

vignette

Introduction to causalglm

```
#devtools::install_github("tlverse/causalglm")
library(causalglm)

## Loading required package: sl3
## Loading required package: hal9001
## Loading required package: Rcpp
## hal9001 v0.4.0: The Scalable Highly Adaptive Lasso
## note: fit_hal defaults have changed. See ?fit_hal for details
## Loading required package: data.table
## Loading required package: R6
## Loading required package: tmle3
##
## Attaching package: 'tmle3'
## The following object is masked from 'package:sl3':
##
##     loss_loglik_binomial
```

causalglm is an R package for robust generalized linear models and interpretable causal inference for heterogeneous (or conditional) treatment effects. Specifically, causalglm very significantly relaxes the assumptions needed for useful causal estimates and correct inference by employing semi and nonparametric models and adaptive machine-learning through targeted maximum likelihood estimation (TMLE). See the writeup [causalglm.pdf](#) for a more theoretical overview of the methods implemented in this package.

The statistical data-structure used throughout this package is $O = (W, A, Y)$ where W represents a random vector of baseline (pretreatment) covariates/confounders, A is a usually binary treatment assignment with values in $c(0, 1)$, and Y is some outcome variable. For marginal structural models, we also consider a subvector $V \subset W$ that represents a subset of baseline variables that are of interest.

The estimands supported by causalglm are

1. Conditional average treatment effect (CATE) for arbitrary outcomes: $E[Y|A=1, W] - E[Y|A=0, W]$
2. Conditional odds ratio (OR) for binary outcomes: $\frac{P(Y=1|A=1, W)/P(Y=0|A=1, W)}{P(Y=1|A=0, W)/P(Y=0|A=0, W)}$
3. Conditional relative risk (RR) for binary, count or nonnegative outcomes: $E[Y|A=1, W]/E[Y|A=0, W]$
4. Conditional treatment-specific mean (TSM) : $E[Y|A = a, W]$
5. Conditional average treatment effect among the treated (CATT) : the best approximation of $E[Y|A=1, W] - E[Y|A=0, W]$ based on a user-specified formula/parametric model among the treated (i.e. observations with $A = 1$)

causalglm also supports the following marginal structural model estimands:

1. Marginal structural models for the CATE: $E[CATE(W)|V] := E[E[Y|A=1,W] - E[Y|A=0,W]|V]$
2. Marginal structural models for the RR: $E[E[Y|A=1,W]|V]/E[E[Y|A=0,W]|V]$
3. Marginal structural models for the TSM : $E[E[Y|A = a, W]|V]$
4. Marginal structural models for the CATT : $E[CATE(W)|V, A = 1] := E[E[Y|A = 1, W] - E[Y|A = 0, W]|V, A = 1]$

causalglm consists of four main functions:

1. spglm for semiparametric estimation of correctly specified parametric models for the CATE, RR and OR
2. npglm for robust nonparametric estimation for user-specified approximation models for the CATE, CATT, TSM, RR or OR
3. msmglm for robust nonparametric estimation for user-specified marginal structural models for the CATE, CATT, TSM or RR
4. causalglmnet for high dimensional confounders W (a custom wrapper function for spglm focused on big data where standard ML may struggle)

spglm is a semiparametric method which means that it assumes the user-specified parametric model is correct for inference. This method should be used if you are very confident in your parametric model. npglm is a nonparametric method that views the user-specified parametric model as an approximation or working-model for the true nonparametric estimand. The estimands are the best causal approximation of the true conditional estimand (i.e. projections). Because of this model agnostic view, npglm provides interpretable estimates and correct inference under no conditions. The user-specified parametric model need not be correct or even a good approximation for inference! npglm should be used if you believe your parametric model is a good approximation but are not very confident that it is correct. Also, it never hurts to run both spglm and npglm for robustness! If the parametric model is close to correct then the two methods should give similar estimates. Finally, msmglm deals with marginal structural models for the conditional treatment effect estimands. This method is useful if you are only interested in modeling the causal treatment effect as a function of a subset of variables V adjusting for all the available confounders W that remain. This allows for parsimonious causal modeling, still maximally adjusting for confounding. This function can be used to understand the causal variable importance of individual variables (by having V be a single variable) and allows for nice plots (see `plot_msm`).

Overview of features using `estimand = "CATE"` as an example

We will begin with the conditional average treatment effect estimand (CATE) and use it to illustrate the features of causalglm. Afterwards, we will go through all the other available estimands.

We will use the following simulated data throughout this part.

```
n <- 250
W1 <- runif(n, min = -1, max = 1)
W2 <- runif(n, min = -1, max = 1)
A <- rbinom(n, size = 1, prob = plogis((W1 + W2)/3))
Y <- rnorm(n, mean = A * (1 + W1 + 2*W1^2) + sin(4 * W2) + sin(4 * W1), sd = 0.3)
data <- data.frame(W1, W2, A, Y)
```

spglm with CATE

All methods in causalglm have a similar argument setup. Mainly, they require a formula that specifies a parametric form for the conditional estimand, a data.frame with the data, and character vectors containing the names of the variables W , A and Y . The estimand is specified with the argument `estimand` and the learning method is specified with the `learning_method` argument.

```
formula <- ~ poly(W1, degree = 2, raw = T) # A correctly specified polynomial model of degree 2
output <- spglm(formula,
  data,
  W = c("W1", "W2"), A = "A", Y = "Y",
  estimand = "CATE", # Options are CATE, RR, OR
  learning_method = "HAL" # A bunch of options. Default is a custom semiparametric Highly Adaptive
)
```

```
## (max) epsilon: -1.685808e-04 max(abs(ED)): 5.411470e-17
```

output contains a `spglm` fit object. It contains estimates information and `tlverse/tmle3` objects that store the fit likelihood, `tmle_tasks`, and target parameter objects. There are a number of extractor functions that should suffice for almost everyone. The `summary`, `coefs`, `print` and `predict` functions should be useful. They work as follows.

```
# Print tells you the object, estimand, and a fit formula/equation for the estimand
print(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
```

```
## CATE(W) = 0.938 * (Intercept) + 1.02 * poly(W1, degree = 2, raw = T)1 + 2.27 * poly(W1, degree = 2, raw = T)2
```

Summary provides the coefficient estimates (`tmle_est`), 95% confidence intervals (lower, upper), and p-values (`p_value`). The `coef` function provides pretty much the same thing as `summary`.

```
summary(output) # Summary gives you the estimates and inference
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
```

```
## CATE(W) = 0.938 * (Intercept) + 1.02 * poly(W1, degree = 2, raw = T)1 + 2.27 * poly(W1, degree = 2, raw = T)2
```

```
##
```

```
## Coefficient estimates and inference:
```

```
##      type                param tmle_est      se      lower      upper
## 1: CATE                (Intercept) 0.938008 0.06056681 0.8192993 1.056717
## 2: CATE poly(W1, degree = 2, raw = T)1 1.019331 0.06715733 0.8877051 1.150957
## 3: CATE poly(W1, degree = 2, raw = T)2 2.274394 0.13999045 2.0000180 2.548770
##      Z_score p_value
## 1: 244.8736      0
## 2: 239.9893      0
## 3: 256.8842      0
```

