

## Assessing Covariate Balance with Small Sample Sizes

George Hripcsak, MD<sup>1,2</sup>, Linying Zhang, PhD<sup>2,3</sup>, Kelly Li<sup>2,4</sup>, Marc A. Suchard, MD, PhD<sup>2,4,5</sup>,  
Patrick B. Ryan, PhD<sup>2,6</sup>, Martijn J. Schuemie, PhD<sup>2,6</sup>

<sup>1</sup> Department of Biomedical Informatics, Columbia University Irving Medical Center, New York, NY, USA

<sup>2</sup> Observational Health Data Science and Informatics, New York, NY, USA

<sup>3</sup> Institute for Informatics, Data Science and Biostatistics, Washington University in St. Louis, St. Louis, MO, USA

<sup>4</sup> Department of Biostatistics, University of California, Los Angeles, Los Angeles, CA, USA

<sup>5</sup> VA Informatics and Computing Infrastructure, US Department of Veterans Affairs, Salt Lake City, UT, USA

<sup>6</sup> Global Epidemiology Organization, Johnson & Johnson, Titusville, NJ, USA

**Funding statement:** This work is partially supported through US National Institutes of Health grants (T15 LM007079, R01 LM006910, and R01 HL169954).

**Conflict of interest statement:** PBR and MJS are employees of Johnson & Johnson. MAS receives a contract from Johnson & Johnson to support methods research not directly related to this study. Johnson & Johnson and Janssen did not have input in the design, execution, interpretation of results or decision to publish. All other authors declare no competing interests.

**Ethics of approval statement:** The research was approved by the Columbia University Institutional Review Board as an OHDSI network study.

**Corresponding author:** George Hripcsak, Department of Biomedical Informatics, Columbia University Irving Medical Center, 622 W. 168th Street, PH20, New York, 10032, NY, USA; email: [hripcsak@columbia.edu](mailto:hripcsak@columbia.edu)

**Keywords:** confounding, covariate balance, propensity score, meta-analysis

## ABSTRACT

Propensity score adjustment addresses confounding by balancing covariates in subject treatment groups through matching, stratification, inverse probability weighting, etc. Diagnostics ensure that the adjustment has been effective. A common technique is to check whether the standardized mean difference for each relevant covariate is less than a threshold like 0.1. For small sample sizes, the probability of falsely rejecting the validity of a study because of chance imbalance when no underlying balance exists approaches 1. We propose an alternative diagnostic that checks whether the standardized mean difference statistically significantly exceeds the threshold. Through simulation and real-world data, we find that this diagnostic achieves a better trade-off of type 1 error rate and power than standard nominal threshold tests and not testing for sample sizes from 250 to 4000 and for 20 to 100,000 covariates. In network studies, meta-analysis of effect estimates must be accompanied by meta-analysis of the diagnostics or else systematic confounding may overwhelm the estimated effect. Our procedure for statistically testing balance at both the database level and the meta-analysis level achieves the best balance of type-1 error rate and power. Our procedure supports the review of large numbers of covariates, enabling more rigorous diagnostics.

## INTRODUCTION

One of the major challenges facing observational research is the risk of producing a biased estimate due to confounding. Propensity score adjustment, invented 40 years ago [1], is a commonly used solution for measured confounders. It is a balancing score in the sense that matching subject treatment groups based on the score tends to balance all of the covariates used to estimate the score. Applying the score to a causal analysis can be done in several ways—matching, stratification, inverse probability treatment weighting, etc.—but the result is the same: removing the effect of the confounder on the causal estimate. Causal inference methods require diagnostics to ensure that any attempted adjustment has been successful. Propensity score adjustment is often assessed [2] using the standardized mean difference of suspected confounders among the treatment groups [3]. A high standardized mean difference reflects imbalance among the treatment groups and potentially ineffective adjustment for that confounder.

It is known that standardized mean difference can falsely reject studies when sample size is too small or too large. With small sample sizes, chance imbalance can cause large deviations of the standardized mean difference from zero. Austin [3] suggests measuring the empirical distribution of the standardized mean difference to account for chance imbalance, but this is rarely done in practice. Additional improvements such as comparing moments beyond the mean have been suggested [3] but are essentially never used. With large sample sizes, small degrees of systematic imbalance can be detected even though imbalance at that level is unlikely to cause an appreciable change in the effect estimate [3]. Researchers often pick a threshold to which a nominal estimate of the standardized mean difference can be compared; 0.1 is chosen most often [3-15] but 0.25 has also been used [11,12]. These thresholds are large enough to accommodate chance imbalance in moderate to large studies and to accommodate real but likely unimportant deviations in balance.

The problem of chance imbalance grows as sample size decreases, either due to a small observational data source or an uncommon treatment. The problem also grows with the number of covariates. Traditional manual selection of confounders leads to 5 to 20 or more covariates. High-dimensional propensity score adjustment [16] leads to hundreds of covariates. Large-scale propensity score adjustment [17,18] leads to tens of thousands of covariates. The probability of falsely rejecting a study that has no confounding or instruments can be calculated. Assume the diagnostic checks for standardized mean difference with a threshold of 0.1, and assume  $J$  independent binary covariates each of prevalence 0.5. The variance,  $\sigma_d^2$ , of the standardized mean difference,  $d$ , can be approximated as follows [19]:

$$\sigma_d^2 = \frac{n_1 + n_0}{n_1 n_0} + \frac{d^2}{2(n_1 + n_0 - 2)}$$

where  $n_1$  and  $n_0$  are the sizes of the two treatment groups and  $d$  is the standardized mean difference. Given  $d$  is small and assuming two equally sized treatment groups, the probability of false rejection is then given by:

$$P(\text{false rejection}) = 1 - \left( 2\Phi\left(\frac{\sqrt{N}}{20}\right) - 1 \right)^J$$

where  $\Phi$  is the cumulative distribution function of the standard normal distribution and  $N$  is the total sample size,  $n_1+n_0$ . With a sample size of 250, it takes only 5 covariates to reject 90% of studies by chance, and with a sample size of 1000, it takes 20 covariates. These numbers are well within the number of covariates often adjusted for in traditional observational research. With techniques like large-scale propensity score adjustment, it takes 8000 subjects to reach reasonable acceptance rates.

A further challenge arises with the growth of networks of observational databases and federated analyses across those databases. Combining the effect estimates from the databases using meta-analysis can result in more precise estimates. The increasing precision of the effect estimates should be matched with increased precision of diagnostics. That is, if there is systematic imbalance across the databases, then the level of imbalance will have a proportionately larger effect compared to the increased precision of the effect estimate.

