional error-in-variables regression. *Annals of Statistics*, 45(6):2400–2426.

[Deaner, 2018] Deaner, B. (2018). Proxy controls and panel data. *arXiv:1810.00283*.

[Dwork and Lei, 2009] Dwork, C. and Lei, J. (2009). Differential privacy and robust statistics. In *Symposium on Theory of Computing*, pages 371–380.

[Dwork et al., 2006] Dwork, C., McSherry, F., Nissim, K., and Smith, A. (2006). Calibrating noise to sensitivity in private data analysis. *Theory of Cryptography*, pages 265–284.

[Feng, 2020] Feng, Y. (2020). Causal inference in possibly nonlinear factor models. *arXiv:2008.13651*.

[Fernández-Val et al., 2021] Fernández-Val, I., Freeman, H., and Weidner, M. (2021). Low-rank approximations of nonseparable panel models. *The Econometrics Journal*, 24(2):C40–C77.

[Hahn and Ridder, 2013] Hahn, J. and Ridder, G. (2013). Asymptotic variance of semi-parametric estimators with generated regressors. *Econometrica*, 81(1):315–340.

[Hasminskii and Ibragimov, 1979] Hasminskii, R. Z. and Ibragimov, I. A. (1979). On the nonparametric estimation of functionals. In *Prague Symposium on Asymptotic Statistics*, volume 473, pages 474–482.

[Hastie et al., 2015] Hastie, T., Mazumder, R., Lee, J. D., and Zadeh, R. (2015). Matrix completion and low-rank SVD via fast alternating least squares. *Journal of Machine Learning Research*, 16(104):3367–3402.

[Hausman et al., 1991] Hausman, J. A., Newey, W. K., Ichimura, H., and Powell, J. L. (1991). Identification and estimation of polynomial errors-in-variables models. *Journal of Econometrics*, 50(3):273–295.

[Hawes, 2021] Hawes, M. (2021). The Census Bureau's simulated reconstruction-abetted re-identification attack on the 2010 Census. `https://www.census.gov/data/academy/webinars/2021/disclosure-avoidance-series/simulated-reconstruction-abetted-re-identification-attack-on-the-2010-census.html`.

[Hirshberg and Wager, 2021] Hirshberg, D. A. and Wager, S. (2021). Augmented minimax linear estimation. *The Annals of Statistics*, 49(6):3206–3227.

[Hotz et al., 2022] Hotz, V. J., Bollinger, C. R., Komarova, T., Manski, C. F., Moffitt, R. A., Nekipelov, D., Sojourner, A., and Spencer, B. D. (2022). Balancing data privacy

and usability in the federal statistical system. *Proceedings of the National Academy of Sciences*, 119(31):e2104906119.

[Hu and Schennach, 2008] Hu, Y. and Schennach, S. M. (2008). Instrumental variable treatment of nonclassical measurement error models. *Econometrica*, 76(1):195–216.

[Kallus et al., 2018] Kallus, N., Mao, X., and Udell, M. (2018). Causal inference with noisy and missing covariates via matrix factorization. In *Advances in Neural Information Processing Systems*, volume 31.

[Keshavan et al., 2009] Keshavan, R., Montanari, A., and Oh, S. (2009). Matrix completion from noisy entries. In *Advances in Neural Information Processing Systems*, volume 22.

[Klaassen, 1987] Klaassen, C. A. J. (1987). Consistent estimation of the influence function of locally asymptotically linear estimators. *The Annals of Statistics*, 15(4):1548–1562.

[Komarova and Nekipelov, 2020] Komarova, T. and Nekipelov, D. (2020). Identification and formal privacy guarantees. *arXiv:2006.14732*.

[Komarova et al., 2018] Komarova, T., Nekipelov, D., and Yakovlev, E. (2018). Identification, data combination, and the risk of disclosure. *Quantitative Economics*, 9(1):395–440.

[Li and Vuong, 1998] Li, T. and Vuong, Q. (1998). Nonparametric estimation of the measurement error model using multiple indicators. *Journal of Multivariate Analysis*, 65(2):139–165.

[Loh and Wainwright, 2012] Loh, P.-L. and Wainwright, M. J. (2012). High-dimensional regression with noisy and missing data: Provable guarantees with nonconvexity. *The Annals of Statistics*, 40(3):1637–1664.

[Miao et al., 2018] Miao, W., Geng, Z., and Tchetgen Tchetgen, E. J. (2018). Identifying causal effects with proxy variables of an unmeasured confounder. *Biometrika*, 105(4):987–993.

[Newey, 1994] Newey, W. K. (1994). The asymptotic variance of semiparametric estimators. *Econometrica*, 62(6):1349–1382.

[Onatski, 2009] Onatski, A. (2009). Testing hypotheses about the number of factors in large factor models. *Econometrica*, 77(5):1447–1480.

[Robins and Rotnitzky, 1995] Robins, J. M. and Rotnitzky, A. (1995). Semiparametric efficiency in multivariate regression models with missing data. *Journal of the American Statistical Association*, 90(429):122–129.

[Robinson, 1988] Robinson, P. M. (1988). Root-$n$-consistent semiparametric regression. *Econometrica*, 56(4):931–954.

[Rosenbaum and Tsybakov, 2013] Rosenbaum, M. and Tsybakov, A. B. (2013). Improved matrix uncertainty selector. In *From Probability to Statistics and Back: High-Dimensional Models and Processes–A Festschrift in Honor of Jon A. Wellner*, volume 9, pages 276–290. Institute of Mathematical Statistics Collections.

[Rotnitzky et al., 2021] Rotnitzky, A., Smucler, E., and Robins, J. M. (2021). Characterization of parameters with a mixed bias property. *Biometrika*, 108(1):231–238.

[Rubin, 1976] Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63(3):581–592.

[Schennach, 2004] Schennach, S. M. (2004). Estimation of nonlinear models with measurement error. *Econometrica*, 72(1):33–75.

[Schennach, 2007] Schennach, S. M. (2007). Instrumental variable estimation of nonlinear errors-in-variables models. *Econometrica*, 75(1):201–239.

[Shevtsova, 2010] Shevtsova, I. G. (2010). An improvement of convergence rate estimates in the Lyapunov theorem. *Doklady Mathematics*, 82(3):862–864.

[Smith, 2011] Smith, A. (2011). Privacy-preserving statistical estimation with optimal convergence rates. In *Symposium on Theory of Computing*, pages 813–822.

[Steed et al., 2022] Steed, R., Liu, T., Wu, Z. S., and Acquisti, A. (2022). Policy impacts of statistical uncertainty and privacy. *Science*, 377(6609):928–931.

[Stock and Watson, 2002a] Stock, J. H. and Watson, M. W. (2002a). Forecasting using principal components from a large number of predictors. *Journal of the American Statistical Association*, 97(460):1167–1179.

[Stock and Watson, 2002b] Stock, J. H. and Watson, M. W. (2002b). Macroeconomic forecasting using diffusion indexes. *Journal of Business & Economic Statistics*, 20(2):147–162.

[Van der Laan and Rose, 2018] Van der Laan, M. J. and Rose, S. (2018). *Targeted Learning in Data Science.* Springer.

[Van der Laan and Rubin, 2006] Van der Laan, M. J. and Rubin, D. (2006). Targeted maximum likelihood learning. *The International Journal of Biostatistics*, 2(1).

[Vershynin, 2018] Vershynin, R. (2018). *High-Dimensional Probability: An Introduction with Applications in Data Science*, volume 47 of *Cambridge Series in Statistical and Probabilistic Mathematics.* Cambridge University Press.

[Xiong and Pelger, 2023] Xiong, R. and Pelger, M. (2023). Large dimensional latent factor modeling with missing observations and applications to causal inference. *Journal of Econometrics*, 233(1):271–301.

[Zheng and Van der Laan, 2011] Zheng, W. and Van der Laan, M. J. (2011). *Cross-Validated Targeted Minimum-Loss-Based Estimation*, pages 459–474. Springer.

# Appendix

## Table of Contents

# A    Data cleaning

In this appendix, we replace the symbol $X_{i,\cdot}$ with the symbol $A_{i,\cdot}$. We suppress indexing by the folds (TRAIN, TEST) to lighten notation. As in Assumption 5.3, we identify `NA` with $0$ in $\boldsymbol{Z}$ hereafter. We slightly abuse notation by letting $n$ be the number of observations in TRAIN, departing from the notation of the main text. The entire section is conditional on $\boldsymbol{A}$, which we omit. We write $\|\cdot\| = \|\cdot\|_{op}$ and let $C$ be an absolute constant.

Recall $\boldsymbol{A} = \boldsymbol{A}^{(\text{LR})} + \boldsymbol{E}^{(\text{LR})}$ and $r = rank\{\boldsymbol{A}^{(\text{LR})}\}$. We denote the SVDs $\boldsymbol{A}^{(\text{LR})} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^T$, $\hat{\boldsymbol{A}} = \hat{\boldsymbol{U}}_k\hat{\boldsymbol{\Sigma}}_k\hat{\boldsymbol{V}}_k^T$, and $\boldsymbol{Z}\,\hat{\boldsymbol{\rho}}^{-1} = \hat{\boldsymbol{U}}\hat{\boldsymbol{\Sigma}}\hat{\boldsymbol{V}}^T$. The first $k$ left singular vectors of $\boldsymbol{A}^{(\text{LR})}$ are $\boldsymbol{U}_k$. We denote $s_k = \Sigma_{kk}$ and $\hat{s}_k = \hat{\Sigma}_{kk}$.

Recall that $\hat{\boldsymbol{A}}^{\text{TRAIN}}$ is constructed by taking TRAIN covariates then filling and cleaning them using TRAIN alone. As a theoretical device we also study $\hat{\boldsymbol{A}}^{\text{TEST}}$, obtained by taking TEST covariates, filling them using TRAIN, and cleaning them using TEST. The analysis does not depend on whether $\hat{\boldsymbol{\rho}}^{\text{TRAIN}}$ or $\hat{\boldsymbol{\rho}}^{\text{TEST}}$ is used when filling in missing values.

Consider a matrix $\boldsymbol{B} \in \mathbb{R}^{n \times p}$ with SVD $\boldsymbol{B} = \sum_{i=1}^{n \wedge p} \sigma_i u_i v_i^T$. We define the linear function $\varphi_\lambda^{\boldsymbol{B}} : \mathbb{R}^n \to \mathbb{R}^n$ as $\varphi_\lambda^{\boldsymbol{B}}(w) = \sum_{i=1}^{n \wedge p} 1(\sigma_i \geq \lambda)u_i u_i^T w$. We use the shorthand $\varphi^{\boldsymbol{B}} = \varphi_0^{\boldsymbol{B}}$.

Define the events:

$$\mathcal{E}_1 = \left\{ \|\boldsymbol{Z} - \boldsymbol{A}\boldsymbol{\rho}\| \leq (\sqrt{n} + \sqrt{p})\Delta_{H,op} \right\}, \quad \mathcal{E}_2 = \left\{ \max_{j \in [p]} \|Z_{\cdot,j} - \rho_j A_{\cdot,j}\|_2^2 \leq n\Delta_H \right\},$$

$$\mathcal{E}_3 = \left\{ \max_{j \in [p]} \|\boldsymbol{U}_k\boldsymbol{U}_k^T(Z_{\cdot,j} - \rho_j A_{\cdot,j})\|_2^2 \leq k\Delta_H \right\}, \quad \mathcal{E}_4 = \left\{ \forall j \in [p], \frac{1}{\delta}\rho_j \leq \hat{\rho}_j \leq \delta\rho_j \right\},$$

$$\mathcal{E}_5 = \left\{ \max_{j \in [p]} |\hat{\rho}_j - \rho_j| \leq C\sqrt{\frac{\ln(np)}{n}} \right\},$$

where

$$\Delta_{H,op} = C\bar{A}(\kappa + K_a + \bar{K})\ln^{\frac{3}{2}}(np), \quad \Delta_H = C(K_a + \bar{A}\bar{K})^2\ln^2(np), \quad \delta = \frac{1}{1 - \sqrt{\frac{22\ln(np)}{n\rho_{\min}}}}.$$

The Online Appendix shows that under Assumptions 5.1, 5.2, and 5.3, $\mathbb{P}(\mathcal{E}^c) \leq \frac{10}{n^{10}p^{10}}$, where $\mathcal{E} := \cap_{k=1}^5 \mathcal{E}_k$.

39

**Lemma A.1.** *Set $k = r$. Then,*

$$\|\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1} - \boldsymbol{A}^{(\text{LR})}\| \mid \{\mathcal{E}_1, \mathcal{E}_4, \mathcal{E}_5\} \leq C\frac{\delta}{\rho_{\min}}\left((\sqrt{n} + \sqrt{p})\Delta_{H,op} + \|\boldsymbol{E}^{(\text{LR})}\| + \sqrt{\frac{\ln(np)}{n}}\|\boldsymbol{A}^{(\text{LR})}\|\right);$$

$$\|\boldsymbol{U}\boldsymbol{U}^T - \hat{\boldsymbol{U}}_r\hat{\boldsymbol{U}}_r^T\| \mid \{\mathcal{E}_1, \mathcal{E}_4, \mathcal{E}_5\} \leq C\frac{\delta}{\rho_{\min}s_r}\left((\sqrt{n} + \sqrt{p})\Delta_{H,op} + \|\boldsymbol{E}^{(\text{LR})}\| + \sqrt{\frac{\ln(np)}{n}}\|\boldsymbol{A}^{(\text{LR})}\|\right);$$

$$\|\boldsymbol{V}\boldsymbol{V}^T - \hat{\boldsymbol{V}}_r\hat{\boldsymbol{V}}_r^T\| \mid \{\mathcal{E}_1, \mathcal{E}_4, \mathcal{E}_5\} \leq C\frac{\delta}{\rho_{\min}s_r}\left\{(\sqrt{n} + \sqrt{p})\Delta_{H,op} + \|\boldsymbol{E}^{(\text{LR})}\| + \sqrt{\frac{\ln(np)}{n}}\|\boldsymbol{A}^{(\text{LR})}\|\right\};$$

$$|s_r - \hat{s}_r| \mid \{\mathcal{E}_1, \mathcal{E}_4, \mathcal{E}_5\} \leq C\frac{\delta}{\rho_{\min}}\left\{(\sqrt{n} + \sqrt{p})\Delta_{H,op} + \|\boldsymbol{E}^{(\text{LR})}\| + \sqrt{\frac{\ln(np)}{n}}\|\boldsymbol{A}^{(\text{LR})}\|\right\}.$$

*Proof.* To begin, write $\|\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1} - \boldsymbol{A}^{(\text{LR})}\| \leq \|\hat{\boldsymbol{\rho}}^{-1}\|\|\boldsymbol{Z} - \boldsymbol{A}^{(\text{LR})}\hat{\boldsymbol{\rho}}\| = \frac{\|\boldsymbol{Z}-\boldsymbol{A}^{(\text{LR})}\hat{\boldsymbol{\rho}}\|}{\min_j \hat{\rho}_j}$. By triangle inequality $\|\boldsymbol{Z} - \boldsymbol{A}^{(\text{LR})}\hat{\boldsymbol{\rho}}\| \leq \|\boldsymbol{Z} - \boldsymbol{A}^{(\text{LR})}\boldsymbol{\rho}\| + \|\boldsymbol{A}^{(\text{LR})}\|\|\boldsymbol{\rho} - \hat{\boldsymbol{\rho}}\|$. Applying $\mathcal{E}_1$ to the first term, $\mathcal{E}_5$ to the second term, and $\mathcal{E}_4$ to the denominator yields the first result. From the first result, Wedin's $\sin\Theta$ Theorem yields the second and third results while Weyl's inequality yields the fourth result. $\qquad\square$

**Lemma A.2** (Eq. 43 of [Agarwal et al., 2021])**.** *Take $\lambda^* = \hat{s}_k$. Then $\hat{A}_{\cdot,j} = \frac{1}{\hat{\rho}_j}\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(Z_{\cdot,j})$.*

**Lemma A.3.** *Suppose $k = r$, Assumptions 5.1 and 5.4 hold, and $\rho_{\min} > \frac{23\ln(np)}{n}$. Then*

$$\left\|\hat{\boldsymbol{A}} - \boldsymbol{A}\right\|_{2,\infty}^2 \mid \mathcal{E} \leq C\bar{A}^4(K_a + \bar{K})^2(\kappa + \bar{K} + K_a)^2 \cdot \frac{r\ln^5(np)}{\rho_{\min}^4}\left(1 + \frac{n}{p} + n\Delta_E^2\right).$$

*Proof.* Fix a column index $j \in [p]$. Observe that $\hat{A}_{\cdot,j} - A_{\cdot,j} = \left\{\hat{A}_{\cdot,j} - \varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(A_{\cdot,j})\right\} + \left\{\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(A_{\cdot,j}) - A_{\cdot,j}\right\}$. Recall that $\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}$ projects onto the span of the top $r$ left singular vectors $\{\hat{u}_1, ..., \hat{u}_r\}$ of $\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}$, which are are also the top $r$ left singular vectors of $\boldsymbol{Z}$ since $\hat{\boldsymbol{\rho}}^{-1}$ is diagonal. Hence $\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(A_{\cdot,j}) - A_{\cdot,j} \in span\{\hat{u}_1, \ldots, \hat{u}_r\}^\perp$. By Lemma A.2, $\hat{A}_{\cdot,j} - \varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(A_{\cdot,j}) = \frac{1}{\hat{\rho}_j}\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(Z_{\cdot,j}) - \varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(A_{\cdot,j}) \in span\{\hat{u}_1, \ldots, \hat{u}_r\}$. Therefore $\left\|\hat{A}_{\cdot,j} - A_{\cdot,j}\right\|_2^2 = \left\|\hat{A}_{\cdot,j} - \varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(A_{\cdot,j})\right\|_2^2 + \left\|\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(A_{\cdot,j}) - A_{\cdot,j}\right\|_2^2$. Again applying Lemma A.2,

$$\left\|\hat{A}_{\cdot,j} - \varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(A_{\cdot,j})\right\|_2^2 = \left\|\frac{1}{\hat{\rho}_j}\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(Z_{\cdot,j} - \rho_j A_{\cdot,j}) + \frac{\rho_j - \hat{\rho}_j}{\hat{\rho}_j}\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(A_{\cdot,j})\right\|_2^2$$

$$\leq 2\left\|\frac{1}{\hat{\rho}_j}\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(Z_{\cdot,j} - \rho_j A_{\cdot,j})\right\|_2^2 + 2\left\|\frac{\rho_j - \hat{\rho}_j}{\hat{\rho}_j}\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(A_{\cdot,j})\right\|_2^2$$

$$\leq \frac{2\delta^2}{\rho_j^2}\left\|\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(Z_{\cdot,j} - \rho_j A_{\cdot,j})\right\|_2^2 + C\frac{\delta^2}{\rho_j^2}\frac{\ln(np)}{n}\|A_{\cdot,j}\|_2^2$$

where the last line uses $\mathcal{E}_4$ and $\mathcal{E}_5$. Note that $\left\|\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(Z_{\cdot,j} - \rho_j A_{\cdot,j})\right\|_2^2 \le 2\left\|\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(Z_{\cdot,j} - \rho_j A_{\cdot,j}) - \varphi^{\boldsymbol{A}^{(\mathrm{LR})}}(Z_{\cdot,j} - \rho_j A_{\cdot,j})\right\|_2^2 + 2\left\|\varphi^{\boldsymbol{A}^{(\mathrm{LR})}}(Z_{\cdot,j} - \rho_j A_{\cdot,j})\right\|_2^2$. Since $k = r$, $\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(w) = \hat{\boldsymbol{U}}_r\hat{\boldsymbol{U}}_r^T w$ and $\varphi^{\boldsymbol{A}^{(\mathrm{LR})}}(w) = \boldsymbol{U}\boldsymbol{U}^T w$ for $w \in \mathbb{R}^n$. By Lemma A.1,

$$\left\|\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(Z_{\cdot,j} - \rho_j A_{\cdot,j}) - \varphi^{\boldsymbol{A}^{(\mathrm{LR})}}(Z_{\cdot,j} - \rho_j A_{\cdot,j})\right\|_2 \le \|\boldsymbol{U}\boldsymbol{U}^T - \hat{\boldsymbol{U}}_r\hat{\boldsymbol{U}}_r^T\|\|Z_{\cdot,j} - \rho_j A_{\cdot,j}\|_2$$
$$\le C\frac{\delta}{\rho_{\min}s_r}\left((\sqrt{n} + \sqrt{p})\Delta_{H,op} + \|\boldsymbol{E}^{(\mathrm{LR})}\| + \sqrt{\frac{\ln(np)}{n}}\|\boldsymbol{A}^{(\mathrm{LR})}\|\right)\|Z_{\cdot,j} - \rho_j A_{\cdot,j}\|_2.$$

Combining the inequalities together, we have $\left\|\hat{A}_{\cdot,j} - \varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(A_{\cdot,j})\right\|_2^2$ is bounded by

$$\frac{C\delta^4}{\rho_{\min}^2}\left(\frac{(\sqrt{n} + \sqrt{p})\Delta_{H,op}}{\rho_{\min}s_r} + \frac{\|\boldsymbol{E}^{(\mathrm{LR})}\|}{\rho_{\min}s_r} + \frac{\sqrt{\ln(np)/n}\|\boldsymbol{A}^{(\mathrm{LR})}\|}{\rho_{\min}s_r}\right)^2\|Z_{\cdot,j} - \rho_j A_{\cdot,j}\|_2^2$$
$$+ \frac{4\delta^2}{\rho_{\min}^2}\left\|\varphi^{\boldsymbol{A}^{(\mathrm{LR})}}(Z_{\cdot,j} - \rho_j A_{\cdot,j})\right\|_2^2 + C\frac{\delta^2}{\rho_{\min}^2}\frac{\ln(np)}{n}\|A_{\cdot,j}\|_2^2.$$

Since $\boldsymbol{A} = \boldsymbol{A}^{(\mathrm{LR})} + \boldsymbol{E}^{(\mathrm{LR})}$,

$$\left\|\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(A_{\cdot,j}) - A_{\cdot,j}\right\|_2^2 \le 2\left\|\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(A_{\cdot,j}^{(\mathrm{LR})}) - A_{\cdot,j}^{(\mathrm{LR})}\right\|_2^2 + 2\left\|\varphi_{\lambda^*}^{\boldsymbol{Z}\hat{\boldsymbol{\rho}}^{-1}}(E_{\cdot,j}^{(\mathrm{LR})}) - E_{\cdot,j}^{(\mathrm{LR})}\right\|_2^2$$
$$\le 2\|\boldsymbol{U}\boldsymbol{U}^T - \hat{\boldsymbol{U}}_r\hat{\boldsymbol{U}}_r^T\|^2\left\|A_{\cdot,j}^{(\mathrm{LR})}\right\|_2^2 + 2\left\|E_{\cdot,j}^{(\mathrm{LR})}\right\|_2^2$$
$$\le C\delta^2\left(\frac{(\sqrt{n} + \sqrt{p})\Delta_{H,op}}{\rho_{\min}s_r} + \frac{\|\boldsymbol{E}^{(\mathrm{LR})}\|}{\rho_{\min}s_r} + \frac{\sqrt{\ln(np)/n}\|\boldsymbol{A}^{(\mathrm{LR})}\|}{\rho_{\min}s_r}\right)^2\left\|A_{\cdot,j}^{(\mathrm{LR})}\right\|_2^2 + 2\left\|E_{\cdot,j}^{(\mathrm{LR})}\right\|_2^2,$$

where the final inequality appeals to Lemma A.1. To conclude, substitute the bounds on each term, appeal to $\mathcal{E}_2$ and $\mathcal{E}_3$, then simplify using Assumptions 5.1 and 5.4. $\qquad\square$

**Lemma A.4.** *Suppose Assumptions 5.1, 5.2, and 5.3 hold. Then* $\mathbb{E}\left[\|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{2,\infty}^2 \mathbb{1}\{\mathcal{E}^c\}\right] \le \Delta_{adv}\frac{1}{n^2p^5}$, *where* $\Delta_{adv} := C\left\{\bar{A}^2 + K_a^2\ln^2(np)\right\}$.

*Proof.* By Cauchy-Schwarz, $\mathbb{E}\left[\|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{2,\infty}^2 \mathbb{1}\{\mathcal{E}^c\}\right] \le \sqrt{\mathbb{E}\left[\|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{2,\infty}^4\right]}\sqrt{\mathbb{E}\left[\mathbb{1}^2\{\mathcal{E}^c\}\right]}$. Within the first factor, $\max_{j\in[p]}\|\hat{A}_{\cdot,j}\|_2 \le \max_{j\in[p]}\frac{1}{\hat{\rho}_j}\|Z_{\cdot,j}\|_2 \le n \cdot \sqrt{n}(\bar{A} + \max_{i,j}|H_{ij}|)$, so

$$\mathbb{E}\left[\|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{2,\infty}^4\right] \le \mathbb{E}[\{n^{\frac{3}{2}}(\bar{A} + \max_{i,j}|H_{ij}|) + \sqrt{n}\bar{A}\}^4] \le Cn^6\{\bar{A}^4 + K_a^4\ln^4(np)\}.$$

The final inequality holds since $\mathbb{E}[\max_{i,j}|H_{ij}|^4] \le CK_a^4\ln^{\frac{4}{a}}(np)$. Since $\mathbb{P}(\mathcal{E}^c) \le \frac{C}{n^{10}p^{10}}$, we conclude that

$$\mathbb{E}\left[\|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{2,\infty}^2 \mathbb{1}\{\mathcal{E}^c\}\right] \le C\sqrt{n^6(\bar{A}^4 + K_a^4\ln^4(np))}\sqrt{\frac{1}{n^{10}p^{10}}}.$$

$\qquad\square$

*Proof of Theorem 5.1.* The result follows from Lemmas A.3 and A.4. By the law of iterated expectations, results conditional on $\boldsymbol{A}$ imply the same unconditional on $\boldsymbol{A}$. $\qquad\square$

# B  Error-in-variable regression

We write the proofs without nonlinear dictionaries for clarity.