The `predict` function allows you get individual-level treatment effect predictions and 95% prediction (confidence) intervals. Specifically, for each observation, the individual CATE estimate derived from the coefficient estimates is given and a 95% confidence interval + p-values for it.

```
preds <- predict(output, data = data)
preds <- predict(output) # By default, training data is used.
head(preds)
```

```
##      (Intercept) poly(W1, degree = 2, raw = T)1 poly(W1, degree = 2, raw = T)2
## 1              1              0.2317417              0.05370422
## 2              1              0.3261620              0.10638163
## 3              1              0.5575424              0.31085354
## 4              1              0.8219950              0.67567575
## 5              1             -0.5678458              0.32244887
## 6              1              0.4941714              0.24420533
##      CATE(W)      se  CI_left CI_right Z-score p-value
## 1 1.296374 0.9332725 1.1806844 1.412064 21.96301      0
## 2 1.512429 0.9089525 1.3997539 1.625104 26.30896      0
```

```
## 3 2.213332 0.9445086 2.0962493 2.330414 37.05191      0
## 4 3.312646 1.4232223 3.1362216 3.489070 36.80207      0
## 5 1.092561 0.7565141 0.9987826 1.186339 22.83488      0
## 6 1.997151 0.9091497 1.8844520 2.109851 34.73326      0
```

It is common to want to obtain multiple fits using multiple formulas. We recommend doing this with `npglm` since it always provides correct interpretable inference even when these models are wrong. It is computationally expensive to recall `spglm` for each formula since the machine-learning is redone. Instead, we can reuse the machine-learning fits from previous calls to `spglm`. Due to the semiparametric nature of `spglm`, the way this works for `spglm` differs from `npglm` and `msmgm`. For `spglm`, you can pass a previous `spglm` fit object through the `data` argument with a new formula. The previous fits will then automatically be reused. The catch for `spglm` is that the new formula must be a subset of the original formula from the previous fit. Thus, one should first fit the most complex formula that contains all terms of interest and then call `spglm` with the desired subformulas. Lets see how this works. Fortunately, `npglm` and `msmgm` also allow for reusing fits and they even work across estimands and for arbitrary formulas (not just subformulas).

```
# Start with big formula
```

```
formula_full <- ~ poly(W1, degree = 3, raw = T)
output_full <- spglm(formula_full,
  data,
  W = c("W1", "W2"), A = "A", Y = "Y",
  estimand = "CATE",
  learning_method = "HAL"
)
```

```
## (max) epsilon: 2.392373e-03 max(abs(ED)): 2.335632e-16
```

```
summary(output_full)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
```

```
## CATE(W) = 0.94 * (Intercept) + 1.09 * poly(W1, degree = 3, raw = T)1 + 2.27 * poly(W1, degree = 3, raw = T)2
```

```
##
```

```
## Coefficient estimates and inference:
```

```
##      type                param    tmle_est      se      lower
## 1: CATE                (Intercept)  0.9403136 0.06159109 0.8195973
## 2: CATE poly(W1, degree = 3, raw = T)1  1.0859235 0.16861231 0.7554494
## 3: CATE poly(W1, degree = 3, raw = T)2  2.2685579 0.14098354 1.9922353
## 4: CATE poly(W1, degree = 3, raw = T)3 -0.1217830 0.27573122 -0.6622063
##      upper    Z_score    p_value
## 1: 1.0610299 241.393105 0.0000e+00
## 2: 1.4163975 101.830984 0.0000e+00
## 3: 2.5448806 254.420123 0.0000e+00
## 4: 0.4186403  6.983462 2.8799e-12
```

```
# This will give a warning since the term names for `poly(W1, degree = 2, raw = T)` are not a subset of
# Use argument warn = FALSE to turn this off.
```

```
subformula <- ~ poly(W1, degree = 2, raw = T) # one less degree
output<- spglm(subformula,
  data = output_full, # replace data with output_full
  estimand = "CATE" # No need to specify the variables again.
)
```

```
## Warning in spglm(subformula, data = output_full, estimand = "CATE"): Terms of
## new formula could not be confirmed as subsets of original formula. Make sure
## this formula is truly a subformula or else the results may be unreliable..
```

```
## (max) epsilon: -2.307535e-04 max(abs(ED)): 2.075076e-17
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
## CATE(W) = 0.94 * (Intercept) + 1.02 * poly(W1, degree = 2, raw = T)1 + 2.26 * poly(W1, degree = 2, raw = T)2
##
## Coefficient estimates and inference:
##   type                param tmle_est      se      lower      upper
## 1: CATE                (Intercept) 0.9401686 0.06059758 0.8213995 1.058938
## 2: CATE poly(W1, degree = 2, raw = T)1 1.0164289 0.06736025 0.8844053 1.148453
## 3: CATE poly(W1, degree = 2, raw = T)2 2.2588624 0.14027185 1.9839346 2.533790
##   Z_score p_value
## 1: 245.3129      0
## 2: 238.5851      0
## 3: 254.6181      0
```

```
subformula <- ~ 1 + W1 # one less degree
```

```
output<- spglm(subformula,
  data = output_full, # replace data with output_full
  estimand = "CATE", warn = FALSE # No need to specify the variables again.
)
```

```
## (max) epsilon: -5.941307e-05 max(abs(ED)): 4.780898e-17
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
## CATE(W) = 1.72 * (Intercept) + 1.22 * W1
##
## Coefficient estimates and inference:
##   type      param tmle_est      se      lower      upper Z_score p_value
## 1: CATE (Intercept) 1.716769 0.03894028 1.640447 1.793091 697.0803      0
## 2: CATE      W1 1.215516 0.06728152 1.083646 1.347385 285.6504      0
```

```
subformula <- ~ 1 # one less degree
```

```
output<- spglm(subformula,
  data = output_full, # replace data with output_full
  estimand = "CATE", warn = FALSE # No need to specify the variables again.
)
```

```
## (max) epsilon: -5.479033e-05 max(abs(ED)): 1.515281e-17
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
## CATE(W) = 1.75 * (Intercept)
##
```

```
## Coefficient estimates and inference:
```

```
##   type      param tmle_est      se      lower      upper Z_score p_value
## 1: CATE (Intercept) 1.747629 0.03795186 1.673244 1.822013 728.0916      0
```

```
# That was fast! Look how different the estimates are when the model is misspecified! (npglm would do b
```

