

# Toolkit for Weighting and Analysis of Nonequivalent Groups: A guide to the `twang` package

Greg Ridgeway, Dan McCaffrey, Andrew Morral, Matthew Cefalu, Lane Burgette, Joseph Pane, and Beth Ann Griffin\*  
RAND

December 22, 2025

## 1 Introduction

The Toolkit for Weighting and Analysis of Nonequivalent Groups, `twang`, contains a set of functions and procedures to support causal modeling of observational data through the estimation and evaluation of propensity scores and associated weights for binary, multinomial, and time-varying treatments. This package was developed in 2004. The `twang` package underwent extensive revisions in 2012 and 2020. Starting with version 2.0, the `twang` package includes revisions that improve computational efficiency. Code written for prior versions of `twang` should still run without modification. However, in certain situations, results from the updated version of `twang` will no longer replicate previous results. If a user wishes to reproduce results from a prior version of `twang`, a new option (`version="legacy"`) must be specified.

This tutorial provides an introduction to `twang` and demonstrates its use through illustrative examples. Interested readers can review the related tutorials for more than two treatment groups and time-varying treatments at <https://www.rand.org/statistics/twang/tutorials.html>.

The foundation of the methods supported by `twang` is the propensity score. The propensity score is the probability that a particular case would be assigned or exposed to a treatment condition. Rosenbaum & Rubin (1983) showed that knowing the propensity score is sufficient to separate the effect of a treatment on an outcome from observed confounding factors that influence both treatment assignment and outcomes, provided the necessary conditions hold. The propensity score has the balancing property that given the propensity score the distribution of features for the treatment cases is the same as that for the control cases. While the treatment selection probabilities are generally not known, good estimates of them can be effective at diminishing or eliminating confounds between pretreatment group differences and treatment outcomes in the estimation of treatment effects.

There are now numerous propensity scoring methods in the literature. They differ in how they estimate the propensity score (e.g. logistic regression, CART), the target estimand (e.g. treatment effect on the treated, population treatment effect), and how they utilize the resulting estimated propensity scores (e.g. stratification, matching, weighting, doubly robust estimators). We originally developed the `twang` package with a particular process in mind, namely, generalized boosted regression to estimate the propensity scores and weighting of the comparison cases to estimate the average treatment effect on the treated (ATT). However, we have updated the package to also meaningfully handle the case where interest lies in using the population weights (e.g., weighting of comparison and treatment cases to estimate the population average treatment effect, ATE).

---

\*The `twang` package and this tutorial were developed under NIDA grants R01 DA017507, R01 DA015697-03, and R01 DA034065

The main workhorse of **twang** is the `ps()` function, which implements gradient boosted models with either the **gbm** package or the **xgboost** package to estimate the propensity scores. However, the framework of the package is flexible enough to allow the user to use propensity score estimates from other methods and to assess the usefulness of those estimates for ensuring equivalence (or “balance”) in the pretreatment covariate distributions of treatment and control groups using tools from the package. The same set of functions is also useful for other tasks such as non-response weighting, as discussed in Section 5.

The **twang** package aims to (i) compute from the data estimates of the propensity scores which yield accurate causal effect estimates, (ii) check the quality of the resulting propensity score weights by assessing whether or not they have the balancing properties that we expect in theory, and (iii) use them in computing treatment effect estimates. Users who are more comfortable with Stata or SAS than R are encouraged to visit [www.rand.org/statistics/twang/](http://www.rand.org/statistics/twang/) for information on Stata commands and SAS macros that implement these methods. Additionally, we now have Shiny apps available that are menu-driven and user-friendly versions of the R tools.

## 2 Principles of twang

The **twang** package utilizes gradient boosted models to estimate propensity scores. These models are flexible, machine learning approaches that can automatically incorporate nonlinearities and interactions among the covariates. In general, the algorithms implemented in the **twang** package optimize the similarity between treatment and comparison observations by selecting the iteration (complexity) of the gradient boosted model that achieves the best balance. The user simply selects the set of covariates they wish to balance, and the algorithm automatically incorporates nonlinearity and interactions to achieve the optimal balance. This vignette focuses on binary treatments using the `ps` function of the **twang** package, but many of the functions and arguments are shared across **twang** functions developed for other treatment types (e.g., multinomial treatment using `mnp`). See [www.rand.org/statistics/twang/](http://www.rand.org/statistics/twang/) for more information.

The arguments to the `ps` can be broken up into three categories: (1) those that specify the model; (2) those that control the gradient boosting; and (3) those that improve computational efficiency. We will describe these arguments in detail now, and then highlight their use in the next section.

### 2.1 Arguments that specify the propensity score model

`formula` is a symbolic description of the propensity score model using standard R syntax with the treatment indicator on the left side of the formula and the potential confounding variables on the right side. In contrast to the `lm()` function, there is no need to specify interaction terms in the formula as the gradient boosted model automatically considers interactions. There is also no need — and it can be counterproductive — to create indicator, or “dummy coded,” variables to represent categorical covariates, provided the categorical variables are stored as a `factor` or as `ordered` (see `help(factor)` for more details).

`data` indicates the dataset to be used, which should include a binary treatment indicator and the covariates to be balanced.

`estimand` is the causal effect of interest, which is either “ATE” (average treatment effect) or “ATT” (average treatment effect on the treated). ATE estimates the change in the outcome if the treatment were applied to the entire population versus if the control were applied to the entire population. ATT estimates the analogous effect, averaging only over the treated population. The `estimand` argument was added to the 2012 revision of the package which integrated ATE weighting into the package and the `ps` function estimate of the propensity score. The default value is “ATE”.

`sampw` are optional sampling weights. If specified, the sampling weights are automatically incorporated into the derivation of the propensity score weights.

## 2.2 Arguments that control the gradient boosting

`stop.method` is a method or set of methods for measuring and summarizing balance across pretreatment variables. The `ps` function selects the optimal number of gradient boosting iterations to minimize the differences between the treatment and control groups as measured by the rules of the given `stop.method`. Current options are `ks.mean`, `ks.max`, `es.mean`, and `es.max`. `ks` refers to the Kolmogorov-Smirnov statistic and `es` refers to standardized effect size. These are summarized across the pretreatment variables by either the maximum (`.max`) or the mean (`.mean`). The default value is to utilize both `ks.mean` and `es.mean`.

`n.trees` is the maximum number of gradient boosting iterations to be considered. The more iterations allows for more nonlinearity and interactions to be considered. The default value is 10,000.

`interaction.depth` is a positive integer denoting the tree depth used in gradient boosting. This can be loosely be interpreted as the maximum number of variables that can be included in an interaction. Higher values allow for greater model complexity. The default value is 3.

`shrinkage` is a numeric value between 0 and 1 denoting the learning rate. Smaller values restrict the complexity that is added at each iteration of the gradient boosting algorithm. A smaller learning rate requires more iterations (`n.trees`), but adds some protection against model overfit. The default value is 0.01.

`n.minobsinnode` is the minimum number of observations in the terminal nodes of the trees used in the gradient boosting algorithm. Smaller values allow for more model complexity, while larger values provide some protection against model overfit. The default value is 10.

## 2.3 Arguments that improve computational efficiency

`n.keep` is a positive integer indicating that the algorithm should only consider every `n.keep`-th iteration of the propensity score model and optimize balance over this set instead of all iterations. Larger values of `n.keep` improve computational efficiency by only assessing balance on a subset of the gradient boosting iterations, but is likely to achieve worse balance than considering every iteration. The default value is 1, which is to optimize over all iterations.

`n.grid` is a positive integer that sets the grid size for an initial search of the region most likely to minimize the `stop.method`. `n.grid` corresponds to a grid of points on the kept iterations as defined by `n.keep`. For example, `n.keep=20` with `n.trees=5000` will keep the iterations (20,40,...,5000). The option `n.grid=10` then splits this vector into 10 points. Thus, `n.grid*n.keep` must be less than or equal to `n.trees`. The default grid size is 25.

`ks.exact` improves the speed when calculating the p-value for the weighted two-sample Kolmogorov-Smirnov (KS) statistic. The default behavior when calculating the weighted two-sample KS p-value differs from implementations of `twang` prior to Version 2.0. Specifically, if `ks.exact=NULL` and the product of the effective sample sizes is less than 10,000, then an approximation based on the exact distribution of the unweighted KS statistic is used. This approximation via the exact distribution can also be requested directly by `ks.exact=TRUE`. Otherwise, an approximation based on the asymptotic distribution of the unweighted KS statistic is used. The default is `ks.exact=NULL`. We note that `ks.exact=TRUE` replicates the calculation of KS p-values from version of `twang` prior to Version 2.0 but adds substantial computation time for larger datasets. We recommend using the default value.

`version` allows the user to specify which gradient boosting package to use for the propensity score model. Current options are `"gbm"`, `"xgboost"`, and `"legacy"`. The default value is `version="gbm"` which uses gradient boosting from the `gbm` package and provides a faster implementation of key components of the `twang` package to improve speed when calculating and optimizing the balance statistics.

`version="xgboost"` uses gradient boosting from the `xgboost` package. This provides users with access to a cutting-edge implementation of gradient boosting, while potentially providing speed improvements in larger datasets. The behavior of `xgboost` can be controlled using the the boosting parameters of `gbm` discussed earlier, but can also be controlled directly using the options of `xgboost`. See <https://xgboost.readthedocs.io/en/latest/parameter.html> for a description of the different options. `version="legacy"` uses the implementation of the `ps` function prior to Version 2.0. To replicate results from analyses utilizing `twang` prior to Version 2.0, specify the `version="legacy"` option.

## 2.4 Typical workflow with `twang`

Here, we outline the usual steps for a propensity score analysis using `twang`. While meant to be illustrative, these steps should not be considered comprehensive. Every analysis is unique and requires careful consideration. These steps will be expanded upon in the subsequent examples.

1. Determine the estimand of interest, either ATE or ATT.
2. Determine the observed confounding factors to be balanced.
3. Fit the propensity score model using the `ps()` function.
  - Evaluate the convergence of the algorithm.
  - Assess the balance of the confounding factors before and after applying the propensity score weights.
  - If needed, rerun the algorithm adjusting the gradient boosting parameters to achieve convergence or improve balance.
4. Extract the propensity score weights and fit a weighted model to estimate the treatment effect.

## 3 An ATT example

### 3.1 Installing `twang` for first time users

The package is installed by typing `install.packages("twang", dependencies=TRUE)` into the R console. You will only need to do this step once. In the future, running `update.packages()` regularly will ensure that you have the latest versions of the packages, including bug fixes and new features.

### 3.2 Load the data

To start using `twang`, first load the package using `library(twang)`. This is required at the beginning of each R session. We will highlight features of `twang` using data from Lalonde's National Supported Work Demonstration analysis (Lalonde 1986, Dehejia & Wahba 1999, <http://users.nber.org/~rdehejia/nswwdata2.html>). This dataset is provided with the `twang` package and can be loaded using `data(lalonde)`.

```
> library(twang)
> data(lalonde)
```

R can read data from many other sources. The manual "R Data Import/Export," available at <http://cran.r-project.org/doc/manuals/R-data.pdf>, describes that process in detail.

### 3.3 Fit a gradient boosted model with ps

For the `lalonge` dataset, the variable `treat` is the treatment indicator, where 1 indicates being part of the National Supported Work Demonstration (“treatment”) and 0 indicates cases drawn from the Current Population Survey (“comparison”). In order to estimate a treatment effect for this demonstration program that is unbiased by pretreatment group differences on other observed covariates, we include the following covariates in a propensity score model of treatment: age, education, black, Hispanic, having no degree, married, earnings in 1974 (pretreatment), and earnings in 1975 (pretreatment). Note that we specify no outcome variables at this time. The `ps()` function is the primary method in `twang` for estimating propensity scores for binary treatments.

```
> ps.lalonge.gbm = ps(treat ~ age + educ + black + hispan + nodegree +
+                     married + re74 + re75,
+                     data = lalonge,
+                     n.trees=5000,
+                     interaction.depth=2,
+                     shrinkage=0.01,
+                     estimand = "ATT",
+                     stop.method=c("es.mean", "ks.max"),
+                     n.minobsinnode = 10,
+                     n.keep = 1,
+                     n.grid = 25,
+                     ks.exact = NULL,
+                     verbose=FALSE)
```

In this model, we have specified `estimand = "ATT"` indicating that we will estimate the average treatment effect among those who participated in the National Supported Work Demonstration. We also specified two stopping criteria: `es.mean` and `ks.max`. This indicates that the algorithm should separately determine the iteration that minimizes the average absolute standardized effect size (`es.mean`) and the maximum KS statistic (`ks.max`). We will compare the results of these stopping rules to evaluate the sensitivity of the results to this choice.

