

REVIEW TOPIC OF THE WEEK

Comparison of Propensity Score Methods and Covariate Adjustment



Evaluation in 4 Cardiovascular Studies

Markus C. Elze, PhD,^{a,b} John Gregson, PhD,^a Usman Baber, MD, PhD,^c Elizabeth Williamson, PhD,^a Samantha Sartori, PhD,^c Roxana Mehran, MD, PhD,^c Melissa Nichols, PhD,^{d,e} Gregg W. Stone, MD, PhD,^{d,e} Stuart J. Pocock, PhD^a

ABSTRACT

Propensity scores (PS) are an increasingly popular method to adjust for confounding in observational studies. Propensity score methods have theoretical advantages over conventional covariate adjustment, but their relative performance in real-world scenarios is poorly characterized. We used datasets from 4 large-scale cardiovascular observational studies (PROMETHEUS, ADAPT-DES [the Assessment of Dual AntiPlatelet Therapy with Drug-Eluting Stents], THIN [The Health Improvement Network], and CHARM [Candesartan in Heart Failure-Assessment of Reduction in Mortality and Morbidity]) to compare the performance of conventional covariate adjustment with 4 common PS methods: matching, stratification, inverse probability weighting, and use of PS as a covariate. We found that stratification performed poorly with few outcome events, and inverse probability weighting gave imprecise estimates of treatment effect and undue influence to a small number of observations when substantial confounding was present. Covariate adjustment and matching performed well in all of our examples, although matching tended to give less precise estimates in some cases. PS methods are not necessarily superior to conventional covariate adjustment, and care should be taken to select the most suitable method. (J Am Coll Cardiol 2017;69:345-57) © 2017 by the American College of Cardiology Foundation.

Evaluations of therapeutic interventions generally fall into 2 categories, observational studies and randomized controlled trials (RCTs). The choice of treatment in observational studies may be influenced by patient characteristics, for example, higher-risk patients may be more or less likely to receive the intervention. Some of these

differences are collected in standard databases, whereas others are not (e.g., frailty). In contrast, when studying the effect of an intervention in RCTs, confounding from both measured and unmeasured variables is avoided, and RCTs are thus generally considered the highest form of scientific investigation. Nonetheless, accurate treatment effect



Listen to this manuscript's
audio summary by
JACC Editor-in-Chief
Dr. Valentin Fuster.



From the ^aDepartment of Biostatistics, London School of Hygiene and Tropical Medicine, London, United Kingdom; ^bInnovative Pediatric Oncology Drug Development, F. Hoffmann-La Roche AG, Basel, Switzerland; ^cCardiovascular Institute, Icahn School of Medicine at Mount Sinai, New York, New York; ^dDivision of Cardiology, Columbia University Medical Center, New York-Presbyterian Hospital, New York, New York; and the ^eCardiovascular Research Foundation, New York, New York. Dr. Elze is an employee of and owns stock in F. Hoffmann-La Roche AG. Dr. Mehran has received research support from Eli Lilly/Daiichi-Sankyo Inc., Bristol-Myers Squibb, AstraZeneca, The Medicines Company, OrbusNeich, Bayer, CSL Behring, Abbott Laboratories, Watermark Research Partners, Novartis Pharmaceuticals, Medtronic, AUM Cardiovascular Inc., Cardiovascular Inc., and Beth Israel Deaconess Medical Center; is a compensated member of the Janssen Pharmaceuticals and Osprey Medical Executive Committees; has financial relationships with Watermark Research Partners; is a consultant for Medscape, The Medicines Company, Boston Scientific, Merck & Company, Cardiovascular Systems Inc., Sanofi USA, Shanghai BraccoSine Pharmaceutical, and AstraZeneca; holds equity in Claret Medical and Elixir Medical; is a consultant for Medscape, The Medicines Company, Boston Scientific, Merck & Company, Cardiovascular Systems Inc., Sanofi USA, Shanghai BraccoSine Pharmaceutical, and AstraZeneca; and holds equity in Claret Medical and Elixir Medical. All other authors have reported that they have no relationships relevant to the contents of this paper to disclose.

Manuscript received September 16, 2016; accepted October 19, 2016.

ABBREVIATIONS AND ACRONYMS

BMI = body mass index

HPR = high platelet reactivity

IPW = inverse probability
weighting

MACE = major adverse
cardiovascular event(s)

PS = propensity score(s)

RCT = randomized controlled
trial

estimates from observational databases can provide complementary value to RCTs. This is particularly true when RCTs enroll highly selected patients (yielding results not generalizable to all real-world scenarios), are small (because of their greater complexity and cost), or are not feasible to conduct (1).

The conventional method used to adjust for baseline differences between treatment groups in observational databases is covariate adjustment, where all relevant patient characteristics are included in a regression model relating the outcome of interest to the alternative treatments. A commonly cited concern is that such models might be overfitted when the number of covariates is large compared with the number of patients or outcome events. Although a rule of thumb is to have at least 10 events per covariate included in the model (2), more recent opinions favor relaxing this rule (3).

Propensity score (PS) methods are increasingly being used in observational studies of cardiovascular interventions as an alternative to conventional covariate adjustment; many such examples can be found published in the *Journal* (4-7). A PS is defined as the probability of a patient being assigned to an intervention, given a set of covariates (8). As the PS summarizes all patient characteristics into a single covariate, it reduces (although does not eliminate [9]) the potential for overfitting. PS methods aim to achieve some of the characteristics of RCTs by compensating for different patients having different probabilities of being assigned to the exposures under investigation. Thus, the aim of these methods is to attenuate problems of confounding of patient characteristics and assignment to an intervention typically found in observational studies.

Popular PS methods include stratification, matching, inverse probability weighting (IPW), and use of the PS as a covariate in a conventional regression model (10-12). However, there is lack of clear guidance as to how to make a sensible choice from among these various PS methods or conventional covariate adjustment for any given database. We therefore applied several PS methods to 4 large-scale observational cardiovascular datasets to critically examine the specific advantages and pitfalls of the different methods and to compare their results with those using classic covariate adjustment.

METHODS

DATASETS. We analyzed data from the CHARM (Candesartan in Heart Failure-Assessment of Reduction in Mortality and Morbidity) program (13), the

ADAPT-DES study (14), the THIN study (15), and the PROMETHEUS study (16). For each dataset, we focused on 1 “treatment” comparison and 1 outcome of prime interest. The overall goal was to produce relevant PS models across a range of different settings, so for some cases these choices differed from the primary objectives of the original publications. The terms *treatment* and *control* are used throughout to simplify the language, even though 1 study (14) performed comparisons for platelet reactivity. All outcomes studied were time-to-event, with censoring occurring at the end of planned follow-up, or at the time of patient withdrawal or lost to follow-up.

The CHARM program (13) randomized 7,599 patients with chronic heart failure to candesartan or placebo therapy, with a median follow-up of 3.1 years. We investigated the association between treatment with beta-blockers at baseline (3,396 untreated, 4,203 treated) and all-cause death (1,831 events). That is, we used the CHARM program as an observational database for making inferences about the association between use of beta-blockers and risk of mortality. Our PS model contains cardiovascular risk factors (age, sex, body mass index [BMI], smoking, diabetes), as well as prior cardiovascular events and hospitalizations (18 variables in all).

