# vignette

## Introduction to causalglm

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: 5.128246e-04 max(abs(ED)): 1.548831e-16
```

`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) = 1.02 * (Intercept) + 0.985 * poly(W1, degree = 2, raw = T)1 + 1.98 * poly(W1, degree = 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) = 1.02 * (Intercept) + 0.985 * poly(W1, degree = 2, raw = T)1 + 1.98 * poly(W1, degree = 2,
##
## Coefficient estimates and inference:
##    type                        param  tmle_est        se     lower    upper
## 1: CATE                  (Intercept) 1.0167804 0.06359281 0.8921408 1.141420
## 2: CATE poly(W1, degree = 2, raw = T)1 0.9848061 0.07054760 0.8465353 1.123077
## 3: CATE poly(W1, degree = 2, raw = T)2 1.9776380 0.14922063 1.6851709 2.270105
##    Z_score p_value
## 1: 252.8071       0
## 2: 220.7184       0
## 3: 209.5501       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                      0.6750553                     0.45569965
## 2           1                     -0.2698063                     0.07279542
## 3           1                      0.4441243                     0.19724636
## 4           1                      0.4435261                     0.19671542
## 5           1                     -0.9318924                     0.86842352
## 6           1                     -0.8833139                     0.78024338
##       CATE(W)        se   CI_left CI_right  Z-score p-value
## 1 2.5827879 0.9430547 2.4658856 2.699690 43.30339       0
## 2 0.8950365 0.9691898 0.7748945 1.015179 14.60165       0
## 3 1.8442386 0.7950148 1.7456875 1.942790 36.67853       0
## 4 1.8425995 0.7951861 1.7440272 1.941172 36.63803       0
## 5 1.8164744 1.8340782 1.5891197 2.043829 15.65963       0
## 6 1.6899265 1.6598760 1.4841661 1.895687 16.09764       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 msmglm. 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 msmglm 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: 1.105621e-02 max(abs(ED)): 2.766953e-16
```

```r
summary(output_full)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
## CATE(W) = 1.01 * (Intercept) + 1.1 * poly(W1, degree = 3, raw = T)1 + 1.99 * poly(W1, degree = 3, ra
##
## Coefficient estimates and inference:
##     type                          param    tmle_est         se        lower
## 1: CATE                     (Intercept)   1.0132674  0.0632651    0.8892701
## 2: CATE poly(W1, degree = 3, raw = T)1   1.0961918  0.1627495    0.7772087
## 3: CATE poly(W1, degree = 3, raw = T)2   1.9892254  0.1483258    1.6985122
## 4: CATE poly(W1, degree = 3, raw = T)3  -0.1948914  0.2648984   -0.7140828
##        upper   Z_score p_value
## 1: 1.1372647 253.23857       0
## 2: 1.4151749 106.49690       0
## 3: 2.2799386 212.04955       0
## 4: 0.3242999  11.63278       0
```

```r
# 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: 1.661585e-03 max(abs(ED)): 3.094053e-16
```

```r
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
## CATE(W) = 1.01 * (Intercept) + 0.987 * poly(W1, degree = 2, raw = T)1 + 2.01 * poly(W1, degree = 2,
##
## Coefficient estimates and inference:
##     type                          param   tmle_est          se       lower      upper
## 1: CATE                     (Intercept)  1.0082855  0.06338630   0.8840506   1.132520
## 2: CATE poly(W1, degree = 2, raw = T)1  0.9866073  0.06988868   0.8496280   1.123587
## 3: CATE poly(W1, degree = 2, raw = T)2  2.0081064  0.14855881   1.7169365   2.299276
##     Z_score p_value
## 1: 251.5116       0
## 2: 223.2068       0
## 3: 213.7265       0
```

```r
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: 1.263999e-04 max(abs(ED)): 1.970125e-17
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
## CATE(W) = 1.66 * (Intercept) + 1.01 * W1
##
## Coefficient estimates and inference:
##     type         param tmle_est         se     lower      upper  Z_score p_value
## 1: CATE (Intercept) 1.655086 0.04038273 1.5759373 1.734235 648.0297       0
## 2: CATE           W1 1.006137 0.07024163 0.8684663 1.143808 226.4815       0
```

```r
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: 1.145922e-05 max(abs(ED)): 2.874784e-17
```

```r
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
## CATE(W) = 1.68 * (Intercept)
##
## Coefficient estimates and inference:
##     type         param tmle_est         se     lower      upper Z_score p_value
## 1: CATE (Intercept) 1.681462 0.03982041 1.603416 1.759509 667.654       0
```

```r
# 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 conditinal mean outcome). For spglm, $E[Y|A, W]$ is learned in a semiparametric way. By default, the learning algorithm is provided the design matrix $cbind(W, A \cdot formula(W))$ where $W$ is a matrix with columns being the baseline variable observations and $A \cdot formula(W)$ is a matrix with columns being the treatment interaction observations specified by the formula argument. Specifically, the design matrix is constructed as follows:

```r
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  0.6750553 -0.8973945           1  0.6750553
## 2 -0.2698063 -0.3427651           1 -0.2698063
## 3  0.4441243  0.3357459           1  0.4441243
## 4  0.4435261  0.9856179           1  0.4435261
## 5 -0.9318924 -0.9657694           1 -0.9318924
## 6 -0.8833139  0.1192606           1 -0.8833139
```

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 $A \cdot formula(W)$. 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: 9.140583e-02 max(abs(ED)): 5.832557e-16
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
## CATE(W) = 0.773 * (Intercept) + 0.899 * poly(W1, degree = 2, raw = T)1 + 2.3 * poly(W1, degree = 2,
##
## Coefficient estimates and inference:
##     type                          param  tmle_est        se     lower     upper
## 1: CATE                      (Intercept) 0.7727220 0.1919245 0.396557 1.148887
## 2: CATE poly(W1, degree = 2, raw = T)1 0.8990937 0.2253121 0.457490 1.340697
## 3: CATE poly(W1, degree = 2, raw = T)2 2.3042587 0.4913236 1.341282 3.267235
##     Z_score p_value
## 1: 63.65946       0
## 2: 63.09434       0
## 3: 74.15383       0
```