The saved object `ps.lalonge.gbm` is of class `ps` and contains the propensity score model fit information. We will work with this object to assess model convergence, to assess the balance of pretreatment characteristics, and to extract the propensity score weights needed for subsequent analysis. By default, the `ps()` function uses the `gbm` package for gradient boosting. A user can request use of `xgboost` by specifying the option `version="xgboost"`. The use of `xgboost` will not be discussed further in this tutorial. See the documentation for additional details.

#### 3.3.1 Convergence diagnostic checks

Before continuing to outcome analyses, an analyst should perform diagnostic checks assessing the convergence of the gradient boosting algorithm. The specified value of `n.trees` should be large enough to allow the gradient boosted model to have explored sufficiently complicated models. We can do this quickly with the `plot()` function.<sup>1</sup> As a default, the `plot()` function applied to a `ps` object gives the balance measures as a function of the number of iterations in the gradient boosting algorithm, with higher iterations corresponding to more complicated fitted models. In the example below, 2127 iterations minimized the average absolute standardized effect size and 1751 iterations minimized the largest of the eight Kolmogorov-Smirnov (KS) statistics computed for the covariates.

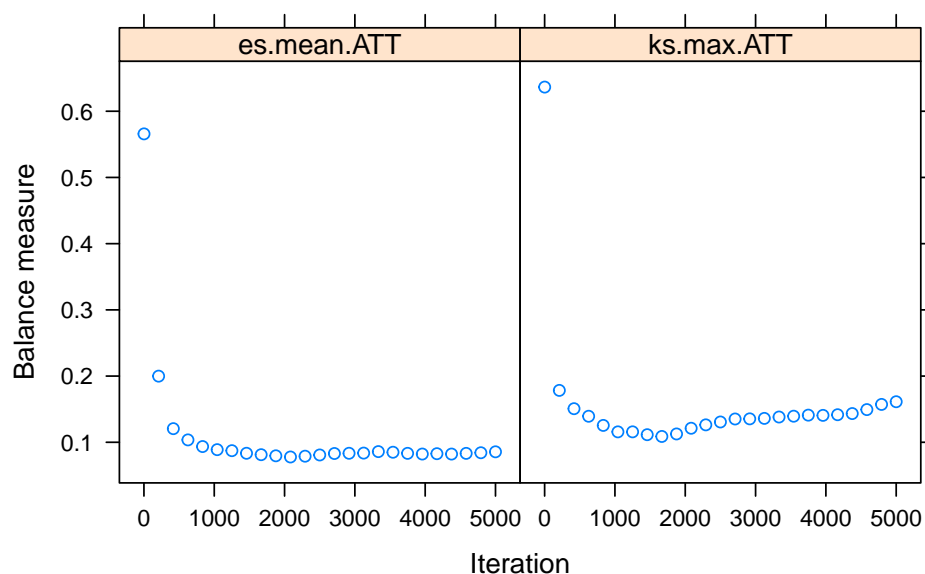
If it appears that additional iterations would be likely to result in lower values of the balance statistic, `n.trees` should be increased. However, after a point, additional complexity typically makes the balance

---

<sup>1</sup>In versions 1.0.x of the `twang` package, the `ps` function itself included some plotting functions. This is no longer the case (and the function no longer includes a `plots` argument); these functions have been moved to the generic `plot()` function.

worse, as in the example below. This figure also gives information on how compatible two or more stopping rules are: if the minima for multiple stopping rules under consideration are near one another, the results should not be sensitive to which stopping rule one uses for the final analysis. See Section 6.3 for a discussion of these and other balance measures.

```
> plot(ps.lalonde.gbm)
```



### 3.3.2 Assessing “balance” using balance tables

Having estimated the propensity scores, `bal.table()` produces a table that shows how well the resulting weights succeed in manipulating the control group so that its weighted pretreatment characteristics match, or balance, those of the unweighted treatment group when `estimand = "ATT"`. If `estimand = "ATE"`, both the control and treatment groups are weighted so that the weighted pretreatment characteristics match, or balance, with one another.

By default, the `bal.table()` function uses the value of `estimand` set with the `ps()` function call. For example, in the analysis we set `estimand = "ATT"` when calling `ps()` to estimate the propensity scores. The function `bal.table()` automatically uses the correct weights when checking balance and comparing the distributions of pre-treatment variables for the weighted control group with those from the unweighted treatment group.

```
> lalonde.balance <- bal.table(ps.lalonde.gbm)
> lalonde.balance
```

\$unw	tx.mn	tx.sd	ct.mn	ct.sd	std.eff.sz	stat	p	ks	ks.pval
age	25.816	7.155	28.030	10.787	-0.309	-2.994	0.003	0.158	0.003
educ	10.346	2.011	10.235	2.855	0.055	0.547	0.584	0.111	0.081
black	0.843	0.365	0.203	0.403	1.757	19.371	0.000	0.640	0.000
hispan	0.059	0.237	0.142	0.350	-0.349	-3.413	0.001	0.083	0.339

nodegree	0.708	0.456	0.597	0.491	0.244	2.716	0.007	0.111	0.081
married	0.189	0.393	0.513	0.500	-0.824	-8.607	0.000	0.324	0.000
re74	2095.574	4886.620	5619.237	6788.751	-0.721	-7.254	0.000	0.447	0.000
re75	1532.055	3219.251	2466.484	3291.996	-0.290	-3.282	0.001	0.288	0.000

\$es.mean.ATT

	tx.mn	tx.sd	ct.mn	ct.sd	std.eff.sz	stat	p	ks	ks.pval
age	25.816	7.155	25.802	7.279	0.002	0.015	0.988	0.122	0.892
educ	10.346	2.011	10.573	2.089	-0.113	-0.706	0.480	0.099	0.977
black	0.843	0.365	0.842	0.365	0.003	0.027	0.978	0.001	1.000
hispan	0.059	0.237	0.042	0.202	0.072	0.804	0.421	0.017	1.000
nodegree	0.708	0.456	0.609	0.489	0.218	0.967	0.334	0.099	0.977
married	0.189	0.393	0.189	0.392	0.002	0.012	0.990	0.001	1.000
re74	2095.574	4886.620	1556.930	3801.566	0.110	1.027	0.305	0.066	1.000
re75	1532.055	3219.251	1211.575	2647.615	0.100	0.833	0.405	0.103	0.969

\$ks.max.ATT

	tx.mn	tx.sd	ct.mn	ct.sd	std.eff.sz	stat	p	ks	ks.pval
age	25.816	7.155	25.760	7.412	0.008	0.059	0.953	0.106	0.923
educ	10.346	2.011	10.571	2.141	-0.112	-0.709	0.478	0.107	0.919
black	0.843	0.365	0.835	0.372	0.023	0.191	0.849	0.008	1.000
hispan	0.059	0.237	0.043	0.203	0.069	0.776	0.438	0.016	1.000
nodegree	0.708	0.456	0.601	0.490	0.235	1.101	0.271	0.107	0.919
married	0.189	0.393	0.199	0.400	-0.025	-0.171	0.864	0.010	1.000
re74	2095.574	4886.620	1678.117	3948.049	0.085	0.791	0.429	0.053	1.000
re75	1532.055	3219.251	1259.757	2677.096	0.085	0.715	0.475	0.093	0.973

`bal.table()` returns information on the pretreatment covariates before and after weighting. The object is a list with named components, one for an unweighted analysis (named `unw`) and one for each `stop.method` specified, here `es.mean` and `ks.max`. McCaffrey et al (2004) essentially used `es.mean` for the analyses, but our more recent work has sometimes used `ks.max`. See McCaffrey et al. (2013) for a more detailed description of these choices.

If there are missing values (represented as NA) in the covariates, `twang` will attempt to construct weights that also balance rates of missingness in the treatment and control arms. In this case, the `bal.table()` will have an extra row for each variable that has missing entries. User should note that missing data in `xgboost` is handled differently than in `gbm`. In `gbm`, missing values are placed in their own node, while in `xgboost` missing values are placed in the left or right node based on minimizing the objective function.

The columns of the table consist of the following items:

**tx.mn, ct.mn** The treatment means and the control means for each of the variables. The unweighted table (`unw`) shows the unweighted means. For each stopping rule the means are weighted using weights corresponding to the `gbm` model selected by `ps()` using the stopping rule. When `estimand = "ATT"` the weights for the treatment group always equal 1 for all cases and there is no difference between unweighted and propensity score weighted `tx.mn`.

**tx.sd, ct.sd** The propensity score weighted treatment and control groups' standard deviations for each of the variables. The unweighted table (`unw`) shows the unweighted standard deviations.

**std.eff.sz** The standardized effect size, defined as the treatment group mean minus the control group mean divided by the treatment group standard deviation if `estimand = "ATT"` or divided by the pooled sample (treatment and control) standard deviation if `estimand = "ATE"`. (In discussions of propensity scores this value is sometimes referred to as "standardized bias".) Occasionally, lack of treatment group

or pooled sample variance on a covariate results in very large (or infinite) standardized effect sizes. For purposes of analyzing mean effect sizes across multiple covariates, we set all standardized effect sizes larger than 500 to NA (missing values).

**stat**, **p** Depending on whether the variable is continuous or categorical, **stat** is a t-statistic or a  $\chi^2$  statistic. **p** is the associated p-value

**ks**, **ks.pval** The Kolmogorov-Smirnov test statistic and its associated p-value. P-values for the KS statistics are either derived from Monte Carlo simulations or analytic approximations, depending on the specifications made in the **perm.test.iters** argument of the **ps** function. For categorical variables this is just the  $\chi^2$  test p-value.

Components of these tables are useful for demonstrating that pretreatment differences between groups on observed variables have been eliminated using the weights.

The **summary()** method for **ps** objects offers a compact summary of the sample sizes of the groups and the balance measures. If **perm.test.iters**>0 was used to create the **ps** object, then Monte Carlo simulation is used to estimate p-values for the maximum KS statistic that would be expected across the covariates, had individuals with the same covariate values been assigned to groups randomly. Thus, a p-value of 0.04 for **max.ks.p** indicates that the largest KS statistic found across the covariates is larger than would be expected in 96% of trials in which the same cases were randomly assigned to groups. Otherwise, **max.ks.p** will be NA.

```
> summary(ps.lalonde.gbm)
```

	n.treat	n.ctrl	ess.treat	ess.ctrl	max.es	mean.es	max.ks
unw	185	429	185	429.00000	1.7567745	0.56872589	0.6404460
es.mean.ATT	185	429	185	22.96430	0.2177817	0.07746175	0.1223384
ks.max.ATT	185	429	185	27.18469	0.2346990	0.08007702	0.1069914

	max.ks.p	mean.ks	iter
unw	NA	0.27024507	NA
es.mean.ATT	NA	0.06361021	2127
ks.max.ATT	NA	0.06259804	1751

In general, weighted means can have greater sampling variance than unweighted means from a sample of equal size. The effective sample size (ESS) of the weighted comparison group captures this increase in variance as

$$ESS = \frac{(\sum_{i \in C} w_i)^2}{\sum_{i \in C} w_i^2}. \quad (1)$$

The ESS is approximately the number of observations from a simple random sample that yields an estimate with sampling variation equal to the sampling variation obtained with the weighted comparison observations. Therefore, the ESS will give an estimate of the number of comparison participants that are comparable to the treatment group when **estimand** = "ATT". The ESS is an accurate measure of the relative size of the variance of means when the weights are fixed or they are uncorrelated with outcomes. Otherwise the ESS underestimates the effective sample size (Little & Vartivarian, 2004). With propensity score weights, it is rare that weights are uncorrelated with outcomes. Hence the ESS typically gives a lower bound on the effective sample size, but it still serves as a useful measure for choosing among alternative models and assessing the overall quality of a model, even if it provides a possibly conservative picture of the loss in precision due to weighting.