Currently all learning was done with HAL (default and recommended in most cases). There are a number of other options. All methods in this package require machine-learning of $P(A = 1|W)$ (the propensity score) and $E[Y|A, W]$ (the conditional mean outcome). For `spglm`, $E[Y|A, W]$ is learned in a semiparametric way. By default, the learning algorithm is provided the design matrix `cbind(W, A · formula(W))` where W is a matrix with columns being the baseline variable observations and $A · formula(W)$ is a matrix with columns being the treatment interaction observations specified by the formula argument. Specifically, the design

matrix is constructed as follows:

```
formula <- ~ 1 + W1
AW <- model.matrix(formula, data)
design_mat_sp_Y <- as.matrix(cbind(data[,c("W1", "W2")], AW))
head(as.data.frame(design_mat_sp_Y))
```

```
##           W1           W2 (Intercept)           W1
## 1  0.2317417 -0.3183951           1  0.2317417
## 2  0.3261620 -0.6161039           1  0.3261620
## 3  0.5575424  0.9674222           1  0.5575424
## 4  0.8219950  0.6710243           1  0.8219950
## 5 -0.5678458 -0.1877137           1 -0.5678458
## 6  0.4941714  0.4327743           1  0.4941714
```

Since the design matrix automatically contains the treatment interaction terms, additive learners like glm, glmnet or gam can in principle perform well (since they will model treatment interactions). Note that the final regression fit based on this design matrix will be projected onto the semiparametric model using glm.fit to ensure all model constraints are satisfied (this is not important and happens behind the scenes).

This learning method corresponds with the default argument specification `append_design_matrix = TRUE`. The other option `append_design_matrix = FALSE` performs treatment-stratified estimation. Specifically, the machine-learning algorithm is used to learn the placebo conditional mean $E[Y|A = 0, W]$ by performing the regression of Y on W using only the observations with $A = 0$. Next, this initial estimator of $E[Y|A = 0, W]$ is used as an offset in a glm-type regression of Y on $A \cdot \text{formula}(W)$. This two-stage approach does not pool data across treatment arms and is thus not preferred.

Now that we got the nitty and gritty details out of the way. Lets use some different algorithms. We see that glm and glmnet perform very badly because of model misspecification. (The true model is quite nonlinear in the noninteraction terms). This motivates using causalglm over conventional methods like glm.

```
formula <- ~ poly(W1, degree = 2, raw = T)
output <- spglm(formula,
  data,
  W = c("W1", "W2"), A = "A", Y = "Y",
  estimand = "CATE",
  learning_method = "glm"
)
```

```
## (max) epsilon: 2.881782e-01 max(abs(ED)): 3.272799e-16
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
```

```
## CATE(W) = 0.627 * (Intercept) + 0.802 * poly(W1, degree = 2, raw = T)1 + 2.83 * poly(W1, degree = 2,
##
```

```
## Coefficient estimates and inference:
```

```
##   type                param  tmle_est      se    lower    upper
## 1: CATE                (Intercept) 0.6274789 0.1848293 0.2652201 0.9897377
## 2: CATE poly(W1, degree = 2, raw = T)1 0.8022313 0.1980889 0.4139842 1.1904784
## 3: CATE poly(W1, degree = 2, raw = T)2 2.8290055 0.4256897 1.9946689 3.6633420
##      Z_score p_value
## 1:   53.67824      0
## 2:   64.03383      0
## 3:  105.07771      0
```

```
output <- spglm(formula,
  data,
```

```
W = c("W1", "W2"), A = "A", Y = "Y",
estimand = "CATE",
learning_method = "glmnet"
)
```

```
## (max) epsilon: 2.761770e-01 max(abs(ED)): 2.108903e-16
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
```

```
## CATE(W) = 0.635 * (Intercept) + 0.804 * poly(W1, degree = 2, raw = T)1 + 2.81 * poly(W1, degree = 2,
##
```

```
## Coefficient estimates and inference:
```

| ## | type | param | tmle_est | se | lower | upper |
|-------|-------------------------------------|-------------|-----------|-----------|-----------|-----------|
| ## 1: | CATE | (Intercept) | 0.6349186 | 0.1850922 | 0.2721445 | 0.9976927 |
| ## 2: | CATE poly(W1, degree = 2, raw = T)1 | 0.8043955 | 0.1974651 | 0.4173711 | 1.1914200 | |
| ## 3: | CATE poly(W1, degree = 2, raw = T)2 | 2.8069908 | 0.4255645 | 1.9728996 | 3.6410820 | |

| ## | Z_score | p_value |
|-------|-----------|---------|
| ## 1: | 54.23752 | 0 |
| ## 2: | 64.40942 | 0 |
| ## 3: | 104.29069 | 0 |

```
output <- spglm(formula,
data,
W = c("W1", "W2"), A = "A", Y = "Y",
estimand = "CATE",
learning_method = "gam"
)
```

```
## (max) epsilon: 1.242178e-03 max(abs(ED)): 1.036879e-16
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
```

```
## CATE(W) = 0.943 * (Intercept) + 1.03 * poly(W1, degree = 2, raw = T)1 + 2.24 * poly(W1, degree = 2,
##
```

```
## Coefficient estimates and inference:
```

| ## | type | param | tmle_est | se | lower | upper |
|-------|-------------------------------------|-------------|------------|------------|-----------|----------|
| ## 1: | CATE | (Intercept) | 0.9428427 | 0.05842236 | 0.8283370 | 1.057348 |
| ## 2: | CATE poly(W1, degree = 2, raw = T)1 | 1.0250157 | 0.06743076 | 0.8928539 | 1.157178 | |
| ## 3: | CATE poly(W1, degree = 2, raw = T)2 | 2.2445624 | 0.13838356 | 1.9733356 | 2.515789 | |

| ## | Z_score | p_value |
|-------|----------|---------|
| ## 1: | 255.1703 | 0 |
| ## 2: | 240.3491 | 0 |
| ## 3: | 256.4586 | 0 |

```
output <- spglm(formula,
data,
W = c("W1", "W2"), A = "A", Y = "Y",
estimand = "CATE",
learning_method = "mars"
)
```

```
## (max) epsilon: 3.442839e-03 max(abs(ED)): 5.835306e-17
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
```

```
## CATE(W) = 0.955 * (Intercept) + 1.02 * poly(W1, degree = 2, raw = T)1 + 2.23 * poly(W1, degree = 2,
##
## Coefficient estimates and inference:
##   type                param  tmle_est      se    lower    upper
## 1: CATE                (Intercept) 0.9554359 0.06154026 0.8348192 1.076053
## 2: CATE poly(W1, degree = 2, raw = T)1 1.0180194 0.07362810 0.8737110 1.162328
## 3: CATE poly(W1, degree = 2, raw = T)2 2.2279090 0.14535863 1.9430113 2.512807
##   Z_score p_value
## 1: 245.4778      0
## 2: 218.6163      0
## 3: 242.3409      0

output <- spglm(formula,
  data,
  W = c("W1", "W2"), A = "A", Y = "Y",
  estimand = "CATE",
  learning_method = "xgboost"
)
```

```
## (max) epsilon: 3.595954e-02 max(abs(ED)): 8.312968e-17
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
## CATE(W) = 0.853 * (Intercept) + 1.02 * poly(W1, degree = 2, raw = T)1 + 2.37 * poly(W1, degree = 2,
##
## Coefficient estimates and inference:
##   type                param  tmle_est      se    lower    upper
## 1: CATE                (Intercept) 0.8528911 0.08551666 0.6852815 1.020501
## 2: CATE poly(W1, degree = 2, raw = T)1 1.0171690 0.11806008 0.7857755 1.248563
## 3: CATE poly(W1, degree = 2, raw = T)2 2.3674294 0.23494939 1.9069371 2.827922
##   Z_score p_value
## 1: 157.6932      0
## 2: 136.2260      0
## 3: 159.3209      0
```

