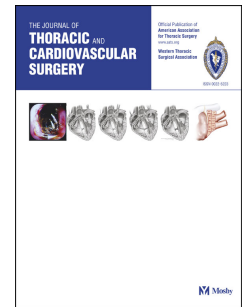# Accepted Manuscript

Propensity Scores: Methods, Considerations, and Applications in the Journal of Thoracic and Cardiovascular Surgery

Timothy L. McMurry, PhD, Yinin Hu, MD, Eugene H. Blackstone, MD, Benjamin D. Kozower, MD, MPH

Please cite this article as: McMurry TL, Hu Y, Blackstone EH, Kozower BD, Propensity Scores: Methods, Considerations, and Applications in the Journal of Thoracic and Cardiovascular Surgery, *The Journal of Thoracic and Cardiovascular Surgery* (2015), doi: 10.1016/j.jtcvs.2015.03.057.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

1 <u>Title</u>: Propensity Scores: Methods, Considerations, and Applications in the Journal of Thoracic and

2 Cardiovascular Surgery

3

4 <u>Date</u>: March 12, 2015

5

6 <u>Authors</u>:

7 Timothy L. McMurry, PhD[1]; Yinin Hu, MD[2]; Eugene H. Blackstone, MD[3]; Benjamin D. Kozower, MD,

8 MPH[1,2]

9

10

11 <u>Affiliations</u>:

12 [1]Department of Public Heath Sciences, University of Virginia, Charlottesville, VA 22908.

13 [2]Department of Surgery/Division of Thoracic Surgery, University of Virginia Health System,

14 Charlottesville, VA 22908.

15 [3]Department of Thoracic and Cardiovascular Surgery and Department of Quantitative Health Sciences,

16 Research Institute, Cleveland Clinic, Cleveland, OH 44195.

17

19

20 <u>Word count</u>: 3570

21

23

25

26 <u>Corresponding author</u>:

27 Benjamin D. Kozower, MD, MPH

28 University of Virginia Health System

29 General Thoracic Surgery

30 P.O. Box 800679

31 Charlottesville, VA 22908-0679

32 E: bdk8g@virginia.edu

33 P: 1(434) 924-2145

34 F: 1(434) 244-9429

35

36    **Structured Abstract**

37    **Objective:** We review the published literature using propensity scoring, describe

38    shortcomings in the use of this technique, and provide conceptual background for

39    understanding and correctly implementing propensity matched studies.

40    **Methods:** We survey the published statistical literature and make recommendations for a set

41    of standard criteria for propensity matching studies. We then applied these criterion to recent

42    publications in the Journal of Thoracic and Cardiovascular Surgery and determined how well

43    the standards were applied.

44    **Results:** We found that propensity matched studies are rarely documented well enough to be

45    convincing in their results.   When documentation is available, statistical shortcomings are

46    common.

47    **Conclusions:** Improved statistical practice is needed when using propensity scoring.  This

48    article creates standard criterion for using this method in JTCVS publications.

49

50    **Central Picture**

51        Timothy L. McMurry, PhD

52

53    **Central Message**

54        We provide conceptual background and good practice recommendations for the use of

55    propensity matching in surgical outcomes research.

56

57    **Clinical Perspective**

58        The frequent use of propensity matching requires that it be understood by clinicians

59    and researchers. Retrospective studies comparing two treatments are hampered by selection

60    bias; propensity scoring reduces the bias by selecting similar groups for comparison. Review

61    of recent JTCVS publications suggests that the methods are not consistently applied. This

62    article highlights the important concepts.

63


# 64    1.  Introduction

65    Propensity scoring is a powerful tool to strengthen causal inferences drawn from observational

66    studies. The motivation is simple: in order to compare the effects of two treatment options, which we

67    generically refer to as "A" and "B" with B being the more common, we want to compare the outcomes of

68  similar groups of patients receiving each treatment. Propensity scoring helps select similar patient groups

69  for comparison.

