

# Implementing CTA from Within Stata: Using CTA to Generate Propensity Score Weights (*Invited*)

Ariel Linden, Dr.P.H.  
Linden Consulting Group, LLC

In contrast to randomized studies in which individuals have no control over their treatment assignment, participants in observational studies self-select into the treatment arm and are therefore likely to differ in their characteristics compared to those who elect not to participate. Analytic approaches using the propensity score to adjust for differences between study groups are thus popular among investigators of observational data. In this paper, I describe how the new Stata package for implementing CTA can be used to generate propensity score weights.

Prior papers<sup>1,2</sup> introduced the new Stata package called **cta**<sup>3</sup> for implementing CTA from within the Stata environment. This package is a wrapper for the CTA software<sup>4</sup>, thus the CTA64.exe file must be loaded on the computer for the **cta** package to work (CTA software is available at <https://odajournal.com/resources/>). To download the **cta** package, at the Stata command line type: “ssc install cta” (without the quotation marks).

This paper demonstrates how the **cta** package can be used in observational studies to compute propensity score<sup>5-9</sup> weights that will thereafter be used in the outcomes (treatment effects) model.

Computing propensity score weights in **cta** is a two-step process. First a CTA model is generated by specifying the treatment assign-

ment indicator as the *class* variable and all the observed pre-intervention covariates as *attributes*. Next, weights are computed separately for each strata (i.e., endpoint) in the model.

## Methods

### Data

This paper uses data from a prior evaluation of a health plan-based program intended to reduce 30-day readmission rates for patients hospitalized with one or more chronic illnesses. The intervention was modeled after that described in Linden and Butterworth,<sup>10</sup> which focused on behavioral change to help patients actively

engage in their own health care, which in turn was expected to reduce the likelihood of readmission.<sup>11-15</sup> This subset of the retrospectively collected data consists of observations for 1398 participants and 7957 nonparticipants.

Ten pre-intervention characteristics available for every observation included demographic variables ([Age] and [Gender]), health services use in the 12 months prior to the index hospitalization (office visits [Office], emergency department [ED] visits, hospitalizations [Admits]), length of stay for the index [Index] hospitalization, indicator variables for whether the patient had congestive heart failure [CHF] and/or chronic obstructive pulmonary disease [COPD], the patient's Charlson comorbidity index score [CCI],<sup>16</sup> and a diagnosis-based risk adjustment score [Riskscore]. The outcome was the number of days post-discharge from the index hospitalization: patients were classified as censored if they were lost to follow-up prior to 30 days, or if they did not experience a readmission within 30 days.<sup>17</sup>

Analytic process

Generating a CTA model

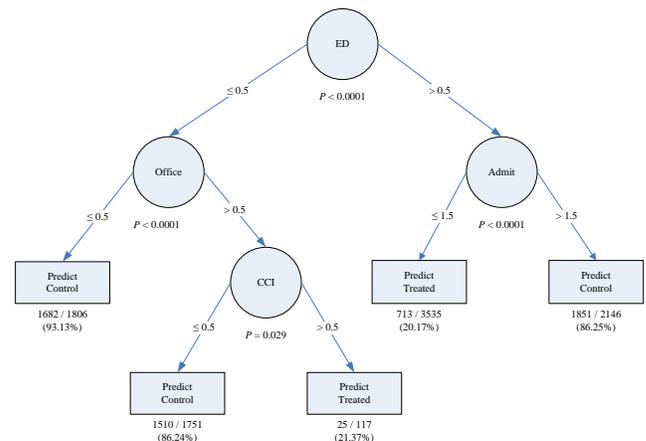
For evaluating self-selection in an observational study, **cta** is implemented with the following syntax (see the help file for **cta** for a complete description of the syntax options):

```
cta treat age gender riskscore admits er
office index cci chf copd , pathcta("C:\
CTA\") store("C:\CTA") cat( gender chf
copd) iter(10000) loo(stable) prune(0.05)
enumerate
```

The above syntax is explained as follows: The outcome variable “treat” is the *class* variable; the 10 variables listed until the comma are covariates specified as the *attributes*; the directory path where the CTA64.exe file is located on my computer is “C:\CTA\”; the directory path where the output and other files

generated during the analysis should be stored is “C:\CTA\output”; the *cat()* option indicates which attributes are categorical; the number of iterations (repetitions) for computing a permutation *P*-value is 10,000; leave-one-out analysis is used and attributes are only retained if they are stable; the tree is pruned with a *P*-value of 0.05 used as the cutpoint for inclusion; and an enumerated model (which enumerates the first three nodes) is conducted. (Yarnold and Soltysik<sup>4</sup> provide a complete description of the CTA modeling process and interpretation of results).