npglm with CATE

npglm is a model-robust version of spglm that we personally recommend (at least as a robustness check). npglm works similarly to spglm. Fitting and extractor functions are pretty much the same.

```
formula <- ~ poly(W1, degree = 2, raw = T)
output <- npglm(formula,
  data,
  W = c("W1", "W2"), A = "A", Y = "Y",
  estimand = "CATE",
  learning_method = "HAL"
)
```

```
## (max) epsilon: 2.089282e-02 max(abs(ED)): 6.914955e-17
```

```
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand CATE with formula:
## CATE(W) = 0.966 * (Intercept) + 1.03 * poly(W1, degree = 2, raw = T)1 + 2.22 * poly(W1, degree = 2,
##
## Coefficient estimates and inference:
```



```
##      type                param  tmle_est      se      lower      upper
## 1: CATE                (Intercept) 0.9662744 0.05249831 0.8633797 1.069169
## 2: CATE poly(W1, degree = 2, raw = T)1 1.0285348 0.06739285 0.8964473 1.160622
## 3: CATE poly(W1, degree = 2, raw = T)2 2.2230771 0.12688418 1.9743887 2.471766
##      Z_score p_value
## 1: 291.0216      0
## 2: 241.3099      0
## 3: 277.0238      0
```

```
head(predict(output))
```

```
##      (Intercept) poly(W1, degree = 2, raw = T)1 poly(W1, degree = 2, raw = T)2
## 1              1              0.2317417              0.05370422
## 2              1              0.3261620              0.10638163
## 3              1              0.5575424              0.31085354
## 4              1              0.8219950              0.67567575
## 5              1             -0.5678458              0.32244887
## 6              1              0.4941714              0.24420533
##      CATE(W)      se CI_left CI_right Z-score p-value
## 1 1.324017 0.7978590 1.225114 1.422921 26.23841      0
## 2 1.538238 0.7900863 1.440298 1.636178 30.78357      0
## 3 2.230778 0.9211842 2.116586 2.344969 38.28951      0
## 4 3.313804 1.4940020 3.128606 3.499003 35.07080      0
## 5 1.099054 0.7373495 1.007651 1.190457 23.56762      0
## 6 2.017434 0.8544971 1.911510 2.123359 37.33007      0
```

npglm can reuse fits across both formulas and estimands with no restrictions. This is because the conditional mean and propensity score are learned fully nonparametrically (the previous semiparametric learning method no longer applies). The nice thing about npglm is that all models are viewed as approximations and thus each model below is interpretable as the best approximation. The intercept model is actually a nonparametric estimate for the marginal ATE! (See writeup.) Additionally, the inference for each model is correct (we don't require correctly specified parametric models!).

```
formula <- ~ 1 # We can start with simplest model. npglm does not care.
output_full <- npglm(formula,
  data,
  W = c("W1", "W2"), A = "A", Y = "Y",
  estimand = "CATE",
  learning_method = "HAL"
)
```

```
## (max) epsilon: 1.166403e-03 max(abs(ED)): 1.720846e-17
```

```
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand CATE with formula:
```

```
## CATE(W) = 0.966 * (Intercept) + 1.03 * poly(W1, degree = 2, raw = T)1 + 2.22 * poly(W1, degree = 2, raw = T)2
```

```
##
```

```
## Coefficient estimates and inference:
```

```
##      type                param  tmle_est      se      lower      upper
## 1: CATE                (Intercept) 0.9662744 0.05249831 0.8633797 1.069169
## 2: CATE poly(W1, degree = 2, raw = T)1 1.0285348 0.06739285 0.8964473 1.160622
## 3: CATE poly(W1, degree = 2, raw = T)2 2.2230771 0.12688418 1.9743887 2.471766
##      Z_score p_value
## 1: 291.0216      0
## 2: 241.3099      0
## 3: 277.0238      0
```

```
formula <- ~ 1 + W1
output <- npglm(formula,
  output_full,
  estimand = "CATE"
)
```

```
## [1] "Reusing previous fit..."
## (max) epsilon: 1.644325e-03 max(abs(ED)): 3.717165e-17
```

```
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand CATE with formula:
## CATE(W) = 1.73 * (Intercept) + 1.21 * W1
```

```
##
```

```
## Coefficient estimates and inference:
```

| | type | param | tmle_est | se | lower | upper | Z_score | p_value |
|-------|------|-------------|----------|------------|-----------|----------|----------|---------|
| ## 1: | CATE | (Intercept) | 1.726024 | 0.05478257 | 1.6186526 | 1.833396 | 498.1665 | 0 |
| ## 2: | CATE | W1 | 1.209410 | 0.11419634 | 0.9855889 | 1.433230 | 167.4523 | 0 |

```
formula <- ~ poly(W1, degree = 2, raw = T)
output <- npglm(formula,
  output_full,
  estimand = "CATE"
)
```

```
## [1] "Reusing previous fit..."
## (max) epsilon: 2.118516e-02 max(abs(ED)): 1.711773e-16
```

```
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand CATE with formula:
```

```
## CATE(W) = 0.967 * (Intercept) + 1.02 * poly(W1, degree = 2, raw = T)1 + 2.22 * poly(W1, degree = 2, raw = T)2
```

```
##
```

```
## Coefficient estimates and inference:
```