70      Propensity scoring is common in the literature and the methodology is widely discussed[1-8]. Despite

71  propensity scoring's popularity, we are concerned that its use is conceptually more intricate than many

72  investigators realize. The consequence can results that are misleading or difficult for readers, referees, and

73  investigators to evaluate objectively. These concerns persist despite having been raised previously in the

74  cardiothoracic surgery literature.[9] The problem is compounded because recommendations in the

75  methodological literature are not consistent (see Section 8).

76      In the present article, we review the basics of propensity scoring, highlight areas where practical

77  recommendations are in general agreement or disagreement, discuss choices available to the investigator,

78  examine how this family of techniques has been applied in recent Journal of Thoracic and Cardiovascular

79  Surgery articles, and establish guidelines for future articles. Sections 2–9 provide conceptual background,

80  Section 10 presents general guidelines, and Section 11 examines recent articles published in the Journal of

81  Thoracic and Cardiovascular Surgery (JTCVS).

## 82  2.  The Conceptual Framework for Propensity Matching

83      Many medical studies are designed to evaluate the effectiveness of a treatment in comparison to an

84  alternative. The goal is to estimate how much better (or worse) the outcomes are for patients in one

85  treatment group as compared with what would have happened had they received the other treatment.

86  Importantly, these comparisons need to account for differences in patients that may have contributed to

87  their allocation to one of the treatments.

88      The randomized controlled trial, where patients are randomly assigned to treatment groups, is the gold

89  standard for this comparison. The importance of random assignment is that the patient characteristics

90  affecting outcomes (e.g., age, sex, comorbidities, mental disposition) tend to be equally distributed across

91  groups. Thus, significant differences in outcomes can be attributed to the treatment.

92      However, there are many questions that a randomized trial cannot address due to cost, ethical or

93  practical considerations, or timeliness. For example, in a study of smoking related healthcare costs,

94  patients cannot ethically (or practically) be randomly assigned to smoke or not smoke. [6] Moreover, a

95  randomized study would take decades to run. Other variables of interest cannot be assigned; Koch et al.

96  considered whether men and women fare differently after coronary artery bypass grafting. [10] Many of

97  these clinical questions can be addressed with robust observational data.

98      Unfortunately, in observational data, patients receiving treatment A typically differ systematically

99  from those receiving B, leaving direct comparisons of outcomes heavily confounded. For example,

100  patients who receive surgery are deemed well enough before treatment to survive the surgery, while

101    extremely frail patients may be deemed inoperable. It is not fair or appropriate to compare outcomes

102    between such dissimilar patients.

103        Historically, investigators tried to account for differences between groups by using multiple regression

104    to adjust for confounding characteristics. However, because the groups of patients may differ

105    systematically, using a regression to estimate the potential effect of treatment A on a patient who received

106    B can be an unreliable extrapolation. Therefore, it is imperative to only compare patients who are

107    legitimate candidates for either procedure.

108        An intuitively appealing approach would be to match patients from one treatment group with patients

109    from the other on a number of important characteristics, and then compare their outcomes. Unfortunately,

110    matching on even a modest number of criteria often leaves a large majority of patients unmatched and

111    unavailable for analysis, making the results less reliable. [11]

112        Propensity scores solve the problem of matching on multiple covariates by reducing them to a single

113    quantity, the propensity score. A patient's propensity score is defined as the probability the patient

114    receives treatment A (instead of B) given all relevant conditions, comorbidities, and other characteristics

115    at the time the treatment decision is made. What makes propensity scores so powerful is that, under some

116    conditions, patients with the same propensity score have the same probabilistic distribution of other

117    covariates regardless of whether they received A or B. [12] As a result, it can be sufficient to compare the

118    outcomes across treatment groups of pairs or pools of patients with similar propensity scores.

## 3. Steps in a propensity score analysis

120        There are four main steps in an analysis using propensity scores. First, the propensity scores must be

121    estimated (Section 5). Second, the data need to be matched or grouped based on the estimated propensity

122    scores (Section 6). Third, balance must assessed to ensure that the grouping produced similar pools of

123    patients receiving treatments A and B (Section 7). Finally, data can be analyzed to estimate the treatment

124    effect size and its clinical and statistical significance (Section 8). The first three of these steps are "design"

125    steps used to frame a comparison around similar groups of patients; they must be performed without

126    looking at the outcomes data. None of the steps can be adequately performed by following a simple

127    recipe. Importantly, even before tackling these technical issues, two crucial assumptions must be met in

128    order for propensity matching to provide useful results.