The covariates that are tested for imbalance are usually the covariates adjusted for. Most commonly, researchers select a small number of covariates, say 4 to 20, that are suspected to be potential confounders and adjust for them. The process is unreliable, as illustrated by a sample of hypertension studies [20-24], each of which claims to have adequately addressed confounding but for which there is only moderate overlap of the confounders across studies. An alternative, which has empirical backing, is to adjust for all observed covariates—potentially tens of thousands—that are not eliminated as potential mediators, colliders, or instruments [17,18]. In comparisons, adjusting for all covariates has outperformed manual confounder selection [18,25,26] and empirical confounder selection [17].

Even if one chooses to adjust for a small number of covariates, the question remains how many covariates to check for balance. If the problem of false rejections of the study can be addressed, then it would seem that the covariates include information about the study even if they are not adjusted for. Rather than taking a head-in-the-sand approach of ignoring such imbalance, it should be measured and explained. If the covariates are not plausible instruments, then they should be adjusted for or the study should be discarded.

In this paper, we propose a simple approach to address false rejection using a statistical test for exceeding the 0.1 threshold and with Bonferroni correction where appropriate. We do this at both the single database level and the network level. The approach is easy to understand and easy to implement. We recognize that testing for the presence of any statistically significant imbalance among the covariates is appropriately frowned upon [11,27], but we believe that testing for exceeding the 0.1 threshold is justified and argue the case further in the discussion. We study the operating characteristics of this diagnostic procedure in simulation and on real-world data, and we report on the implications for single database studies and for meta-analyses.

## METHODS

### *Simulation*

Our goal for the simulation was to create a data set that would mimic the balance characteristics of a data set whose treatments groups had been completely or imperfectly adjusted for confounding, for example using 1-to-1 matching based on a propensity score. We used a binary treatment and a binary outcome, with a set of binary covariates that could include a confounder or other correlations that were not sufficiently balanced by the adjustment procedure.

For the base case, we varied database sample size, effect size, and degree of confounding, holding outcome prevalence, covariate prevalence, and aggregate sample size constant and using homogeneous confounding across databases. We used an aggregate sample size of 20,000 subjects, which were apportioned among databases with 4000, 2000, 1000, 500, and 250 subjects, producing 5, 10, 20, 40, and 80 databases, respectively. Each experiment was carried out 200 times to estimate error rates. Each subject,  $i$ , was defined as follows:

$$x_{i,j} \sim \text{Bernoulli}(0.5)$$

$$t_i \sim \text{Bernoulli} \left( 0.5 + c_t(1 - 2x_{i,1}) \right)$$

$$y_i \sim \text{Bernoulli} \left( 0.25 + c_e(1 - 2t_i) + c_y(1 - 2x_{i,1}) + c_x(1 - 2x_{i,2}) + \dots + c_y(1 - 2x_{i,10}) \right)$$

where the  $x_{ij}$  were the non-treatment covariates for subject  $i$  and covariate  $j$ ,  $t_i$  was the treatment, and  $y_i$  was the outcome. Index  $j$  varied from 1 to 1000 with  $x_{i,1}$  being a potential confounder, with  $x_{i,2} \dots x_{i,10}$  being causally linked to the outcome but not the treatment, and with  $x_{i,11} \dots x_{i,1000}$  not being linked to treatment or outcome. Constant  $c_e$  determined the treatment effect,  $c_t$  was the link from the confounder  $x_{i,1}$  to the treatment  $t_i$ ,  $c_y$  was the link from the confounder  $x_{i,1}$  to the outcome  $y_i$  and  $c_x$  was the link from covariates  $x_{i,2}$  to  $x_{i,10}$  to the outcome  $y_i$ . We used 1000 covariates to simulate the probability of detecting imbalance among covariates by chance in an analysis with large-scale covariate adjustment such as is used in large-scale propensity score adjustment [17,18] or high-dimensional propensity score adjustment [16]. For the base case,  $c_e$  varied from 0 to 0.1,  $c_t$  varied from 0 to 0.3,  $c_y$  was held at 0.1, and  $c_x$  was held at 0.1. The parameters were chosen such that they best illustrated the range of performance, for example from an effect that was undetectable to an effect that was detectable by all studied approaches, and they were iterated upon to find weaknesses in the approaches. When  $c_e$  and  $c_t$  both equaled zero, there was no systematic source of imbalance between the two groups and any detected imbalance was due to chance.

To estimate an effect size, we used function `glm` (family = binomial) in R to carry out a simple logistic regression using  $t_i$  to predict  $y_i$ , ignoring  $x_{i,j}$ , as one might do if one assumes a previous adjustment procedure succeeded in achieving balance. We studied our ability to detect imbalance among covariates in the treated versus untreated group, correlating those results with measured type 1 error and power based on the true effect  $c_e$  and the effect estimates and their variances from the model.

To quantify imbalance, we used the standardized mean difference (SMD),  $smd_j$ , for covariates  $x_{i,j}$ , defining it as follows:

$$n_1 = \sum_i t_i$$

$$n_0 = \sum_i 1 - t_i$$

$$s_{1,j} = \sum_i t_i x_{i,j}$$

$$s_{0,j} = \sum_i (1 - t_i) x_{i,j}$$

$$sd_j = \sqrt{\frac{\left(\frac{s_{1,j}}{n_1}\right)\left(\frac{1 - s_{1,j}}{n_1}\right) + \left(\frac{s_{0,j}}{n_0}\right)\left(\frac{1 - s_{0,j}}{n_0}\right)}{2}}$$

$$smd_j = \frac{\frac{s_{1,j}}{n_1} - \frac{s_{0,j}}{n_0}}{sd_j}$$

$$\text{var} \text{smd}_j = \frac{n_1 + n_0}{n_1 n_0} + \frac{\text{smd}_j^2}{2(n_1 + n_0 - 2)}$$

where  $\text{smd}_j$  is the pooled sample standard deviation for covariates  $x_{i,j}$  and  $\text{var} \text{smd}_j$  is the variance of  $\text{smd}_j$ . We use a large-sample estimate of the variance instead of its empirical distribution to ensure that the resulting procedure is simple enough for actual adoption.

We estimated the type 1 error rate as the proportion of 200 study iterations where the effect coefficient estimate differed statistically significantly from 0 when  $c_e=0$ . We estimated the power as the proportion of the 200 iterations where the effect coefficient estimate differed statistically significantly from 0 when  $c_e \neq 0$ .