The ADAPT-DES study (14) investigated the relationship between high platelet reactivity (HPR) in patients taking clopidogrel (HPR: $n = 4,930$; no HPR: $n = 3,650$) and stent thrombosis and other cardiovascular events at 12 months’ follow-up in a prospective, multicenter registry of patients receiving drug-eluting stents. Herein we focused on stent thrombosis (56 events). The study authors reported an adjusted hazard ratio (HR) of 2.49 for HPR compared to patients without HPR. Our PS model contained information about age, sex, medication, diabetes, ethnicity, smoking, renal function, and other cardiovascular risk factors (39 variables in all).

The THIN population-based cohort study (15) compared 30,811 statin users with 60,921 patients not using statins, treated by the same general practitioners (total: $n = 91,732$) for several outcome events, including all-cause mortality (17,296 events, HR: 0.79). The inclusion criteria required at least 12 months of follow-up; thus, the first year must be excluded due to so-called immortality bias. Herein we investigated the effect of statin use on all-cause mortality. The study authors reported an adjusted HR of 0.78, comparing statin users with nonusers. Previously, a large RCT (16) in a similar patient population found an HR of 0.87. Our PS model contained cardiovascular risk factors, age, sex, BMI, smoking, drinking, other medications, and other diseases (48 variables in all).

The PROMETHEUS cohort study (17) compared prasugrel (treatment) therapy with clopidogrel (control) therapy for major adverse cardiovascular event (MACE) outcomes (death, myocardial infarction, stroke, or unplanned revascularization) at 90 days (1,580 events) in 19,914 patients (4,017 prasugrel; 15,587 clopidogrel) using databases from 8 U.S. hospitals. The investigators reported an unadjusted HR of 0.58 and an adjusted HR using a PS model of 0.89. Our PS model contained cardiovascular risk factors, age, sex, BMI, smoking, and prior cardiovascular events, as well as details about the implanted stent and an indicator for study center (35 variables in all).

PROPENSITY SCORES: A BRIEF OVERVIEW. The PS for an individual is defined as the probability of being assigned to “treatment” given all relevant covariates (8). The PS is typically estimated using a logistic regression model that incorporates all variables that may be related to the outcome and/or the treatment decision. All such variables should be included in the logistic model, regardless of their statistical significance or collinearity with other variables in the model. However, variables that are exclusively associated with the treatment decision but not the outcome should not be incorporated (18). As in any predictive regression model, any variable collected after the treatment decision should not be used. As far as possible, covariates identified as relevant in the 4 original studies will be incorporated in the PS models used here. Note that all relevant variables remain in the model regardless of their statistical significance.

For each covariate, individuals with the same PS should have, on average, the same distribution of that covariate irrespective of treatment decision (*covariate balance*). This can be checked using plots of the covariate balance or several diagnostic tests.

After the PS has been calculated, there are several options for how to use them to estimate “treatment” effects. Note throughout that, although PS methods strive to estimate the true “treatment effect,” the usual caveats for observational studies apply, such as the inability to include all relevant confounders (especially those unmeasured). As described later, popular PS methods include matching or stratifying observations on the basis of the PS, IPW applied to each observation, or simply including the PS as an additional variable in a regression model. The more conventional covariate adjustment offers an alternative to PS techniques by simply incorporating all relevant covariates into the final model (19).

FOUR PS METHODS. For each dataset, the goal was to estimate the “treatment” effect on a time-to-event outcome, using Cox proportional hazards models.

After creating the PS for each individual, there are several ways to adjust for confounding.

PS stratification. PS stratification splits the dataset into several strata on the basis of the individual’s PS alone, without reference to treatment (exposure) group. A treatment effect is then estimated within each stratum, and an overall estimated treatment effect is calculated by taking a (weighted) average across strata. Here, 5 and 10 strata, with an equal number of individuals in each stratum, are used. An alternative is to split the range of possible PS into equal parts, which usually results in fewer individuals in the more extreme strata. Stratification has the additional advantage that effect estimates are available for each stratum, which may reveal potential heterogeneity of “treatment” effects across strata.

PS matching. PS matching tries to find 1 (or more) individual(s) with similar PS in the treatment and control groups. There are various methods to match individuals, but here we use 1:1 nearest-neighbor matching, with an added constraint that the difference between the PS (*caliper width*) may be at most 0.1 to avoid pairing dissimilar individuals. We chose this method for its computational simplicity. Following matching, the treatment effect is calculated by applying either a conventional (unmatched) regression model or a matched pair analysis to the set of patients who are successfully matched (20). We opt for an unpaired analysis here due to its greater simplicity, noting that in our examples, a paired analysis gave almost identical results (Online Table 1).

The matching process results in an analysis based upon only those patients who are successfully matched. Therefore, if the treatment effect varies according to patients’ characteristics and their likelihood of receiving treatment, the treatment effect estimated from this subset of patients may differ from the effect in the original study population. This issue is covered in greater detail in the discussion.

Inverse probability weighting. Inverse probability weighting uses the whole dataset but reweights individuals to increase the weights of those who received unexpected exposures. This procedure can be thought of as producing additional observations for those parts of the target population from which there were few observations. It effectively generates a pseudopopulation with near-perfect covariate balance between treatment groups (12). IPW applies weights corresponding to $1/PS$ for patients in the treated cohort and $[1/(1 - PS)]$ for those in the control cohort. Due to the large weight assigned to

these observations, PS close to 0 (for the treated) or 1 (for the control) may be problematic for IPW. We discuss methods to resolve this issue, such as trimming or truncating large weights, later.

Although these 3 PS methods aim to balance all covariates between the treatment and control groups, the more conventional *covariate adjustment* aims to control for covariate effects (confounding) using a prediction model for the outcome event (in our case a proportional hazards model for a time-to-event outcome). Care must be taken to specify the correct functional form for any covariates that may have nonlinear effects. Covariate adjustment has its critics, but there is little practical evidence that it gives misleading results. For comparison, we provide the crude effect estimate, as well as the covariate-adjusted effect estimate using all covariates from the PS models.

Including the PS as an additional covariate. Including the PS as an additional covariate in the regression model represents the fourth PS method examined. Alternatively, one could have only the PS and treatment in a model of the outcome of interest.

VARIATIONS ON PS METHODS. There are several variations on the 4 PS methods presented in the preceding section. Using a doubly robust approach can compensate for a lack of covariate balance. The dataset can be pre-processed by “trimming” away (removing) individuals with extreme PS. Alternatively, large IPW weights can be avoided by truncating the weights.