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

```
## (max) epsilon: 1.396560e-01 max(abs(ED)): 1.634942e-15
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
## CATE(W) = 0.768 * (Intercept) + 0.905 * poly(W1, degree = 2, raw = T)1 + 2.32 * poly(W1, degree = 2,
##
## Coefficient estimates and inference:
##     type                          param  tmle_est        se     lower     upper
## 1: CATE                      (Intercept) 0.7677723 0.1922913 0.3908882 1.144656
## 2: CATE poly(W1, degree = 2, raw = T)1 0.9045716 0.2241210 0.4653025 1.343841
## 3: CATE poly(W1, degree = 2, raw = T)2 2.3217030 0.4885867 1.3640906 3.279315
##     Z_score p_value
## 1: 63.13101       0
## 2: 63.81612       0
## 3: 75.13374       0
```

6

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

## (max) epsilon: 3.189288e-03 max(abs(ED)): 1.870136e-16

```r
summary(output)
```

## A causalglm fit object obtained from spglm for the estimand CATE with formula:
## CATE(W) = 0.991 * (Intercept) + 0.983 * poly(W1, degree = 2, raw = T)1 + 2.01 * poly(W1, degree = 2,
##
## Coefficient estimates and inference:
##      type                            param  tmle_est         se      lower     upper
## 1: CATE                        (Intercept) 0.9913251 0.06244799 0.8689293 1.113721
## 2: CATE poly(W1, degree = 2, raw = T)1 0.9828660 0.06892129 0.8477828 1.117949
## 3: CATE poly(W1, degree = 2, raw = T)2 2.0149649 0.14049922 1.7395914 2.290338
##      Z_score p_value
## 1: 250.9965       0
## 2: 225.4815       0
## 3: 226.7585       0

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

## (max) epsilon: -8.069873e-03 max(abs(ED)): 3.693053e-17

```r
summary(output)
```

## A causalglm fit object obtained from spglm for the estimand CATE with formula:
## CATE(W) = 1 * (Intercept) + 0.994 * poly(W1, degree = 2, raw = T)1 + 1.99 * poly(W1, degree = 2, raw
##
## Coefficient estimates and inference:
##      type                            param  tmle_est         se      lower     upper
## 1: CATE                        (Intercept) 1.0010540 0.06433902 0.8749518 1.127156
## 2: CATE poly(W1, degree = 2, raw = T)1 0.9940742 0.06873645 0.8593532 1.128795
## 3: CATE poly(W1, degree = 2, raw = T)2 1.9879028 0.14684470 1.7000925 2.275713
##      Z_score p_value
## 1: 246.0102       0
## 2: 228.6660       0
## 3: 214.0459       0

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

## (max) epsilon: 2.118945e-02 max(abs(ED)): 1.710576e-16

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand CATE with formula:
## CATE(W) = 0.982 * (Intercept) + 1.01 * poly(W1, degree = 2, raw = T)1 + 2.13 * poly(W1, degree = 2, :
##
## Coefficient estimates and inference:
##     type                       param  tmle_est        se      lower    upper
## 1: CATE                 (Intercept) 0.9823545 0.08320182 0.8192820 1.145427
## 2: CATE poly(W1, degree = 2, raw = T)1 1.0083446 0.09591448 0.8203557 1.196334
## 3: CATE poly(W1, degree = 2, raw = T)2 2.1266219 0.19898182 1.7366247 2.516619
##     Z_score p_value
## 1: 186.6833       0
## 2: 166.2244       0
## 3: 168.9845       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: 1.990756e-02 max(abs(ED)): 4.202923e-16
```

```
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand CATE with formula:
## CATE(W) = 1.01 * (Intercept) + 1.02 * poly(W1, degree = 2, raw = T)1 + 2 * poly(W1, degree = 2, raw =
##
## Coefficient estimates and inference:
##     type                       param  tmle_est        se      lower    upper
## 1: CATE                 (Intercept) 1.010602 0.06114292 0.8907645 1.130440
## 2: CATE poly(W1, degree = 2, raw = T)1 1.021836 0.06482964 0.8947727 1.148900
## 3: CATE poly(W1, degree = 2, raw = T)2 2.002873 0.13839630 1.7316210 2.274125
##     Z_score p_value
## 1: 261.3390       0
## 2: 249.2171       0
## 3: 228.8226       0
```

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

```
##   (Intercept) poly(W1, degree = 2, raw = T)1 poly(W1, degree = 2, raw = T)2
## 1           1                     0.6750553                     0.45569965
## 2           1                    -0.2698063                     0.07279542
## 3           1                     0.4441243                     0.19724636
## 4           1                     0.4435261                     0.19671542
## 5           1                    -0.9318924                     0.86842352
## 6           1                    -0.8833139                     0.78024338
##     CATE(W)        se  CI_left  CI_right  Z-score p-value
## 1 2.6131070 0.8281288 2.5104510 2.7157629 49.89181       0
```

```
## 2 0.8807045 0.9233819 0.7662409 0.9951681 15.08061        0
## 3 1.8594842 0.7528583 1.7661589 1.9528094 39.05254        0
## 4 1.8578095 0.7531303 1.7644505 1.9511685 39.00328        0
## 5 1.7977026 1.7299800 1.5832521 2.0121531 16.43035        0
## 6 1.6707284 1.5660240 1.4766020 1.8648547 16.86854        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: 3.383414e-04 max(abs(ED)): 4.054396e-17
```

```
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand CATE with formula:
## CATE(W) = 1.01 * (Intercept) + 1.02 * poly(W1, degree = 2, raw = T)1 + 2 * poly(W1, degree = 2, raw =
##
## Coefficient estimates and inference:
##     type                         param tmle_est          se      lower    upper
## 1: CATE                     (Intercept) 1.010602 0.06114292 0.8907645 1.130440
## 2: CATE poly(W1, degree = 2, raw = T)1 1.021836 0.06482964 0.8947727 1.148900
## 3: CATE poly(W1, degree = 2, raw = T)2 2.002873 0.13839630 1.7316210 2.274125
##     Z_score p_value
## 1: 261.3390       0
## 2: 249.2171       0
## 3: 228.8226       0
```

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