Given these definitions, we studied alternative decision rules to determine whether to reject a network study based on covariate imbalance, and we compared their type 1 error rate and power after rejecting the imbalanced studies and assigning them a status of not statistically significant regardless of the effect estimate. Our decision rules operated at two levels, at the single database level and at the network level, and there were three types of rules: accept all studies (i.e., ignore imbalance), reject studies with any  $\text{smd}_j \geq 0.1$  (following Austin [3]), and reject studies where the  $\text{smd}_j$  is statistically significantly greater than 0.1 using  $\text{var} \text{smd}_j$  and a Bonferroni correction for the number of covariates. For the network study, we employed a random effects meta-analysis using the R function `rma` both to determine the overall effect size and to determine the overall standardized mean difference for each covariate. The result was nine total rules; we use the shorthand of “all” for accepting all studies, “nominal” for checking for any  $\text{smd}_j \geq 0.1$ , and “signif” for checking for any  $\text{smd}_j$  being statistically significantly greater than 0.1 after Bonferroni correction:

- AllOnAll** (all network, all database): Accept each database regardless of imbalance, and then accept the network meta-analysis regardless of overall imbalance
- AllOnNominal** (all network, nominal database): Reject databases with any  $\text{smd}_j \geq 0.1$ , and then accept the network meta-analysis of the remaining databases regardless of overall imbalance
- AllOnSignif** (all network, statistical database): Reject databases with any  $\text{smd}_j$  being statistically significantly greater than 0.1 after Bonferroni correction for the number of hypotheses, and then accept the network meta-analysis of the remaining databases regardless of overall imbalance
- NominalOnAll** (nominal network, all database): Accept each database regardless of imbalance, and then reject the network meta-analysis if the meta-analytic overall imbalance was greater than or equal to 0.1 for any covariate
- NominalOnNominal** (nominal network, nominal database): Reject databases with any  $\text{smd}_j \geq 0.1$ , and then reject the network meta-analysis of the remaining databases if the meta-analytic overall imbalance was greater than or equal to 0.1 for any covariate
- NominalOnSignif** (nominal network, statistical database): Reject databases with any  $\text{smd}_j$  being statistically significantly greater than 0.1 after Bonferroni correction for the number of hypotheses, and then reject the network meta-analysis of the remaining databases if the meta-analytic overall imbalance was greater than or equal to 0.1 for any covariate
- SignifOnAll** (statistical network, all database): Accept each database regardless of imbalance, and then reject the network meta-analysis if the meta-analytic overall imbalance was statistically significantly greater than 0.1 for any covariate after Bonferroni correction for the number of covariates
- SignifOnNominal** (statistical network, nominal database): Reject databases with any  $\text{smd}_j \geq 0.1$ , and then reject the network meta-analysis of the remaining databases if the meta-analytic overall imbalance was statistically significantly greater than 0.1 for any covariate after Bonferroni correction for the number of covariates

**SignifOnSignif** (statistical network, statistical database): Reject databases with any  $smd_j$  being statistically significantly greater than 0.1 after Bonferroni correction for the number of hypotheses, and then reject the network meta-analysis of the remaining databases if the meta-analytic overall imbalance was statistically significantly greater than 0.1 for any covariate after Bonferroni correction for the number of covariates

Of note, rule AllOnAll ignores imbalance and accepts all studies regardless of the detection of a potential for confounding. Rules NominalOnAll and SigniOnAll ignore database-level imbalance and rely solely on a network-level detection. Rules AllOnNominal and AllOnSignif ignore any network-level imbalance but exploit database-level detection. In general, database-level imbalance detection is expected to be severely affected by sample size, either missing imbalance or declaring false-positive imbalance on small databases. Network-level imbalance is expected to be more likely to detect imbalance that is shared among databases but may miss an aberrant database.

We carried out several simulation experiments in addition to the base case. We reran the study under the conditions of low covariate prevalence (10% instead of 50%), low outcome prevalence (1% instead of 25%), heterogeneous confounding (varying  $c_t$  from  $-0.3$  to  $0.3$  for each database instead of holding it constant), and fewer databases (capped at 5 databases instead of holding the aggregate total constant).

### *Real-World Data*

We exploited a data set and protocol used in two previous studies to illustrate the effect of choice of covariate balance detection on real-world covariates and confounding. The Observational Health Data Sciences and Informatics (OHDSI) [28,29] LEGEND hypertension [30] and type-2 diabetes studies [31] comprehensively evaluated the comparative effects of all pharmaceutical treatments for their respective disease areas. Here we select two comparisons from each, including the negative control outcomes used in these studies.

We used three databases shown in Table 1. Each database uses OHDSI's Observational Medical Outcome Partnership (OMOP) common data model [32] populated with patient characteristics, health care visits, diseases, medications, procedures, and, optionally, other data types such as laboratory tests. Data elements were translated to standard terminologies [33,34] such as Systematized Nomenclature of Medicine (diseases, procedures), RxNorm (medications), and Logical Observation Identifiers Names and Codes (laboratory tests).

The gold standard was the collection of 110 real negative controls used in the original studies as well as synthetic positive controls generated from them. Each negative control contained a target drug, a comparator drug, and an outcome that was determined through review of the literature, product labels, spontaneous reports, and clinical experts not to be causally associated with either drug [35]. Therefore, the hazard ratio of the appearance of the outcome between the two drugs should be 1, indicating no difference. Synthetic positive controls were generated from the negative controls by fitting an L1-regularized Poisson regression model [36] on presence of the outcome based on all pre-treatment covariates and then using that model to insert additional simulated events into the data set to achieve hazard ratios of 1.5, 2, and 4. More details are available [30,31].

For each of the three large databases, we created several sets of smaller databases to illustrate the effect of a network study across databases with small sample sizes. We kept the total number of cases constant at 20,000: 5 databases with 4000 cases each, 10 databases with 2000 cases each, 20 databases with 1000 cases each, 40 databases with 500 cases each, and 80 databases with 250 cases each. For each of these sets, we calculated large-scale propensity scores [18] for four treatment comparisons—lisinopril versus hydrochlorothiazide, lisinopril versus metoprolol, sitagliptin versus liraglutide, and sitagliptin versus

glimepiride—on 98,681 covariates found in the three databases. Only pretreatment variables were used to minimize the probability of mediators or colliders, suspected instruments were removed, and lack of strong instruments was confirmed by measuring equipoise. We calculated propensity scores two ways: on each small database and on a pooled sample of all 20,000 cases. We used three methods to apply the propensity score to a causal analysis: no adjustment (crude), 1-to-1 matching, and stratification. We then used a Cox proportionate hazards model to estimate the hazard ratio for each hypothesis, calculated a confidence interval, and tested for significance based the confidence interval excluding no effect (hazard ratio of 1).

The result was a set of analyses characterized by which of three large databases it came from, by the four hypotheses, by the sample size of the generated database (250, 500, 1000, 2000, 4000), by the underlying hazard ratio (1, 1.5, 2, 24), and by the analysis method (unadjusted=crude, matched on the database-level propensity score, stratified on the database-level propensity score, matched on the pooled propensity score, stratified on the pooled propensity score). Each analysis had an effect estimate and a standardized mean difference for every covariate.

We then further carried out a meta-analysis as described in the Simulation methods section across the 80, 40, 20, 10, or 5 generated data sets (for 250, 500, 1000, 2000, and 4000 sample sizes, respectively) getting meta-analytic estimates for the effect and all the covariate SMDs. We applied the nine covariate balance rules described above in the Simulation methods section. We calculated the proportion of studies that passed the covariate balance rule and report the type 1 error rate and power for each rule for each sample size and for two groups of analyses: unadjusted=crude versus all propensity-adjusted methods. The unadjusted analyses should display greater confounding than the adjusted ones.