The **ess.treat** and **ess.ctrl** columns in the summary results shows the ESS for the estimated propensity scores. Note that although the original comparison group had 429 cases, the propensity score estimates effectively utilize only 24 or 36.4 of the comparison cases, depending on the rules and measures used to



estimate the propensity scores. While this may seem like a large loss of sample size, this indicates that many of the original cases were unlike the treatment cases and, hence, were not useful for isolating the treatment effect. Moreover, similar or even greater reductions in ESS would be expected from alternative approaches to using propensity scores, such as matching or stratification. Since the estimand of interest in this example is ATT, `ess.treat = n.treat` throughout (i.e., all treatment cases have a weight of 1).

### 3.3.3 Graphical assessments of balance

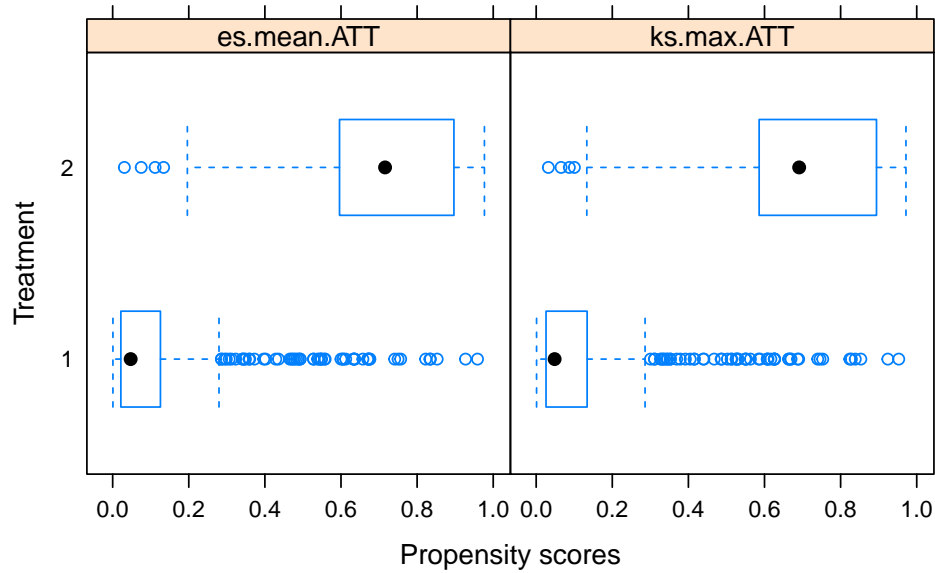
The `plot()` function can generate useful diagnostic plots from the propensity score objects. The full set of plots available in `twang` and the argument value of `plot` to produce each one are given in Table 1. The convergence plot — the default — was discussed above.

Descriptive argument	Numeric argument	Description
"optimize"	1	Balance measure as a function of GBM iterations
"boxplot"	2	Boxplot of treatment/control propensity scores
"es"	3	Standardized effect size of pretreatment variables
"t"	4	<i>t</i> -test <i>p</i> -values for weighted pretreatment variables
"ks"	5	Kolmogorov-Smirnov <i>p</i> -values for weighted pretreatment variables
"histogram"	6	Histogram of weights for treatment/control

Table 1: Available options for `plots` argument to `plot()` function.

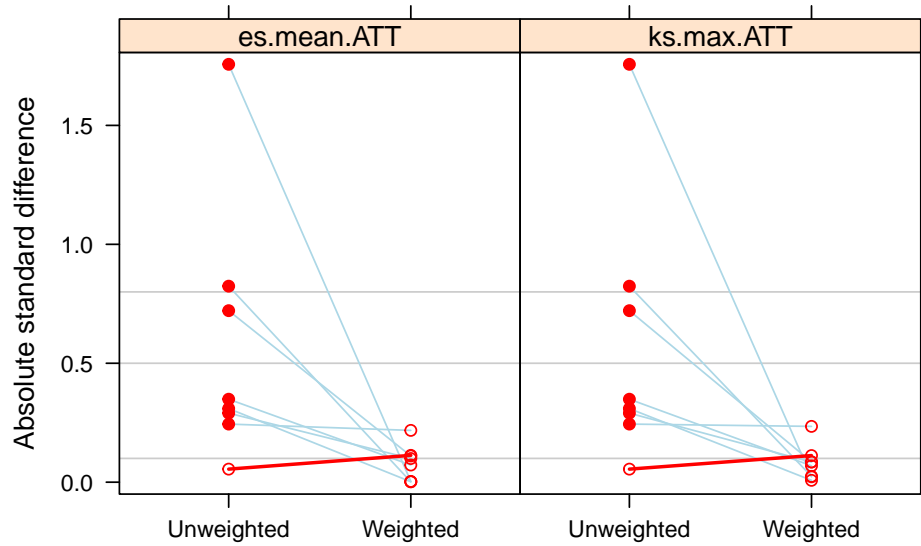
The `plot()` function takes a `plots` argument in order to produce other diagnostic plots. For example, specifying `plots = 2` or `plots = "boxplot"` produces boxplots illustrating the spread of the estimated propensity scores in the treatment and comparison groups. Whereas propensity score stratification requires considerable overlap in these spreads, excellent covariate balance can often be achieved with weights, even when the propensity scores estimated for the treatment and control groups show little overlap.

```
> plot(ps.lalonde.gbm, plots=2)
```



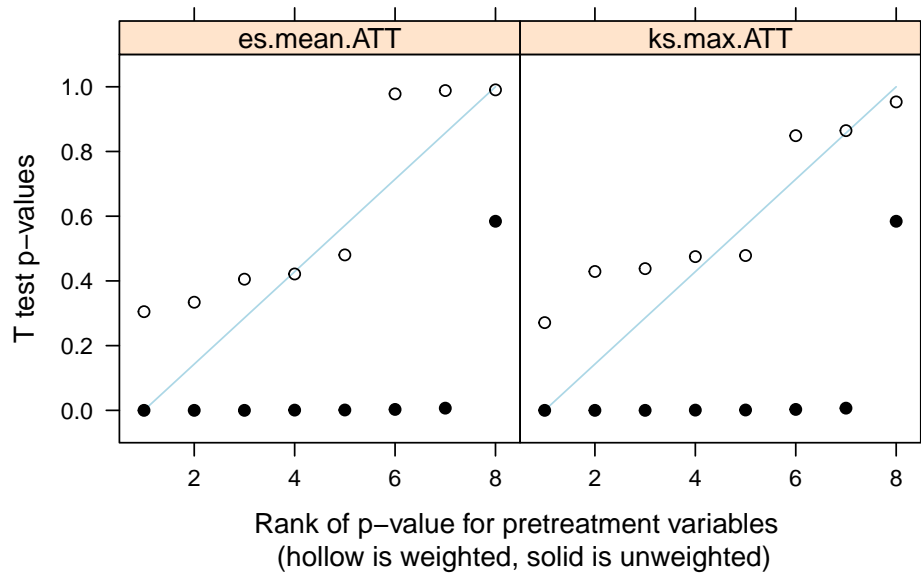
The effect size plot illustrates the effect of weights on the magnitude of differences between groups on each pretreatment covariate. These magnitudes are standardized using the standardized effect size described earlier. In these plots, substantial reductions in effect sizes are observed for most variables (blue lines), with only one variable showing an increase in effect size (red lines), but only a seemingly trivial increase. Closed red circles indicate a statistically significant difference, many of which occur before weighting, none after. In some analyses variables can have very little variance in the treatment group sample or the entire sample and group differences can be very large relative to the standard deviations. In these situations, the user is warned that some effect sizes are too large to plot.

```
> plot(ps.lalonde.gbm, plots=3)
```



When many of the p-values testing individual covariates for balance are very small, the groups are clearly imbalanced and inconsistent with what we would expect had the groups been formed by random assignment. After weighting we would expect the p-values to be larger if balance had been achieved. We use a QQ plot comparing the quantiles of the observed p-values to the quantiles of the uniform distribution (45 degree line) to conduct this check of balance. Ideally, the p-values from independent tests in which the null hypothesis is true will have a uniform distribution. Although the ideal is unlikely to hold even if we had random assignment (Bland, 2013), severe deviation of the p-values below the diagonal suggests lack of balance and p-values running at or above the diagonal suggests balance might have been achieved. The p-value plot (`plots=4` or `plots="t"`) allows users to visually inspect the p-values of the t-tests for group differences in the covariate means.

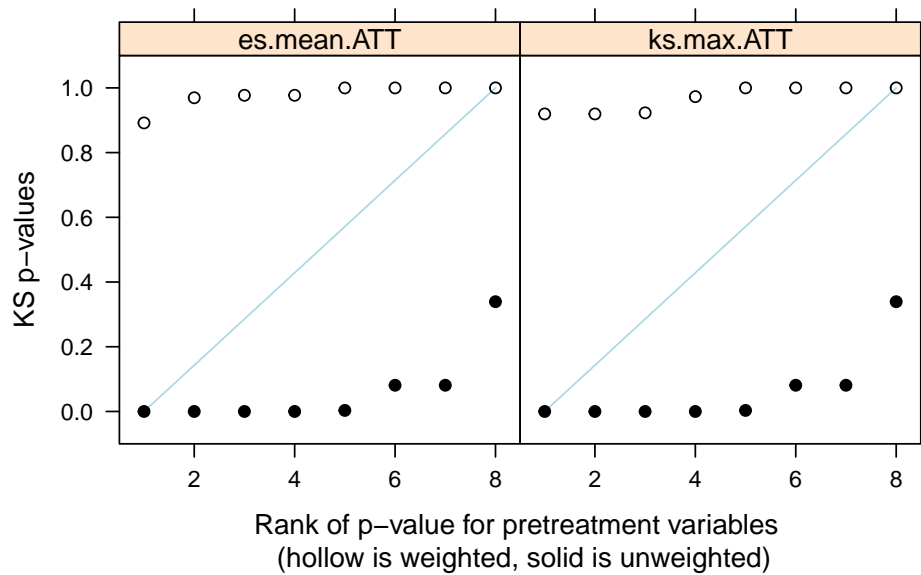
```
> plot(ps.lalonde.gbm, plots = 4)
```



Before weighting (closed circles), the groups have statistically significant differences on many variables (i.e., p-values are near zero). After weighting (open circles) the p-values are generally above the 45-degree line, which represents the cumulative distribution of a uniform variable on  $[0,1]$ . This indicates that the p-values are even larger than would be expected in a randomized study.