Recall that the FILL operator rescales using $\hat{\boldsymbol{\rho}}$ calculated from TRAIN. Denote the SVDs $\boldsymbol{A}^{(\text{LR}),\text{TRAIN}} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^T$, $\text{FILL}(\boldsymbol{Z}^{\text{TRAIN}}) = \boldsymbol{Z}^{\text{TRAIN}}\hat{\boldsymbol{\rho}}^{-1} = \hat{\boldsymbol{U}}\hat{\boldsymbol{\Sigma}}\hat{\boldsymbol{V}}^T$, $\hat{\boldsymbol{A}}^{\text{TRAIN}} = \hat{\boldsymbol{U}}_k\hat{\boldsymbol{\Sigma}}_k\hat{\boldsymbol{V}}_k^T$. In this notation, $\boldsymbol{V}$ is an orthonormal basis for $\text{ROW}\{\boldsymbol{A}^{(\text{LR}),\text{TRAIN}}\}$. Let $\boldsymbol{V}_\perp$ be an orthonormal basis for its orthogonal complement. Likewise we define $\hat{\boldsymbol{V}}_{k,\perp}$. Define $s_k$ and $\hat{s}_k$ as the $k$-th singular values of $\boldsymbol{A}^{(\text{LR}),\text{TRAIN}}$ and $\hat{\boldsymbol{A}}^{\text{TRAIN}}$, respectively. Next, denote the SVDs $\boldsymbol{A}^{(\text{LR}),\text{TEST}} = \boldsymbol{U}'\boldsymbol{\Sigma}'(\boldsymbol{V}')^T$, $\text{FILL}(\boldsymbol{Z}^{\text{TEST}}) = \boldsymbol{Z}^{\text{TEST}}\hat{\boldsymbol{\rho}}^{-1} = \hat{\boldsymbol{U}}'\hat{\boldsymbol{\Sigma}}'(\hat{\boldsymbol{V}}')^T$, $\hat{\boldsymbol{A}}^{\text{TEST}} = \hat{\boldsymbol{U}}'_k\hat{\boldsymbol{\Sigma}}'_k(\hat{\boldsymbol{V}}'_k)^T$. We define $\boldsymbol{V}'_\perp$ and $\hat{\boldsymbol{V}}'_\perp$ analogously to $\boldsymbol{V}_\perp$. Define $s'_k$ and $\hat{s}'_k$ as the $k$-th singular values of $\boldsymbol{A}^{(\text{LR}),\text{TEST}}$ and $\hat{\boldsymbol{A}}^{\text{TEST}}$, respectively. Finally, denote the SVD of the row-wise concatenation of $\boldsymbol{A}^{(\text{LR}),\text{TRAIN}}$ and $\boldsymbol{A}^{(\text{LR}),\text{TEST}}$ as $\tilde{\boldsymbol{U}}\tilde{\boldsymbol{\Sigma}}\tilde{\boldsymbol{V}}^T$. We define $\tilde{\boldsymbol{V}}_\perp$ analogously to $\boldsymbol{V}_\perp$ but with respect to the row-wise concatenation of $\boldsymbol{A}^{(\text{LR}),\text{TRAIN}}$ and $\boldsymbol{A}^{(\text{LR}),\text{TEST}}$.

We define $\beta^* \in \mathbb{R}^p$ as the unique solution to the following optimization problem across TRAIN and TEST: $\min_{\beta\in\mathbb{R}^p}\|\beta\|_2$ such that $\beta \in \text{argmin}\left\|\begin{bmatrix}\gamma_0(\boldsymbol{A}^{\text{TRAIN}}) \\ \gamma_0(\boldsymbol{A}^{\text{TEST}})\end{bmatrix} - \begin{bmatrix}\boldsymbol{A}^{(\text{LR}),\text{TRAIN}} \\ \boldsymbol{A}^{(\text{LR}),\text{TEST}}\end{bmatrix}\beta\right\|_2^2$.
$\beta^*$ is not the quantity of interest, but rather a theoretical device. It defines the unique, minimal-norm, low-rank, linear approximation to the regression $\gamma_0$.

Recall $Y_i = A_{i,:}^{(\text{LR})}\beta^* + \phi_i^{(\text{LR})} + \varepsilon_i$. Denote by $Y^{\text{TRAIN}} \in \mathbb{R}^n$ the concatenation of $(Y_i)_{i\in\text{TRAIN}}$. Likewise for $\varepsilon^{\text{TRAIN}}$ and $\phi^{(\text{LR}),\text{TRAIN}}$. In the argument for TRAIN ERROR, all objects correspond to TRAIN. For this reason, we suppress superscipt TRAIN when possible. For $i \in$ TEST, let $\hat{\gamma}_i = Z_{i,:}\hat{\boldsymbol{\rho}}^{-1}\hat{\beta}$ and $\gamma_i = \gamma_0(A_{i,:})$ which form the vectors $\hat{\boldsymbol{\gamma}}, \boldsymbol{\gamma}_0 \in \mathbb{R}^n$.

Define the events:

$$\tilde{\mathcal{E}}_1 := \left\{\|\hat{\boldsymbol{A}}^{\text{TEST}} - \boldsymbol{A}^{\text{TEST}}\|_{2,\infty}^2, \|\hat{\boldsymbol{A}}^{\text{TRAIN}} - \boldsymbol{A}^{(\text{LR}),\text{TRAIN}}\|_{2,\infty}^2, \leq \tilde{\Delta}_1\right\},$$

$$\tilde{\Delta}_1 := C_1 \cdot \frac{r\ln^5(np)}{\rho_{\min}^4}\left(1 + \frac{n}{p} + n\Delta_E^2\right) \quad \text{and} \quad C_1 = C\bar{A}^4(K_a + \bar{K})^2(\kappa + \bar{K} + K_a)^2;$$

$$\tilde{\mathcal{E}}_2 := \left\{\|\boldsymbol{Z}^{\text{TEST}}\hat{\boldsymbol{\rho}}^{-1} - \boldsymbol{A}^{(\text{LR}),\text{TEST}}\|^2 \leq \tilde{\Delta}_2\right\}, \quad \tilde{\Delta}_2 := C\bar{A}^2(\kappa + \bar{K} + K_a)^2\frac{\ln^3(np)}{\rho_{\min}^2}\left(n + p + np\Delta_E^2\right);$$

$$\tilde{\mathcal{E}}_3 := \left\{ \|\boldsymbol{V}\boldsymbol{V}^T - \hat{\boldsymbol{V}}_r\hat{\boldsymbol{V}}_r^T\|^2, \|\boldsymbol{V}'(\boldsymbol{V}')^T - \hat{\boldsymbol{V}}_r'(\hat{\boldsymbol{V}}_r')^T\|^2 \le \tilde{\Delta}_3 \right\}, \quad \tilde{\Delta}_3 := \frac{r}{np}\tilde{\Delta}_2;$$

$$\tilde{\mathcal{E}}_4 := \{\hat{s}_r \gtrsim s_r\};$$

$$\tilde{\mathcal{E}}_5 := \left\{ \langle \hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \varepsilon \rangle \le \tilde{\Delta}_5 \right\}, \quad \tilde{\Delta}_5 := C\bar{\sigma}^2 \ln(np)\left\{ r + \|\phi^{(\text{LR})}\|_2 + \|\beta^*\|_1(\sqrt{n}\bar{A} + \tilde{\Delta}_1^{1/2}) \right\}.$$

The Online Appendix uses the results in Appendix A to show that if the conditions of Theorem 5.1 hold and $\rho_{\min} \gg \tilde{C}\sqrt{r}\ln^{\frac{3}{2}}(np)\left(\frac{1}{\sqrt{p}} \vee \frac{1}{\sqrt{n}} \vee \Delta_E\right)$, where $\tilde{C} := C\bar{A}\left(\kappa + \bar{K} + K_a\right)$, then $\mathbb{P}(\tilde{\mathcal{E}}^c) \le \frac{C}{n^{10}p^{10}}$ where $\tilde{\mathcal{E}} := \cap_{k=1}^5 \tilde{\mathcal{E}}_k$.

**Lemma B.1.** *If Assumption 5.6 holds then $\hat{\boldsymbol{V}}_{k,\perp}^T\hat{\beta} = 0$ and $\boldsymbol{V}_\perp^T\beta^* = (\boldsymbol{V}_\perp')^T\beta^* = 0$.*

*Proof.* We generalize [Agarwal et al., 2020a, Property 4.1], using our new definition of $\beta^*$ and noting that $\text{ROW}(\boldsymbol{V}^T) = \text{ROW}\{(\boldsymbol{V}')^T\} = \text{ROW}(\tilde{\boldsymbol{V}}^T)$. $\qquad\square$

**Lemma B.2.** *Deterministically, $\|\hat{\boldsymbol{A}}\hat{\beta} - \boldsymbol{A}^{(\text{LR})}\beta^*\|_2^2 \le C\{\|\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})}\|_{2,\infty}^2\|\beta^*\|_1^2 \vee \|\phi^{(\text{LR})}\|_2^2 \vee \langle\hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \varepsilon\rangle\}$.*

*Proof.* To begin, write $\|\hat{\boldsymbol{A}}\hat{\beta} - Y\|_2^2 = \|\hat{\boldsymbol{A}}\hat{\beta} - \boldsymbol{A}^{(\text{LR})}\beta^* - \phi^{(\text{LR})}\|_2^2 + \|\varepsilon\|_2^2 - 2\langle\hat{\boldsymbol{A}}\hat{\beta} - \boldsymbol{A}^{(\text{LR})}\beta^*, \varepsilon\rangle + 2\langle\phi^{(\text{LR})}, \varepsilon\rangle$. By optimality of $\hat{\beta}$, we have $\|\hat{\boldsymbol{A}}\hat{\beta} - Y\|_2^2 \le \|\hat{\boldsymbol{A}}\beta^* - Y\|_2^2 = \|(\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})})\beta^* - \phi^{(\text{LR})}\|_2^2 + \|\varepsilon\|_2^2 - 2\langle(\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})})\beta^*, \varepsilon\rangle + 2\langle\phi^{(\text{LR})}, \varepsilon\rangle$. Combining these results, $\|\hat{\boldsymbol{A}}\hat{\beta} - \boldsymbol{A}^{(\text{LR})}\beta^* - \phi^{(\text{LR})}\|_2^2 \le \|(\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})})\beta^* - \phi^{(\text{LR})}\|_2^2 + 2\langle\hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \varepsilon\rangle$. Moreover, since $\|\hat{\boldsymbol{A}}\hat{\beta} - \boldsymbol{A}^{(\text{LR})}\beta^* - \phi^{(\text{LR})}\|_2^2 = \|\hat{\boldsymbol{A}}\hat{\beta} - \boldsymbol{A}^{(\text{LR})}\beta^*\|_2^2 + \|\phi^{(\text{LR})}\|_2^2 - 2\langle\hat{\boldsymbol{A}}\hat{\beta} - \boldsymbol{A}^{(\text{LR})}\beta^*, \phi^{(\text{LR})}\rangle$ and $\|(\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})})\beta^* - \phi^{(\text{LR})}\|_2^2 = \|(\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})})\beta^*\|_2^2 + \|\phi^{(\text{LR})}\|_2^2 - 2\langle(\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})})\beta^*, \phi^{(\text{LR})}\rangle$ we conclude that $\|\hat{\boldsymbol{A}}\hat{\beta} - \boldsymbol{A}^{(\text{LR})}\beta^*\|_2^2 \le \|(\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})})\beta^*\|_2^2 + 2\langle\hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \phi^{(\text{LR})}\rangle + 2\langle\hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \varepsilon\rangle$. By Cauchy-Schwarz and triangle inequalities, $\langle\hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \phi^{(\text{LR})}\rangle \le (\|\hat{\boldsymbol{A}}\hat{\beta} - \boldsymbol{A}^{(\text{LR})}\beta^*\|_2 + \|\hat{\boldsymbol{A}}\beta^* - \boldsymbol{A}^{(\text{LR})}\beta^*\|_2) \cdot \|\phi^{(\text{LR})}\|_2$.

In summary, for $a = \|\hat{\boldsymbol{A}}\hat{\beta} - \boldsymbol{A}^{(\text{LR})}\beta^*\|_2^2$, $b = \|\hat{\boldsymbol{A}}\beta^* - \boldsymbol{A}^{(\text{LR})}\beta^*\|_2^2$, and $c = b + 2\sqrt{b}\|\phi^{(\text{LR})}\|_2 + 2\langle\hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \varepsilon\rangle$, we have shown $a \le 2\sqrt{a}\|\phi^{(\text{LR})}\|_2 + c$, which we now analyze. Since $a \ge 0$, there are three possible cases: (i) $c \ge 0$, $2\sqrt{a}\|\phi^{(\text{LR})}\|_2 \ge c$, so $a \le 4\sqrt{a}\|\phi^{(\text{LR})}\|_2$ implies $a \le 16\|\phi^{(\text{LR})}\|_2^2$; (ii) $c \ge 0$, $2\sqrt{a}\|\phi^{(\text{LR})}\|_2 < c$, so $a \le 2c$; (iii) $c < 0$, so $a < 2\sqrt{a}\|\phi^{(\text{LR})}\|_2$ implies $a < 4\|\phi^{(\text{LR})}\|_2^2$. The three cases imply $a \le 2c \vee 16\|\phi^{(\text{LR})}\|_2^2$.

Finally, let $d := 2b + 4\sqrt{b}\|\phi^{(\text{LR})}\|_2 + 2\|\phi^{(\text{LR})}\|_2^2$. Then $d = 2\{\sqrt{b} + \|\phi^{(\text{LR})}\|_2\}^2 \le 4\{b + \|\phi^{(\text{LR})}\|_2^2\}$. Note $2c \le d + 4\langle\hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \varepsilon\rangle$. Together with the earlier results, this implies $a \le C\{b \vee \|\phi^{(\text{LR})}\|_2^2 \vee \langle\hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \varepsilon\rangle\}$. Finally note $b \le \|\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})}\|_{2,\infty}^2\|\beta^*\|_1^2$. $\qquad\square$

**Proposition B.1** (Projected TRAIN ERROR). *Suppose conditions of Theorem 5.1 hold. Further suppose Assumptions 5.5 and 5.6 hold. Let $k = r$ and $\rho_{\min} \gg \tilde{C}\sqrt{r}\ln^{\frac{3}{2}}(np)\left(\frac{1}{\sqrt{p}} \vee \frac{1}{\sqrt{n}} \vee \Delta_E\right)$. Then with probability at least $1 - O\{(np)^{-10}\}$, $\|\hat{\boldsymbol{V}}_r\hat{\boldsymbol{V}}_r^T(\hat{\beta} - \beta^*)\|_2^2$ is bounded by*

$$C\bar{A}^4(K_a + \bar{K})^2(\kappa + \bar{K} + K_a)^2\frac{\bar{\sigma}^2}{\rho_{\min}^4} \cdot r\ln^6(np) \cdot \left\{ \frac{1}{np}\|\phi^{(\text{LR})}\|_2^2 + \|\beta^*\|_1^2\left(\frac{\sqrt{n}}{\|\beta^*\|_1 np} + \frac{r}{np} + \frac{r}{p^2} + \frac{r}{p}\Delta_E^2\right)\right\}.$$

*Proof.* We show that for any $k$, $\|\hat{\boldsymbol{V}}_k\hat{\boldsymbol{V}}_k^T(\hat{\beta} - \beta^*)\|_2^2 \leq \frac{C}{\hat{s}_k^2}\left\{\|\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})}\|_{2,\infty}^2\|\beta^*\|_1^2 \vee \|\phi^{(\text{LR})}\|_2^2 \vee \langle \hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \varepsilon\rangle\right\}$. Appealing to $\mathbb{P}(\tilde{\mathcal{E}}^c) \leq \frac{C}{n^{10}p^{10}}$ yields the result. Since $\hat{\boldsymbol{V}}_k$ is an isometry, $\|\hat{\boldsymbol{V}}_k\hat{\boldsymbol{V}}_k^T(\hat{\beta} - \beta^*)\|_2^2 = \|\hat{\boldsymbol{V}}_k^T(\hat{\beta} - \beta^*)\|_2^2$. Therefore $\|\hat{\boldsymbol{A}}(\hat{\beta} - \beta^*)\|_2^2 = (\hat{\beta} - \beta^*)^T\hat{\boldsymbol{V}}_k\hat{\boldsymbol{\Sigma}}_k^2\hat{\boldsymbol{V}}_k^T(\hat{\beta} - \beta^*) \geq \hat{s}_k^2\|\hat{\boldsymbol{V}}_k^T(\hat{\beta} - \beta^*)\|_2^2$. Next, consider $\|\hat{\boldsymbol{A}}(\hat{\beta} - \beta^*)\|_2^2 \leq 2\|\hat{\boldsymbol{A}}\hat{\beta} - \boldsymbol{A}^{(\text{LR})}\beta^*\|_2^2 + 2\|\boldsymbol{A}^{(\text{LR})} - \hat{\boldsymbol{A}}\|_{2,\infty}^2\|\beta^*\|_1^2$. Combining, $\|\hat{\boldsymbol{V}}_k\hat{\boldsymbol{V}}_k^T(\hat{\beta} - \beta^*)\|_2^2 \leq \frac{2}{\hat{s}_k^2}\left\{\|\hat{\boldsymbol{A}}\hat{\beta} - \boldsymbol{A}^{(\text{LR})}\beta^*\|_2^2 + \|\boldsymbol{A}^{(\text{LR})} - \hat{\boldsymbol{A}}\|_{2,\infty}^2\|\beta^*\|_1^2\right\}$. Bound $\|\hat{\boldsymbol{A}}\hat{\beta} - \boldsymbol{A}^{(\text{LR})}\beta^*\|_2^2$ by Lemma B.2. $\qquad\square$

**Proposition B.2** (TRAIN ERROR). *Suppose conditions of Proposition B.1 hold. Then with probability at least $1 - O\{(np)^{-10}\}$, $\|\hat{\beta} - \beta^*\|_2^2$ is bounded by*

$$C\bar{A}^4(K_a + \bar{K})^2(\kappa + \bar{K} + K_a)^2\frac{\bar{\sigma}^2}{\rho_{\min}^4} \cdot r\ln^6(np) \cdot \left\{ \frac{1}{np}\|\phi^{(\text{LR})}\|_2^2 + \|\beta^*\|_2^2\left(\frac{r}{n} + \frac{r}{p} + r\Delta_E^2\right)\right\}.$$

*Proof.* We show $\|\hat{\beta} - \beta^*\|_2^2 \leq C\left[\|\boldsymbol{V}\boldsymbol{V}^T - \hat{\boldsymbol{V}}_k\hat{\boldsymbol{V}}_k^T\|^2\|\beta^*\|_2^2 + \frac{1}{\hat{s}_k^2}\left\{\|\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})}\|_{2,\infty}^2\|\beta^*\|_1^2 \vee \|\phi^{(\text{LR})}\|_2^2 \vee \langle\hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \varepsilon\rangle\right\}\right]$. Appealing to $\mathbb{P}(\tilde{\mathcal{E}}^c) \leq \frac{C}{n^{10}p^{10}}$ yields the result. Write $\|\hat{\beta} - \beta^*\|_2^2 = \|\hat{\boldsymbol{V}}_k\hat{\boldsymbol{V}}_k^T(\hat{\beta} - \beta^*)\|_2^2 + \|\hat{\boldsymbol{V}}_{k,\perp}\hat{\boldsymbol{V}}_{k,\perp}^T(\hat{\beta} - \beta^*)\|_2^2$. Proposition B.1 bounds the former. By Lemma B.1, the latter equals $\|\hat{\boldsymbol{V}}_{k,\perp}\hat{\boldsymbol{V}}_{k,\perp}^T\beta^*\|_2^2 = \|(\hat{\boldsymbol{V}}_{k,\perp}\hat{\boldsymbol{V}}_{k,\perp}^T\beta^* - \boldsymbol{V}_\perp\boldsymbol{V}_\perp^T)\beta^*\|_2^2$, which we bound by $\|\hat{\boldsymbol{V}}_{k,\perp}\hat{\boldsymbol{V}}_{k,\perp} - \boldsymbol{V}_\perp\boldsymbol{V}_\perp^T\|^2\|\beta^*\|_2^2 = \|\boldsymbol{V}\boldsymbol{V}^T - \hat{\boldsymbol{V}}_k\hat{\boldsymbol{V}}_k^T\|^2\|\beta^*\|_2^2$. $\qquad\square$

**Proposition B.3** (TEST ERROR). *Let the conditions of Theorem 5.2 hold. Then $\mathbb{E}[\|\hat{\boldsymbol{A}}^{\text{TEST}}\hat{\beta} - \boldsymbol{A}^{\text{TEST}}\beta^*\|_2^2 \mathbb{1}\{\tilde{\mathcal{E}}\}]$ is bounded by*

$$C_1 C_2 \cdot \bar{\sigma}^2 \cdot \frac{r^3\ln^8(np)}{\rho_{\min}^6}\|\beta^*\|_1^2\left\{1 + \frac{p}{n} + \frac{n}{p} + (n+p)\Delta_E^2 + np\Delta_E^4\right\} + C_2 \cdot \frac{r^2\ln^3(np)}{\rho_{\min}^2}\left(1 + \Delta_E^2\right)\|\phi^{\text{TRAIN}}\|_2^2.$$

*Proof.* We show $\|\hat{\boldsymbol{A}}^{\text{TEST}}\hat{\beta} - \boldsymbol{A}^{\text{TEST}}\beta^*\|_2^2 \leq C\sum_{m=1}^3 \Delta_m$ where $\Delta_1 := \left\{\|\boldsymbol{Z}^{\text{TEST}}\hat{\boldsymbol{\rho}}^{-1} - \boldsymbol{A}^{(\text{LR}),\text{TEST}}\|^2 + \|\boldsymbol{A}^{(\text{LR}),\text{TEST}}\|^2\|\boldsymbol{V}\boldsymbol{V}^T - \hat{\boldsymbol{V}}_r\hat{\boldsymbol{V}}_r^T\|^2\right\}\|\hat{\beta} - \beta^*\|_2^2$, $\Delta_2 := \frac{\|\boldsymbol{A}^{(\text{LR}),\text{TEST}}\|^2}{\hat{s}_r^2}\left\{\|\hat{\boldsymbol{A}}^{\text{TRAIN}} - \boldsymbol{A}^{(\text{LR}),\text{TRAIN}}\|_{2,\infty}^2\|\beta^*\|_1^2 \vee \|\phi^{(\text{LR}),\text{TRAIN}}\|_2^2 \vee \langle\hat{\boldsymbol{A}}^{\text{TRAIN}}(\hat{\beta} - \beta^*), \varepsilon\rangle\right\}$, and $\Delta_3 := \|\hat{\boldsymbol{A}}^{\text{TEST}} - \boldsymbol{A}^{\text{TEST}}\|_{2,\infty}^2\|\beta^*\|_1^2$. Appealing to $\tilde{\mathcal{E}}$ yields the result. To begin, write $\|\hat{\boldsymbol{A}}^{\text{TEST}}\hat{\beta} - \boldsymbol{A}^{\text{TEST}}\beta^*\|_2^2 \leq 2\|\hat{\boldsymbol{A}}^{\text{TEST}}(\hat{\beta} - \beta^*)\|_2^2 + 2\|(\hat{\boldsymbol{A}}^{\text{TEST}} - \boldsymbol{A}^{\text{TEST}})\beta^*\|_2^2$. We bound the latter term by matrix Hölder: $\|(\hat{\boldsymbol{A}}^{\text{TEST}} - \boldsymbol{A}^{\text{TEST}})\beta^*\|_2^2 \leq \|\hat{\boldsymbol{A}}^{\text{TEST}} - \boldsymbol{A}^{\text{TEST}}\|_{2,\infty}^2\|\beta^*\|_1^2$. In what remains, we analyze the former term using $\|\hat{\boldsymbol{A}}^{\text{TEST}}(\hat{\beta} - \beta^*)\|_2^2 \leq 2\|\{\hat{\boldsymbol{A}}^{\text{TEST}} - \boldsymbol{A}^{(\text{LR}),\text{TEST}}\}(\hat{\beta} - \beta^*)\|_2^2 + 2\|\boldsymbol{A}^{(\text{LR}),\text{TEST}}(\hat{\beta} - \beta^*)\|_2^2$.

By Weyl's inequality, we have $\|\hat{A}^{\text{TEST}} - Z^{\text{TEST}}\hat{\rho}^{-1}\| = \hat{s}'_{r+1} = \hat{s}'_{r+1} - s'_{r+1} \le \|Z^{\text{TEST}}\hat{\rho}^{-1} - A^{(\text{LR}),\text{TEST}}\|$. In turn, this gives $\|\hat{A}^{\text{TEST}} - A^{(\text{LR}),\text{TEST}}\| \le 2\|Z^{\text{TEST}}\hat{\rho}^{-1} - A^{(\text{LR}),\text{TEST}}\|$ and hence $\|\{\hat{A}^{\text{TEST}} - A^{(\text{LR}),\text{TEST}}\}(\hat{\beta} - \beta^*)\|_2^2 \le 4\|Z^{\text{TEST}}\hat{\rho}^{-1} - A^{(\text{LR}),\text{TEST}}\|^2 \cdot \|\hat{\beta} - \beta^*\|_2^2$.