```
## [1] "Reusing previous fit..."
## (max) epsilon: -2.252898e-03 max(abs(ED)): 4.220929e-17
```

```
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand CATE with formula:
## CATE(W) = 1.66 * (Intercept) + 1.05 * W1
##
## Coefficient estimates and inference:
##     type       param tmle_est        se      lower    upper  Z_score p_value
## 1: CATE (Intercept) 1.660375 0.0522615 1.5579448 1.762806 502.3362       0
## 2: CATE          W1 1.049280 0.1034683 0.8464861 1.252074 160.3446       0
```

9

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

```
## [1] "Reusing previous fit..."
## (max) epsilon: 2.363362e-02 max(abs(ED)): 1.936541e-16
```

```r
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand CATE with formula:
## CATE(W) = 1.01 * (Intercept) + 1.02 * poly(W1, degree = 2, raw = T)1 + 2 * poly(W1, degree = 2, raw =
##
## Coefficient estimates and inference:
##    type                          param tmle_est        se    lower    upper
## 1: CATE                     (Intercept) 1.010267 0.0615818 0.8895692 1.130965
## 2: CATE poly(W1, degree = 2, raw = T)1 1.020919 0.0648660 0.8937835 1.148054
## 3: CATE poly(W1, degree = 2, raw = T)2 2.001053 0.1382532 1.7300822 2.272025
##    Z_score p_value
## 1: 259.3904       0
## 2: 248.8536       0
## 3: 228.8514       0
```

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

```
## [1] "Reusing previous fit..."
## (max) epsilon: 2.387915e-02 max(abs(ED)): 2.431735e-16
```

```r
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand CATE with formula:
## CATE(W) = 1.01 * (Intercept) + 1.14 * poly(W1, degree = 3, raw = T)1 + 2 * poly(W1, degree = 3, raw =
##
## Coefficient estimates and inference:
##    type                          param tmle_est        se       lower
## 1: CATE                     (Intercept)  1.009075 0.06171169   0.8881222
## 2: CATE poly(W1, degree = 3, raw = T)1  1.144151 0.15991535   0.8307223
## 3: CATE poly(W1, degree = 3, raw = T)2  2.004282 0.13802363   1.7337604
## 4: CATE poly(W1, degree = 3, raw = T)3 -0.219686 0.25703948  -0.7234742
##       upper   Z_score p_value
## 1: 1.1300276 258.53892       0
## 2: 1.4575789 113.12616       0
## 3: 2.2748031 229.60182       0
## 4: 0.2841021  13.51365       0
```

### 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.

```r
formula <- ~ poly(W1, degree = 3, raw = T)
output <- causalglmnet(formula,
```

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

## (max) epsilon: 3.351427e+00 max(abs(ED)): 2.205069e-15

```
summary(output)
```

```
## A causalglm fit object obtained from causalglmnet for the estimand CATE with formula:
## CATE(W) = 0.849 * (Intercept) + 0.964 * poly(W1, degree = 3, raw = T)1 + 2.04 * poly(W1, degree = 3,
##
## Coefficient estimates and inference:
##     type                          param    tmle_est          se        lower
## 1: CATE                       (Intercept)  0.84854081 0.1622606   0.53051579
## 2: CATE poly(W1, degree = 3, raw = T)1   0.96356900 0.4592845   0.06338783
## 3: CATE poly(W1, degree = 3, raw = T)2   2.04090701 0.4006807   1.25558726
## 4: CATE poly(W1, degree = 3, raw = T)3  -0.03108612 0.7633613  -1.52724668
##        upper    Z_score p_value
## 1: 1.166566 82.6855365 0.00000
## 2: 1.863750 33.1719493 0.00000
## 3: 2.826227 80.5368783 0.00000
## 4: 1.465074  0.6438822 0.51965
```

### msmglm with CATE

msmglm 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 <- msmglm(formula,
    data,
    V = "W1",
    W = c("W1", "W2"), A = "A", Y = "Y",
    estimand = "CATE",
    learning_method = "HAL"
    )
```

## (max) epsilon: -4.094664e-02 max(abs(ED)): 6.199347e-16

```
summary(output)
```

```
## A causalglm fit object obtained from msmglm for the estimand CATE with formula:
## E[CATE(W)|V] = 0.927 * (Intercept) + 1.09 * poly(W1, degree = 3, raw = T)1 + 2.03 * poly(W1, degree =
##
## Coefficient estimates and inference:
##     type                          param   tmle_est          se       lower
## 1: CATE                       (Intercept)  0.9272988 0.1398557   0.6531867
## 2: CATE poly(W1, degree = 3, raw = T)1   1.0882385 0.3790807   0.3452540
## 3: CATE poly(W1, degree = 3, raw = T)2   2.0273469 0.3222212   1.3958049
## 4: CATE poly(W1, degree = 3, raw = T)3  -0.4562005 0.6034215  -1.6388849
```

```
##        upper   Z_score p_value
## 1: 1.2014109 104.83580       0
## 2: 1.8312231  45.39023       0
## 3: 2.6588889  99.48186       0
## 4: 0.7264838  11.95377       0
```

```r
plot_msm(output)
```

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

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

E[CATE(W)|V] = 0.927 * (Intercept) + 1.09 * poly(W1, degree = 3, raw = T)1 + 2.03 * poly(W1, degree = 3, raw = T)2 + −0.45(



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

```
## (max) epsilon: -2.745625e-03 max(abs(ED)): 8.104628e-18
```

```r
summary(output)
```

```
## A causalglm fit object obtained from msmglm for the estimand CATE with formula:
## E[CATE(W)|V] = 1.59 * (Intercept) + 0.863 * W1
##
## Coefficient estimates and inference:
##    type       param  tmle_est       se     lower     upper   Z_score p_value
```

```
## 1: CATE (Intercept) 1.5858907 0.09771712 1.3943686 1.777413 256.60941          0
## 2: CATE         W1 0.8633534 0.18134452 0.5079247 1.218782  75.27559          0
```

```
plot_msm(output)
```