| | type | param | tmle_est | se | lower | upper | Z_score | p_value |
|-------|------|--------------------------------|-----------|------------|-----------|----------|----------|---------|
| ## 1: | CATE | (Intercept) | 0.9668825 | 0.05235138 | 0.8642757 | 1.069489 | 292.0220 | 0 |
| ## 2: | CATE | poly(W1, degree = 2, raw = T)1 | 1.0247047 | 0.06447282 | 0.8983403 | 1.151069 | 251.2998 | 0 |
| ## 3: | CATE | poly(W1, degree = 2, raw = T)2 | 2.2166762 | 0.12181121 | 1.9779306 | 2.455422 | 287.7299 | 0 |

```
formula <- ~ poly(W1, degree = 3, raw = T)
output <- npglm(formula,
  output_full,
  estimand = "CATE"
)
```

```
## [1] "Reusing previous fit..."
## (max) epsilon: 2.136724e-02 max(abs(ED)): 3.334485e-16
```

```
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand CATE with formula:
```

```
## CATE(W) = 0.968 * (Intercept) + 1.07 * poly(W1, degree = 3, raw = T)1 + 2.22 * poly(W1, degree = 3, raw = T)2
```

```
##
```

```
## Coefficient estimates and inference:
```

```
##      type                param    tmle_est      se      lower
## 1: CATE (Intercept)  0.96759949 0.05212237 0.8654415
## 2: CATE poly(W1, degree = 3, raw = T)1  1.07455363 0.15184635 0.7769403
## 3: CATE poly(W1, degree = 3, raw = T)2  2.22041970 0.12167284 1.9819453
## 4: CATE poly(W1, degree = 3, raw = T)3 -0.08771426 0.23252500 -0.5434549
##      upper      Z_score    p_value
## 1: 1.0697575 293.522560 0.0000e+00
## 2: 1.3721670 111.890639 0.0000e+00
## 3: 2.4588941 288.543584 0.0000e+00
## 4: 0.3680264  5.964452 2.4546e-09
```

causalglmnet with CATE

causalglmnet is a wrapper for spglm that uses the LASSO with glmnet for all estimation. This is made for high dimensional settings. It is used in the same way as spglm.

```
formula <- ~ poly(W1, degree = 3, raw = T)
output <- causalglmnet(formula,
  data,
  W = c("W1", "W2"), A = "A", Y = "Y",
  estimand = "CATE"
)
```

```
## (max) epsilon: 2.563074e+00 max(abs(ED)): 1.381742e-15
```

```
summary(output)
```

```
## A causalglm fit object obtained from causalglmnet for the estimand CATE with formula:
```

```
## CATE(W) = 0.711 * (Intercept) + 0.668 * poly(W1, degree = 3, raw = T)1 + 2.73 * poly(W1, degree = 3,
##
```

```
## Coefficient estimates and inference:
```

```
##      type                param    tmle_est      se      lower      upper
## 1: CATE (Intercept)  0.7105508 0.1706801 0.3760240 1.045078
## 2: CATE poly(W1, degree = 3, raw = T)1  0.6680459 0.4472515 -0.2085509 1.544643
## 3: CATE poly(W1, degree = 3, raw = T)2  2.7325617 0.3676210 2.0120379 3.453086
## 4: CATE poly(W1, degree = 3, raw = T)3  0.3115097 0.7045428 -1.0693688 1.692388
##      Z_score    p_value
## 1:  65.823708 0.0000e+00
## 2:  23.616989 0.0000e+00
## 3: 117.527559 0.0000e+00
## 4:   6.990918 2.7309e-12
```

msmgglm with CATE

msmgglm is for learning marginal structural models (e.g. marginal estimands like the ATE, ATT, and marginal relative risk). It operates in the same way as npglm. It is also a nonparametrically robust method that does not require correct model specification and estimates the best approximation. The only difference is that the marginal covariate(s) of interest V need to be specified. It also has a useful plotting feature that displays 95% confidence bands (only if V is one-dimensional). This method is used if you have many confounders W for which to adjust but only care about the treatment effect association with a subset of variables V . This can be used to build causal predictors that only utilize a handful of variables.

```
formula <- ~ poly(W1, degree = 3, raw = T)
output <- msmgglm(formula,
  data,
  V = "W1",
  W = c("W1", "W2"), A = "A", Y = "Y",
```

```

estimand = "CATE",
learning_method = "HAL"
)

```

```
## (max) epsilon: 2.711186e-02 max(abs(ED)): 3.334624e-16
```

```
summary(output)
```

```
## A causalglm fit object obtained from msglm for the estimand CATE with formula:
```

```
## E[CATE(W)|V] = 0.971 * (Intercept) + 1.07 * poly(W1, degree = 3, raw = T)1 + 2.21 * poly(W1, degree = 3, raw = T)2 + -0.09 * poly(W1, degree = 3, raw = T)3
```

```
##
```

```
## Coefficient estimates and inference:
```

| | type | param | tmle_est | se | lower |
|-------|-------------------------------------|-------------|-------------|------------|------------|
| ## 1: | CATE | (Intercept) | 0.97126183 | 0.05302221 | 0.8673402 |
| ## 2: | CATE poly(W1, degree = 3, raw = T)1 | | 1.07254487 | 0.15330403 | 0.7720745 |
| ## 3: | CATE poly(W1, degree = 3, raw = T)2 | | 2.20969784 | 0.12523093 | 1.9642497 |
| ## 4: | CATE poly(W1, degree = 3, raw = T)3 | | -0.09467703 | 0.23700798 | -0.5592041 |

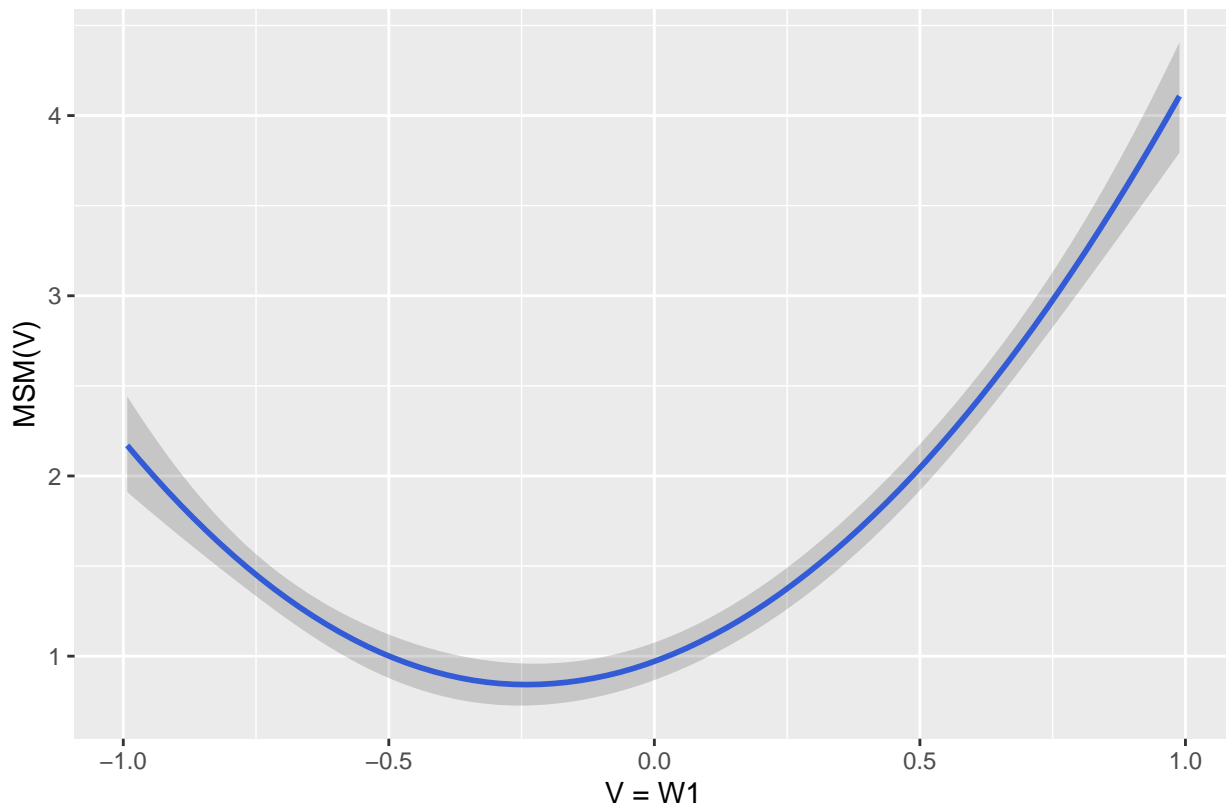
| | upper | Z_score | p_value |
|-------|-----------|------------|------------|
| ## 1: | 1.0751835 | 289.633303 | 0.0000e+00 |
| ## 2: | 1.3730152 | 110.619554 | 0.0000e+00 |
| ## 3: | 2.4551459 | 278.991710 | 0.0000e+00 |
| ## 4: | 0.3698501 | 6.316139 | 2.6818e-10 |

```
plot_msm(output)
```

```
## Loading required package: ggplot2
```

```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```