Assumption 5.6 implies $(V')^T V_\perp = 0$ and hence $A^{(\text{LR}),\text{TEST}} V_\perp V_\perp^T = 0$. As a result,

$$\|A^{(\text{LR}),\text{TEST}}(\hat{\beta} - \beta^*)\|_2^2 = \|A^{(\text{LR}),\text{TEST}}(VV^T + V_\perp V_\perp^T)(\hat{\beta} - \beta^*)\|_2^2$$

$$= \|A^{(\text{LR}),\text{TEST}} VV^T(\hat{\beta} - \beta^*)\|_2^2 \le \|A^{(\text{LR}),\text{TEST}}\|^2 \|VV^T(\hat{\beta} - \beta^*)\|_2^2$$

where $\|VV^T(\hat{\beta} - \beta^*)\|_2^2 \le 2\|VV^T - \hat{V}_r \hat{V}_r^T\|^2 \|\hat{\beta} - \beta^*\|_2^2 + 2\|\hat{V}_r \hat{V}_r^T(\hat{\beta} - \beta^*)\|_2^2$. Finally appeal to the proof of Proposition B.1. $\square$

**Proposition B.4** (Implicit cleaning). *Let the conditions of Theorem 5.2 hold. Then* $\mathbb{E}[\|Z^{\text{TEST}}\hat{\rho}^{-1}\hat{\beta} - \hat{A}^{\text{TEST}}\hat{\beta}\|_2^2 \mathbb{1}\{\tilde{\mathcal{E}}\}]$ *has the same bound as Proposition B.3.*

*Proof.* We show $\|Z^{\text{TEST}}\hat{\rho}^{-1}\hat{\beta} - \hat{A}^{\text{TEST}}\hat{\beta}\|_2^2 \le C\|Z^{\text{TEST}}\hat{\rho}^{-1} - A^{(\text{LR}),\text{TEST}}\|^2 \cdot \{\|\hat{\beta} - \beta^*\|_2^2 + \|\hat{V}'_r(\hat{V}'_r)^T - V'(V')^T\|^2 \|\beta^*\|_2^2\}$. Appealing to $\tilde{\mathcal{E}}$ yields the result. Using the definitions $Z^{\text{TEST}}\hat{\rho}^{-1} = \hat{U}'\hat{\Sigma}'(\hat{V}')^T$, $\hat{A}^{\text{TEST}} = \hat{U}'_r \hat{\Sigma}'_r (\hat{V}'_r)^T$, and $\hat{A}_\perp^{\text{TEST}} = \hat{U}'_{r,\perp} \hat{\Sigma}'_{r,\perp} (\hat{V}'_{r,\perp})^T$, write $Z^{\text{TEST}}\hat{\rho}^{-1} = \hat{A}^{\text{TEST}} + \hat{A}_\perp^{\text{TEST}}$. Therefore $\|Z^{\text{TEST}}\hat{\rho}^{-1}\hat{\beta} - \hat{A}^{\text{TEST}}\hat{\beta}\|_2 \le \|\hat{U}'_{r,\perp}\| \cdot \|\hat{\Sigma}'_{r,\perp}\| \cdot \|(\hat{V}'_{r,\perp})^T\hat{\beta}\|_2 = \|\hat{\Sigma}'_{r,\perp}\| \cdot \|(\hat{V}'_{r,\perp})^T\hat{\beta}\|_2$. By Weyl's inequality, $\|\hat{\Sigma}'_{r,\perp}\| = \hat{s}'_{r+1} = \hat{s}'_{r+1} - s'_{r+1} \le \|Z^{\text{TEST}}\hat{\rho}^{-1} - A^{(\text{LR}),\text{TEST}}\|$. Moreover, $\|(\hat{V}'_{r,\perp})^T\hat{\beta}\|_2 = \|\hat{V}'_{r,\perp}(\hat{V}'_{r,\perp})^T\hat{\beta}\|_2 \le \|\hat{V}'_{r,\perp}(\hat{V}'_{r,\perp})^T(\hat{\beta} - \beta^*)\|_2 + \|\hat{V}'_{r,\perp}(\hat{V}'_{r,\perp})^T\beta^*\|_2$. Focusing on the former term, $\|\hat{V}'_{r,\perp}(\hat{V}'_{r,\perp})^T(\hat{\beta} - \beta^*)\|_2 \le \|\hat{\beta} - \beta^*\|_2$. By Lemma B.1, the latter term equals

$$\|\hat{V}'_{r,\perp}(\hat{V}'_{r,\perp})^T V'(V')^T\beta^*\|_2 \le \|\{\hat{V}'_{r,\perp}(\hat{V}'_{r,\perp})^T - V'_\perp(V'_\perp)^T\}V'(V')^T\beta^*\|_2 + \|V'_\perp(V'_\perp)^T V'(V')^T\beta^*\|_2$$

$$= \|\{\hat{V}'_{r,\perp}(\hat{V}'_{r,\perp})^T - V'_\perp(V'_\perp)^T\}V'(V')^T\beta^*\|_2 = \|\{\hat{V}'_{r,\perp}(\hat{V}'_{r,\perp})^T - V'_\perp(V'_\perp)^T\}\beta^*\|_2$$

$$= \|\{\hat{V}'_r(\hat{V}'_r)^T - V'(V')^T\}\beta^*\|_2 \le \|\hat{V}'_r(\hat{V}'_r)^T - V'(V')^T\| \|\beta^*\|_2,$$

which is dominated by the former term by the proof of Proposition B.2. $\square$

*Proof of Theorem 5.2.* To begin, write $\mathbb{E}\|\hat{\gamma} - \gamma_0\|_2^2 \le 2\mathbb{E}\big[\|Z^{\text{TEST}}\hat{\rho}^{-1}\hat{\beta} - A^{\text{TEST}}\beta^*\|_2^2 \mathbb{1}\{\tilde{\mathcal{E}}\}\big] + 2\mathbb{E}\big[\|Z^{\text{TEST}}\hat{\rho}^{-1}\hat{\beta} - A^{\text{TEST}}\beta^*\|_2^2 \mathbb{1}\{\tilde{\mathcal{E}}^c\}\big] + 2\|\phi^{\text{TEST}}\|_2^2$. Consider the first term. Using the bound

$$\mathbb{E}\big[\|Z^{\text{TEST}}\hat{\rho}^{-1}\hat{\beta} - A^{\text{TEST}}\beta^*\|_2^2 \mathbb{1}\{\tilde{\mathcal{E}}\}\big] \le 2\mathbb{E}\big[\|Z^{\text{TEST}}\hat{\rho}^{-1}\hat{\beta} - \hat{A}^{\text{TEST}}\hat{\beta}\|_2^2 \mathbb{1}\{\tilde{\mathcal{E}}\}\big]$$

$$+ 2\mathbb{E}\big[\|\hat{A}^{\text{TEST}}\hat{\beta} - A^{\text{TEST}}\beta^*\|_2^2 \mathbb{1}\{\tilde{\mathcal{E}}\}\big],$$

we may appeal to Propositions B.3 and B.4. Consider the second term. Write

$$\mathbb{E}\Big[\|\boldsymbol{Z}^{\text{TEST}}\,\hat{\boldsymbol{\rho}}^{-1}\,\hat{\beta} - \boldsymbol{A}^{\text{TEST}}\beta^*\|_2^2\,\mathbb{1}\{\tilde{\mathcal{E}}^c\}\Big] \leq 2\mathbb{E}\Big[\|\boldsymbol{Z}^{\text{TEST}}\,\hat{\boldsymbol{\rho}}^{-1}\,\hat{\beta}\|_2^2\,\mathbb{1}\{\tilde{\mathcal{E}}^c\}\Big] + 2\mathbb{E}\Big[\|\boldsymbol{A}^{\text{TEST}}\beta^*\|_2^2\,\mathbb{1}\{\tilde{\mathcal{E}}^c\}\Big].$$

Since $\|\boldsymbol{A}^{\text{TEST}}\beta^*\|_2^2 \leq \|\boldsymbol{A}^{\text{TEST}}\|_{2,\infty}^2\|\beta^*\|_1^2 \leq n\bar{A}^2\|\beta^*\|_1^2$, we have that $\mathbb{E}\Big[\|\boldsymbol{A}^{\text{TEST}}\beta^*\|_2^2\,\mathbb{1}\{\tilde{\mathcal{E}}^c\}\Big] \leq n\bar{A}^2\|\beta^*\|_1^2\mathbb{P}(\tilde{\mathcal{E}}^c)$. By Cauchy-Schwarz inequality,

$$\mathbb{E}\big[\|\boldsymbol{Z}^{\text{TEST}}\,\hat{\boldsymbol{\rho}}^{-1}\,\hat{\beta}\|_2^2\,\mathbb{1}\{\tilde{\mathcal{E}}^c\}\big] \leq \sqrt{\mathbb{E}\big[\|\boldsymbol{Z}^{\text{TEST}}\,\hat{\boldsymbol{\rho}}^{-1}\,\hat{\beta}\|_2^4\big]}\,\sqrt{\mathbb{P}(\tilde{\mathcal{E}}^c)}.$$

These expressions are dominated by the bound from Proposition B.3 under Assumption 5.7.

$\square$

# C   Error-in-variable balancing weight

As before, we write the proofs without nonlinear dictionaries for clarity.

We define $\eta^* \in \mathbb{R}^p$ as the unique solution to the following optimization problem across TRAIN and TEST: $\min_{\eta\in\mathbb{R}^p}\|\eta\|_2$ such that $\eta \in \text{argmin}\left\|\begin{bmatrix}\alpha_0(\boldsymbol{W}^{\text{TRAIN}})\\\alpha_0(\boldsymbol{W}^{\text{TEST}})\end{bmatrix} - \begin{bmatrix}\boldsymbol{A}^{(\text{LR}),\text{TRAIN}}\\\boldsymbol{A}^{(\text{LR}),\text{TEST}}\end{bmatrix}\eta\right\|_2^2$. $\eta^*$ is not the quantity of interest, but rather a theoretical device. It defines the unique, minimal-norm, low-rank, linear approximation to the balancing weight $\alpha_0$.

$\hat{M}$ is the counterfactual moment. Write $\hat{\boldsymbol{G}} = \frac{1}{n}(\hat{\boldsymbol{A}}^{\text{TRAIN}})^T\hat{\boldsymbol{A}}^{\text{TRAIN}}$ as the covariance matrix after data cleaning. In this notation, $\hat{\boldsymbol{G}}\hat{\eta} = \hat{M}^T$, and these feasible objects are computed from TRAIN. We analogously define $M^*$ and $\boldsymbol{G}^*$ using the low rank approximation to the signal. In this notation, $\boldsymbol{G}^*\eta^* = (M^*)^T$, and these infeasible objects are defined from TRAIN and TEST. See the Online Appendix for a more formal statement. Finally, let $\Delta_{RR} := n \cdot \Big\{\|\hat{M}^T - (M^*)^T\|_{\max} + \|\boldsymbol{G}^* - \hat{\boldsymbol{G}}\|_{\max}\|\eta^*\|_1\Big\}$.

Define the event: $\tilde{\mathcal{E}}_5 := \Big\{\Delta_{RR} \leq \tilde{\Delta}_5\Big\}$ where $\tilde{\Delta}_5 := C\bar{A}^5(\sqrt{C_m'}+C_m''+\bar{\alpha}+\bar{A})\frac{(K_a+\bar{K})^2(\kappa+\bar{K}+K_a)^2}{\rho_{\min}^4}r\cdot \ln^5(np) \cdot n\|\eta^*\|_1\Big(\frac{1}{n} + \frac{1}{p} + \frac{n}{p^2} + \Delta_E^2 + n\Delta_E^4\Big)^{\frac{1}{2}}$. Set $\tilde{\mathcal{E}} := \cap_{k=1}^5\tilde{\mathcal{E}}_k$ where the remaining events are defined in Appendix B. The Online Appendix shows that if the conditions of Theorem 5.3 hold, then $\mathbb{P}(\tilde{\mathcal{E}}^c) \leq \frac{C}{n^{10}p^{10}}$.

**Lemma C.1.** *If Assumptions 5.6 and 5.8 hold, $\hat{\boldsymbol{V}}_{k,\perp}^T\hat{\eta} = 0$ and $\boldsymbol{V}_\perp^T\eta^* = (\boldsymbol{V}_\perp')^T\eta^* = 0$.*

*Proof.* For the former result, $\hat{\eta}$ is the unique solution to the program $\min_{\eta\in\mathbb{R}^p}\|\eta\|_2$ such that $\eta \in \text{argmin} -2\hat{M}\eta + \eta^T\hat{\boldsymbol{G}}\eta$ where $\hat{M} \in \text{ROW}(\hat{\boldsymbol{A}}^{\text{TRAIN}})$ by Assumption 5.8 and $\text{ROW}(\hat{\boldsymbol{G}}) =$

$\text{ROW}\{(\hat{\boldsymbol{A}}^{\text{TRAIN}})^T \hat{\boldsymbol{A}}^{\text{TRAIN}}\} = \text{ROW}(\hat{\boldsymbol{A}}^{\text{TRAIN}})$. Therefore $\hat{\eta} \in \text{ROW}(\hat{\boldsymbol{A}}^{\text{TRAIN}})$, so we can appeal to the same reasoning as Lemma B.1. The latter result is similar. $\qquad\square$

**Lemma C.2.** $\|\hat{\boldsymbol{A}}\hat{\eta} - \boldsymbol{A}^{(\text{LR})}\eta^*\|_2^2 \leq C\big\{\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta} - \eta^*)\|_1 \cdot \Delta_{RR} \vee \|\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})}\|_{2,\infty}^2 \|\eta^*\|_1^2\big\}.$

*Proof.* To begin, write $\|\hat{\boldsymbol{A}}\hat{\eta} - \boldsymbol{A}^{(\text{LR})}\eta^*\|_2^2 \leq 2\|\hat{\boldsymbol{A}}(\hat{\eta} - \eta^*)\|_2^2 + 2\|(\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})})\eta^*\|_2^2$. We bound the latter term by $\|\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})}\|_{2,\infty}^2 \|\eta^*\|_1^2$. Hereafter, we focus on the former term:

$$\frac{1}{n}\|\hat{\boldsymbol{A}}(\hat{\eta} - \eta^*)\|_2^2 = \frac{1}{n}(\hat{\eta} - \eta^*)^T \hat{\boldsymbol{A}}^T \hat{\boldsymbol{A}}(\hat{\eta} - \eta^*) = (\hat{\eta} - \eta^*)^T \hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T \hat{\boldsymbol{G}}(\hat{\eta} - \eta^*)$$

which is bounded by $\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta} - \eta^*)\|_1 \cdot \|\hat{\boldsymbol{G}}(\hat{\eta} - \eta^*)\|_{\max}$. Using the first order conditions,

$$\|\hat{\boldsymbol{G}}(\hat{\eta} - \eta^*)\|_{\max} \leq \|\hat{\boldsymbol{G}}\hat{\eta} - \hat{M}^T\|_{\max} + \|\hat{M}^T - (M^*)^T\|_{\max} + \|(M^*)^T - \boldsymbol{G}^*\eta^*\|_{\max} + \|\boldsymbol{G}^*\eta^* - \hat{\boldsymbol{G}}\eta^*\|_{\max}$$

$$= 0 + \|\hat{M}^T - (M^*)^T\|_{\max} + 0 + \|\boldsymbol{G}^* - \hat{\boldsymbol{G}}\|_{\max}\|\eta^*\|_1.\square$$

**Proposition C.1** (Projected TRAIN ERROR). *Suppose the conditions of Theorem 5.1 hold. Further suppose Assumptions 5.6, 5.8, 5.9, and 5.10 hold, and $\|\alpha_0\|_\infty \leq \bar{\alpha}$. Let $k = r$ and $\rho_{\min} \gg \tilde{C}\sqrt{r}\ln^{\frac{3}{2}}(np)\big(\frac{1}{\sqrt{p}} \vee \frac{1}{\sqrt{n}} \vee \Delta_E\big)$, where $\tilde{C} := C\bar{A}\big(\kappa + \bar{K} + K_a\big)$. Then with probability at least $1 - O\{(np)^{-10}\}$, $\|\hat{\boldsymbol{V}}_r \hat{\boldsymbol{V}}_r^T(\hat{\eta} - \eta^*)\|_2^2$ is bounded by*

$$C\bar{A}^{10}(\sqrt{C_m'} + C_m'' + \bar{\alpha} + \bar{A})^2(K_a + \bar{K})^4(\kappa + \bar{K} + K_a)^4 \cdot r^4 \cdot \ln^{10}(np) \cdot \frac{\|\eta^*\|_1^2}{\rho_{\min}^8}\Big(\frac{1}{np} + \frac{1}{p^2} + \frac{n}{p^3} + \frac{1}{p}\Delta_E^2 + \frac{n}{p}\Delta_E^4\Big).$$

*Proof.* We show that for any $k$, $\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta}-\eta^*)\|_2^2 \leq C\big\{\frac{1}{\hat{s}_k^2}\|\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})}\|_{2,\infty}^2\|\eta^*\|_1^2 + \frac{1}{\hat{s}_k^4}p \cdot \Delta_{RR}^2\big\}$. Appealing to $\mathbb{P}(\tilde{\mathcal{E}}^c) \leq \frac{C}{n^{10}p^{10}}$ yields the result. As in the proof of Proposition B.1, $\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta} - \eta^*)\|_2^2 \leq \frac{2}{\hat{s}_k^2}\big\{\|\hat{\boldsymbol{A}}\hat{\eta} - \boldsymbol{A}^{(\text{LR})}\eta^*\|_2^2 + \|\boldsymbol{A}^{(\text{LR})} - \hat{\boldsymbol{A}}\|_{2,\infty}^2\|\eta^*\|_1^2\big\}$. Using Lemma C.2, we conclude that $\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta} - \eta^*)\|_2^2 \leq \frac{C}{\hat{s}_k^2}\big\{\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta} - \eta^*)\|_1 \cdot \Delta_{RR} \vee \|\boldsymbol{A}^{(\text{LR})} - \hat{\boldsymbol{A}}\|_{2,\infty}^2\|\eta^*\|_1^2\big\}$. There are two cases. In the first case, $\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta} - \eta^*)\|_2^2 \leq C\frac{1}{\hat{s}_k^2}\|\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})}\|_{2,\infty}^2\|\eta^*\|_1^2$, giving the first term. In the second case, $\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta} - \eta^*)\|_2^2 \leq C\frac{1}{\hat{s}_k^2}\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta} - \eta^*)\|_1 \cdot \Delta_{RR}$. Bounding $\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta} - \eta^*)\|_1 \leq \sqrt{p}\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta} - \eta^*)\|_2$, dividing both sides by $\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta} - \eta^*)\|_2$, and squaring gives the second term. $\qquad\square$

**Proposition C.2** (TRAIN ERROR). *Suppose the conditions of Proposition C.1 hold. Then with probability at least $1 - O\{(np)^{-10}\}$, $\|\hat{\eta} - \eta^*\|_2^2$ is bounded by*

$$C\bar{A}^{10}(\sqrt{C_m'} + C_m'' + \bar{\alpha} + \bar{A})^2(K_a + \bar{K})^4(\kappa + \bar{K} + K_a)^4 \cdot r^4 \cdot \ln^{10}(np) \cdot \frac{\|\eta^*\|_2^2}{\rho_{\min}^8}\Big(\frac{1}{n} + \frac{1}{p} + \frac{n}{p^2} + \Delta_E^2 + n\Delta_E^4\Big).$$

*Proof.* We show $\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta} - \eta^*)\|_2^2 \le C\left\{\frac{1}{\hat{s}_k^2}\|\hat{\boldsymbol{A}} - \boldsymbol{A}^{(\text{LR})}\|_{2,\infty}^2\|\eta^*\|_1^2 + \frac{1}{\hat{s}_k^4}p \cdot \Delta_{RR}^2\right\}$. Appealing to $\mathbb{P}(\tilde{\mathcal{E}}^c) \le \frac{C}{n^{10}p^{10}}$ yields the result. The argument is similar to Proposition B.2, replacing Lemma B.1 with Lemma C.1 and Proposition B.1 with Proposition C.1. $\qquad\square$

**Proposition C.3** (TEST ERROR). *Let the conditions of Theorem 5.3 hold. Then* $\mathbb{E}[\|\hat{\boldsymbol{A}}^{\text{TEST}}\hat{\eta} - \boldsymbol{A}^{\text{TEST}}\eta^*\|_2^2\, \mathbb{1}\{\tilde{\mathcal{E}}\}]$ *is bounded by*

$$C_3 \cdot \frac{r^5 \ln^{13}(np)}{\rho_{\min}^{10}} \cdot \|\eta^*\|_1^2 \left\{1 + \frac{p}{n} + \frac{n}{p} + \frac{n^2}{p^2} + \left(n + p + \frac{n^2}{p}\right)\Delta_E^2 + (np + n^2)\Delta_E^4 + n^2 p \Delta_E^6\right\}.$$

*Proof.* The proof is similar to Proposition B.3, updating $\Delta_2 = \frac{\|\boldsymbol{A}^{(\text{LR}),\text{TEST}}\|^2}{\hat{s}_r^2}\left\{\|\hat{\boldsymbol{V}}_k \hat{\boldsymbol{V}}_k^T(\hat{\eta} - \eta^*)\|_1 \cdot \Delta_{RR} \vee \|\boldsymbol{A}^{(\text{LR}),\text{TRAIN}} - \hat{\boldsymbol{A}}^{\text{TRAIN}}\|_{2,\infty}^2 \|\eta^*\|_1^2\right\}$. In particular, we bound $\|\hat{\boldsymbol{V}}_r \hat{\boldsymbol{V}}_r^T(\hat{\eta} - \eta^*)\|_2^2$ using Proposition C.1 instead of Proposition B.1. $\qquad\square$

**Proposition C.4** (Implicit cleaning). *Let the conditions of Theorem 5.3 hold. Then* $\mathbb{E}[\|\boldsymbol{Z}^{\text{TEST}}\hat{\boldsymbol{\rho}}^{-1}\hat{\eta} - \hat{\boldsymbol{A}}^{\text{TEST}}\hat{\eta}\|_2^2\, \mathbb{1}\{\tilde{\mathcal{E}}\}]$ *has the same bound as Proposition C.3.*

*Proof.* The proof is analogous to Proposition B.4. $\qquad\square$

*Proof of Theorem 5.3.* The proof is analogous to Theorem 5.2, instead appealing to Propositions C.3 and C.4. $\qquad\square$

# D  Data cleaning-adjusted confidence intervals

Let $\psi_i = \psi(W_{i,\cdot}, \theta_i, \gamma_0, \alpha_0)$ where $\psi(w, \theta, \gamma, \alpha) = m(w, \gamma) + \alpha(w)\{y - \gamma(w)\} - \theta$ and $\gamma \mapsto m(w, \gamma)$ is linear. We take as given that $(\gamma_0, \alpha_0)$ exist, though the latter is implied by Assumption 5.10. The Gateaux derivative of $\psi(w, \theta, \gamma, \alpha)$ with respect to its argument $\gamma$ in the direction $u$ is $\{\partial_\gamma \psi(w, \theta, \gamma, \alpha)\}(u) = \frac{\partial}{\partial \tau}\psi(w, \theta, \gamma + \tau u, \alpha)\Big|_{\tau=0}$.

Let $L$ be the number of folds, with fold $\ell$ indexed by $I_\ell$. Train $(\hat{\gamma}_\ell, \hat{\alpha}_\ell)$ on observations in $I_\ell^c$, which serves as TRAIN. Let $m = |I_\ell| = n/L$ be the number of observations in $I_\ell$, which serves as TEST. Denote by $\mathbb{E}_\ell[\cdot]$ the average over observations in $I_\ell$. This generalized notation allows us to reverse the roles of TRAIN and TEST, and to allow for more than two folds. The target and oracle are $\hat{\theta} = L^{-1}\sum_{\ell=1}^L \mathbb{E}_\ell[m(W_{i,\cdot}, \hat{\gamma}_\ell) + \hat{\alpha}_\ell(W_{i,\cdot})\{Y_i - \hat{\gamma}_\ell(W_{i,\cdot})\}]$ and $\bar{\theta} = L^{-1}\sum_{\ell=1}^L \mathbb{E}_\ell[m(W_{i,\cdot}, \gamma_0) + \alpha_0(W_{i,\cdot})\{Y_i - \gamma_0(W_{i,\cdot})\}]$. For $i \in I_\ell$, let $\hat{\psi}_i = \psi(W_{i,\cdot}, \hat{\theta}, \hat{\gamma}_\ell, \hat{\alpha}_\ell)$.

**Lemma D.1.** *Suppose Assumptions 5.9 and 5.10 hold, $\mathbb{E}[\varepsilon_i^2|W_{i,\cdot}] \leq \bar{\sigma}^2$, $\|\alpha_0\|_\infty \leq \bar{\alpha}$, and that for $(i,j) \in I_\ell$, $\hat{\gamma}_\ell(W_{i,\cdot}) \perp\!\!\!\perp \hat{\gamma}_\ell(W_{j,\cdot})|I_\ell^c$ and $\hat{\alpha}_\ell(W_{i,\cdot}) \perp\!\!\!\perp \hat{\alpha}_\ell(W_{j,\cdot})|I_\ell^c$. Then with probability $1 - \epsilon$, $\frac{n^{1/2}}{\sigma}|\hat{\theta} - \bar{\theta}| \leq \Delta = \frac{3L}{\epsilon\sigma}\big[(\bar{Q}^{1/2} + \bar{\alpha})\{\mathcal{R}(\hat{\gamma}_\ell)\}^{\bar{q}/2} + \bar{\sigma}\{\mathcal{R}(\hat{\alpha}_\ell)\}^{1/2} + \{n\mathcal{R}(\hat{\gamma}_\ell)\mathcal{R}(\hat{\alpha}_\ell)\}^{1/2}\big].*

*Proof.* We generalize [Chernozhukov et al., 2023, Proposition S6] to the new norm, noting that Jensen's inequality and $\bar{q} \in (0,1]$ imply that $\mathbb{E}_\ell[\mathbb{E}[m(W_{i,\cdot},u)^2|I_\ell^c]] \leq \bar{Q}\mathbb{E}_\ell[\{\mathbb{E}[u(W_i)^2|I_\ell^c]\}^{\bar{q}}] \leq \bar{Q}\{\mathbb{E}_\ell[\mathbb{E}[u(W_i)^2|I_\ell^c]]\}^{\bar{q}}$. Moreover, using the shorthand $u_i = u(W_{i,\cdot})$ and $v_i = v(W_{i,\cdot})$, $\mathbb{E}[\mathbb{E}_\ell\{|u(W_{i,\cdot})v(W_{i,\cdot})|\}] = \frac{1}{m}\mathbb{E}[u^Tv] \leq \frac{1}{m}(\mathbb{E}[u^Tu])^{1/2}(\mathbb{E}[v^Tv])^{1/2} = \big(\frac{1}{m}\mathbb{E}[\mathbb{E}_\ell[u_i^2]]\big)^{1/2}\big(\frac{1}{m}\mathbb{E}[\mathbb{E}_\ell[v_i^2]]\big)^{1/2} = \sqrt{\mathcal{R}(\hat{\gamma}_\ell)}\sqrt{\mathcal{R}(\hat{\alpha}_\ell)}$. □

*Proof of Theorem 5.4.* We generalize [Chernozhukov et al., 2023, Theorem 1] using Lemma D.1 and an i.n.i.d. Berry Esseen lemma for $\bar{\theta} - \theta_0 = \mathbb{E}_n\psi_i$ [Shevtsova, 2010]. □

**Lemma D.2.** *Suppose Assumptions 5.9 and 5.10 hold, $\mathbb{E}[\varepsilon_i^2|W_{i,\cdot}] \leq \bar{\sigma}^2$, and $\|\hat{\alpha}_\ell\|_\infty \leq \bar{\alpha}'$. Then with probability $1 - \epsilon'/2$, $\mathbb{E}_n\{(\hat{\psi}_i - \psi_i + \theta_0 - \theta_i)^2\} \leq \Delta' = 4(\hat{\theta} - \theta_0)^2 + \frac{24L}{\epsilon'}\big[\{\bar{Q} + (\bar{\alpha}')^2\}\mathcal{R}(\hat{\gamma}_\ell)^{\bar{q}} + \bar{\sigma}^2\mathcal{R}(\hat{\alpha}_\ell)\big]$.*

*Proof.* We generalize [Chernozhukov et al., 2023, Proposition S10] to the new norm, appealing to Jensen's inequality and $\bar{q} \in (0,1]$ as in the proof of Lemma D.1. □

**Lemma D.3.** *With probability $1 - \epsilon'/2$, $|\mathbb{E}_n(\psi_i^2) - \sigma^2| \leq \Delta'' = \big(\frac{2}{\epsilon'}\big)^{1/2}\frac{\chi^2}{n^{1/2}}$.*

*Proof.* Let $B_i = \psi_i^2$ and $\bar{B} = \mathbb{E}_n[B_i]$ so that $\mathbb{E}[\bar{B}] = \frac{1}{n}\sum_{i=1}^n \mathbb{E}[B_i] = \frac{1}{n}\sum_{i=1}^n \sigma_i^2 = \sigma^2$ and $\mathbb{V}[\bar{B}] = \frac{\sum_{i=1}^n \mathbb{V}(B_i)}{n^2} \leq \frac{\sum_{i=1}^n \mathbb{E}[B_i^2]}{n^2} = \frac{\sum_{i=1}^n \chi_i^4}{n^2} = \frac{\chi^4}{n}$. By Markov inequality $\mathbb{P}(|\bar{B} - \mathbb{E}[\bar{B}]| > t) \leq \frac{\mathbb{V}[\bar{B}]}{t^2} = \frac{\epsilon'}{2}$. Solving gives $t = \Delta''$. □