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

E[CATE(W)|V] = 1.59 * (Intercept) + 0.863 * W1



```
# This gives a nonparametric estimate for the marginal ATE
formula <-  ~ 1
output <- msmglm(formula,
    data,
    V = "W1",
    W = c("W1", "W2"), A = "A", Y = "Y",
    estimand = "CATE",
    learning_method = "HAL"
    )
```

```
## (max) epsilon: 4.882896e-04 max(abs(ED)): 6.694645e-17
```

```
summary(output)
```

```
## A causalglm fit object obtained from msmglm for the estimand CATE with formula:
## E[CATE(W)|V] = 1.61 * (Intercept)
##
## Coefficient estimates and inference:
##    type      param tmle_est        se    lower    upper  Z_score p_value
## 1: CATE (Intercept) 1.605885 0.1027628 1.404474 1.807296 247.0863       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: 5.143039e-03 max(abs(ED)): 1.647328e-16
```

```
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand CATE with formula:
## CATE(W) = 0.881 * (Intercept) + 0.969 * poly(W1, degree = 2, raw = TRUE)1 + 2.12 * poly(W1, degree =
##
## Coefficient estimates and inference:
##     type                              param  tmle_est         se     lower
## 1: CATE                          (Intercept) 0.8812678 0.04732328 0.7885158
## 2: CATE poly(W1, degree = 2, raw = TRUE)1 0.9689233 0.06002482 0.8512768
## 3: CATE poly(W1, degree = 2, raw = TRUE)2 2.1221705 0.10444566 1.9174608
##        upper  Z_score p_value
## 1: 0.9740197 294.4442       0
## 2: 1.0865698 255.2281       0
## 3: 2.3268802 321.2624       0
```

```
# CATT, lets reuse fit
output <- npglm(formula,
      output,
      estimand = "CATT")
```

```
## [1] "Reusing previous fit..."
## (max) epsilon: 2.138162e-03 max(abs(ED)): 2.354505e-16
```

```
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand CATT with formula:
## CATT(W) = 0.886 * (Intercept) + 0.977 * poly(W1, degree = 2, raw = TRUE)1 + 2.12 * poly(W1, degree =
##
## Coefficient estimates and inference:
##     type                              param  tmle_est         se     lower
## 1: CATT                          (Intercept) 0.8862398 0.04658125 0.7949422
## 2: CATT poly(W1, degree = 2, raw = TRUE)1 0.9767066 0.06343350 0.8523792
## 3: CATT poly(W1, degree = 2, raw = TRUE)2 2.1160669 0.10433146 1.9115810
##        upper  Z_score p_value
## 1: 0.9775374 300.8224       0
```

```
## 2: 1.1010339 243.4532          0
## 3: 2.3205528 320.6890          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: 1.431745e-02 max(abs(ED)): 9.187096e-17
```

```
summary(outputs[[1]])
```

```
## A causalglm fit object obtained from npglm for the estimand TSM with formula:
## TSM(W) = 0.98 * E[Y_{A=1}]: (Intercept) + 1.31 * E[Y_{A=1}]: poly(W1, degree = 2, raw = TRUE)1 + 1.89
##
## Coefficient estimates and inference:
##     type                                    param  tmle_est          se
## 1:   TSM                     E[Y_{A=1}]: (Intercept) 0.979965 0.08997471
## 2:   TSM E[Y_{A=1}]: poly(W1, degree = 2, raw = TRUE)1 1.308504 0.10570176
## 3:   TSM E[Y_{A=1}]: poly(W1, degree = 2, raw = TRUE)2 1.893816 0.21344042
##        lower    upper  Z_score p_value
## 1: 0.8036179 1.156312 172.2107       0
## 2: 1.1013321 1.515675 195.7324       0
## 3: 1.4754800 2.312151 140.2914       0
```

```
summary(outputs[[2]])
```

```
## A causalglm fit object obtained from npglm for the estimand TSM with formula:
## TSM(W) = 0.098 * E[Y_{A=0}]: (Intercept) + 0.34 * E[Y_{A=0}]: poly(W1, degree = 2, raw = TRUE)1 + -0
##
## Coefficient estimates and inference:
##     type                                    param    tmle_est          se
## 1:   TSM                     E[Y_{A=0}]: (Intercept)  0.09796239 0.09430772
## 2:   TSM E[Y_{A=0}]: poly(W1, degree = 2, raw = TRUE)1  0.34016312 0.11665214
## 3:   TSM E[Y_{A=0}]: poly(W1, degree = 2, raw = TRUE)2 -0.22736046 0.22424548
##           lower      upper  Z_score p_value
## 1: -0.08687735 0.2828021 16.42412       0
## 2:  0.11152913 0.5687971 46.10675       0
## 3: -0.66687353 0.2121526 16.03102       0
```

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: -3.839044e-06 (max) epsilon: 8.765768e-03 max(abs(ED)): 1.727728e-03

```
summary(output)
```

```
## A causalglm fit object obtained from npglm for the estimand OR with formula:
## log OR(W) = 1.15 * (Intercept) + 1.97 * W
##
## Coefficient estimates and inference:
##    type        param tmle_est        se     lower     upper  psi_exp lower_exp
## 1:  OR (Intercept) 1.154000 0.3213434 0.5241784 1.783821 3.170850  1.689070
## 2:  OR           W 1.966543 0.5033627 0.9799700 2.953115 7.145928  2.664376
##    upper_exp  Z_score p_value
## 1:  5.952559 56.78144       0
## 2: 19.165568 61.77210       0
```

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

## risk_change: -3.416995e-05 (max) epsilon: 3.867929e-02 max(abs(ED)): 1.712181e-03

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand OR with formula:
## log OR(W) = 1.21 * (Intercept) + 2.01 * W
##
## Coefficient estimates and inference:
##    type        param tmle_est        se     lower     upper  psi_exp lower_exp
## 1:  OR (Intercept) 1.211124 0.3105147 0.6025263 1.819721 3.357256  1.826728
## 2:  OR           W 2.005836 0.5102673 1.0057308 3.005942 7.432308  2.733904
##    upper_exp  Z_score p_value
## 1:  6.17014 61.67035       0
## 2: 20.20524 62.15381       0
```