Doubly robust methods incorporate relevant covariates in both the PS model and the outcome regression model for the treatment effect: this can compensate for insufficient covariate balance (21). As the name implies, this approach offers some robustness to model misspecification, either in the PS or the outcome regression model. It is recommended (8,22) when using the PS as a covariate to also include individual covariates in the outcome regression model. When the method is used in this way, as we do throughout this report, it is doubly robust. Doubly robust methods are also commonly used with IPW but less frequently with matching or stratification, possibly due to the reduced sample size when using these methods. However, the doubly robust approach removes a key advantage of PS methods: having only 1 covariate in the final model.

Trimming can be performed after calculation of the PS. This involves dropping the individuals with most extreme PS values in both the treatment and control groups, as they may lack a match in the other group and can be predisposed to residual

confounding. This can help avoid extreme weights in IPW, improve comparability between the exposures, and remove unusual “outlying” patients for whom the expected treatment (or control) was not chosen. Typical trimming methods might remove the most extreme 1% or 5% of all observations.

Weight truncation reduces any “large” weight down to a maximum weight. There is no standard definition of a large weight for IPW. Here, we considered any weight above 10 to be large, and reduced any weight >10 down to this threshold. Removal of large weights is sometimes recommended for sensitivity analyses. However, complete removal of all individuals with weights larger than 10 may increase the imbalance between treatment and control groups.

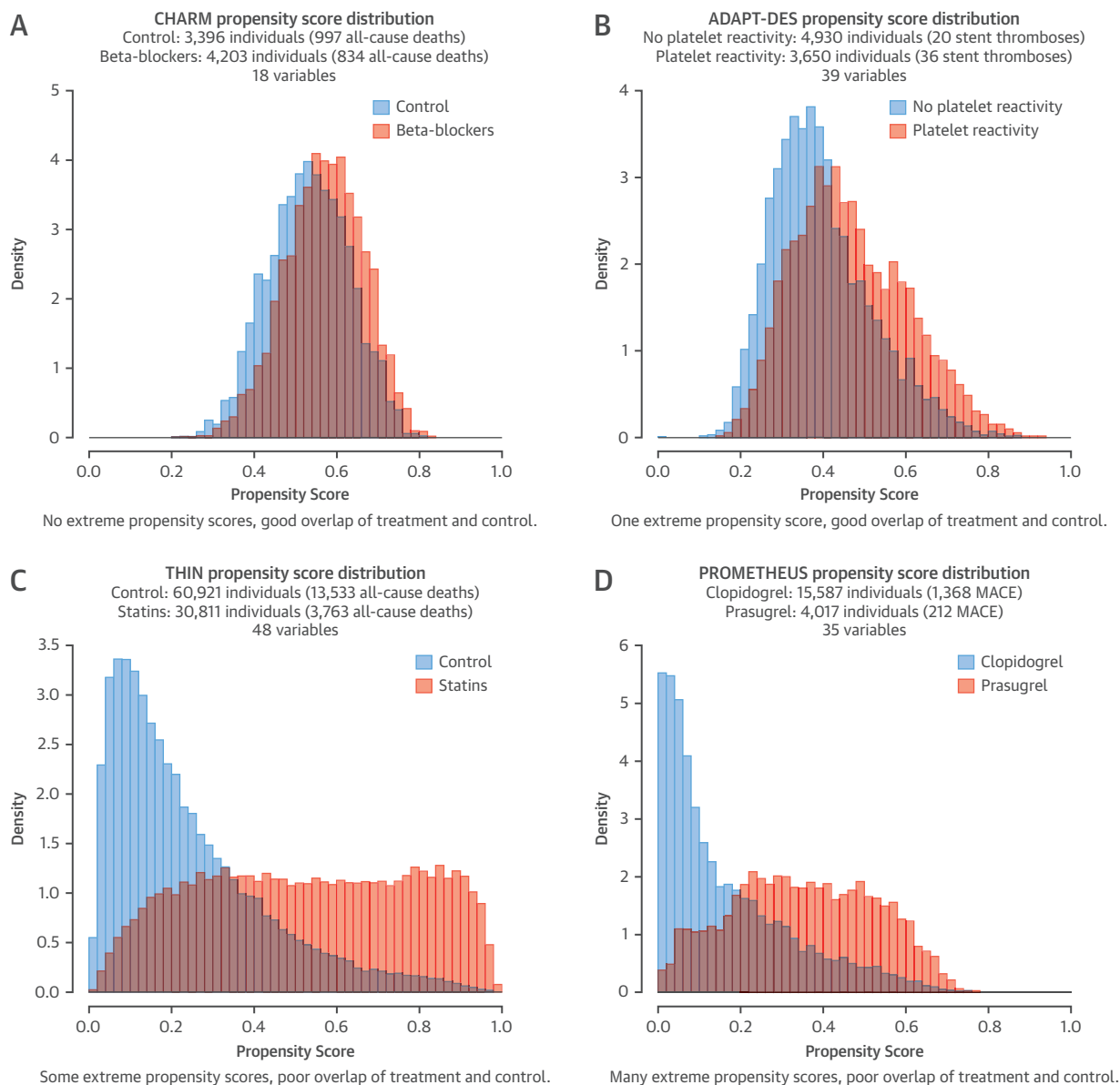
STANDARD ERRORS AND P VALUES. All SE reported here are given on the effect (i.e., log HR) scale, and we used the usual sandwich variance estimator when using IPW (23). The calculation of p values can then be done in the usual fashion. A special case is stratification, where it is necessary to aggregate SEs and p values from multiple models. This is done by calculating the overall variance for a particular parameter as the weighted average of the variances for that parameter from each stratum and dividing by the number of strata (24). Assuming asymptotic normality on the overall effect, p values can then be calculated.

RESULTS

Propensity scores for the CHARM, ADAPT-DES, THIN, and PROMETHEUS studies showed a range of different distributions (Figure 1). Full PS models are given in Online Tables 2 to 5, and for comparison, covariate-adjusted models are given in Online Tables 6 to 9. Both CHARM and ADAPT-DES exhibited good overlap between the PS for the treatment and control groups. A single individual in the control group for ADAPT-DES had a PS close to 0 and could be considered an outlier.

In contrast, the THIN and PROMETHEUS studies showed markedly different PS distributions for the treatment and control groups. This indicates that it may be difficult to provide valid comparisons between the 2 groups. THIN had a substantial number of PS close to 0 or 1. There were 1,134 treated patients (4% of all treated) and 15,514 control patients (25% of all control) with a PS <0.1 . Conversely, there were 2,235 treated individuals (7% of all treated) and 173 control patients (0.3% of all control) above a PS of 0.9. Clearly, there are key variables in the PS that played an important role in who did (and did not)

FIGURE 1 Overview of the PS Distribution



PS distributions for the control (blue) and treatment (orange) groups for (A) CHARM (Candesartan in Heart Failure-Assessment of Reduction in Mortality and Morbidity), (B) ADAPT-DES (the Assessment of Dual AntiPlatelet Therapy with Drug-Eluting Stents), (C) THIN (The Health Improvement Network), and (D) PROMETHEUS. Greater overlap in PS indicates a lesser risk of confounding by indication. PS = propensity score.

receive a statin. Patients where PS and chosen exposure strongly disagreed (high PS but received control; low PS but received treatment) may be atypical but received large IPW weights.