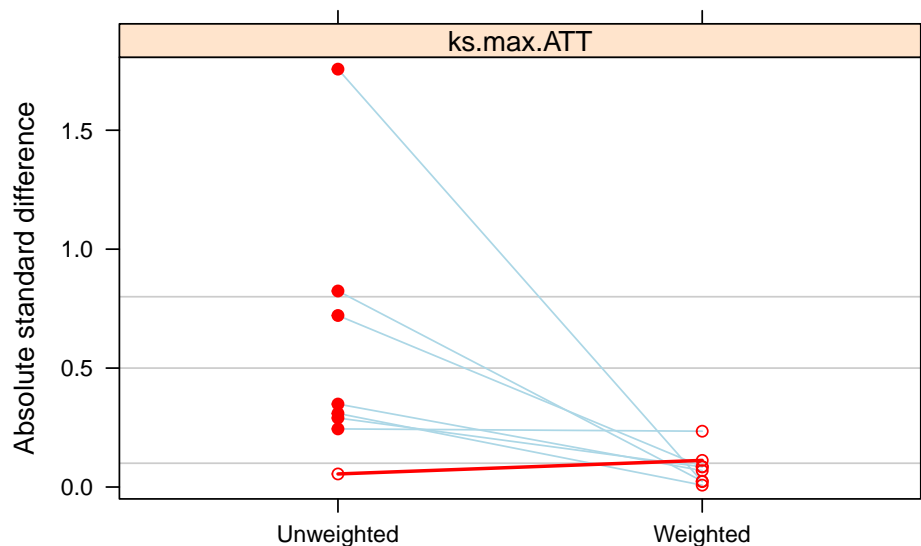
One can inspect similar plots for the KS statistic with the argument `plots = "ks"` or `plots = 5`.

```
> plot(ps.lalonde.gbm, plots = 5)
```



In all cases, the `subset` argument can be used if we wish to focus on results from one stopping rule.

```
> plot(ps.lalonde.gbm, plots = 3, subset = 2)
```



### 3.3.4 Understanding the relationship between the covariates and the treatment assignment

The `gbm` package has various tools for exploring the relationship between the covariates and the treatment assignment indicator if these are of interest. These tools can be used with the `ps` object to extract useful information. In particular, the `summary()` function applied to the underlying `gbm` object computes the relative influence of each variable for estimating the probability of treatment assignment. The relative influence of each variable changes as more iteration are added to the gradient boosted model. In this example, we choose the number of iterations to be the optimal number for minimizing the largest of the KS statistics. This value can be found in the `ps.lalonde.gbm$desc$ks.max.ATT$n.trees`. Figure 1 shows the barchart of the relative influence and is produced when `plot=TRUE` in the call to `summary()`.

```
> summary(ps.lalonde.gbm$gbm.obj,
+         n.trees=ps.lalonde.gbm$desc$ks.max.ATT$n.trees,
+         plot=FALSE)
```

	var	rel.inf
black	black	57.84196723
age	age	16.51419940
re74	re74	15.61206620
re75	re75	3.57534275
married	married	3.02199483
educ	educ	2.91568989
nodegree	nodegree	0.43801453
hispan	hispan	0.08072517

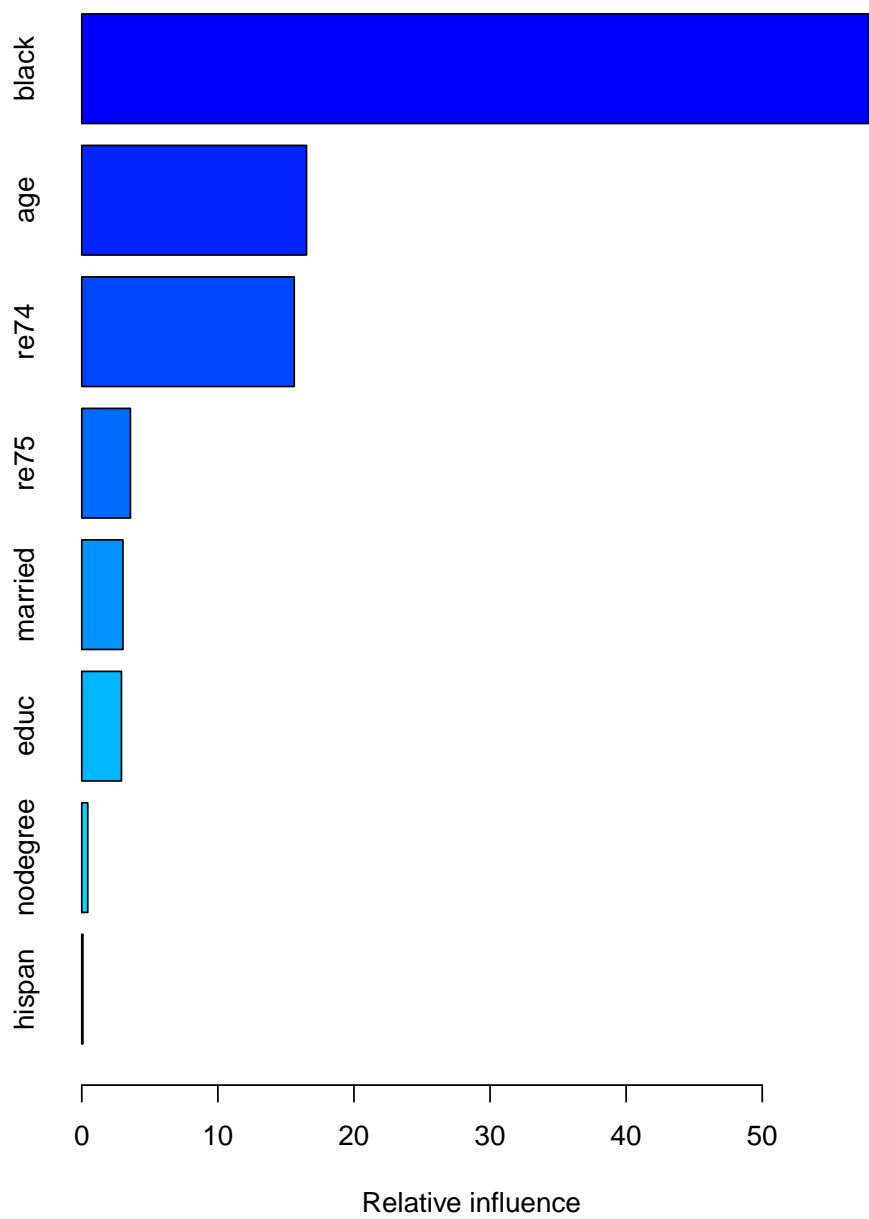


Figure 1: Relative influence of the covariates on the estimated propensity score.

Additionally, our team has developed the **SBdecomp** package for R which can be used to quantify the proportion of the estimated selection bias explained by each observed confounder when estimating causal effects using propensity score weights. It includes two approaches to quantify the proportion of the selection bias explained by each observed confounder: a single confounder removal approach and a single confounder inclusion approach. This tool can help analyze data where there is a substantive interest in identifying the variable or variables that explains the largest proportion of the estimated selection bias. Interested users are encourage to review the tutorial for that package at <https://www.rand.org/pubs/tools/TLA570-3.html>.

## 3.4 Analysis of outcomes

### 3.4.1 Propensity scores from `ps()`

A separate R package, the **survey** package, is useful for performing the outcomes analyses using weights. Its statistical methods account for the weights when computing standard error estimates. It is not a part of the standard R installation but installing **twang** should automatically install **survey** as well.

```
> library(survey)
```

The `get.weights()` function extracts the propensity score weights from a **ps** object. Those weights may then be used as case weights in a **svydesign** object. By default, it returns weights corresponding to the estimand (ATE or ATT) that was specified in the original call to `ps()`. If needed, the user can override the default via the optional `estimand` argument.

```
> lalonde$w <- get.weights(ps.lalonde.gbm, stop.method="es.mean")
> design.ps <- svydesign(ids=~1, weights=~w, data=lalonde)
```

The `stop.method` argument specifies which GBM model, and consequently which weights, to utilize. The **svydesign** function from the **survey** package creates an object that stores the dataset along with design information needed for analyses. See `help(svydesign)` for more details on setting up **svydesign** objects.

The aim of the National Supported Work Demonstration analysis is to determine whether the program was effective at increasing earnings in 1978. The propensity score adjusted test can be computed with **svyglm**.

```
> glm1 <- svyglm(re78 ~ treat, design=design.ps)
> summary(glm1)
```

Call:

```
svyglm(formula = re78 ~ treat, design = design.ps)
```

Survey design:

```
svydesign(ids = ~1, weights = ~w, data = lalonde)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	5616.6	884.9	6.347	4.28e-10 ***
treat	732.5	1056.6	0.693	0.488

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for gaussian family taken to be 49804197)

Number of Fisher Scoring iterations: 2

The analysis estimates an increase in earnings of \$733 for those that participated in the NSW compared with similarly situated people observed in the CPS. The effect, however, does not appear to be statistically significant.

Some authors have recommended utilizing both propensity score adjustment and additional covariate adjustment to minimize mean square error or to obtain “doubly robust” estimates of the treatment effect (Huppler-Hullsiek & Louis 2002, Bang & Robins 2005). These estimators are consistent if either the propensity scores are estimated correctly *or* the regression model is specified correctly. For example, note that the balance table for `es.mean.ATT` made the two groups more similar on `nodegree`, but still some differences remained, 70.8% of the treatment group had no degree while 60.9% of the comparison group had no degree. While linear regression is sensitive to model misspecification when the treatment and comparison groups are dissimilar, the propensity score weighting has made them more similar, perhaps enough so that additional modeling with covariates can adjust for any remaining differences. In addition to potential bias reduction, the inclusion of additional covariates can reduce the standard error of the treatment effect if some of the covariates are strongly related to the outcome.

```
> glm2 <- svyglm(re78 ~ treat + nodegree, design=design.ps)
> summary(glm2)
```

Call:

```
svyglm(formula = re78 ~ treat + nodegree, design = design.ps)
```

Survey design:

```
svydesign(ids = ~1, weights = ~w, data = lalonde)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	6768.4	1471.0	4.601	5.11e-06 ***
treat	920.3	1082.8	0.850	0.396
nodegree	-1891.8	1261.9	-1.499	0.134

---

Signif. codes: 0 ‘\*\*\*’ 0.001 ‘\*\*’ 0.01 ‘\*’ 0.05 ‘.’ 0.1 ‘ ’ 1

(Dispersion parameter for gaussian family taken to be 49013778)

Number of Fisher Scoring iterations: 2

Adjusting for the remaining group difference in the `nodegree` variable slightly increased the estimate of the program’s effect to \$920, but the difference is still not statistically significant. We can further adjust for the other covariates, but that too in this case has little effect on the estimated program effect.

```
> glm3 <- svyglm(re78 ~ treat + age + educ + black + hispan + nodegree +
+               married + re74 + re75,
+               design=design.ps)
> summary(glm3)
```

Call:

```
svyglm(formula = re78 ~ treat + age + educ + black + hispan +
       nodegree + married + re74 + re75, design = design.ps)
```

Survey design:

```
svydesign(ids = ~1, weights = ~w, data = lalonde)
```



Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-2.459e+03	4.289e+03	-0.573	0.56671
treat	7.585e+02	1.019e+03	0.745	0.45674
age	3.005e+00	5.558e+01	0.054	0.95691
educ	7.488e+02	2.596e+02	2.884	0.00406 **
black	-7.627e+02	1.012e+03	-0.753	0.45153
hispan	6.106e+02	1.711e+03	0.357	0.72123
nodegree	5.350e+02	1.626e+03	0.329	0.74227
married	4.918e+02	1.072e+03	0.459	0.64660
re74	5.699e-02	1.801e-01	0.316	0.75176
re75	1.568e-01	1.946e-01	0.806	0.42076

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for gaussian family taken to be 47150852)

Number of Fisher Scoring iterations: 2

### 3.4.2 Estimating the program effect using linear regression

The classical regression approach to estimating the program effect would fit a linear model with a treatment indicator and linear terms for each of the covariates.

```
> glm4 <- lm(re78 ~ treat + age + educ + black + hispan + nodegree +
+           married + re74 + re75,
+           data=lalonde)
> summary(glm4)
```

Call:

```
lm(formula = re78 ~ treat + age + educ + black + hispan + nodegree +
    married + re74 + re75, data = lalonde)
```

Residuals:

Min	1Q	Median	3Q	Max
-13595	-4894	-1662	3929	54570

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	6.651e+01	2.437e+03	0.027	0.9782
treat	1.548e+03	7.813e+02	1.982	0.0480 *
age	1.298e+01	3.249e+01	0.399	0.6897
educ	4.039e+02	1.589e+02	2.542	0.0113 *
black	-1.241e+03	7.688e+02	-1.614	0.1071
hispan	4.989e+02	9.419e+02	0.530	0.5966
nodegree	2.598e+02	8.474e+02	0.307	0.7593
married	4.066e+02	6.955e+02	0.585	0.5590
re74	2.964e-01	5.827e-02	5.086	4.89e-07 ***
re75	2.315e-01	1.046e-01	2.213	0.0273 *

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6948 on 604 degrees of freedom

Multiple R-squared: 0.1478, Adjusted R-squared: 0.1351

F-statistic: 11.64 on 9 and 604 DF, p-value: < 2.2e-16

This model estimates a rather strong treatment effect, estimating a program effect of \$1548 with a p-value=0.048. Several variations of this regression approach also estimate strong program effects. For example using square root transforms on the earnings variables yields a p-value=0.016. These estimates, however, are very sensitive to the model structure since the treatment and control subjects differ greatly as seen in the unweighted balance comparison from `bal.table(ps.lalonde.gbm)`.