```
E[CATE(W)|V] = 0.971 * (Intercept) + 1.07 * poly(W1, degree = 3, raw = T)1 + 2.21 * poly(W1, degree = 3, raw = T)2 + -0.09 * poly(W1, degree = 3, raw = T)3
```



```

formula <- ~ 1 + W1 # Best linear approximation
output <- msglm(formula,
  data,
  V = "W1",
  W = c("W1", "W2"), A = "A", Y = "Y",
  estimand = "CATE",
  learning_method = "HAL"
)

```

```
## (max) epsilon: 1.804041e-03 max(abs(ED)): 3.808065e-17
```

```
summary(output)
```

```

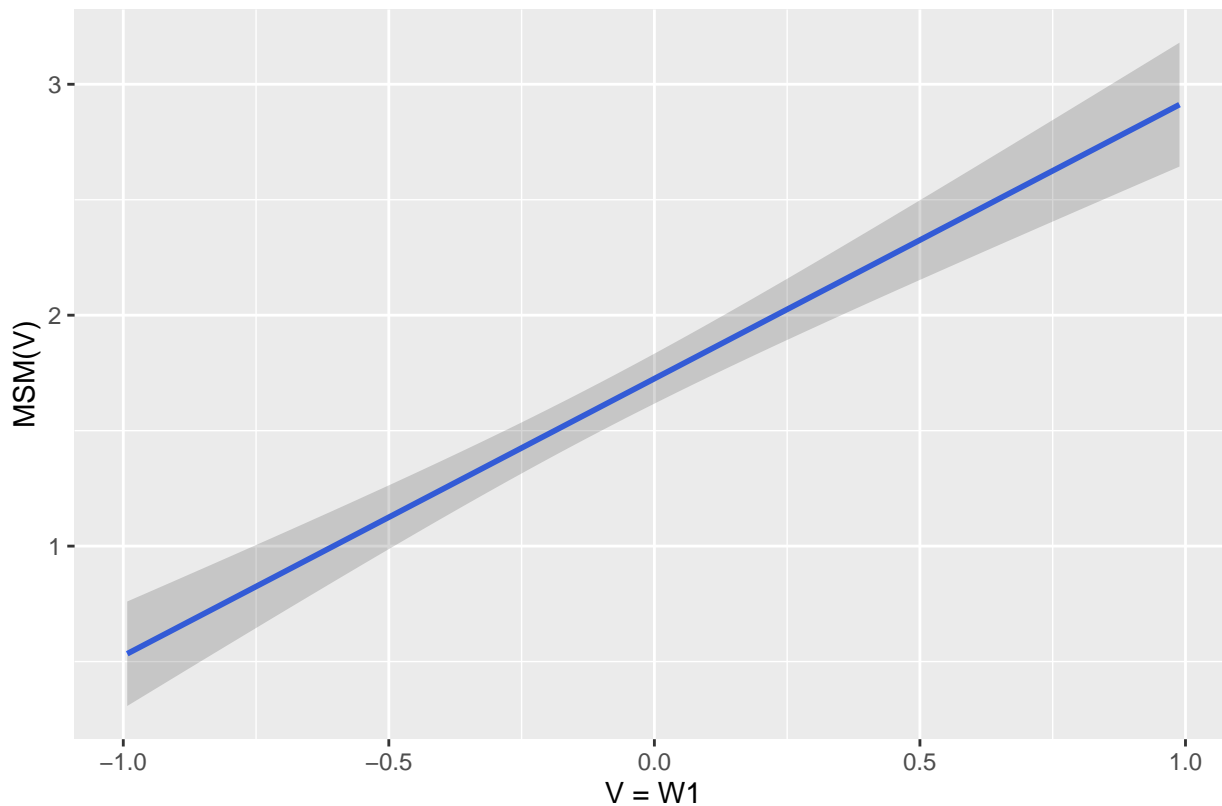
## A causalglm fit object obtained from msglm for the estimand CATE with formula:
## E[CATE(W)|V] = 1.73 * (Intercept) + 1.2 * W1
##
## Coefficient estimates and inference:
##   type      param tmle_est      se    lower    upper  Z_score p_value
## 1: CATE (Intercept) 1.725167 0.05520031 1.6169760 1.833357 494.1508      0
## 2: CATE              W1 1.200215 0.11506330 0.9746951 1.425735 164.9272      0

```

```
plot_msm(output)
```

```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```

$E[CATE(W)|V] = 1.73 * (\text{Intercept}) + 1.2 * W1$



```
# This gives a nonparametric estimate for the marginal ATE
```

```

formula <- ~ 1
output <- msglm(formula,
  data,

```

```

V = "W1",
W = c("W1", "W2"), A = "A", Y = "Y",
estimand = "CATE",
learning_method = "HAL"
)

## (max) epsilon: 1.084738e-03 max(abs(ED)): 3.275158e-17

summary(output)

## A causalglm fit object obtained from msmglm for the estimand CATE with formula:
## E[CATE(W)|V] = 1.75 * (Intercept)
##
## Coefficient estimates and inference:
##   type      param tml_e_st      se    lower    upper  Z_score p_value
## 1: CATE (Intercept) 1.754421 0.06374637 1.629481 1.879362 435.1595      0

```

Learning other estimands.

All of the vignette discussed so far can be applied to other estimands by specifying a different “estimand” argument.