## RESULTS

### *Simulation*

Figure 1 shows the performance of the three types of rules at the database level, with type 1 error rate in the first column of graphs (no true effect with  $c_e=0$ ) and power in the rest of the columns (increasing effect with  $c_e>0$ ). The first row, labeled All, ignores imbalance. It has good power but increasing type 1 error with increasing confounding, reaching 1 for high confounding. The second row, labeled Nominal, uses the standard practice of rejecting studies where at least one standardized mean difference equals or exceeds 0.1. This rule rejects all studies with fewer than 4000 subjects even when confounding is zero due to chance imbalance, producing little power. The third row, labeled Signif, rejects studies where at least one standardized mean difference is statistically significantly great than 0.1 after Bonferroni correction. Type 1 error rate is relatively controlled with an average of 0.054, with the highest rate of 0.14 at intermediate levels of confounding ( $c_i=-0.06$ ), where the confounding did not trigger a rejection but caused a false positive result. Power rises with effect size, reaching near one except where confounding is most strong, leading to rejection of those highly confounded studies.

Figure 2 shows the performance of the nine rules defined in Methods for the base case of the network study. The first row, AllOnAll, illustrates that completely ignoring imbalance produces increasingly poor type 1 error rate with increasing confounding, getting to 1 for the highest confounding. The second row, AllOnNominal, illustrates that using the nominal rule of checking for standardized mean difference greater than or equal to 0.1 at the database level leads to rejecting all databases smaller than 4000 even when there is no true imbalance, resulting in a power of zero. For the same reason, rules NominalOnNominal and SignifOnNominal, which use the nominal test at the database level, fail with low power. The third row, rule AllOnSignif, illustrates that checking for statistically significant imbalance only at the database level leads to low type 1 error rate when there is no confounding and when there is high confounding (when it can be detected) but high type 1 error rate near 1 with intermediate



confounding ( $c_i=0.1$  or  $-0.1$ ) that is strong enough to cause a false positive result but not strong enough to trigger the rule to drop the database. The rows for rules NominalOnAll and NominalOnSignif illustrate that checking for nominally reaching the 0.1 threshold on the meta-analytic standardized mean difference estimates produces low type 1 error and high power if the confounding is low. The rows for rules SignifOnAll and SignifOnSignif illustrate that checking whether the meta-analytic standardized mean difference estimate is statistically significantly greater than or equal to 0.1 produces a generally reasonable type 1 error rate and high power for low confounding. We were able to find a combination that produced a high type 1 error rate of 0.295 when  $c_i=0.03$ , representing an intermediate level of confounding.

Thus, four rules—NominalOnAll, NominalOnSignif, SignifOnAll, and SignifOnSignif—were viable according to the simulation, with the latter two showing some higher type 1 error. The four rules differentiated when the number of databases was limited to five (Figure S1). Using a nominal test for the meta-analytic standardized mean difference reaching the 0.1 threshold (rules NominalOnAll and NominalOnSignif) resulted in zero power for databases smaller than 1000 because chance imbalance always disqualified the study even when there was no confounding. Testing for statistical significance (rules SignifOnAll and SignifOnSignif) avoided discarding non-confounded small studies. In other words, the nominal threshold test of the meta-analytic threshold only worked when the network was large enough, but the statistical test behaved gracefully with smaller network sizes.

We also performed a number of further sensitivity analyses, shown in the Supplement. When outcome prevalence is low (Figure S2), the rules perform similarly to each other, and they all lose some power compared to the base case, where outcome prevalence is higher. When covariate prevalence is low (Figure S3), the four rules perform similarly to the base case. When confounding is heterogeneous (Figure S4), the four rules produce similar results; rules NominalOnAll and NominalOnSignif have less power, but because the study is confounded, greater power is not the goal. We tried the Signif rules without the Bonferroni correction (Figure S5); the type 1 error rate remained high at 0.260 for the  $c_i=0.03$  case, and it reduced power near 0 in the smallest studies.

We studied the effect of using fewer covariates. When the covariate count is 20 (Figure S6), we see similar results to the base case except that the curves are shifted to the left, with the same issues occurring across the nine rules, but with one quarter the sample size.

To provide context for the type 1 error rates at small sample sizes, we reproduced a more typical study with sample size increased to 20,000 and keeping the covariates at 20 (Figure S7). We show that type 1 error rates can get over 0.57 for all nine rules if we select confounding carefully. We also note the similarity of the results for the Nominal rules versus the Signif rules with this large sample size (Figures S7 and S16).

### *Real-World Data*

Figure 3 shows the type 1 error ( $RR=1$ ) and power ( $RR>1$ ) for all nine rules on the real-world data for different sample sizes and levels of adjustment (adjusted versus unadjusted, likely reflecting lower and higher residual confounding). When covariate imbalance was ignored (rule AllOnAll), the type 1 error rate was high at 0.09 for the unadjusted study (no covariate adjustment and therefore presumably higher confounding) when sample size was 250. As was observed in the simulation, when covariate balance was checked at the database level using a threshold of the standardized mean difference nominally exceeding 0.1 (rules AllOnNominal, NominalOnNominal, and SignifOnNominal), all studies were rejected resulting in no power. Checking at the network level for the meta-analytic standardized mean difference nominally exceeding 0.1 (rules NominalOnAll, NominalOnNominal, and NominalOnSignif) also resulted in rejecting most studies with reduced power. This left only rules that checked for statistically significant

covariate imbalance. Of note, for rule AllOnSignif, which only checks for covariate imbalance at the database level, the type 1 error rate was high at 0.07 when sample size was 250 with the crude analysis, reflecting the consequence of ignoring network-level estimates of standardized mean difference. Rules SignifOnAll and SignifOnSignif, which test for statistically significant imbalance at the network level, both produced type 1 error less than 0.05 and showed power similar to no checking (rule AllOnAll) in the adjusted studies (when confounding was likely to be low) and showed low power due to appropriately rejecting unadjusted studies (when confounding is likely to be high).

## DISCUSSION

Our study illustrates two principles and identifies an algorithm that appears to best address the tradeoff between type 1 error and power in the simulation and real-world data. First, as sample size decreases, using a nominal test of imbalance such as described by Austin [3] will result in near-certain rejection of the study even with no confounding due to chance imbalance, and this effect occurs with few covariates as well as many covariates. Using a statistical test of imbalance exceeding a threshold like 0.1 will maintain power without substantially raising the type 1 error rate. When doing a study with small sample sizes, it may be hard to detect small-to-moderate confounding, but our results illustrate that the small study's wide confidence intervals will avert a high type 1 error rate. That is, it takes significant confounding to cause high type 1 error rate. Second, when doing a network study, it is important to carry out a meta-analysis not just of the effect estimate but also of the diagnostics such as standardized mean difference. The meta-analysis of the effect estimate may potentially produce a more precise effect estimate and a narrow confidence interval so that small-to-moderate confounding can yield too many false positive results. The meta-analysis of the standardized mean differences, however, permits the detection of lower levels of confounding despite the sample sizes being individually small.