PROMETHEUS had a very large number of PS close to 0, especially in the control group receiving clopidogrel (7,282 control individuals below PS 0.1, 47% of the control group). This indicates that key variables in

the PS had a marked influence on physician choice of clopidogrel rather than prasugrel. There were also a considerable number of patients in the treatment group receiving prasugrel with a PS close to 0 (330 treated individuals below PS of 0.1, 8% of the treated group). These individuals may be unusual and may not offer a representative comparison with the other group.

CHARM. Results for CHARM, a nonrandomized comparison of the effect of beta-blocker use versus control on all-cause death, showed excellent agreement across all PS methods and covariate adjustment (Figure 2A). As expected, the crude estimate (Figure 2A, first row) was different from the covariate-adjusted estimate (Figure 2A, second row) or estimates provided by the different PS methods (Figure 2A, other rows). The adjusted HRs were all ~ 0.73 , with 95% confidence intervals (CIs) of ~ 0.65 to 0.81. SEs were very similar across all methods, and p values were highly significant for all methods.

ADAPT-DES. The ADAPT-DES study, which investigated the relationship between HPR and the risk of stent thrombosis, produced similar HRs for most methods (Figure 2B). Covariate adjustment, matching, IPW, and use of the PS as a covariate all arrived at an HR of ~ 2.2 , comparing patients with HPR with those without HPR. A notable exception is stratification, which showed a wider CI with 5 strata and an unstable results with 10 strata. Otherwise, SEs and p values were comparable for all methods, although matching had slightly poorer precision.

An investigation of the strata for ADAPT-DES revealed that the relatively low number of 56 events in the dataset were divided unevenly in the 10 strata (Online Table 10). Two strata received only a single event, making precise estimation of the treatment effect within those strata impossible. These findings strongly suggest that stratification with this many strata should not be used when the number of events is sparse.

THIN. The different PS methods and covariate adjustment mostly produced similar results for the THIN study, arriving at HR of ~ 0.85 and a highly significant mortality reduction for those taking a statin (Figure 2C). The exception was IPW, which estimated a smaller treatment effect with a wider CI. Trimming individuals with extreme PS from IPW gave similar results, whereas truncating large weights in IPW brought the HR in line with the other methods. Similarly, a doubly robust approach of including all covariates in the final regression model also brought the HR into agreement with the other approaches. Additionally, a strong influence of confounders in this database was noted: the crude HR of 0.55 greatly exaggerated the treatment effect due to individuals on statins tending to be at lower mortality risk. Note that from RCTs, an HR of approximately 0.87 is expected.

A plot of the IPW weights revealed very large weights for some individuals (Figure 3A), which may be why IPW produced different results from

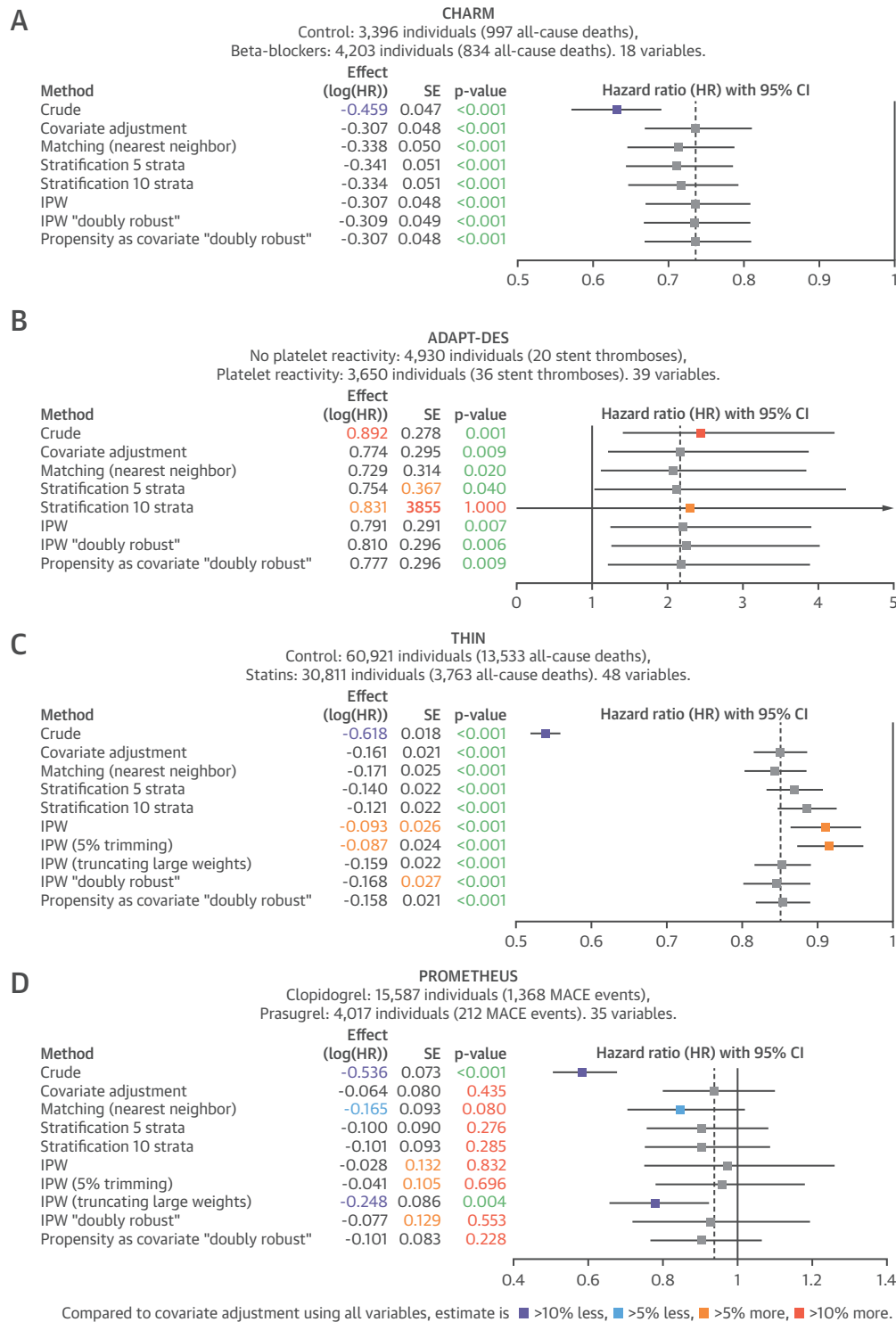
other methods. For 1,307 patients, weights exceeded 10. This group of patients, 1.4%, had the same total weight as the 22% of patients with the lowest weights, which may have given undue influence to very few observations. This is especially problematic considering that those large weights were given to the most unusual individuals, as most were patients taking statin treatment who the PS model strongly predicted would be controls (i.e., not taking a statin).