### 3.5 Propensity scores from logistic regression

Propensity score analysis is intended to avoid problems associated with the misspecification of covariate adjusted models of outcomes, but the quality of the balance and the treatment effect estimates can be sensitive to the method used to estimate the propensity scores. Consider estimating the propensity scores using logistic regression instead of `ps()`.

```
> ps.logit <- glm(treat ~ age + educ + black + hispan + nodegree +
+                 married + re74 + re75,
+                 data = lalonde,
+                 family = binomial)
> lalonde$w.logit <- rep(1,nrow(lalonde))
> lalonde$w.logit[lalonde$treat==0] <- exp(predict(ps.logit,subset(lalonde,treat==0)))
```

`predict()` for logistic regression model produces estimates on the log-odds scale by default. Exponentiating those predictions for the comparison subjects gives the ATT weights  $p/(1-p)$ . `dx.wts()` from the `twang` package diagnoses the balance for an arbitrary set of weights producing a balance table. This function requires the user to specify the estimand argument in order to perform the appropriate calculations relative to the target group on which we are drawing inferences. The function `dx.wts` has not been updated in Version 2.0 and still relies on the older version of the balance calculations.

```
> bal.logit <- dx.wts(x = lalonde$w.logit,
+                    data=lalonde,
+                    vars=c("age", "educ", "black", "hispan", "nodegree",
+                          "married", "re74", "re75"),
+                    treat.var="treat",
+                    perm.test.iters=0, estimand = "ATT")
> bal.logit
```

	type	n.treat	n.ctrl	ess.treat	ess.ctrl	max.es	mean.es	max.ks
1	unw	185	429	185	429.00000	1.7567745	0.5687259	0.6404460
2		185	429	185	99.81539	0.1188496	0.0318841	0.3078039

```

      mean.ks iter
1 0.27024507  NA
2 0.09302319  NA
```

Applying the `bal.table()` function to this object returns a variable-by-variable summary of balance, just like it did for the `ps` object.

```
> bal.table(bal.logit)
```

```
$unw
```

	tx.mn	tx.sd	ct.mn	ct.sd	std.eff.sz	stat	p	ks
age	25.816	7.155	28.030	10.787	-0.309	-2.994	0.003	0.158
educ	10.346	2.011	10.235	2.855	0.055	0.547	0.584	0.111
black	0.843	0.365	0.203	0.403	1.757	19.371	0.000	0.640
hispan	0.059	0.237	0.142	0.350	-0.349	-3.413	0.001	0.083
nodegree	0.708	0.456	0.597	0.491	0.244	2.716	0.007	0.111
married	0.189	0.393	0.513	0.500	-0.824	-8.607	0.000	0.324
re74	2095.574	4886.620	5619.237	6788.751	-0.721	-7.254	0.000	0.447
re75	1532.055	3219.251	2466.484	3291.996	-0.290	-3.282	0.001	0.288

```
ks.pval
```

```

age      0.003
educ     0.075
black    0.000
hispan   0.319
nodegree 0.075
married  0.000
re74     0.000
re75     0.000

```

```
[[2]]
```

```

      tx.mn   tx.sd   ct.mn   ct.sd std.eff.sz   stat    p    ks
age      25.816   7.155   24.966  10.535    0.119  0.739 0.460 0.308
educ     10.346   2.011   10.403   2.459   -0.028 -0.219 0.827 0.036
black     0.843   0.365    0.845   0.362   -0.006 -0.069 0.945 0.002
hispan    0.059   0.237    0.059   0.236    0.001  0.008 0.993 0.000
nodegree  0.708   0.456    0.690   0.463    0.040  0.332 0.740 0.018
married   0.189   0.393    0.171   0.377    0.047  0.456 0.649 0.019
re74     2095.574 4886.620 2106.045 4235.832   -0.002 -0.022 0.983 0.228
re75     1532.055 3219.251 1496.541 2716.258    0.011  0.107 0.915 0.133

      ks.pval
age      0.000
educ     1.000
black    1.000
hispan   1.000
nodegree 1.000
married  1.000
re74     0.002
re75     0.185

```

For weights estimated with logistic regression, the largest KS statistic was reduced from the unweighted sample's largest KS of 0.64 to 0.31, which is still quite a large KS statistic.

Table 2 and 3 compares the balancing quality of the weights directly with one another using the standardize effect size and the KS statistic, respectively. The standardized effect sizes for both sets of propensity score weights are improved compared to the unweighted analysis. The KS statistics for **age** and **re74** are larger for the logistic regression model than the GBM-based propensity score model. Combining the results of Table 2 and 3 we can conclude that logistic regression is achieving mean balance between the groups (Table 2), but it does not balance the full distribution (Table 3).

Covariate	Unweighted	Using <b>ps</b>	Using logistic regression
age	-0.309	0.002	0.119
educ	0.055	-0.113	-0.028
black	1.757	0.003	-0.006
hispan	-0.349	0.072	0.001
nodegree	0.244	0.218	0.04
married	-0.824	0.002	0.047
re74	-0.721	0.11	-0.002
re75	-0.29	0.1	0.011

Table 2: Standardized effect size of covariates using GBM and logistic regression.

Covariate	Unweighted	Using ps	Using logistic regression
age	0.158	0.122	0.308
educ	0.111	0.099	0.036
black	0.64	0.001	0.002
hispan	0.083	0.017	0
nodegree	0.111	0.099	0.018
married	0.324	0.001	0.019
re74	0.447	0.066	0.228
re75	0.288	0.103	0.133

Table 3: KS statistic of covariates using GBM and logistic regression.

```
> design.logit <- svydesign(ids=~1, weights=~w.logit, data=lalonde)
> glm6 <- svyglm(re78 ~ treat, design=design.logit)
> summary(glm6)
```

Call:

```
svyglm(formula = re78 ~ treat, design = design.logit)
```

Survey design:

```
svydesign(ids = ~1, weights = ~w.logit, data = lalonde)
```

Coefficients:

```
      Estimate Std. Error t value Pr(>|t|)
(Intercept)  5135.1      588.9   8.719  <2e-16 ***
treat        1214.1      824.7   1.472    0.142
```

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for gaussian family taken to be 49598072)

Number of Fisher Scoring iterations: 2

The analysis estimates an increase in earnings of \$1214 for those that participated in the NSW compared with similarly situated people observed in the CPS. Table 4 compares all of the treatment effect estimates.

Treatment effect	PS estimate	Linear adjustment
\$733	GBM, minimize ES	none
\$920	GBM, minimize ES	nodegree
\$758	GBM, minimize ES	all
\$1548	None	all
\$1214	Logistic regression	none
\$1237	Logistic regression	all

Table 4: Treatment effect estimates by various methods

## 4 An ATE example

In the analysis of Section 2, we focused on estimating ATT for the `lalonge` dataset. In this situation, the ATE is not of great substantive interest because not all people who are offered entrance into the program could be expected to take advantage of the opportunity. Further, there is some evidence that the treated subjects were drawn from a subset of the covariate space. In particular, in an ATE analysis, we see that we are unable to achieve balance, especially for the “black” indicator.

We now turn to an ATE analysis that is feasible and meaningful. We focus on the `lindner` dataset, which was included in the `USPS` package (Obenchain 2011), and is now included in `twang` for convenience. A tutorial by Helmreich and Pruzek (2009; HP) for the `PSAgraphics` package also uses propensity scores to analyze a portion of these data. HP describe the data as follows on p. 3 with our minor recodings in square braces:

The `lindner` data contain data on 996 patients treated at the Lindner Center, Christ Hospital, Cincinnati in 1997. Patients received a Percutaneous Coronary Intervention (PCI). The data consists of 10 variables. Two are outcomes: `[sixMonthSurvive]` ranges over two values... depending on whether patients survived to six months post treatment [denoted by `TRUE`] or did not survive to six months [`FALSE`]... Secondly, `cardbill` contains the costs in 1998 dollars for the first six months (or less if the patient did not survive) after treatment... The treatment variable is `abcix`, where 0 indicates PCI treatment and 1 indicates standard PCI treatment and additional treatment in some form with abciximab. Covariates include `acutemi`, 1 indicating a recent acute myocardial infarction and 0 not; `ejecfrac` for the left ventricle ejection fraction, a percentage from 0 to 90; `ves1proc` giving the number of vessels (0 to 5) involved in the initial PCI; `stent` with 1 indicating coronary stent inserted, 0 not; `diabetic` where 1 indicates that the patient has been diagnosed with diabetes, 0 not; `height` in centimeters and `female` coding the sex of the patient, 1 for female, 0 for male.

HP focus on `cardbill` — the cost for the first months after treatment — as their outcome of interest. However, since not all patients survived to six months, it is not clear whether a lower value of `cardbill` is good or not. For this reason, we choose six-month survival (`sixMonthSurvive`) as our outcome of interest.

Ignoring pre-treatment variables, we see that `abcix` is associated with lower rates of 6-month mortality:

```
> data(lindner)
> table(lindner$sixMonthSurvive, lindner$abcix)

      0    1
FALSE  15   11
TRUE  283 687

> chisq.test(table(lindner$sixMonthSurvive, lindner$abcix))

Pearson's Chi-squared test with Yates' continuity correction

data:  table(lindner$sixMonthSurvive, lindner$abcix)
X-squared = 8.5077, df = 1, p-value = 0.003536
```

The question is whether this association is causal. If health care policies were to be made on the basis of these data, we would wish to elicit expert opinion as to whether there are likely to be other confounding pretreatment variables. For this tutorial, we simply follow HP in choosing the pre-treatment covariates. The `twang` model is fit as follows

```

> set.seed(1)
> ps.lindner <- ps(abcix ~ stent + height + female + diabetic +
+                 acutemi + ejecfrac + ves1proc,
+                 data = lindner,
+                 estimand = "ATE",
+                 verbose = FALSE)

```

We set `estimand = "ATE"` because we are interested in the effects of abciximab on everyone in the population. We do not specify the stopping rules. Consequently `ps()` uses the defaults: `es.mean` and `ks.mean`. We then inspect pre- and post-weighting balance with the command

```

> bal.table(ps.lindner)

```

```

$unw
      tx.mn tx.sd ct.mn ct.sd std.eff.sz stat      p      ks ks.pval
stent    0.705 0.456 0.584 0.494      0.257 3.624 0.000 0.121 0.004
height  171.443 10.695 171.446 10.589      0.000 -0.005 0.996 0.025 0.999
female    0.331 0.471 0.386 0.488     -0.115 -1.647 0.100 0.055 0.554
diabetic  0.205 0.404 0.268 0.444     -0.152 -2.127 0.034 0.064 0.367
acutemi   0.179 0.384 0.060 0.239      0.338 5.923 0.000 0.119 0.006
ejecfrac  50.403 10.419 52.289 10.297     -0.181 -2.640 0.008 0.114 0.009
ves1proc  1.463 0.706 1.205 0.480      0.393 6.693 0.000 0.188 0.000