*Proof of Theorem 5.5.* To begin, write $\hat{\sigma}^2 - (\sigma^2 + \text{BIAS}) = \{\hat{\sigma}^2 - \mathbb{E}_n(\psi_i^2) - \text{BIAS}\} + \{\mathbb{E}_n(\psi_i^2) - \sigma^2\} \leq \{\hat{\sigma}^2 - \mathbb{E}_n(\psi_i^2) - \text{BIAS}\} + \Delta''$ where the inequality holds with probability $1 - \epsilon'/2$ by Lemma D.3. In what follows, we focus on the former term. In particular, we write $\hat{\sigma}^2 = \mathbb{E}_n\{(\hat{\psi}_i - \psi_i)^2\} + 2\mathbb{E}_n\{(\hat{\psi}_i - \psi_i)\psi_i\} + \mathbb{E}_n(\psi_i^2)$, then solve for $\text{BIAS} = \text{BIAS}_1 + \text{BIAS}_2$ as a function of $\Delta_{\text{OUT}}$ in the decomposition $\hat{\sigma}^2 - \mathbb{E}_n(\psi_i^2) - \text{BIAS} = \mathbb{E}_n\{(\hat{\psi}_i - \psi_i)^2\} - \text{BIAS}_1 + 2\mathbb{E}_n\{(\hat{\psi}_i - \psi_i)\psi_i\} - \text{BIAS}_2$. To derive $\text{BIAS}_1$, open the square and write

$$\mathbb{E}_n\{(\hat{\psi}_i - \psi_i)^2\} = \mathbb{E}_n\{(\hat{\psi}_i - \psi_i + \theta_0 - \theta_i)^2\} + \mathbb{E}_n\{(\theta_i - \theta_0)^2\} + 2\mathbb{E}_n\{(\hat{\psi}_i - \psi_i + \theta_0 - \theta_i)(\theta_i - \theta_0)\}$$

$$\leq \mathbb{E}_n\{(\hat{\psi}_i - \psi_i + \theta_0 - \theta_i)^2\} + \mathbb{E}_n\{(\theta_i - \theta_0)^2\} + 2[\mathbb{E}_n\{(\hat{\psi}_i - \psi_i + \theta_0 - \theta_i)^2\}]^{1/2}[\mathbb{E}_n\{(\theta_i - \theta_0)^2\}]^{1/2}$$

$$\leq \Delta' + \Delta_{\text{OUT}} + 2(\Delta')^{1/2}\Delta_{\text{OUT}}^{1/2}$$

49

where the last line holds with probability $1 - \epsilon'/2$ by Lemma D.2. Taking $\text{BIAS}_1 = \Delta_{\text{OUT}}$, we have shown $\mathbb{E}_n\{(\hat{\psi}_i - \psi_i)^2\} - \text{BIAS}_1 \leq \Delta' + 2(\Delta')^{1/2}\Delta_{\text{OUT}}^{1/2}$. To derive $\text{BIAS}_2$, write

$$\mathbb{E}_n\{(\hat{\psi}_i - \psi_i)\psi_i\} \leq \left[\mathbb{E}_n\{(\hat{\psi}_i - \psi_i)^2\}\right]^{1/2}\{|\mathbb{E}_n(\psi_i^2) - \sigma^2| + \sigma^2\}^{1/2}$$

$$\leq \{\Delta' + \Delta_{\text{OUT}} + 2(\Delta')^{1/2}\Delta_{\text{OUT}}^{1/2}\}^{1/2} \cdot \{\Delta'' + \sigma^2\}^{1/2}$$

where the last line holds with probability $1 - \epsilon'$ appealing to Lemmas D.2 and D.3 as well as the analysis for $\text{BIAS}_1$. In summary,

$$2\mathbb{E}_n\{(\hat{\psi}_i - \psi_i)\psi_i\} \leq 2\{\Delta' + \Delta_{\text{OUT}} + 2(\Delta')^{1/2}\Delta_{\text{OUT}}^{1/2}\}^{1/2} \cdot \{\Delta'' + \sigma^2\}^{1/2}$$

$$\leq 2\{(\Delta')^{1/2} + \Delta_{\text{OUT}}^{1/2} + 2^{1/2}(\Delta')^{1/4}\Delta_{\text{OUT}}^{1/4}\} \cdot \{(\Delta'')^{1/2} + \sigma\}.$$

Taking $\text{BIAS}_2 = 2\Delta_{\text{OUT}}^{1/2}\sigma$, we have shown

$$2\mathbb{E}_n\{(\hat{\psi}_i - \psi_i)\psi_i\} - \text{BIAS}_2 \leq 2(\Delta')^{1/2}\{(\Delta'')^{1/2} + \sigma\} + 2\Delta_{\text{OUT}}^{1/2}(\Delta'')^{1/2} + 2^{3/2}(\Delta')^{1/4}\Delta_{\text{OUT}}^{1/4}\{(\Delta'')^{1/2} + \sigma\}.$$

Thus with probability $1 - \epsilon'$, $\hat{\sigma}^2 - (\sigma^2 + \text{BIAS}) \leq \{\hat{\sigma}^2 - \mathbb{E}_n(\psi_i^2) - \text{BIAS}\} + \Delta''$, equalling

$$\mathbb{E}_n\{(\hat{\psi}_i - \psi_i)^2\} - \text{BIAS}_1 + 2\mathbb{E}_n\{(\hat{\psi}_i - \psi_i)\psi_i\} - \text{BIAS}_2 + \Delta''$$

$$\leq \Delta' + 2(\Delta')^{1/2}\Delta_{\text{OUT}}^{1/2} + 2(\Delta')^{1/2}\{(\Delta'')^{1/2} + \sigma\} + 2\Delta_{\text{OUT}}^{1/2}(\Delta'')^{1/2} + 2^{3/2}(\Delta')^{1/4}\Delta_{\text{OUT}}^{1/4}\{(\Delta'')^{1/2} + \sigma\} + \Delta''.$$

Combining terms yields the desired result. $\qquad\square$

# E   Additional examples

**Semiparametric estimands**. We consider causal parameters of the form $\theta_0 = \frac{1}{n}\sum_{i=1}^{n}\theta_i$, where $\theta_i = \mathbb{E}[m(W_{i,\cdot}, \gamma_0)]$, in an i.n.i.d. data generating process where $Y_i = \gamma_0(D_i, X_{i,\cdot}) + \varepsilon_i$, $Z_{i,\cdot} = (X_{i,\cdot} + H_{i,\cdot}) \odot \pi_{i,\cdot}$, and $W_{i,\cdot} = (D_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})$. $(D_i, X_{i,\cdot})$ concatenate the various arguments of $\gamma_0$, which we hereby call regressors. This model includes the scenario in which some variables are corrupted and other are not. Which regressors are corrupted or uncorrupted constrains the construction of technical regressors; see Appendix F. We concatentate signal and noise as $W_{i,\cdot}$. Appendix J generalizes Assumption 5.9 to impose invariance of the regression $\gamma_0$ and generalized balancing weight $\alpha_0$ across observations.

**Example E.1** (Average treatment effect). *Let $(D_i, X_{i,\cdot})$ concatenate treatment $D_i \in \{0, 1\}$ and covariates $X_{i,\cdot} \in \mathbb{R}^p$. Denote $\gamma_0(D_i, X_{i,\cdot}) := \mathbb{E}[Y_i | D_i, X_{i,\cdot}]$. Under the assumption of*

selection on $X_{i,\cdot}$, the average treatment effect is given by $\theta_i = \mathbb{E}[\gamma_0(1, X_{i,\cdot}) - \gamma_0(0, X_{i,\cdot})]$. With uncorrupted treatment and corrupted covariates, $W_{i,\cdot} = (D_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})$ where $(H_{i,\cdot}, \pi_{i,\cdot})$ are measurement error and missingness for the covariates.[9]

While the true regression $\gamma_0(D_i, X_{i,\cdot})$ is only a function of signal $(D_i, X_{i,\cdot})$, our regression estimator $\hat{\gamma}(D_i, Z_{i,\cdot})$ is a function of both signal and noise $W_{i,\cdot}$. In other words, the hypothesis space for estimation is the extended space of functions $\mathbb{L}_2(\mathcal{W})$, and we must define an extended functional over $\mathbb{L}_2(\mathcal{W})$. In Example E.1, the extended functional is $\gamma \mapsto \mathbb{E}[\gamma(1, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) - \gamma(0, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})]$.

**Example E.2** (Local average treatment effect)**.** *Let $(U_i, X_{i,\cdot})$ concatenate instrument $U_i \in \{0, 1\}$ and covariates $X_{i,\cdot} \in \mathbb{R}^p$. Denote $\gamma_0(U_i, X_{i,\cdot}) := \mathbb{E}[Y_i | U_i, X_{i,\cdot}]$ and $\delta_0(U_i, X_{i,\cdot}) := \mathbb{E}[D_i | U_i, X_{i,\cdot}]$. Under standard instrumental variable assumptions, the local average treatment effect for the subpopulation of compliers is given by $\beta_0 = \frac{\theta_0}{\theta_0'}$ where $\theta_i = \mathbb{E}[\gamma_0(1, X_{i,\cdot}) - \gamma_0(0, X_{i,\cdot})]$ and $\theta_i' = \mathbb{E}[\delta_0(1, X_{i,\cdot}) - \delta_0(0, X_{i,\cdot})]$. With uncorrupted instrument and corrupted covariates, $W_{i,\cdot} = (U_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})$ where $(H_{i,\cdot}, \pi_{i,\cdot})$ are measurement error and missingness for the covariates.*

**Example E.3** (Average policy effect)**.** *Let $X_{i,\cdot} \in \mathbb{R}^p$ be the covariates. Consider the counterfactual transportation of covariates $x_{i,\cdot} \mapsto t(x_{i,\cdot})$. Denote $\gamma_0(X_{i,\cdot}) := \mathbb{E}[Y_i | X_{i,\cdot}]$. The average policy effect is $\theta_i = \mathbb{E}[\gamma_0\{t(X_{i,\cdot})\} - \gamma_0(X_{i,\cdot})]$. With corrupted covariates, $W_{i,\cdot} = (X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})$ where $(H_{i,\cdot}, \pi_{i,\cdot})$ are measurement error and missingness for covariates.*

**Example E.4** (Price elasticity of demand)**.** *Let $Y_i$ be price of a particular good. Let $(D_i, X_{i,\cdot})$ concatenate quantities sold of the particular good $D_i$ and other goods $X_{i,\cdot} \in \mathbb{R}^p$. Denote $\gamma_0(D_i, X_{i,\cdot}) = \mathbb{E}[Y_i | D_i, X_{i,\cdot}]$. The average price elasticity of demand is $\theta_i = \mathbb{E}[\nabla_d \gamma_0(D_i, X_{i,\cdot})]$. With uncorrupted quantity for the particular good and corrupted quantities for the other goods, $W_{i,\cdot} = (D_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})$ where $(H_{i,\cdot}, \pi_{i,\cdot})$ are measurement error and missingness for the other goods.*

**Weighted estimands**. In empirical economic research with aggregate units, it is common to weight units by their size. It is also common to consider partially linear models.

---

[9]More generally, treatment observations may be corrupted as well. For readability, we exposit the simpler and plausible case that treatment is uncorrupted.

For example, the estimand of [Autor et al., 2013] may be viewed as a weighted partially linear instrumental variable regression. To bridge theory with practice, we provide these examples next. A weighted functional $\theta_0 \in \mathbb{R}$ is a scalar that takes the form $\theta_0 = \frac{1}{n} \sum_{i=1}^n \theta_i$ where $\theta_i = \mathbb{E}[\ell_i m(W_{i,\cdot}, \gamma_0)]$ and $\ell_i$ is the weight for aggregate unit $i$. For simplicity, we take the weights $\ell_i$ to be known, but their uncertainty can be incorporated as well.

**Example E.5** (Weighted partially linear regression). *Let $(D_i, X_i)$ concatenate treatment $D \in \mathbb{R}$ and covariates $X_{i,\cdot} \in \mathbb{R}^p$. Denote $\gamma_0(D_i, X_{i,\cdot}) = \mathbb{E}[Y_i | D_i, X_{i,\cdot}]$. The weighted partially regression coefficient is given by $\theta_i = \mathbb{E}[\ell_i \{\gamma_0(d+1, X_{i,\cdot}) - \gamma_0(d, X_{i,\cdot})\}]$. With uncorrupted treatment and corrupted covariates, $W_{i,\cdot} = (D_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})$ where $(H_{i,\cdot}, \pi_{i,\cdot})$ are measurement error and missingness for the covariates.*

**Example E.6** (Weighted partially linear instrumental variable regression). *Let $(U_i, X_{i,\cdot})$ concatenate instrument $U_i \in \mathbb{R}$ and covariates $X_{i,\cdot} \in \mathbb{R}^p$. Denote $\gamma_0(U_i, X_{i,\cdot}) := \mathbb{E}[Y_i | U_i, X_{i,\cdot}]$ and $\delta_0(U_i, X_{i,\cdot}) := \mathbb{E}[D_i | U_i, X_{i,\cdot}]$. Under standard instrumental variable assumptions, the weighted partially linear instrumental variable regression coefficient is given by $\beta_0 = \frac{\theta_0}{\theta_0'}$, where $\theta_i = \mathbb{E}[\ell_i \{\{\gamma_0(u+1, X_{i,\cdot}) - \gamma_0(u, X_{i,\cdot})\}]$ and $\theta_i' = \mathbb{E}[\ell_i \{\delta_0(u+1, X_{i,\cdot}) - \delta_0(u, X_{i,\cdot})\}]$. With uncorrupted instrument and corrupted covariates, $W_{i,\cdot} = (U_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})$ where $(H_{i,\cdot}, \pi_{i,\cdot})$ are measurement error and missingness for the covariates.*

**Nonparametric estimands**. A local functional $\theta_0^{\lim} \in \mathbb{R}$ is a scalar that takes the form $\theta_0^{\lim} = \lim_{h \to 0} \theta_0^h$, where $\theta_0^h = \frac{1}{n} \sum_{i=1}^n \theta_i^h$, $\theta_i^h = \mathbb{E}[m_h(W_{i,\cdot}, \gamma_0)] = \mathbb{E}[\ell_h(W_{ij}) m(W_{i,\cdot}, \gamma_0)]$. Here, $\ell_h$ is a Nadaraya Watson weighting with bandwidth $h$ and $W_{ij}$ is a scalar component of $W_{i,\cdot}$. $\theta_0^{\lim}$ is a nonparametric quantity. However, it can be approximated by the sequence $\{\theta_0^h\}$. Each $\theta_0^h$ can be analyzed like a weighted functional as long as we keep track of how certain quantities depend on $h$. By this logic, finite sample semiparametric theory for $\theta_0^h$ translates to finite sample nonparametric theory for $\theta_0^{\lim}$ up to some approximation error. In this sense, our analysis encompasses both semiparametric and nonparametric estimands.

**Example E.7** (Heterogeneous treatment effect). *Let $(D_i, V_i, X_{i,\cdot})$ concatenate treatment $D_i \in \{0, 1\}$, covariate of interest $V_i \in \mathbb{R}$, and other covariates $X_{i,\cdot} \in \mathbb{R}^p$. Denote $\gamma_0(D_i, V_i, X_{i,\cdot}) := \mathbb{E}[Y_i | D_i, V_i, X_{i,\cdot}]$. Under the assumption of selection on $(V_i, X_{i,\cdot})$ and identical distribution of $V_i$, the heterogeneous treatment effect for the subpopulation with subcovariate value $v$ is given by $\theta_i = \mathbb{E}[\gamma_0(1, V_i, X_{i,\cdot}) - \gamma_0(0, V_i, X_{i,\cdot}) | V_i = v] = \lim_{h \to 0} \mathbb{E}[\ell_h(V_i) \{\gamma_0(1, V_i, X_{i,\cdot}) -$*

$\gamma_0(0, V_i, X_{i,\cdot})\}]$ *where* $\ell_h(V_i) = \frac{K\{(V_i - v)/h\}}{\omega}$, $\omega = \mathbb{E}[K\{(V_i - v)/h\}]$, *and* $K$ *is a standard kernel function. With uncorrupted treatment, uncorrupted covariate of interest, and corrupted other covariates,* $W_{i,\cdot} = (D_i, V_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})$ *where* $(H_{i,\cdot}, \pi_{i,\cdot})$ *are measurement error and missingness for the other covariates.*

Appendix J formally defines our general class of semiparametric and nonparametric estimands. Each example belongs to the class under generalizations of Assumption 5.10.

**Missing outcomes**. So far, $Y_i$ has been uncorrupted. Measurement error and differential privacy of $Y_i$ are allowed by response noise $\varepsilon_i$. An important additional issue is outcome attrition: for some observations, $Y_i$ is missing in a way that may depend on the true regressors. The enriched observation model is $Y_i = \gamma_0(D_i, X_{i,\cdot}, S_i) + \varepsilon_i$, $Z_{i,\cdot} = [X_{i,\cdot} + H_{i,\cdot}] \odot \pi_{i,\cdot}$, and $\tilde{Y}_i = Y_i \cdot S_i$ with $S_i \in \{1, \mathtt{NA}\}$. Instead of $(Y_i, D_i, X_{i,\cdot})$, the analyst observes $(\tilde{Y}_i, D_i, Z_{i,\cdot})$. Outcome $Y_i$ may be missing at random conditional on true regressors $(D_i, X_{i,\cdot})$, of which $X_{i,\cdot}$ may be corrupted. The extended semiparametric model is $\mathbb{E}[\tilde{Y}_i | D_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}, S_i = 1] = \mathbb{E}[Y_i | D_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}, S_i = 1] = \mathbb{E}[Y_i | D_i, X_{i,\cdot}, S_i = 1] = \gamma_0(D_i, X_{i,\cdot}, S_i = 1)$. For this extension, replace $Y_i$ with $\tilde{Y}_i$ and replace $(D, X_{i,\cdot})$ with $(D_i, X_{i,\cdot}, S_i)$.

# F   Nonlinearity

We characterize the class of nonlinear dictionaries $b : \mathbb{R}^p \to \mathbb{R}^{p'}$ for which our main results go through. We discuss two classes of dictionaries and delay proofs to the end.

**Polynomial dictionary**. We refer to the following three simple properties as dictionary continuity, since they imply that the data cleaning results for original regressors imply similar data cleaning results for technical regressors constructed from the dictionary. We state the properties then verify them for the polynomial dictionary of degree $d_{\max}$.

**Assumption F.1** (Dictionary continuity). *(i) For any two matrices* $\boldsymbol{M}^{(1)}, \boldsymbol{M}^{(2)} \in \mathbb{R}^{n \times p}$, $\|b(\boldsymbol{M}^{(1)}) - b(\boldsymbol{M}^{(2)})\|_{2,\infty}^2 \leq C_b' \|\boldsymbol{M}^{(1)} - \boldsymbol{M}^{(2)}\|_{2,\infty}^2$; *(ii) for any* $\boldsymbol{M} \in \mathbb{R}^{n \times p}$, $rank\{b(\boldsymbol{M})\} \leq \{rank(\boldsymbol{M})\}^{C_b''}$; *(iii) for any* $v \in \mathbb{R}^p$, $\|b(v)\|_{\max} \leq (\|v\|_{\max})^{C_b'''}$.

For much of our argument to go through, it suffices that the dictionary exhibits three simple properties: clean original regressors imply clean technical regressors; low rank

original regressors imply low rank technical regressors; and bounded original regressors imply bounded technical regressors. Polynomial dictionaries have these properties.

**Definition F.1** (Polynomial dictionary). *Let $v = (v_1, v_2, \ldots, v_p) \in \mathbb{R}^p$. Consider the dictionary $b^{\text{POLY}}$, where for $k \in [p']$, $b_k^{\text{POLY}}(v) = \prod_{\ell=1}^{d(k)} v_\ell$ with $v_\ell \in \{v_1, \ldots, v_p\}$.*

That is, each basis function $b_k^{\text{POLY}}(v)$ in the dictionary is a polynomial of degree $d(k) \leq d_{\max}$ constructed from coordinates of $v$, allowing for repeats. This class of dictionaries is commonly used in empirical economic research. It nests as a special case the interacted dictionary studied in the main text, which permits a rich model of heterogeneous treatment effects. Pleasingly, for this class, the dictionary constants $(C_b', C_b'', C_b''')$ do not depend on $p'$. Rather, $(C_b', C_b'', C_b''')$ depend on the maximum degree $d_{\max}$ of the polynomial dictionary.

**Proposition F.1** (Verifying dictionary continuity). *$b^{\text{POLY}}$ of degree $d_{\max}$ satisfies Assumption F.1 with $C_b' \leq 2^{d_{\max}} \cdot \|\boldsymbol{M}^{(1)}\|_{\max}^{2d_{\max}} \cdot \|\boldsymbol{M}^{(2)}\|_{\max}^{2d_{\max}}$, $C_b'' \leq d_{\max}$, and $C_b''' \leq d_{\max}$.*

This class of dictionaries preserves the low rank approximation in the following sense.

**Proposition F.2** (Low rank approximation is preserved). *Suppose Assumption 5.1 holds and the true covariates satisfy $\boldsymbol{X} = \boldsymbol{X}^{(LR)} + \boldsymbol{E}^{(LR)}$ where $r = rank\{\boldsymbol{X}^{(LR)}\}$ and $\Delta_E = \|\boldsymbol{E}^{(LR)}\|_{\max}$. Consider $b = b^{\text{POLY}}$ of degree $d_{\max}$. Then $r' := rank\{b(\boldsymbol{X}^{(LR)})\} \leq r^{d_{\max}}$ and $\Delta_E' := \|b(\boldsymbol{X}) - b(\boldsymbol{X}^{(LR)})\|_{\max} \leq C \bar{A}^{d_{\max}} \cdot d_{\max} \Delta_E$.*

The same logic applies for dictionaries applied to $(D_i, X_{i,.})$ rather than $X_{i,.}$. The generalization of Appendix A with nonlinear dictionaries is immediate from these results.

**Polynomial dictionary with uncorrupted nonlinearity**. Assumption F.1 suffices to generalize our data cleaning results. For analysis of the error-in-variable estimators, we impose a further assumption, which constrains which kinds of terms can appear as technical regressors. Consider the polynomial dictionary of degree $d_{\max}$, where the only source of nonlinearity is powers and interactions with regressors known to be uncorrupted.

**Definition F.2** (Polynomial dictionary with uncorrupted nonlinearity). *Suppose the observed regressors consist of one uncorrupted regressor $D_i$ and several corrupted regressors $X_{i,.}$. Consider a polynomial dictionary $b^{\text{POLY}}$ of degree $d_{\max}$ such that each basis function $b_k^{\text{POLY}}$ is at most linear in the corrupted regressors. By definition, $p' \leq C \cdot d_{\max} p$.*

For example, in Example E.1 where $D_i$ is uncorrupted, the interacted dictionary $b :$ $(D_i, X_{i,\cdot}) \mapsto \{D_i X_{i,\cdot}, (1 - D_i) X_{i,\cdot}\}$ satisfies this property. In Example E.4 where $D_i$ is uncorrupted, the nonlinear dictionary $b : (D_i, X_{i,\cdot}) \mapsto (1, D_i, X_{i,\cdot}, D_i X_{i,\cdot}, D_i^2)$ satisfies this property as well since it contains $D_i^2$ but does not contain $X_{ij}^2$. Intuitively, this family of dictionaries avoids compounding measurement error because the corrupted regressors are not multiplied with each other. For readability, we focus on the case of one uncorrupted regressor, which can be conceptualized as $b : (D_i, X_{i,\cdot}) \mapsto (1, D_i, ..., D_i^{d_{\max}}, X_{i,\cdot}, D_i X_{i,\cdot}, ..., D_i^{d_{\max}-1} X_{i,\cdot})$ where $D_i$ is uncorrupted and $X_{i,\cdot}$ are corrupted. Definition F.2 naturally generalizes to the case of multiple uncorrupted regressors. We require three properties to hold after the dictionary is applied to the data.

**Assumption F.2** (Dictionary is non-collapsing)**.** *The dictionary does not collapse in the following sense. (i) Recall that we set $k := rank(\hat{\boldsymbol{X}})$ equal to $r := rank\{\boldsymbol{X}^{(\text{LR})}\}$. We further assume $k' := rank\{b(D, \hat{\boldsymbol{X}})\}$ is equal to $r' := rank[b\{D, \boldsymbol{X}^{(\text{LR})}\}]$. (ii) Assumption 5.4 posits that the smallest singular value of $\boldsymbol{X}^{(\text{LR})}$ is $s_r \geq C\sqrt{\frac{np}{r}}$. We further posit that the smallest singular value of $b\{D, \boldsymbol{X}^{(\text{LR})}\}$ is $s'_{r'} \geq C\sqrt{\frac{np}{r'}}$. (iii) Using the notation of one uncorrupted regressor, the technical regressors $(1, D_i, ...D_i^{d_{\max}})$ are full rank.*

The first property in Assumption F.2 ensures two matrices of equal rank get mapped to two new matrices of equal rank. The second property imposes that singular values, after dictionary mapping, remain well balanced. In particular, we allow for a weaker signal to noise ratio for technical regressors since $r' \geq r$. We do *not* impose $s'_{r'} \geq C\sqrt{\frac{np'}{r'}}$, which is a stronger and less plausible requirement since it implies that the signal to noise ratio increases with the dictionary dimension $p'$. The third property is a technical assumption which allows the theory of implicit data cleaning to generalize.

Appendices B and C generalize to accommodate nonlinear dictionaries under this additional assumption. See the previous draft for explicit algebra. We turn to proofs.