Best overall performance appeared to be achieved by testing for statistically significant imbalance defined as standardized mean difference over 0.1 and using a Bonferroni correction across covariates. For network analyses, rule SignifOnSignif, which tested for statistically significant imbalance at both the single database level and across the database network using a meta-analysis of standardized mean differences, worked best. It produced power at all studied sample sizes, and its power was generally not too much less than ignoring imbalance, yet it achieved substantially lower type 1 error rates. In the simulation, we used iterative parameter testing to find a combination that produced a type 1 error rate approaching 0.3, but we emphasize that we could achieve such high error rates even with standard practice. For example, even in a much larger study with a sample size of 20,000, using only 20 covariates, setting  $c_t$  to 0.02,  $c_y$  to 0.2, and  $c_e$  to 0, and using a nominal threshold of 0.1, the type 1 error rate was 0.57. Given any threshold, it is possible to design a simulation that thwarts it. The most important question is what happens with real-world data, and in fact, rule SignifOnSignif produced low type 1 error rate. Rule SignifOnSignif had other good properties. If the confounding was heterogeneous, doing the standardized mean difference test at both levels worked best: drop the databases that fail at the individual level and do a meta-analysis on the rest. As the number of databases varied from 1 to 5 up to 80, rule SignifOnSignif produced a stable output, and it worked whether the number of covariates was 20, 1000, or 98,681. When covariate prevalence or outcome prevalence was low, the results were similar.

Rules that appeared to work on simulation failed on the real-world data. For example, rule NominalOnAll, which tested for any meta-analytic standardized mean difference to be nominally over 0.1, had too-low power in real-world data compared to testing for statistical significance: even with 4000 cases in a study with low confounding, most studies were rejected, rendering the rule ineffective for small studies. It also failed in the simulation when the number of databases in the network was only 5. The two techniques that most clearly failed represent standard practice: doing a nominal test for imbalance at the database level (AllOnNominal) rejected all databases regardless of confounding as sample size fell; and doing no test for imbalance (AllOnAll) led to very poor type 1 error as confounding rose.

Applying meta-analysis to study diagnostics will clearly be difficult given that very few studies share the details of their balance results (although it has been suggested [37]). The approach is still relevant to distributed networks that use meta-analysis to combine their results; sharing of study diagnostics can be incorporated into the study protocol [31]. It is important to at least recognize this limitation of meta-analysis on observational research without study diagnostics.

Showing power for studies in which confounding is non-zero can be seen as misguided. Should not the metric instead be the ability to detect confounded studies? We believe that power is the right metric. As Austin [3] showed, the goal is not to detect increasingly minute imbalance, but to detect imbalance that might reflect confounding that matters. As sample size falls, confidence intervals widen, and it takes more confounding to alter the result appreciably. If important confounding is slipping through the imbalance test, then that will be reflected in the type 1 error rate. If the type 1 error rate stays near its nominal value (say 0.05), then we argue that power reflects acceptance of studies that while they might be confounded to some degree, that confounding is insufficient to cause a frequent false result. All observational studies in fields like medicine have some small degree of confounding, so it is important that the goal be recognized as achieving best power given a nominal type 1 error rate, not eliminating all studies.

This study potentially informs other observational research as well. The current custom is to check balance only on the covariates suspected to be confounders and adjusted for using propensity scores, ignoring potentially useful information about other covariates. Given our results, it may be more informative to incorporate all available covariate information. If a covariate is unbalanced, that should be explained as an instrument if domain knowledge is available or the study should be suspected to be biased. Our work demonstrates that it is possible to test for imbalance without triggering too many false positives and also without missing confounding that could substantially affect the results.

As authors have pointed out [3,11,27,38], using a statistical test to detect the presence of imbalance (difference from 0) performs poorly and does not achieve the proper goal, yet we believe that this test for exceeding a threshold (exceed 0.1) is in fact useful. Imai et al. [27] argue three points: that a statistical test is dependent on sample size yet the actual imbalance is not, that any threshold like a p-value is arbitrary, and that the target of analysis is the sample itself and not some underlying population. On the first, we too are dependent on sample size but we believe that that is appropriate. As the size of the sample shrinks, the effect estimate will become less precise with a larger confidence interval and any given level of confounding will have a proportionately smaller influence on the conclusion of the study. Therefore, the threshold for imbalance ought to become less stringent as sample size falls. On the second, we agree that a threshold is arbitrary at first, but the observational research field including most of these authors have come to relative agreement on a 0.1 constant threshold [3-15] despite the fact that no threshold can guarantee immunity from important bias. The appropriateness of a threshold is decided slowly over time as real study results are compared to baseline knowledge and validated in later randomized experiments. Our use of statistical comparison to 0.1 is no more arbitrary than current practice. On the third, we acknowledge that we are concerned with the current sample, but we use the statistical test to effectively adjust the level of the threshold to the sample size, in effect accounting for the size of the effect estimate sample size. Our test may be better seen as a heuristic that uses a threshold whose level varies with sample size.

Our use of a Bonferroni correction may be questioned, but we argue it is necessary and appropriate. With increasing numbers of independent covariates, the likelihood of chance imbalance will rise, so as we have demonstrated (Figure S5), correction for multiple hypotheses is needed to avoid rejecting too many studies. One may question why existing confounders should be allowed more imbalance just by adding new independent covariates. First, and most important, we include many covariates because small sets of manually chosen or empirically selected confounders are likely missing confounding (both directly and

indirectly measured [18]), which is likely why large-scale propensity adjustment appears to perform better than other methods [17,18,25,26]. That is, we agree that domain knowledge is effective for identifying important confounders but we believe it has little ability to rule out other confounders. Our further simulations (see Supplement section III) demonstrate that if confounders are distributed among the covariates, then adding covariates actually increases our ability to detect confounding despite the Bonferroni correction. We also point out that for the range of standard errors of standardized mean difference that we saw in the real-world data (0.044) adding a Bonferroni correction for 100,000 covariates when testing for exceeding at threshold of 0.1 is still stricter than not using Bonferroni correction but using a threshold of 0.25, which is a previously accepted alternative threshold [11,12]. In the end, we believe that the question of how many covariates to include is an empirical one, balancing the benefit of covering more confounders with the risk of including inappropriate variables, and our experience so far is that more has been better [17,18,25,26], and this experiment shows that on real-world data, chance imbalance can be addressed by rule SignifOnSignif even with 98,861 covariates. Nevertheless, whether a researcher selects 20 covariates or 100,000 covariates, the results of this experiment remain relevant.

