

# Semiparametric Least Squares Inference for Causal Effects with R

Pierre Chausse\*, Mihai Giurcanu†, Marinela Capanu‡, George Luta§

## Abstract

This vignette explains how to use the `causalSLSE` package to estimate causal effects using the semiparametric least squares methods developed by Giurcanu et al. (2023). We describe the classes and methods implemented in the package as well as how they can be used to analyze synthetic and real data.

## 1 Introduction

This document presents the `causalSLSE` package describing the functions implemented in the package. It is intended for users interested in the details about the methods presented in Giurcanu et al. (2023) and how they are implemented.

The general semiparametric additive regression model is

$$\begin{aligned} Y &= \beta_0(1 - Z) + \beta_1 Z + \sum_{l=1}^q f_{l,0}(X_l)(1 - Z) + \sum_{l=1}^q f_{l,1}(X_l)Z + \xi \\ &\equiv \beta_0(1 - Z) + \beta_1 Z + f_0(X)(1 - Z) + f_1(X)Z + \xi, \end{aligned}$$

where  $Y \in \mathbb{R}$  is the response variable,  $Z$  is the treatment indicator defined as  $Z = 1$  for the treated and  $Z = 0$  for the nontreated, and  $X \in \mathbb{R}^q$  is a  $q$ -vector of confounders. We approximate this model by the following regression model:

$$\begin{aligned} Y &= \beta_0(1 - Z) + \beta_1 Z + \sum_{l=1}^q \psi_{l,0}^T U_{l,0}(1 - Z) + \sum_{l=1}^q \psi_{l,1}^T U_{l,1}Z + \zeta \\ &\equiv \beta_0(1 - Z) + \beta_1 Z + \psi_0^T U_0(1 - Z) + \psi_1^T U_1Z + \zeta, \end{aligned}$$

where  $U_{l,k} = u_{l,k}(X_l) = (u_{j,l,k}(X_l) : 1 \leq j \leq p_{l,k}) \in \mathbb{R}^{p_{l,k}}$  is a vector of basis functions corresponding to the  $l^{\text{th}}$  nonparametric component of the  $k^{\text{th}}$  group  $f_{l,k}(X_l)$ ,  $\psi_{l,k} \in \mathbb{R}^{p_{l,k}}$  is an unknown vector of regression coefficients,  $U_k = u_k(X) = (u_{l,k}(X_l) : 1 \leq l \leq q) \in \mathbb{R}^{p_k}$  and  $\psi_k = (\psi_{l,k} : 1 \leq l \leq q) \in \mathbb{R}^{p_k}$ , with  $p_k = \sum_{l=1}^q p_{l,k}$ . In this paper, we propose a data-driven method for selecting the vectors of basis functions  $u_0(X)$  and  $u_1(X)$ . Note that we allow the number of basis functions ( $p_{l,k}$ ) to differ across confounders and groups.

To understand the package, it is important to know how the  $u_{l,k}(X_l)$ 's are defined. For clarity, let's write  $U_{l,k} = u_{l,k}(X_l)$  as  $U = u(X) = (u_j(X) : 1 \leq j \leq p) \in \mathbb{R}^p$ . We just need to keep in mind that it is different for the treated and nontreated groups and also for different confounders. We describe here how to construct the local linear splines for a given confounder  $X$ . To this end, let  $\{\kappa_1, \dots, \kappa_{p-1}\}$  be a set of  $p - 1$  knots

---

\*University of Waterloo, pchausse@uwaterloo.ca

†University of Chicago, giurcanu@uchicago.edu

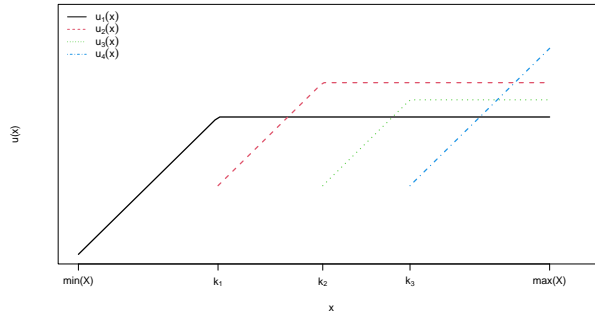
‡Memorial Sloan Kettering Cancer Center, capanu@mskcc.org

§Georgetown University, George.Luta@georgetown.edu

strictly inside the support of  $X$  satisfying  $\kappa_1 < \kappa_2 < \dots < \kappa_{p-1}$ . In the case of local linear splines described in the paper, we have:

$$\begin{aligned} u_1(x) &= xI(x \leq \kappa_1) + \kappa_1 I(x > \kappa_1) \\ u_j(x) &= (x - \kappa_{j-1})I(\kappa_{j-1} \leq x \leq \kappa_j) + (\kappa_j - \kappa_{j-1})I(x > \kappa_j), \quad 2 \leq j \leq p-1 \\ u_p(x) &= (x - \kappa_{p-1})I(x > \kappa_{p-1}) \end{aligned}$$

Therefore, if the number of knots is equal to 1, we only have two local linear splines. Since the knots must be strictly inside the support of  $X$ , for any categorical variable with two levels, the number of knots must be equal to zero. In this case,  $u(x) = x$ . For general ordinal variables, the number of knots cannot exceed the number of levels minus two. The following illustrates the local splines when the number of knots is equal to 3:



Note that for the sample regression, the knots of  $X_l$  for group  $k$ ,  $l = 1, \dots, q$ , must be strictly inside the sample range of  $(X_{i,l} : 1 \leq i \leq n, Z_i = k) \in \mathbb{R}^{n_k}$ , where  $n_k$  is the sample size in group  $k$ , instead of inside the support of  $X_l$ .

## 2 The causalSLSE package

### 2.1 Setting up the Model

The first step in using the package is to define the causal model. The model contains the information about the outcome ( $Y$ ), the treatment indicator ( $Z$ ), the confounders ( $X$ ) and their knots  $(\kappa_{l,k})$ . This is the starting point before applying any basis selection method. To illustrate how to use the package, we are using the dataset from Lalonde (1986). The dataset, called **nsw**, contains some continuous and categorical variables, so we can illustrate how knots are selected initially. The dataset is included in the **causalSLSE** package.

```
library(causalSLSE)
data(nsw)
```

The outcome is the real income in 1978 (**re78**) and the purpose is to estimate the causal effect of a training program (**treat**) on **re78**. The dataset includes the continuous covariates age (**age**), education (**ed**), the 1975 real income (**re75**), and binary variables (**black**, **hisp**, **married** and **nodeg**). We start by considering the variables **age**, **re75**, **ed**, and **married**. To setup the model, we simply run the following command:

```
model1 <- cslseModel(re78 ~ treat | ~ age + re75 + ed + married, data = nsw)
```

The left of  $|$  is designated for the formula linking the outcome (**re78**) and the treatment indicator (**treat**). The confounders are entered after  $|$  as a formula without a dependent variable. This formula works similarly to formulas in the **lm** function. For example, we can add interactions, transformations of the variables, etc. The following is an example:

```
model0 <- cslseModel(re78 ~ treat | ~ age + I(age^2) + re75 + ed * married,
                    data = nsw)
```

This will create the vector of covariates {age, age<sup>2</sup>, re75, ed, married, ed×married}. Note that adding age<sup>2</sup> is not recommended since we already model nonlinearities via basis functions. This is presented to give an example of what can be added to the formula. The function `cslseModel` creates an object of class `cslseModel` with its own print method, which will be presented later.

The following subsections explain the arguments of `cslseModel`.

### 2.1.1 The starting knots

By default, the function automatically generates knots for each variable based on the following procedure applied separately for the treated and nontreated. The term **sample size** refers either to the number of observations in the treated or nontreated group.

1. The starting number of knots is a function of the sample size and is determined by the argument **nbasis**, a function of one argument, the sample size. The floor value of what the function returns is the number of basis functions. The starting number of knots is therefore equal to the **floor** of what the function returns minus 1 (or 0 if the function returns a value strictly less than 1). The default function is `function(n) n^0.3`. For example, if the total sample size is 500, with 200 treated and 300 nontreated, the starting number of knots in the treated and nontreated groups are respectively equal to `3=floor(200^0.3)-1` and `4=floor(300^0.3)-1`, respectively. It is possible to have a number of knots that does not depend on the sample size. All we need is to set the argument **nbasis** to a function that returns an integer, e.g., `nbasis=function(n) 4` for 4 basis functions or 3 knots.
2. Let  $(p - 1)$  be the number of knots determined in the previous step. The default knots are obtained by computing  $p + 1$  quantiles of  $X$  for equally spaced probabilities from 0 to 1, and by dropping the first and last quantiles. For example, if the number of knots is 3, then the initial knots are given by quantiles for the probabilities 0.25, 0.5 and 0.75.
3. We drop any duplicated knots and any knots equal to either the max or the min of  $X$ . If the resulting number of knots is equal to 0, the vector of knots is set to `NULL`. When the knots is equal to `NULL` for a variable  $X$ , it means that  $u(x) = x$ .

The last step implies that the number of knots for all categorical variables with two levels is equal to 0. For nominal variables with a small number of levels, the number of knots, a subset of the levels, may be smaller than the ones defined by **nbasis**. For example, when the number of levels for a nominal variable is 3, the number of knots cannot exceed 1.

We can inspect the knots of the current model as follows. Note that each object in the package is S3-class, so the elements can be accessed using the operator `$`. The elements `knots0` and `knots1` are the list of knots for the nontreated and treated groups, respectively. For example, in our case the initial knots for the treated are:

```
model1$knots1

## $age
## 20% 40% 60% 80%
## 19 22 25 28
##
## $re75
##      40%      60%      80%
## 357.9499 1961.8640 5588.6640
##
## $ed
## 20% 40% 60% 80%
## 9 10 11 12
##
```

```
## $married
## NULL
```

We see that it is set to NULL for `married`, because it is a binary variable. The sample size for the treated is 297. Given the default `nbasis`, it implies a number of starting knots equal to  $4 = \text{floor}(297^{0.3}) - 1$ . This is the number of knots we have for `ed` and `age`. However, the number of knots for `re75` is 3. The reason is that `re75` contains a large fraction of zeros. Since the 20<sup>th</sup> percentile is equal to 0 and 0 is also the minimum value of `ed75`, it is dropped (the `type` argument of the `quantile` function is the same as it is implemented in the package). This can be seen as follows:

```
quantile(nsw[nsw$treat==1, 're75'], c(.2, .4, .6, .8), type=1)
```

```
##          20%          40%          60%          80%
##    0.0000  357.9499 1961.8640 5588.6640
```

By printing the object, we see a description of the model. It includes the list of variables with a positive number of knots and with no knots.

```
model1

## Semiparametric LSE Model
## *****
##
## Number of treated:  297
## Number of nontreated: 425
## Number of missing values:  0
## Selection method: Default
## Covariates approximated by SLSE:
##  age, re75, ed
## Covariates not approximated by SLSE:
##  married
```

Note that the selection method is set to `Default`. We refer to this method when the knots are automatically selected by the method described above. Later in the document, we will present methods for selecting a subset of `Default` using `SLSE`, which stands for Semiparametric Least Squares Estimator (Giurcanu (2016)).