Let us begin with npglm (msmglm acts in the same exact way). Both npglm and msmglm support the CATE, OR, RR, CATT and TSM

```

n <- 250
W1 <- runif(n, min = -1, max = 1)
W2 <- runif(n, min = -1, max = 1)
A <- rbinom(n, size = 1, prob = plogis((W1 + W2 )/3))
Y <- rnorm(n, mean = A * (1 + W1 + 2*W1^2) + sin(4 * W2) + sin(4 * W1), sd = 0.3)
data <- data.frame(W1, W2,A,Y)
# CATE
formula = ~ poly(W1, degree = 2, raw = TRUE)
output <- npglm(formula,
  data,
  W = c("W1", "W2"), A = "A", Y = "Y",
  estimand = "CATE")

```

```
## (max) epsilon: 1.278125e-02 max(abs(ED)): 9.580184e-17
```

```
summary(output)
```

```

## A causalglm fit object obtained from npglm for the estimand CATE with formula:
## CATE(W) = 0.989 * (Intercept) + 1.05 * poly(W1, degree = 2, raw = TRUE)1 + 2.06 * poly(W1, degree = 2, raw = TRUE)2
##
## Coefficient estimates and inference:
##   type      param tml_e_st      se    lower    upper  Z_score p_value
## 1: CATE (Intercept) 0.9890949 0.05135432 0.8884423
## 2: CATE poly(W1, degree = 2, raw = TRUE)1 1.0482672 0.06064779 0.9293997
## 3: CATE poly(W1, degree = 2, raw = TRUE)2 2.0563742 0.11627220 1.8284849
##   upper  Z_score p_value
## 1: 1.089747 304.5306      0
## 2: 1.167135 273.2921      0
## 3: 2.284264 279.6380      0

# CATT, lets reuse fit
output <- npglm(formula,

```

```

    output,
    estimand = "CATT")

## [1] "Reusing previous fit..."
## (max) epsilon: 1.613398e-02 max(abs(ED)): 4.648781e-16

summary(output)

## A causalglm fit object obtained from npglm for the estimand CATT with formula:
## CATT(W) = 0.995 * (Intercept) + 1.05 * poly(W1, degree = 2, raw = TRUE)1 + 2.04 * poly(W1, degree = 2, raw = TRUE)2
##
## Coefficient estimates and inference:
##      type                param tmle_est      se      lower
## 1: CATT                (Intercept) 0.9947908 0.05462320 0.8877313
## 2: CATT poly(W1, degree = 2, raw = TRUE)1 1.0527126 0.05867836 0.9377051
## 3: CATT poly(W1, degree = 2, raw = TRUE)2 2.0383799 0.11876642 1.8056020
##      lower      upper  Z_score p_value
## 1: 1.101850 287.9550      0
## 2: 1.167720 283.6625      0
## 3: 2.271158 271.3698      0

# TSM, note this provides a list of npglm objects for each level of `A`.
outputs <- npglm(formula,
  output,
  estimand = "TSM")

## [1] "Reusing previous fit..."
## (max) epsilon: 2.053621e-02 max(abs(ED)): 1.842138e-16
## [1] 0 1

summary(outputs[[1]])

## A causalglm fit object obtained from npglm for the estimand TSM with formula:
## TSM(W) = 1.08 * E[Y_{A=1}]: (Intercept) + 1.14 * E[Y_{A=1}]: poly(W1, degree = 2, raw = TRUE)1 + 2.04 * E[Y_{A=1}]: poly(W1, degree = 2, raw = TRUE)2
##
## Coefficient estimates and inference:
##      type                param tmle_est      se
## 1: TSM                E[Y_{A=1}]: (Intercept) 1.077243 0.08917413
## 2: TSM E[Y_{A=1}]: poly(W1, degree = 2, raw = TRUE)1 1.144763 0.11697361
## 3: TSM E[Y_{A=1}]: poly(W1, degree = 2, raw = TRUE)2 2.011043 0.20929915
##      lower      upper  Z_score p_value
## 1: 0.9024652 1.252021 191.0051      0
## 2: 0.9154989 1.374027 154.7383      0
## 3: 1.6008242 2.421262 151.9231      0

summary(outputs[[2]])

## A causalglm fit object obtained from npglm for the estimand TSM with formula:
## TSM(W) = 0.0912 * E[Y_{A=0}]: (Intercept) + 0.107 * E[Y_{A=0}]: poly(W1, degree = 2, raw = TRUE)1 + 0.04465915 * E[Y_{A=0}]: poly(W1, degree = 2, raw = TRUE)2
##
## Coefficient estimates and inference:
##      type                param tmle_est      se
## 1: TSM                E[Y_{A=0}]: (Intercept) 0.09116781 0.09139883
## 2: TSM E[Y_{A=0}]: poly(W1, degree = 2, raw = TRUE)1 0.10702472 0.12196780
## 3: TSM E[Y_{A=0}]: poly(W1, degree = 2, raw = TRUE)2 -0.04465915 0.22152713
##      lower      upper  Z_score p_value
## 1: -0.08797061 0.2703062 15.771423 0.000000

```

```
## 2: -0.13202778 0.3460772 13.874230 0.000000
## 3: -0.47884435 0.3895260 3.187525 0.001435
```