```

```

$ks.mean.ATE
      tx.mn tx.sd ct.mn ct.sd std.eff.sz stat      p      ks ks.pval
stent    0.683 0.466 0.657 0.475      0.054 0.720 0.472 0.026 1.000
height  171.470 10.549 171.589 10.594     -0.011 -0.153 0.879 0.015 1.000
female    0.338 0.473 0.345 0.476     -0.015 -0.202 0.840 0.007 1.000
diabetic  0.215 0.411 0.229 0.421     -0.033 -0.430 0.667 0.014 1.000
acutemi   0.148 0.355 0.107 0.310      0.114 1.333 0.183 0.040 0.947
ejecfrac  51.051 10.334 51.604 9.112     -0.053 -0.799 0.425 0.027 1.000
ves1proc  1.395 0.666 1.337 0.573      0.089 1.199 0.231 0.027 1.000

```

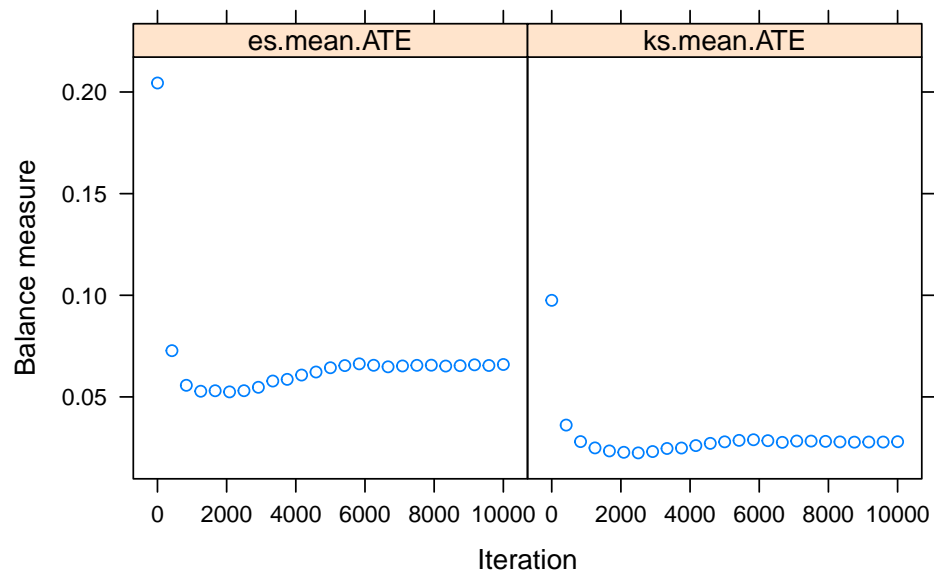
```

$es.mean.ATE
      tx.mn tx.sd ct.mn ct.sd std.eff.sz stat      p      ks ks.pval
stent    0.683 0.466 0.656 0.476      0.056 0.751 0.453 0.027 1.00
height  171.467 10.542 171.586 10.660     -0.011 -0.151 0.880 0.016 1.00
female    0.338 0.473 0.345 0.476     -0.015 -0.206 0.837 0.007 1.00
diabetic  0.215 0.411 0.231 0.422     -0.039 -0.506 0.613 0.016 1.00
acutemi   0.148 0.355 0.108 0.311      0.113 1.327 0.185 0.040 0.95
ejecfrac  51.037 10.348 51.546 9.171     -0.049 -0.733 0.464 0.027 1.00
ves1proc  1.396 0.666 1.342 0.579      0.082 1.087 0.277 0.025 1.00

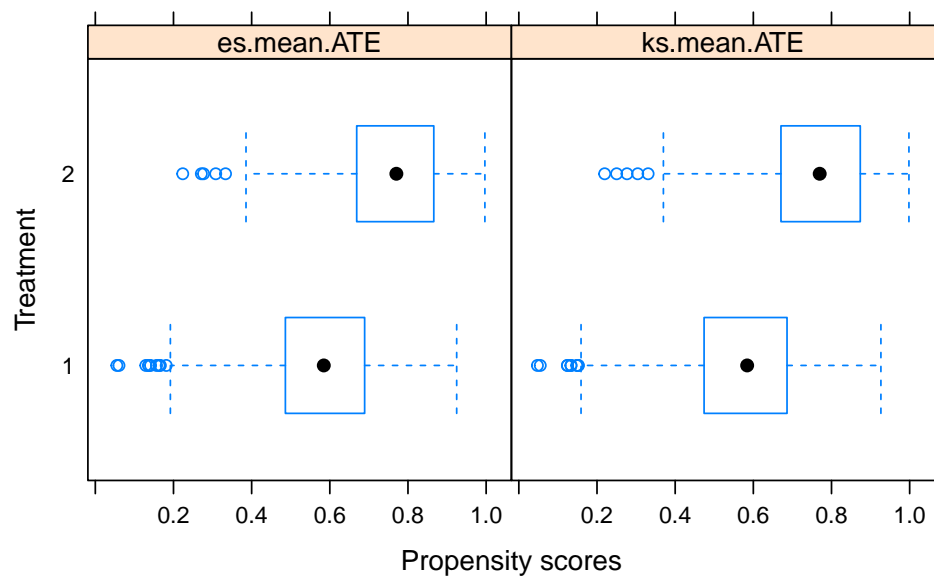
```

This balance table shows that `stent`, `acutemi`, and `ves1proc` were all significantly imbalanced before weighting. After weighting (using either `stop.method` considered) we do not see problems in this regard. Examining `plot(ps.lindner, plots = x)` for `x` running from 1 to 5 does not reveal problems, either. In regard to the optimize plot, we note that the scales of the KS and ES statistics presented in the optimize plots are not necessarily comparable. The fact that the KS values are lower than the ES values in the optimize plot does not suggest that the KS stopping rule is finding superior models. Each panel of the optimize plot indicates the gbm model that minimizes each stopping rule. The panels should not be compared other than to compare the number of iterations selected by each rule.

```
> plot(ps.lindner, plots = 1)
```

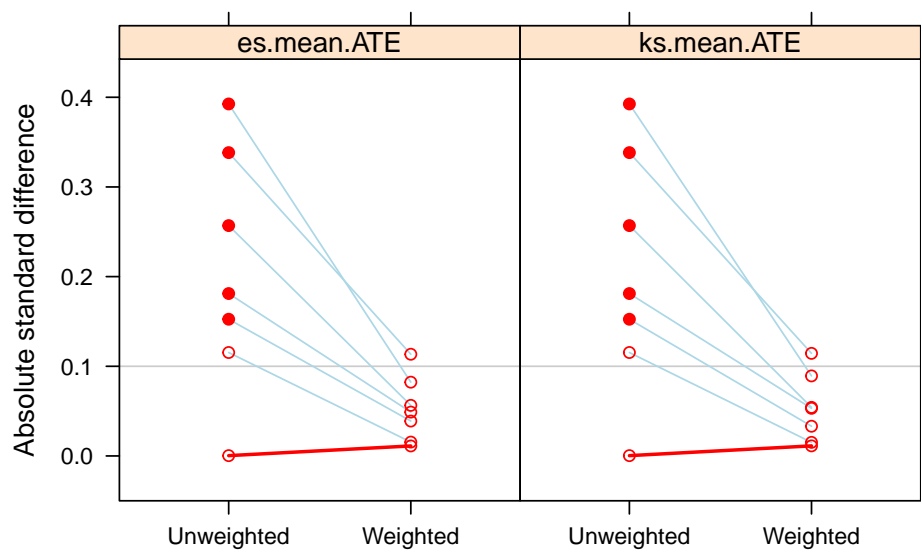


```
> plot(ps.lindner, plots = 2)
```

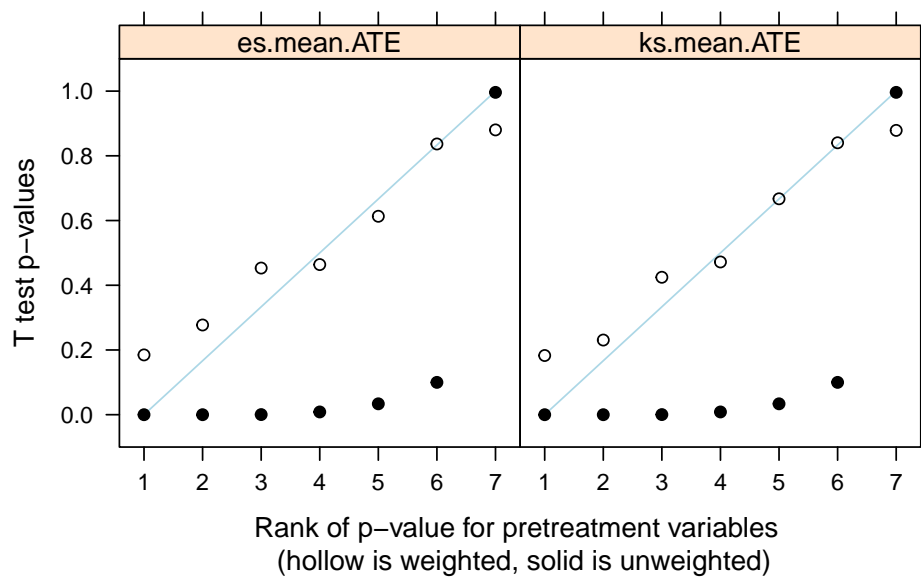




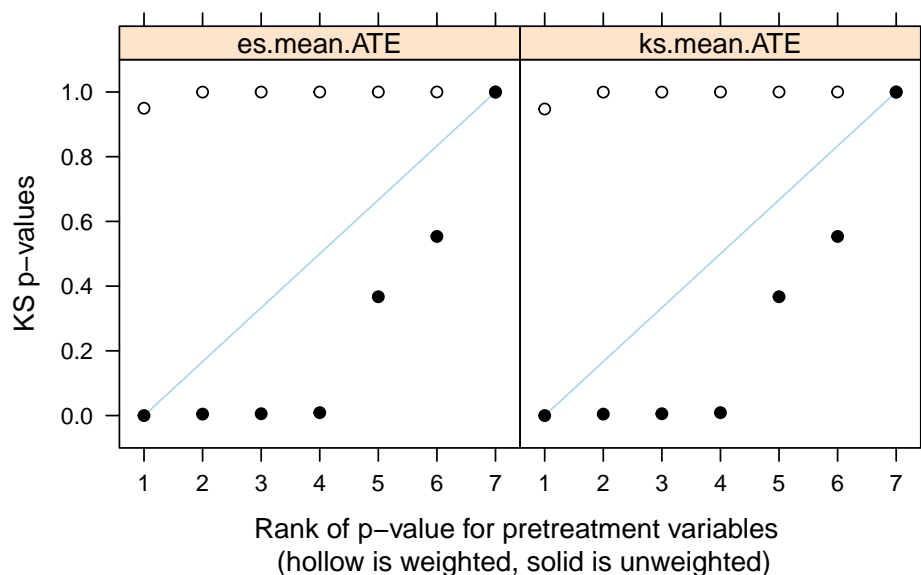
```
> plot(ps.lindner, plots = 3)
```



```
> plot(ps.lindner, plots = 4)
```



```
> plot(ps.lindner, plots = 5)
```



From a call to `summary()`, we see that the `es.mean.ATE` stopping rule results in a slightly higher ESS with comparable balance measures, so we proceed with those weights. Also, we note that `ess.treat` is no longer equal to `n.treat` since we are focusing on ATE rather than ATT.

```
> summary(ps.lindner)
```

	n.treat	n.ctrl	ess.treat	ess.ctrl	max.es	mean.es	max.ks
unw	698	298	698.0000	298.0000	0.3925637	0.20528943	0.18841945
ks.mean.ATE	698	298	655.7339	229.0034	0.1144059	0.05292428	0.04013745
es.mean.ATE	698	298	658.3496	230.5940	0.1134167	0.05236774	0.03979038

	max.ks.p	mean.ks	iter
unw	NA	0.09791845	NA
ks.mean.ATE	NA	0.02233864	2582
es.mean.ATE	NA	0.02262830	2116

As before, we use the `survey` package to reweight our sample and perform the analysis.

```
> lindner$w <- get.weights(ps.lindner, stop.method = "es.mean")
> design.ps <- svydesign(ids=~1, weights = ~w, data = lindner)
> svychisq(~sixMonthSurvive + abcix, design = design.ps)
```

Pearson's  $X^2$ : Rao & Scott adjustment

```
data: svychisq(~sixMonthSurvive + abcix, design = design.ps)
F = 9.3574, ndf = 1, ddf = 995, p-value = 0.00228
```

The reweighting does not diminish the association between the treatment and the outcome. Indeed, it is marginally more significant after the reweighting. Alternatively, we can run a generalised linear model.