One can also argue whether observational studies with small sample sizes can be trusted. They are clearly important because even in large databases, uncommon treatments can result in small cohorts. As we demonstrate in our experiment, such studies can achieve type 1 error rates and power comparable to larger studies. With modern regularized regression [36], even propensity models with 100,000 covariates can be stably estimated with sample sizes down to 250 [39]. Therefore, we see no need to exclude them.

One limitation of our design is that type 1 error rate and power were estimated by counting rejected studies as if their effect estimates were not distinguishable from the null hypothesis. One can recover the type 1 error rate and power within only the non-rejected databases by dividing the type 1 error or power by the proportion valid (Figures S8-S17). The relative effectiveness of the rules remains the same, however.

## REFERENCES

1. Rosenbaum PR, Rubin DB. "The central role of the propensity score in observational studies for causal effects." *Biometrika* 1983; 70:41–55.
2. Granger E, Watkins T, Sergeant JC, Lunt M. A review of the use of propensity score diagnostics in papers published in high-ranking medical journals. *BMC Med Res Methodol* 2020 May 27;20(1):132.
3. Austin PC. Balance diagnostics for comparing the distribution of baseline covariates between treatment groups in propensity-score matched samples. *Stat Med* 2009;28:3083-107.
4. Normand SL T, Landrum MB, Guadagnoli E, Ayanian JZ, Ryan TJ, Cleary PD, McNeil BJ. 5 Validating recommendations for coronary angiography following an acute myocardial infarction in the elderly: A matched analysis using propensity scores. *Journal of Clinical Epidemiology*. 2001;54:387–398.
5. Mamdani M, Sykora K, Li P, et al. Reader's guide to critical appraisal of cohort studies: 2. Assessing potential for confounding. *Br Med J*. 2005;330(7497):960–962.
6. Markoulidakis A, Taiyari K, Holmans P, Pallmann P, Busse M, Godley MD, Griffin BA. A tutorial comparing different covariate balancing methods with an application evaluating the causal effects of substance use treatment programs for adolescents. *Health Serv Outcomes Res Methodol* 2023;23(2):115-148.
7. Austin PC. A critical appraisal of propensity-score matching in the medical literature between 1996 and 2003. *Stat. Med.* 2008;27(12):2037–2049.
8. Austin PC, Grootendorst P, Anderson GM. A comparison of the ability of different propensity score models to balance measured variables between treated and untreated subjects: a Monte Carlo study. *Stat. Med.* 2007;26(4):734–753.

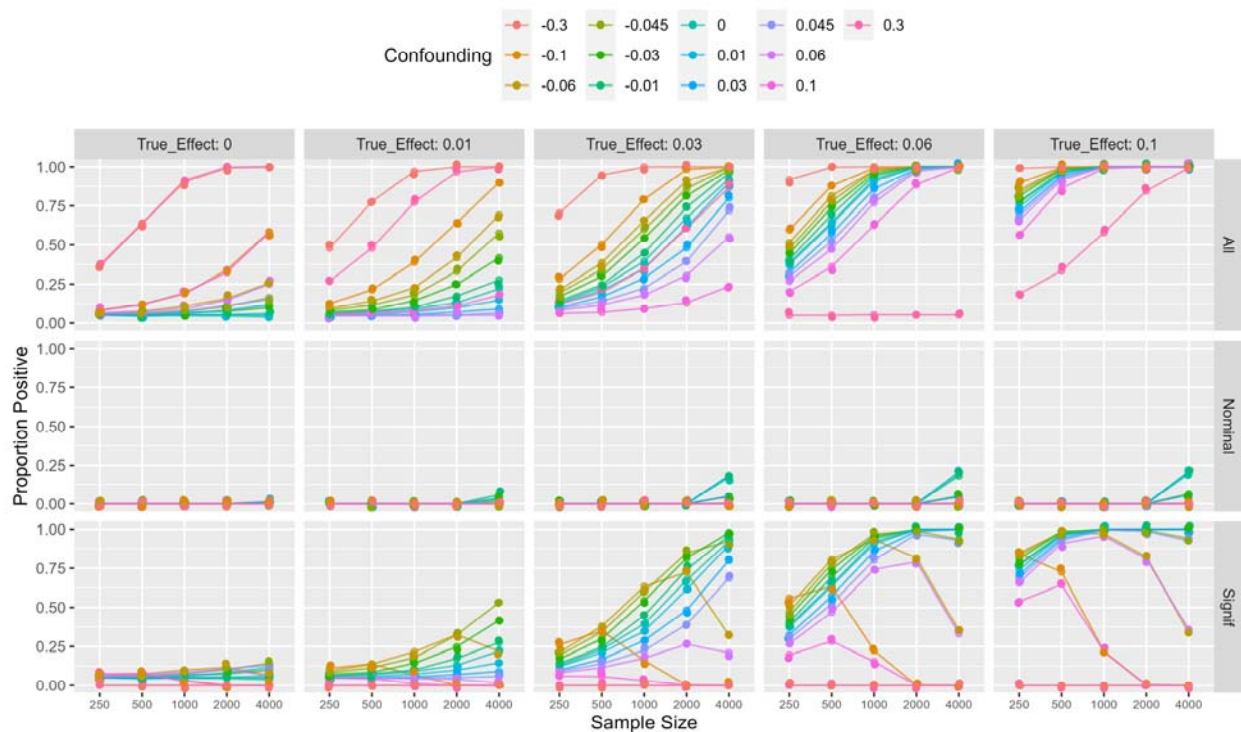
9. Griffin BA, McCaffrey DF, Almirall D, Burgette LF, Setodji CM. Chasing balance and other recommendations for improving nonparametric propensity score models. *J Causal Inference* 2017;5(2).
10. Griffin BA, Ramchand R, Almirall D, Slaughter ME, Burgette LF, McCaffery DF. Estimating the causal effects of cumulative treatment episodes for adolescents using marginal structural models and inverse probability of treatment weighting. *Drug Alcohol Depend.* 2014;136:69–78.
11. Ho DE, Imai K, King G, Stuart EA. Matching as nonparametric preprocessing for reducing model dependence in parametric causal inference. *Polit. Anal.* 2007;15(3):199–236.
12. Stuart EA, Lee BK, Leacy FP. Prognostic score-based balance measures can be a useful diagnostic for propensity score methods in comparative effectiveness research. *J. Clin. Epidemiol.* 2013;66(8):S84–S90.
13. Zhang Z, Kim HJ, Lonjon G, Zhu Y, et al. Balance diagnostics after propensity score matching. *Ann Trans Med* 2019; 7(1).
14. Medaglio D, Stephens-Shields AJ, Leonard CE. Research and scholarly methods: Propensity scores. *J Am Coll Clin Pharm* 2022 Apr;5(4):467-475.
15. Chang TH, Nguyen TQ, Lee Y, Jackson JW, Stuart EA. Flexible propensity score estimation strategies for clustered data in observational studies. *Stat Med* 2022 Nov 10;41(25):5016-5032.
16. Schneeweiss S, Rassen JA, Glynn RJ, Avorn J, Mogun H, Brookhart MA. High-dimensional propensity score adjustment in studies of treatment effects using health care claims data. *Epidemiology.* 2009; 20(4): 512-522.
17. Tian Y, Schuemie MJ, Suchard MA. Evaluating large-scale propensity score performance through real-world and synthetic data experiments. *Int J Epidemiol* 2018; 47:2005-14.
18. Zhang L, Wang Y, Schuemie M, Blei D, Hripcsak G. Adjusting for indirectly measured confounding using large-scale propensity score. *Journal of Biomedical Informatics.* 2022 Oct;134:104204.
19. Zejnnullahi R, Hedges LV. Robust variance estimation in small meta-analysis with the standardized mean difference. *Res Syn Meth* 2023;1-17.
20. Chien SC, Ou SM, Shih CJ et al. Comparative Effectiveness of Angiotensin-Converting Enzyme Inhibitors and Angiotensin II Receptor Blockers in Terms of Major Cardiovascular Disease Outcomes in Elderly Patients. *Medicine* 2015; 94(43): e1751.
21. Hicks BM, Fillion KB, Yin H et al. Angiotensin converting enzyme inhibitors and risk of lung cancer: population based cohort study. *BMJ (Clinical Research Ed)* 2018; 363: k4209.
22. Ku E, McCulloch CE, Vittinghoff E et al. Use of Antihypertensive Agents and Association With Risk of Adverse Outcomes in Chronic Kidney Disease: Focus on Angiotensin-Converting Enzyme Inhibitors and Angiotensin Receptor Blockers. *Journal of the American Heart Association* 2018; 7(19): e009992.
23. Magid DJ, Shetterly SM, Margolis KL et al. Comparative Effectiveness of Angiotensin-Converting Enzyme Inhibitors Versus Beta-Blockers as Second-Line Therapy for Hypertension. *Circulation: Cardiovascular Quality and Outcomes* 2010; 3(5): 453–458.
24. Hasvold LP, Bodeg<sup>o</sup>ard J, Thuresson M et al. Diabetes and CVD risk during angiotensin-converting enzyme inhibitor or angiotensin II receptor blocker treatment in hypertension: a study of 15 990 patients. *Journal of Human Hypertension* 2014; 28(11): 663–669.
25. Weinstein RB, Ryan P, Berlin J A, Matcho A, Schuemie M, Swerdel J, Patel K, Fife D. Channeling in the use of nonprescription paracetamol and ibuprofen in an electronic medical records database: evidence and implications. *Drug Safety* 2017;40:1279–92.
26. Weinstein RB, Ryan PB, Berlin JA, Schuemie MJ, Swerdel J, Fife D. Channeling bias in the analysis of risk of myocardial infarction, stroke, gastrointestinal bleeding, and acute renal failure with the use of paracetamol compared with ibuprofen. *Drug Safety* 2020;43:927-942.
27. Imai K, King G, Stuart E. Misunderstandings between experimentalists and observationalists about causal inference. *Journal of the Royal Statistical Society, Series A* 2008;171:481-502.
28. Hripcsak G, Duke JD, Shah NH, Reich CG, Huser V, Schuemie MJ, Suchard MA, Park RW, Wong ICK, Rijnbeek PR, van der Lei J, Pratt N, Norén GN, Lim YC, Stang PE, Madigan D, Ryan PB. Observational Health Data Sciences and Informatics (OHDSI): opportunities for observational researchers. *MEDINFO'15; 2015 August 19 - 23, São Paulo, Brazil; 2015.*