```
summary(output)
```

```
## A causalglm fit object obtained from spglm for the estimand OR with formula:
## log OR(W) = 1.21 * (Intercept) + 2.01 * W
##
## Coefficient estimates and inference:
##    type        param tmle_est        se     lower     upper  psi_exp lower_exp
## 1:  OR (Intercept) 1.211124 0.3105147 0.6025263 1.819721 3.357256  1.826728
## 2:  OR           W 2.005836 0.5102673 1.0057308 3.005942 7.432308  2.733904
##    upper_exp  Z_score p_value
## 1:  6.17014 61.67035       0
## 2: 20.20524 62.15381       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) = 0.56 * (Intercept) + 1.44 * poly(W, degree = 2, raw = TRUE)1 + 2.94 * poly(W, degree = 2
##
## Coefficient estimates and inference:
##     type                         param  tmle_est        se     lower
## 1:    RR                     (Intercept) 0.5599643 0.2186846 0.1313503
## 2:    RR poly(W, degree = 2, raw = TRUE)1 1.4365641 0.2271766 0.9913061
## 3:    RR poly(W, degree = 2, raw = TRUE)2 2.9362518 0.4616893 2.0313574
##        upper   psi_exp lower_exp upper_exp   Z_score p_value
## 1: 0.9885784  1.750610  1.140367  2.687411  40.48667       0
## 2: 1.8818222  4.206219  2.694752  6.565457  99.98419       0
## 3: 3.8411462 18.845078  7.624429 46.578831 100.55727       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) = 0.948 * (Intercept) + 1.19 * poly(W, degree = 2, raw = TRUE)1 + 2.13 * poly(W, degree = 2
##
## Coefficient estimates and inference:
##     type                         param  tmle_est        se     lower     upper
## 1:    RR                     (Intercept) 0.9477762 0.1788572 0.5972224 1.298330
## 2:    RR poly(W, degree = 2, raw = TRUE)1 1.1893842 0.2121553 0.7735675 1.605201
## 3:    RR poly(W, degree = 2, raw = TRUE)2 2.1296411 0.5242031 1.1022218 3.157060
##     psi_exp lower_exp upper_exp  Z_score p_value
## 1: 2.579966  1.817065  3.663174 83.78558       0
## 2: 3.285058  2.167485  4.978860 88.64174       0
## 3: 8.411847  3.010848 23.501409 64.23575       0
```

```
output <-
  msmglm(
    formula,
    data,
    V = "W",
```

```
    W = "W", A = "A", Y = "Y",
    estimand = "RR",
    verbose = FALSE
  )
summary(output)
```

```
## A causalglm fit object obtained from msmglm for the estimand RR with formula:
## log E[RR(W)|V] = 0.449 * (Intercept) + 1.48 * poly(W, degree = 2, raw = TRUE)1 + 3.23 * poly(W, degre
##
## Coefficient estimates and inference:
##    type                          param  tmle_est        se       lower
## 1:   RR                     (Intercept) 0.4492428 0.2527816 -0.04620004
## 2:   RR poly(W, degree = 2, raw = TRUE)1 1.4760137 0.2528665  0.98040449
## 3:   RR poly(W, degree = 2, raw = TRUE)2 3.2313368 0.5614794  2.13085738
##        upper   psi_exp lower_exp upper_exp  Z_score p_value
## 1: 0.9446857  1.567125 0.9548509  2.572005 28.09996       0
## 2: 1.9716229  4.375469 2.6655342  7.182324 92.29307       0
## 3: 4.3318162 25.313473 8.4220846 76.082340 90.99518       0
```

## Custom learners with sl3

We refer to the documentation of the tlverse/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 msmglm 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 msmglm and npglm, make sure the learner models interactions, specifically treatment interactions, as these are crucial for fitting the conditional treatment effect estimands well.

```
library(sl3)
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: 5.519004e-02 max(abs(ED)): 3.876899e-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.

# Effects of categorical treatments with npglm and msmglm

For `msmglm` and `npglm`, the CATE, CATT, TSM and RR can be learned for categorical treatments relative to a control treatment. To do this, you need to specify the arguments treatment_level and control_level. The estimands are then user-specified parametric models in $W$ for