## 4. Crucial assumptions

130        There are two conditions on the data which must be met for analyses based on propensity scores to

131    provide valid results. The most important condition is known as "strong ignorability," which technically

132 says a patient's treatment assignment (A or B) is independent of his/her potential outcomes under the two

133 treatment scenarios given the covariates. In other words, the observed covariates contain all the

134 information about the patient's condition that is relevant to the patient's potential outcomes. Strong

135 ignorability makes intuitive sense. If the goal is to compare similar groups of patients receiving different

136 treatments, we need to know all the factors that determine whether patients are comparable at the time of

137 treatment allocation.

138 The second important condition states that, given the covariates, the patient needs to have a positive

139 probability of receiving either of the treatments. Intuitively, there is no gain in asking what the potential

140 benefits of surgery are for a patient whose comorbidities preclude surviving an operation. The only

141 interest is in comparing patients for whom both treatments are realistic.

# 5. Constructing the propensity score model

143 The first step is to estimate the propensity scores for each patient. The most common approach is to use

144 logistic regression, but there are other regression models that can estimate classification probabilities.[13,14]

## Which variables to include

146 Guidelines for constructing and evaluating a regression model depend on its intended application.[15]

147 For propensity scores, the "strong ignorability" condition necessitates inclusion of covariates that predict

148 potential outcomes under either treatment scenario as well as any covariates that predict treatment

149 assignment, although these two criteria are typically related. From a practical point of view, the second of

150 these requirements deserves particular emphasis: if the data do not contain the information used to make

151 the treatment decision (or are systematically missing in one of the two patient groups), the propensity

152 model will be inadequate, and all subsequent analyses will be suspect.

153 The consensus is that if the sample size is too small for the propensity score model to include all

154 variables of interest, it is most important to include the variables that are strongly related to outcome.[16]

155 These should be selected *a priori* based on scientific understanding and previous literature, and without

156 reference to the outcomes within the data set.[17] It is probably better to err on the side of too many

157 predictors rather than too few,[7] and when data sizes are large, good propensity models can contain many

158 predictors.[6]

## Logistic regression model diagnostics

160 The goal of the propensity score model is to create balanced groups of patients receiving each

161 treatment. Therefore, some model evaluation tools, such as those evaluating discriminative ability (e.g.,

162 the *c*-statistic), multicollinearity, and model selection are only of secondary importance.[18] The crucial

163 diagnostic step is to compare covariate balance between the resulting groups of patients receiving the two

164 treatments; see Section 7.

165 A model which accurately estimates the likelihood of treatment allocation is the key to achieving this

166 balance. Nonetheless, some metrics which are common in many regression applications are of diminished

167 importance in the present context. For example, multicollinearity occurs when highly correlated predictors

168 produce instability in their corresponding coefficients. Fortunately, multicollinearity does not affect the

169 resulting fitted values, in this case, the propensity scores. However, if the sample size is limited, it may

170 still be advantageous to remove highly correlated variables in order to include less correlated covariates.

171 There is also concern about traditional model selection strategies such as stepwise variable selection.

172 These approaches are designed for prediction rather than covariate balance. The concern is that these

173 selection methods might remove variables that are weakly related to treatment assignment but strongly

174 related to outcome[16], while variables related to outcome are thought to be at least as important.

175 The commonly used $c$-statistic also requires nuanced interpretation in this setting. In most

176 applications, a predictive model with a low $c$-statistic is useless. A propensity model with a low $c$-statistic

177 could be caused by poor construction. However, it could also be indicative of differences in practice that

178 are not related to patient condition. The first of these problems invalidates subsequent analyses, while the

179 latter can be beneficial. To illustrate: imagine trying to estimate propensity scores for a randomized trial.