PROMETHEUS. Results for the PROMETHEUS study (Figure 2D), comparing prasugrel with clopidogrel for risk of MACE, showed substantial disagreement between the methods, although the results were nonsignificant for almost all methods. Covariate adjustment, stratification, IPW, and use of the PS as a covariate all produced HRs of ~ 0.94 . Matching showed a lower HR of ~ 0.85 . Inverse probability weighting without any modification had a much higher SE than the other methods. Investigating the IPW weight distribution revealed very large weights for 8% (330 individuals) of the treatment group (Figure 3B), which may explain the stark change in HR seen when truncating large weights. The crude estimate of treatment (HR: 0.59) is attributable to the marked confounding present, whereby patients chosen to receive prasugrel tended to be at lower risk of MACE.

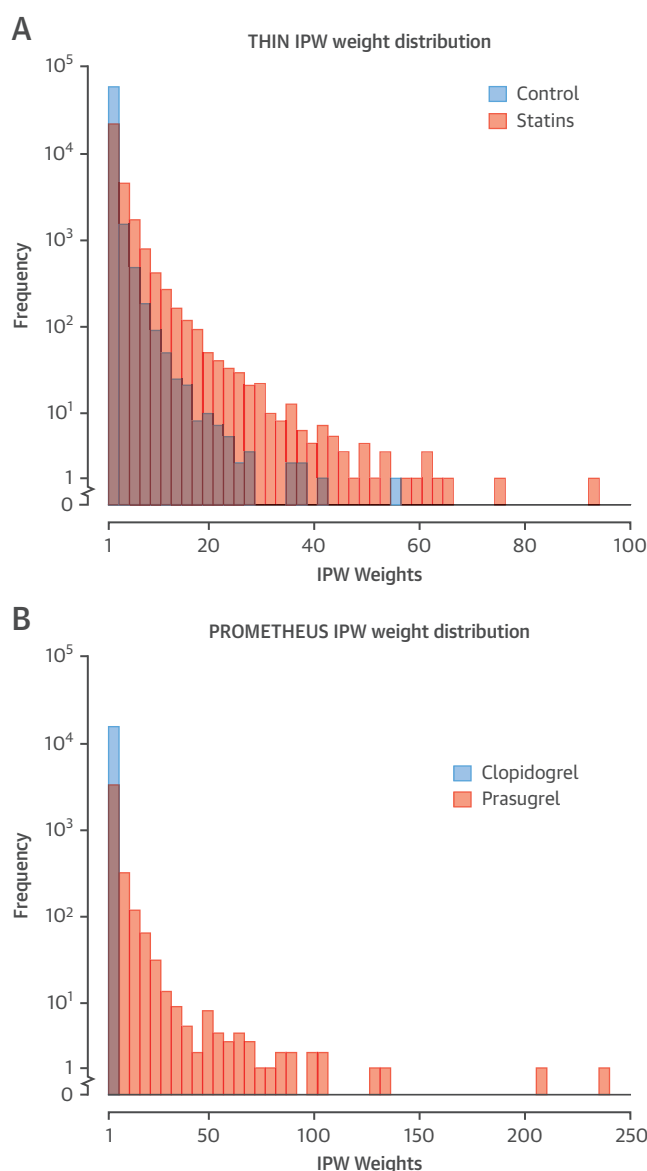
Further examination showed that covariate balance is insufficient for some methods. Figure 4 compares the covariate balance values for matching, stratification, and IPW by using the absolute standardized difference between the treatment and control groups. Without use of PS methods, covariate balance was insufficient for almost all variables. Matching produced excellent balance for all variables. Stratification mostly achieved satisfactory covariate balance, except for previous percutaneous coronary intervention, with age, hypertension, and prior congestive heart failure as borderline cases. Inverse probability weighting showed very poor covariate balance for previous percutaneous coronary intervention and poor or borderline balance for hypertension, previous myocardial infarction, and prior peripheral artery disease. Due to the lack of covariate balance, the results for stratification and IPW may be considered unreliable.

EFFECT OF TRIMMING AND TRUNCATION OF IPW WEIGHTS. We used PS trimming in the THIN and PROMETHEUS studies to attempt to reduce the impact of large weights in IPW. However, for both studies, even 5% trimming was not sufficient to fully remove all large IPW weights from the datasets. Consequently, the SEs for the estimated treatment effect remained

FIGURE 2 Comparison of Hazard Ratios From Different PS Methods and Covariate Adjustment



HRs and covariate adjustments for (A) CHARM, (B) ADAPT-DES, (C) THIN, and (D) PROMETHEUS. In the plot, covariate adjustment is used as the basis for the comparison (dashed line). Colors are used if results for the other methods differed by more than 5%. CI = confidence interval; HR = hazard ratio; IPW = inverse probability weighting; other abbreviations as in Figure 1.

FIGURE 3 Distribution of the Weights for IPW in the THIN and PROMETHEUS Studies

Distributions in (A) THIN and (B) PROMETHEUS. To facilitate display, the vertical axis is on a logarithmic scale. In (B) PROMETHEUS, there are no patients with extreme weights in the clopidogrel group; hence, all patients treated with clopidogrel will appear in the bar with the smallest inverse probability weight. Abbreviations as in Figures 1 and 2.

large relative to other methods after trimming, particularly in PROMETHEUS (Figures 2C and 2D). We additionally applied 1% and 5% trimming and compared findings for each PS method and covariate adjustment on the trimmed data sets. However, trimming did little to reconcile differences in the estimates produced (Online Figures 1 to 4). Finally, we truncated large weights in the THIN and PROMETHEUS studies

to a maximum of 10. In both examples, this resulted in a large reduction in SE and an estimated HR much closer to the crude estimate (Figures 2C and 2D). However, this brought the estimates closer to the other methods in THIN (Figure 2C), whereas it took estimates farther away from other methods in PROMETHEUS (Figure 2D).