$$W \mapsto E[Y|A = a, W] - E[Y|A = 0, W]$$

$$W \mapsto E[Y|A = a, W]$$

$$W \mapsto E[Y|A = a, W]/E[Y|A = 0, W]$$

where $a$ is the specified treatment level.

```
n <- 250
V <- runif(n, min = -1, max = 1)
W <- runif(n, min = -1, max = 1)
A <- rbinom(n, size = 1, prob = 0.66*plogis(W))
A[A==1] <- 2
A[A==0] <- rbinom(n, size = 1, prob = plogis(W))
```

```
## Warning in A[A == 0] <- rbinom(n, size = 1, prob = plogis(W)): number of items
## to replace is not a multiple of replacement length
```

```
table(A)
```

```
## A
##  0  1  2
## 77 85 88
```

```
Y <- rnorm(n, mean = A * (1 + W  ) + W , sd = 0.5)
data <- data.table(W,A,Y)
```

```
output_init <- npglm(~1+W, data, W = "W", A = "A", Y = "Y", estimand = "CATE", learning_method = "mars"
```

```
## (max) epsilon: -8.145000e-03 max(abs(ED)): 2.571034e-17
```

```
summary(output_init)
```

```
## A causalglm fit object obtained from npglm for the estimand CATE with formula:
## CATE(W) = 0.956 * (Intercept) + 0.983 * W
##
## Coefficient estimates and inference:
##    type      param  tmle_est         se     lower     upper   Z_score p_value
## 1: CATE (Intercept) 0.9563525 0.08718246 0.7854780 1.127227 173.44383       0
## 2: CATE           W 0.9833113 0.15563476 0.6782728 1.288350  99.89746       0
```

```
output <- msmglm(~1+W, data, V = "W", W = "W", A = "A", Y = "Y", estimand = "CATE", learning_method = "
```

```
## (max) epsilon: -8.145000e-03 max(abs(ED)): 2.571034e-17
```

```
summary(output)
```

```
## A causalglm fit object obtained from msmglm for the estimand CATE with formula:
## E[CATE(W)|V] = 0.956 * (Intercept) + 0.983 * W
```

```
## 
## Coefficient estimates and inference:
##    type         param tmle_est        se    lower    upper   Z_score p_value
## 1: CATE (Intercept) 0.9563525 0.08718246 0.7854780 1.127227 173.44383       0
## 2: CATE           W 0.9833113 0.15563476 0.6782728 1.288350  99.89746       0
```
```r
# Reuse fits
output <- npglm(~1+W, output_init , estimand = "CATT",   treatment_level = 2, control_level = 0)
```
```
## [1] "Reusing previous fit..."
## (max) epsilon: -3.353597e-02 max(abs(ED)): 1.241436e-16
```
```r
summary(output)
```
```
## A causalglm fit object obtained from npglm for the estimand CATT with formula:
## CATT(W) = 1.97 * (Intercept) + 1.85 * W
## 
## Coefficient estimates and inference:
##    type         param tmle_est        se    lower    upper  Z_score p_value
## 1: CATT (Intercept) 1.967138 0.08502992 1.800483 2.133794 365.7911       0
## 2: CATT           W 1.851514 0.16534631 1.527441 2.175587 177.0527       0
```
```r
output <- npglm(~1+W, output_init , estimand = "TSM",   treatment_level = c(0,1,2))
```
```
## [1] "Reusing previous fit..."
## (max) epsilon: 1.332798e-02 max(abs(ED)): 4.880818e-17
```
```r
lapply(output, summary)
```
```
## A causalglm fit object obtained from npglm for the estimand TSM with formula:
## TSM(W) = 0.0176 * E[Y_{A=0}]: (Intercept) + 1.06 * E[Y_{A=0}]: W
## 
## Coefficient estimates and inference:
##    type               param   tmle_est        se     lower     upper
## 1:  TSM E[Y_{A=0}]: (Intercept) 0.01756881 0.06739586 -0.1145247 0.1496623
## 2:  TSM         E[Y_{A=0}]: W 1.05611065 0.12345656  0.8141402 1.2980811
##       Z_score    p_value
## 1:   4.121725 3.7605e-05
## 2: 135.258716 0.0000e+00
## A causalglm fit object obtained from npglm for the estimand TSM with formula:
## TSM(W) = 0.977 * E[Y_{A=1}]: (Intercept) + 2.04 * E[Y_{A=1}]: W
## 
## Coefficient estimates and inference:
##    type               param  tmle_est        se     lower    upper
## 1:  TSM E[Y_{A=1}]: (Intercept) 0.9766055 0.05490313 0.8689973 1.084214
## 2:  TSM         E[Y_{A=1}]: W 2.0430308 0.09359917 1.8595798 2.226482
##     Z_score p_value
## 1: 281.2497       0
## 2: 345.1222       0
## A causalglm fit object obtained from npglm for the estimand TSM with formula:
## TSM(W) = 1.99 * E[Y_{A=2}]: (Intercept) + 2.98 * E[Y_{A=2}]: W
## 
## Coefficient estimates and inference:
##    type               param tmle_est        se    lower    upper  Z_score
## 1:  TSM E[Y_{A=2}]: (Intercept) 1.987028 0.04649532 1.895899 2.078157 675.7170
## 2:  TSM         E[Y_{A=2}]: W 2.980155 0.07845877 2.826378 3.133931 600.5751
##     p_value
```

```
## 1:        0
## 2:        0
```

```
## $`E[Y_{A=0}]`
##     type                   param   tmle_est        se       lower     upper
## 1:  TSM E[Y_{A=0}]: (Intercept) 0.01756881 0.06739586 -0.1145247 0.1496623
## 2:  TSM          E[Y_{A=0}]: W 1.05611065 0.12345656  0.8141402 1.2980811
##       Z_score    p_value
## 1:   4.121725 3.7605e-05
## 2: 135.258716 0.0000e+00
##
## $`E[Y_{A=1}]`
##     type                   param  tmle_est         se     lower     upper
## 1:  TSM E[Y_{A=1}]: (Intercept) 0.9766055 0.05490313 0.8689973 1.084214
## 2:  TSM          E[Y_{A=1}]: W 2.0430308 0.09359917 1.8595798 2.226482
##      Z_score p_value
## 1: 281.2497       0
## 2: 345.1222       0
##
## $`E[Y_{A=2}]`
##     type                   param tmle_est         se     lower    upper  Z_score
## 1:  TSM E[Y_{A=2}]: (Intercept) 1.987028 0.04649532 1.895899 2.078157 675.7170
## 2:  TSM          E[Y_{A=2}]: W 2.980155 0.07845877 2.826378 3.133931 600.5751
##     p_value
## 1:        0
## 2:        0
##
## $estimand
##     Length     Class      Mode
##          1 character character
##
## $levels_A
##     Min. 1st Qu.  Median    Mean 3rd Qu.     Max.
##      0.0     0.5     1.0     1.0     1.5      2.0
```

```r
n <- 250
V <- runif(n, min = -1, max = 1)
W <- runif(n, min = -1, max = 1)
A <- rbinom(n, size = 1, prob = 0.66*plogis(W))
A[A==1] <- 2
A[A==0] <- rbinom(n, size = 1, prob = plogis(W))
```

```
## Warning in A[A == 0] <- rbinom(n, size = 1, prob = plogis(W)): number of items
## to replace is not a multiple of replacement length
```

```r
table(A)
```

```
## A
##  0  1  2
## 93 74 83
```

```r
Y <- rpois(n, lambda = exp( A * (1 + W)  + sin(5 * W)))
data <- data.table(W,A,Y)


output_init <- npglm(~1+W, data, W = "W", A = "A", Y = "Y", estimand = "RR", learning_method = "gam", t:
```

```
## risk_change: -4.685731e-02 (max) epsilon: 2.499999e-02 max(abs(ED)): 1.354800e+00
```

```
## risk_change: -3.259055e-02 (max) epsilon: 2.499999e-02 max(abs(ED)): 9.589864e-01
## risk_change: -2.127771e-02 (max) epsilon: 2.499999e-02 max(abs(ED)): 5.522427e-01
## risk_change: -7.464972e-03 (max) epsilon: 2.470098e-02 max(abs(ED)): 1.142749e-01
## risk_change: -1.309352e-04 (max) epsilon: 3.331534e-03 max(abs(ED)): 1.691141e-02
## risk_change: -6.660274e-06 (max) epsilon: 4.231507e-04 max(abs(ED)): 4.530236e-03
## risk_change: -7.131212e-07 (max) epsilon: 3.222550e-04 max(abs(ED)): 2.361127e-03
```