180 A well constructed model accounting for all relevant clinical covariates will have a $c$-statistic around 0.5,

181 and all patients should have similar propensity scores. Nonetheless, a randomized trial is ideally suited for

182 causal inference. At the other extreme, a $c$-statistic close to 1 indicates that the regression model is able to

183 differentiate patients receiving A from those receiving B. This may be an indication that the these two

184 groups are so different that their outcomes will be difficult to meaningfully compare.

# 6. Grouping the data

186 Once propensity scores have been estimated, the data are typically grouped by either subclassification

187 (sometimes called stratification) or matching. Both of these methods prune the original data set down to

188 groups or sets of patients with similar propensity scores. While there are other approaches, we focus on

189 these two for their simplicity and frequency of use.

190 It may be reasonable to remove patients receiving one treatment who have propensity scores either

191 much larger or much smaller than any patient receiving the other treatment, the "oranges" as discussed in

192 Blackstone.[2] The rationale for exclusion is that these patients do not appear to have been candidates for

193 the alternative treatment. Nonetheless, excluded patients should be examined carefully. If there are many

194 unmatchable patients, the propensity score model may include a variable that is a strong surrogate for

195  treatment assignment, which may be removed. [2] Evaluation of the "oranges" will also help reveal the

196  limits within which a valid comparison of the two treatments is possible.

### Subclassification

198      Subclassification is frequently suggested in the methodological literature but less frequently

199  applied. The idea is simple: propensity scores are grouped into (e.g.) quintiles or deciles (5 to 10 groups is

200  typical). [19] Within each group, the propensity scores are similar, so grouped patients should have similar

201  covariate distributions, and thus can be compared to each other. An analysis is performed in each group,

202  and then results are aggregated. Subclassification has intuitive appeal because it focuses comparisons on

203  pools of patients with similar propensity scores. In contrast, if patients are matched, there may be many

204  suitable matches for a group A patient, with some potential matches arbitrarily excluded from final

205  comparisons.

### Matching

207      The more common approach is to match individual patients receiving one treatment to patients with

208  similar propensity scores receiving the other. While conceptually simple, the details lead to different

209  algorithms which can affect subsequent analyses. These variations include the methods for measuring

210  distances between propensity scores, the threshold for what constitutes matching scores, how one match is

211  chosen from many candidates, the number of patients in group B (the larger group) matched to each

212  patient in A, and whether or not a single patient in group B can be matched to more than one individual in

213  group A.

214      Intuition suggests that the distance between propensity scores should be measured by the simple

215  difference between estimated probabilities of treatment. This approach is commonly used, however there

216  is evidence that it is more effective to match on the "linear propensity score," or the difference between

217  propensity scores on the logit scale.[20]

218      One needs to decide how close two propensity scores need to be before they can be potential matches;

219  this threshold is known as a "caliper." A narrow caliper can prevent inaccurate matching, but if too many

220  patients go unmatched the results can become uninterpretable. The appropriate caliper size depends on the

221  relative variations of propensity scores in the two treatment groups.[20,21]

222      Next, the user must decide how many patients in group B should be matched to each patient in A. The

223  most common approach is to match each patient in A to a single patient in B. If group B is much larger

224  than A, it may be advantageous to match a larger number of patients in B, but the benefits are reduced if

225  the extra matches are of poor quality. Finally, it is possible to re-use patients in group B. This makes the

226  matching process independent of the order in which the matches are selected and may improve the overall

7

227  match quality. However, without adjustment, reused patients have too much weight in the final
228  comparison of outcomes.

229      Once these decisions have been made, pairing is often done by a "greedy" algorithm. The group A
230  patients are randomly ordered. The first of these randomly ordered patients is then matched to their best
231  group B counterpart. The group B patient is removed from the set of potential future matches, and the
232  process is iterated. An alternative to "greedy" matching is "optimal" matching, [4] which seems to produce
233  better matched pairs, but does not substantially improve the balance of the matched groups as a whole. [22]
234  Stuart maintains a web page describing available software. [23]

## 7.  Assessment of covariate balance in matched groups

236      Since the goal in using a propensity scores is to create pools of similar patients for comparison, it is
237  extremely important to assess the post-matching similarity across groups before any assessment of
238  outcomes. In particular, all covariates affecting patients' prognoses prior to treatment and indications for
239  treatment need to be compared. If clinically-relevant differences remain post-match, subsequent analyses
240  are unreliable.