We have also included, in the package, the simulated dataset `simDat4`, which contains special types of covariates. It helps to further illustrate how the knots are determined. The dataset contains a continuous variable `X1` with a large proportion of zeros, the categorical variable `X2` with 3 levels, an ordinal variable `X3` with 3 levels, and a binary variable `X4`. The levels for `X2` are {"first", "second", "third"} and for `X3` the levels are {1, 2, 3}.

```
data(simDat4)
model2 <- cslseModel(Y ~ Z | ~ X1 + X2 + X3 + X4, data = simDat4)
model2$knots0
```

```
## $X1
##          40%          60%          80%
##  0.2531388  2.9118507 12.1110772
##
## $X2second
## NULL
##
## $X2third
## NULL
##
## $X3
## 40%
```

```
## 2
##
## $X4
## NULL
```

Character-type variables are automatically converted into factors. It is also possible to define a numerical variable like `X3` as a factor by using the function `as.factor` in the formula. We see that the 2 binary variables `X2second` and `X2third` are created and `X2first` is omitted to avoid multicollinearity. For the binary variable `X4`, the number of knots is set to 0. For the ordinal variable `X3`, the number of knots is set to 1 because the min and max values 1 and 3 cannot be selected.

### 2.1.2 Setting the knots manually

The user has control over the selection of knots through the arguments `knots0` and `knots1`. When the arguments are missing (the default), all knots are set automatically as described above. One way to set the number of knots to 0 for all variables in a given group is to set the argument to `NULL`. For example, the number of knots is equal to 0 for all variables of the treated group using the following command:

```
cslseModel(re78 ~ treat | ~ age + re75 + ed + married, data = nsw,
           knots1 = NULL)
```

```
## Semiparametric LSE Model
## *****
##
## Number of treated: 297
## Number of nontreated: 425
## Number of missing values: 0
## Selection method: User Based
## Covariates approximated by SLSE:
## Treated: None
## Nontreated: age, re75, ed
## Covariates not approximated by SLSE:
## Treated: age, re75, ed, married
## Nontreated: married
```

Notice that the selection method is defined as “User Based” whenever the knots are provided manually by the user. Also, the print method shows the lists of covariates by group only when they differ, which is the case here. The other option is to provide a list of knots. For each element, we have three options:

- NA: The knots are set automatically for this variable only.
- NULL: The number of knots is set to 0 for this variable only.
- A numeric vector: The vector cannot contain missing or duplicated values and must be strictly inside the range of the variable for the group.

In the following, we describe all possible formats for the list of knots.

1. An unnamed list of length equal to the number of covariates. In that case, the knots must be defined in the same order of covariates implied by the formula.

Suppose we want to set for the nontreated group an automatic selection for `age`, no knots for `ed`, the knots `{1000,5000,10000}` for `re75`, and the knots to be automatically selected for the treated group. We proceed as follows. Note that setting the value to `NA` or `NULL` has the same effect for the binary variable `married`. In the following, the argument `knots=TRUE` is added to the `print` method to only print the knots.

```
model <- cslseModel(re78 ~ treat | ~ age + re75 + ed + married, data = nsw,
                   knots0 = list(NA, c(1000,5000,10000), NULL, NA))
```

```
print(model, knots = TRUE)
```

```
## Semiparametric LSE Model
## *****
##
## Selection method: User Based
## Lists of knots for the treated
## *****
## age:
## 20% 40% 60% 80%
## 19 22 25 28
## re75:
## 40% 60% 80%
## 357.9499 1961.8640 5588.6640
## ed:
## 20% 40% 60% 80%
## 9 10 11 12
## married:
## None
##
## Lists of knots for the nontreated
## *****
## age:
## 16.66667% 33.33333% 50% 66.66667% 83.33333%
## 18 20 23 26 30
## re75:
## k1 k2 k3
## 1000 5000 10000
## ed:
## None
## married:
## None
```

2. A named list of length equal to the number of covariates. In that case, the order of the list of variables does not matter. The `cslseModel` function will automatically reorder the variables to match the order implied by the formula. The names must match perfectly the covariate names generated by R.

In the following example, we want to add the interaction between `ed` and `age`. We want the same set of knots as in the previous example and no knots for the interaction term. The name of the interaction depends on how we enter it in the formula. For example, it is “age:ed” if we enter `age*ed` in the formula and “ed:age” if we enter `ed*age`. For factors, the names depend on which binary variable is omitted. Using the above example with the `simDat4` model, if we interact `X2` and `X4` by adding `X2*X4` to the formula, the names of the interaction terms are “X2second:X4” and “X2third:X4”. When we are uncertain about the names, we can print the knots of a model with the default sets of knots. In the following, we change the order of variables to show that the order does not matter.

```
knots <- list(married = NA, ed = NULL, 'age:ed' = NULL,
             re75 = c(1000,5000,10000), age = NA)
model <- cslseModel(re78 ~ treat | ~ age * ed + re75 + married, data = nsw,
                  knots0 = knots)
model$knots0
```

```
## $age
## 16.66667% 33.33333% 50% 66.66667% 83.33333%
## 18 20 23 26 30
```

```
##
## $ed
## NULL
##
## $re75
##      k1      k2      k3
## 1000  5000 10000
##
## $married
## NULL
##
## $`age:ed`
## NULL
```

3. A named list of length strictly less than the number of covariates. The names of the selected covariates must match perfectly the names generated by R and the order does not matter. This is particularly useful when the number of covariates is large.

If we consider the previous example, the knots are set manually only for `ed`, `ed:age` and `re75`. By default, all names not included in the list of knots are set to NA. Therefore, we can create the same model from the previous example as follows:

```
knots <- list(ed = NULL, 'age:ed' = NULL, re75 = c(1000,5000,10000))
model <- cslseModel(re78 ~ treat | ~ age * ed + re75 + married, data = nsw,
                    knots0 = knots)
model$knots0
```

```
## $age
## 16.66667% 33.33333%      50% 66.66667% 83.33333%
##      18      20      23      26      30
##
## $ed
## NULL
##
## $re75
##      k1      k2      k3
## 1000  5000 10000
##
## $married
## NULL
##
## $`age:ed`
## NULL
```

Note that the previous case offers an easy way of setting the number of knots to 0 for a subset of covariates. For example, suppose we want to add more interaction terms and set the knots to 0 for all of them. We can proceed as follows.

```
knots <- list('age:ed' = NULL, 'ed:re75' = NULL, 'ed:married' = NULL)
model <- cslseModel(re78 ~ treat | ~ age * ed + re75 * ed + married * ed, data = nsw,
                    knots0 = knots, knots1 = knots)
model
```

```
## Semiparametric LSE Model
## *****
##
## Number of treated: 297
```

```
## Number of nontreated: 425
## Number of missing values: 0
## Selection method: User Based
## Covariates approximated by SLSE:
## age, ed, re75
## Covariates not approximated by SLSE:
## married, age:ed, ed:re75, ed:married
```

Note also that `cslseModel` deals with interaction terms as any other variable. For example, `ed:black` is like a continuous variable with a large proportion of zeros. The following shows the default selected knots for `ed:black`.

```
model <- cslseModel(re78 ~ treat | ~ age + ed * black, data = nsw)
model$knots0[["ed:black"]]
```

```
## 33.33333%      50% 66.66667%
##          9       10       11
```

## 2.2 Estimating the model

Given the set of knots from the model object, the estimation is just a least squares method applied to the extended set of covariates defined as the local linear splines corresponding to the set of knots. The regression model is given by:

$$Y = \beta_0(1 - Z) + \beta_1 Z + \psi_0^T U_0(1 - Z) + \psi_1^T U_1 Z + \zeta,$$

where  $U_0 = u_0(X)$  and  $U_1 = u_1(X)$  are defined above (which depend on the knots of the model). The method that estimates the model is `estSLSE` which has three arguments, but two of them are mainly used internally by other functions. We present them in case they are needed. The arguments are:

- `model`: A model created by the function `cslseModel`.
- `w0` and `w1`: lists of integers to select knots for the nontreated and treated respectively. For example, suppose we have 2 covariates with 5 knots each. If we want to estimate the model with only the first knot for the first covariate and knots 3 and 5 for the second, we set `w0` to `list(1L, c(3L, 5L))`. By default they are set to `NULL` and all the knots from the model are used.

We illustrate the use of `estSLSE` with a simple model containing 2 covariates and one knot per variable.

```
model <- cslseModel(re78 ~ treat | ~ age + married, data = nsw,
                  nbasis = function(n) 2)
print(model, knots = TRUE)
```

```
## Semiparametric LSE Model
## *****
##
## Selection method: Default
## Lists of knots for the treated
## *****
## age:
## 50%
## 23
## married:
## None
##
## Lists of knots for the nontreated
## *****
```



```
## age:
## 50%
## 23
## married:
## None

fit <- estSLSE(model)
fit

## Semiparametric LSE
## *****
## Selection method: Default
##
## factor(treat)0 factor(treat)1 Xf0age_1 Xf0age_2 Xf0married
## 4558.28061 3754.98326 27.79868 -12.51415 -115.81593
## Xf1age_1 Xf1age_2 Xf1married
## 89.25358 22.22331 1435.28205
```

The object of class `slseFit` returned by `estSLSE` has its own print method that returns the coefficient estimates. A more detailed presentation of the results can be obtained using the `summary` method. The following is an example with one knot per eligible variable.

```
summary(fit)

## Semiparametric LSE
## *****
## Selection method: Default
##
## Estimate Std. Error t value Pr(>|t|)
## factor(treat)0 4558.28 2739.40 1.664 0.0961 .
## factor(treat)1 3754.98 3704.37 1.014 0.3107
## Xf0age_1 27.80 136.61 0.203 0.8387
## Xf0age_2 -12.51 56.06 -0.223 0.8234
## Xf0married -115.82 795.35 -0.146 0.8842
## Xf1age_1 89.25 185.53 0.481 0.6305
## Xf1age_2 22.22 82.52 0.269 0.7877
## Xf1married 1435.28 1226.68 1.170 0.2420
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Multiple R-squared: 0.009618, Adjusted R-squared: -9.119e-05
```

For example, the coefficient of `Xf0age_1` is the effect of age for the control on `re78` when `age`  $\leq$  23 and `Xf0age_2` is the effect when `age`  $>$  23. Note that the  $R^2$  and adjusted  $R^2$  are different from what we obtain using the summary of the `lm` object:

```
summary(fit$lm.out)[c("r.squared", "adj.r.squared")]

## $r.squared
## [1] 0.4379272
##
## $adj.r.squared
## [1] 0.4316295
```