The **cta** package produces an extract of the total output produced by CTA software (the complete output is stored in the specified directory with the extension “.out”). Here we include a diagram of the pruned model, which achieved overall ESS of 16.17 (a weak effect)—slightly less than the enumerated model (ESS=18.35), but much more parsimonious (5 vs. 31 nodes).



In reviewing this diagram, it is evident that those patients predicted to participate in the intervention follow a different pathway than those predicted to serve as controls. That is, a patient is predicted to participate in the intervention if (1) they had less than one ED visit, more than one office visit, and a CCI level of at least 1 in the past year, or if (2) they had at least one ED visit and less than two hospital admissions in the prior year. While the accuracy of these predictions is low (21.37% and 20.17%

for pathway 1 and 2, respectively), these pathways were statistically significant ( $P < 0.0001$ ).

### Computing Propensity Score Weights

For CTA models, a stratified weight is generated for each individual based on both their actual treatment assignment and their specific stratum (model endpoint): observations have identical weights if they are classified into the same endpoint and they have the same actual treatment assignment (i.e., treated or non-treated). CTA model-based stratified weights are computed using the following formula:

$$\frac{n_s \times \Pr(Z = z)}{n_{z = z, s}}$$

where  $n_s$  is the total number of individuals in a given stratum  $s$ ,  $\Pr(Z = z)$  is the estimated probability of assignment to treatment group  $z$  (i.e., the proportion of individuals actually receiving treatment  $z$  in the sample), and  $n_{z = z, s}$  is the total number of individuals in stratum  $s$  who were actually assigned to treatment  $z$ . Thus, the weight is proportional to the ratio of the number of individuals in a given stratum relative to the number of individuals within that stratum who do (not) receive treatment. Taken together, the stratification reduces bias in the observed covariates used to create the propensity score, and the weighting standardizes each treatment group to the target population. We developed this stratified weighting approach for the CTA models to ensure that weights conform exactly to the underlying geometry and findings of the CTA model.<sup>18,19</sup>

To demonstrate the computation of the stratified weight we take the first stratum (far left end-point in the diagram) as an example. For controls, the weight is:

$$\frac{1806 \times 0.8506}{1682} = 0.9133$$

where 1806 is the total number of individuals in that stratum, 0.8506 is the proportion of controls in the entire sample, and 1682 is the number of controls in that stratum. Thus, every control in that stratum will be assigned a weight of 0.9133. More specifically, the following Stata code would generate a variable named *cta1* to represent all controls in the pathway to the first end-point:

```
gen cta1 = 0.9133 if er <= 0.5 & office <= 0.5 & treat == 0
```

For treated individuals in that stratum, the computation is:

$$\frac{1806 \times 0.1494}{124} = 2.1759$$

where 1806 is the total number of individuals in that stratum, 0.1494 is the proportion of treated individuals in the entire sample, and 124 is the number of treated individuals in that stratum. Thus, every treated individual in that stratum will be assigned a weight of 2.1759. We can then use the following Stata code to replace values of *cta1* to represent treated patients in the pathway to the first end-point:

```
replace cta1 = 2.1759 if er <= 0.5 & office <= 0.5 & treat == 1
```

Treated individuals are weighted substantially more than controls in this stratum because the end-point predicts controls. Thus, treated individuals who are in this stratum are wrongly predicted to be there and therefore must be weighted more in order to achieve balance with controls in the covariate pathway.

Upon computing the weights for all strata and assigning them to each individual in the sample, a treatment effects analysis can be conducted in which these weights are specified in the *wt()* option (a detailed discussion on how to estimate treatment effects for observational data using **oda** is presented elsewhere<sup>20,21</sup>).