241      The types of comparisons differ depending on how the data are grouped. If the data are subclassified,
242  then one should perform diagnostics within each subclass. If the data are matched, then typically
243  comparison is between the matched pools of patients.

244      Most investigators assess covariate balance using hypothesis tests. For example, an investigator might
245  test the hypothesis that the average age of patients in group A is the same as their matched counterparts in
246  B. Unfortunately, hypothesis tests answer the wrong question. The *p*-value from a hypothesis test depends
247  on the difference between the two groups and their sample sizes. However, only the difference between
248  the two groups is relevant to covariate balance. [24]

249      A better metric for continuous covariates would be to use a measure that does not depend on sample
250  size, such as the standardized difference in means: $(\bar{X}_A - \bar{X}_B)/\sigma_A$ [20], which expresses the difference
251  between the two groups in standard deviations. The improvement in balance achieved by matching can be
252  demonstrated by comparing standardized differences in means before and after matching (using the same
253  estimate for $\sigma_A$ in both quantities). Binary covariates can be compared with a simple difference in
254  proportions or by a similar standardized difference. [25] Alternatively, Rubin suggests a set of powerful but
255  less intuitive diagnostics. [17]

## 8. Analysis of the matched data

The final step in a propensity score analysis is to estimate the treatment effect size and its clinical and statistical significance. Literature on the proper analysis of matched data is sparse and occasionally in conflict. For example, Rosenbaum recommends analyzing the data with permutation tests in the same way one would analyze an unmatched observational trial. [4] Austin argues that propensity matched data should be analyzed using procedures for matched analyses, such as paired *t*-tests, and McNemar's test. [26] Stuart replies that matched analyses are not necessary, and that the data can be analyzed using a standard regression that includes a treatment indicator and the variables used in the matching. [27] A recent article by Li and Greene suggests that a weighting method is optimal. [28] Many articles make almost no mention of statistical inference.

This confusion has resulted because statistical understanding is still evolving, and assumptions made about the data and matching process can alter the estimates' derived properties. In most cases relevant to surgical outcomes, regression is defensible and even recommendable. Propensity scores provide an objective way to restrict the domain of analysis to patients who are legitimate candidates for either procedure. Outcomes in the two groups are then compared using a regression model that controls for all covariates used in matching, plus a treatment indicator variable. The coefficient associated with this indicator is interpreted as the treatment effect. An advantage of regression is that it provides some level of "double robustness" by adjusting for any remaining small covariate imbalances. [29] It is for this reason that even randomized trials are sometimes analyzed with regression models.

Regression is more important following subclassification because, within subclasses, meaningful covariate imbalances may still remain. In this setting, recommendations for the regression remain similar. If enough data are available, a regression model containing the treatment indicator and all covariates can be fit in each subclass, and the results combined. If data are more limited, a single regression model may be fit containing subclass indicators and subclass by treatment interactions along with the other covariates. This keeps the covariate relationships fixed but allows different size treatment effects across the subclasses. After regression modeling, subclass-specific treatment effects are then combined by a weighted average of the treatment effects in each subclass, where the effects are typically weighted by the number of group A individuals in the subclass.

## 9. Interpretation

For matched data, patients receiving A have been grouped with a probabilistically similar pool of group B patients. Therefore, the estimated effect size represents the average improvement of the group A patients relative to similar patients in group B. This quantity is traditionally described in the literature as

9

288 the Average Treatment Effect in the Treated (ATT). ATT is not the same as the average effect of
289 treatment across the entire population, referred to in the literature as Average Treatment Effect (ATE). In
290 most cases, we suspect ATT is the desired quantity, as it describes the benefits/risks of A relative to
291 similar patients receiving B, rather than as a potential benefit averaged across all patients.

292    Both ATT and ATE assume that all group A patients in the initial data set were included in the final
293 analyzed groups. If many patients have been excluded, the interpretation may change or results may
294 become uninterpretable; see Section 6.