DISCUSSION

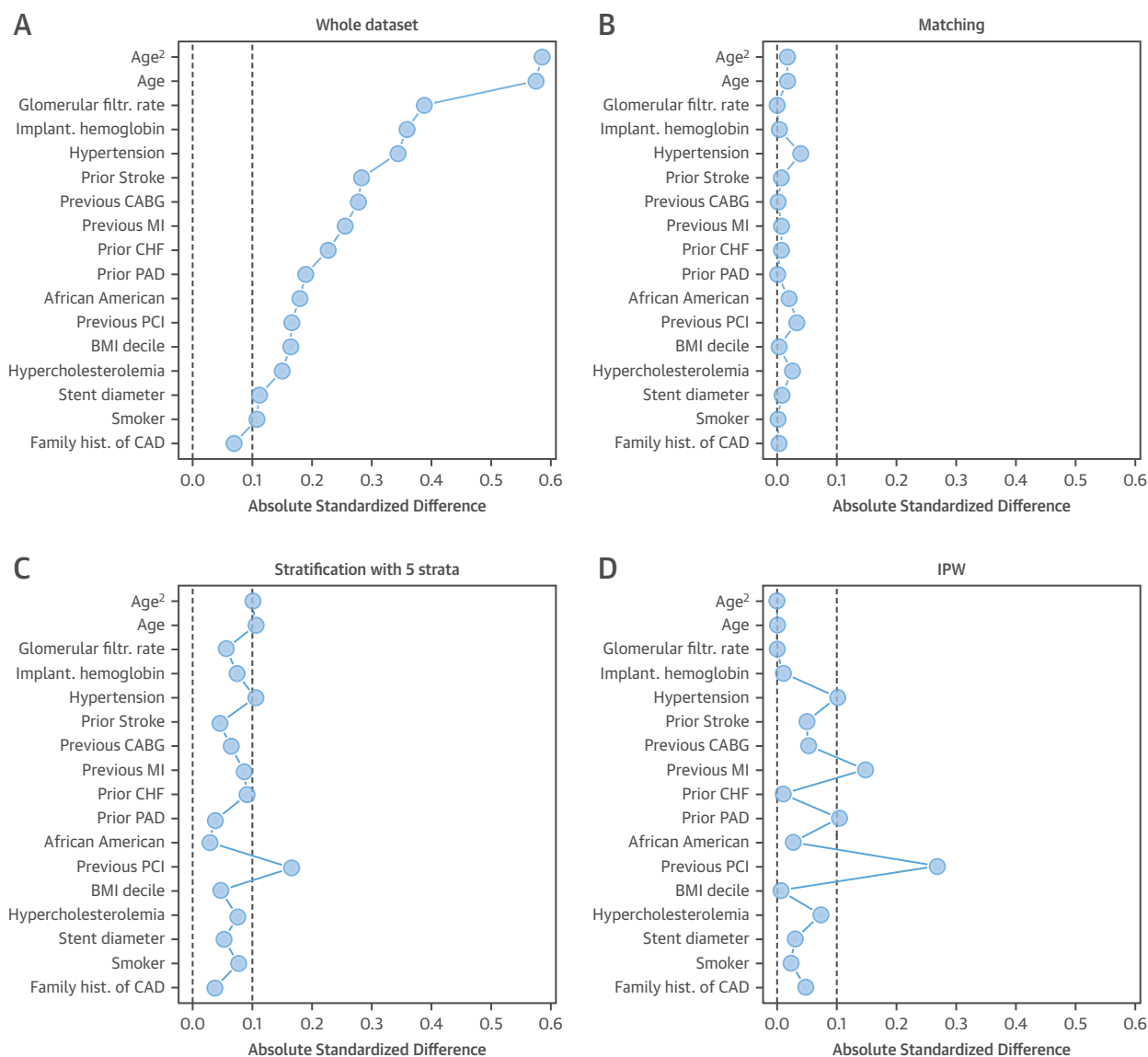
For observational cohort studies that compare alternative treatments (and other exposures), it has become standard practice to use PS methods to correct for selection biases and potential confounding when examining the relative risks (hazards) of event outcomes. Although the principles of PS methods are clear, there exists a diversity of alternative approaches (e.g., propensity matching, stratification, IPW with or without trimming) alongside the more conventional method of covariate adjustment. Although there is a substantial body of methodologic literature on PS approaches with some limited guidance on which specific methods may be preferable (10,11), there is no general agreement as to the choice of PS method that is best suited to any particular scenario. Thus, researchers may choose a suboptimal method that preserves more bias and/or imprecision than necessary.

To provide insight into this common problem, we have here undertaken an in-depth assessment of many of the available PS and covariate-adjustment approaches as applied to 4 large-scale cardiovascular studies. The present analysis illustrates the challenges faced in determining which methods actually produce the most valid results in different settings.

Our first example, the CHARM study examining the impact of beta-blocker use at baseline on mortality in heart failure patients, is the most straightforward. The PS distributions for the 17 chosen baseline variables showed considerable overlap between the 2 groups with no extreme values. In addition, the study was large, with the 2 groups being of similar size. The results showed a consistency across all PS methods and also covariate adjustment. Note, the crude estimate produced an exaggerated treatment effect, indicating the importance of taking confounding into account by using any of these methods. However, the extent of confounding is less extreme than in several of the other studies.

The next example, ADAPT-DES, comparing the risk of stent thrombosis in acute coronary syndrome patients with and without HPR, has some methodologic similarities to CHARM but also the complication of

FIGURE 4 Comparison of the Extent of Covariate Imbalance in PROMETHEUS



Comparisons using (A) crude comparisons, (B) propensity matching, (C) propensity stratification, and (D) IPW. Graphs show the absolute standardized difference between treatment and control; values <0.1 are conventionally considered acceptable. BMI = body mass index; CABG = coronary artery bypass graft; CAD = coronary artery disease; CHF = congestive heart failure; MI = myocardial infarction; PAD = peripheral arterial disease; PCI = percutaneous coronary intervention; other abbreviations as in Figures 1 to 3.

having fewer outcome events (only 56 stent thromboses). Here, PS stratification performed badly, with too few events per stratum. PS matching and IPW showed good agreement, although the former had less precise estimates due to not using all patients in the matched analysis. Surprisingly, covariate adjustment with 39 covariates and only 56 events held up well, producing very similar estimates to IPW.

In the last 2 examples, both the THIN study, comparing the mortality of individuals on and off statins, and the PROMETHEUS study, comparing prasugrel versus clopidogrel for the risk of MACE in acute coronary syndrome patients, presented more of a challenge in choosing a robust PS method. The reason for this in both cases was the marked separation of PS probability distributions between the 2 groups: statin

versus no statin and prasugrel versus clopidogrel, respectively. In particular, a substantial number of PS were close to 0 and 1 in THIN and close to 0 in PROMETHEUS. As a consequence, IPW included more than a few very influential individuals with very large weights in the IPW analysis. This, in turn, led to imprecise estimates of treatment effect and a worrisome lack of covariate balance for some potentially important confounders. Additionally, in both these examples, IPW analyses estimated HRs closer to the null.

In both of these examples, the use of IPW in a doubly robust fashion (i.e., also including all covariates in the final analysis) induced compatibility with other methods but did not reduce the SE, thereby leaving the 95% CIs unduly wide. The use of trimming (e.g., removing the 5% of individuals with the most extreme PS) was somewhat helpful, but the imprecision of estimates remained greater than for other methods. The use of PS stratification also has its problems when there is marked selection bias, as seen in THIN and PROMETHEUS. This is because using only 5 (or 10) strata does not wholly correct for covariate imbalance.

What can we learn from these experiences in order to make recommendations for the future use of PS methods and covariate adjustment? As in all studies, the primary analysis strategy should be prespecified. Post hoc selection of a preferred method after data exploration introduces bias and should only be considered for exploratory or sensitivity analysis.

One useful approach is to examine the baseline covariates before accessing any outcome data in order to determine which PS method (or covariate adjustment) may be most suitable given the characteristics of the PS, such as the degree of overlap in PS between treatment and control groups. Even so, relying on 1 method of analysis (which may have its flaws) may be too restrictive, and it is wise to predefine a number of secondary sensitivity analyses using alternative approaches. This enables one to determine whether there is a consistency of the findings regarding the estimated treatment “effect,” which, if present, instills confidence in the primary results.