29. Hripcsak G, Schuemie MJ, Madigan D, Ryan PB, Suchard MA. Drawing reproducible conclusions from observational clinical data with OHDSI. *Yearb Med Inform.* 2021 Apr 21.
30. Suchard MA, Schuemie MJ, Krumholz HM, You SC, Chen R, Pratt N, Reich CG, Duke J, Madigan D, Hripcsak G, Ryan PB. Comprehensive comparative effectiveness and safety of first-line antihypertensive drug classes. *Lancet* 2019;394(10211):1816-26.
31. Khera R, Schuemie MJ, Lu Y, Ostropolets A, Chen R, Hripcsak G, Ryan PB, Krumholz HM, Suchard MA. Large-scale evidence generation and evaluation across a network of databases for type 2 diabetes mellitus (LEGEND-T2DM): Protocol for a series of multinational, real-world comparative cardiovascular effectiveness and safety studies. *BMJ Open* 2022 Jun 9;12(6):e057977.
32. Overhage JM, Ryan PB, Reich CG, Hartzema AG, Stang PE. Validation of a common data model for active safety surveillance research. *J Am Med Inform Assoc.* 2012 Jan-Feb;19(1):54-60.
33. Reich C, Ostropolets A, Ryan P, Rijnbeek P, Schuemie M, Davydov A, Dymshyts D, Hripcsak G. OHDSI Standardized Vocabularies—a large-scale centralized reference ontology for international data harmonization, *Journal of the American Medical Informatics Association*, 2024; ocad247.
34. Bodenreider O, Cornet R, Vreeman DJ. Recent Developments in Clinical Terminologies – SNOMED CT, LOINC, and RxNorm. *Yearb Med Inform.* 2018 Aug;27(1):129-139.
35. Voss EA, Boyce RD, Ryan PB, van der Lei J, Rijnbeek PR, Schuemie MJ. Accuracy of an automated knowledge base for identifying drug adverse reactions. *J Biomed Inform* 2017 Feb;66:72-81.
36. Suchard MA, Simpson SE, Zorych I, Ryan PB, Madigan D. Massive parallelization of serial inference algorithms for a complex generalized linear model. *ACM Transactions on Modeling and Computer Simulation* 2013;23(1):1-17.
37. Schuemie MJ, Ryan PB, Pratt N et al. Principles of Largescale Evidence Generation and Evaluation across a Network of Databases (LEGEND). *Journal of the American Medical Informatics Association* 2020; 27(8): 1331–1337.
38. Garrido MM, Kelley AS, Paris J, Roza K, Meier DE, Morrison RS, Aldridge MD. Methods for constructing and assessing propensity scores. *Health Serv Res* 2014 Oct;49(5):1701-20.
39. Schuemie M, Suchard MA, Nishimura A, Zhang L, Hripcsak G. Evaluating confounding adjustment when sample size is small. *Observational Health Data Sciences and Informatics (OHDSI) Symposium*; 2023 October 20; 2023.