This is because our model does not contain an intercept and the  $R^2$  is computed differently for models without an intercept. The definition of the  $R^2$  used by R is the following (RSS means residual sum of squares):

$$R^2 = 1 - \frac{\text{RSS for the model with the regressors}}{\text{RSS for the model without the regressors}}.$$

In a model with an intercept, the residual of the model without the regressors is  $Y_i - \bar{Y}$ , but it is equal to  $Y_i$  when the model does not have an intercept. As a result, the  $R^2$  with and without an intercept are respectively

$$R^2_{with} = 1 - \frac{\sum_{i=1}^n \hat{e}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

and

$$R^2_{without} = 1 - \frac{\sum_{i=1}^n \hat{e}_i^2}{\sum_{i=1}^n Y_i^2}.$$

## 2.3 The predict and plot methods

The `predict` method is very similar to the `predict.lm` method. We use the same arguments: `object`, `interval`, `se.fit`, `newdata` and `level`. The difference is that it returns the predicted outcome for the treated and nontreated separately, and the argument `vcov` provides a way of changing how the least squares covariance matrix is computed. By default, it is computed using `vcovHC` from the `sandwich` package (Zeileis (2006)). The function returns a list of 2 elements, `treated` and `nontreated`. By default (`se.fit=FALSE` and `interval="none"`), each element contains a vector of predictions. Here is an example with the previously fitted model `fit`:

```
predict(fit,
        newdata = data.frame(treat = c(1,1,0,0), age = 20:23, married = 1))

## $treated
## [1] 6975.337 7064.591
##
## $nontreated
## [1] 5054.036 5081.834
```

If `interval` is set to “confidence”, but `$se.fit` remains equal to `FALSE`, each element contains a matrix containing the prediction, and the lower and upper confidence limits, with the confidence level determined by the argument `level` (set to 0.95 by default). Here is an example with the same fitted model:

```
predict(fit,
        newdata = data.frame(treat = c(1,1,0,0), age = 20:23, married = 1),
        interval = "confidence")

## $treated
##      fit      lower      upper
## 1 6975.337 4646.673 9304.001
## 2 7064.591 4741.653 9387.528
##
## $nontreated
##      fit      lower      upper
## 3 5054.036 3574.096 6533.975
## 4 5081.834 3544.849 6618.820
```

If `se.fit` is set to `TRUE`, each element, `treated` or `nontreated`, is a list with the elements `pr`, containing the predictions, and `se.fit`, containing the standard errors. In the following, we only show the result for the `treated`:

```
predict(fit,
        newdata = data.frame(treat = c(1,1,0,0), age = 20:23, married = 1),
        se.fit = TRUE)$treated
```

```
## $fit
## [1] 6975.337 7064.591
##
## $se.fit
##      1      2
## 1188.116 1185.194
```

The `predict` method is called by the `plot` method to visually assess the predicted outcome for the treated and nontreated with respect to a given covariate, controlling for the other variables in the model. The arguments of the `plot` method are:

- **x**: An object of class `slseFit`.
- **y**: An alias for **which** for compatibility with the generic `plot` function.
- **which**: covariate to plot against the outcome variable. It could be an integer (the position of the covariate) or a character (the name of the covariate)
- **interval**: The type of confidence interval to display. The default is “none”. The alternative is “confidence”.
- **level**: The confidence level when **interval**=“confidence”. The default is 0.95.
- **fixedCov0** and **fixedCov1**: Optional named lists of fixed values for some or all other covariates in each group. The values of the covariates not specified are determined by the argument **FUN**. By default, **fixedCov1** is equal to **fixedCov0**, so it is not necessary to set a value for **fixedCov1** if we want the same covariates to be fixed to the same values in both groups.
- **legendPos**: The position of the legend. The default is “topright”.
- **vcov.**: An optional function to compute the estimated matrix of covariance of the least squares estimators. This argument only affects the confidence intervals. The default is `vcovHC` with **type**=“HC3”.
- **col0**, **col1**, **lty0**, **lty1**: The line colors and shapes for the nontreated and treated. The defaults are **col0**=1 (black), **col1**=2 (red), **lty0**=1 (solid) and **lty1**=2 (dashed).
- **add.**: Should the curves be added to an existing plot? The default is **FALSE**.
- **addToLegend**: An optional character string to add to the legend next to “treated” and “nontreated”.
- **cex**: The font size for the legend. The default is 1.
- **ylim**, **xlim**: optional ranges for the y-axis and x-axis.
- **addPoints**: Should we include the scatterplot of the outcome and covariate to the graph? The default is **FALSE**.
- **FUN**: A function to determine how the other covariates are fixed. The default is **mean**. Note that the function is applied to each group separately.
- **main**: An optional title to replace the default one.
- **plot**: By default, the method produces a graph. Alternatively, we can set this argument to **FALSE** and it returns one **data.frame** per group with the variable selected by **which** and the prediction.
- **...**: Other arguments are passed to the **vcov.** function. For example, it is possible to change the type of `vcovHC` from the default HC3 to any available methods included in the **sandwich** package.

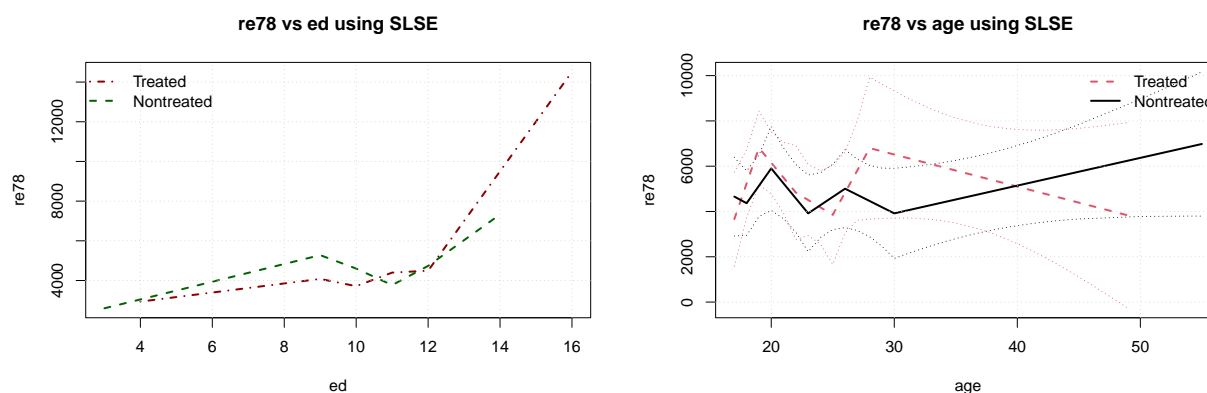
In the following, we illustrate some examples.

### 2.3.0.1 Example 1: Consider the model:

```
model1 <- cslseModel(re78 ~ treat | ~ age + re75 + ed + married, data = nsw)
fit1 <- estSLSE(model1)
```

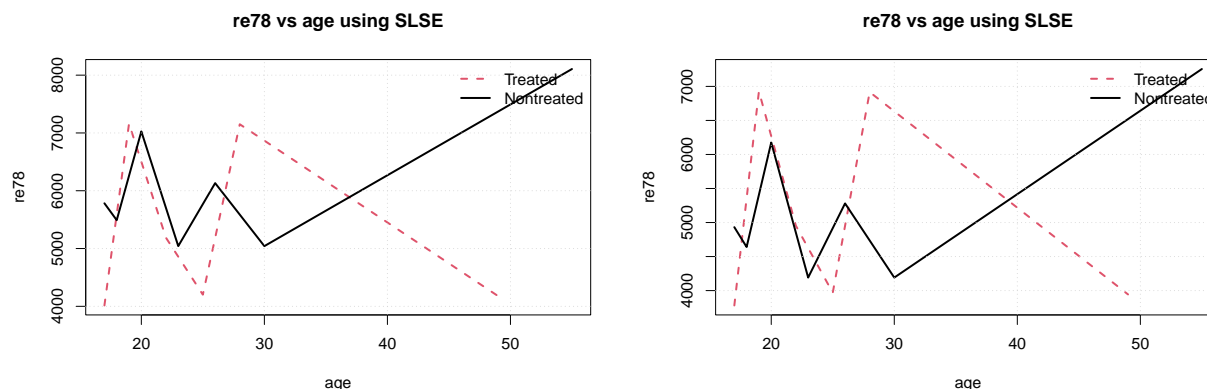
Suppose we want to compare the predicted income between the two treatment groups with respect to age or education, holding the other variables fixed to their group means (the default). The following are two examples with some of the default arguments modified. Note that `vcov.lm` is used in the first plot function and `vcovHC` (the default) of type `HC1` in the second plot.

```
library(sandwich)
plot(fit1, "ed", col0 = "darkgreen", col1 = "darkred", lty0 = 2, lty1 = 4,
     legendPos = "topleft", vcov. = vcov)
plot(fit1, "age", interval = 'confidence', level = 0.9, type = "HC1")
```



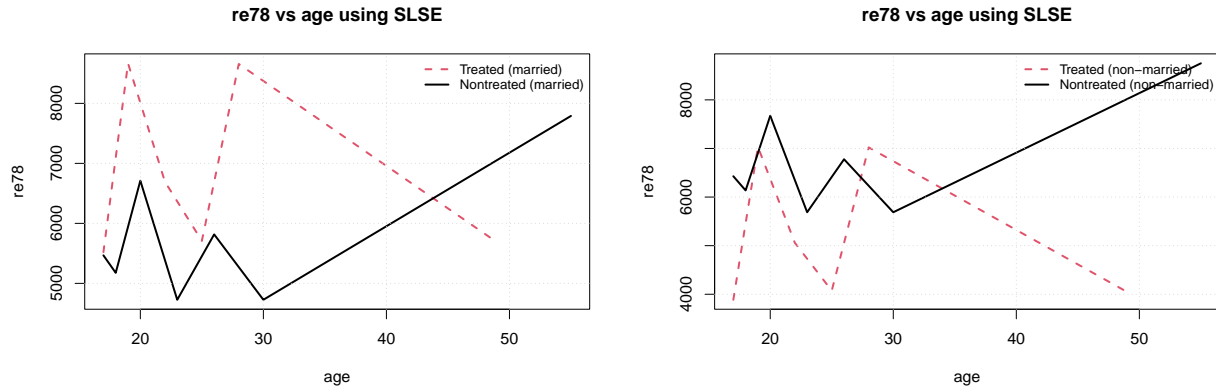
**2.3.0.2 Example 2:** If we want to fix the other covariates using another function, we can change the argument `FUN`. The new function must be a function of one argument. For example, if we want to fix the other covariates to their group medians, we set `FUN` to `median` (no quotes). We proceed the same way for any function that requires only one argument. If the function requires more than one argument, we have to create a new function. For example, if we want to fix them to their 20% group empirical quantiles, we can set the argument to `function(x) quantile(x, .20)`. The following illustrates the two cases:

```
plot(fit1, "age", FUN = median)
plot(fit1, "age", FUN = function(x) quantile(x, 0.20))
```