However, for any specific study, what should be chosen as the primary analysis method? We see no single “right answer” to meet all circumstances, but the following insights should help in making the choice:

1. PS matching appears to be a reliable method, in that it provides excellent covariate balance in most circumstances. It has the advantage of being simple to analyze, present, and interpret. Its main disadvantage is that some individuals end up not

matched and hence excluded from the analysis, resulting in a loss of both precision and generalizability. In our examples, up to 60% of patients were excluded by matching, although it should be noted that some of these patients could have been successfully matched by using more sophisticated matching algorithms. Finally, whatever the choice of matching algorithm, it is important to predefine the precise algorithm to be used.

2. PS stratification tends to work well when covariate imbalance is not very marked. It has the merit of keeping all individuals in the analysis and also provides the opportunity to explore potential interactions between treatment and the PS on outcome risk. Stratification tends to perform less well in datasets with few outcomes, particularly when the number of strata is large. When choosing the number of strata, one needs to trade off the need for accurate control of confounding with the requirement of having a sufficient number of events in each strata. Previous research shows that 5 strata may reduce confounding bias by up to 90%, so a modest number of strata should suffice in studies with few outcomes and/or only moderate confounding bias (25). However, in studies with many outcome events, using more strata will further reduce confounding bias, which may be important if covariate imbalance is marked (26,27). Beyond these recommendations, further research is needed to determine the best strategy to define the number and size distribution of strata. Are strata of equal size preferred, or is it better to have larger numbers in the middle of the PS distribution, therefore enabling a more detailed exploration of the tails?
3. Inverse probability weighting offers a conceptually simple method that is easy to implement in practice and retains all study participants. Some have advocated it as a preferred method (28,29). However, when there is marked covariate imbalance, PS scores close to the extreme probabilities of 0 and 1 occur, with some individuals ending up with very large weights. This seems intuitively inappropriate because these influential data points occur in individuals who represent a small proportion in their chosen treatment group. In our 2 examples with such extreme weights, THIN and PROMETHEUS, use of IPW produced less precise estimates than other methods and notable covariate imbalance.
4. Trimming has been recommended as an appropriate way of limiting the influence of heavily weighted individuals. The difficulty here is to define in advance what level of trimming is

CENTRAL ILLUSTRATION Comparison of Propensity Score Methods and Covariate Adjustment

Primary study analysis method	 Pros	 Cons
Traditional covariate adjustment	<ul style="list-style-type: none"> Performed well Provides prognostic model for outcome of interest 	<ul style="list-style-type: none"> May not be suitable with many covariates in smaller studies
Propensity score (PS) stratification	<ul style="list-style-type: none"> Retains data from all study participants Opportunity to explore interactions between treatment and PS on outcome risk Provides effect estimates for every stratum 	<ul style="list-style-type: none"> Performs less well in datasets with few outcomes, particularly when the number of strata is large May not account for strong confounding
PS matching	<ul style="list-style-type: none"> Reliable; provides excellent covariate balance in most circumstances Simple to analyze, present and interpret 	<ul style="list-style-type: none"> Some patients are unmatched leading to information excluded from the analysis Less precise
PS inverse probability weighting	<ul style="list-style-type: none"> Retains data from all study participants Easy to implement Creates a pseudo population with perfect covariate balance 	<ul style="list-style-type: none"> Can be unstable when extreme weights occur
PS covariate adjustment (use of PS as a covariate)	<ul style="list-style-type: none"> Performed well 	<ul style="list-style-type: none"> Adding the PS as an additional covariate produced results very similar (and not necessarily superior) to traditional covariate adjustment

Elze, M.C. et al. *J Am Coll Cardiol.* 2017;69(3):345-57.

An overview of the pros and cons of covariate adjustment and various propensity score methods.

desirable: should we exclude just 1% or up to 5% of extreme weights? In our examples with major covariate imbalance, trimming increased the precision of our estimates, but did not alter the estimated associations. Truncating large weights resulted in more precise estimates and had large effects on the estimated associations. In both of our examples, the estimated associations moved closer to the crude HR following truncation, perhaps suggesting that this method may lead to inadequate adjustment for covariate imbalance. Given the difficulty of limiting the influence of heavily weighted individuals, IPW may be best confined to datasets for which extremes of the PS distribution do not occur, such as in the CHARM and ADAPT-DES studies, although this will generally not be known in advance of examining the data distribution.

5. Use of a doubly robust IPW approach, where covariates are also included in the outcome

regression model, appeared to produce results similar to conventional coverage adjustment, but with notably wider CIs. They also added a level of complexity to the analysis, and the inclusion of covariates in the outcome regression model removes a key advantage of the PS methods. They may therefore be unattractive as a primary method of analysis, and be best reserved for sensitivity explorations.

6. Covariate adjustment is the conventional method for correcting for covariate imbalance, selection bias, and potential confounding and existed long before PS methods were developed. Recently, some have argued that PS methods may be more robust or offer more complete adjustment for confounders (28). However, although there are theoretical grounds on which to favor PS adjustment (30), we see little practical evidence to justify such negative claims and, in our examples, covariate adjustment provided reliable and statistically

efficient estimates (**Central Illustration**). One issue for datasets with few event outcomes is that the number of covariates considered for inclusion in the model may be limited, whereas many more covariates may be included in a PS model without raising concerns of overfitting or lack of model convergence. However, in ADAPT-DES, including 39 covariates with only 56 events still produced reliable results. This example demonstrates that having fewer than 10 events per covariate does not necessarily preclude using covariate adjustment (3), although it does not allay concerns of overfitting in all similar scenarios. A further advantage of covariate adjustment is that it provides a predictive model (including treatment) for the risk (hazard) of the event outcome, which gives insight as to which covariates have the strongest influence on risk. Perhaps it is time that old-fashioned covariate adjustment deserves a revival in its use. Finally, adding the PS as an additional covariate produced results very similar to covariate adjustment, with similar estimates and SEs across all examples.

STUDY LIMITATIONS. First, with just 4 datasets explored in depth, caution is needed in drawing any generalizable conclusions. This is particularly the case for small studies, which are not examined here, although it should be noted that ADAPT-DES is small in terms of the number of outcome events included. Despite these limitations, we believe that the diversity of our examples facilitate a practical debate on the basis of real experiences, which is better than relying on purely theoretical arguments. Second, our study assumes throughout that the effect of treatment on outcome does not differ by the likelihood that a patient is treated. When treatment effects do differ, as can be detected by comparing estimated HRs across strata of the PS, some PS methods will produce results that are systematically different

from covariate adjustment, even when both methods provide adequate adjustment for confounders (31). This is because certain PS methods estimate the treatment effects relating to certain sections of the study population, such as only the treated or only the control patients. In these scenarios, investigators need to select an appropriate method to estimate the treatment effect in the set of patients in whom they most want to understand the impact of treatment: usually this is either the treated patients, the control patients, or the entire study population (11). In addition, a further technical detail is that both IPW and PS matching (using an unpaired analysis) estimate a marginal treatment effect, whereas multivariate regression, stratification, and doubly robust methods all estimate a conditional HR.