**Lemma F.1.** *For $b^{\text{POLY}}$, $C'_b \leq 2^{d_{\max}} \cdot \|\boldsymbol{M}^{(1)}\|_{\max}^{2d_{\max}} \cdot \|\boldsymbol{M}^{(2)}\|_{\max}^{2d_{\max}}$.*

*Proof.* We introduce the notation $[b^{\text{POLY}}(\boldsymbol{M})]_{ik} = \prod_{\{j(k)\}} M_{ij(k)}$, where $j(k) \in [p]$, $M_{ij(k)} \in \{M_{i1}, ..., M_{ip}\}$, and $|\{j(k)\}| = d(k)$. We will simplify notation in the following way. Fix $k$. Let $M_{i\ell}$ refer to the $\ell$-th element of the product, where $\ell \in [d(k)]$. Therefore

$[b^{\text{POLY}}(\boldsymbol{M})]_{ik} = \prod_{\{j(k)\}} M_{ij(k)} = \prod_{\ell=1}^{d(k)} M_{i\ell}$. Then for any column $k \in [p']$,

$$\|b(\boldsymbol{M}^{(1)})_{\cdot,k} - b(\boldsymbol{M}^{(2)})_{\cdot,k}\|_2^2 = \sum_{i=1}^n \left( \prod_{\ell=1}^{d(k)} M_{i\ell}^{(1)} - \prod_{\ell=1}^{d(k)} M_{i\ell}^{(2)} \right)^2$$

$$\leq 2 \sum_{i=1}^n \left( \prod_{\ell=1}^{d(k)} M_{i\ell}^{(1)} - M_{i1}^{(2)} \prod_{\ell=2}^{d(k)} M_{i\ell}^{(1)} \right)^2 + 2 \sum_{i=1}^n \left( \prod_{\ell=1}^{d(k)} M_{i\ell}^{(2)} - M_{i1}^{(2)} \prod_{\ell=2}^{d(k)} M_{i\ell}^{(1)} \right)^2.$$

The first term equals $\sum_{i=1}^n \left( M_{i1}^{(1)} - M_{i1}^{(2)} \right)^2 \left( \prod_{\ell=2}^{d(k)} M_{i\ell}^{(1)} \right)^2 \leq \|\boldsymbol{M}^{(1)}\|_{\max}^{2d_{\max}} \sum_{i=1}^n \left( M_{i1}^{(1)} - M_{i1}^{(2)} \right)^2 \leq$ $\|\boldsymbol{M}^{(1)}\|_{\max}^{2d_{\max}} \|\boldsymbol{M}^{(1)} - \boldsymbol{M}^{(2)}\|_{2,\infty}^2$. The second term equals $\sum_{i=1}^n \left( M_{i1}^{(2)} \left( \prod_{\ell=2}^{d(k)} M_{i\ell}^{(2)} - \prod_{\ell=2}^{d(k)} M_{i\ell}^{(1)} \right) \right)^2 \leq$ $\|\boldsymbol{M}^{(2)}\|_{\max}^2 \sum_{i=1}^n \left( \prod_{\ell=2}^{d(k)} M_{i\ell}^{(2)} - \prod_{\ell=2}^{d(k)} M_{i\ell}^{(1)} \right)^2$. Recursing with $\sum_{i=1}^n \left( \prod_{\ell=2}^{d(k)} M_{i\ell}^{(2)} - \prod_{\ell=2}^{d(k)} M_{i\ell}^{(1)} \right)^2$ gives the desired result. $\qquad\square$

**Lemma F.2.** *For $b^{\text{POLY}}$, $C_b'' \leq d_{\max}$.*

*Proof.* Fix $\boldsymbol{M} \in \mathbb{R}^{n \times p}$ with rank $r$. For notational simplicity, let $M_{i\ell}$ refer to the $\ell$-th element of the product in $[b^{\text{POLY}}(\boldsymbol{M})]_{ik}$. Observe that $b^{\text{POLY}}(\boldsymbol{M})$ can be equivalently represented as $b^{\text{POLY}}(\boldsymbol{M}) = \boldsymbol{B}^{(1)} \odot, ..., \odot \boldsymbol{B}^{(d_{\max})}$, where $\odot$ means Hadamard product, $\boldsymbol{B}^{(\ell)} \in \mathbb{R}^{n \times p'}$, and for $\ell \in [d_{\max}], i \in [n], k \in [p'], [\boldsymbol{B}^{(\ell)}]_{ik} = M_{i\ell}$ if $\ell \leq d(k)$ and $[\boldsymbol{B}^{(\ell)}]_{ik} = 1$ if $\ell > d(k)$. Since each column of each $\boldsymbol{B}^{(\ell)}$ is either a column of $\boldsymbol{M}$ or a column of ones, it has rank at most $r$. The rank of a Hadamard product is bounded by the product of ranks and so $rank\{b^{\text{POLY}}(\boldsymbol{M})\} \leq \prod_{\ell=1}^{d_{\max}} r = r^{d_{\max}}$. $\qquad\square$

**Lemma F.3.** *For $b^{\text{POLY}}$, $C_b''' \leq d_{\max}$.*

*Proof.* Denote $v \in \mathbb{R}^p$ with $\|v\|_\infty \leq \bar{A}$. Then $b_k^{\text{POLY}}(v) = \prod_{\ell=1}^{d(k)} v_\ell \leq \bar{A}^{d_{\max}}$. $\qquad\square$

*Proof of Proposition F.1.* Immediate from Lemmas F.1, F.2, and F.3. $\qquad\square$

**Lemma F.4.** *If Assumption 5.1 holds, then $\|\boldsymbol{X}^{(\text{LR})}\|_{\max} \leq 3\bar{A}$.*

*Proof.* Suppose we have $\boldsymbol{X}^{(\text{LR})}$ with rank $r$ such that $\|\boldsymbol{X}^{(\text{LR})}\|_{\max} > 3\bar{A}$. By reverse triangle inequality $\Delta_{E,\boldsymbol{X}^{(\text{LR})}} = \|\boldsymbol{X}^{(\text{LR})} - \boldsymbol{X}\|_{\max} \geq \|\boldsymbol{X}^{(\text{LR})}\|_{\max} - \|\boldsymbol{X}\|_{\max} > 2\bar{A}$. We construct $\boldsymbol{B}^{(\text{LR})}$ with rank $r$ such that $\|\boldsymbol{B}^{(\text{LR})}\|_{\max} \leq 3\bar{A}$ and $\Delta_{E,\boldsymbol{B}^{(\text{LR})}} < \Delta_{E,\boldsymbol{X}^{(\text{LR})}}$. Set $\boldsymbol{B}^{(\text{LR})} = \frac{\bar{A}}{\|\boldsymbol{X}^{(\text{LR})}\|_{\max}} \cdot \boldsymbol{X}^{(\text{LR})}$. Clearly, $rank\{\boldsymbol{B}^{(\text{LR})}\} = rank\{\boldsymbol{X}^{(\text{LR})}\}$. By construction $\|\boldsymbol{B}^{(\text{LR})}\|_{\max} \leq \bar{A}$, so $\Delta_{E,\boldsymbol{B}^{(\text{LR})}} = \|\boldsymbol{B}^{(\text{LR})} - \boldsymbol{X}\|_{\max} \leq \|\boldsymbol{B}^{(\text{LR})}\|_{\max} + \|\boldsymbol{X}\|_{\max} \leq 2\bar{A}$. $\qquad\square$

*Proof of Proposition F.2.* By definition, $r = rank\{\boldsymbol{X}^{(\text{LR})}\}$. The first result follows directly from Proposition F.1. To see the second result, consider the case where $d_{\max} = 2$. Then any higher order entry of $b(\boldsymbol{X}) - b(\boldsymbol{X}^{(\text{LR})})$ is of the form $|X_{ij}X_{ik} - X_{ij}^{(\text{LR})}X_{ik}^{(\text{LR})}| \leq |X_{ij}X_{ik} - X_{ij}^{(\text{LR})}X_{ik}| + |X_{ij}^{(\text{LR})}X_{ik} - X_{ij}^{(\text{LR})}X_{ik}^{(\text{LR})}| \leq \bar{A}\Delta_E + 3\bar{A}\Delta_E$ by Lemma F.4. More generally, there are $d_{\max}$ such terms, and the largest is of the form $(3\bar{A})^{d_{\max}}\Delta_E$. $\qquad\square$

# G  Data cleaning supporting details

Define the unit ball $\mathbb{B}^p = \{v \in \mathbb{R}^p : \|v\|_2 \leq 1\}$ and sphere $\mathbb{S}^{p-1} = \{v \in \mathbb{R}^p : \|v\|_2 = 1\}$.

*Proof of Proposition 4.1.* Immediate from the law of iterated expectations. $\qquad\square$

**Proposition G.1** (Bound on $\|\hat{\boldsymbol{A}}\|_{\max}$). *Suppose $k = r$ and $\hat{s}_1, ..., \hat{s}_r \leq C\sqrt{\frac{np}{r}}$. Assume the following incoherence conditions for the corrupted singular vectors: $\|\hat{\boldsymbol{U}}_r\|_{\max} \leq Cn^{-1/2}$ and $\|\hat{\boldsymbol{V}}_r\|_{\max} \leq Cp^{-1/2}$. Then $\|\hat{\boldsymbol{A}}\|_{\max} \leq Cr^{1/2}$.*

The condition $\hat{s}_1, ..., \hat{s}_r \leq C\sqrt{\frac{np}{r}}$ holds with high probability under $s_1, ..., s_r \leq C\sqrt{\frac{np}{r}}$ by Weyl's inequality, similar to Proposition H.2. The condition $s_1, ..., s_r \leq C\sqrt{\frac{np}{r}}$ complements Assumption 5.4. To interpret the incoherence conditions, note that $U_{.,j} \in \mathbb{R}^n$ and $V_{.,j} \in \mathbb{R}^p$.

*Proof.* Write $\hat{A}_{ij} = \sum_{\ell=1}^{r} \hat{U}_{i\ell}\hat{s}_\ell\hat{V}_{j\ell}$. Hence $|\hat{A}_{ij}| \leq \sum_{\ell=1}^{r} |\hat{U}_{i\ell}| \cdot |\hat{s}_\ell| \cdot |\hat{V}_{j\ell}|$. $\qquad\square$

**Lemma G.1.** *Under Assumption 5.3, $\|\mathbb{E}[(\boldsymbol{Z}-\boldsymbol{A}\boldsymbol{\rho})^T(\boldsymbol{Z}-\boldsymbol{A}\boldsymbol{\rho})]\| \leq \rho_{\max}(1-\rho_{\min})\big(\max_{j\in[p]}\|A_{.,j}\|_2^2 + \|diag(\mathbb{E}[\boldsymbol{H}^T\boldsymbol{H}])\|\big) + \rho_{\max}\|\mathbb{E}[\boldsymbol{H}^T\boldsymbol{H}]\|$, where $\rho_{\max} := \max_{j\in[p]}\rho_j \leq 1$.*

*Proof.* Write $\mathbb{E}[(\boldsymbol{Z} - \boldsymbol{A}\boldsymbol{\rho})^T(\boldsymbol{Z} - \boldsymbol{A}\boldsymbol{\rho})] = \sum_{\ell=1}^{n}\mathbb{E}[(Z_{\ell,.} - A_{\ell,.}\boldsymbol{\rho}) \otimes (Z_{\ell,.} - A_{\ell,.}\boldsymbol{\rho})]$. Let $\boldsymbol{X} = \boldsymbol{A} + \boldsymbol{H}$. For any $(\ell, j) \in [n] \times [p]$, $\mathbb{E}[Z_{\ell j}] = \rho_j A_{\ell j}$ and $\mathbb{E}[Z_{\ell j}^2] = \rho_j\mathbb{E}[X_{\ell j}^2]$. Fix a row $\ell \in [n]$ and denote $\boldsymbol{W}^{(\ell)} = (Z_{\ell,.} - A_{\ell,.}\boldsymbol{\rho}) \otimes (Z_{\ell,.} - A_{\ell,.}\boldsymbol{\rho})$. By linearity of expectations, $\mathbb{E}[W_{ij}^{(\ell)}] = \mathbb{E}[Z_{\ell i}Z_{\ell j}] - \rho_j\mathbb{E}[Z_{\ell i}]A_{\ell j} - \rho_i\mathbb{E}[Z_{\ell j}]A_{\ell i} + \rho_i\rho_j A_{\ell i}A_{\ell j}$. Suppose $i = j$, then $\mathbb{E}[W_{ii}^{(\ell)}] = \rho_i\mathbb{E}[X_{\ell i}^2] - \rho_i^2 A_{\ell i}^2 = \rho_i(1 - \rho_i)\mathbb{E}[X_{\ell i}^2] + \rho_i^2\mathbb{E}[(X_{\ell i} - A_{\ell i})^2]$. On the other hand, if $i \neq j$, $\mathbb{E}[W_{ij}^{(\ell)}] \leq \sqrt{\rho_i\rho_j}\mathbb{E}[(X_{\ell i} - A_{\ell i})(X_{\ell j} - A_{\ell j})]$ since $\mathbb{E}[Z_{\ell i}Z_{\ell j}] = \mathbb{E}[\pi_{i\ell}\pi_{\ell j}]\mathbb{E}[X_{\ell i}X_{\ell j}] \leq \sqrt{\mathbb{E}[\pi_{\ell i}^2]}\sqrt{\mathbb{E}[\pi_{\ell j}^2]}\mathbb{E}[X_{\ell i}X_{\ell j}] = \sqrt{\rho_i\rho_j}\mathbb{E}[X_{\ell i}X_{\ell j}]$. Therefore, we can bound $\boldsymbol{W}^{(\ell)}$ as the sum of two matrices and hence $\mathbb{E}[\boldsymbol{W}^{(\ell)}] \leq \rho_{\max}(1 - \rho_{\min})\mathbb{E}[diag(X_{\ell,.} \otimes X_{\ell,.})] + \rho_{\max}\mathbb{E}[H_{\ell,.} \otimes H_{\ell,.}]$. Summing over all rows $\ell \in [n]$ yields $\mathbb{E}[(\boldsymbol{Z}-\boldsymbol{A}\boldsymbol{\rho})^T(\boldsymbol{Z}-\boldsymbol{A}\boldsymbol{\rho})] \leq \rho_{\max}(1-\rho_{\min})diag(\mathbb{E}[\boldsymbol{X}^T\boldsymbol{X}]) + \rho_{\max}\mathbb{E}[\boldsymbol{H}^T\boldsymbol{H}]$. To complete the proof, we apply triangle inequality: $\|\mathbb{E}[(\boldsymbol{Z} - \boldsymbol{A}\boldsymbol{\rho})^T(\boldsymbol{Z} - \boldsymbol{A}\boldsymbol{\rho})]\| \leq$

$\rho_{\max}(1-\rho_{\min})\big\|diag(\mathbb{E}[\boldsymbol{X}^T\boldsymbol{X}])\big\|+\rho_{\max}\big\|\mathbb{E}[\boldsymbol{H}^T\boldsymbol{H}]\big\|$ and since $\boldsymbol{H}$ is zero mean, $\big\|diag(\mathbb{E}[\boldsymbol{X}^T\boldsymbol{X}])\big\| \leq \big\|diag(\boldsymbol{A}^T\boldsymbol{A})\big\| + \big\|diag(\mathbb{E}[\boldsymbol{H}^T\boldsymbol{H}])\big\|$. $\qquad\qquad\square$

**Lemma G.2** (Lemma H.2 of [Agarwal et al., 2021]). *Suppose that $X \in \mathbb{R}^n$ and $P \in \{0,1\}^n$ are random vectors. Then for any $a \geq 1$, $\|X \odot P\|_{\psi_a} \leq \|X\|_{\psi_a}$.*

**Lemma G.3.** *Under Assumptions 5.1, 5.2, and 5.3, $\|Z_{i,\cdot} - A_{i,\cdot}\boldsymbol{\rho}\|_{\psi_a} \leq K_a + \bar{A}\bar{K}$.*

*Proof.* To begin, write $\|Z_{i,\cdot} - A_{i,\cdot}\boldsymbol{\rho}\|_{\psi_a} \leq \|(X_{i,\cdot} - A_{i,\cdot}) \odot \pi_{i,\cdot}\|_{\psi_a} + \|A_{i,\cdot} \odot \pi_{i,\cdot} - A_{i,\cdot}\boldsymbol{\rho}\|_{\psi_a}$. By Lemma G.2 and Assumption 5.2, $\|(X_{i,\cdot} - A_{i,\cdot}) \odot \pi_{i,\cdot}\|_{\psi_a} \leq \|(X_{i,\cdot} - A_{i,\cdot})\|_{\psi_a} = \|H_{i,\cdot}\|_{\psi_a} \leq K_a$. By the definition of $\|\cdot\|_{\psi_a}$ and Assumption 5.1, $\|A_{i,\cdot}\odot\pi_{i,\cdot}-A_{i,\cdot}\boldsymbol{\rho}\|_{\psi_a} = \sup_{u\in\mathbb{B}^p}\big\|\sum_{j=1}^p u_j A_{ij}(\pi_{ij}-\rho_j)\big\|_{\psi_a}$ $= \bar{A}\sup_{u\in\mathbb{B}^p}\big\|\sum_{j=1}^p u_j\frac{A_{ij}}{\bar{A}}(\pi_{ij}-\rho_j)\big\|_{\psi_a}$. Let $v_j = u_j\frac{A_{ij}}{\bar{A}}$. Since $v \in \mathbb{B}^p$, we have

$$\sup_{u\in\mathbb{B}^p}\left\|\sum_{j=1}^p u_j\frac{A_{ij}}{\bar{A}}(\pi_{ij}-\rho_j)\right\|_{\psi_a} \leq \sup_{v\in\mathbb{B}^p}\left\|\sum_{j=1}^p v_j(\pi_{ij}-\rho_j)\right\|_{\psi_a} = \|\pi_{i,\cdot} - (\rho_1,...,\rho_p)\|_{\psi_a} \leq \bar{K},$$

using Assumption 5.3. $\qquad\qquad\square$

**Lemma G.4** (Proposition H.1 of [Agarwal et al., 2021]). *Let $\boldsymbol{W} \in \mathbb{R}^{n\times p}$ be a random matrix whose rows $\boldsymbol{W}_{i,\cdot}$ are independent $\psi_a$-random vectors for some $a \geq 1$. Then for any $\tau > 0$, with probability at least $1-\frac{2}{n^{1+\tau}p^\tau}$, $\|\boldsymbol{W}\| \leq \big\|\mathbb{E}\boldsymbol{W}^T\boldsymbol{W}\big\|^{1/2} + \sqrt{(1+\tau)p}\max_{i\in[n]}\|\boldsymbol{W}_{i,\cdot}\|_{\psi_a}\Big\{1+(2+\tau)\ln(np)\Big\}^{\frac{1}{a}}\sqrt{\ln(np)}$.*

**Proposition G.2.** *Under Assumptions 5.1, 5.2, and 5.3 $\mathbb{P}(\mathcal{E}_1^c) \leq \frac{2}{n^{11}p^{10}} < \frac{2}{n^{10}p^{10}}$.*

*Proof.* We show that for all $\tau > 0$, with probability $1-\frac{2}{n^{1+\tau}p^\tau}$, $\|\boldsymbol{Z}-\boldsymbol{A}\boldsymbol{\rho}\| \leq C\sqrt{n}(\bar{A}+\kappa+K_a)+\sqrt{1+\tau}\sqrt{p}(K_a+\bar{A}\bar{K})\{1+(2+\tau)\ln(np)\}^{\frac{1}{a}}\sqrt{\ln(np)}$. Setting $\tau = 10$ and simplifying the bound yields the result. By Lemma G.1, $\|\mathbb{E}[(\boldsymbol{Z}-\boldsymbol{A}\boldsymbol{\rho})^T(\boldsymbol{Z}-\boldsymbol{A}\boldsymbol{\rho})]\| \leq \max_{j\in[p]}\|A_{\cdot,j}\|_2^2 + \|diag(\mathbb{E}[\boldsymbol{H}^T\boldsymbol{H}])\|+\|\mathbb{E}[\boldsymbol{H}^T\boldsymbol{H}]\|$. We bound these terms as $n\bar{A}^2$, $nCK_a$, and $n\kappa^2$, respectively, then plug them and Lemma G.3 into Lemma G.4. $\qquad\qquad\square$

**Lemma G.5** (Lemma H.4 of [Agarwal et al., 2021]). *Let $X_1,\ldots,X_n$ be independent random variables with mean zero. For $a \geq 1$, $\|\sum_{i=1}^n X_i\|_{\psi_a} \leq C\left(\sum_{i=1}^n\|X_i\|_{\psi_a}^2\right)^{1/2}$.*

**Lemma G.6.** *Under Assumptions 5.1, 5.2, and 5.3, $\|Z_{\cdot,j} - \rho_j A_{\cdot,j}\|_{\psi_a} \leq C(K_a + \bar{A}\bar{K})$.*

*Proof.* Write $\|Z_{\cdot,j} - \rho_j A_{\cdot,j}\|_{\psi_a} = \sup_{u\in\mathbb{S}^{n-1}}\|u^T(\boldsymbol{Z}-\boldsymbol{A}\boldsymbol{\rho})e_j\|_{\psi_a} = \sup_{u\in\mathbb{S}^{n-1}}\big\|\sum_{i=1}^n u_i(Z_{i,\cdot} - A_{i,\cdot}\boldsymbol{\rho})e_j\big\|_{\psi_a}$. By Lemma G.5, its bound is $C\sup_{u\in\mathbb{S}^{n-1}}\left(\sum_{i=1}^n u_i^2\|(Z_{i,\cdot} - A_{i,\cdot}\boldsymbol{\rho})e_j\|_{\psi_a}^2\right)^{1/2} \leq C\max_{i\in[n]}\|(Z_{i,\cdot} - A_{i,\cdot}\boldsymbol{\rho})e_j\|_{\psi_a}$. The conclusion follows from Lemmas G.2 and G.3. $\qquad\qquad\square$

**Lemma G.7** (Lemma I.7 of [Agarwal et al., 2021])**.** *Let* $W_1, \ldots, W_n$ *be a sequence of* $\psi_a$-*random variables for some* $a \geq 1$. *For any* $t \geq 0$, $\mathbb{P}(\sum_{i=1}^n W_i^2 > t) \leq 2 \sum_{i=1}^n \exp\left\{-\left(\frac{t}{n\|W_i\|_{\psi_a}^2}\right)^{a/2}\right\}$.

**Proposition G.3.** *Under Assumptions 5.1, 5.2, and 5.3,* $\mathbb{P}(\mathcal{E}_2^c) \leq \frac{2}{n^{10}p^{10}}$

*Proof.* Fix $j$. Write $\|Z_{\cdot,j} - \rho_j A_{\cdot,j}\|_2^2 = \sum_{i=1}^n W_i^2$, where $W_i = e_i^T(Z_{\cdot,j} - \rho_j A_{\cdot,j})$. By Lemmas G.2 and G.6, $\|W_i\|_{\psi_a} \leq \|Z_{\cdot,j} - \rho_j A_{\cdot,j}\|_{\psi_a} \leq C(K_a + \bar{K}\bar{A})$. By Lemma G.7 and the union bound, we arrive at the conclusion. □

**Proposition G.4.** *Under Assumptions 5.1, 5.2, and 5.3,* $\mathbb{P}(\mathcal{E}_3^c) \leq \frac{2}{n^{10}p^{10}}$.

*Proof.* The key equality is $\|\boldsymbol{U}_k\boldsymbol{U}_k^T(Z_{\cdot,j} - \rho_j A_{\cdot,j})\|_2^2 = \sum_{i=1}^k W_i^2$, where $W_i = u_i^T(Z_{\cdot,j} - \rho_j A_{\cdot,j})$. To see that it holds, set $v = Z_{\cdot,j} - \rho_j A_{\cdot,j}$. Then $\|\boldsymbol{U}_k\boldsymbol{U}_k^T v\|_2^2 = v^T\boldsymbol{U}_k\boldsymbol{U}_k^T\boldsymbol{U}_k\boldsymbol{U}_k^T v = v^T\boldsymbol{U}_k\boldsymbol{U}_k^T v = W^TW$. The rest is analogous to Proposition G.3. □

**Proposition G.5.** *Under Assumption 5.3,* $\mathbb{P}(\mathcal{E}_4^c) \leq \frac{2}{n^{10}p^{10}}$.

*Proof.* Fix $\delta > 1$. Define the event $\mathcal{E}_{(j)} = \{\frac{1}{\delta}\rho_j \leq \hat{\rho}_j \leq \delta\rho_j\}$. By the Chernoff bound for binary variables, $\mathbb{P}(\mathcal{E}_{(j)}^c) \leq 2\exp\left(-\frac{(\delta-1)^2}{2\delta^2}n\rho_j\right) \leq 2\exp\left(-\frac{(\delta-1)^2}{2\delta^2}n\rho_{\min}\right)$. Hence by De Morgan's law and the union bound $\mathbb{P}(\mathcal{E}_4^c) = \mathbb{P}\left(\left\{\bigcap_{j\in[p]}\mathcal{E}_{(j)}\right\}^c\right) = \mathbb{P}\left(\bigcup_{j\in[p]}\mathcal{E}_{(j)}^c\right) \leq 2p\exp\left(-\frac{(\delta-1)^2}{2\delta^2}n\rho_{\min}\right)$. Solve $\frac{2}{n^{10}p^{10}} \geq \frac{2}{n^{11}p^{10}} = 2p\exp\left(-\frac{(\delta-1)^2}{2\delta^2}n\rho_{\min}\right)$ for $\delta$. □

**Proposition G.6.** *Under Assumption 5.3,* $\mathbb{P}(\mathcal{E}_5^c) \leq \frac{2}{n^{10}p^{10}}$.

*Proof.* Define the event $\mathcal{E}_{(j)} = \{|\hat{\rho}_j - \rho_j| \leq t\}$. By Hoeffding's inequality for bounded variables, $\mathbb{P}(\mathcal{E}_{(j)}^c) \leq 2\exp(-2nt^2)$. By De Morgan's law and the union bound, $\mathbb{P}(\mathcal{E}_5^c) = \mathbb{P}\left(\left\{\bigcap_{j\in[p]}\mathcal{E}_{(j)}\right\}^c\right) = \mathbb{P}\left(\bigcup_{j\in[p]}\mathcal{E}_{(j)}^c\right) \leq 2p\exp(-2nt^2)$. To arrive at the desired result, solve $\frac{2}{n^{10}p^{10}} \geq \frac{2}{n^{11}p^{10}} = 2p\exp(-2nt^2)$ for $t$. □

# H   Error-in-variable regression supporting details

We propose a new norm for error-in-variable regression analysis $\mathcal{R}(\hat{\gamma})$ (Theorem 5.2). Our norm is essentially an on-average generalization error, and we demonstrate its compatibility with semiparametric theory (Theorem 5.4). To prove the guarantee, we combine an analysis of TEST ERROR (Proposition B.3) with a new theory of implicit data cleaning (Proposition B.4). The former builds on an analysis of TRAIN ERROR (Proposition B.2).

In the PCR literature, the norm called TRAIN ERROR is standard. The norm called TEST ERROR is similar in spirit to [Agarwal et al., 2020a]. Our norm $\mathcal{R}(\hat{\gamma})$ extends these ideas.

We also propose a new norm for error-in-variable balancing weight analysis $\mathcal{R}(\hat{\alpha})$ (Theorem 5.3) that is compatible with semiparametric theory (Theorem 5.4). Our definitions of balancing weight TRAIN ERROR (Proposition C.2), TEST ERROR (Proposition C.3), and implicit data cleaning (Proposition C.4) norms all appear to be conceptual innovations.