**summary**(output_init)

```
## A causalglm fit object obtained from npglm for the estimand RR with formula:
## log RR(W) = 0.897 * (Intercept) + 1.67 * W
##
## Coefficient estimates and inference:
##     type        param tmle_est        se      lower     upper  psi_exp lower_exp
## 1:   RR (Intercept) 0.8970385 0.1770838 0.5499606 1.244116 2.452330  1.733185
## 2:   RR            W 1.6716038 0.5186227 0.6551220 2.688086 5.320694  1.925377
##     upper_exp Z_score p_value
## 1:   3.469867 80.09442       0
## 2: 14.703501 50.96263       0
```

output <- **npglm**(~1+W, output_init , estimand = "RR",  treatment_level = 2, control_level = 0)

```
## [1] "Reusing previous fit..."
## risk_change: -1.613633e-01 (max) epsilon: 2.499999e-02 max(abs(ED)): 9.469931e-01
## risk_change: -3.114208e-02 (max) epsilon: 1.693221e-02 max(abs(ED)): 1.235503e-01
## risk_change: -7.645272e-04 (max) epsilon: 8.478769e-04 max(abs(ED)): 1.168460e-01
## risk_change: -2.572193e-04 (max) epsilon: 1.924575e-03 max(abs(ED)): 4.757747e-02
## risk_change: -4.406177e-05 (max) epsilon: 3.096226e-04 max(abs(ED)): 1.673728e-02
## risk_change: -4.462983e-06 (max) epsilon: 2.476927e-04 max(abs(ED)): 5.928181e-03
## risk_change: -8.914715e-07 (max) epsilon: 4.112759e-05 max(abs(ED)): 2.760402e-03
```

**summary**(output)

```
## A causalglm fit object obtained from npglm for the estimand RR with formula:
## log RR(W) = 1.92 * (Intercept) + 2.07 * W
##
## Coefficient estimates and inference:
##     type        param tmle_est        se      lower     upper  psi_exp lower_exp
## 1:   RR (Intercept) 1.916972 0.1713000 1.5812306 2.252714 6.800339  4.860934
## 2:   RR            W 2.068841 0.6375705 0.8192259 3.318456 7.915644  2.268743
##     upper_exp  Z_score p_value
## 1:   9.513522 176.94105       0
## 2: 27.617683  51.30609       0
```

## Effects of a continuous treatment with contglm

The function `contglm` supports treatment effects for continuous treatments. Currently, the CATE, OR and RR estimands are supported. Specifically, `contglm` computes estimates and nonparametric inference for the best approximation of the true CATE $E[Y|A = a, W] - E[Y|A = 0, W]$ ( for instance ) with respect to the parametric working model $E[Y|A = a, W] - E[Y|A = 0, W] = 1(a > 0) \cdot \beta^T \underline{f}(W) + a \cdot \beta^T \underline{g}(W)$ where $\underline{f}(W)$ and $\underline{g}(W)$ are user-specified parametric models. $\underline{f}(W)$ is specified with the argument `formula\_binary` and captures the treatment effect caused by being treated or not treated $(1(A > 0))$. $\underline{g}(W)$ is specified with the argument `formula_continuous` and captures the treatment effect caused by dosage of continuous effects in the treatment $A$. Note $A$ should be a nonnegative treatment value with $A = 0$ being the placebo group and $A > 0$ being a continuous or ordered numeric dose value.

Thus, unlike other functions, both the argument `formula\_continuous` and `formula\_binary` need to be specified.

For the OR and RR, the models are

$$\log OR(a, W) := \log P(Y = 1|A = a, W)/P(Y = 0|A = a, W) - \log P(Y = 1|A = 0, W)/P(Y = 0|A = 0, W)$$