## 5 Non-response weights

The `twang` package was designed to estimate propensity score weights for the evaluation of treatment effects in observational or quasi-experimental studies. However, we find that the package includes functions and diagnostic tools that are highly valuable for other applications, such as for generating and diagnosing non-response weights for survey nonresponse or study attrition. We now present an example that uses the tools in `twang`. This example uses the subset of the US Sustaining Effects Study data distributed with the HLM software (Bryk, Raudenbush, Congdon, 1996) and also available in the R package `mlmRev`. The data include mathematics test scores for 1721 students in kindergarten to fourth grade. They also include student race (black, Hispanic, or other), gender, an indicator for whether or not the student had been retained in grade, the percent low income students at the school, the school size, the percent of mobile students, the students' grade-levels, student and school IDs, and grades converted to year by centering. The study analysis plans to analyze growth in math achievement from grade 1 to grade 4 using only students with complete data. However, the students with complete data differ from other students. To reduce bias that could potentially result from excluding incomplete cases, our analysis plan is to weight complete cases with nonresponse weights.

The goal of nonresponse weighting is to develop weights for the respondents that make them look like the entire sample — both the respondents and nonrespondents. Since the respondents already look like themselves, the hard part is to figure out how well each respondent represents the nonrespondents. Nonresponse weights equal the reciprocal of the probability of response and are applied only to respondents.

Note that the the probability of response is equivalent to the propensity score if we consider subjects with an observed outcome to be the “treated” group, and those with an unobserved outcome to be the “controls”. We wish to reweight the sample to make it equivalent to the population from which the sample was drawn, so ATE weights are more appropriate in this case. Further, recall that the weights for the treated subjects are  $1/p$  in an ATE analysis. Therefore we can reweight the sample of respondents using the `get.weights()` function.

Before we can generate nonresponse weights, we need to prepare the data using the following commands. First we load the data.

```
> data(egsingle)
```

Next we create a response indicator variable that we can merge onto the student by test score level data. We want to include only students with scores from all of grades 1 to 4. The data include scores from kindergarten (grade = 0) to grade 5 with some students having multiple scores from the same grade. First we keep the unique grades for each student:

```
> tmp <- tapply(egsingle$grade, egsingle$childid, unique)
```

Because students do not all have the same number of score, `tapply()` returns a list with one element per student. Each element contains the unique set of grades observed for each student. We now check for grades in 1 to 4:

```
> tmp <- lapply(tmp, function(x){return(x %in% 1:4)})
```

The list `tmp` now contains a boolean vector for each student, where “TRUE” indicates the grade took on a value in 1 to 4. The sum of this vector for each student determines how many of grades 1 to 4 we observed for him or her.

```
> tmp <- lapply(tmp, sum)
```

A student is a respondent if he or she has scores from all four of grades 1 to 4 or if the value of `tmp` is 4.

```
> tmp <- sapply(tmp, function(x){as.numeric(x == 4)})
```

We create a data frame of response indicators so we can merge them onto the test scores data:

```
> tmp <- data.frame(tmp)
> names(tmp) <- "resp"
> tmp$childid <- row.names(tmp)
```

and merge this back to create a single data frame

```
> egsingle <- merge(egsingle, tmp)
```

Because nonresponse is a student-level variable rather than a student-by-year-level variable we create one record per student.

```
> egsingle.one <- unique(egsingle[, -c(3:6)])
```

We also create a race variable

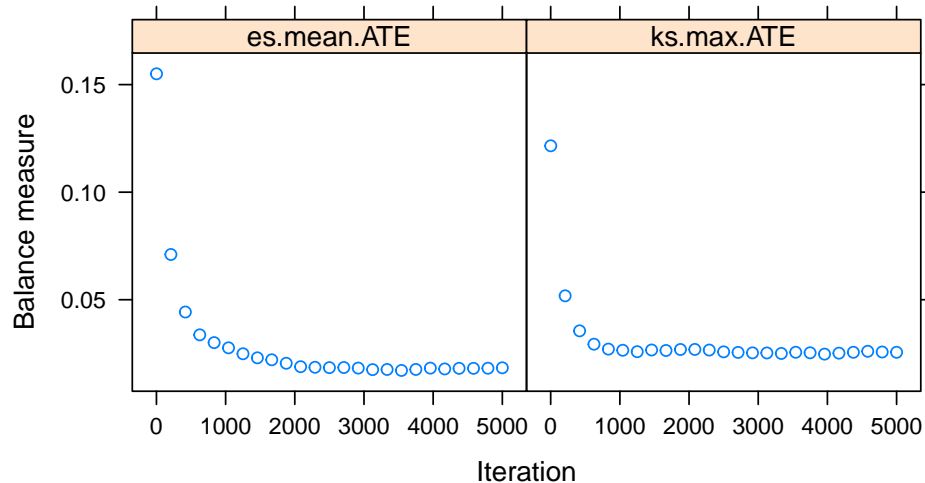
```
> egsingle.one$race <- as.factor(race <- ifelse(egsingle.one$black==1, 1,
+                                             ifelse(egsingle.one$hispanic==1, 2, 3)))
```

As discussed above, to use `ps()` to estimate nonresponse, we need to let respondents be the treatment group by modeling an indicator of response.

```
> egsingle.ps <- ps(resp ~ race + female + size + lowinc + mobility,
+ data=egsingle.one,
+ stop.method=c("es.mean", "ks.max"),
+ n.trees=5000,
+ verbose=FALSE,
+ estimand = "ATE")
```

As in standard propensity score applications, we should check that `n.trees` was set large enough, so that balance would seem not to be improved if more complex models were considered. Recall that for `es.mean.ATE` the measure is the average effect size difference between the two groups and for `ks.max.ATE` the measure is the largest of the KS statistics.

```
> plot(egsingle.ps)
```



By default the balance table generated by `ps()` compares the weighted treatment group (respondents) to the weighted comparison group (nonrespondents) – both groups weighted to equal the overall population. However, the goal is to weight the respondents to match the population not to compare the weighted respondents and nonrespondents. The default balance table may be useful for evaluating the propensity scores, but it does not directly assess the quality of the weights for balancing the weighted respondents with the overall population.

We can “trick” the `dx.wts()` function in the `twang` package into making the desired comparison. We want to compare the weighted respondents to the unweighted full sample. When evaluating ATT weights we compare the weighted comparison group with the unweighted treatment group. If we apply `dx.wts()` to a data set where the “treatment” group is the entire `esingle.one` sample and the “control” group is the `esingle.one` respondents and the weights equal one for every student in the pseudo-treatment group and equal the weights from `ps()` for every student in the pseudo-control group, we can obtain the balance statistics we want.

We begin by setting up the data with the pseudo-treatment and control groups. We add ATE weights from the “ks.max” stopping rule as our nonresponse weights.

```
> egsingle.one$wgt <- get.weights(egsingle.ps, stop.method="ks.max")
```

We now stack the full sample and the respondents. The variable “nr2” is the pseudo-treatment indicator. We set it equal to one for the full sample and 0 for the respondents. Similarly, “wgt2” is the pseudo-ATT weight which is set equal to one for the full sample and equal to the nonresponse weights for the respondents.

```
> egtmp <- rbind(data.frame(egsingle.one, nr2=1, wgt2=1),
+               data.frame(egsingle.one, nr2=0, wgt2=egsingle.one$wgt)[egsingle.one$resp==1,])
```

We now run `dx.wts()` to obtain the balance statistics. Switching to ATT from ATE yields a warning that can be ignored in this case.

```
> egdxwts <- dx.wts(x=egtmp$wgt2,
+                  data=egtmp,
+                  estimand="ATT",
+                  vars=c("race", "female", "size", "lowinc", "mobility"),
+                  treat.var="nr2")

> # pretty.tab<-bal.table(egdxwts)[[2]][,c("tx.mn","ct.mn","std.eff.sz","ks")]
> # names(pretty.tab) <- c("OverallS Sample","Weighted responders","Std ES","KS")
> # xtable(pretty.tab,
> #       caption = "Balance of the nonrespondents and respondents",
> #       label = "tab:balance2",
> #       digits = c(0, 2, 2, 2, 2),
> #       align=c("l","r","r","r","r"))
> bal.table(egdxwts)[[2]]
```

```
tx.mn tx.sd ct.mn ct.sd std.eff.sz stat p ks race:1 0.694 0.461 0.690 0.463 0.011 0.147 0.862 0.005 race:2 0.145
0.352 0.142 0.349 0.010 NA NA 0.003 race:3 0.160 0.367 0.169 0.374 -0.022 NA NA 0.008 female:Female
0.492 0.500 0.486 0.500 0.013 0.092 0.761 0.007 female:Male 0.508 0.500 0.514 0.500 -0.013 NA NA 0.007 size
755.888 314.292 756.535 312.935 -0.002 -0.048 0.961 0.018 lowinc 78.172 26.505 78.507 27.214 -0.013 -0.297
0.767 0.032 mobility 34.588 13.994 34.222 13.697 0.026 0.626 0.532 0.018 ks.pval race:1 0.862 race:2 0.862
race:3 0.862 female:Female 0.761 female:Male 0.761 size 0.991 lowinc 0.631 mobility 0.991
```

The balance is very good. We can use the weighted respondents for analysis. We select only the records with an observed outcome. This will be our analysis sample and the variable “wgt” will contains the nonresponse weights.

```
> egsingle.resp <- merge(subset(egsingle, subset=resp==1),
+                         subset(egsingle.one, subset=resp==1,
+                               select=c(childid, wgt)) )
```

## 6 The details of twang

### 6.1 Propensity scores and weighting

Propensity scores can be used to reweight comparison cases so that the distribution of their features match the distribution of features of the treatment cases, for ATT, or cases from both treatment and control groups to match each other, for ATE (Rosenbaum 1987, Wooldridge 2002, Hirano and Imbens 2001, McCaffrey *et al.* 2004) Let  $f(\mathbf{x}|t = 1)$  be the distribution of features for the treatment cases and  $f(\mathbf{x}|t = 0)$  be the distribution of features for the comparison cases. If treatments were randomized then we would expect these two distributions to be similar. When they differ for ATT we will construct a weight,  $w(\mathbf{x})$ , so that

$$f(\mathbf{x}|t = 1) = w(\mathbf{x})f(\mathbf{x}|t = 0). \quad (2)$$

For example, if  $f(\text{age}=65, \text{sex}=F|t = 1) = 0.10$  and  $f(\text{age}=65, \text{sex}=F|t = 0) = 0.05$  (i.e. 10% of the treatment cases and 5% of the comparison cases are 65 year old females) then we need to give a weight of 2.0 to every 65 year old female in the comparison group so that they have the same representation as in the treatment group. More generally, we can solve (2) for  $w(\mathbf{x})$  and apply Bayes Theorem to the numerator and the denominator to give an expression for the propensity score weight for comparison cases,

$$w(\mathbf{x}) = K \frac{f(t = 1|\mathbf{x})}{f(t = 0|\mathbf{x})} = K \frac{P(t = 1|\mathbf{x})}{1 - P(t = 1|\mathbf{x})}, \quad (3)$$

where  $K$  is a normalization constant that will cancel out in the outcomes analysis. Equation (3) indicates that if we assign a weight to comparison case  $i$  equal to the odds that a case with features  $\mathbf{x}_i$  would be exposed to the treatment, then the distribution of their features would balance. Note that for comparison cases with features that are atypical of treatment cases, the propensity score  $P(t = 1|\mathbf{x})$  would be near 0 and would produce a weight near 0. On the other hand, comparison cases with features typical of the treatment cases would receive larger weights.