**2.3.0.3 Example 3:** It is also possible to set some of the other covariates to a specific value by changing the argument `fixedCov0` and `fixedCov1`. By default, `fixedCov1` is equal to `fixedCov0`, so if we want to fixed the same covariates to the same values in both groups, we only need to set `fixedCov0`. The argument must be a named list with the names corresponding to the variables you want to fix. You can also add a description to the legend with the argument `addToLegend`.

```
plot(fit1, "age", fixedCov0 = list(married = 1, re75 = 10000),
     addToLegend = "married", cex = 0.8)
plot(fit1, "age", fixedCov0 = list(married = 0, re75 = 10000),
     addToLegend = "non-married", cex = 0.8)
```



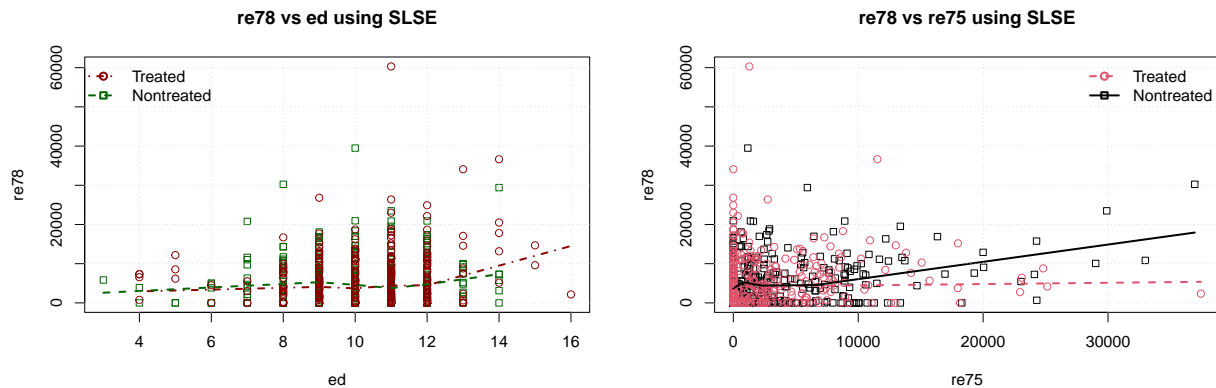
**2.3.0.4 Example 4:** To better compare the two groups, it is also possible to have them plotted on the same graph by setting the argument `add.` to `TRUE`. We just need to adjust some of the arguments to better distinguish the different curves. In the following example, we set the colors and line shapes to different values and change the position of the legend in the second plot function.

```
plot(fit1, "age", fixedCov0 = list(married = 1, re75 = 10000),
     addToLegend = "married", cex = 0.8, ylim. = c(3000,10000))
plot(fit1, "age", fixedCov0 = list(married = 0, re75 = 10000),
     addToLegend = "non-married", cex = 0.8, legendPos = 'topleft',
     col0 = "darkgreen", col1 = "darkred", lty0 = 4, lty1 = 5,
     add. = TRUE)
```



**2.3.0.5 Example 5:** Finally, it is also possible to add the observed points to the graph.

```
plot(fit1, "ed", col0 = "darkgreen", col1 = "darkred", lty0 = 2, lty1 = 4,
     legendPos = "topleft", addPoints = TRUE)
plot(fit1, "re75", addPoints = TRUE)
```



### 2.3.1 Factors, interactions and functions of covariates

The package allows some of the covariates to be factors, functions of other covariates or interactions. For example, the dataset `simDat4` includes one factor, `X2`, with levels equal to “first”, “second” and “third”. We can include this covariate directly to the list of covariates. For example,

```
data(simDat4)
mod <- cslseModel(Y ~ Z | ~ X1 + X2 + X4, data = simDat4)
mod
```

```
## Semiparametric LSE Model
## *****
##
## Number of treated: 246
## Number of nontreated: 254
## Number of missing values: 0
## Selection method: Default
## Covariates approximated by SLSE:
## X1
## Covariates not approximated by SLSE:
## X2second, X2third, X4
```

We see that R has created 2 binary variables, one for `X2="second"` and one for `X2="third"`. These two variables are automatically included in the group of covariates not approximated by SLSE because they are binary variables like `X4`. If we want to plot `Y` against `X1`, the binary variables `X2second`, `X2third` and `X4` are fixed to their group averages which, in case of binary variables, represent the proportions of ones in each group.

For interaction terms or functions of covariates, `FUN` is applied to the functions of covariates. This is how we have to proceed to obtain the average prediction in regression models. For example, if we interact `X2` and `X4`, we obtain:

```
data(simDat4)
mod <- cslseModel(Y ~ Z | ~ X1 + X2 * X4, data = simDat4)
mod
```

```
## Semiparametric LSE Model
## *****
```

```
##
## Number of treated: 246
## Number of nontreated: 254
## Number of missing values: 0
## Selection method: Default
## Covariates approximated by SLSE:
## X1
## Covariates not approximated by SLSE:
## X2second, X2third, X4, X2second:X4, X2third:X4
```

In this case, when `FUN=mean`, `X2second:X4` is replaced by the proportion of ones in `X2second×X4` for each group. It is not replaced by the proportion of ones in `X2second` times the proportion of ones in `X4`. The same applies to functions of covariates. For functions of covariates, which can be defined in the formula using a built-in function like `log` or using the identity function `I()` (e.g. we can interact `X1` and `X4` by using `I(X1*X4)`), `FUN` is applied to the function (e.g. the average `log(X)` or the average `I(X1*X4)`).

To fix a factor to a specific level, we just set its value to the `fixedCov0` and `fixedCov1` arguments. In the following example, we fix `X2` to “first”, so `X2second` and `X2third` are set to 0.

```
fit <- estSLSE(mod)
plot(fit, "X1", fixedCov0 = list(X2 = "first"))
```

Note that if a function of covariates (or an interaction) involves the covariate we want to plot the outcome against, we factorize the covariate out, apply `FUN` to the remaining of the function and add the covariate back. For example, if we interact `X1` with `X4` and `FUN=mean`, `X1:X4` is replaced by `X1` times the proportion of ones in `X4` for each group.

## 2.4 Optimal selection of the knots

We have implemented two methods for selecting the knots: the backward SLSE (BSLSE) and the forward SLSE (FSLSE) methods. For each method, we have 3 criteria: the p-value threshold (PVT), the Akaike Information criterion (AIC), and the Bayesian Information criterion (BIC). The two selection methods can be summarized as follows:

### BSLSE:

1. We estimate the model with all knots included in the model.
2. For each knot, we test if the slopes of the basis functions adjacent to the knot are the same, and return the p-value.
3. The knots are selected using one of the following criteria
  - **PVT**: We remove all knots with a p-value greater than a specified threshold.
  - **AIC** or **BIC**: We order the p-values in descending order. Then, going from the largest to the smallest, we remove the knot associated with the p-value one by one, estimate the model and return the information criterion. We keep the model with the smallest information criterion.

### FSLSE:

1. We estimate the model by including a subset of the knots, one variable at the time. When we test a knot for one covariate, the number of knots is set to 0 for all other variables.
2. For each knot, we test if the adjacent slopes to the knot is the same, and return the p-value. The set of knots used for each test depends on the following:
  - Variables with 1 knot: we return the p-value of the test of equality of the slopes adjacent to the knot.

- Variables with 2 knots: we include the two knots and return the p-values of the test of equality of the slopes adjacent to each knot.
  - Variables with  $p$  knots ( $p > 2$ ): We test the equality of the slopes adjacent to knot  $i$ , for  $i = 1, \dots, p$ , using the sets of knots  $\{1, 2\}$ ,  $\{1, 2, 3\}$ ,  $\{2, 3, 4\}$ ,  $\dots$ ,  $\{p-2, p-1, p\}$  and  $\{p-1, p\}$  respectively.
3. The knots are selected using one of the following criteria
- **PVT**: We remove all knots with a p-value greater than a specified threshold.
  - **AIC** or **BIC**: We order the p-values in ascending order. Then, starting with a model with no knots and going from the smallest to the highest highest p-value, we add the knot associated with the smallest remaining p-value one by one, estimate the model and return the information criterion. We keep the model with the smallest information criterion.

The knot selection is done using the `selSLSE` method. The arguments are:

- **model**: An object of class `cslseModel`.
- **selType**: This is the selection method. We have the choice between “FSLSE” and “BSLSE” (the default).
- **selCrit**: This is the criterion used by the selection method. We have the choice between “AIC” (the default), “BIC” or “PVT”.
- **pvalT**: This is a function that returns the p-value threshold. It is a function of one argument, the average number of basis functions per covariate. The default is `function(p) 1/log(p)` and it is applied to each group separately. Therefore, the threshold may be different for the treated and non-treated. It is also possible to set it to a fix threshold. For example, `function(p) 0.20` sets the threshold to 0.2. This argument affects the result only when **method** is set to “PVT”.
- **vcov.**: An optional function to compute the least squares standard errors. By default, the p-values are computed using the `vcovHC` method from the `sandwich` package with `type="HC3"`.
- **...**: This is used to pass arguments to the `vcov.` function.