$$= 1(a > 0) * formula\_binary(W) + a * formula\_continuous(W)$$

and

$$\log RR(a, W) := log E[Y|A = a, W] - \log E[Y|A = 0, W]$$

$$= 1(a > 0) * formula\_binary(W) + a * formula\_continuous(W)$$

```
# CATE
n <- 500
W <- runif(n, min = -1, max = 1)
Abinary <- rbinom(n , size = 1, plogis(W))
A <- rgamma(n, shape = 1, rate = exp(W))
A <- A * Abinary
Y <- rnorm(n, mean =   (A > 0) + A * (1 + W) + W , sd = 0.5)
data <- data.table(W, A, Y)

# Model is CATE(A,W) = formula_binary(W) 1(A > 0) + A * formula_continuous(W)

out <- contglm(
  formula_continuous = ~ 1 + W,
  formula_binary = ~ 1,
  data = data,
  W = "W", A = "A", Y = "Y",
  estimand = "CATE"
)
```

```
## risk_change: -1.438542e-04 (max) epsilon: 4.902379e-03 max(abs(ED)): 5.493324e-03
```

```
summary(out)
```

```
## A causalglm fit object obtained from contglm for the estimand CATE with formula:
## contCATE(A,W) = 0.964 * 1(A>0)*(Intercept) + 1.03 * A*(Intercept) + 1.01 * A*W
##
## Coefficient estimates and inference:
##         type              param    tmle_est         se     lower    upper  Z_score
## 1: contCATE 1(A>0)*(Intercept) 0.9643142 0.05713009 0.8523413 1.076287 377.4319
## 2: contCATE      A*(Intercept) 1.0323259 0.03819698 0.9574612 1.107191 604.3281
## 3: contCATE                A*W 1.0071813 0.04531184 0.9183718 1.095991 497.0282
##    p_value
## 1:       0
## 2:       0
## 3:       0
```

```
# The CATE predictions are now a function of `A`
### head(predict(out))
```

```
# OR
# Model is log OR(a,W) =
# log P(Y=1|A=a,W)/P(Y=0|A=a,W) - log P(Y=1|A=0,W)/P(Y=0|A=0,W)
```

```
# ~ 1(a>0) * formula_binary(W) + a * formula_continuous(W)
n <- 5000
W <- runif(n, min = -1, max = 1)
Abinary <- rbinom(n ,size = 1, plogis(W))
A <- pmin(rgamma(n, shape = 1, rate = exp(W)),1)
A <- A * Abinary
quantile(A)
```

```
##         0%        25%        50%        75%       100%
## 0.00000000 0.00000000 0.00197522 0.56165206 1.00000000
```

```
Y <- rbinom(n, size = 1,  plogis((A>0) + A * (1 + W  ) + W))
data <- data.table(W,A,Y)
out <- contglm(formula_continuous = ~1+W, formula_binary = ~1, estimand = "OR", data =data, W = "W", A =
```

```
## risk_change: -3.536438e-05 (max) epsilon: 2.725785e-03 max(abs(ED)): 8.863804e-02
## risk_change: -1.044470e-05 (max) epsilon: 1.824511e-03 max(abs(ED)): 2.408413e-02
## risk_change: -3.921032e-06 (max) epsilon: 1.014529e-03 max(abs(ED)): 2.850798e-02
```

```
summary(out)
```

```
## A causalglm fit object obtained from contglm for the estimand OR with formula:
## log contCATE(W) = 1.18 * 1(A>0)*(Intercept) + 1.07 * A*(Intercept) + 1.04 * A*W
##
## Coefficient estimates and inference:
##         type              param tmle_est        se    lower    upper  psi_exp
## 1: contCATE 1(A>0)*(Intercept) 1.184363 0.1357195 0.9183580 1.450369 1.184363
## 2: contCATE      A*(Intercept) 1.073405 0.2091382 0.6635021 1.483309 1.073405
## 3: contCATE                A*W 1.043236 0.2618311 0.5300560 1.556415 1.043236
##    lower_exp upper_exp  Z_score p_value
## 1: 0.9183580  1.450369 617.0605       0
## 2: 0.6635021  1.483309 362.9238       0
## 3: 0.5300560  1.556415 281.7385       0
```

```
# The OR predictions are now a function of `A`
#head(predict(out))
```

```
# RR
# Model is log RR(a,W) =
# log E[Y|A=a,W]    - log E[Y|A=0,W]
# ~ 1(a>0) * formula_binary(W) + a * formula_continuous(W)
n <- 1000
W <- runif(n, min = -1, max = 1)
Abinary <- rbinom(n ,size = 1, plogis(W))
A <- pmin(rgamma(n, shape = 1, rate = exp(W)), 1)
A <- A * Abinary
quantile(A)
```

```
##         0%        25%        50%        75%       100%
## 0.00000000 0.00000000 0.03129013 0.58957793 1.00000000
```

```
Y <- rpois(n,  exp((A>0) + A * (1 + W  ) + W))
table(Y)
```

```
## Y
##   0   1   2   3   4   5   6   7   8   9  10  11  12  13  14  15  16  17  18  19
## 233 202 104  77  60  54  37  32  19  16  20  20  14   9  11  12   5   6   6   9
```

```
## 20 21 22 23 24 25 27 28 29 30 31 33 34 35 36 37 38 39 42 43
##  6  4  5  3  3  2  2  1  4  1  2  3  1  2  1  1  2  2  2  1
## 44 46 49 50 56
##  1  2  1  1  1
```

```r
data <- data.table(W,A,Y)
out <- contglm(formula_continuous = ~1+W, formula_binary = ~1, data =data, W = "W", A = "A", Y = "Y",
               estimand = "RR")
```

```
## risk_change: -2.584808e-04 (max) epsilon: 4.999999e-02 max(abs(ED)): 4.088323e-02
## risk_change: -1.406292e-04 (max) epsilon: 4.999999e-02 max(abs(ED)): 3.350878e-02
## risk_change: -7.338708e-05 (max) epsilon: 2.286296e-02 max(abs(ED)): 2.112006e-02
## risk_change: -3.903614e-05 (max) epsilon: 2.653255e-02 max(abs(ED)): 1.716340e-02
## risk_change: -2.032749e-05 (max) epsilon: 1.169375e-02 max(abs(ED)): 1.106802e-02
## risk_change: -1.089249e-05 (max) epsilon: 1.451904e-02 max(abs(ED)): 9.262725e-03
```

```r
summary(out)
```

```
## A causalglm fit object obtained from contglm for the estimand RR with formula:
## log contCATE(W) = 1.08 * 1(A>0)*(Intercept) + 1.03 * A*(Intercept) + 0.766 * A*W
##
## Coefficient estimates and inference:
##        type              param  tmle_est         se     lower     upper
## 1: contCATE 1(A>0)*(Intercept) 1.0777291 0.06294287 0.9543633 1.2010949
## 2: contCATE      A*(Intercept) 1.0277804 0.06444655 0.9014675 1.1540933
## 3: contCATE                A*W 0.7657016 0.11785871 0.5347028 0.9967005
##     psi_exp lower_exp upper_exp  Z_score p_value
## 1: 1.0777291 0.9543633 1.2010949 541.4558       0
## 2: 1.0277804 0.9014675 1.1540933 504.3136       0
## 3: 0.7657016 0.5347028 0.9967005 205.4461       0
```

```r
# The RR predictions are now a function of `A`
#head(predict(out))
```