## 10.  Recommendations for published literature

296    While we recognize the importance of brevity, it is important that propensity scoring methods be
297 described well enough that results can be evaluated and replicated. Most of our recommendations can be
298 implemented with one or two paragraphs. In some cases, additional tables may be provided in on-line
299 appendices.

300    While different analyses are appropriate for different data sets and clinical questions, we propose that
301 articles utilizing propensity matching should include the following:

302    1. The original sample sizes for the pools of patients in each group.

303    2. The sample sizes available after matching.

304    3. The type of regression model used to estimate the propensity scores.

305    4. The variables considered for inclusion in the propensity model, the variables included in the final
306       model, and the inclusion criteria.

307    5. The type of matching algorithm used.

308    6. Diagnostics demonstrating the quality of the resulting matches.

309    7. Characterization of the unmatched patients.

310    8. An indication of the statistical procedures used for analyses.

## 11.  JTCVS literature review

312    We reviewed all publications in JTCVS from 2013 and 2014 using propensity score matching.
313 We found 25 such articles in 2013 and 64 in 2014.  While many of these papers were well done, some
314 exhibited substantial statistical shortcomings, and many did not provide enough detail for objective
315 evaluation. Our results are summarized in Table 1. Notably, many papers showed evidence of inadequate
316 covariate balance after matching, and no paper carefully evaluated the excluded patients.

## 12. Conclusions

Propensity matching is a powerful tool for observational data analyses because it facilitates the comparison of outcomes between similar groups of patients. Although propensity matching has become a popular technique, the methodology is acqtually quite complex. This review is intended to help surgeons understand the concepts behind propensity matching which may influence their own research and/or help them critically evaluate the published literature. We identified eight criteria which we feel should be reported in any manuscript using propensity matching. When we applied these criteria to the publications in JTCVS from 2013 and 2014, concerns were raised about the use of this methodology and appropriateness of the applications. We recommend that the Journal adopt these criteria to create a standard for future articles submitted to JTCVS using propensity matching.

## References

1. Blackstone EH. Breaking down barriers: Helpful breakthrough statistical methods you need to understand better. *J Thorac Cardiovasc Surg*. 2001;122(3):430-439.

2. Blackstone EH. Comparing apples and oranges. *J Thorac Cardiovasc Surg*. 2002;123(1):8-15.

3. Brookhart MA, Wyss R, Layton JB, Sturmer T. Propensity score methods for confounding control in nonexperimental research. *Circ Cardiovasc Qual Outcomes*. 2013;6(5):604-611. doi: 10.1161/CIRCOUTCOMES.113.000359 [doi].

4. Rosenbaum PR. *Observational studies.* 2nd ed. New York: Springer; 2002.

5. Rubin D. *Matched sampling for casual effects.* Vol 10. New York: Cambridge University Press; 2006:12.

6. Rubin DB. The design versus the analysis of observational studies for causal effects: Parallels with the design of randomized trials. *Stat Med*. 2007;26(1):20-36.

7. Stuart EA. Matching methods for causal inference: A review and a look forward. *Stat Sci*. 2010;25(1):1-21.

8. d'Agostino RB. Tutorial in biostatistics: Propensity score methods for bias reduction in the comparison of a treatment to a non-randomized control group. *Stat Med*. 1998;17(19):2265-2281.

341   9. Austin PC. Propensity-score matching in the cardiovascular surgery literature from 2004 to 2006: A systematic

342   review and suggestions for improvement. *J Thorac Cardiovasc Surg*. 2007;134(5):1128-1135. e3.

343   10. Koch CG, Khandwala F, Nussmeier N, Blackstone EH. Gender and outcomes after coronary artery bypass

344   grafting: A propensity-matched comparison. *J Thorac Cardiovasc Surg*. 2003;126(6):2032-2043.

345   11. Rosenbaum PR, Rubin DB. The bias due to incomplete matching. *Biometrics*. 1985:103-116.

346   12. Rosenbaum PR, Rubin DB. The central role of the propensity score in observational studies for causal effects.

347   *Biometrika*. 1983;70(1):41-55.

348   13. Lee BK, Lessler J, Stuart EA. Improving propensity score weighting using machine learning. *Stat Med*.