*Proof of Proposition 4.2.* For $i \in$ TEST, $\hat{\gamma}(D_i, Z_{i,\cdot}) = b(D_i, Z_{i,\cdot} \hat{\boldsymbol{\rho}}^{-1})\hat{\beta}$, which equals

$$\begin{bmatrix} D_i Z_{i,\cdot} \hat{\boldsymbol{\rho}}^{-1} & (1 - D_i)Z_{i,\cdot} \hat{\boldsymbol{\rho}}^{-1} \end{bmatrix} \begin{bmatrix} \hat{\beta}^{\text{TREAT}} \\ \hat{\beta}^{\text{UNTREAT}} \end{bmatrix} = \begin{bmatrix} D_i Z_{i,\cdot} & (1 - D_i)Z_{i,\cdot} \end{bmatrix} \begin{bmatrix} \hat{\boldsymbol{\rho}}^{-1} \hat{\beta}^{\text{TREAT}} \\ \hat{\boldsymbol{\rho}}^{-1} \hat{\beta}^{\text{UNTREAT}} \end{bmatrix}.$$

Both $(\hat{\boldsymbol{\rho}}, \hat{\beta})$ are calculated from TRAIN, while $(D_i, Z_{i,\cdot})$ and $(D_j, Z_{j,\cdot})$ are i.n.i.d. $\qquad\square$

**Lemma H.1** (Theorem 4.6.1 of [Vershynin, 2018])**.** *Let $\boldsymbol{U} \in \mathbb{R}^{m \times r}$ whose rows are independent, mean zero, subGaussian, and isotropic with $\|U_{i,\cdot}\|_{\psi_2} \leq K_u$. Then for any $t \geq 0$, with probability $1 - 2e^{-t^2}$, $\sqrt{m} - CK_u^2(\sqrt{r} + t) \leq s_r(\boldsymbol{U}) \leq s_1(\boldsymbol{U}) \leq \sqrt{m} + CK_u^2(\sqrt{r} + t)$.*

**Proposition H.1** (Verifying row space inclusion)**.** *By hypothesis, $\text{rank}\{\boldsymbol{A}^{(\text{LR})}\} = r$, so it admits a representation $A_{ij}^{(\text{LR})} = \langle u_i, v_j, \rangle$ where $u_i, v_j \in \mathbb{R}^r$. Suppose $\{u_i\}$ are independent, mean zero, subGaussian, and isotropic with $\|u_i\|_{\psi_2} \leq K_u$. If $m \gg K_u^4 \cdot r \ln(mp)$ then with probability $1 - O\{(mp)^{-10}\}$, $\text{ROW}\{\boldsymbol{A}^{(\text{LR}),\text{TRAIN}}\} = \text{ROW}\{\boldsymbol{A}^{(\text{LR}),\text{TEST}}\}$.*

*Proof.* Consider $\boldsymbol{A}^{(\text{LR}),\text{TRAIN}}$. Let $\boldsymbol{U}$ have rows $\{U_{i,\cdot}\}$. By Lemma H.1 with $t = \ln^{\frac{1}{2}}(mp)$, $s_r(\boldsymbol{U}) \geq \sqrt{m} - CK_u^2\{\sqrt{r} + \ln^{\frac{1}{2}}(mp)\} \gg 0$. With high probability, $s_r(\boldsymbol{U}) \gg 0$, implying that $\{U_{i,\cdot}\}$ are full rank: $\text{ROW}(\boldsymbol{U}) = \mathbb{R}^r$. Consider $\boldsymbol{A}^{(\text{LR}),\text{TEST}}$. Let $\boldsymbol{U}'$ have rows $\{U'_{i,\cdot}\}$. Fix $i \in$ TEST. Since $U'_{i,\cdot} \in \mathbb{R}^r = \text{ROW}(\boldsymbol{U})$, there exists some $\lambda \in \mathbb{R}^r$ such that $U'_{i,\cdot} = \sum_{k=1}^{r} \lambda_k U_{k,\cdot}$. Therefore $A_{ij}^{(\text{LR}),\text{TEST}} = \langle U'_{i,\cdot}, V_{\cdot,j} \rangle = \langle \sum_{k=1}^{r} \lambda_k U_{k,\cdot}, V_{\cdot,j} \rangle = \sum_{k=1}^{r} \lambda_k \langle U_{k,\cdot}, V_{\cdot,j} \rangle = \sum_{k=1}^{r} \lambda_k A_{kj}^{(\text{LR}),\text{TRAIN}}$. Thus for any $i \in$ TEST, $A_{i,\cdot}^{(\text{LR}),\text{TEST}} \in \text{ROW}\{\boldsymbol{A}^{(\text{LR}),\text{TRAIN}}\}$. Therefore $\text{ROW}\{\boldsymbol{A}^{(\text{LR}),\text{TEST}}\} \subset \text{ROW}\{\boldsymbol{A}^{(\text{LR}),\text{TRAIN}}\}$. Likewise for the other direction. $\qquad\square$

The results for $\tilde{\mathcal{E}}_1, \tilde{\mathcal{E}}_2, \tilde{\mathcal{E}}_3$ follow from the results for $\mathcal{E}$. We focus on $\tilde{\mathcal{E}}_4$ and $\tilde{\mathcal{E}}_5$.

**Proposition H.2.** *If Assumptions 5.1, 5.2, 5.3, and 5.4 hold, $k = r$, and $\rho_{\min} \gg \tilde{C}\sqrt{r}\ln^{\frac{3}{2}}(np)\left(\frac{1}{\sqrt{p}} \vee \frac{1}{\sqrt{n}} \vee \Delta_E\right)$, where $\tilde{C} := C\bar{A}\left(\kappa + \bar{K} + K_a\right)$, then $\mathbb{P}(\tilde{\mathcal{E}}_4^c) \leq \frac{C}{n^{10}p^{10}}$.*

*Proof.* By Lemma A.1, with probability at least $1 - O\{1/(np)^{10}\}$, $|\hat{s}_r - s_r| \leq \Delta$. Hence $\hat{s}_r \geq s_r - \Delta$. We want to show $\Delta = o(s_r)$, i.e. $\Delta \leq c_n s_r$ where $c_n \to 0$. It suffices to show $\frac{\Delta}{s_r} \to 0$. In such case, $\hat{s}_r \geq s_r - \Delta \geq s_r - c_n s_r = (1 - c_n)s_r$, i.e. $\hat{s}_r \gtrsim s_r$ as desired. We upper bound $\Delta$ using Lemma A.1 and lower bound $s_r \geq C\sqrt{\frac{np}{r}}$ using Assumption 5.4 to derive the stated sufficient condition on $\rho_{\min}$. $\qquad\square$

**Lemma H.2.** *If Assumption 5.5 holds then* $\mathbb{E}[\langle \hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \varepsilon \rangle | \boldsymbol{A}] \leq \bar{\sigma}^2 k$.

*Proof.* Note that $\hat{\beta} = \hat{\boldsymbol{A}}^{\dagger} Y = \hat{\boldsymbol{A}}^{\dagger} \{\boldsymbol{A}^{(\mathrm{LR})} \beta^* + \varepsilon + \phi^{(\mathrm{LR})}\}$. Since $\varepsilon$ is conditionally independent of $\hat{\boldsymbol{A}}$, $\boldsymbol{A}^{(\mathrm{LR})}$, $\beta^*$, and $\phi^{(\mathrm{LR})}$ we have $\mathbb{E}[\langle \hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \varepsilon \rangle | \boldsymbol{A}] = \mathbb{E}[\langle \hat{\boldsymbol{A}}\hat{\boldsymbol{A}}^{\dagger}\varepsilon, \varepsilon \rangle \boldsymbol{A}]$. By properties of trace algebra, conditional independence of $\varepsilon$ from $\hat{\boldsymbol{A}}$, Assumption 5.5, and the fact that $rank(\hat{\boldsymbol{A}}) = k$,

$$\mathbb{E}[\langle \hat{\boldsymbol{A}}\hat{\boldsymbol{A}}^{\dagger}\varepsilon, \varepsilon \rangle | \boldsymbol{A}] = \mathbb{E}\left[trace\left(\hat{\boldsymbol{A}}\hat{\boldsymbol{A}}^{\dagger}\varepsilon\varepsilon^T\right) \mid \boldsymbol{A}\right] = trace\left(\mathbb{E}\left[\hat{\boldsymbol{A}}\hat{\boldsymbol{A}}^{\dagger} \mid \boldsymbol{A}\right]\mathbb{E}\left[\varepsilon\varepsilon^T \mid \boldsymbol{A}\right]\right)$$
$$\leq \bar{\sigma}^2 trace\left(\mathbb{E}\left[\hat{\boldsymbol{A}}\hat{\boldsymbol{A}}^{\dagger} \mid \boldsymbol{A}\right]\right) = \bar{\sigma}^2 k.$$

$\qquad\square$

**Lemma H.3** (Lemma A.3 of [Agarwal et al., 2020a]). *Let $X \in \mathbb{R}^n$ be random vector with independent mean zero subGaussian random coordinates with $\|X_i\|_{\psi_2} \leq K$. Let $a \in \mathbb{R}^n$ be another random vector that satisfies $\|a\|_2 \leq b$ almost surely for some constant $b \geq 0$. Then for all $t \geq 0$, $\mathbb{P}\left(\left|\sum_{i=1}^n a_i X_i\right| \geq t\right) \leq 2\exp\left(-\frac{ct^2}{K^2 b^2}\right)$, where $c > 0$ is a universal constant.*

**Lemma H.4** (Lemma A.4 of [Agarwal et al., 2020a]). *Let $X \in \mathbb{R}^n$ be a random vector with independent mean zero subGaussian coordinates where $\|X_i\|_{\psi_2} \leq K$. Let $\boldsymbol{B} \in \mathbb{R}^{n \times n}$ be a random matrix satisfying $\|\boldsymbol{B}\| \leq a$ and $\|\boldsymbol{B}\|_{Fr}^2 \leq b$ almost surely for some $a, b \geq 0$. Then for any $t \geq 0$, $\mathbb{P}\left(|X^T \boldsymbol{B} X - \mathbb{E}[X^T \boldsymbol{B} X]| \geq t\right) \leq 2 \cdot \exp\left\{-c \min\left(\frac{t^2}{K^4 b}, \frac{t}{K^2 a}\right)\right\}$.*

**Proposition H.3.** *If Theorem 5.1 conditions and Assumption 5.5 hold, $\mathbb{P}(\tilde{\mathcal{E}}_5^c) \leq \frac{C}{n^{10}p^{10}}$.*

*Proof.* We show that under Assumptions 5.1 and 5.5, with $k = r$, the following holds with probability at least $1 - O\{1/(np)^{10}\}$ with respect to randomness in $\varepsilon$: $\langle \hat{\boldsymbol{A}}(\hat{\beta} - \beta^*), \varepsilon \rangle \leq C\bar{\sigma}^2 \ln(np)\left\{r + \|\phi^{(\mathrm{LR})}\|_2 + \|\beta^*\|_1(\sqrt{n}\bar{A} + \|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{2,\infty})\right\}$. Simplifying yields the desired result.

Recall that $\hat{\beta} = \hat{V}_k \hat{\Sigma}_k^{-1} \hat{U}_k^T Y$, $\hat{A} = \hat{U}_k \hat{\Sigma}_k \hat{V}_k^T$, and $Y = A^{(\text{LR})} \beta^* + \phi^{(\text{LR})} + \varepsilon$. Thus, $\hat{A}\hat{\beta} = \hat{U}_k \hat{\Sigma}_k \hat{V}_k^T \hat{V}_k \hat{\Sigma}_k^{-1} \hat{U}_k^T Y = \hat{U}_k \hat{U}_k^T A^{(\text{LR})} \beta^* + \hat{U}_k \hat{U}_k^T \phi^{(\text{LR})} + \hat{U}_k \hat{U}_k^T \varepsilon$. Therefore, $\langle \hat{A}(\hat{\beta} - \beta^*), \varepsilon \rangle = \langle \hat{U}_k \hat{U}_k^T A^{(\text{LR})} \beta^*, \varepsilon \rangle + \langle \hat{U}_k \hat{U}_k^T \phi^{(\text{LR})}, \varepsilon \rangle + \langle \hat{U}_k \hat{U}_k^T \varepsilon, \varepsilon \rangle - \langle \hat{A}\beta^*, \varepsilon \rangle$.

For the first, second, and fourth terms, we use Lemma H.3. Note that $\|\hat{U}_k \hat{U}_k^T A^{(\text{LR})} \beta^*\|_2 \leq \|A^{(\text{LR})} \beta^*\|_2 \leq 3\sqrt{n}\bar{A}\|\beta^*\|_1$ since $\hat{U}_k \hat{U}_k^T$ is a projection matrix and $\|A^{(\text{LR})}\|_{2,\infty} \leq 3\sqrt{n}\bar{A}$ due to Lemma F.4. It follows that $\mathbb{P}\Big( \langle \hat{U}_k \hat{U}_k^T A^{(\text{LR})} \beta^*, \varepsilon \rangle \geq t \Big) \leq \exp\Big( -\frac{ct^2}{n\bar{A}^2 \|\beta^*\|_1^2 \bar{\sigma}^2} \Big)$. Similarly, $\|\hat{U}_k \hat{U}_k^T \phi^{(\text{LR})}\|_2 \leq \|\phi^{(\text{LR})}\|_2$ since $\hat{U}_k \hat{U}_k^T$ is a projection matrix. Hence $\mathbb{P}\Big( \langle \hat{U}_k \hat{U}_k^T \phi^{(\text{LR})}, \varepsilon \rangle \geq t \Big) \leq \exp\Big( -\frac{ct^2}{\|\phi^{(\text{LR})}\|_2^2 \bar{\sigma}^2} \Big)$. Finally, $\|\hat{A}\beta^*\|_2 \leq \big( \|\hat{A} - A\|_{2,\infty} + \|A\|_{2,\infty} \big)\|\beta^*\|_1$. Therefore $\mathbb{P}\Big( \langle \hat{A}\beta^*, \varepsilon \rangle \geq t \Big) \leq \exp\Big( -\frac{ct^2}{\bar{\sigma}^2 (n\bar{A}^2 + \|\hat{A} - A\|_{2,\infty}^2)\|\beta^*\|_1^2} \Big)$.

For the third term, we use Lemma H.4. Recall that $\varepsilon$ is conditionally independent of $\hat{U}_k, \hat{\Sigma}_k, \hat{V}_k$ since $\hat{A}$ is determined by $Z$. Hence $\mathbb{E}\big[ \langle \hat{U}_k \hat{U}_k^T \varepsilon, \varepsilon \rangle | A \big] = \mathbb{E}\big[ \varepsilon^T \hat{U}_k \hat{U}_k^T \varepsilon | A \big] = \mathbb{E}\big[ trace(\varepsilon \varepsilon^T \hat{U}_k \hat{U}_k^T) | A \big] = trace(\mathbb{E}[\varepsilon \varepsilon^T | A] \mathbb{E}[\hat{U}_k \hat{U}_k^T | A]) \leq \bar{\sigma}^2 trace(\mathbb{E}[\hat{U}_k^T \hat{U}_k | A]) = \bar{\sigma}^2 k$. Since $\hat{U}_k \hat{U}_k^T$ is a projection matrix, $\|\hat{U}_k \hat{U}_k^T\| \leq 1$ and $\|\hat{U}_k \hat{U}_k^T\|_{Fr}^2 = trace(\hat{U}_k \hat{U}_k^T \hat{U}_k \hat{U}_k^T) = trace(\hat{U}_k^T \hat{U}_k) = k$. Using Lemma H.4, it follows that for any $t > 0$ $\mathbb{P}\Big( \langle \hat{U}_k \hat{U}_k^T \varepsilon, \varepsilon \rangle \geq \bar{\sigma}^2 k + t | A \Big) \leq \exp\Big\{ -c\min\Big( \frac{t^2}{k\bar{\sigma}^4}, \frac{t}{\bar{\sigma}^2} \Big) \Big\}$. Hence it also holds unconditionally.

Set each probability equal to $1/(np)^{10}$, solve for $t$, then combine terms. $\qquad\square$

# I   Error-in-variable balancing weight supporting details

**Counterfactual moments**. We describe the counterfactual moments for general parameters and general dictionaries. In this appendix, we consider causal parameters of the form $\theta_0 = \frac{1}{2n} \sum_{i \in \text{TRAIN,TEST}} \theta_i$, where $\theta_i = \mathbb{E}[m(W_{i,\cdot}, \gamma_0)]$, $W_{i,\cdot} = (A_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})$, and we slightly abuse sample size notation to avoid overloading $m$. Given a dictionary $b : \mathbb{R}^p \to \mathbb{R}^{p'}$, define $b^{\text{SIGNAL}}(W_{i,\cdot}) = b(A_{i,\cdot})$ and $b^{\text{NOISE}}(W_{i,\cdot}) = b(Z_{i,\cdot})$.

**Algorithm I.1** (Counterfactual moment with data cleaning). *Given corrupted training covariates $Z^{TRAIN} \in \mathbb{R}^{n \times p}$, the dictionary $b : \mathbb{R}^p \to \mathbb{R}^{p'}$, and the formula $m : \mathcal{W} \times \mathbb{L}_2 \to \mathbb{R}$: (i) perform data cleaning on $Z^{TRAIN}$ to obtain $\hat{A}^{TRAIN} \in \mathbb{R}^{n \times p}$; (ii) for $i \in TRAIN$ calculate $m_{i,\cdot} = m(W_{i,\cdot}, b^{NOISE}) \in \mathbb{R}^{p'}$; (iii) for $i \in TRAIN$, calculate $\hat{m}_{i,\cdot}$ from $m_{i,\cdot}$ by overwriting $Z_{i,\cdot}$ with $\hat{A}_{i,\cdot}$; (iv) calculate $\hat{M} = \frac{1}{n} \sum_{i \in TRAIN} \hat{m}_{i,\cdot}$.*

To specialize this procedure, it suffices to describe $\hat{m}_i$. We provide the explicit expressions

for $\hat{m}_i$ for each leading example in the proof of Proposition I.1 below.

As theoretical devices, we introduce several related objects. First, we define the counterfactual vectors $\tilde{m}_{i,\cdot}, \hat{m}_{i,\cdot} \in \mathbb{R}^{p'}$ for observation $i$. The former vector uses clean data, while the latter uses cleaned data. In particular, $\tilde{m}_{i,\cdot} = m(W_{i,\cdot}, b^{\text{SIGNAL}})$, and $\hat{m}_{i,\cdot} = m(W_{i,\cdot}, b^{\text{NOISE}})$ overwriting $Z_{i,\cdot}$ with $\hat{A}_{i,\cdot}$. We concatenate the vectors $\tilde{m}_{i,\cdot}$ as rows in the matrix $\tilde{\boldsymbol{M}}$. We concatenate the vectors $\hat{m}_{i,\cdot}$ as rows in the matrix $\hat{\boldsymbol{M}}$. We refer to these objects as the counterfactual matrices. We also use the counterfactual vectors to define the counterfactual moments $M^*, \hat{M} \in \mathbb{R}^{p'}$: $M^* = \frac{1}{2n} \sum_{i \in \text{TRAIN,TEST}} \alpha_0(W_{i,\cdot}) b\{A_{i,\cdot}^{(\text{LR})}\}$ and $\hat{M} = \frac{1}{n} \sum_{i \in \text{TRAIN}} \hat{m}_{i,\cdot}$. Finally, we introduce notation for the covariance matrices $\boldsymbol{G}^*, \hat{\boldsymbol{G}} \in \mathbb{R}^{p' \times p'}$: $\boldsymbol{G}^* = \frac{1}{2n} \sum_{i \in \text{TRAIN,TEST}} b\{A_{i,\cdot}^{(\text{LR})}\}^T b\{A_{i,\cdot}^{(\text{LR})}\}$ and $\hat{\boldsymbol{G}} = \frac{1}{n} b(\hat{\boldsymbol{A}}^{\text{TRAIN}})^T b(\hat{\boldsymbol{A}}^{\text{TRAIN}})$.

A desirable property is that data cleaning guarantees for the corrupted regressors imply data cleaning guarantees of the counterfactual moments. We refer to this property as data cleaning continuity, and verify that it holds for the leading examples.

**Assumption I.1** (Data cleaning continuity). *There exist $C_m', C_m'' < \infty$ such that (i) $\|\hat{\boldsymbol{M}} - \tilde{\boldsymbol{M}}\|_{2,\infty}^2 \leq C_m' \|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{2,\infty}^2$; (ii) $\max_{j \in [p']} |\tilde{m}_{ij}| \leq C_m''$.*

**Proposition I.1** (Verifying data cleaning continuity). *Suppose Assumption 5.1 holds. In Example E.1 with the interacted dictionary, $C_m' = 1$ and $C_m'' = \bar{A}$. In Example E.2 with the interacted dictionary, $C_m' = 1$ and $C_m'' = \bar{A}$ in the numerator and denominator. In Example E.3 with the identity dictionary, suppose the counterfactual policy is of the form $t : A_{i,\cdot} \mapsto t_1 \odot A_{i,\cdot} + t_2$ where $t_1, t_2 \in \mathbb{R}^p$. Then $C_m' = (\|t_1\|_{\max} + 1)^2$ and $C_m'' = (\|t_1\|_{\max} + 1)\bar{A} + \|t_2\|_{\max}$. In Example E.4 with the interacted quadratic dictionary, $C_m' = 4\bar{A}^2$ and $C_m'' = 2\bar{A}^2$.[10] In Example E.5 with the partially linear dictionary, $C_m' = 0$ and $C_m'' = 1$.[11] In Example E.6 with the partially linear dictionary, $C_m' = 0$ and $C_m'' = 1$ in the numerator and denominator. In Example E.7 with the interacted dictionary, $C_m' = 1$ and $C_m'' = \bar{A}$.[12]*

*Proof.* In Example E.1, write $m_{i,\cdot} = b(1, Z_{i,\cdot}) - b(0, Z_{i,\cdot}) = \{Z_{i,\cdot}, 0\} - \{0, Z_{i,\cdot}\} = (Z_{i,\cdot}, -Z_{i,\cdot})$.

---

[10]Likewise for any polynomial of $D_i$ interacted with $Z_{i,\cdot}$.

[11]Recall that, to estimate a weighted balancing weight, we propose estimating an unweighted balancing weight then applying the weighting. The verification here is for the unweighted balancing weight that will be weighted.

[12]Recall that, to estimate a local balancing weight, we propose estimating a global balancing weight then applying the localization. The verification here is for the global balancing weight that will be localized.

Hence $\hat{m}_{i,\cdot} = (\hat{A}_{i,\cdot}, -\hat{A}_{i,\cdot})$ and $\tilde{m}_{i,\cdot} = (A_{i,\cdot}, -A_{i,\cdot})$. Example E.2 is analogous. In Example E.3, write $m_{i,\cdot} = b\{t(Z_{i,\cdot})\} - b(Z_{i,\cdot}) = t(Z_{i,\cdot}) - Z_{i,\cdot} = t_1 \odot Z_{i,\cdot} + t_2 - Z_{i,\cdot} = [(t_1 - \mathbb{1}^T) \odot Z_{i,\cdot}] + t_2$. Hence $\hat{m}_{i,\cdot} = [(t_1 - \mathbb{1}^T) \odot \hat{A}_{i,\cdot}] + t_2$ and $\tilde{m}_{i,\cdot} = [(t_1 - \mathbb{1}^T) \odot A_{i,\cdot}] + t_2$. In Example E.4, write $m_{i,\cdot} = \nabla_d b(D_i, Z_{i,\cdot}) = \nabla_d (1, D_i, D_i^2, Z_{i,\cdot}, D_i Z_{i,\cdot}, D_i^2 Z_{i,\cdot}) = (0, 1, 2D_i, 0, Z_{i,\cdot}, 2D_i Z_{i,\cdot})$. Hence $\hat{m}_{i,\cdot} = (0, 1, 2D_i, 0, \hat{A}_{i,\cdot}, 2D_i \hat{A}_{i,\cdot})$ and $\tilde{m}_{i,\cdot} = (0, 1, 2D_i, 0, A_{i,\cdot}, 2D_i A_{i,\cdot})$. In Example E.5, let $b(D_i, Z_{i,\cdot}) = \{D_i, \tilde{b}(Z_{i,\cdot})\}$. Write $m_{i,\cdot} = b(1, Z_{i,\cdot}) - b(0, Z_{i,\cdot}) = \{1, \tilde{b}(Z_{i,\cdot})\} - \{0, \tilde{b}(Z_{i,\cdot})\} = (1, 0, ..., 0)$. Hence $\hat{m}_{i,\cdot} = (1, 0, ..., 0)$ and $\tilde{m}_{i,\cdot} = (1, 0, ..., 0)$. Example E.6 is analogous to Example E.5. Example E.7 is analogous to Example E.1. $\qquad\square$

**Properties**. The error-in-variable balancing weight confers balance and equivalence.

**Proposition I.2** (Finite sample balance). *For any finite training sample size $n$, and any dictionary $b$, the coefficient $\hat{\eta}$ balances the cleaned actual regressors with the corresponding cleaned counterfactuals in the sense that $\frac{1}{n} \sum_{i \in \text{TRAIN}} b(\hat{A}_{i,\cdot}) \cdot \hat{\omega}_i = \frac{1}{n} \sum_{i \in \text{TRAIN}} \hat{m}_{i,\cdot}$ where $\hat{\omega}_i \in \mathbb{R}$ are balancing weights computed from $\hat{\eta}$: for each $i \in \text{TRAIN}$, $\hat{\omega}_i = b(\hat{A}_{i,\cdot})\hat{\eta}$.*

*Proof.* $\frac{1}{n} \sum_{i \in \text{TRAIN}} b(\hat{A}_{i,\cdot})^T b(\hat{A}_{i,\cdot})\hat{\eta} = \hat{\boldsymbol{G}}\hat{\eta} = \hat{M}^T = \frac{1}{n} \sum_{i \in \text{TRAIN}} (\hat{m}_{i,\cdot})^T$. $\qquad\square$

In words, $\hat{\eta}$ serves to balance actual observations with counterfactual queries.

*Proof of Proposition 4.3.* By Proposition I.2, $\frac{1}{n} \sum_{i \in \text{TRAIN}} b(D_i, \hat{A}_{i,\cdot}) \cdot \hat{\omega}_i = \frac{1}{n} \sum_{i \in \text{TRAIN}} \hat{m}_{i,\cdot}$. By Proposition I.1, $\hat{m}_{i,\cdot} = (\hat{A}_{i,\cdot}, -\hat{A}_{i,\cdot})$. Notice that $b(D_i, \hat{A}_{i,\cdot}) = \{D_i \hat{A}_{i,\cdot}, (1 - D_i)\hat{A}_{i,\cdot}\}$ and $\hat{\omega}_i = b(D_i, \hat{A}_{i,\cdot})\hat{\eta} = D_i \cdot \hat{\omega}_i^{\text{TREAT}} - (1 - D_i)\hat{\omega}_i^{\text{UNTREAT}}$. Therefore $b(D_i, \hat{A}_{i,\cdot}) \cdot \hat{\omega}_i = \{D_i \hat{A}_{i,\cdot} \cdot \hat{\omega}_i^{\text{TREAT}}, (1 - D_i)\hat{A}_{i,\cdot} \cdot (-\hat{\omega}_i^{\text{UNTREAT}})\}$. In summary, $\frac{1}{n} \sum_{i \in \text{TRAIN}} D_i \hat{A}_{i,\cdot} \cdot \hat{\omega}_i^{\text{TREAT}} = \frac{1}{n} \sum_{i \in \text{TRAIN}} \hat{A}_{i,\cdot}$ and $\frac{1}{n} \sum_{i \in \text{TRAIN}} (1 - D_i)\hat{A}_{i,\cdot} \cdot (-\hat{\omega}_i^{\text{UNTREAT}}) = \frac{1}{n} \sum_{i \in \text{TRAIN}} (-\hat{A}_{i,\cdot})$. $\qquad\square$

A well-known equivalence holds for treatment effects when using OLS with the interacted dictionary (without data cleaning).[13] We generalize it in three ways: (i) for our entire class of semiparametric and nonparametric estimands, (ii) for any square integrable dictionary, (iii) for estimation with or without data cleaning. We define, for $i \in \text{TRAIN}$, $\tilde{\gamma}(D_i, Z_{i,\cdot}) = b(D, \hat{X}_{i,\cdot})\hat{\beta}$ and $\tilde{\alpha}(D_i, Z_{i,\cdot}) = b(D, \hat{X}_{i,\cdot})\hat{\eta}$. We also define $\mathbb{E}_{\text{TRAIN}}[\cdot] = \frac{1}{m} \sum_{i \in \text{TRAIN}}[\cdot]$.