## Discussion

This paper demonstrates how to compute propensity score weights after generating a CTA model using the new Stata package `cta`. CTA provides accurate, parsimonious decision rules that are easy to visually display and interpret, while reporting  $P$  values derived via permutation tests at every node, in addition to corresponding partial ESS statistics. CTA is also insensitive to skewed data or outliers, and has the ability to handle any variable metric including categorical, Likert-type integer, and real number measurement scales. Moreover, CTA also has the distinct ability to ascertain where optimal (maximum-accuracy) cutpoints are on each variable, which in turn, facilitates the use of measures of predictive accuracy. Moreover, CTA can perform cross-validation using LOO which allows for assessing the cross-generalizability of the model to potentially new study participants or non-participants.<sup>22</sup>

Finally, the findings continue to support our recommendation to employ the ODA and CTA frameworks to evaluate the efficacy of health-improvement interventions and policy initiatives.<sup>23-38</sup>

## References

- <sup>1</sup>Linden A (2020). Implementing CTA from Within Stata: Implementing CTA from Within Stata: Assessing the Quality of the Randomization Process in Randomized Controlled Trials (*Invited*). *Optimal Data Analysis*, 9, 57-62.
- <sup>2</sup>Linden A (2020). Implementing CTA from Within Stata: Characterizing Participation in Observational Studies (*Invited*). *Optimal Data Analysis* 2020;9:63-67.
- <sup>3</sup>Linden A. (2020). CTA: Stata module for conducting Classification Tree Analysis. *Statistical Software Components S458729*, Boston College Department of Economics.
- <sup>4</sup>Yarnold PR, Soltysik RC (2016). *Maximizing Predictive Accuracy*. Chicago, IL: ODA Books. DOI: 10.13140/RG.2.1.1368.3286
- <sup>5</sup>Linden A, Adams J (2006). Evaluating disease management program effectiveness: an introduction to instrumental variables. *Journal of Evaluation in Clinical Practice*, 12, 148-154.
- <sup>6</sup>Linden A, Adams JL (2010). Using propensity score-based weighting in the evaluation of health management programme effectiveness. *Journal of Evaluation in Clinical Practice*, 16, 175-179.
- <sup>7</sup>Linden A, Adams JL (2010). Evaluating health management programmes over time. Application of propensity score-based weighting to longitudinal data. *Journal of Evaluation in Clinical Practice*, 16, 180-185.
- <sup>8</sup>Linden A (2014). Combining propensity score-based stratification and weighting to improve causal inference in the evaluation of health care interventions. *Journal of Evaluation in Clinical Practice*, 20, 1065-1071.
- <sup>9</sup>Linden A, Uysal SD, Ryan A, Adams JL (2016). Estimating causal effects for multivalued treatments: A comparison of approaches. *Statistics in Medicine*, 35, 534-552.
- <sup>10</sup>Linden A, Butterworth SW (2014). A comprehensive hospital-based Intervention to reduce readmissions for chronically ill patients: A randomized controlled trial. *American Journal of Managed Care*, 20, 783-792.
- <sup>11</sup>Linden A, Roberts N (2004). Disease management interventions: What's in the black box? *Disease Management*, 7, 275-291.
- <sup>12</sup>Linden A, Butterworth S, Roberts N (2006). Disease management interventions II: what else is in the black box? *Disease Management*, 9, 73-85.