349   2010;29(3):337-346.

350   14. Westreich D, Lessler J, Funk MJ. Propensity score estimation: Neural networks, support vector machines,

351   decision trees (CART), and meta-classifiers as alternatives to logistic regression. *J Clin Epidemiol*. 2010;63(8):826-

352   833.

353   15. Vittinghoff E, Shiboski S, McCulloch CE. *Regression methods in biostatistics*. 2nd ed. New York: Springer;

354   2010.

355   16. Brookhart MA, Schneeweiss S, Rothman KJ, Glynn RJ, Avorn J, Sturmer T. Variable selection for propensity

356   score models. *Am J Epidemiol*. 2006;163(12):1149-1156.

357   17. Rubin DB. Using propensity scores to help design observational studies: Application to the tobacco litigation.

358   *Health Serv Outcomes Res*. 2001;2(3-4):169-188.

359   18. Rubin DB. On principles for modeling propensity scores in medical research. *Pharmacoepidemiol Drug Saf*.

360   2004;13(12):855-857.

361   19. Rosenbaum PR, Rubin DB. Reducing bias in observational studies using subclassification on the propensity

362   score. *Journal of the American Statistical Association*. 1984;79(387):516-524.

363     20. Rosenbaum PR, Rubin DB. Constructing a control group using multivariate matched sampling methods that

364     incorporate the propensity score. *The American Statistician*. 1985;39(1):33-38.

365     21. Cochran WG, Rubin DB. Controlling bias in observational studies: A review. *Sankhyā: The Indian Journal of*

366     *Statistics, Series A*. 1973:417-446.

367     22. Gu XS, Rosenbaum PR. Comparison of multivariate matching methods: Structures, distances, and algorithms.

368     *Journal of Computational and Graphical Statistics*. 1993;2(4):405-420.

369     23. Stuart E. Software for implementing matching methods and propensity scores. Elizabeth Stuart's Propensity

370     Score Software Page Web site. http://www.biostat.jhsph.edu/~estuart/propensityscoresoftware.html. Updated 2012.

371     Accessed 9/5, 2014.

372     24. Imai K, King G, Stuart EA. Misunderstandings between experimentalists and observationalists about causal

373     inference. *Journal of the Royal Statistical Society: Series A*. 2008;171(2):481-502.

374     25. Austin PC. Using the standardized difference to compare the prevalence of a binary variable between two groups

375     in observational research. *Communications in Statistics-Simulation and Computation*. 2009;38(6):1228-1234.

376     26. Austin PC. A critical appraisal of propensity-score matching in the medical literature between 1996 and 2003.

377     *Stat Med*. 2008;27(12):2037-2049.

378     27. Stuart EA. Developing practical recommendations for the use of propensity scores: Discussion of 'A critical

379     appraisal of propensity score matching in the medical literature between 1996 and 2003' by peter austin, statistics in

380     medicine. *Stat Med*. 2008;27(12):2062-2065.

381     28. Li L, Greene T. A weighting analogue to pair matching in propensity score analysis. *The International Journal of*

382     *Biostatistics*. 2013;9(2):215-234.

383     29. Rubin DB. The use of matched sampling and regression adjustment to remove bias in observational studies.

384     *Biometrics*. 1973:185-203.

13

387

388    Table 1: Characteristics of recent (2013-2014) JTCVS papers using propensity scores.

| Criteria | Number (%) of Papers Providing Information |
|---|---|
| Sample size for original data set | 87/89 (98%) |
| Matched sample size | 81/89 (91%) |
| Type of regression model used to estimate the propensity score | 79/89 (89%) |
| Matching algorithm | 60/89 (67%) |
| Analysis of covariate balance | 66/89 (74%) |
|     Evidence of inadequate covariate balance | 17/66 (26%) |
| Comparison of matched to unmatched patients | 0/89 (0%) |
| Type of statistical procedure | Number of papers |
|     Univariate, independent samples | 59/89 (66%) |
|     Univariate, paired | 11/89 (12%) |
|     Regression after matching | 31/89 (35%) |
|     Regression including the propensity score as a covariate | 10/89 (11%) |

389