---

[13]We thank David Bruns-Smith and Avi Feller for suggesting this connection. See e.g. [Ben-Michael et al., 2021] for a recent summary, and references therein.

**Proposition I.3** (Equivalence in TRAIN). *If Assumption J.1 holds and b is square integrable, then the outcome, balancing weight, and doubly robust estimators coincide on the training set:*
$$\mathbb{E}_{\text{TRAIN}}[m(W_{i,\cdot}, \tilde{\gamma})] = \mathbb{E}_{\text{TRAIN}}[Y_i \tilde{\alpha}(D_i, Z_{i,\cdot})] = \mathbb{E}_{\text{TRAIN}}[m(W_{i,\cdot}, \tilde{\gamma}) + \tilde{\alpha}(D_i, Z_{i,\cdot})\{Y_i - \tilde{\gamma}(D_i, Z_{i,\cdot})\}].$$
*The same result holds without data cleaning.*

*Proof.* To prove the second equality, we appeal to the first order condition for $\hat{\eta}$: $\hat{\eta}^T \hat{\boldsymbol{G}} = \hat{M}$. Multiplying by $\hat{\beta}$, we have $\hat{\eta}^T \hat{\boldsymbol{G}} \hat{\beta} = \hat{\eta}^T \mathbb{E}_{\text{TRAIN}}[b(D, \hat{X}_{i,\cdot})^T b(D, \hat{X}_{i,\cdot})]\hat{\beta} = \mathbb{E}_{\text{TRAIN}}[\tilde{\alpha}(D_i, Z_{i,\cdot})\tilde{\gamma}(D_i, Z_{i,\cdot})]$ and $\hat{M}\hat{\beta} = \mathbb{E}_{\text{TRAIN}}[\hat{m}_{i,\cdot}]\hat{\beta} = \mathbb{E}_{\text{TRAIN}}[m(W_{i,\cdot}, \tilde{\gamma})]$. In summary, $\mathbb{E}_{\text{TRAIN}}[\tilde{\alpha}(D_i, Z_{i,\cdot})\tilde{\gamma}(D_i, Z_{i,\cdot})] = \mathbb{E}_{\text{TRAIN}}[m(W_{i,\cdot}, \tilde{\gamma})]$ which implies the result. To prove the first equality, we appeal to the first order condition for $\hat{\beta}$: $\hat{\beta}^T \hat{\boldsymbol{G}} = \mathbb{E}_{\text{TRAIN}}[Y_i b(D_i, \hat{X}_{i,\cdot})]$. Multiplying by $\hat{\eta}$, we have $\mathbb{E}_{\text{TRAIN}}[Y_i b(D_i, \hat{X}_{i,\cdot})]\hat{\eta} = \mathbb{E}_{\text{TRAIN}}[Y_i \tilde{\alpha}(D_i, Z_{i,\cdot})]$ and, appealing to the previous result, $\hat{\beta}^T \hat{\boldsymbol{G}} \hat{\eta} = \mathbb{E}_{\text{TRAIN}}[\tilde{\alpha}(D_i, Z_{i,\cdot})\tilde{\gamma}(D_i, Z_{i,\cdot})] = \mathbb{E}_{\text{TRAIN}}[m(W_{i,\cdot}, \tilde{\gamma})]$. $\square$

However, our estimator involves sample splitting and implicit data cleaning to break dependence, motivated by our goal of inference after data cleaning.

**Proposition I.4** (Non-equivalence in TEST). *If Assumption J.1 holds and b is square integrable, then the outcome, balancing weight, and doubly robust estimators generically do not coincide on the test set:* $\mathbb{E}_{\text{TEST}}[m(Z_{i,\cdot}, \hat{\gamma})] \neq \mathbb{E}_{\text{TEST}}[Y_i \hat{\alpha}(D_i, Z_{i,\cdot})] \neq \mathbb{E}_{\text{TEST}}[m(Z_{i,\cdot}, \hat{\gamma}) + \hat{\alpha}(D_i, Z_{i,\cdot})\{Y_i - \hat{\gamma}(D_i, Z_{i,\cdot})\}].$

*Proof.* The first order conditions for $(\hat{\beta}, \hat{\eta})$ hold for TRAIN after data cleaning. They do not hold for TEST, especially since we do not clean the test covariates. $\square$

**High probability events.** Define the events:

$$\mathcal{E}_6 = \left\{ \max_{j \in [p]} \left| \frac{1}{n} \sum_{i \in \text{TRAIN}} \{\tilde{m}_{ij} - \mathbb{E}[\tilde{m}_{ij}]\} \right| \leq C \cdot C''_m \sqrt{\frac{\log(np)}{n}} \right\};$$

$$\mathcal{E}_7 = \left\{ \max_{j \in [p]} \left| \frac{1}{2n} \sum_{i \in \text{TRAIN,TEST}} \{\alpha_0(W_{i,\cdot})A_{ij} - \mathbb{E}[\alpha_0(W_{i,\cdot})A_{ij}]\} \right| \leq C \cdot \bar{\alpha}\bar{A} \sqrt{\frac{\log(np)}{n}} \right\};$$

$$\mathcal{E}_8 = \left\{ \max_{j,k \in [p]} \left| \frac{1}{n} \sum_{i \in \text{TRAIN}} \{A_{ij}A_{ik} - \mathbb{E}[A_{ij}A_{ik}]\} \right| \leq C \cdot \bar{A}^2 \sqrt{\frac{\log(np)}{n}} \right\};$$

$$\mathcal{E}_9 = \left\{ \max_{j,k \in [p]} \left| \frac{1}{2n} \sum_{i \in \text{TRAIN,TEST}} \{A_{ij}A_{ik} - \mathbb{E}[A_{ij}A_{ik}]\} \right| \leq C \cdot \bar{A}^2 \sqrt{\frac{\log(np)}{n}} \right\}.$$

**Lemma I.1.** *Under Assumption I.1,* $\mathbb{P}(\mathcal{E}_6^c) \leq \frac{2}{n^{10}p^{10}}$. *Under Assumption 5.1 and* $\|\alpha_0\|_\infty \leq \bar{\alpha}$,
$\mathbb{P}(\mathcal{E}_7^c) \leq \frac{2}{n^{10}p^{10}}$. *Under Assumption 5.1,* $\mathbb{P}(\mathcal{E}_8^c) \leq \frac{2}{n^{10}p^{10}}$ *and* $\mathbb{P}(\mathcal{E}_9^c) \leq \frac{2}{n^{10}p^{10}}$.

*Proof.* By Assumption I.1, $\tilde{m}_{ij} \leq C_m''$. By Assumption 5.1, $|\alpha_0(W_{i,\cdot})A_{ij}| \leq \bar{\alpha}\bar{A}$ and $|A_{ij}A_{ik}| \leq \bar{A}^2$. For $\mathcal{E}_6, \mathcal{E}_7$ we appeal to Hoeffding for any $j \in [p]$, then take the union bound. For $\mathcal{E}_8, \mathcal{E}_9$ we appeal to Hoeffding for any $j, k \in [p]$, then take the union bound. $\square$

**Lemma I.2.** *Suppose Assumptions 5.1, I.1, J.1, and J.2 hold, and* $\|\alpha_0\|_\infty \leq \bar{\alpha}$. *Then*
$\|\hat{M} - M^*\|_{\max} | \{\mathcal{E}_6, \mathcal{E}_7\} \leq \Delta_M = \sqrt{C_m'}\frac{1}{\sqrt{n}}\|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{2,\infty} + C \cdot (C_m'' + \bar{\alpha}\bar{A})\sqrt{\frac{\ln(np)}{n}} + \bar{\alpha} \cdot \Delta_E$.

*Proof.* Write $\hat{M} - M^* = \sum_{k=1}^5 R^{(k)}$, where $\{R^{(k)}\}$ are below. By triangle inequality, it suffices to bound $R_j^{(k)}$. Write

$$\{R_j^{(1)}\}^2 = \left\{ \frac{1}{n}\sum_{i \in \text{TRAIN}} (\hat{m}_{ij} - \tilde{m}_{ij}) \right\}^2 \leq \frac{1}{n}\sum_{i \in \text{TRAIN}} (\hat{m}_{ij} - \tilde{m}_{ij})^2 \leq \frac{1}{n}\|\hat{\boldsymbol{M}} - \tilde{\boldsymbol{M}}\|_{2,\infty}^2 \leq \frac{1}{n}C_m'\|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{2,\infty}^2,$$

appealing to Assumption I.1. Write $R_j^{(2)} = \frac{1}{n}\sum_{i \in \text{TRAIN}}\{\tilde{m}_{ij} - \mathbb{E}[\tilde{m}_{ij}]\}$, then appeal to $\mathcal{E}_6$. Write $R_j^{(3)} = \frac{1}{n}\sum_{i \in \text{TRAIN}}\mathbb{E}[\tilde{m}_{ij}] - \frac{1}{2n}\sum_{i \in \text{TRAIN,TEST}}\mathbb{E}[\alpha_0(W_{i,\cdot})A_{ij}] = 0$ by Riesz representation and ex ante identical distribution of folds in the random partition (TRAIN, TEST). In particular, since $b_j^{\text{SIGNAL}} \in \mathbb{L}_2(\mathcal{W})$, $\mathbb{E}[\tilde{m}_{ij}] = \mathbb{E}[m(W_{i,\cdot}, b_j^{\text{SIGNAL}})] = \mathbb{E}[\alpha_0(W_{i,\cdot})b_j^{\text{SIGNAL}}(W_{i,\cdot})] = \mathbb{E}[\alpha_0(W_{i,\cdot})b_j(A_{i,\cdot})]$. Write $-R_j^{(4)} = \frac{1}{2n}\sum_{i \in \text{TRAIN,TEST}}\{\alpha_0(W_{i,\cdot})A_{ij} - \mathbb{E}[\alpha_0(W_{i,\cdot})A_{ij}]\}$ then appeal to $\mathcal{E}_7$. Write $|R_j^{(5)}| = \left| \frac{1}{2n}\sum_{i \in \text{TRAIN,TEST}}\alpha_0(W_{i,\cdot})E_{ij}^{(\text{LR})} \right| \leq \bar{\alpha} \cdot \Delta_E$ where $\alpha_0(W_{i,\cdot}) \leq \bar{\alpha}$. $\square$

**Lemma I.3.** *Suppose Assumptions 5.1 holds. Then* $\|\hat{\boldsymbol{G}} - \boldsymbol{G}^*\|_{\max} | \{\mathcal{E}_8, \mathcal{E}_9\} \leq \Delta_G$ *where*
$\Delta_G = (\bar{A} + \|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{2,\infty})\frac{1}{\sqrt{n}}\|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{2,\infty} + C \cdot \bar{A}^2\sqrt{\frac{\ln(np)}{n}} + C \cdot \bar{A}\Delta_E$.

*Proof.* Write $\hat{\boldsymbol{G}} - \boldsymbol{G}^* = \sum_{\ell=1}^7 S^{(\ell)}$, where $\{S^{(\ell)}\}$ are below. By triangle inequality, it suffices to bound $S_{jk}^{(\ell)}$. Write $S_{jk}^{(1)} = \frac{1}{n}\sum_{i \in \text{TRAIN}}\hat{A}_{ij}(\hat{A}_{ik} - A_{ik}) \leq \|\hat{\boldsymbol{A}}\|_{\max} \cdot \frac{1}{n}\sum_{i \in \text{TRAIN}}(\hat{A}_{ik} - A_{ik})$ hence $\{S_{jk}^{(1)}\}^2 \leq \|\hat{\boldsymbol{A}}\|_{\max}^2 \cdot \left\{ \frac{1}{n}\sum_{i \in \text{TRAIN}}(\hat{A}_{ik} - A_{ik}) \right\}^2 \leq \|\hat{\boldsymbol{A}}\|_{\max}^2 \cdot \frac{1}{n}\sum_{i \in \text{TRAIN}}(\hat{A}_{ik} - A_{ik})^2 \leq \|\hat{\boldsymbol{A}}\|_{\max}^2\frac{1}{n}\|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{2,\infty}^2$. Then use $\|\hat{\boldsymbol{A}}\|_{\max} \leq \|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{\max} + \|\boldsymbol{A}\|_{\max}$. Write $S_{jk}^{(2)} = \frac{1}{n}\sum_{i \in \text{TRAIN}}(\hat{A}_{ij} - A_{ij})A_{ik} \leq \bar{A} \cdot \frac{1}{n}\sum_{i \in \text{TRAIN}}(\hat{A}_{ij} - A_{ij})$. By a similar argument, $\{S_{jk}^{(2)}\}^2 \leq \bar{A}^2\frac{1}{n}\|\hat{\boldsymbol{A}} - \boldsymbol{A}\|_{2,\infty}^2$. Write $S_{jk}^{(3)} = \frac{1}{n}\sum_{i \in \text{TRAIN}}\{A_{ij}A_{ik} - \mathbb{E}[A_{ij}A_{ik}]\}$ then appeal to $\mathcal{E}_8$. Write $S_{jk}^{(4)} = \frac{1}{n}\sum_{i \in \text{TRAIN}}\mathbb{E}[A_{ij}A_{ik}] - \frac{1}{2n}\sum_{i \in \text{TRAIN,TEST}}\mathbb{E}[A_{ij}A_{ik}] = 0$ by ex ante identical distribution of folds in the random partition (TRAIN, TEST). Write $-S_{jk}^{(5)} = \frac{1}{2n}\sum_{i \in \text{TRAIN,TEST}}\{A_{ij}A_{ik} - \mathbb{E}[A_{ij}A_{ik}]\}$ then appeal to $\mathcal{E}_9$. By Assumption 5.1, $S_{jk}^{(6)} = \frac{1}{2n}\sum_{i \in \text{TRAIN,TEST}}A_{ij}E_{ik}^{(\text{LR})} \leq \bar{A}\Delta_E$.

By Assumption 5.1 and Lemma F.4. $S_{jk}^{(7)} = \frac{1}{2n} \sum_{i \in \text{TRAIN,TEST}} E^{(\text{LR})} A_{ik}^{(\text{LR})} \leq \|\boldsymbol{E}^{(\text{LR})}\|_{\max} \|\boldsymbol{A}^{(\text{LR})}\|_{\max} \leq 3\bar{A}\Delta_E$. $\qquad\square$

**Proposition I.5.** *If the conditions of Proposition C.2 hold, then* $\mathbb{P}(\tilde{\mathcal{E}}_5^c) \leq \frac{C}{n^{10}p^{10}}$.

*Proof.* The result follows from Lemmas I.2 and I.3. $\qquad\square$

# J   Data cleaning-adjusted confidence intervals details

**Riesz representation**. We generalize Assumptions 5.9 and 5.10 from the ATE example to the general case. In doing so, we also generalize the balancing weight to a Riesz representer.

**Assumption J.1** (Linearity and mean square continuity). *(i) The functional* $\gamma \mapsto \mathbb{E}[m(W_{i,\cdot}, \gamma)]$ *is linear. (ii) There exists* $\bar{Q} < \infty$ *and* $\bar{q} \in (0,1]$ *such that for all* $\gamma \in \Gamma$, $\mathbb{E}[m(W_{i,\cdot}, \gamma)^2] \leq \bar{Q} \cdot \{\mathbb{E}[\gamma(W_{i,\cdot})^2]\}^{\bar{q}}$.

These restrictions generalize the usual propensity score assumptions; Assumption 5.10 is a special case of Assumption J.1. Assumption J.1 implies that the balancing weight exists.

**Lemma J.1** ([Chernozhukov et al., 2022b]). *Suppose Assumption J.1 holds and* $\gamma_0 \in \Gamma$, *which may be imposed in estimation. Then there exists a Riesz representer* $\alpha_0 \in \mathbb{L}_2(\mathcal{W})$ *such that for all* $\gamma \in \Gamma$, $\mathbb{E}[m(W_{i,\cdot}, \gamma)] = \mathbb{E}[\alpha_0(W_{i,\cdot})\gamma(W_{i,\cdot})]$. *There exists a unique minimal Riesz representer* $\alpha_0^{\min} \in \Gamma$ *satisfying this equation. Moreover, denoting by* $\bar{M}$ *the operator norm of* $\gamma \mapsto \mathbb{E}[m(W_{i,\cdot}, \gamma)]$, $\{\mathbb{E}[\alpha_0^{\min}(W_{i,\cdot})^2]\}^{\frac{1}{2}} = \bar{M} \leq \bar{Q}^{\frac{1}{2}} < \infty$.

The balancing weight is a special case of a Riesz representer. Hereafter, we refer to the Riesz representer as a balancing weight nonetheless, since our estimator $\hat{\alpha}$ achieves balance across examples; see Proposition I.2, which generalizes Proposition 4.3. To lighten notation, we will typically consider the case where $\Gamma = \mathbb{L}_2(\mathcal{W})$ and $\alpha_0^{\min} = \alpha_0$. When we consider the more general case, as in Example E.4, we will use the richer notation.

We impose that $(\gamma_0, \alpha_0)$ do not vary across observations, generalizing familiar distribution shift assumptions. Assumption 5.9 is a special case of Assumption J.2, which we now state.

**Assumption J.2** (Marginal distribution shift). *For all observations* $i \in [n]$, *(i) the regression* $\gamma_0$ *does not vary:* $\mathbb{E}[\gamma_0(W_{i,\cdot})v(W_{i,\cdot})] = \mathbb{E}[Y_i v(W_{i,\cdot})]$ *for all* $v \in \mathbb{L}_2(\mathcal{W})$; *(ii) the Riesz representer* $\alpha_0$ *does not vary:* $\mathbb{E}[\alpha_0(W_{i,\cdot})u(W_{i,\cdot})] = \mathbb{E}[m(W_{i,\cdot}, u)]$ *for all* $u \in \mathbb{L}_2(\mathcal{W})$.

**Proposition J.1** (Verifying Assumptions J.1 and J.2). *Assumptions J.1 and J.2 hold under simple and interpretable conditions for the leading examples. Recall that $\|\alpha_0\|_\infty \leq \bar{\alpha}$, while $(\bar{Q}, \bar{q})$ are defined in Assumption J.1. In Example E.1, $\alpha_0(W_{i,\cdot}) = \frac{D_i}{\phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})} -$*
*$\frac{1-D_i}{1-\phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})}$, where $\phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) := \mathbb{E}[D_i | X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}]$. Suppose $0 < \underline{\phi} \leq \phi_0(X_{i,\cdot}, H_{i,\cdot}, \phi_{i,\cdot}) \leq$*
*$\bar{\phi} < 1$. Then $\bar{\alpha} = \frac{1}{\underline{\phi}} + \frac{1}{1-\bar{\phi}}$, $\bar{Q} = \frac{2}{\underline{\phi}} + \frac{2}{1-\bar{\phi}}$, and $\bar{q} = 1$. We impose that the outcome regression and treatment propensity score do not vary. In Example E.2, $\alpha_0(W_{i,\cdot}) =$*
*$\frac{U_i}{\phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})} - \frac{1-U_i}{1-\phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})}$, where $\phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) := \mathbb{E}[U_i | X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}]$. Suppose*
*$0 < \underline{\phi} \leq \phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) \leq \bar{\phi} < 1$. Then $\bar{\alpha} = \frac{1}{\underline{\phi}} + \frac{1}{1-\bar{\phi}}$, $\bar{Q} = \frac{2}{\underline{\phi}} + \frac{2}{1-\bar{\phi}}$, and $\bar{q} = 1$*
*. We impose that the outcome regression and instrument propensity score do not vary. In Example E.3, $\alpha_0(W_{i,\cdot}) = \omega(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) - 1$, where $\omega(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) = \frac{f\{t(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})\}}{f(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})}$.*
*Suppose $\omega(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) \leq \bar{\omega} < \infty$. Then $\bar{\alpha} = \bar{\omega} + 1$, $\bar{Q} = 2\bar{\omega} + 2$, and $\bar{q} = 1$. We impose that the outcome regression and covariate density ratio do not vary. In Example E.4,*
*$\alpha_0(W_{i,\cdot}) = -\nabla_d \ln f(D_i | X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})$. Suppose $-\nabla_d \ln f(D_i | X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) \leq \bar{f} < \infty$.*
*Then $\bar{\alpha} = \bar{f}$, $\bar{Q} = \bar{f}(\bar{\gamma} + \bar{\gamma}')$, and $\bar{q} = 1/2$ for $\Gamma$ that satisfies a Sobolev condition:*
*$\mathbb{E}[\{\nabla_d \gamma(D_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})\}^2] \leq \bar{\gamma}^2 < \infty$ and $\mathbb{E}[\{\partial_d^2 \gamma(D_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})\}^2] \leq (\bar{\gamma}')^2 < \infty$. We*
*impose that the outcome regression and conditional density of goods do not vary. In Example E.5, $\alpha_0(W_{i,\cdot}) = \ell_i \frac{D_i - \phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})}{\mathbb{E}[\{D_i - \phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})\}^2]}$, where $\phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) := \mathbb{E}[D_i | X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}]$.*
*Suppose $\mathbb{E}[\{D_i - \phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})\}^2] > \underline{\phi}$ and $|\ell_i| \leq \bar{\ell}$. Then $\bar{\alpha} = \frac{2\bar{\ell}\bar{A}}{\underline{\phi}}$, $\bar{Q} = \frac{4\bar{\ell}^2\bar{A}^2}{\underline{\phi}^2}$, and $\bar{q} = 1$*
*. We impose that the outcome regression and treatment regression do not vary. In Example E.6, $\alpha_0(W_{i,\cdot}) = \ell_i \frac{U_i - \phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})}{\mathbb{E}[\{U_i - \phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})\}^2]}$, where $\phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) := \mathbb{E}[U_i | X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}]$*
*. Suppose $\mathbb{E}[\{D_i - \phi_0(X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})\}^2] > \underline{\phi}$ and $|\ell_i| \leq \bar{\ell}$. Then $\bar{\alpha} = \frac{2\bar{\ell}\bar{A}}{\underline{\phi}}$, $\bar{Q} = \frac{4\bar{\ell}^2\bar{A}^2}{\underline{\phi}^2}$,*
*and $\bar{q} = 1$. We impose that the outcome regression and instrument regression do not vary. In Example E.7, $\alpha_0(W_{i,\cdot}) = \ell_h(V_i)\left\{\frac{D_i}{\phi_0(V_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})} - \frac{1-D_i}{1-\phi_0(V_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot})}\right\}$. Suppose*
*$0 < \underline{\phi} \leq \phi_0(V_i, X_{i,\cdot}, H_{i,\cdot}, \pi_{i,\cdot}) \leq \bar{\phi} < 1$ and other regularity conditions hold given in Lemma J.2 below. Then $\bar{\alpha}_h \leq C \cdot \frac{1}{h}\left(\frac{1}{\underline{\phi}} + \frac{1}{1-\bar{\phi}}\right)$, $\bar{Q}_h \leq C \cdot \frac{1}{h^2}\left(\frac{2}{\underline{\phi}} + \frac{2}{1-\bar{\phi}}\right)$, $\bar{q} = 1$. We impose that the outcome regression and treatment propensity score do not vary.*

*Proof of Proposition J.1.* The results follow from the law of iterated expectations and integration by parts. See [Chernozhukov et al., 2023, Lemmas S3 and S4] for mean square continuity of Example E.4 and [Chernozhukov et al., 2023, Theorem 2] for the characterization of $(\bar{\alpha}_h, \bar{Q}_h)$ with localization. □

**Nonparametrics**. A local functional $\theta_0^{\lim} \in \mathbb{R}$ is a scalar that takes the form $\theta_0^{\lim} = \lim_{h\to 0} \theta_0^h$, where $\theta_0^h = \frac{1}{n}\sum_{i=1}^n \theta_i^h$ and $\theta_i^h = \mathbb{E}[m_h(W_{i,\cdot}, \gamma_0)] = \mathbb{E}[\ell_h(W_{ij})m(W_{i,\cdot}, \gamma_0)]$. Here, $\ell_h$ is a Nadaraya Watson weighting with bandwidth $h$ and $W_{ij}$ is a scalar component of $W_{i,\cdot}$. $\theta_0^{\lim}$ is a nonparametric quantity that we approximate by the sequence $\{\theta_0^h\}$, incurring the nonparametric approximation error $\Delta_h = n^{1/2}\sigma^{-1}|\theta_0^h - \theta_0^{\lim}|$. Each $\theta_0^h$ can be analyzed like $\theta_0$ above as long as we keep track of how certain quantities depend on $h$

**Lemma J.2** (Theorem 2 of [Chernozhukov et al., 2023]). *If response noise has finite variance then $\bar{\sigma}^2 < \infty$. Suppose bounded moment and heteroscedasticity conditions hold. Then for global functionals $\xi/\sigma \lesssim \sigma \asymp \bar{M} < \infty$, $\xi, \chi \lesssim \bar{M}^2 \leq \bar{Q} < \infty$, and $\bar{\alpha} < \infty$. Suppose bounded moment, heteroscedasticity, density, and derivative conditions hold. Then for local functionals $\xi_h/\sigma_h \lesssim h^{-1/6}$, $\sigma_h \asymp \bar{M}_h \asymp h^{-1/2}$, $\xi_h \lesssim h^{-2/3}$, $\chi_h \lesssim h^{-3/4}$, $\bar{\alpha}_h \lesssim h^{-1}$, $\bar{Q}_h \lesssim h^{-2}$, and $\Delta_h \lesssim n^{1/2}h^{\nu+1/2}$ where $\nu$ is the order of differentiability of the kernel $K$.*

Equipped with this lemma, we prove validity of the data cleaning-adjusted confidence interval for nonparametric quantities.