The function returns a model of class `cslseModel` with the optimal selection of knots. For example, we can compare the starting knots of `model1`, with the model selected by the default arguments.

```
print(model1, knots = TRUE)

## Semiparametric LSE Model
## *****
##
## Selection method: Default
## Lists of knots for the treated
## *****
## age:
## 20% 40% 60% 80%
## 19 22 25 28
## re75:
## 40% 60% 80%
## 357.9499 1961.8640 5588.6640
## ed:
## 20% 40% 60% 80%
## 9 10 11 12
## married:
## None
##
## Lists of knots for the nontreated
## *****
## age:
## 16.66667% 33.33333% 50% 66.66667% 83.33333%
## 18 20 23 26 30
```



```
## re75:
##      50% 66.66667% 83.33333%
## 823.2544 2292.1710 6567.3290
## ed:
## 16.66667% 33.33333% 66.66667% 83.33333%
##      9      10      11      12
## married:
## None

model2 <- selSLSE(model1)
print(model2, knots = TRUE)

## Semiparametric LSE Model
## *****
##
## Selection method: BSLSE
## Criterion: AIC
##
## Lists of knots for the treated
## *****
## age:
## 20% 60% 80%
## 19 25 28
## re75:
## None
## ed:
## 80%
## 12
## married:
## None
##
## Lists of knots for the nontreated
## *****
## age:
## 33.33333%      50% 66.66667% 83.33333%
##      20      23      26      30
## re75:
##      50% 83.33333%
## 823.2544 6567.3290
## ed:
## 16.66667% 66.66667%
##      9      11
## married:
## None
```

For example, the BSLSE-AIC method has removed all knots from `re75` for the treated and kept two knots for the nontreated. The print method indicates which method was used to select the knots. In the following example, we see BSLSE as selection method and BIC as criterion. Note that the BIC selects 0 knots for all covariates.

```
model3 <- selSLSE(model1, selType = "BSLSE", selCrit = "BIC")
model3

## Semiparametric LSE Model
## *****
##
## Number of treated: 297
## Number of nontreated: 425
## Number of missing values: 0
## Selection method: BSLSE
## Criterion: BIC
##
## Covariates approximated by SLSE:
## None
## Covariates not approximated by SLSE:
```

```
## age, re75, ed, married
```

Since the `selSLSE` method returns a new model, we can apply the `estSLSE` to it:

```
estSLSE(selSLSE(model1, selType = "FSLSE", selCrit = "BIC"))

## Semiparametric LSE
## *****
## Selection method: FSLSE
## Criterion: BIC
##
## factor(treat)0 factor(treat)1 Xf0age Xf0re75 Xf0ed
## 4.825878e+03 -3.889679e+02 -2.010566e+01 2.982477e-01 2.500219e+00
## Xf0married Xf1age Xf1re75 Xf1ed Xf1married
## -1.094084e+03 4.105403e+01 2.676162e-02 4.849161e+02 1.417291e+03
```

## 2.5 The causalSLSE method for slseFit objects

The regression model estimated by `estSLSE`, as described in the introduction, can be written as

$$Y_i = \beta_0(1 - Z_i) + \beta_1 Z_i + \psi'_0 U_{i,0}(1 - Z_i) + \psi'_1 U_{i,1} Z_i + \zeta_i \text{ for } i = 1, \dots, n.$$

Let  $\hat{\beta}_0$ ,  $\hat{\beta}_1$ ,  $\hat{\psi}_0$  and  $\hat{\psi}_1$  be the least squares estimates of the regression parameters. Then, the SLSE average causal effect (ACE), causal effect on the treated (ACT) and causal effect on the non-treated (ACN) are defined respectively as follows:

$$\begin{aligned} \text{ACE} &= \hat{\beta}_1 - \hat{\beta}_0 + \hat{\psi}'_1 \overline{U_1} - \hat{\psi}'_0 \overline{U_0} \\ \text{ACT} &= \hat{\beta}_1 - \hat{\beta}_0 + \hat{\psi}'_1 \overline{U_1 Z} - \hat{\psi}'_0 \overline{U_0 Z} \\ \text{ACN} &= \hat{\beta}_1 - \hat{\beta}_0 + \hat{\psi}'_1 \overline{U_1(1 - Z)} - \hat{\psi}'_0 \overline{U_0(1 - Z)}, \end{aligned}$$

where

$$\begin{aligned} \overline{U_j} &= \frac{1}{n} \sum_{i=1}^n u_j(X_i), \text{ for } j=0,1 \\ \overline{U_j Z} &= \frac{1}{n_1} \sum_{i=1}^n u_j(X_i) Z_i, \text{ for } j=0,1 \\ \overline{U_j(1 - Z)} &= \frac{1}{n_0} \sum_{i=1}^n u_j(X_i) (1 - Z_i), \text{ for } j=0,1 \end{aligned}$$

and  $n_0$  and  $n_1$  are the sample size in the nontreated and treated groups. The method `causalSLSE` estimates the causal effects from `slseFit` objects using the knots included in the estimated model. The arguments of the method are:

- **object**: An object of class `slseFit`.
- **seType**: The method to compute the standard errors of the causality measures. By default, they are computed using an analytic expression derived in the paper. Alternatively, we can set the argument to “lm” and use the least squares standard errors based on the asymptotic properties.
- **causal**: What causality measure should the function compute? We have the choice between “All” (the default), “ACT”, “ACE” or “ACN”.

- **vcov.**: An alternative function used to compute the covariance matrix of the least squares estimates. By default, **vcovHC** is used with **type="HC3"**.
- **...**: This is used to pass arguments to the **vcov.** function.

In the following example, we estimate the causal effect with the initial knots (without selection).

```
model1 <- cslseModel(re78 ~ treat | ~ age + re75 + ed + married, data=nsw)
fit1 <- estSLSE(model1)
causalSLSE(fit1)
```

```
## Causal Effect using Semiparametric LSE
## *****
## Selection method: Default
##
## ACE = 814.3083
## ACT = 831.8856
## ACN = 802.0249
```

We see that the selection method used to select the knots are set to SLSE. This is explained in the section “Setting up the Model”. The method returns an object of class **cslse**. We see above what its **print** method returns. The following shows its **summary** method:

```
ce <- causalSLSE(fit1)
summary(ce)
```

```
## Causal Effect using Semiparametric LSE
## *****
## Selection method: Default
##      Estimate Std. Error t value Pr(>|t|)
## ACE      814.3      482.1   1.689   0.0912 .
## ACT      831.9      499.5   1.665   0.0958 .
## ACN      802.0      498.9   1.608   0.1079
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

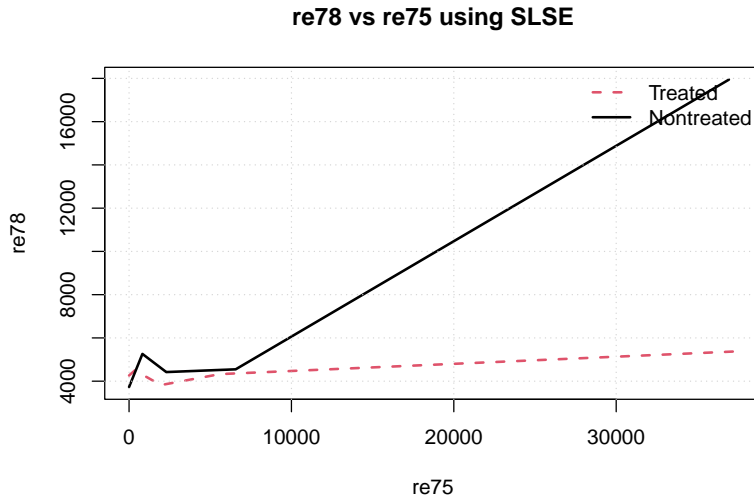
By default, the standard errors are computed using an analytic expression derived in the paper. In the following, we estimate the standard errors using the HC3 type of heteroskedasticity robust standard errors, which is the default when **seType="lm"**.

```
ce2 <- causalSLSE(fit1, seType="lm")
summary(ce2)
```

```
## Causal Effect using Semiparametric LSE
## *****
## Selection method: Default
##      Estimate Std. Error t value Pr(>|t|)
## ACE      814.3      506.1   1.609   0.108
## ACT      831.9      527.4   1.577   0.115
## ACN      802.0      514.2   1.560   0.119
```

The object **cslse** inherits from the class **slseFit**, so we can apply the **plot** (or the **predict**) method directly on this object.

```
plot(ce2, "re75")
```



### 2.5.1 The extract method

The package comes with an `extract` method for objects of class `cslse`, which is a required method for creating Latex tables using the `texreg` package (Leifeld (2013)). For example, we can compare different methods in a single table.

```
library(texreg)
c1 <- causalSLSE(fit1)
fit2 <- estSLSE(selSLSE(model1, selType="BSLSE"))
fit3 <- estSLSE(selSLSE(model1, selType="FSLSE"))
c2 <- causalSLSE(fit2)
c3 <- causalSLSE(fit3)
texreg(list(SLSE=c1, BSLSE=c2, FSLSE=c3), table=FALSE, digits=4)
```

	SLSE	BSLSE	FSLSE
ACE	814.3083 (482.1393)	818.1598 (482.8785)	824.4901 (481.8267)
ACT	831.8856 (499.4948)	837.0768 (501.3497)	852.4659 (496.6795)
ACN	802.0249 (498.8671)	804.9401 (491.0229)	804.9401 (490.4101)
Num. knots (Nontreated)	12	8	6
Num. knots (Treated)	11	4	4
Num. covariates	4	4	4
Num. obs. (Nontreated)	425	425	425
Num. obs. (Treated)	297	297	297
$R^2$	0.0869	0.0852	0.0840
$R^2_{adj}$	0.0445	0.0577	0.0592

\*\*\* $p < 0.001$ ; \*\* $p < 0.01$ ; \* $p < 0.05$

The option `table=FALSE`, from the `texreg` package, is used to remove the Latex floating table environment. With this option, the table appears right after the code instead of being placed somewhere else by Latex. The arguments of the `extract` methods, which control what is printed and can be modified through the `texreg` function, are:

- **include.nobs**: Include the number of observations. The default is `TRUE`.
- **include.nknots**: Include the number of knots. The default is `TRUE`.
- **include.rsquared**: Include the  $R^2$ . The default is `TRUE`.
- **include.adjrsquared**: Include the adjusted  $R^2$ . The default is `TRUE`.

- **which:** Which causal effects should be printed? The options are “ALL” (the default), “ACE”, “ACT”, “ACN”, “ACE-ACT”, “ACE-ACN” or “ACT-ACN”.

Here is one example on how to change some arguments:

```
texreg(list(SLSE=c1, BSLSE=c2, FSLSE=c3), table=FALSE,
        which="ACE-ACT", include.adjrsquared=FALSE)
```

	SLSE	BSLSE	FSLSE
ACE	814.31 (482.14)	818.16 (482.88)	824.49 (481.83)
ACT	831.89 (499.49)	837.08 (501.35)	852.47 (496.68)
Num. knots (Nontreated)	12	8	6
Num. knots (Treated)	11	4	4
Num. covariates	4	4	4
Num. obs. (Nontreated)	425	425	425
Num. obs. (Treated)	297	297	297
R <sup>2</sup>	0.09	0.09	0.08

\*\*\* $p < 0.001$ ; \*\* $p < 0.01$ ; \* $p < 0.05$

## 2.6 The causalSLSE method for cslseModel objects

When applied directly to `cslseModel` objects, the `causalSLSE` method offers the possibility to select the knots and estimate the causal effects all at once. The method also returns an object of class `cslse`. The arguments are the same as the method for `slseFit` objects, plus the necessary arguments for the knots selection. The following are the arguments not already defined for objects of class `slseFit`. The details of these arguments are presented in the section Optimal selection of knots.

- **object:** An object of class `cslseModel`.
- **selType:** This is the selection method. We have the choice between “SLSE” (the default), “FSLSE” and “BSLSE”. The SLSE method performs no selection, so all knots from the model are kept.
- **selCrit:** This is the criterion used by the selection method when `selType` is set to “FSLSE” or “BSLSE”. The default is “AIC”.
- **pvalT:** This is a function that returns the p-value threshold. We explained this argument when we presented the `selSLSE` method.

For example, we can generate the previous table as follows.