CONCLUSIONS

In the present detailed examination of alternative PS methods and covariate adjustment in several topical cardiovascular studies, covariate adjustment and matching performed well in all of our examples, although matching tended to give less precise estimates in some cases. Propensity score methods are not necessarily superior to conventional covariate adjustment, and care should be taken to select the most suitable method. We hope these insights will guide others to make wise choices in their use of PS methods, and to rekindle interest in old-fashioned covariate adjustment, which may be viewed as a suitable primary analysis method in many cases.

REPRINT REQUESTS AND CORRESPONDENCE: Prof. Stuart J. Pocock, London School of Hygiene and Tropical Medicine, Department of Medical Statistics, Keppel Street, London WC1E 7HT, United Kingdom. E-mail: stuart.pocock@lshtm.ac.uk.

REFERENCES

1. Rothwell PM. External validity of randomised controlled trials: "to whom do the results of this trial apply?" *Lancet* 2005;365:82-93.
2. Harrell F Jr., Lee KL, Califf RM, et al. Regression modelling strategies for improved prognostic prediction. *Stat Med* 1984;3:143-52.
3. Vittinghoff E, McCulloch CE. Relaxing the rule of ten events per variable in logistic and Cox regression. *Am J Epidemiol* 2007;165:710-8.
4. Park DW, Seung KB, Kim YH, et al. Long-term safety and efficacy of stenting versus coronary artery bypass grafting for unprotected left main coronary artery disease: 5-year results from the MAIN-COMPARE (Revascularization for Unprotected Left Main Coronary Artery Stenosis: Comparison of Percutaneous Coronary Angioplasty Versus Surgical Revascularization) registry. *J Am Coll Cardiol* 2010;56:117-24.
5. Ramos R, Garcá-Gil M, Comas-Cuf M, et al. Statins for prevention of cardiovascular events in a low-risk population with low ankle brachial index. *J Am Coll Cardiol* 2016;67:630-40.
6. Tamburino C, Barbanti M, D'Errigo P, et al., for the OBSERVANT Research Group. 1-Year outcomes after transfemoral transcatheter or surgical aortic valve replacement: results from the Italian OBSERVANT Study. *J Am Coll Cardiol* 2015;66:804-12.
7. Solomon MD, Go AS, Shilane D, et al. Comparative effectiveness of clopidogrel in medically managed patients with unstable angina and non-ST-segment elevation myocardial infarction. *J Am Coll Cardiol* 2014;63:2249-57.
8. Rosenbaum PR, Rubin DB. The central role of the propensity score in observational studies for causal effects. *Biometrika* 1983;70:41-55.
9. Senn S, Graf E, Caputo A. Stratification for the propensity score compared with linear regression

techniques to assess the effect of treatment or exposure. *Stat Med* 2007;26:5529–44.

10. Austin PC. An introduction to propensity score methods for reducing the effects of confounding in observational studies. *Multivariate Behav Res* 2011;46:399–424.
11. Williamson E, Morley R, Lucas A, et al. Propensity scores: from naive enthusiasm to intuitive understanding. *Stat Methods Med Res* 2012;21:273–93.
12. Heinze G, Jüni P. An overview of the objectives of and the approaches to propensity score analyses. *Eur Heart J* 2011;32:1704–8.
13. Swedberg K, Pfeffer M, Granger C, et al., for the Charm-Programme Investigators. Candesartan in heart failure—assessment of reduction in mortality and morbidity (CHARM): rationale and design. *J Card Fail* 1999;5:276–82.
14. Stone GW, Witzenbichler B, Weisz G, et al., for the ADAPT-DES Investigators. Platelet reactivity and clinical outcomes after coronary artery implantation of drug-eluting stents (ADAPT-DES): a prospective multicentre registry study. *Lancet* 2013;382:614–23.
15. Smeeth L, Douglas I, Hall AJ, et al. Effect of statins on a wide range of health outcomes: a cohort study validated by comparison with randomized trials. *Br J Clin Pharmacol* 2009;67:99–109.
16. Heart Protective Study Collaborative Group. MRC/BHF Heart Protection Study of cholesterol lowering with simvastatin in 20,536 high-risk individuals: a randomised placebo-controlled trial. *Lancet* 2002;360:7–22.
17. Wayangankar SA, Baber U, Poddar K, et al. Predictors of 1 year net adverse cardiovascular

events (NACE) among ACS patients undergoing PCI with clopidogrel or prasugrel: analysis from the PROMETHEUS registry. *J Am Coll Cardiol* 2016;67:562.

18. Myers JA, Rassen JA, Gagne JJ, et al. Effects of adjusting for instrumental variables on bias and precision of effect estimates. *Am J Epidemiol* 2011;174:1213–22.
19. Fisher R. *The Design of Experiments*. 9th edition. London, UK: Macmillan, 1971.
20. Stuart EA. Developing practical recommendations for the use of propensity scores: discussion of “A critical appraisal of propensity score matching in the medical literature between 1996 and 2003” by Peter Austin. *Stat Med* 2008;27:2062–5; discussion 2066–9.
21. Kang JDY, Schafer JL. Demystifying double robustness: a comparison of alternative strategies for estimating a population mean from incomplete data. *Stat Sci* 2007;22:523–39.
22. D’Agostino RB Jr. Propensity score methods for bias reduction in the comparison of a treatment to a non-randomized control group. *Stat Med* 1998;17:2265–81.
23. Joffe MM, Ten Have TR, Feldman HI, et al. Model selection, confounder control, and marginal structural models. *Am Stat* 2004;58:272–9.
24. Mosteller F, Tukey JW. *Data Analysis and Regression: A Second Course in Statistics*. London, UK: Pearson, 1977.
25. Rosenbaum PR, Rubin DB. Reducing bias in observational studies using subclassification on the propensity score. *J Am Stat Assoc* 1984;79:516–24.

26. Cochran WG. The effectiveness of adjustment by subclassification in removing bias in observational studies. *Biometrics* 1968;24:295–313.

27. Hullsiek KH, Louis TA. Propensity score modeling strategies for the causal analysis of observational data. *Biostatistics* 2002;2:179–93.
28. Austin PC. The relative ability of different propensity score methods to balance measured covariates between treated and untreated subjects in observational studies. *Med Decis Making* 2009;29:661–77.
29. Lunceford JK, Davidian M. Stratification and weighting via the propensity score in estimation of causal treatment effects: a comparative study. *Stat Med* 2004;23:2937–60.
30. Glynn RJ, Schneeweiss S, Stürmer T. Indications for propensity scores and review of their use in pharmacoepidemiology. *Basic Clin Pharmacol Toxicol* 2006;98:253–9.
31. Stürmer T, Rothman KJ, Glynn RJ. Insights into different results from different causal contrasts in the presence of effect-measure modification. *Pharmacoepidemiol Drug Saf* 2006;15:698–709.

KEY WORDS bias, comparison of methods, observational studies, regression

APPENDIX For supplemental tables and figures, please see the online version of this article.