**Corollary J.1** (Confidence interval coverage). *Suppose the conditions of Corollary 5.3 and Lemma J.2 hold. Update the rate conditions to be (i) bandwidth regularity: $n^{-1/2}h^{-3/2} \to 0$ and $\Delta_h \to 0$; (ii) error-in-variable regression rate: $(h^{-1} + \bar{\alpha}')\{\mathcal{R}(\hat{\gamma})\}^{\bar{q}/2} \to 0$; (iii) error-in-variable balancing weight rate: $\bar{\sigma}h^{-1}\{\mathcal{R}(\hat{\alpha})\}^{1/2} \to 0$; (iv) product of rates is fast: $h^{-1/2}\{n\mathcal{R}(\hat{\gamma})\mathcal{R}(\hat{\alpha})\}^{1/2} \to 0$. Then the conclusions of Corollary 5.3 hold, replacing $(\hat{\theta}, \theta_0)$ with $(\hat{\theta}^h, \theta_0^{\lim})$.*

*Proof of Corollary J.1.* By Lemma J.2, the regularity condition on moments is $\{(\kappa/\sigma)^3 + \zeta^2\}n^{-1/2} \lesssim \{(h^{-1/6})^3 + (h^{-3/4})^2\}n^{-1/2} \lesssim h^{-3/2}n^{-1/2}$. By Lemma J.2, the first learning rate condition is $(\bar{Q}^{1/2} + \bar{\alpha}/\sigma + \bar{\alpha}')\{\mathcal{R}(\hat{\gamma})\}^{1/2} \lesssim (h^{-1} + h^{-1}/h^{-1/2} + \bar{\alpha}')\{\mathcal{R}(\hat{\gamma})\}^{1/2} \lesssim (h^{-1} + \bar{\alpha}')\{\mathcal{R}(\hat{\gamma})\}^{1/2}$. By [Chernozhukov et al., 2023, Lemma S.9], the second learning rate condition is $\bar{\sigma}\{\mathcal{R}(\hat{\alpha}^h)\}^{1/2} \lesssim \bar{\sigma}h^{-1}\{\mathcal{R}(\hat{\alpha})\}^{1/2}$. By Lemma J.2 and [Chernozhukov et al., 2023, Lemma S.9], the third learning rate condition is $\{n\mathcal{R}(\hat{\gamma})\mathcal{R}(\hat{\alpha}^h)\}^{1/2}/\sigma \lesssim \{n\mathcal{R}(\hat{\gamma})\mathcal{R}(\hat{\alpha})\}^{1/2}h^{-1}/h^{-1/2} = h^{-1/2}\{n\mathcal{R}(\hat{\gamma})\mathcal{R}(\hat{\alpha})\}^{1/2}$.

$\square$

# K  Nonlinear factor model

We denote $\mathcal{R}_\gamma = \mathcal{R}(\hat{\gamma}_\ell)$ and $\mathcal{R}_\alpha = \mathcal{R}(\hat{\alpha}_\ell)$ to lighten notation. The distinction between $n$ and $m = \frac{n}{2}$ is irrelevant in the context of $(\mathcal{R}_\gamma, \mathcal{R}_\alpha)$ due to the absolute constant $C$.

**Lemma K.1** ([Agarwal et al., 2021]). *Suppose Assumption 5.11 holds for some fixed* $\mathcal{H}(q, S, C_H)$. *Then for any small* $\delta > 0$, *there exists* $\boldsymbol{A}^{(lr)}$ *such that* $r = rank(\boldsymbol{A}^{(lr)}) \leq C \cdot \delta^{-q}$ *and* $\Delta_E = \|\boldsymbol{A} - \boldsymbol{A}^{(lr)}\|_{\max} \leq C_H \cdot \delta^S$, *where* $C$ *may depend on* $(q, S)$.

*Proof of Corollary 5.4.* From Lemma K.1, $r \leq C \cdot \delta^{-q}$ and $\Delta_E \leq C \cdot \delta^S$. The conditions of Corollary 5.4 imply $(\sigma, \bar{\sigma}, \bar{\alpha}, \bar{\alpha}', \bar{Q})$ are irrelevant. We verify simplified rate conditions from Corollary 5.3: $\mathcal{R}_\gamma \to 0$, $\mathcal{R}_\alpha \to 0$, $\sqrt{n\mathcal{R}_\gamma\mathcal{R}_\alpha} \to 0$. The relevant terms in $(\mathcal{R}_\gamma, \mathcal{R}_\alpha)$ simplify as well. From Theorem 5.2, these are $\mathcal{R}_\gamma \leq Cr^3\left\{\frac{1}{n} + \frac{p}{n^2} + \frac{1}{p} + \left(1 + \frac{p}{n}\right)\Delta_E^2 + p\Delta_E^4\right\}$. From Theorem 5.3, these are $\mathcal{R}_\alpha \leq Cr^5\left\{\frac{1}{n} + \frac{1}{p} + \frac{p}{n^2} + \frac{n}{p^2} + \left(1 + \frac{p}{n} + \frac{n}{p}\right)\Delta_E^2 + (n+p)\Delta_E^4 + np\Delta_E^6\right\}$.

Suppose $n = p^\upsilon$ with $\upsilon \geq 1$. Then $\mathcal{R}_\gamma \leq Cr^3\left(\frac{1}{p} + \Delta_E^2 + p\Delta_E^4\right) \leq C\delta^{-3q}\left(\frac{1}{p} + \delta^{2S} + p\delta^{4S}\right)$. The three terms are equalized with $\delta^{2S} = p^{-1}$. Hence $\mathcal{R}_\gamma \leq C\delta^{-3q}\frac{1}{p} = Cp^{\frac{3q}{2S}}\frac{1}{p} = Cp^{\frac{3q}{2S}-1}$. Similarly $\mathcal{R}_\alpha \leq Cr^5\left(\frac{n}{p^2} + \frac{n}{p}\Delta_E^2 + n\Delta_E^4 + np\Delta_E^6\right) \leq C\delta^{-5q}\left(\frac{n}{p^2} + \frac{n}{p}\delta^{2S} + n\delta^{4S} + np\delta^{6S}\right)$. The four terms are equalized with $\delta^{2S} = p^{-1}$. Hence $\mathcal{R}_\alpha \leq C\delta^{-5q}\frac{n}{p^2} = Cp^{\frac{5q}{2S}}\frac{n}{p^2} = Cp^{\frac{5q}{2S}-2}n$. To satisfy $\mathcal{R}_\gamma \leq \mathcal{R}_\alpha \to 0$, it suffices that $\frac{q}{S} < \frac{2}{5}(2 - \upsilon)$. To satisfy $\sqrt{n\mathcal{R}_\gamma\mathcal{R}_\alpha} \to 0$, it suffices that $\frac{q}{S} < \frac{1}{2}\left(\frac{3}{2} - \upsilon\right)$. In summary, a sufficient generalized factor model is one in which $\frac{q}{S} < \frac{2}{5}(2 - \upsilon) \wedge \frac{1}{2}\left(\frac{3}{2} - \upsilon\right)$ where $\upsilon \leq \frac{3}{2}$. The latter condition binds for $1 \leq \upsilon \leq \frac{3}{2}$.

If instead $n = p^\upsilon$ with $\upsilon \geq 1$, then a similar argument arrives at the same condition. $\square$

When using a polynomial dictionary, the relevant terms in $(\mathcal{R}_\gamma, \mathcal{R}_\alpha)$ are as above, instead using $(r', \Delta_E')$. Let $q' = d_{\max} \cdot q$. Then $r' \leq C \cdot r^{d_{\max}} \leq C \cdot \delta^{-qd_{\max}} = C \cdot \delta^{-q'}$ and $\Delta_E' \leq C\bar{A}^{d_{\max}} \cdot d_{\max}\Delta_E \leq C \cdot \Delta_E \leq C \cdot \delta^S$. Hence the proof of Corollary 5.4 remains the same.

# L  Simulation and application

**Simulation design**. We focus on average treatment effect (ATE) with corrupted covariates (Example E.1). A single observation is a triple $(Y_i, D_i, Z_{i,\cdot})$ for outcome, treatment, and corrupted covariates where $Y \in \mathbb{R}$, $D_i \in \{0, 1\}$, and $Z_{i,\cdot} \in \mathbb{R}^p$ are generated as follows.

First, we generate signal from a factor model. Sample $\boldsymbol{U} \sim \mathcal{N}(0, \boldsymbol{I}_{n \times r})$ and $\boldsymbol{V} \sim \mathcal{N}(0, \boldsymbol{I}_{p \times r})$. Then set $\boldsymbol{X} = \boldsymbol{U}\boldsymbol{V}^T$. By construction,

$$\mathbb{E}[X_{ij}] = \mathbb{E}\left[\sum_{s=1}^{r} U_{is} V_{sj}\right] = \sum_{s=1}^{r} \mathbb{E}[U_{is}]\mathbb{E}[V_{sj}] = 0$$

and $\mathbb{V}[X_{ij}] = \mathbb{V}[\sum_{s=1}^{r} U_{is} V_{sj}] = \sum_{s=1}^{r} \mathbb{V}[U_{is}]\mathbb{V}[V_{sj}] = r$.

Draw response noise as $\varepsilon_i \overset{i.i.d.}{\sim} \mathcal{N}(0, 1)$. Define the vector $\beta \in \mathbb{R}^p$ by $\beta_j = j^{-2}$. Then set $D_i \sim \text{Bernoulli}\{\Lambda(0.25 X^T \beta)\}$ and $Y_i = 2.2 D_i + 1.2 X_{i,\cdot}\beta + D_i X_{i1} + \varepsilon_i$ where $\Lambda(t) = (0.95 - 0.05)\frac{\exp(t)}{1+\exp(t)} + 0.05$ is the truncated logistic function. The ATE is $\theta_0 = 2.2$.

We observe the corrupted covariate $Z_{i,\cdot} = [X_{i,\cdot} + H_{i,\cdot}] \odot \pi_{i,\cdot}$. $H_{ij} \overset{i.i.d.}{\sim} F_H$ is drawn i.i.d. with mean zero and variance $\sigma_H^2$. $\pi_{ij}$ is 1 with probability $\rho$ and NA with probability $1 - \rho$. We consider different choices of the measurement error distribution $F_H$ to corresponding to classical measurement error, discretization, and differential privacy. In summary, the three data corruption parameters are $(F_H, \sigma_H, \rho)$. The remaining design parameters are $(n, p, r)$ corresponding to the sample size, dimension of covariates, and rank of the signal.

For classical measurement error, $F_H = \mathcal{N}(0, \sigma_H^2)$. For discretization, we generate $Z_{ij} = sign(X_{ij}) \cdot Poisson(|X_{ij}|)$ and implicitly define $F_H$ by $H_{ij} = Z_{ij} - X_{ij}$. Note that $\mathbb{E}[Z_{ij}|X_{ij}] = sign(X_{ij})\mathbb{E}[Poisson(|X_{ij}|)|X_{ij}] = X_{ij}$ as desired. Below, we show that $\sigma_H^2 = \mathbb{V}[H_{ij}] = 1.7$ in this construction. For differential privacy, $F_H = Laplace(0, \frac{\sigma_H}{\sqrt{2}})$.

**Proposition L.1** (Discretization variance). *Given some random variable $X$, define $P = Poisson(|X|)$, $Z = sign(X) \cdot P$, and $H = Z - X$. Then $\mathbb{E}[H] = 0$ and $\mathbb{V}[H] = \mathbb{E}[|X|]$.*

*Proof.* To begin, write $\mathbb{E}[Z|X] = sign(X) \cdot \mathbb{E}[P|X] = X$. By the law of total variance, $\mathbb{V}[H] = \mathbb{E}[\mathbb{V}[H|X]] + \mathbb{V}[\mathbb{E}[H|X]]$. In the latter term, $\mathbb{E}[H|X] = \mathbb{E}[Z - X|X] = 0$. In the former term, $\mathbb{V}[H|X] = \mathbb{V}[Z|X] = \mathbb{E}[Z^2|X] - \{\mathbb{E}[Z|X]\}^2$. Moreover, $\mathbb{E}[Z^2|X] = \mathbb{E}[P^2|X] = \mathbb{V}[P|X] + \{\mathbb{E}[P|X]\}^2 = |X| + X^2$. In summary, $\mathbb{V}[H|X] = |X|$. $\square$

**Robustness to data dimensions**. In the main text, each sample from the simulated data generating process produces a matrix of covariates $\boldsymbol{X} \in \mathbb{R}^{100 \times 100}$ with rank $r = 5$. How robust is our end-to-end procedure across realistic dimensions of economic data? We consider the following variations: $\boldsymbol{X} \in \mathbb{R}^{50 \times 200}$, $\mathbb{R}^{100 \times 100}$, $\mathbb{R}^{200 \times 50}$, $\mathbb{R}^{500 \times 20}$, and $\mathbb{R}^{1000 \times 10}$. For each choice of $(n, p)$, we set the rank $r = \{\min(n, p)\}^{1/3}$. Across data dimensions, we

introduce measurement error with the fixed noise-to-signal ratio of 20%. We consider the oracle tuning of the PCA hyperparameter $k = r$.

Table 1 quantifies coverage performance. Different rows correspond to different data dimensions. We record the average point estimates, which are close to $\theta_0 = 2.2$. Next, we record the average standard errors, which adaptively decrease in length for larger sample sizes. These confidence intervals are the correct length, since coverage is close to the nominal level.

| Meas. Err. | $n$ | $p$ | ATE | SE | 80% CI | 95% CI |
|---|---|---|---|---|---|---|
| 20% | 50 | 200 | 2.21 | 0.56 | 0.82 | 0.96 |
| 20% | 100 | 100 | 2.21 | 0.35 | 0.82 | 0.95 |
| 20% | 200 | 50 | 2.22 | 0.21 | 0.81 | 0.95 |
| 20% | 500 | 20 | 2.23 | 0.12 | 0.78 | 0.94 |
| 20% | 1000 | 10 | 2.23 | 0.08 | 0.76 | 0.93 |
| 20% | 722 | 30 | 2.26 | 0.12 | 0.78 | 0.92 |

Table 1: Our approach adapts to data shape

We repeat this exercise for the simulated data generating process with $\boldsymbol{X} \in \mathbb{R}^{722 \times 30}$ and rank $r = 5$. Table 1 confirms that our procedure attains nearly nominal coverage.

**Can data corruption flip signs**? In the main text, we show that for the simulated data generating process with $\boldsymbol{X} \in \mathbb{R}^{100 \times 100}$ and rank $r = 5$, OLS performs well with clean data and performs poorly with corrupted data. We investigate two follow-up questions. First, can data corruption flip the sign of OLS estimates, i.e. can it lead to negative point estimates when the average treatment effect is $\theta_0 = 2.2 > 0$? Second, can can data corruption flip the sign of OLS and 2SLS estimates in scenarios more similar to our real world example?

We find that data corruption can flip the sign of OLS estimates some of the time. In particular, measurement error with a 20% noise-to-signal ratio is enough to flip the sign roughly one quarter of the time. We repeat this exercise for the simulated data generating process with $\boldsymbol{X} \in \mathbb{R}^{722 \times 30}$ and rank $r = 5$. Flipping signs requires not only 20% measurement error but also 10% missingness. A similar fraction of OLS estimates have flipped signs.

Finally, we conduct a semi-synthetic sign flipping exercise. We consider the covariates of [Autor et al., 2013] at the commuting zone level. Rather than a synthetic ATE, the estimand is the actual effect of import competition on manufacturing employment in a partially linear instrumental variable model. Flipping signs requires not only 20% measurement error but also 20% missingness. In this thought experiment, we take the reported effect from [Autor et al., 2013] as the ground truth, we take the data set from [Autor et al., 2013] as clean data, and we generate synthetic measurement error and missingness.

We summarize the results of these sign flipping exercises in Table 2. The rows correspond to (i) synthetic data with $\boldsymbol{X} \in \mathbb{R}^{100\times100}$; (ii) synthetic data with $\boldsymbol{X} \in \mathbb{R}^{722\times30}$; and (iii) semi-synthetic data from

| Data | Meas. Err. | Miss. Val. | Sign Flip |
|------|-----------|-----------|-----------|
| $100 \times 100$ | 20% | 0% | 27% |
| $722 \times 30$ | 20% | 10% | 22% |
| Census | 20% | 20% | 9% |

Table 2: Data corruption can flip signs

[Autor et al., 2013]. We interpret the OLS and TSLS results as motivation for data cleaning before data analysis. Our procedure may be viewed as an extension of OLS and TSLS with simple data cleaning that we subsequently account for in our confidence intervals.

**Empirical application**. The variable definitions follow [Autor et al., 2013]. In the authors' original specification [Autor et al., 2013, Table 3, column 6], $X_{i,\cdot} \in \mathbb{R}^{14}$ consists of: a constant, an indicator for the 2000-2007 period, percentage of employment in manufacturing, percentage of college educated population, percentage of foreign-born population, percentage of employment among women, percentage of employment in routine occupations, average offshorability index of occupations, and Census division dummies.

In our augmented specification, $X_{i,\cdot} \in \mathbb{R}^{30}$ consists of variables from the original specification as well as additional variables in [Autor et al., 2013, Appendix Table 2]. These include percentages of the working age population: employed in manufacturing, employed in non-manufacturing, unemployed, not in the labor force, receiving disability benefits; average log weekly wages: manufacturing, non-manufactuing; average benefits per capita: individual transfers, retirement, disability, medical, federal income assistance, unemployment, TAA; and average household income per working age adult: total, wage and salary.

Figure 6 provides analogous results to Figure 5, where now we center and scale the covariates $X_{i,\cdot} \in \mathbb{R}^{30}$ before conducting the exercise. We arrive at similar conclusions.
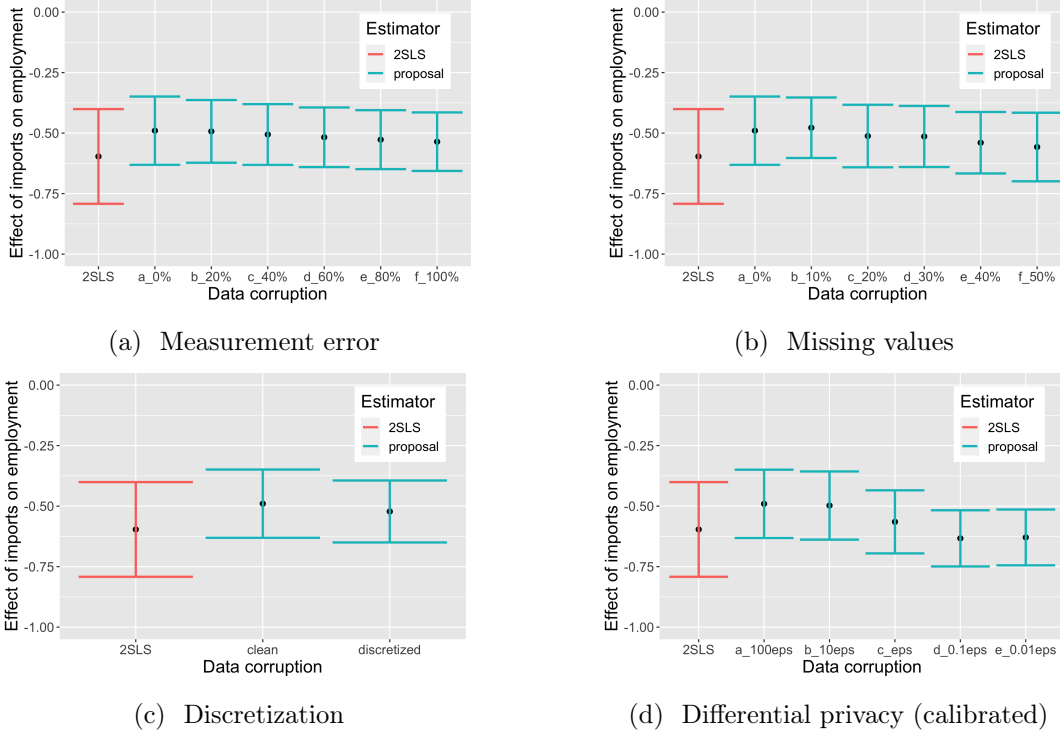
(a) Measurement error      (b) Missing values

(c) Discretization      (d) Differential privacy (calibrated)

Figure 6: Standardizing covariates before synthetic corruption

# M    Formalizing privacy

*Proof of Proposition 6.1.* Fix the commuting zone $i \in [n]$. We refer to the construction of the summary statistic $X_{ij} = f_j(\boldsymbol{M}^{(i)}) = \frac{1}{L_i} \sum_{\ell=1}^{L_i} M_{\ell j}^{(i)}$ as the $j$-th query $f_j$ about $\boldsymbol{M}^{(i)}$, where $j \in [p]$. To ensure privacy level $\epsilon_j$ for query $f_j$, a possible mechanism is, according to [Dwork et al., 2006, Proposition 3.3], $Z_{ij} = X_{ij} + H_{ij}$, where $X_{ij} = f_j(\boldsymbol{M}^{(i)})$ and $H_{ij} \overset{i.i.d.}{\sim} \text{Laplace}(S(f_j)/\epsilon_j)$. $S(f_j)$ is the sensitivity of the query, to which we return below. If no individual appears in two commuting zones, the Bureau can achieve privacy level $\epsilon$ while publishing all $j \in [p]$ variables for this commuting zone by setting $\epsilon_j = \epsilon/p$.

We wish to characterize the resulting subexponential parameters. They are, by independence of the Laplacians, $K_a = \|H_{i,\cdot}\|_{\psi_a} = \max_{j \in [p]} \|H_{ij}\|_{\psi_a} = \max_j \sqrt{2} \cdot S(f_j)/\epsilon_j = \sqrt{2}/\epsilon \cdot p \max_j S(f_j)$ and $\kappa^2 = \|\mathbb{E}[H_{i,\cdot}^T H_{i,\cdot}]\|_{op} = \max_{ij} \mathbb{V}(H_{ij}) = 2 \max_j S(f_j)^2/\epsilon_j^2 = 2/\epsilon^2 \cdot p^2 \max_j S(f_j)^2$. What remains is to define and characterize the the sensitivity $S(f_j)$. The sensitivity of the query $f_j$ is the most that the query may vary if one individual in the microdata were replaced. Formally, $\max_{\boldsymbol{M}^{(i)}, \boldsymbol{M}^{(i')}} |f_j(\boldsymbol{M}^{(i)}) - f_j(\boldsymbol{M}^{(i')})| \leq S(f_j)$ where $\boldsymbol{M}^{(i)}$ and $\boldsymbol{M}^{(i')}$ are two possible data sets of $L_i$ individuals that differ in one individual.

In what follows, we suppress indexing by $i$ to lighten notation. By hypothesis, each entry of microdata is bounded: $|M_{\ell j}| \leq \bar{A}$. This fact, together with the fact that the query $f_j$ is a sample mean, provides a bound on the sensitivity $S(f_j)$. To begin, write $f_j(\boldsymbol{M}) = \frac{1}{L}\left\{\sum_{\ell=1}^{L} M_{\ell j}\right\} = \frac{1}{L}\left\{\sum_{\ell=1}^{L-1} M_{\ell j} + M_{\ell L}\right\}$. Therefore without loss of generality $f_j(\boldsymbol{M}) - f_j(\boldsymbol{M}') = \frac{1}{L}(M_{\ell L} - M'_{\ell L})$ and hence $S(f_j) = \max_{\boldsymbol{M},\boldsymbol{M}'} |f_j(\boldsymbol{M}) - f_j(\boldsymbol{M}')| = \max_{\boldsymbol{M},\boldsymbol{M}'} \left|\frac{1}{L}(M_{\ell L} - M'_{\ell L})\right| \leq \frac{2\bar{A}}{L}$. Lemma M.1 ensures that privacy is preserved. $\qquad\square$

**Lemma M.1.** *Suppose the conditions of Proposition 6.1 hold. If $Z_{i,\cdot} = X_{i,\cdot} + H_{i,\cdot}$ confers $\epsilon$ differential privacy, then $\hat{X}_{i,\cdot}$ remains $\epsilon$ differentially private.*

*Proof.* Extending the notation from the previous proof, let $Z_{i,\cdot} = \mathcal{M}(\boldsymbol{M}^{(i)}) = X_{i,\cdot} + H_{i,\cdot}$ and $Z_{i',\cdot} = \mathcal{M}(\boldsymbol{M}^{(i')}) = X_{i',\cdot} + H_{i,\cdot}$. Recall that $\hat{X}_{i,\cdot}$ is a function of $Z_{i,\cdot}$ and $\{Z_{j,\cdot}\}_{j\neq i}$, i.e. $\hat{X}_{i,\cdot} = \text{CLEAN}[Z_{i,\cdot}; \{Z_{j,\cdot}\}_{j\neq i}]$. Analogously, $\hat{X}_{i',\cdot} = \text{CLEAN}[Z_{i',\cdot}; \{Z_{j,\cdot}\}_{j\neq i}]$.

By hypothesis, for any event $E$, $\frac{\mathbb{P}_{H_{i,\cdot}}(Z_{i,\cdot}\in E)}{\mathbb{P}_{H_{i,\cdot}}(Z_{i',\cdot}\in E)} \leq e^{\epsilon}$ where the subscript emphasizes the source of randomness. We wish to show that, for any event $F$, $\frac{\mathbb{P}_{H_{i,\cdot},\{H_{j,\cdot}\}_{j\neq i}}(\hat{X}_{i,\cdot}\in F)}{\mathbb{P}_{H_{i,\cdot},\{H_{j,\cdot}\}_{j\neq i}}(\hat{X}_{i',\cdot}\in F)} \leq e^{\epsilon}$. Fix $F$ and define $G := \{z \in \mathbb{R} : \text{CLEAN}[z; \{Z_{j,\cdot}\}_{j\neq i}] \in F\}$. Then

$$\mathbb{P}_{H_{i,\cdot},\{H_{j,\cdot}\}_{j\neq i}}(\hat{X}_{i,\cdot}\in F) = \mathbb{P}_{H_{i,\cdot},\{H_{j,\cdot}\}_{j\neq i}}(\text{CLEAN}[Z_{i,\cdot}; \{Z_{j,\cdot}\}_{j\neq i}]\in F)$$

$$= \mathbb{P}_{H_{i,\cdot},\{H_{j,\cdot}\}_{j\neq i}}(Z_{i,\cdot}\in G) = \mathbb{E}_{\{H_{j,\cdot}\}_{j\neq i}}[\mathbb{E}_{H_{i,\cdot}}[\mathbb{1}(Z_{i,\cdot}\in G)|\{H_{j,\cdot}\}_{j\neq i}]].$$

Moreover, $\mathbb{E}_{H_{i,\cdot}}[\mathbb{1}(Z_{i,\cdot}\in G)|\{H_{j,\cdot}\}_{j\neq i}] = \mathbb{P}_{H_{i,\cdot}}(Z_{i,\cdot}\in G) \leq e^{\epsilon} \cdot \mathbb{P}_{H_{i,\cdot}}(Z_{i',\cdot}\in G) = e^{\epsilon} \cdot \mathbb{E}_{H_{i,\cdot}}[\mathbb{1}(Z_{i',\cdot}\in G)|\{H_{j,\cdot}\}_{j\neq i}]$. Reversing the steps above yields the conclusion. $\qquad\square$

**Lemma M.2** ([Bun and Steinke, 2016])**.** *If $\mathcal{M}$ is differentially private with parameter $\epsilon$, then it is zero concentrated differentially private with parameter $\rho = \epsilon^2/2$.*

The Bureau's global privacy loss budget for people, in terms of zero concentrated differential privacy, is 2.56 in 2020 Census redistricting data (P.L 94-171) [Abowd et al., 2022]. Of this budget, $447/4,099$ is for counties. We use these numbers to calibrate a realistic privacy budget of $\rho = 2.56 \cdot 447/4,099$ for a hypothetical data release concerning commuting zones. Lemma M.2 demonstrates that a sufficient degree of differential privacy is $\epsilon = (2\rho)^{1/2}$.