```
c1 <- causalSLSE(model1, selType="SLSE")
c2 <- causalSLSE(model1, selType="BSLSE")
c3 <- causalSLSE(model1, selType="FSLSE")
texreg(list(SLSE=c1, BSLSE=c2, FSLSE=c3), table=FALSE, digits=4)
```

	SLSE	BSLSE	FSLSE
ACE	814.3083 (482.1393)	818.1598 (482.8785)	824.4901 (481.8267)
ACT	831.8856 (499.4948)	837.0768 (501.3497)	852.4659 (496.6795)
ACN	802.0249 (498.8671)	804.9401 (491.0229)	804.9401 (490.4101)
Num. knots (Nontreated)	12	8	6
Num. knots (Treated)	11	4	4
Num. covariates	4	4	4
Num. obs. (Nontreated)	425	425	425
Num. obs. (Treated)	297	297	297
R <sup>2</sup>	0.0869	0.0852	0.0840
R <sup>2</sup> <sub>adj</sub>	0.0445	0.0577	0.0592

\*\*\* $p < 0.001$ ; \*\* $p < 0.01$ ; \* $p < 0.05$

## 2.7 The causalSLSE method for formula objects

This last method, offers an alternative way of estimating the causal effects. It allows the estimation in one step without having to first create a model. The arguments are the same as of the `cslseModel` function and the `causalSLSE` method for `cslseModel` objects. It creates the model, selects the knots and estimates the causal effects in one step. For example, we can create the previous table as follows:

```
c1 <- causalSLSE(re78 ~ treat | ~ age + re75 + ed + married, data=nsw,
  selType="SLSE")
c2 <- causalSLSE(re78 ~ treat | ~ age + re75 + ed + married, data=nsw,
  selType="BSLSE")
c3 <- causalSLSE(re78 ~ treat | ~ age + re75 + ed + married, data=nsw,
  selType="FSLSE")
texreg(list(SLSE=c1, BSLSE=c2, FSLSE=c3), table=FALSE, digits=4)
```

	SLSE	BSLSE	FSLSE
ACE	814.3083 (482.1393)	818.1598 (482.8785)	824.4901 (481.8267)
ACT	831.8856 (499.4948)	837.0768 (501.3497)	852.4659 (496.6795)
ACN	802.0249 (498.8671)	804.9401 (491.0229)	804.9401 (490.4101)
Num. knots (Nontreated)	12	8	6
Num. knots (Treated)	11	4	4
Num. covariates	4	4	4
Num. obs. (Nontreated)	425	425	425
Num. obs. (Treated)	297	297	297
R <sup>2</sup>	0.0869	0.0852	0.0840
R <sub>adj</sub> <sup>2</sup>	0.0445	0.0577	0.0592

\*\*\*  $p < 0.001$ ; \*\*  $p < 0.01$ ; \*  $p < 0.05$

Note that this method calls `cslseModel`, `selSLSE`, `estSLSE` and the method `causalSLSE` for `slseFit` objects sequentially. It is easier to simply work with this method, but manually going through all steps may be beneficial to better understand the procedure. Also, it is more convenient to work with a model when we want to compare the different selection methods, or if we want to compare estimations with different standard errors.

## 2.8 A simulated data set from Model 1

In the package, the data set `datSim1` is generated using the following data generating process with a sample size of 300.

$$\begin{aligned}
 Y(0) &= 1 + X + X^2 + \epsilon(0) \\
 Y(1) &= 1 - 2X + \epsilon(1) \\
 Z &= \text{Bernoulli}[\Lambda(1 + X)] \\
 Y &= Y(1)Z + Y(0)(1 - Z)
 \end{aligned}$$

where  $X$ ,  $\epsilon(0)$  and  $\epsilon(1)$  are independent standard normal and  $\Lambda(x)$  is the CDF of the standard logistic distribution. The causal effects ACE, ACT and ACN are approximately equal to -1, -1.6903 and 0.5867 (estimated using a sample size of  $10^7$ ). We can start by building starting model:

```
data(simDat1)
mod <- cslseModel(Y ~ Z | ~ X, data = simDat1)
```

Then we can compare three different methods:

```
c1 <- causalSLSE(mod, selType = "SLSE")
c2 <- causalSLSE(mod, selType = "BSLSE", selCrit = "BIC")
```

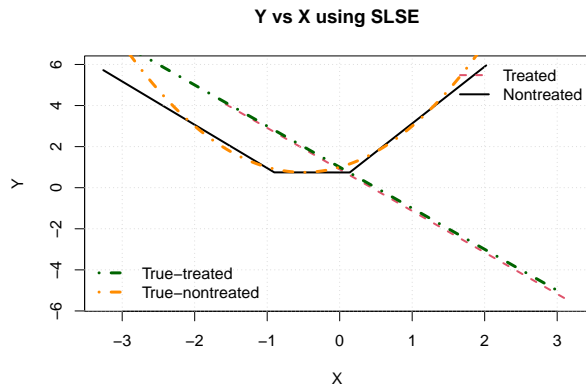
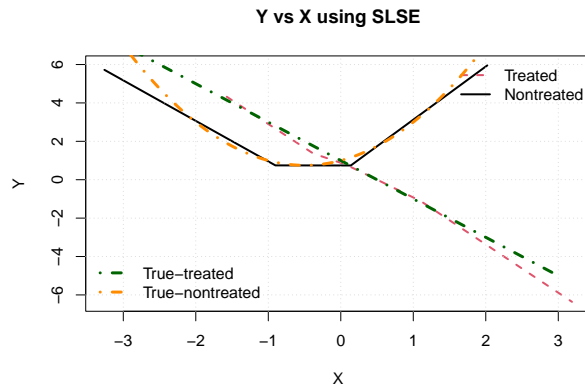
```
c3 <- causalSLSE(mod, selType = "FSLSE", selCrit = "BIC")
texreg(list(SLSE = c1, BSLSE = c2, FSLSE = c3), table = FALSE, digits = 4)
```

	SLSE	BSLSE	FSLSE
ACE	-1.4396*** (0.2614)	-1.4530*** (0.2605)	-1.4533*** (0.2599)
ACT	-1.9316*** (0.3030)	-1.9316*** (0.3024)	-1.9320*** (0.3030)
ACN	-0.0865 (0.3263)	-0.1369 (0.3224)	-0.1369 (0.3224)
Num. knots (Nontreated)	2	2	1
Num. knots (Treated)	4	0	0
Num. covariates	1	1	1
Num. obs. (Nontreated)	80	80	80
Num. obs. (Treated)	220	220	220
R <sup>2</sup>	0.7434	0.7386	0.7303
R <sup>2</sup> <sub>adj</sub>	0.7354	0.7342	0.7266

\*\*\* $p < 0.001$ ; \*\* $p < 0.01$ ; \* $p < 0.05$

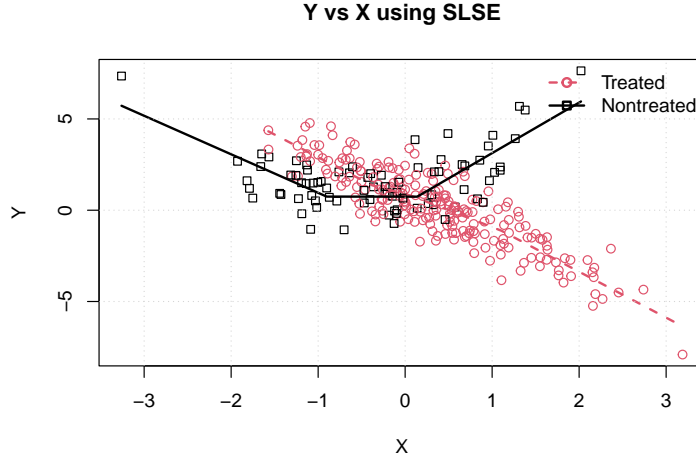
We see that both selection methods choose to assign 0 knots for the treated group, which is not surprising since the true  $f_1(x)$  is linear. We can compare the different fits (we ignore the FSLSE because the selected knots are the same).

```
plot(c1, "X")
curve(1 - 2 * x, -3, 3, col = "darkgreen", lty = 4, lwd = 3, add = TRUE)
curve(1 + x + x^2, -3, 3, col = "darkorange", lty = 4, lwd = 3, add = TRUE)
legend("bottomleft", c("True-treated", "True-nontreated"),
      col=c("darkgreen", "darkorange"), lty = 4, lwd = 3, bty = 'n')
plot(c2, "X")
curve(1 - 2 * x, -3, 3, col="darkgreen", lty = 4, lwd = 3, add = TRUE)
curve(1 + x + x^2, -3, 3, col = "darkorange", lty = 4, lwd = 3, add = TRUE)
legend("bottomleft", c("True-treated", "True-nontreated"),
      col = c("darkgreen", "darkorange"), lty = 4, lwd = 3, bty = 'n')
```



We see that the piecewise polynomials are very close to the true  $f_1(x)$  and  $f_2(x)$ . We can see from the following graph how the lines are fit through the observations by group.

```
plot(c1, "X", addPoints=TRUE)
```



## 2.9 A simulated data set from Model 2

The dataset `datSim2` is a change point regression model (with unknown location of change points) defined as follows:

$$\begin{aligned}
 Y(0) &= (1 + X)I(X \leq -1) + (-1 - X)I(X > -1) + \epsilon(0) \\
 Y(1) &= (1 - 2X)I(X \leq 0) + (1 + 2X)I(X > 0) + \epsilon(1) \\
 Z &= \text{Bernoulli}[\Lambda(1 + X)] \\
 Y &= Y(1)Z + Y(0)(1 - Z)
 \end{aligned}$$

where  $I(A)$  is the indicator function equal to 1 if  $A$  is true, and  $X$ ,  $\epsilon(0)$  and  $\epsilon(1)$  are independent standard normal. The causal effects ACE, ACT and ACN are approximately equal to 3.763, 3.858 and 3.545 (estimated with a sample size of  $10^7$ ). We can compare the SLSE, BSLSE-AIC and BSLSE-BIC.