- <sup>13</sup>Biuso TJ, Butterworth S, Linden A (2007). Targeting prediabetes with lifestyle, clinical and behavioral management interventions. *Disease Management*, 7, 6-15.
- <sup>14</sup>Linden A, Adler-Milstein J (2008). Medicare disease management in policy context. *Health Care Finance Review*, 29, 1-11.
- <sup>15</sup>Linden A, Roberts N (2005). A Users guide to the disease management literature: recommendations for reporting and assessing program outcomes. *American Journal of Managed Care*, 11, 81-90.
- <sup>16</sup>Charlson ME, Pompei P, Ales KL, MacKenzie CR (1987). A new method of classifying prognostic comorbidity in longitudinal studies: development and validation. *Journal of Chronic Disease*, 40, 373-383.
- <sup>17</sup>Linden A, Adams J, Roberts N (2004). Evaluating disease management program effectiveness: an introduction to survival analysis. *Disease Management*, 7, 180-190.
- <sup>18</sup>Linden A, Yarnold PR (2017). Using classification tree analysis to generate propensity score weights. *Journal of Evaluation in Clinical Practice*, 23, 703-712.
- <sup>19</sup>Linden A, Yarnold PR (2018). Estimating causal effects for survival (time-to-event) outcomes by combining classification tree analysis and propensity score weighting. *Journal of Evaluation in Clinical Practice*, 24, 380-387.
- <sup>20</sup>Linden A (2020). Implementing ODA from Within Stata: Implementing ODA from Within Stata: An Application to Estimating Treatment Effects using Observational Data (*Invited*). *Optimal Data Analysis*, 9, 14-20.
- <sup>21</sup>Linden A (2020). Implementing ODA from Within Stata: Evaluating Treatment Effects for Survival (Time-to-Event) Outcomes (*Invited*). *Optimal Data Analysis*, 9, 39-44.
- <sup>22</sup>Linden A, Adams J, Roberts N (2004). The generalizability of disease management program results: getting from here to there. *Managed Care Interface*, 17, 38-45.
- <sup>23</sup>Linden A, Adams J, Roberts N (October, 2003). *Evaluation methods in disease management: determining program effectiveness*. Position Paper for the Disease Management Association of America (DMAA).
- <sup>24</sup>Linden A, Yarnold PR, Nallomothu BK (2016). Using machine learning to model dose-response relationships. *Journal of Evaluation in Clinical Practice*, 22, 860-867.
- <sup>25</sup>Yarnold PR, Linden A (2016). Theoretical aspects of the D statistic. *Optimal Data Analysis*, 22, 171-174.
- <sup>26</sup>Linden A, Yarnold PR (2017). Modeling time-to-event (survival) data using classification tree analysis. *Journal of Evaluation in Clinical Practice*, 23, 1299-1308.
- <sup>27</sup>Linden A, Yarnold PR (2018). Identifying causal mechanisms in health care interventions using classification tree analysis. *Journal of Evaluation in Clinical Practice*, 24, 353-361.
- <sup>28</sup>Linden A, Yarnold PR (2016). Using machine learning to assess covariate balance in matching studies. *Journal of Evaluation in Clinical Practice*, 22, 848-854.
- <sup>29</sup>Linden A, Yarnold PR (2016). Combining machine learning and propensity score weighting to estimate causal effects in multivalued treatments. *Journal of Evaluation in Clinical Practice*, 22, 875-885.
- <sup>30</sup>Linden A, Yarnold PR (2016). Using machine learning to identify structural breaks in single-group interrupted time series designs. *Journal of Evaluation in Clinical Practice*, 22, 855-859.

<sup>31</sup>Linden A, Yarnold PR (2016). Combining machine learning and matching techniques to improve causal inference in program evaluation. *Journal of Evaluation in Clinical Practice*, 22, 868-874.

<sup>32</sup>Linden A, Yarnold PR (2017). Minimizing imbalances on patient characteristics between treatment groups in randomized trials using classification tree analysis. *Journal of Evaluation in Clinical Practice*, 23, 1309-1315.

<sup>33</sup>Linden A, Yarnold PR (2018). Using machine learning to evaluate treatment effects in multiple-group interrupted time series analysis. *Journal of Evaluation in Clinical Practice*, 24, 740-744.

<sup>34</sup>Yarnold PR, Linden A (2017). Computing propensity score weights for CTA models involving perfectly predicted endpoints. *Optimal Data Analysis*, 6, 43-46.

<sup>35</sup>Yarnold PR, Linden A (2016). Using machine learning to model dose-response relationships via ODA: Eliminating response variable baseline variation by ipsative standardization. *Optimal Data Analysis*, 5, 41-52.

<sup>36</sup>Linden A, Yarnold PR (2018). The Australian gun buy-back program and the rate of suicide by firearm. *Optimal Data Analysis*, 7, 28-35.

<sup>37</sup>Linden A, Yarnold PR (2018). Using ODA in the evaluation of randomized controlled trials. *Optimal Data Analysis*, 7, 46-49.

<sup>38</sup>Linden A, Yarnold PR (2018). Using ODA in the evaluation of randomized controlled trials: Application to survival outcomes. *Optimal Data Analysis*, 7, 50-53.

### **Author Notes**

No conflict of interest was reported.