**Table 1. Description and Characteristics of Administrative Claims Data Sources**

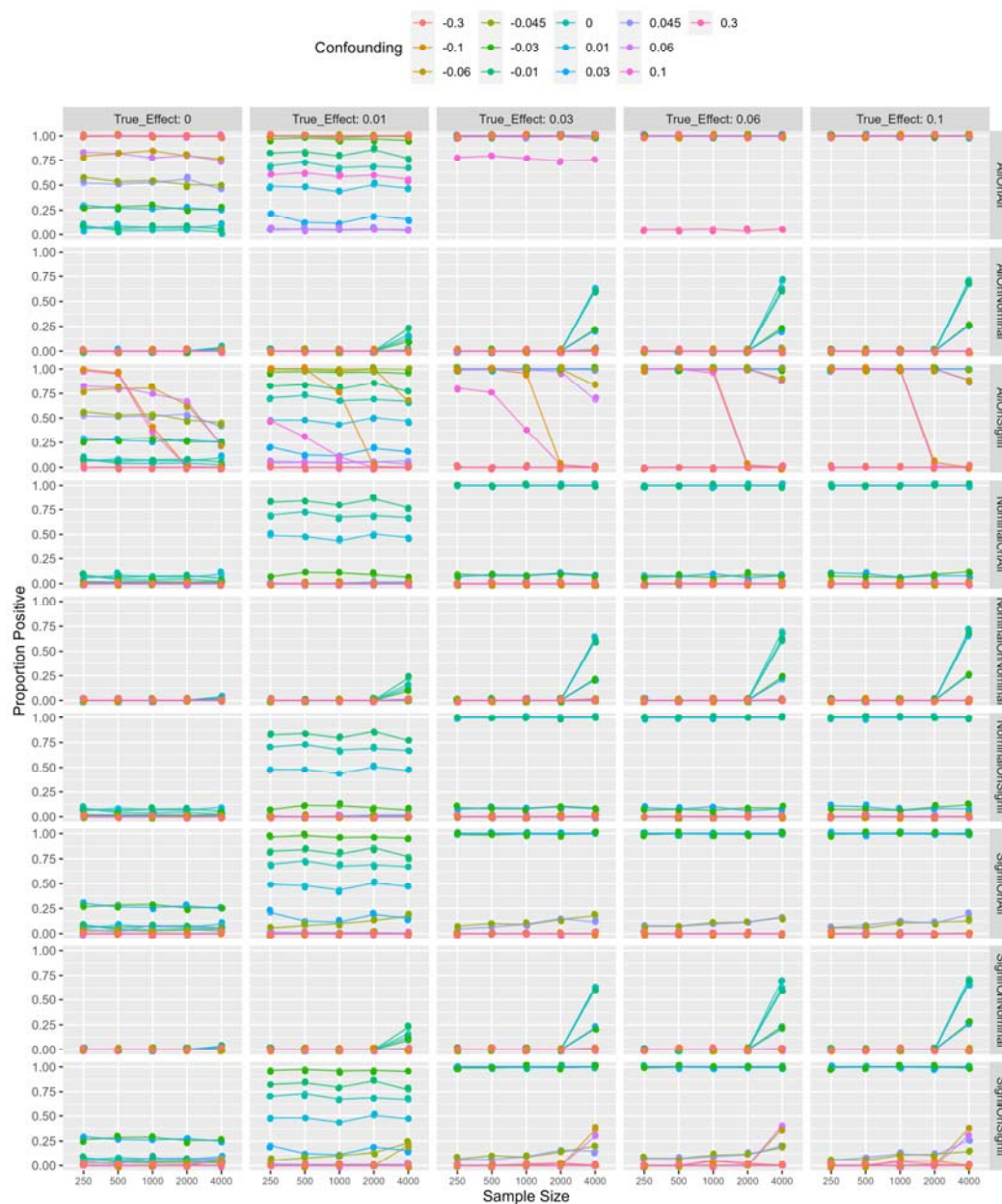
<b>Data Source</b>	<b>Claims Type</b>	<b>Population Enrolled</b>
Merative Medicare Supplemental Database (MDCR)	Adjudicated health insurance claims of retirees with primary or Medicare supplemental coverage through privately insured fee-for-service, point-of-service or capitated health plans.	10M starting 2000 Commercially insured, 65+ years
Merative MarketScan Multi-State Medicaid Database (MDCD)	Adjudicated health insurance claims for Medicaid enrollees from multiple states and includes hospital discharge diagnoses, outpatient diagnoses and procedures, and outpatient pharmacy claims.	26M starting 2006 Medicaid enrollees, racially diverse
Optum® de-identified Electronic Health Record data set (Optum® EHR)	Clinical information, prescriptions, lab results, vital signs, body measurements, diagnoses and procedures derived from clinical notes using natural language processing.	93M starting 2006 US, general

**Figure 1. Rule performance at the database level on simulation.** The proportion of study iterations that were not rejected by the rule and that had effect coefficients that were statistically significantly different from zero is plotted against database sample size. Colored lines represent different levels of confounding with  $c_i$  from  $-0.3$  to  $0.3$ , and graphs from left to right show different values for effect parameter  $c_e$  from  $0$  to  $0.1$ . The rows represent the three types of rules applied only to a single database under study: **All** ignores imbalance, **Nominal** tests for any covariate's standardized mean difference reaching or exceeding  $0.1$ , and **Signif** tests for any covariate's standardized mean difference statistically significantly reaching or exceeding  $0.1$ . The first column, with  $c_e=0$ , represents the type 1 error rate (i.e., higher proportion positive is undesired because it means higher type 1 error rate), and the other columns, with  $c_e>0$ , represent the power with increasing effect size (i.e., higher proportion positive is desired because it means more power in detecting an effect). **All** has unacceptably high type 1 error with high confounding, **Nominal** has unacceptably low power, and **Signif** has moderate type 1 error rate and high power when confounding is low, and when confounding is high, more studies ought to be rejected so low power is expected. (Graph points are jittered to reveal overlapping colors but lines are drawn true.)





**Figure 2. Rule performance at the network level on simulation.** Graphs show the proportion of study iterations that were not rejected by the rule and that had effect coefficients that were statistically significantly different from zero plotted against database sample size. Colored lines represent different levels of confounding  $c_t$  from  $-0.3$  to  $0.3$ , and graphs from left to right show different values for effect parameter  $c_e$  from 0 to 0.1. The nine rows represent the nine rules listed in the Methods section. The first column, where  $c_e=0$ , shows the type 1 error rate, and the other columns, where  $c_e>0$ , show the power with increasing effect size. See text for an explanation of results. (Graph points are jittered to reveal overlapping colors but lines are drawn true.)



**Figure 3. Rule performance at the network level on real-world data.** Graphs show the proportion of study iterations that were not rejected by the rule and that had effect coefficients that were statistically significantly different from zero plotted against database sample size. Colored lines represent adjusted and unadjusted analyses, reflecting lower and higher residual confounding. Graphs from left to right show different true relative risks (RR). The nine rows represent the nine rules listed in the Methods section. The first column, where RR=1, shows the type 1 error rate, and the other columns, where RR>1, show the power with increasing effect size. See text for an explanation of results. (Graph points are jittered to reveal overlapping colors but lines are drawn true.)