```

data(simDat2)
mod <- cslseModel(Y~Z | ~X, data=simDat2)

c1 <- causalSLSE(mod, selType = "SLSE")
c2 <- causalSLSE(mod, selType = "BSLSE", selCrit = "BIC")
c3 <- causalSLSE(mod, selType = "BSLSE", selCrit = "AIC")
texreg(list(SLSE = c1, BSLSE.BIC = c2, BSLSE.AIC = c3), table = FALSE, digits = 4)

```

	SLSE	BSLSE.BIC	BSLSE.AIC
ACE	3.9290*** (0.1703)	3.9201*** (0.1717)	3.9201*** (0.1717)
ACT	3.9552*** (0.1891)	3.9404*** (0.1904)	3.9404*** (0.1904)
ACN	3.8670*** (0.2371)	3.8721*** (0.2362)	3.8721*** (0.2362)
Num. knots (Nontreated)	2	1	1
Num. knots (Treated)	3	2	2
Num. covariates	1	1	1
Num. obs. (Nontreated)	89	89	89
Num. obs. (Treated)	211	211	211
R <sup>2</sup>	0.7833	0.7829	0.7829
R <sup>2</sup> <sub>adj</sub>	0.7774	0.7784	0.7784

\*\*\* $p < 0.001$ ; \*\* $p < 0.01$ ; \* $p < 0.05$

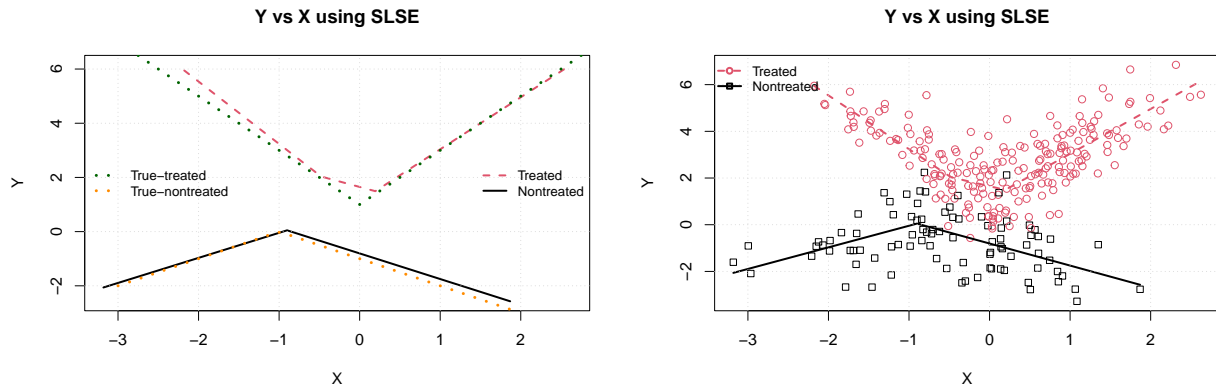
The following shows the fit of BSLSE-AIC with the true  $f_1(x)$  and  $f_0(x)$ , and the observations.



```

plot(c2, "X", legendPos = "right", cex = .8)
curve((1 - 2 * x) * (x <= 0) + (1 + 2 * x) * (x > 0), -3, 3,
      col = "darkgreen", lty = 3, lwd = 3, add = TRUE)
curve((1 + x) * (x <= -1) + (-1 - x) * (x > -1),
      -3, 3, col = "darkorange", lty = 3, lwd = 3, add = TRUE)
legend("left", c("True-treated", "True-nontreated"),
      col = c("darkgreen", "darkorange"), lty = 3, lwd = 3, bty = 'n', cex = .8)
plot(c2, "X", addPoints = TRUE, legendPos = "topleft", cex = .8)

```



## 2.10 A simulated data set from Model 3

The data set `datSim3` is generated from model with multiple confounders defined as follows:

$$\begin{aligned}
 Y(0) &= [1 + X_1 + X_1^2] + [(1 + X_2)I(X_2 \leq -1) + (-1 - X_2)I(X_2 > -1)] + \epsilon(0) \\
 Y(1) &= [1 - 2X_1] + [(1 - 2X_2)I(X_2 \leq 0) + (1 + 2X_2)I(X_2 > 0)] + \epsilon(1) \\
 Z &= \text{Bernoulli}[\Lambda(1 + X_1 + X_2)] \\
 Y &= Y(1)Z + Y(0)(1 - Z),
 \end{aligned}$$

where  $X_1$ ,  $X_2$ ,  $\epsilon(0)$  and  $\epsilon(1)$  are independent standard normal. The causal effects ACE, ACT and ACN are approximately equal to 2.762, 2.204 and 3.922 (estimated with a sample size of  $10^7$ ). We can compare the SLSE, FSLSE with AIC and FSLSE with BIC.

```

data(simDat3)
mod <- cslseModel(Y ~ Z | ~ X1 + X2, data = simDat3)

c1 <- causalSLSE(mod, selType = "SLSE")
c2 <- causalSLSE(mod, selType = "FSLSE", selCrit = "BIC")
c3 <- causalSLSE(mod, selType = "FSLSE", selCrit = "AIC")
texreg(list(SLSE = c1, FSLSE.BIC = c2, FSLSE.AIC = c3), table = FALSE, digits = 4)

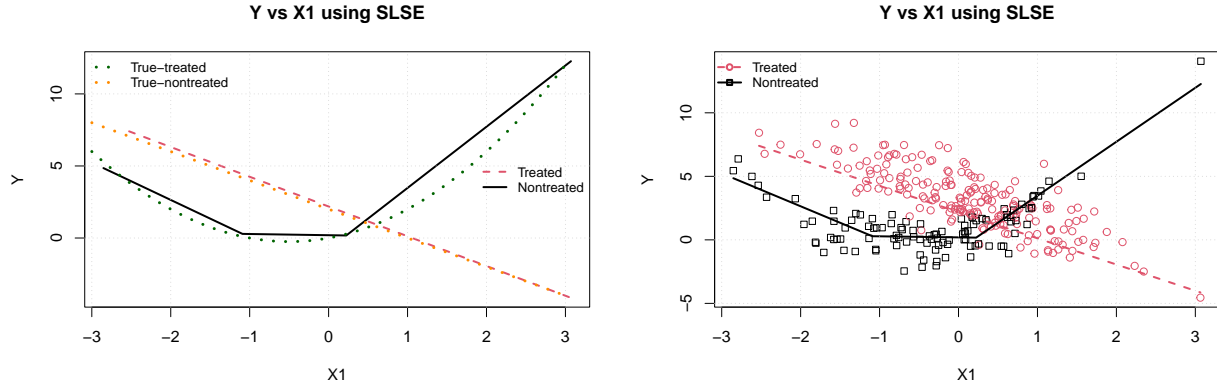
```

	SLSE	FSLSE.BIC	FSLSE.AIC
ACE	2.4699*** (0.2684)	2.4866*** (0.2675)	2.4725*** (0.2684)
ACT	2.0653*** (0.3397)	2.0688*** (0.3380)	2.0688*** (0.3402)
ACN	3.2323*** (0.3445)	3.2739*** (0.3425)	3.2334*** (0.3436)
Num. knots (Nontreated)	6	5	5
Num. knots (Treated)	6	3	4
Num. covariates	2	2	2
Num. obs. (Nontreated)	104	104	104
Num. obs. (Treated)	196	196	196
$R^2$	0.8630	0.8614	0.8625
$R^2_{adj}$	0.8547	0.8551	0.8558

\*\*\* $p < 0.001$ ; \*\* $p < 0.01$ ; \* $p < 0.05$

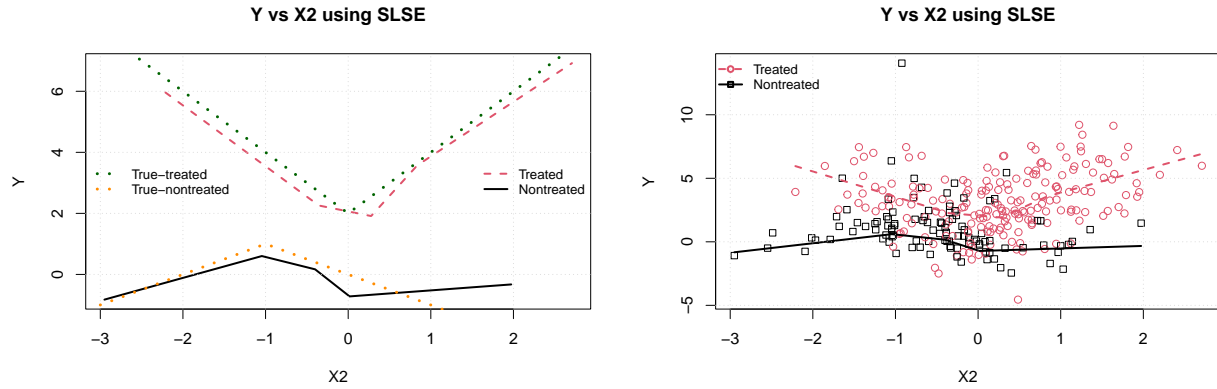
To illustrate the method, since we have two covariates, we need to plot the outcome against one covariate holding the other fixed. The default is to fix it to its sample mean. For the true curve, we fix it to its population mean, which is 0. We first look at the outcome against  $X_1$ . By fixing  $X_2$  to 0, the true curve is  $X_1 + X_1^2$  for the untreated and  $2 - 2X_1$  for the treated. The following graphs show how the FSLSE-BIC method fits the curves.

```
plot(c2, "X1", legendPos = "right", cex = .8)
curve(x + x^2, -3, 3, col = "darkgreen", lty = 3, lwd = 3, add = TRUE)
curve(2 - 2 * x, -3, 3, col = "darkorange", lty = 3, lwd = 3, add = TRUE)
legend("topleft", c("True-treated", "True-nontreated"),
      col = c("darkgreen", "darkorange"), lty = 3, lwd = 3, bty = 'n', cex = .8)
plot(c2, "X1", addPoints = TRUE, legendPos = "topleft", cex = .8)
```



If we fix  $X_1$  to 0, the true curve is  $1 + [(1 + X_2)I(X_2 \leq -1) + (-1 - X_2)I(X_2 > -1)]$  for the nontreated and  $1 + [(1 - 2X_2)I(X_2 \leq 0) + (1 + 2X_2)I(X_2 > 0)]$  for the treated. The following graphs illustrates how these curves are approximated by FSLSE-AIC.