For ATE, each group is weighted to match the population. The weight must satisfy:

$$f(\mathbf{x}|t = 1) = w(\mathbf{x})f(\mathbf{x}), \text{ and} \quad (4)$$

$$f(\mathbf{x}|t = 0) = w(\mathbf{x})f(\mathbf{x}), \text{ and} \quad (5)$$

Again using Bayes Theorem we obtain  $w(\mathbf{x}) \propto 1/f(t = 1|\mathbf{x})$  for the treatment group and  $w(\mathbf{x}) \propto 1/f(t = 0|\mathbf{x})$  for the control group.

### 6.2 Estimating the propensity score

In randomized studies  $P(t = 1|\mathbf{x})$  is known and fixed in the study design. In observational studies the propensity score is unknown and must be estimated, but poor estimation of the propensity scores can cause just as much of a problem for estimating treatment effects as poor regression modeling of the outcome. Linear logistic regression is the common method for estimating propensity scores, and can suffice for many problems. Linear logistic regression for propensity scores estimates the log-odds of a case being in the treatment given  $\mathbf{x}$  as

$$\log \frac{P(t = 1|\mathbf{x})}{1 - P(t = 1|\mathbf{x})} = \beta' \mathbf{x} \quad (6)$$

Usually,  $\beta$  is selected to maximize the logistic log-likelihood

$$\ell\beta = \frac{1}{n} \sum_{i=1}^n t_i \beta' \mathbf{x}_i - \log(1 + \exp(\beta' \mathbf{x}_i)) \quad (7)$$

Maximizing (7) provides the maximum likelihood estimates of  $\beta$ . However, in an attempt to remove as much confounding as possible, observational studies often record data on a large number of potential confounders, many of which can be correlated with one another. Standard methods for fitting logistic regression models to such data with the iteratively reweighted least squares algorithm can be statistically and numerically unstable. To improve the propensity score estimates we might also wish to include non-linear effects and interactions in  $\mathbf{x}$ . The inclusion of such terms only increases the instability of the models.

One increasingly popular method for fitting models with numerous correlated variables is the lasso (least absolute subset selection and shrinkage operator) introduced in statistics in Tibshirani (1996). For logistic regression, lasso estimation replaces (7) with a version that penalizes the absolute magnitude of the coefficients

$$\ell\beta = \frac{1}{n} \sum_{i=1}^n t_i \beta' \mathbf{x}_i - \log(1 + \exp(\beta' \mathbf{x}_i)) - \lambda \sum_{j=1}^J |\beta_j| \quad (8)$$

The second term on the right-hand side of the equation is the penalty term since it decreases the overall of  $\ell\beta$  when there are coefficient that are large in absolute value. Setting  $\lambda = 0$  returns the standard (and potentially unstable) logistic regression estimates of  $\beta$ . Setting  $\lambda$  to be very large essentially forces all of the  $\beta_j$  to be equal to 0 (the penalty excludes  $\beta_0$ ). For a fixed value of  $\lambda$  the estimated  $\hat{\beta}$  can have many coefficients exactly equal to 0, not just extremely small but precisely 0, and only the most powerful predictors of  $t$  will be non-zero. As a result the absolute penalty operates as a variable selection penalty. In practice, if we have several predictors of  $t$  that are highly correlated with each other, the lasso tends to include all of them in the model, shrink their coefficients toward 0, and produce a predictive model that utilizes all of the information in the covariates, producing a model with greater out-of-sample predictive performance than models fit using variable subset selection methods.

Our aim is to include as covariates all piecewise constant functions of the potential confounders and their interactions. That is, in  $\mathbf{x}$  we will include indicator functions for continuous variables like  $I(\text{age} < 15)$ ,  $I(\text{age} < 16)$ ,  $\dots$ ,  $I(\text{age} < 90)$ , etc., for categorical variables like  $I(\text{sex} = \text{male})$ ,  $I(\text{prior MI} = \text{TRUE})$ , and interactions among them like  $I(\text{age} < 16)I(\text{sex} = \text{male})I(\text{prior MI} = \text{TRUE})$ . This collection of basis functions spans a plausible set of propensity score functions, are computationally efficient, and are flat at the extremes of  $\mathbf{x}$  reducing the likelihood of propensity score estimates near 0 and 1 that can occur with linear basis functions of  $\mathbf{x}$ . Theoretically with the lasso we can estimate the model in (8), selecting a  $\lambda$  small enough so that it will eliminate most of the irrelevant terms and yield a sparse model with only the most important main effects and interactions. Boosting (Friedman 2001, 2003, Ridgeway 1999) effectively implements this strategy using a computationally efficient method that Efron *et al.* (2004) showed is equivalent to optimizing (8). With boosting it is possible to maximize (8) for a range of values of  $\lambda$  with no additional computational effort than for a specific value of  $\lambda$ . We use boosted logistic regression as implemented in the generalized boosted modeling (gbm) package in R (Ridgeway 2005).

### 6.3 Evaluating the weights

As with regression analyses, propensity score methods cannot adjust for unmeasured covariates that are uncorrelated with the observed covariates. Nonetheless, the quality of the adjustment for the observed covariates achieved by propensity score weighting is easy to evaluate. The estimated propensity score weights should equalize the distributions of the cases' features as in (2). This implies that weighted statistics of the covariates of the comparison group should equal the same statistics for the treatment group. For example, the weighted average of the age of comparison cases should equal the average age of the treatment cases. To

assess the quality of the propensity score weights one could compare a variety of statistics such as means, medians, variances, and Kolmogorov-Smirnov statistics for each covariate as well as interactions. The **twang** package provides both the standardized effect sizes and KS statistics and p-values testing for differences in the means and distributions of the covariates for analysts to use in assessing balance.

## 6.4 Analysis of outcomes

With propensity score analyses the final outcomes analysis is generally straightforward, while the propensity score estimation may require complex modeling. Once we have weights that equalize the distribution of features of treatment and control cases by reweighting. For ATT, we give each treatment case a weight of 1 and each comparison case a weight  $w_i = p(\mathbf{x}_i)/(1 - p(\mathbf{x}_i))$ . To estimate the ATE, we give control cases weight  $w_i = 1/(1 - p(\mathbf{x}_i))$  and we give the treatment cases  $w_i = 1/p(\mathbf{x}_i)$ . We then estimate the treatment effect estimate with a weighted regression model that contains only a treatment indicator. No additional covariates are needed if the weights account for differences in  $\mathbf{x}$ .

A combination of propensity score weighting and covariate adjustment can be useful for several reasons. First, the propensity scores may not have been able to completely balance all of the covariates. The inclusion of these covariates in addition to the treatment indicator in a weighted regression model may correct this if the imbalance is relatively small. Second, in addition to exposure, the relationship between some of the covariates and the outcome may also be of interest. Their inclusion can provide coefficients that can estimate the direction and magnitude of the relationship. Third, as with randomized trials, stratifying on covariates that are highly correlated with the outcome can improve the precision of estimates. Lastly, the some treatment effect estimators that utilize an outcomes regression model and propensity scores are “doubly robust” in the sense that if either the propensity score model is correct or the regression model is correct then the treatment effect estimator will be unbiased (Bang & Robins 2005).

## About This Tutorial

This tutorial and the R package were supported by funding from grant R01DA017507, R01DA015697, and R01DA034065 from the National Institute on Drug Abuse. The overarching goal of this work is to develop statistical methods and tools that will provide addiction health services researchers and others with the tools and training they need to study the effectiveness of treatments using observational data. For more information about **twang** and other causal tools being developed, see [www.rand.org/statistics/twang](http://www.rand.org/statistics/twang).

RAND Social and Economic Well-Being is a division of the RAND Corporation that seeks to actively improve the health and social and economic well-being of populations and communities throughout the world. This research was conducted in the Social and Behavioral Policy Program within RAND Social and Economic Well-Being. The program focuses on such topics as risk factors and prevention programs, social safety net programs and other social supports, poverty, aging, disability, child and youth health and well-being, and quality of life, as well as other policy concerns that are influenced by social and behavioral actions and systems that affect well-being. For more information, email [sbp@rand.org](mailto:sbp@rand.org).

## References

- [1] Bang H. and J. Robins (2005). “Doubly robust estimation in missing data and causal inference models,” *Biometrics* 61:692–972.
- [2] Bland M. (2013). “Do baseline p-values follow a uniform distribution in randomised trials?” *PLoS ONE* 8(10): e76010: 1–5.
- [3] Dehejia, R.H. and S. Wahba (1999). “Causal effects in nonexperimental studies: re-evaluating the evaluation of training programs,” *Journal of the American Statistical Association* 94:1053–1062.



- [4] Efron, B., T. Hastie, I. Johnstone, R. Tibshirani (2004). “Least angle regression,” *Annals of Statistics* 32(2):407–499.
- [5] Friedman, J.H. (2001). “Greedy function approximation: a gradient boosting machine,” *Annals of Statistics* 29(5):1189–1232.
- [6] Friedman, J.H. (2002). “Stochastic gradient boosting,” *Computational Statistics and Data Analysis* 38(4):367–378.
- [7] Friedman, J.H., T. Hastie, R. Tibshirani (2000). “Additive logistic regression: a statistical view of boosting,” *Annals of Statistics* 28(2):337–374.
- [8] Hastie, T., R. Tibshirani, and J. Friedman (2001). *The Elements of Statistical Learning*. Springer-Verlag, New York.
- [9] Helmreich, J.E., and R.M. Pruzek (2009). “PSAgraphics: An R package to support propensity score analysis,” *Journal of Statistical Software* 29(6):1–23.
- [10] Hirano, K. and G. Imbens (2001). “Estimation of causal effects using propensity score weighting: An application to data on right heart catheterization,” *Health Services and Outcomes Research Methodology* 2:259–278.
- [11] Huppler-Hullsiek, K. and T. Louis (2002) “Propensity score modeling strategies for the causal analysis of observational data,” *Biostatistics* 3:179–193.
- [12] Lalonde, R. (1986). “Evaluating the econometric evaluations of training programs with experimental data,” *American Economic Review* 76:604–620.
- [13] Little, R. J. and S. Vartivarian (2004). “Does weighting for nonresponse increase the variance of survey means?” *ASA Proceedings of the Joint Statistical Meetings*, 3897–3904 American Statistical Association (Alexandria, VA) <http://www.bepress.com/cgi/viewcontent.cgi?article=1034&context=umichbiostat>.
- [14] McCaffrey, D., G. Ridgeway, A. Morral (2004). “Propensity score estimation with boosted regression for evaluating adolescent substance abuse treatment,” *Psychological Methods* 9(4):403–425.
- [15] Obenchain, B. (2011). *USPS 1.2 package manual*. <http://cran.r-project.org/web/packages/USPS/USPS.pdf>
- [16] Ridgeway, G. (1999). “The state of boosting,” *Computing Science and Statistics* 31:172–181.
- [17] Ridgeway, G. (2005). *GBM 1.5 package manual*. <http://cran.r-project.org/doc/packages/gbm.pdf>.
- [18] Ridgeway, G. (2006). “Assessing the effect of race bias in post-traffic stop outcomes using propensity scores.” *Journal of Quantitative Criminology* 22(1):1–29.
- [19] Rosenbaum, P. and D. Rubin (1983). “The Central Role of the Propensity Score in Observational Studies for Causal Effects,” *Biometrika* 70(1):41–55.
- [20] Rosenbaum, P. (1987). “Model-based direct adjustment,” *Journal of the American Statistical Association* 82:387–394.
- [21] Tibshirani, R. (1996). “Regression shrinkage and selection via the lasso,” *Journal of the Royal Statistical Society, Series B* 58(1):267–288.
- [22] Wooldridge, J. (2002). *Econometric analysis of cross section and panel data*, MIT Press, Cambridge.