Both the OR and RR estimands provide the original coefficient estimates and their exponential transforms. This is because the parametric model/formula is actually for the log RR and log OR (that is log-linear models). The predict function gives the exponential of the linear predictor (so actually predicts the OR and RR).

```
# odds ratio
n <- 250
W <- runif(n, min = -1, max = 1)
A <- rbinom(n, size = 1, prob = plogis(W))
Y <- rbinom(n, size = 1, prob = plogis(A + A * W + W + sin(5 * W)))
data <- data.frame(W, A, Y)
output <-
  npglm(
    ~1+W,
    data,
    W = c("W"), A = "A", Y = "Y",
    estimand = "OR"
  )
```

```
## risk_change: -9.444701e-05 (max) epsilon: 2.499999e-02 max(abs(ED)): 5.291292e-02
## risk_change: -2.985830e-05 (max) epsilon: 2.409428e-02 max(abs(ED)): 2.083815e-03
```

```
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand OR with formula:
## log OR(W) = 0.958 * (Intercept) + 1.11 * W
##
## Coefficient estimates and inference:
##      type      param  tmle_est      se      lower      upper  psi_exp lower_exp
## 1:   OR (Intercept) 0.9576024 0.2780915 0.4125532 1.502652 2.605442 1.510670
## 2:   OR              W 1.1058034 0.4893491 0.1466969 2.064910 3.021651 1.158003
##      upper_exp Z_score p_value
## 1:  4.493589 54.44620      0
## 2:  7.884588 35.72969      0
```

```
output <-
  spglm(
    ~1+W,
    data,
    W = c("W"), A = "A", Y = "Y",
    estimand = "OR"
  )
```

```
## risk_change: -2.053396e-05 (max) epsilon: 2.284645e-02 max(abs(ED)): 9.052419e-04
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand OR with formula:
## log OR(W) = 1 * (Intercept) + 1.04 * W
##
## Coefficient estimates and inference:
##      type      param  tmle_est      se      lower      upper  psi_exp lower_exp
## 1:   OR (Intercept) 0.9997166 0.2793692 0.45216298 1.547270 2.717512 1.571708
## 2:   OR              W 1.0425404 0.5061006 0.05060155 2.034479 2.836414 1.051904
##      upper_exp Z_score p_value
```



```
## 1:  4.698627 56.58070      0
## 2:  7.648269 32.57062      0
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand OR with formula:
## log OR(W) = 1 * (Intercept) + 1.04 * W
##
## Coefficient estimates and inference:
##      type      param tmle_est      se      lower      upper psi_exp lower_exp
## 1:   OR (Intercept) 0.9997166 0.2793692 0.45216298 1.547270 2.717512 1.571708
## 2:   OR      W 1.0425404 0.5061006 0.05060155 2.034479 2.836414 1.051904
##      upper_exp Z_score p_value
## 1:  4.698627 56.58070      0
## 2:  7.648269 32.57062      0
```

```
# relative risk
```

```
n <- 250
W <- runif(n, min = -1, max = 1)
A <- rbinom(n, size = 1, prob = plogis(W))
Y <- rpois(n, lambda = exp( A * (1 + W + 2*W^2) + sin(5 * W)))
data <- data.frame(W, A, Y)
formula = ~ poly(W, degree = 2, raw = TRUE)
output <-
  npglm(
    formula,
    data,
    W = "W", A = "A", Y = "Y",
    estimand = "RR",
    verbose = FALSE
  )
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand RR with formula:
## log RR(W) = 0.165 * (Intercept) + 1.82 * poly(W, degree = 2, raw = TRUE)1 + 3.96 * poly(W, degree = 2, raw = TRUE)
##
## Coefficient estimates and inference:
##      type      param tmle_est      se      lower
## 1:   RR      (Intercept) 0.1653645 0.2073348 -0.2410041
## 2:   RR poly(W, degree = 2, raw = TRUE)1 1.8231213 0.2503624 1.3324201
## 3:   RR poly(W, degree = 2, raw = TRUE)2 3.9633509 0.4747933 3.0327730
##      upper psi_exp lower_exp upper_exp Z_score p_value
## 1: 0.5717332 1.179823 0.7858384 1.771334 12.61073      0
## 2: 2.3138225 6.191152 3.7902048 10.113007 115.13743      0
## 3: 4.8939287 52.633398 20.7547057 133.476941 131.98601      0
```

```
output <-
  spglm(
    formula,
    data,
    W = "W", A = "A", Y = "Y",
    estimand = "RR",
    verbose = FALSE
  )
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand RR with formula:
## log RR(W) = 0.776 * (Intercept) + 1.13 * poly(W, degree = 2, raw = TRUE)1 + 2.47 * poly(W, degree = 2, raw = TRUE)2
##
## Coefficient estimates and inference:
##      type                param  tmle_est      se      lower      upper
## 1:   RR      (Intercept) 0.7762781 0.1348308 0.5120145 1.040542
## 2:   RR poly(W, degree = 2, raw = TRUE)1 1.1317651 0.1597535 0.8186540 1.444876
## 3:   RR poly(W, degree = 2, raw = TRUE)2 2.4742653 0.2827654 1.9200553 3.028475
##      psi_exp lower_exp upper_exp  Z_score p_value
## 1:  2.173368  1.668649  2.830750  91.03284      0
## 2:  3.101126  2.267446  4.241327 112.01492      0
## 3: 11.872980  6.821336 20.665697 138.35347      0
```

```
output <-
  msglm(
    formula,
    data,
    V = "W",
    W = "W", A = "A", Y = "Y",
    estimand = "RR",
    verbose = FALSE
  )
summary(output)
```

```
## A causalglm fit object obtained from msglm for the estimand RR with formula:
## log E[RR(W)|V] = 0.162 * (Intercept) + 1.83 * poly(W, degree = 2, raw = TRUE)1 + 3.97 * poly(W, degree = 2, raw = TRUE)2
##
## Coefficient estimates and inference:
##      type                param  tmle_est      se      lower
## 1:   RR      (Intercept) 0.1618289 0.2075373 -0.2449368
## 2:   RR poly(W, degree = 2, raw = TRUE)1 1.8289672 0.2513692  1.3362926
## 3:   RR poly(W, degree = 2, raw = TRUE)2 3.9743707 0.4755755  3.0422598
##      upper  psi_exp lower_exp  upper_exp  Z_score p_value
## 1: 0.5685946  1.175659  0.782754  1.765784 12.32906      0
## 2: 2.3216417  6.227451  3.804911 10.192394 115.04397      0
## 3: 4.9064816 53.216616 20.952538 135.163015 132.13531      0
```

Custom learners with sl3

We refer to the documentation of the `tverse/sl3` package for how learners work. To specify custom learners for the propensity score use the argument `sl3_learner_A` and to specify custom learners for the outcome conditional mean use the argument `sl3_learner_Y`. For `spglm`, keep in mind the argument “`append_design_matrix`” when choosing learners. A good rule of thumb for `spglm` is to think of `sl3_learner_Y` as a learner for $E[Y|A = 0, W]$. For `msglm` and `npglm`, the learning is fully nonparametric and the regression is performed how you would expect (a standard design matrix containing W and A is passed to the learner). For `msglm` and `npglm`, make sure the learner models interactions, specifically treatment interactions, as these are crucial for fitting the conditional treatment effect estimands well.

```
lrnr_A <- Lrnr_gam$new()
lrnr_Y <- Lrnr_xgboost$new(max_depth = 4)
lrnr_Y <- Lrnr_cv$new(lrnr_Y, full_fit = TRUE) #cross-fit xgboost

n <- 250
W1 <- runif(n, min = -1, max = 1)
W2 <- runif(n, min = -1, max = 1)
```

```

A <- rbinom(n, size = 1, prob = plogis((W1 + W2 )/3))
Y <- rnorm(n, mean = A * (1 + W1 + 2*W1^2) + sin(4 * W2) + sin(4 * W1), sd = 0.3)
data <- data.frame(W1, W2,A,Y)
# CATE
formula = ~ poly(W1, degree = 2, raw = TRUE)
output <- npglm(formula,
  data,
  W = c("W1", "W2"), A = "A", Y = "Y",
  estimand = "CATE",
  sl3_Learner_A = lrnr_A,
  sl3_Learner_Y = lrnr_Y)

## (max) epsilon: 1.274702e-01 max(abs(ED)): 1.903148e-16

```

Other arguments

See the documentation for other arguments for all methods. We note that the remaining arguments will likely not be needed for the average user.