```
plot(c2, "X2", legendPos = "right", cex = .8)
curve(1 + (1 - 2 * x) * (x <= 0) + (1 + 2 * x) * (x > 0), -3, 3,
      col = "darkgreen", lty = 3, lwd = 3, add = TRUE)
curve(1 + (1 + x) * (x <= -1) + (-1 - x) * (x > -1),
      -3, 3, col = "darkorange", lty = 3, lwd = 3, add = TRUE)
legend("left", c("True-treated", "True-nontreated"),
      col = c("darkgreen", "darkorange"), lty = 3, lwd = 3, bty = 'n', cex = .8)
plot(c2, "X2", addPoints = TRUE, legendPos = "topleft", cex = .8)
```



## 2.11 A simulated data set with interactions

The data set `datSim5` is generated using the following data generating process with a sample size of 300.

$$\begin{aligned}
 Y(0) &= [1 + X_1 + X_1^2] + [(1 + X_2)I(X_2 \leq -1) + (-1 - X_2)I(X_2 > -1)] \\
 &\quad + [1 + X_1X_2 + (X_1X_2)^2] + \epsilon(0) \\
 Y(1) &= [1 - 2X_1] + [(1 - 2X_2)I(X_2 \leq 0) + (1 + 2X_2)I(X_2 > 0)] \\
 &\quad + [1 - 2X_1X_2] + \epsilon(1) \\
 Z &= \text{Bernoulli}[\Lambda(1 + X_1 + X_2 + X_1X_2)] \\
 Y &= Y(1)Z + Y(0)(1 - Z),
 \end{aligned}$$

where  $X_1$ ,  $X_2$ ,  $e$  and  $u$  are independent standard normal. The causal effects ACE, ACT and ACN are approximately equal to 1.763, 0.998 and 3.194 (estimated with a sample size of  $10^7$ ). We can compare the SLSE, FSLSE-AIC and FSLSE-BIC.

```

data(simDat5)
mod <- cslseModel(Y ~ Z | ~ X1 * X2, data = simDat5)

c1 <- causalSLSE(mod, selType = "SLSE")
c2 <- causalSLSE(mod, selType = "FSLSE", selCrit = "BIC")
c3 <- causalSLSE(mod, selType = "FSLSE", selCrit = "AIC")
texreg(list(SLSE = c1, FSLSE.BIC = c2, FSLSE.AIC = c3), table = FALSE, digits = 4)

```

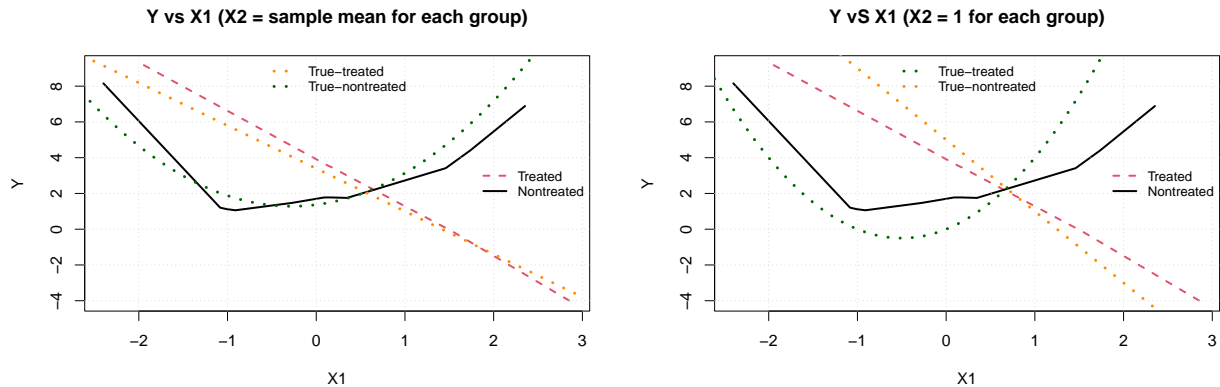
	SLSE	FSLSE.BIC	FSLSE.AIC
ACE	1.7990*** (0.3566)	1.7797*** (0.3615)	1.7744*** (0.3613)
ACT	1.2582** (0.4722)	1.2091* (0.4796)	1.2091* (0.4803)
ACN	2.8183*** (0.4402)	2.8550*** (0.4400)	2.8399*** (0.4378)
Num. knots (Nontreated)	9	8	8
Num. knots (Treated)	9	5	6
Num. covariates	3	3	3
Num. obs. (Nontreated)	104	104	104
Num. obs. (Treated)	196	196	196
R <sup>2</sup>	0.8909	0.8879	0.8894
R <sup>2</sup> <sub>adj</sub>	0.8809	0.8799	0.8811

\*\*\* $p < 0.001$ ; \*\* $p < 0.01$ ; \* $p < 0.05$

In the case of multiple covariates with interactions, the shape of the fitted outcome with respect to one covariate depends on the value of the other covariates. Without interaction, changing the value of the other

covariates only shifts the fitted line without changing its shape. The following graphs compare the estimated relationship between  $Y$  and  $X_1$  for  $X_2$  equal to the group means (left graph) and 1 (right graph). Using a sample of  $10^7$ , we obtain that  $E(X_2|Z = 1)$  and  $E(X_2|Z = 0)$  are approximately equal to 0.1982 and -0.3698, respectively. Therefore, the true curves are  $(1.3698 + 0.6302x + 1.1368x^2)$  for the nontreated and  $(3.3964 - 2.3964x)$  for the treated. If  $X_2 = 1$ , the true curves become  $2x + 2x^2$  for the treated and  $(5 - 4x)$  for the nontreated.

```
x20 <- mean(subset(simDat5, Z == 0)$X2)
x21 <- mean(subset(simDat5, Z == 1)$X2)
plot(c2, "X1", fixedCov0 = list(X2 = x20), fixedCov1 = list(X2 = x21),
     legendPos = "right", cex = .8,
     main="Y vs X1 (X2 = sample mean for each group)")
curve(1.3698 + 0.6302 * x + 1.1368 * x^2, -3, 3,
      col = "darkgreen", lty = 3, lwd = 3, add = TRUE)
curve(3.3964 - 2.3964 * x, -3, 3, col = "darkorange", lty = 3, lwd = 3, add = TRUE)
legend("top", c("True-treated", "True-nontreated"),
      col=c("darkorange", "darkgreen"), lty = 3, lwd = 3, bty = 'n', cex = .8)
plot(c2, "X1", fixedCov0 = list(X2 = 1), legendPos = "right", cex = .8,
     main="Y vs X1 (X2 = 1 for each group)")
curve(2 * x + 2 * x^2, -3, 3, col = "darkgreen", lty = 3, lwd = 3, add = TRUE)
curve(5 - 4 * x, -3, 3, col = "darkorange", lty = 3, lwd = 3, add = TRUE)
legend("top", c("True-treated", "True-nontreated"),
      col = c("darkorange", "darkgreen"), lty = 3, lwd = 3, bty = 'n', cex = .8)
```



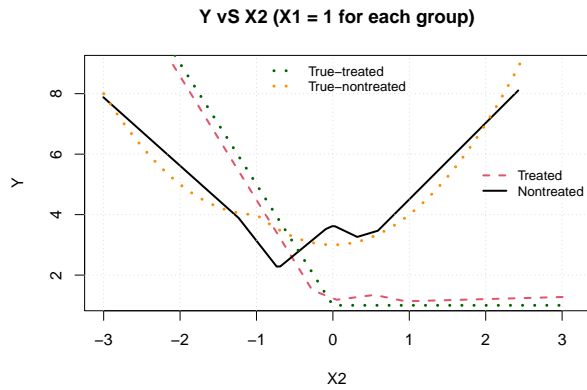
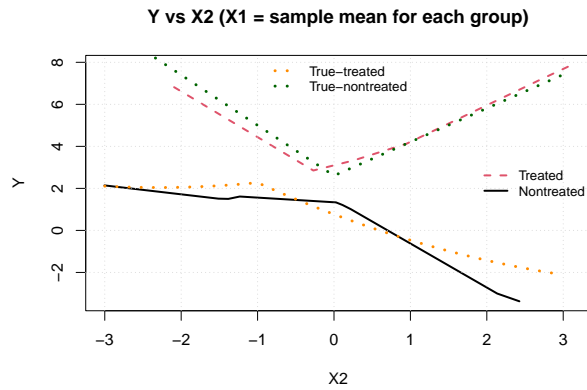
The following graphs illustrate the relationship between  $Y$  and  $X_2$  for a given  $X_1$ . When  $X_1$  is equal its population group means (they are equal to the population means of  $X_2$ ), the true curves are  $[1.6036 - 0.3964x](x \leq 0) + (1 + 2x)(x > 0)$  for the treated and  $[(1.767 - 0.3698x + 0.1368x^2) + (1 + x)(x \leq -1) + (-1 - x)(x > -1)]$  for the nontreated. If  $X_1 = 1$ , the true curves become  $[-2x + (1 - 2x)(x \leq 0) + (1 + 2x)(x > 0)]$  for the treated and  $[(4 + x + x^2) + (1 + x)(x \leq -1) + (-1 - x)(x > -1)]$  for the nontreated.

```
x10 <- mean(subset(simDat5, Z == 0)$X1)
x11 <- mean(subset(simDat5, Z == 1)$X1)
plot(c2, "X2", fixedCov0 = list(X1 = x10), fixedCov1 = list(X1 = x11),
     legendPos = "right", cex = .8,
     main = "Y vs X2 (X1 = sample mean for each group)")
curve(1.603900 - .3964 * x + (1 - 2 * x) * (x <= 0) + (1 + 2 * x) * (x > 0), -3, 3,
      col = "darkgreen", lty = 3, lwd = 3, add = TRUE)
curve(1.767 - 0.3698 * x + 0.1368 * x^2 + (1 + x) * (x <= -1) + (-1 - x) * (x > -1),
      -3, 3, col = "darkorange", lty = 3, lwd = 3, add = TRUE)
legend("top", c("True-treated", "True-nontreated"),
      col = c("darkorange", "darkgreen"), lty = 3, lwd = 3, bty = 'n', cex = .8)
```

```

plot(c2, "X2", fixedCov0 = list(X1 = 1), legendPos = "right", cex = .8,
     main="Y vS X2 (X1 = 1 for each group)")
curve(-2 * x + (1 - 2 * x) * (x <= 0) + (1 + 2 * x) * (x > 0), -3, 3,
      col = "darkgreen", lty = 3, lwd = 3, add = TRUE)
curve(4 + (1 + x) * (x <= -1) + (-1 - x) * (x > -1) + x + x^2,
      -3, 3, col = "darkorange", lty = 3, lwd = 3, add = TRUE)
legend("top", c("True-treated", "True-nontreated"),
      col = c("darkgreen", "darkorange"), lty = 3, lwd = 3, bty = 'n', cex = .8)

```



## References

- Giurcanu, M. 2016. “Thresholding Least-Squares Inference in High-Dimensional Regression Models.” *Electronic Journal of Statistics*. <https://doi.org/10:2124/T1\textendash2156>.
- Giurcanu, M., M. Capanu, P. Chaussé, and G. Luta. 2023. “Semiparametric Thresholding Least Squares Inference for Causal Effects.” *Working Paper*.
- Leifeld, Philip. 2013. “texreg: Conversion of Statistical Model Output in R to LaTeX and HTML Tables.” *Journal of Statistical Software* 55 (8): 1–24. <http://dx.doi.org/10.18637/jss.v055.i08>.
- Zeileis, Achim. 2006. “Object-Oriented Computation of Sandwich Estimators.” *Journal of Statistical Software* 16 (9): 1–16. <https://doi.org/10.18637/jss.v016.i09>.